



# Evaluation Methodology for NLP Tasks



Rohini K. Srihari

Dept. of Computer Science & Engineering  
State University of New York at Buffalo

CEO, Janya Inc.

[www.janyainc.com](http://www.janyainc.com)

# Overview

- Information Extraction
  - InfoXtract
- MUC
- ACE Evaluation Methodology
  - EDR, EMD, EDR with co-reference
  - RDR
  - Events
  - Time Stamping
- Supporting Technology
  - Parsing
  - WSD
- QA (TREC)
- Miscellaneous Issues

Note: Several of the results slides on ACE were taken from the “ACE 2004 NIST Overview” presentation by Mark Przybocki and Audrey Le.

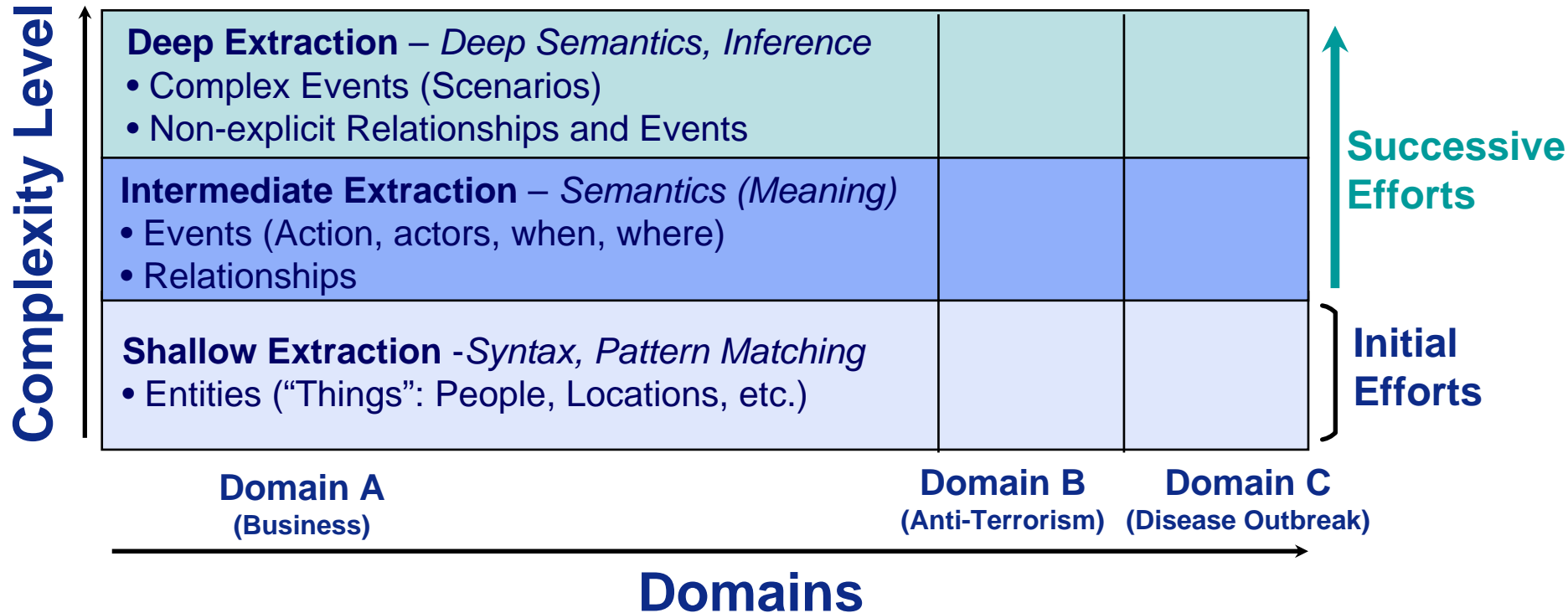
# Motivation

- IE is the most dominant use of NLP technology
  - Counterterrorism
  - Business analytics
  - Biomedical domain
- IE spans virtually every major area of NLP
- Encourage a modular approach to development of IE technology
  - Future NLP tools contests at ICON!

# Information Extraction

- Information extraction is the process of deriving useful, actionable semantics from unstructured text documents
  - Named entity tagging, e.g person/organization names
  - Relationship/Attribute extraction
  - Entity Profiles
    - collection of attributes and relationships associated with an entity
    - Modifiers and descriptors, e.g. “the leading manufacturer of ....”
    - aliases, co-reference
  - Event detection:
    - time stamps, location normalized, numerical unit normalization (e.g. \$ amounts)
    - Permits rich, interactive querying

# Information Extraction: *Domain-Portable Design*



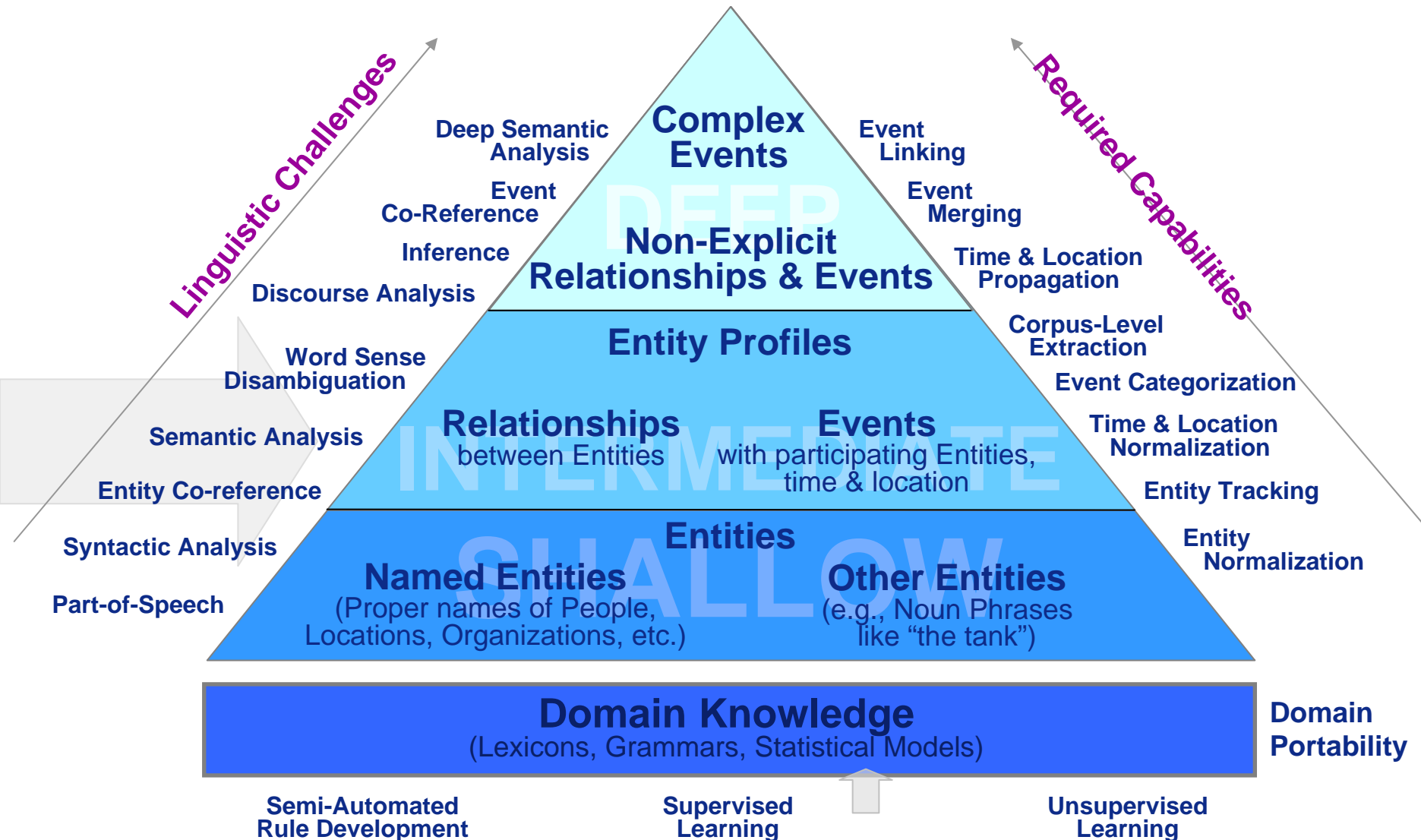
- Establish shallow extraction capable of porting to multiple domains first
- Incrementally evolve more complex extraction capabilities
- Promote plug-and-play architectures to exploit more mature IE technology as it evolves

\* Slide courtesy of Carrie Pine, Air Force Research Laboratories, Rome, NY

# Information Extraction

## What is the State-of-the-Art?

\* Slide courtesy of Carrie Pine, Air Force Research Laboratories, Rome, NY



# InfoXtract

## Generates Entity Profiles & Events from Documents

### Person Profile (EP101)

Profile Name	Herman Cohen
Mentions	official, Mr. Cohen
Aliases	Cohen
Position	assistant secretary of state
Where From	US

---

Assistant secretary of state, **Herman Cohen**,...as the US official currently visiting Khartoum, **he** has been ...

Copyright 1992 Guardian Newspapers Limited  
The Guardian (London)

### Organization Profile (EP402)

Profile Name	al-Barakaat
Descriptors	money transfer terrorist
Staff	Mohamed Barre (EP102)
Founder	Osama bin Laden (EP103)
Events	<i>founding</i> Who: EP103 Org: EP402

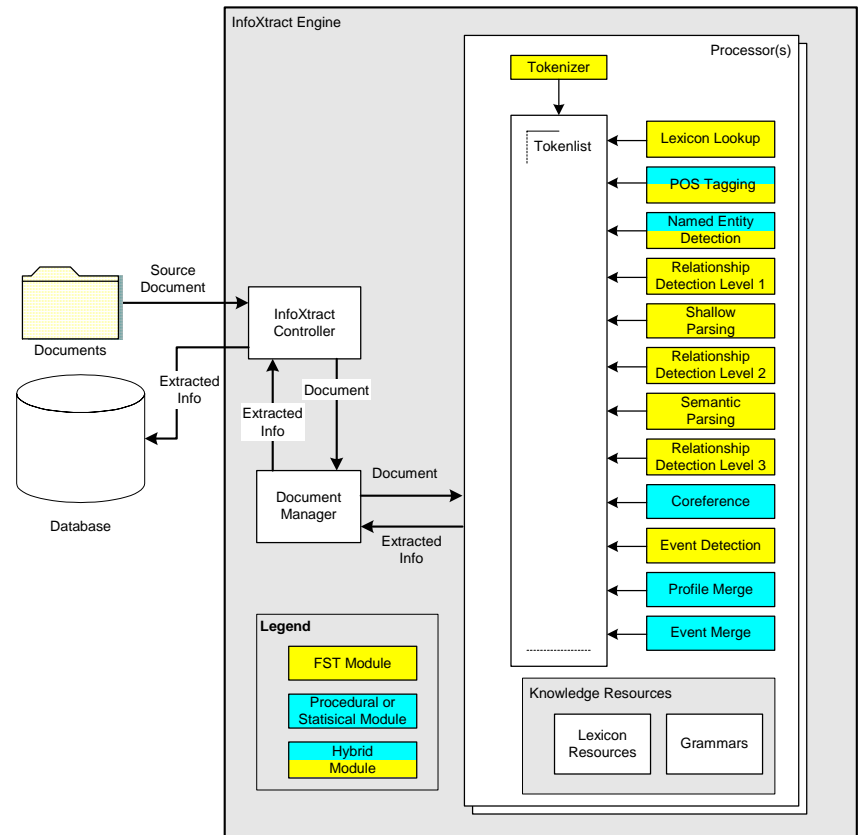
---

According to the Boston Globe, the **al-Barakaat network** was founded by Osama bin Laden...Mohamed Barre has been the **money transfer agency's** broker...

