

# Question Answering

- Not the classic IR scenario
- Text Retrieval Conference TREC-8

*“People have questions and they need answers, not documents. Automatic question answering will definitely be a major advance in the state-of-art information retrieval technology.”*

- Unstructured document corpus
- Answers as portions of surface text
- Ranking of top 5 answers
- 50- and 250-byte versions

# Approaches used in TREC-8

## Information Retrieval

- Application of existing IR technology to the new problem
- Query processed for keywords  
→ *president, USA*
- Passages returned on statistical relevance measure *tf/idf*
- Some performance with 250-byte passages
- Poor performance with 50-byte passages

# Approaches used in TREC-8

## Query Processing & Named Entity Extraction

- Query processed for keywords and required entity class
  - Who  $\rightarrow$  *PERSON*
  - When  $\rightarrow$  *TIME*
  - How long  $\rightarrow$  *DURATION, LENGTH*
- Passages returned using IR techniques
- Entities within passages identified and assigned to semantic classes
  - Bill Clinton  $\rightarrow$  *PERSON*
- Various heuristics for choice of entity
- Most popular and best performing approach

## Problems revealed by TREC-8

**Q 1** *Who is the president of the USA?*

**A 1.1** *Bill Clinton is the president of the USA.*

**A 1.2** *Even if one disapproves of Bill Clinton, it is important to remember that although many people do not agree with his policies, he is one of the most powerful men on earth: he is the leader of the free world and the president of the USA.*

**A 1.3** *Hillary Clinton is the wife of the president of the USA.*

# Approaches used in TREC-8

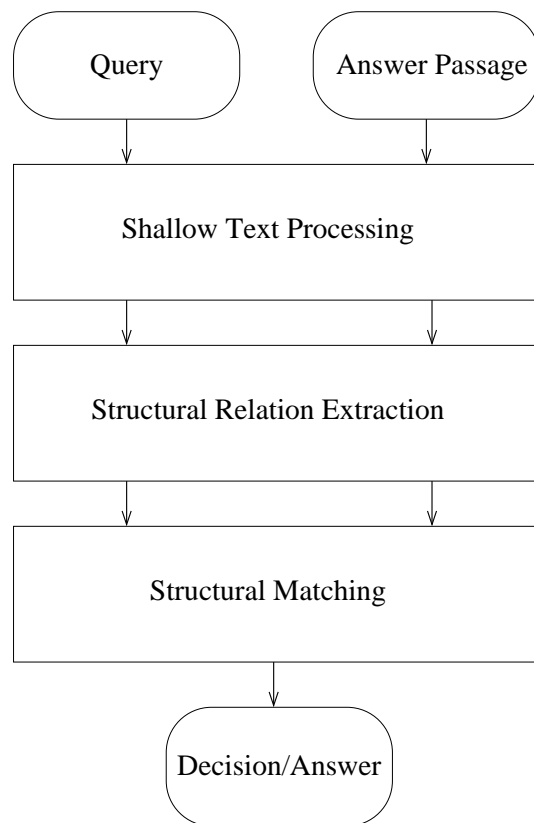
## Sentence Structure

- University of Maryland: dependency parser
  - Dependency tree helps narrow entity choice
  - Not enough to prevent selection of non-answers
- CL Research: semantic triples
  - Semantic roles within sentence
  - Restricted scope
  - Full parser → not robust
  - No coreference resolution

## Aims of Project

- Use sentence structure to select answer entities and prevent
- Requirements:
  - Robust ( $\rightarrow$  no full parse)
  - Select portions of surface text as answer
- Simplifications:
  - Ignored tense, inference, negation
- Assumptions:
  - IR query assumed (only passages containing keywords)
  - NEE assumed (entities added to lexicon)