Copyright 2001 Blethen Maine Newspapers, Inc.  
Portland Press Herald

# InfoXtract™ 3.0

- Given text documents, tags
  - Key *entities* (people, places, organizations,...)
  - *Relationships*
  - *Events*
  - Any other linguistic association that customers need, e.g., messaging
  - Summarizes information in *entity profiles*
- Customization
  - Can be customized by the user to support specific domains, tagging needs
  - Supports variety of data sources
    - **FBIS, USMTF, Lexis-Nexis, HUMINT, Factiva, Dialog, etc.**

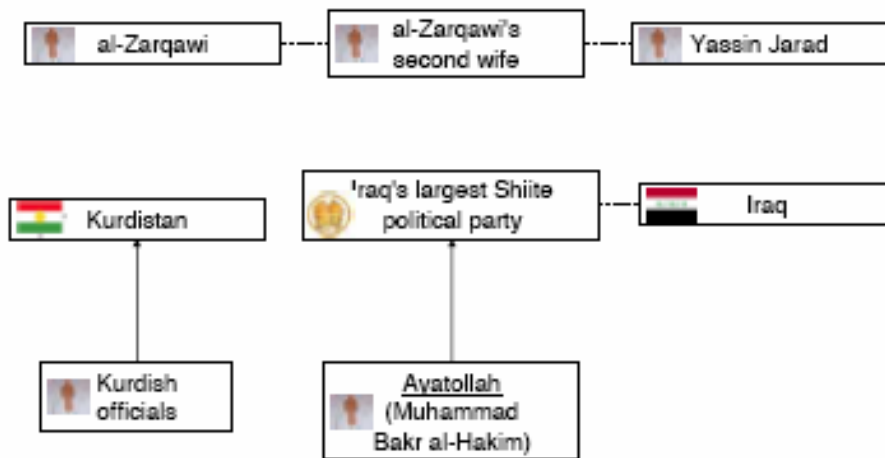


# Use Case 1: Automated Organization Chart Generation from Intelligence

## Message Traffic

End Users: Intelligence analysts at the National Air and Space Intelligence Command (NASIC)

The attack in August 2003 killed more than 85 people, including Ayatollah Muhammad Bakr al-Hakim, who led Iraq's largest Shiite political party. The bombing was carried out with an explosives-laden ambulance driven by Yassin Jarad, the father of al-Zarqawi's second wife, the Kurdish officials said. (*The Seattle Times* -- 08 Feb 2005)



Janya technology used to automatically build and visualize organization charts from HUMINT message traffic.

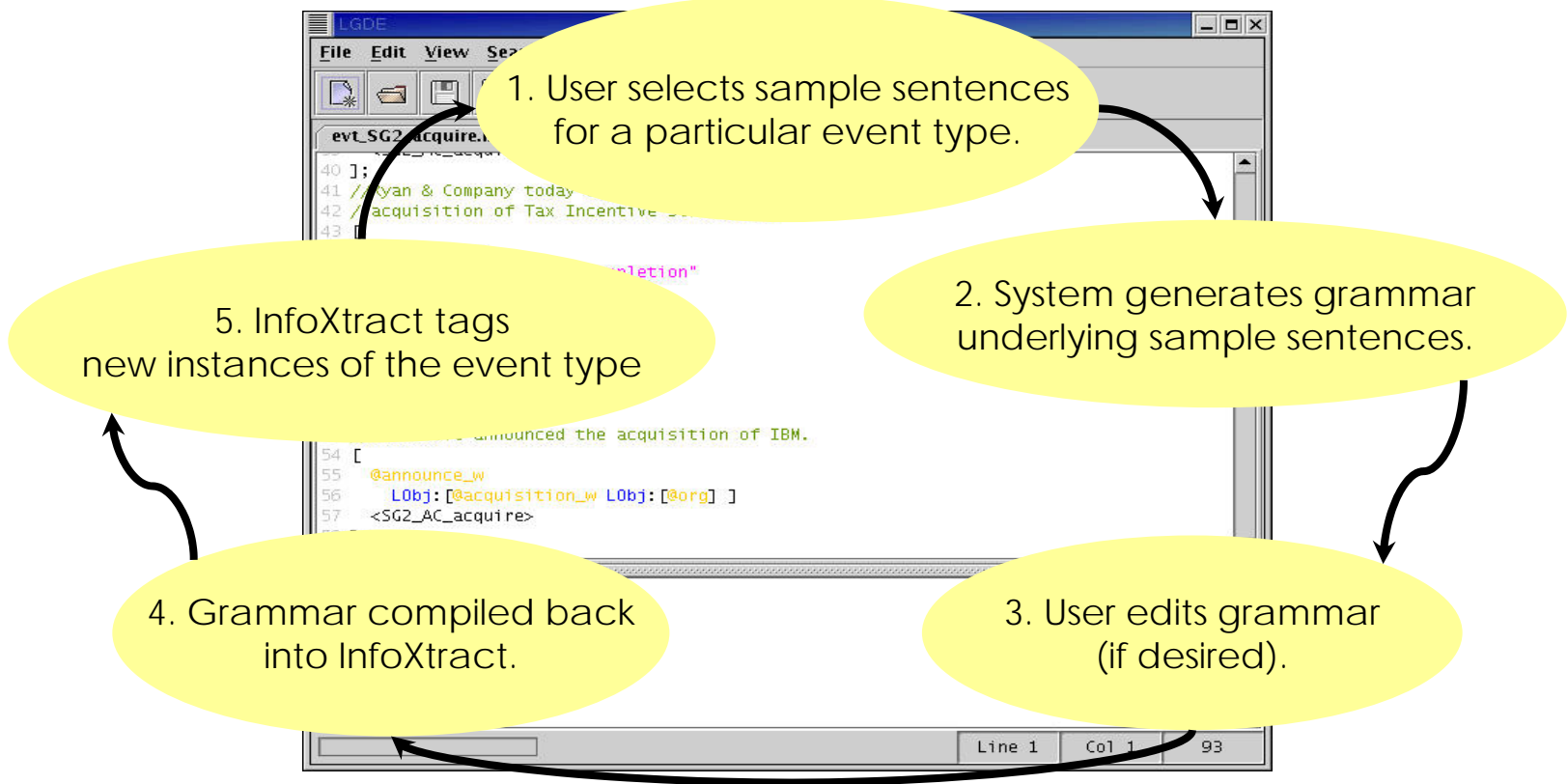
- Funding Agency: Air Force Research Laboratories (AFRL), through the Intermediate Text Exploitation (ITEA) program
- Integrator Partner: General Dynamics Advanced Information Systems

# Use Case 2: Brand Dashboard



- Brand Dashboard™ is powered by InfoXtract™
  - Designed for consumer brand monitoring based on electronic media analysis
  - Marketed by Cymfony Inc.
  - Used by more than 30 Fortune 500 companies, including Pfizer and HP

# InfoXtract Workbench



## Complete IDE for domain customization

- Customization can be performed by non-linguists
- Manual, semi-automatic and automatic tools for domain porting
- Provides a level of customization that is not available in other extraction engines

# Overview

- Information Extraction
  - InfoXtract
- MUC : Message Understanding Conference
- ACE Evaluation Methodology
  - EDR, EMD, EDR with co-reference
  - RDR
  - Events
  - Time Stamping
- Supporting Technology
  - Parsing
  - WSD
- QA (TREC)
- Miscellaneous Issues

# Message Understanding

## • Named entities Conference (MUC) Overview

“The **<ENAMEX TYPE="LOCATION">U.K.</ENAMEX>** satellite television broadcaster said its subscriber base grew **<NUMEX TYPE="PERCENT">17.5 percent </NUMEX>** during **<TIMEX TYPE="DATE">the past year</TIMEX>**.”

## • Equivalence classes

“The **U.K. satellite television broadcaster** said **its** subscriber base grew **17.5 percent** during the past year to **5.35 million.**”

## • Template elements

```
<ENTITY-9602040136-11> :=  
ENT_NAME: "Dennis Gillespie"  
ENT_TYPE: PERSON  
ENT_DESCRIPTOR: "Capt."  
/ "the commander of Carrier Air Wing 11"  
ENT_CATEGORY: PER_MIL
```

# Message Understanding Conference (MUC) Overview (2)

- **Facts**

- relationships between template elements

- <EMPLOYEE\_OF-9602040136-5> :=**

- PERSON: <ENTITY-9602040136-11>**

- ORGANIZATION: <ENTITY-9602040136-1>**

- <ENTITY-9602040136-11> :=**

- ENT\_NAME: "Dennis Gillespie"

- ENT\_TYPE: PERSON

- ENT\_DESCRIPTOR: "Capt."

- / "the commander of Carrier Air Wing 11"

- ENT\_CATEGORY: PER\_MIL

- <ENTITY-9602040136-1> :=**

- ENT\_NAME: "NAVY"

- ENT\_TYPE: ORGANIZATION

- ENT\_CATEGORY: ORG\_GOVT

# Message Understanding Conference (MUC) Overview (3)

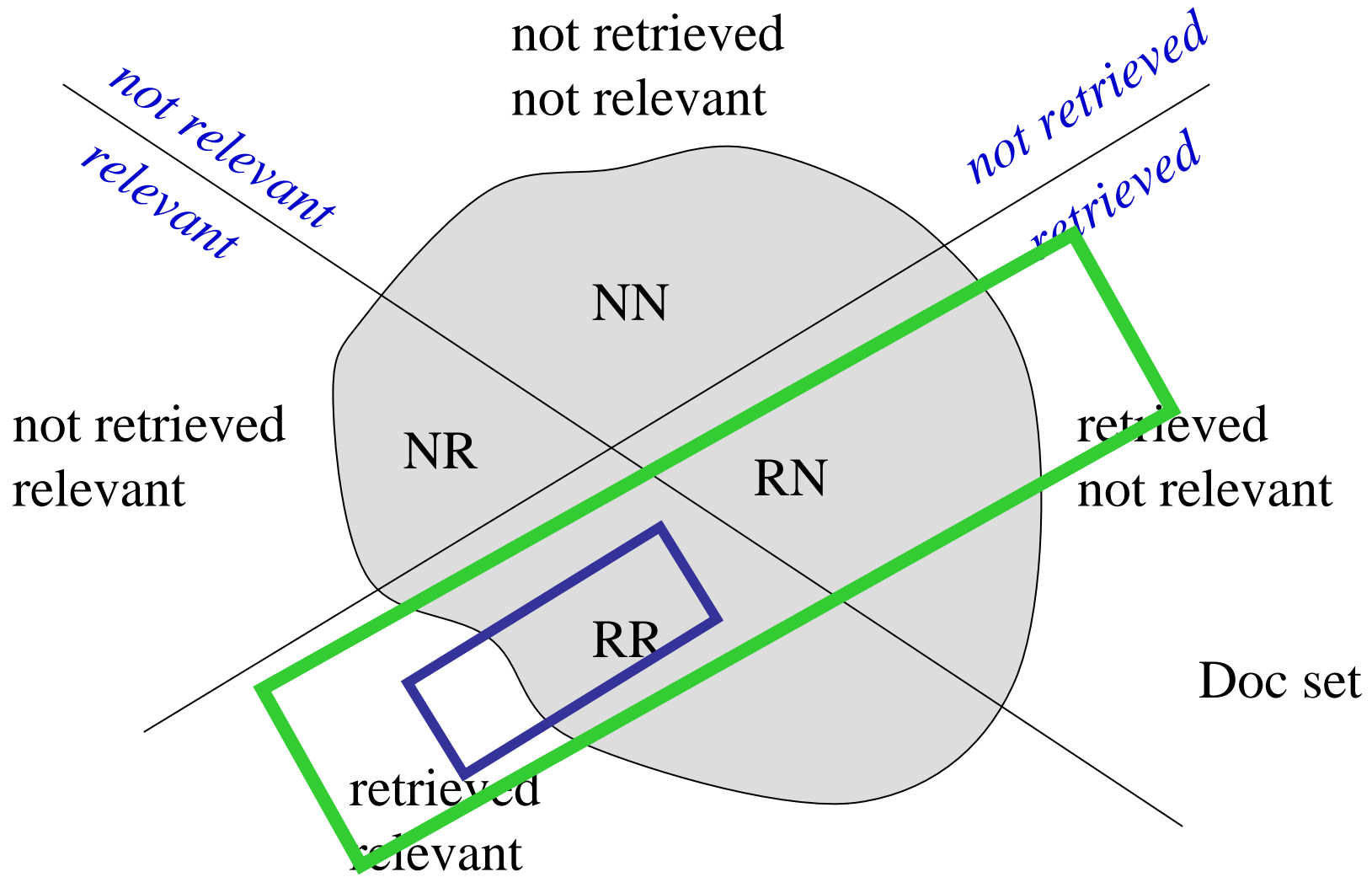
- Scenario Template

- scenario concerns vehicle launch events.
- Template consists of one high-level event object

Chinese rocket carrying a television satellite exploded seconds after launch Wednesday, dealing a potential blow to Rupert Murdoch's ambitions to offer satellite programming in Latin America.

Murdoch's News Corp. is one of four media companies in a partnership that had leased space on the Intelsat satellite to offer the Latin American service. The other partners are Tele-Communications Inc., the nation's largest cable operator; Grupo Televisa SA, the Mexican broadcaster and publisher, and the giant

# Measures based on relevance



# Measures based on relevance

$$\text{recall} = \frac{RR}{RR + NR}$$

$$\text{precision} = \frac{RR}{RR + RN}$$

## **F-measure (harmonic mean)**

$$F(j) = 2 / (1/R(j) + 1/P(j))$$

- Where  $R(j)$  is the recall for  $j$ th doc,  $P(j)$  is precision for  $j$ th doc
- When does  $F$  take on a high value?

# MUC Evaluation/scoring

- NE, Template Elements, Facts, Scenario Template
  - Use the F-Measure as the single measure for systems. The parameters in the F -Measure used are such that recall and precision scores are combined with equal weighting .
  - Scenario Template
    - Systems scored points lower (F-measure) on ST than on Facts.
    - Inter-annotator variability (measured on all articles) was between 85.15 and 96.64 on the F-measures.
- Coreference
  - Objects in the key file are grouped into equivalence classes, as are the objects in the response file. Scores for coreference are based on how well the equivalence classes agree. .

# Overview

- Information Extraction
  - InfoXtract
- MUC : Message Understanding Conference
- ACE Evaluation Methodology
  - EDR, EMD, EDR with co-reference
  - RDR
  - Events
  - Time Stamping
- Supporting Technology
  - Parsing
  - WSD
- QA (TREC)
- Miscellaneous Issues

# ACE: Automated Content Extraction

- Develop technology to extract information from natural language text.
- Primarily interested in
  - Entities
  - Relations
  - Events
  - Time expressions
- Scoring/Evaluation

# ACE Dataset (2004)

Training Data (Oct–Dec 2000)	
English Resources	
BNews	60,291 words
Newsire	59,840 words
Treebank translations Fisher conversations	37,822 words
Arabic Resources	
BNews	63,238 words
Newsire	63,122 words
Treebank	25,010 words
Chinese Resources	
BNews	~67,702 words
Newsire	~60,251 words
Treebank	~25,749 words

Evaluation Data (Jan 2001)	
English Resources	
BNews	25,365 words
Newsire	25,926 words
STT of BNews	~25,000 words
Arabic Resources	
BNews	25,471 words
Newsire	25,056 words
Chinese Resources	
BNews	25,318 words
Newsire	25,379 words

# ACE: Automated Content Extraction

- Develop technology to extract information from natural language text.
- Primarily interested in
  - **Entities**
  - Relations
  - Events
  - Time expressions
- Scoring/Evaluation

# ACE Entity Types

- **PER** Humans, either an individual or a group
- **ORG** Groups defined by an organizational structure
- **VEH** Physical device used for moving an object
- **WEA** Physical device used for harm or destruction
- **GPE** Geographic regions defined by political and/or social groups
- **LOC** Geographic Entities with physical structure
- **FAC** Permanent man-made structures

In decreasing order of evaluation weight

# ACE Entity Subtypes

- **PER** (none)
- **ORG** Government, Commercial, Educational, Non-Profit
- **VEH** Land, Air, Water, Subarea-Vehicle
- **WEA** Blunt, Exploding, Sharp, Chemical, Biological, Shooting, Projectile, Nuclear
- **GPE** Continent, Nation, Sate-or-Province, County-or-District, Population-Center
- **LOC** Address, Boundary, Celestial, Land-Region-Natural, Region-Local, Region-Subnational, Region-National, Region-International, Water-Body
- **FAC** Building, Subarea-Building, Bounded-Area, Conduit, Path, Barrier, Plant,

# ACE Entity Mentions

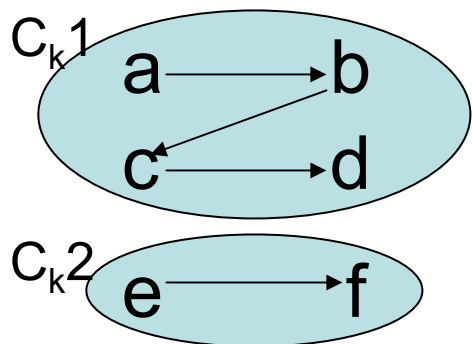
- **Type** Names, Nominals  
Pronouns, Premodifiers
- **GPE Role** PER, LOC, ORG, GPE, FAC
- **Reference Modes (2)** +/-Metonymic

# MUC Scoring[2]

- Most widely used and is more intuitive
- Main Theme:
  - Recall scores based on determining minimum needed to the response classes to map perfectly to key classes
  - Precision scores based on determining minimum needed to the key classes to map perfectly to response classes
- Let key class:  $S_k = \{A, B, C, D\}$ , response class:  $S_r = \{\{A, B\}, \{C\}, \{D\}\}$
- $c(S)$ : #Correct links needed to make key class  $\rightarrow |S_k| - 1$
- $m(S)$ : #Missing links in response class to map perfectly  $\rightarrow |S_r| - 1$
- Thus, Recall = 
$$\frac{c(S) - m(S)}{c(S)}$$
- Similarly calculation for Precision, only switch roles of  $S_r$  and  $S_k$

# MUC Scoring Example

## KEY FILE



Correct Links

$$C_r = \{\{a,b\}, \{c,d\}, \{e,f\}\} , |C_r| = 3$$

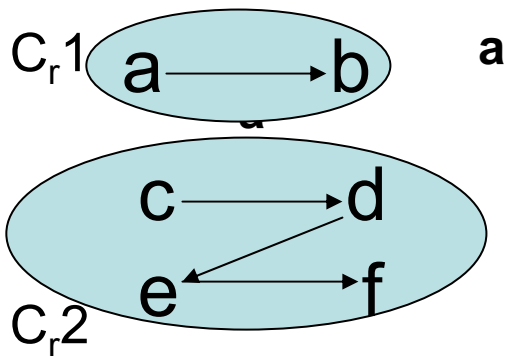
Missing in response

$$M_r = \{b,c\} , |M_r| = 1$$

Missing in key

$$M_k = \{d,e\} , |M_k| = 1$$

## RESPONSE FILE



$$\text{Recall(R): } |C_r| / (|C_r| + |M_k|) = 3 / (3 + 1) = 3/4$$

$$\text{Precision(P): } |C_r| / (|C_r| + |M_k|) = 3 / (3 + 1) = 3/4$$

$$\text{F-Measure: } 2 \times R \times P / (R + P) = 3/4$$

# B-Cubed Scoring[3]

- Tries to eliminate the shortcomings of MUC scoring
  - No credit for separating out singletons from the other co-reference chains
  - Weighting errors. Each entity class given a weight based on its attributes such as Entity Type
- Main Theme:
  - Calculate scores based on Entity Class rather than simply using links
  - Estimate the Spurious links and Missing links in each response class for each element of the response class
- Let Key class:  $S_k = \{A, B, C, D\}$ ,  
Response class:  $S_r = \{\{A, B\}, \{C\}, \{D\}\} \rightarrow \{P_1, P_2, P_3\}$
- $m_j(S)$ : #Missing links in response class partition  $P_j \rightarrow (|S| - |P_j|)$
- $m_j(S)$  calculated for each element of the entity class

# B-Cubed Scoring

- The Recall for a key class  $S_k$  is thus simply:

$$1 - \frac{1}{|S|} \times \sum_j \sum_{e \in P_j} \frac{m_j(S)}{|S|}$$

- Overall Recall of the system is given as weighted recall for individual key classes

- Overall Recall:  $\sum_{i=1}^n w_i \times R_i$

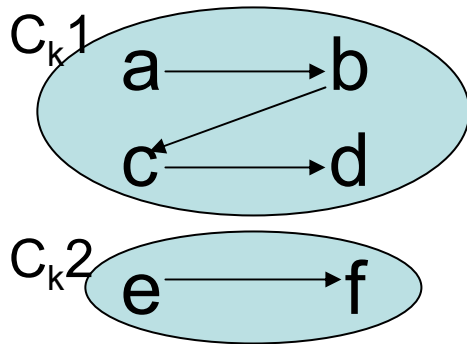
- If the weights are based upon the number of elements in each key class, then Recall:

$$\sum_{i=1}^n \frac{|S_i|}{\sum_{j=1}^n |S_j|} \times R_i$$

- Precision: calculated the same way, only switch roles of  $S_k$  and  $S_r$

# B-Cubed Scoring Example

## KEY FILE



For entity of element 'a':

Part of key class  $\{a,b,c,d\}$

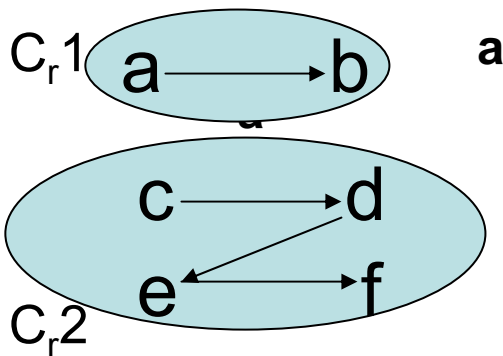
Missing mentions in response entity:  $\{c,d\}$

Recall contribution of 'a' =  $1 - 2/4 = 2/4$

Perform similarly for all elements to get final (un-weighted) Recall:

$$(2/4 + 2/4 + 2/4 + 2/4 + 2/2 + 2/2) * (1/6) \rightarrow 2/3$$

## RESPONSE FILE



Switch roles of key and response classes to get Precision of:

$$(2/2 + 2/2 + 2/4 + 2/4 + 2/4 + 2/4) * (1/6) \rightarrow 2/3$$

F-Measure  $\rightarrow 2/3$

# Co-reference used in InfoXtract

- Aliasing
  - US = America
- Rule-based co-reference
  - [John] nominated [himself]
  - [[President] Bush]
  - [I-90], [the New York State Thruway]
- Statistical-based co-reference

(Features: mention type, number, gender, entity type/subtype string matching, anaphoricity)

- Assad bought a car. He drove to Utica.

# State-of-art Co-reference resolution system

- *IBM's Co-reference resolution based on the Bell Tree* [4]
  - For each mention in document, estimates probability:
    - if it should be linked with one of the existing entities
    - form or create a new entity class
  - Features used:
    - Lexical: {String Matching, Edit Distance, ..}
    - Distance: {Token Distance, Sentence Distance, ..}
    - Syntax: {Mentions in Apposition (Y/N), ..}
    - Count: {Number of times a mentions is seen}
    - Pronoun Features: {Gender Match, Number Match, ..}
  - Feature weights trained using Maximum Entropy Model
- ACE Scores:
  - ACE 2003:            Ace Value: 89.8%            ECM F-measure: 73.2%
  - ACE 2004:            Ace Value: 91.5%
- MUC Scores:
  - MUC 6:            MUC F-measure: 85.7%            ECM F-measure: 76.8%

# ACE: Automated Content Extraction

- Develop technology to extract information from natural language text.
- Primarily interested in
  - Entities
  - **Relations**
  - Events
  - Time expressions
- Scoring/Evaluation

# ACE Relation Types

- **PHYS** Located, Near\*, Part-whole
- **PER-SOC** Business\*, Family\*, Other\*
- **EMP-ORG** Emp-Exec, Employ-Staff, Emp Undet., Member-of-group, Partner\*, Subsidiary
- **ART** User/Owner, Inventor/manufacturer
- **OTHER-AFF** Ethnic, ideology
- **GPE-AFF** Citizen/Resident, Based-in
- **DISC** None (discourse-related)

\* Denotes symmetric relations

# What Do Relations Related Tasks Measure?

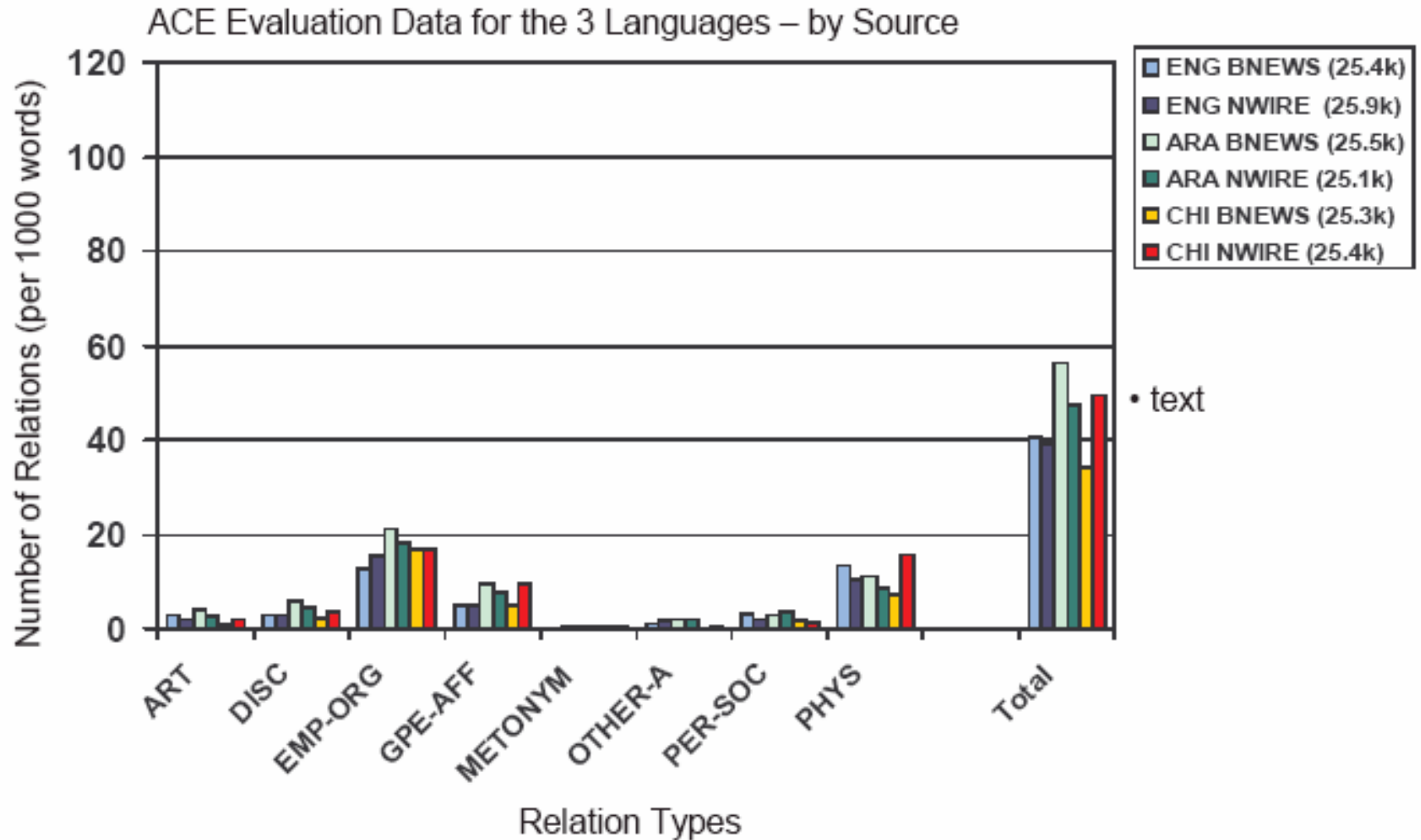
- **Relation Detection and Recognition task (RDR)**
  - detect a **set of specified types of relations**, and
  - recognize **selected information** such as the *type*, *subtype* and *arguments* of each relation.
- **Relation Mention Detection task (RMD)**
  - correctly identify mentions of ACE relations

# Values Assigned To Relations

Relation Types	value weight
ART	1.000
DISC	1.000
EMP-ORG	1.000
GPE-AFF	1.000
METONYMY	1.000
OTHER-AFF	1.000
PER-SOC	1.000
PHYS	1.000

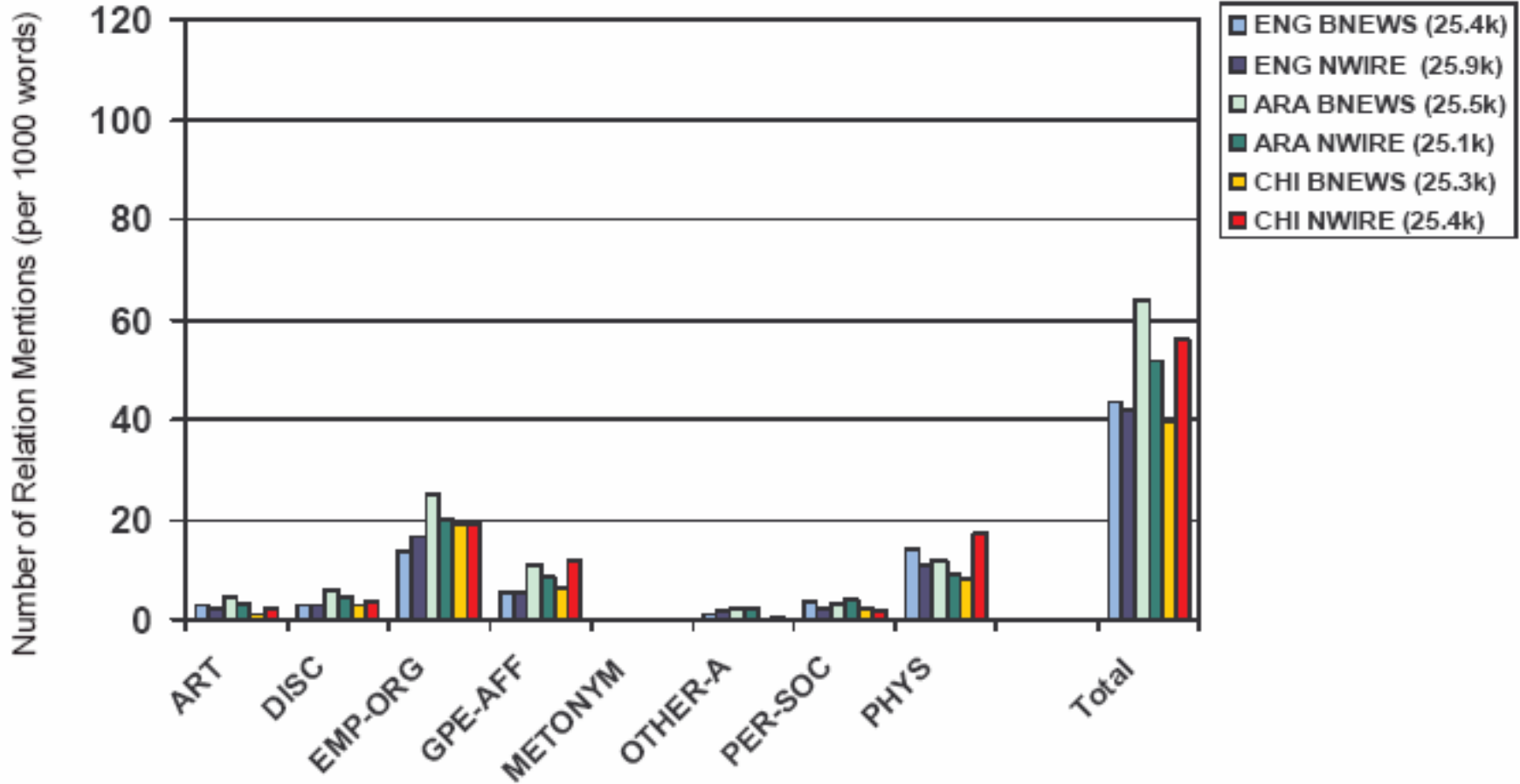
Relation Attribute	value discount
SUBTYPE	0.900
TYPE	0.500
Other Costs	
False Alarm Relation	0.750

# Relation Types In Documents



# Specific Relation Mentions In Documents

ACE Evaluation Data for the 3 Languages – by Source



# RDR Cost Model

The entity evaluation score is the sum of the values of all system ***output relations***

$$\text{RDR\_Value}_{\text{sys}} = \sum_i \text{value\_of\_sys\_entity}_i$$

The overall score of a system is computed as the system output information relative to perfect output

$$\text{system\_value} = \frac{\sum_i \text{value}(\text{sys\_output}_i, \text{reference}_{\text{map}(i)})}{\sum_m \text{value}(\text{reference}_m, \text{reference}_m)}$$

# RDR Cost Model (2)

The value of **each** system **output relation** is the product of an inherent **relation value** and the sum of the values of the **relation's entity arguments**

$$\text{Value}_{\text{sys\_relation}} = \text{Relation\_value}(\text{sys\_entity}) * \sum_a \text{argument\_value}(\text{sys\_argument}_a)$$

# RDR Cost Model (3)

- The ***relation\_value*** of a system output entity is a function of its type
  - if the output entity is mapped, then the minimum value for the system entity and its corresponding reference entity is used
  - if unmapped, it is weighted by a false alarm penalty

$$\text{REL\_VAL} = \left\{ \begin{array}{l} \min \left( \begin{array}{l} R\_TypeVal(\text{sys}), \\ R\_TypeVal(\text{ref}_{\text{sys}}) \end{array} \right) * (W_{R\text{-err-type}} * W_{R\text{-err-subtype}}) \text{ (when mapped)} \\ \\ R\_TypeVal(\text{sys}) * W_{R\text{-FA}} \text{ (when not mapped)} \end{array} \right\}$$

# RDR Cost Model (4)

- The ***argument\_value*** of a system relation argument is the ***entity\_value*** of that entity argument, where the entity argument of the system relation is mapped to the corresponding argument of the reference relation.  
$$\text{Argument\_Value} = \text{Entity\_Value}(\text{sys})$$

# ACE: Automated Content Extraction

- Develop technology to extract information from natural language text.
- Primarily interested in
  - Entities
  - Relations
  - **Events**
  - Time expressions

# What Do Event Related Tasks Measure?

- **Event Detection and Characterization (EDC)**
  - Event types
    - Interaction, Movement, Transfer, Creation, Destruction
  - Event arguments
    - Agent, Object, Source, Target
  - Attributes
    - Temporal, Locative, Instrument, etc.

# Examples of Events

- Quaker Oats rejected PepsiCo's **takeover** offer as too low.
- In April of last year, the CR Company began **bankruptcy** procedures.
- The two were **married** on July 20. **It** was a joyous **event**.

# Event evaluation

- The official ACE metric is a weighted accuracy measure of detected entities called "value score". Events are given different weights determined by a hypothetical application need.
- The value of a system is defined as the weighted value of the detected events compared to the total weighted value of the reference events.
- The "value score" will range from 0-100%.

# ACE: Automated Content Extraction

- Develop technology to extract information from natural language text.
- Primarily interested in
  - Entities
  - Relations
  - Events
  - **Time expressions**

# ISO Standard Time Expressions

## ISO Standard, as used for TIMEX2

Date-Time Strings	$\langle date \rangle \langle time \rangle$ (e.g., 2004-11-04T11:00 or 20041104T1100)
Times	T11:04:33.3-0500 (or T16:04:33.3Z)
Dates (YMD)	2004-11-04 (or 20041104T1100)
Dates (YWD)	2004-W45-4 (or 2004W454)
Dates (YD)	2004-330 (or 2004330)
Periods	$Pn_Y Yn_M Mn_D D$ (e.g., P3Y; PT2H15M)
Intervals	2004-03/2004-08

# TIMEX 2 Tags

- VAL Normalized form of the date/time  
(2004-11-04, P3Y4M, PRESENT\_REF)
- MOD Temporal modifiers (APPROX, EARLY)
- ANCHOR\_VAL Normalized form of an anchoring date/time  
(2004-11-04T13:15)
- ANCHOR\_DIR Relative direction or orientation between  
VAL and ANCHOR\_VAL (AS\_OF)
- SET Expressions denoting sets of time. (YES)

# Sample TIMEX2 Annotation

Assad, 35, has focused on domestic issues since his father, Hafez Assad, died in `<TIMEX2 val="2000-06">June</TIMEX2>` after `<TIMEX2 val="P30Y" anchor val="2000-06" anchor dir="ENDING" comment="context allows more specific anchor val">30 years in power</TIMEX2>`. He has taken steps to modernize his nation's bureaucracy and revive its stagnant economy.

# What Do Entity Related Tasks

## Measure?

- **Entity Detection and Recognition (EDR)**
  - Identify all entities, and
  - recognize selected information: *type*, *subtype*, *class* and *names*
- **Entity Mention Detection (EMD)**
  - Identify all mentions of entities of a certain type
- **EDR with co-reference**
  - link any co-referring mentions together

# Entity Type Example

Person (PER)

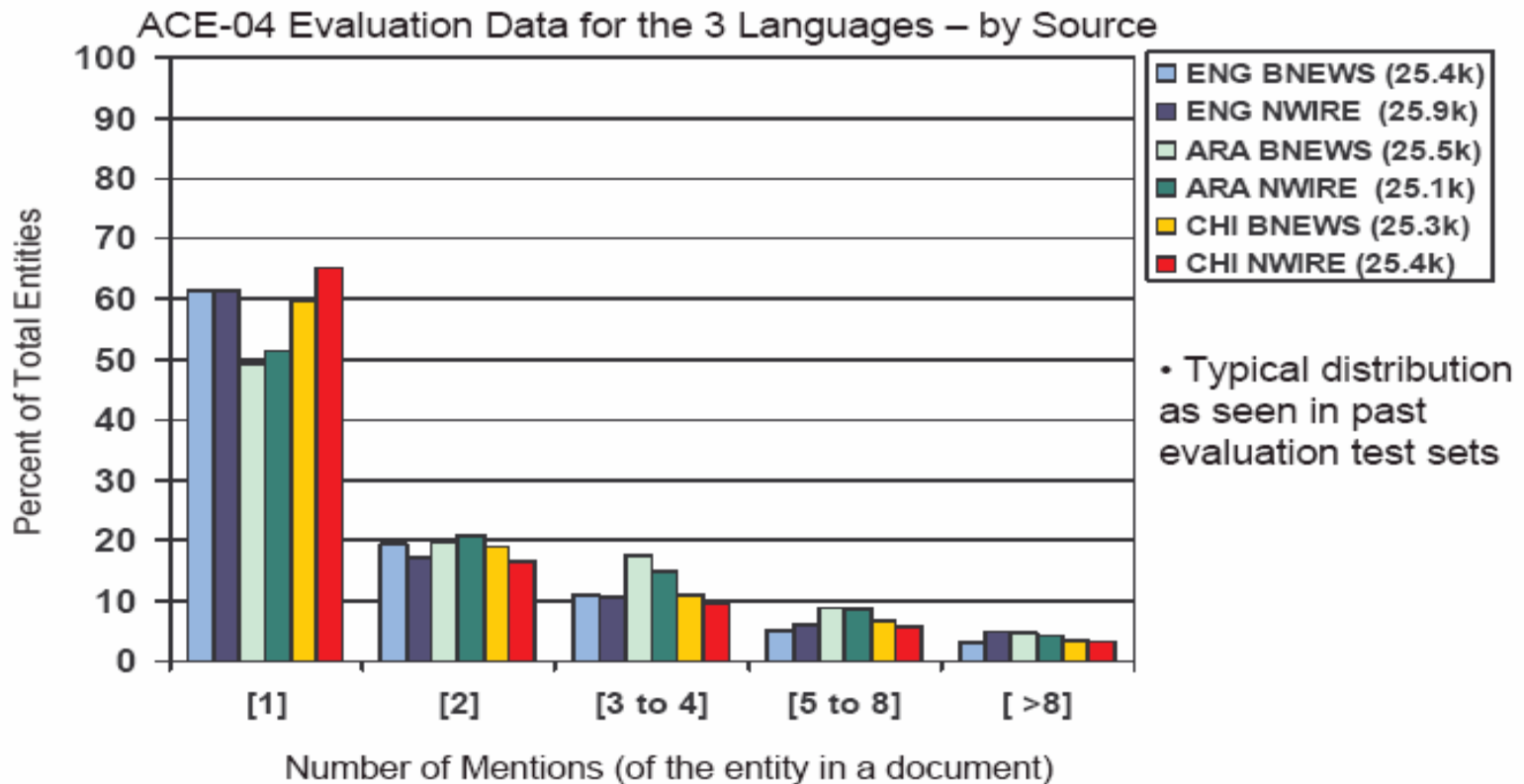
Geo-political entity (GPE)

Location (LOC)

Facility (FAC)

In North Carolina, George Vanderbilt built a house, Biltmore, in the mountains. He later sold the house for a profit. This man builds houses for a living.

# Entities in Documents



# What Do Entity Related Tasks Measure?

- Entity Detection and Recognition (EDR)
  - Identify all entities, and
  - recognize selected information: *type*, *subtype*, *class* and *names*
- **Entity Mention Detection (EMD)**
  - Identify all mentions of entities of a certain type
- EDR with co-reference
  - link any co-referring mentions together

# ACE Entity Mention Example

Proper Name

Nominal

Pronoun

Premodifier

In North Carolina, George Vanderbilt built a house, Biltmore, in the mountains. He later sold the house for a profit. This North American man builds them for a living.

# EMD vs. EDR

## **EMD output:**

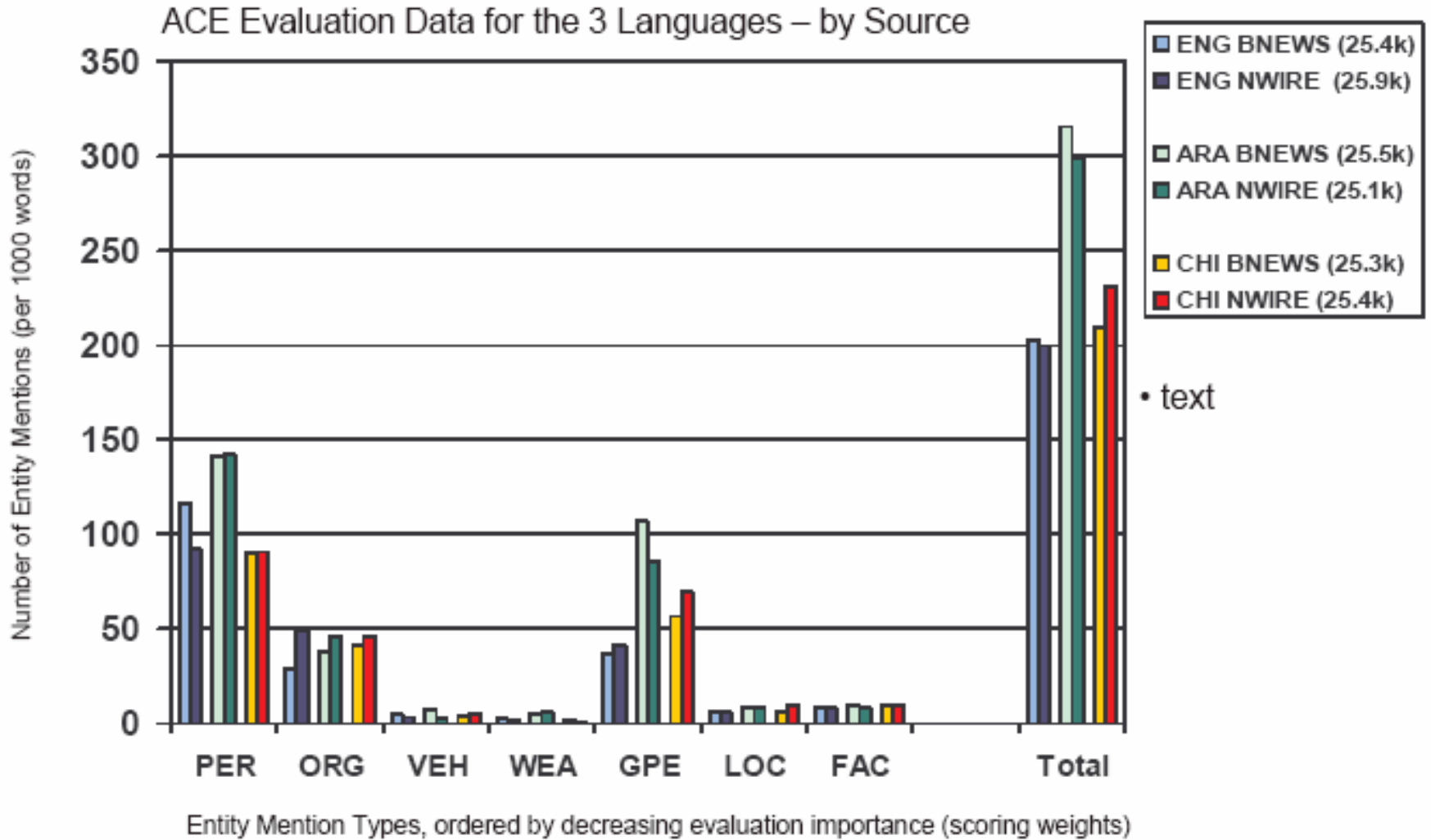
M1 "North Carolina" GPE NAME  
M2 "George Vanderbilt" PER NAME  
M3 "a house" FAC NOM  
M4 "Biltmore" FAC NAME  
M5 "the mountains" LOC NOM  
M6 "He" PER PRO  
M7 "it" FAC PRO  
M8 "This man" PER NOM  
M9 "houses" FAC NOM

## **EDR output:**

E1 GPE: M1  
E2 PER: M2, M6, M8  
E3 FAC: M3, M4, M7  
E4 LOC: M5  
E5 FAC: M9

# Specific Entity Mentions In

ACE Evaluation Data for the 3 Languages – by Source



# ACE Annotation Tool

The screenshot displays the ACE Annotation Tool interface. On the left, a text editor shows XML markup for a news story. On the right, there are three main annotation panels: 'Text', 'Mention', and 'Object'.

**Text Panel:** Shows XML tags: `<DOCID> CNN_ENG_20030630_075848.7 </DOCID>`, `<DOCTYPE SOURCE="broadcast news"> NEWS STORY </DOCID>`, `<DATETIME> 2003-06-30 08:11:12 </DATETIME>`, `<BODY>`, `<TEXT>`, and `<TURN>`.

**Mention Panel:** Contains a table with columns: Text, Coref, Type, Role. A row for 'EM1' has 'sean patrick o'malley' in the Text column, and 'Odd' in the Coref, Type, and Role columns.

	Text	Coref	Type	Role
EM1	sean patrick o'malley	<input type="checkbox"/> Odd	<input type="checkbox"/> Odd	<input type="checkbox"/> Odd

**Object Panel:** Contains a table with columns: Type, Subtype, Class, Gaz., REF, ATR. A row for 'CE1' has 'PER' in Type, 'Individual' in Subtype, 'GPE' in Class, 'Yes!' in Gaz., and 'sean patrick o'malley' in REF.

	Type	Subtype	Class	Gaz.	REF	ATR
CE1	PER	Individual	GPE	Yes!	sean patrick o'malley	

At the bottom, there are dropdown menus for 'PER' (with a list showing Individual, Group, Indeterminate), 'Individual', and 'GEN'. There are also 'Merge' and 'Cross-Typ' buttons.

# ACE Annotation Tool (2)

The screenshot displays the ACE Annotation Tool interface. On the left, a text editor window contains XML metadata and a news snippet. The snippet text is: "pope will reportedly name a new head of the troubled boston cese this week. sean patrick o'malley, the bishop of palm florida, is expected to replace cardinal bernard law. that's ng to the national catholic reporter. law s lepped down in per amid charges of widespread child sexual abuse in the church." The word "church" is underlined in blue. On the right, the "Mention" table shows an annotation for "church" with Type "NAM". Below it, the "Object" table shows two entries: "CE1" for "sean patrick o'malley" (Type: PER, Subtype: Individual, Class: GEN) and "CE2" for "church" (Type: FAC, Subtype: Building-Grounds, Class: GEN). The interface includes a menu bar (File, Edit, Mode, Config, Tools, Window) and a toolbar with icons for file operations.

Text

<DOCID> CNN\_ENG\_20030630\_075848.7 </DOCID>

<DOCTYPE NEWS STORY SOURCE="broadcast news">

<DATETIME> 2003-06-30 08:11:12 </DATETIME>

>

>

>

pope will reportedly name a new head of the troubled boston cese this week. sean patrick o'malley, the bishop of palm florida, is expected to replace cardinal bernard law. that's ng to the national catholic reporter. law s lepped down in per amid charges of widespread child sexual abuse in the church.

>

Mention

Text	Coref	Type	Role	M
church	<input type="checkbox"/> Odd	NAM	<input type="checkbox"/> Odd	<input type="checkbox"/>

NAM <N/A>

Object

	Type	Subtype	Class	Gaz.	REF	ATR	M
CE1	PER	Individual	GEN	<input type="checkbox"/> Yes!	sean patrick o'malley		
CE2	FAC	Building-Grounds	GEN	<input type="checkbox"/> Yes!	church		

FAC Building-Grounds GEN Merge Cross-Type

# Output from the ACE annotation tool

She builds houses in **that town**.

GPE-Population Center, Specific, Nominal, Location-Role

```
<entity ID="E3" TYPE="GPE" CLASS="SPC" SUBTYPE="Population-Center">
  <entity_mention ID="E3m1" TYPE="NOM" ROLE="LOC">
    <extent>
      <charseq START="64" END="72">that town</charseq>
    </extent>
    <head>
      <charseq START="69" END="72">town</charseq>
    </head>
  </entity_mention>
</entity>
```

# Values Assigned To Entities

Entity Mention	value weight
NAM	1.000
NOM	0.200
PRE	0.200
PRO	0.040
Entity Types	value weight
PER	1.000
ORG	0.500
VEH	0.500
WEA	0.500
GPE	0.250
LOC	0.100
FAC	0.050

Entity Classes	value weight
SPECIFIC	1.000
<i>all others</i>	0.000
Entity Attribute	value discount
CLASS	0.750
SUBTYPE	0.900
TYPE	0.500
Mention Attribute	value discount
ROLE	0.900
STYLE	0.900
TYPE	0.900
Other Costs	
False Alarm Entity	0.750
False Alarm Mention	0.750
Discount Incorrect coref	0.000

# ACE Entity Scoring

Each entity's value is the product of its attribute weights and the sum of its mentions' values, each weighted by their attributes. This is primarily the measure of PER/ORG/GPE name detection.

Type: PER >> VEH, ORG, WEA >>  
GPE >> LOC >> FAC

Mention Type: NAM >> NOM, PRE >> PRO

Class: Only specific entities have positive value; generics and negatives worth 0.

# ACE Entity Scoring (2)

- Were all the reference entities correctly recognized?
  - A **MISS\*** occurs whenever a system misses an existing entity e.g., merging two entities into one
- Were all the system output entities valid entities?
  - A **FALSE ALARM\*** occurs whenever a system outputs an entity that doesn't exist e.g., splitting one entity into two.
- Were the valid system output entities correctly recognized?
  - An **ERROR\*** occurs whenever the TYPE, SUBTYPE or CLASS of the system entity does not match that of the reference entity.

# ACE Scoring- Mapping Response File to Key File

- Developed for the Entity Detection and Recognition (EDR) track, of which Co-reference resolution a major task
- Outputs a System\_Value along with Recall, Precision and ECM F-measure values
- Main Theme:
  - Given key files containing the expected co-reference classes and our response files, finds the best mapping from response-classes to key\_file-classes
  - Best mapping → Mapping that maximizes the System\_Value
- $System\_Value = \sum_i (Entity\_Value(sys\_entity_i) \times \sum_j Mention\_Value(sys\_mention_j))$
- *Entity\_Value*: Function of Entity Type
- *Mention\_Value*: Function of Mentions making up the Entity class

# EDR Cost Model

The entity evaluation score is the sum of the values of all system ***output entities***

$$\text{EDR\_Value}_{\text{sys}} = \sum_i \text{value\_of\_sys\_entity}_i$$

The overall score of a system is computed as the system output information relative to perfect output

$$\text{system\_value} = \frac{\sum_i \text{value}(\text{sys\_output}_i, \text{reference}_{\text{map}(i)})}{\sum_m \text{value}(\text{reference}_m, \text{reference}_m)}$$

# EDR Cost Model (2)

The value of **each** system **output entity** is the product of an inherent **entity value** and the sum of the values of the **entity's mentions**

$$\text{Value}_{\text{sys\_entity}} = \text{Entity\_value}(\text{sys\_entity}) * \sum_m \text{Mention\_value}(\text{sys\_mention}_m)$$

# EDR Cost Model (3)

- The ***entity\_value*** of a system output entity is a function of its type

- if the output entity is mapped, then the minimum value for the system entity and its corresponding reference entity is used

– if

$$\text{ENT\_VAL} = \min \left\{ \begin{array}{l} E\_TypeVal(\text{sys}) * \\ E\_ClassVal(\text{sys}), \\ E\_TypeVal(\text{ref}_{\text{sys}}) * \\ E\_ClassVal(\text{ref}_{\text{sys}}) \end{array} \right\} * (W_{E\text{-err-type}} * W_{E\text{-err-subtype}} * W_{E\text{-err-class}})$$

*(when mapped)*

$$E\_TypeVal(\text{sys}) * E\_ClassVal(\text{sys}) * W_{E\text{-FA}}$$

*(when not mapped)*

# EDR Cost Model (4)

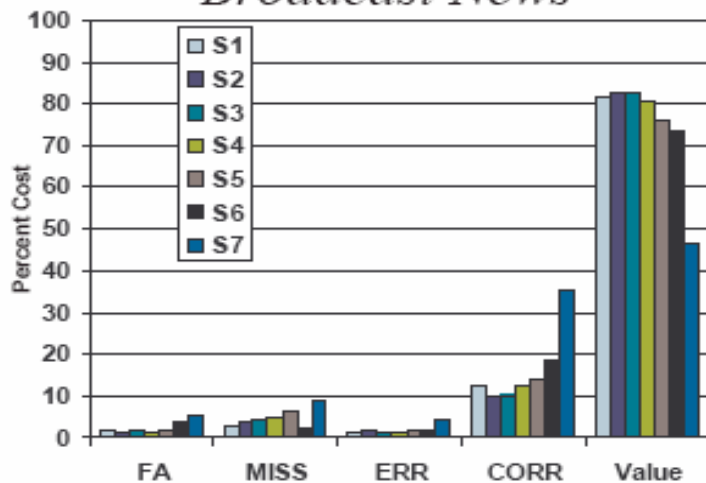
- The ***mention\_value*** of a system entity mention is a function of its type
  - If the mention is mapped, then the minimum value for the sys mention and its corresponding ref mention is used
- If unmapped, it is weighted by a false alarm penalty

$$\text{MEN\_VAL} = \left\{ \begin{array}{l} \min \left( \begin{array}{l} \text{M\_TypeVal}(\text{sys}), \\ \text{M\_TypeVal}(\text{ref}_{\text{sys}}) \end{array} \right) * (W_{\text{M-err-type}} * W_{\text{M-err-role}} * W_{\text{M-err-style}}) \\ \text{M\_TypeVal}(\text{sys}) * (W_{\text{M-FA}} * W_{\text{M-CoRef}}) \end{array} \right. \begin{array}{l} \text{(when mapped)} \\ \text{(when not mapped)} \end{array}$$

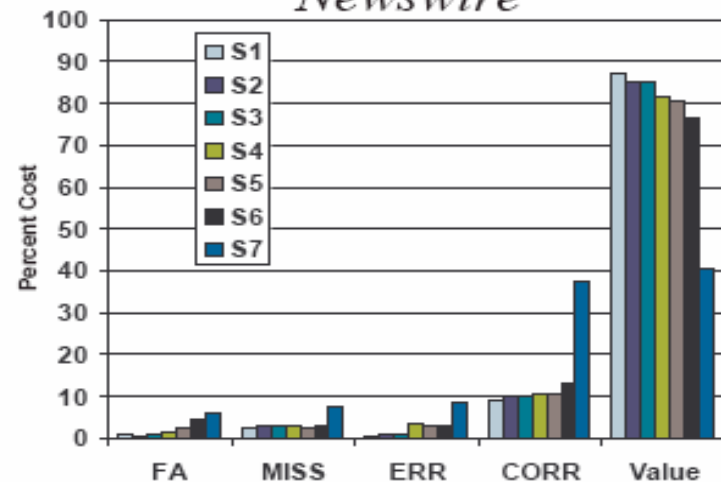
# Performance Of All Systems for Person

PER Value (%)	S1	S2	S3	S4	S5	S6	S7
<b>BNEWS (1222)</b>	81.6	<b>82.8</b>	<b>82.8</b>	80.2	75.8	73.3	46.3
<b>NWIRE (865)</b>	<b>87.4</b>	84.9	84.8	81.3	80.4	76.3	40.5

*Broadcast News*



*Newswire*

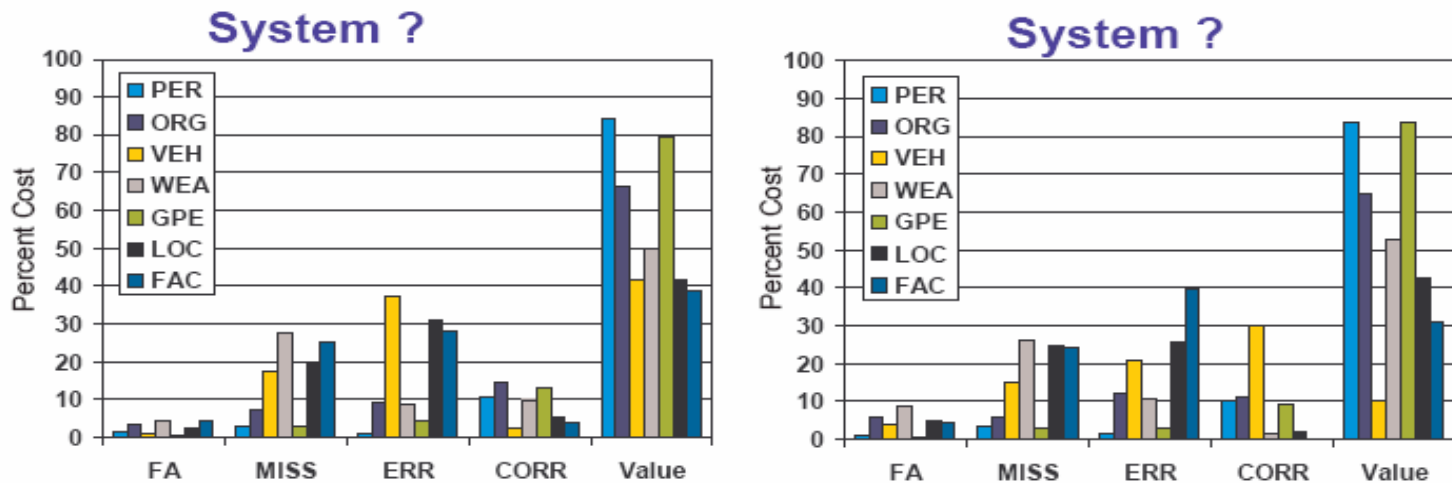


# Sample Performance in EDR

## English EDR

### Percent of Cost by Type

*Combined Sources (Broadcast News & Newswire)*



- Non-zero values in the “corr” column are due to missing or spurious mentions in the system output entity

# Benchmarking Co-reference Systems

Several scoring techniques

- Accuracy Scores

Ratio of Number of Correct links to Number of total links expected

- ACE Scores

Scoring method developed for Automatic Content Extraction program

- MUC Scores

Scoring used by Message Understanding Conferences

Most widely used and more intuitive

- B-Cubed Scores

Attempt to overcome shortcomings of MUC scoring

# Accuracy Scores

- Most straightforward scoring methodology
- For every pair of mentions declared as co-reference check if True\_Positive or False\_Positive
- Accuracy = 
$$\frac{\textit{True\_Positive}}{(\textit{True\_Positive} + \textit{False\_Positive})}$$
- A pair is True\_Positive if placed in the same Co-reference class in key file, else False\_Positive
- Could also be defined as ratio of (number of correctly resolved mentions) to (total number of mentions considered)

# ACE Scoring[1]

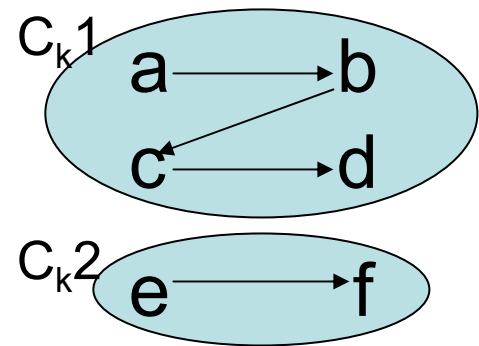
- Developed for the Entity Detection and Recognition (EDR) track, of which Co-reference resolution a major task
- Outputs a System\_Value along with Recall, Precision and ECM F-measure values
- Main Theme:
  - Given key files containing the expected co-reference classes and our response files, finds the best mapping from response-classes to key\_file-classes
  - Best mapping → Mapping that maximizes the System\_Value
- $System\_Value = \sum_i (Entity\_Value(sys\_entity_i) \times \sum_j Mention\_Value(sys\_mention_j))$
- *Entity\_Value*: Function of Entity Type
- *Mention\_Value*: Function of Mentions making up the Entity class

# ACE Scoring

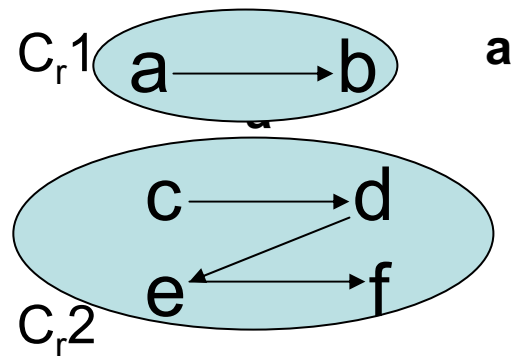
- Performs a brute-force mapping of response classes to key classes based on maximum System\_Value
- After aligning response and key classes, determine:
  - number of mentions in mapped entity class (True Positive)
  - number of spurious mentions in mapped entity class (False Positive)
  - number of missing mentions in mapped entity class (False Negative)
- The F-measure (ECM F-measure) thus obtained measures the percentage of mentions in “right” entities
- F-measure known as (ECM) Entity-Constrained-Mention F-measure since we constrain the mentions to the entities in key files (many-to-one mapping from mentions to entity class in key)

# ACE Scoring Example

## KEY FILE



## RESPONSE FILE



Mapping (Assuming equal weights for all entities):

$$C_{k1} \rightarrow C_{r1} \text{ and } C_{k2} \rightarrow C_{r2}$$

Correct Mentions: (True Positives - TP)

$$\text{Class\#1: } \{a,b\} \rightarrow 2 \quad \text{Class\#2: } \{e, f\} \rightarrow 2$$

Wrong Mentions: (False Positives - FP)

$$\text{Class\#1: } \{\} \rightarrow 0 \quad \text{Class\#2: } \{c,d\} \rightarrow 2$$

Missing Mentions: (False Negatives - FN)

$$\text{Class\#1: } \{c,d\} \rightarrow 2 \quad \text{Class\#2: } \{\} \rightarrow 0$$

$$\text{Recall(R): } TP / (TP + FN) = 4 / (4 + 2) = 2/3$$

$$\text{Precision(P): } TP / (TP + FP) = 4 / (4 + 2) = 2/3$$

$$\text{ECM F-Measure: } 2 \times R \times P / (R + P) = 2/3$$

# Results for Time Expressions

Attribute	pos	act	prec	rec	F	Best P	Best R	Best F
TIMEX2	1828	1382	0.967	0.731	0.832	0.976	0.925	0.926
TEXT	1828	1382	0.901	0.681	0.776	0.901	0.827	0.849
VAL	1329	1336	0.843	0.847	0.845	0.875	0.870	0.872
ANCHOR_DIR	194	73	0.986	0.371	0.539			
ANCHOR_VAL	194	73	0.986	0.371	0.539			
MOD	38	25	0.840	0.553	0.667			
SET	26	0	0.000	0.000	0.000			

# Supporting Technology

# Parsing

- Part of Speech tagging: statistical, rule-based, and lexical lookup
- Shallow parsing of NPs
- Syntactic Parsing
  - Generating parse trees
  - PARSEVAL
    - Uses Penn Treebank
    - Compares brackets in response and key file
- Semantic Parsing
  - Assigning semantic roles
  - PropBank and FrameNet are good resources

# Supporting Technology Through WSD

Word sense disambiguation

– **Interest (noun)**

- Benefit
- Monetary
- Curiosity
- Share

**Senseval:** evaluates the strengths and weaknesses of WSD programs w.r.t. different words and languages.

# Senseval evaluation

- Fine grained
  - assigning probabilities over sense tags output by WSD algorithms.
    - Given a probability distribution over sense tags and a single known correct tag, the algorithm's score should be the probability that the algorithm assigns to the correct tag.
  - Given multiple possible correct tags for a given word token, the algorithm's score should be the sum of ALL probabilities that it assigns to ANY of the correct tags. The premise here is that it is impossible to tell whether a multi-sense annotation was intended as disjunctive or conjunctive, so algorithms should be given the benefit of the doubt.

algorithm's output	score
A	1
.5 A; .5 B	1
.3 A; .7 C	.3
.3 A; .4 B; .3 C	.7
A, B, C	2/3

# Question Answering at TREC

- Question answering competition at TREC consists of answering a set of 500 fact-based questions, e.g.,
  - “*When was Mozart born?*”
- has really pushed the field forward.
- The document set
  - Newswire textual documents from LA Times, San Jose Mercury News, Wall Street Journal, NY Times etcetera: over 1M documents now.
  - Well-formed lexically, syntactically and semantically (were reviewed by professional editors).
- The questions
  - Hundreds of new questions every year, the total is ~2400
- Task
  - Initially extract at most 5 answers: long (250B) and short (50B).
  - Now extract only one exact answer.
  - Several other sub-tasks added later: definition, list, biography.

# TREC Scoring

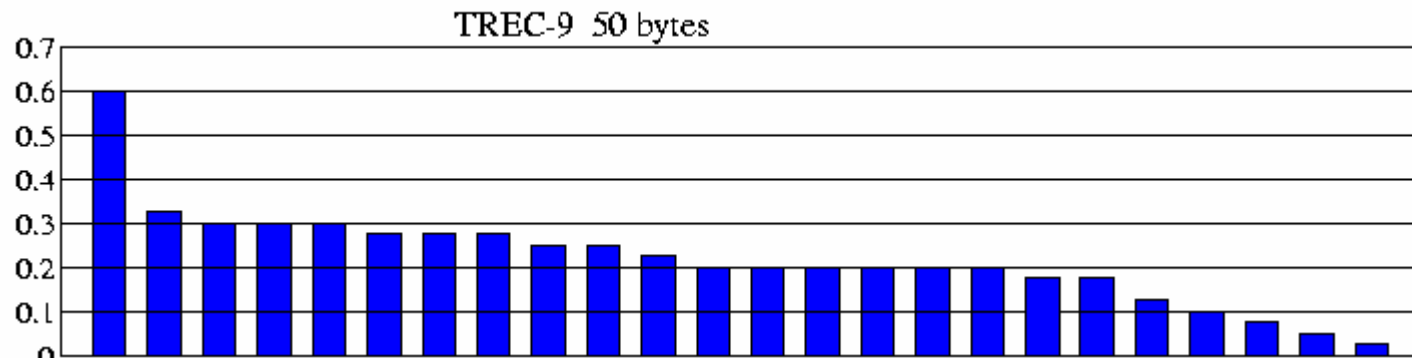
- For the first three years systems were allowed to return 5 ranked answer snippets (50/250 bytes) to each question.
  - Mean Reciprocal Rank Scoring (MRR):
    - Each question assigned the reciprocal rank of the first correct answer. If correct answer at position  $k$ , the score is  $1/k$ .  
1, 0.5, 0.33, 0.25, 0.2, 0 for 1, 2, 3, 4, 5, 6+ position
  - Mainly Named Entity answers (person, place, date, ...)
- From 2002 on, the systems are only allowed to return a single *exact* answer and the notion of confidence has been introduced.

# Sample TREC Questions

1. Who is the author of the book, "The Iron Lady: A Biography of Margaret Thatcher"?
2. What was the monetary value of the Nobel Peace Prize in 1989?
3. What does the Peugeot company manufacture?
4. How much did Mercury spend on advertising in 1993?
5. What is the name of the managing director of Apricot Computer?
6. Why did David Koresh ask the FBI for a word processor?
7. What is the name of the rare neurological disease with symptoms such as: involuntary movements (tics), swearing, and incoherent vocalizations (grunts, shouts, etc.)?

# Top Performing Systems

- In 2003, the best performing systems at TREC can answer approximately 60-70% of the questions
- Approaches and successes have varied a fair deal
  - Knowledge-rich approaches, using a vast array of NLP techniques stole the show in 2000-2003
    - Notably Harabagiu, Moldovan et al. ( SMU/UTD/LCC )
- Statistical systems starting to catch up
  - AskMSR system stressed how much could be achieved by very simple methods with enough text (and now various copycats)
  - People are experimenting with machine learning methods
- Middle ground is to use large collection of surface matching patterns (ISI)



# Benchmarking During Development

- Benchmarking is a key component of the development process
- Regression Testing
  - Many knobs are being tuned simultaneously
  - Important to make sure that nothing breaks due to new grammar rule/feature!
  - Enforce nightly builds
- Benchmarking Accuracy: is the system improving over time?
  - Expensive to manually annotate corpus; only small set can be annotated. May lead to over adaptation.
  - Necessary to benchmark on a much larger, possibly changing corpus

# References

## ACE and MUC

- ACE web site: [www.itl.nist.gov/iad/894.01/tests/ace/](http://www.itl.nist.gov/iad/894.01/tests/ace/)
- MUC:  
[www.itl.nist.gov/iaui/894.02/related\\_projects/muc/proceedings/muc\\_7\\_toc.html](http://www.itl.nist.gov/iaui/894.02/related_projects/muc/proceedings/muc_7_toc.html)
- Przybocki, Mark and Audrey Le. “ACE 2004 NIST Overview”.

## Co-reference

1. The ACE 2004 Evaluation Plan,  
<http://www.nist.gov/speech/tests/ace/ace04/index.htm>
2. A Model-Theoretic Coreference Scoring Scheme, *Marc Vilain, John Burger, John Aberdeen, Dennis Connolly, Lynette Hirschman*
3. Algorithms for scoring coreference chains, *Amit Bagga, Breck Baldwin*
4. A Mention-Synchronous Coreference Resolution Algorithm Based on the Bell Tree, *Xiaoqiang Luo and Abe Ittycheriah, Hongyan Jing and Nanda Kambha Kambhatla and Salim Roukos*

\*\*\* more references forthcoming

*Note: A special thanks to Purnima Shakti of the Linguistics Dept. at SUNY/Buffalo and Cassandre Creswell, Thomas Cornell and Amit Mhatre of Janya Inc. for their help in preparing this material*