# System Overview



# Shallow Text Processing

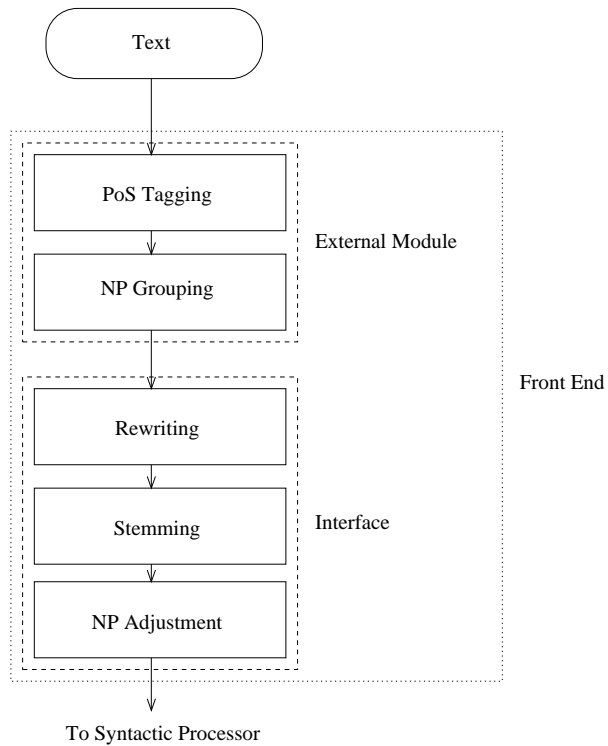


Figure 1: Front End



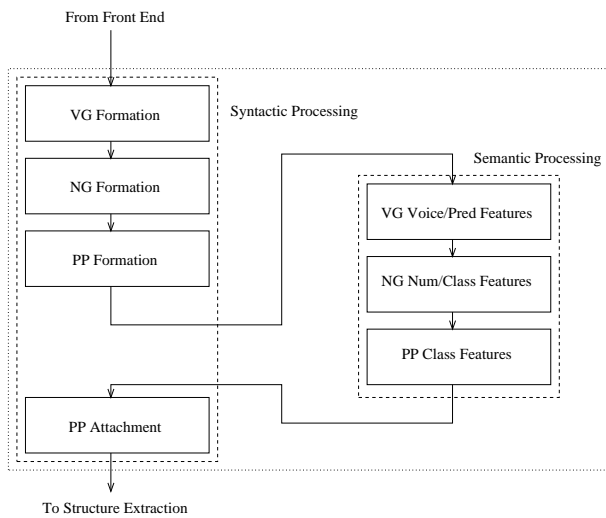


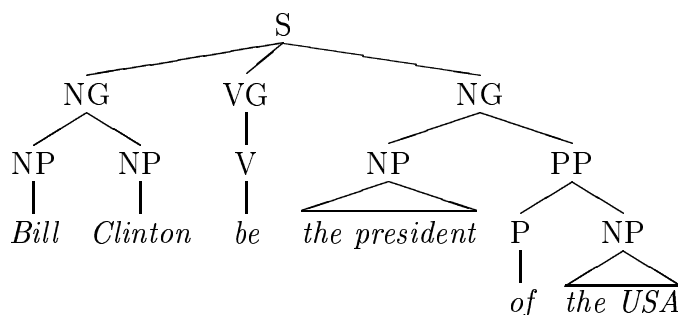
Figure 2: Main Functions

# Syntactic Processing

**Front End** Tagging, stemming:

[np: [bill/nn] , np: [clinton/np] , be/vbz , np: [the  
of/prep , np: [the/det , usa/np] , ' . ' / ' . ' ]

**VG/NG/PP Grouping** Compound NPs, conjunctions,  
aux inversion



# Semantic Processing

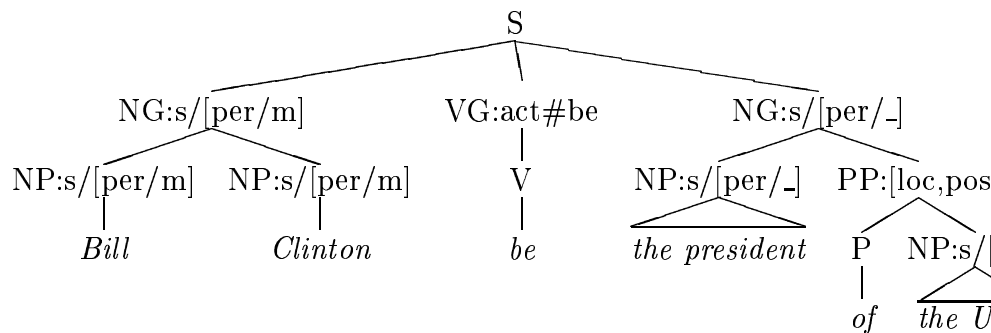
**VG Info** Predicate name, voice

**NG/PP Info** Semantic class, gender, number

Class	Description	Examples
per/f	Person (female)	queen, Mary
per/-	Person (either)	student, candidate
obj	Concrete object	banana, kitchen
abs	Abstract object	election, answer
loc	Location	Wales, Cambridge
org	Organisation	university, Microsoft

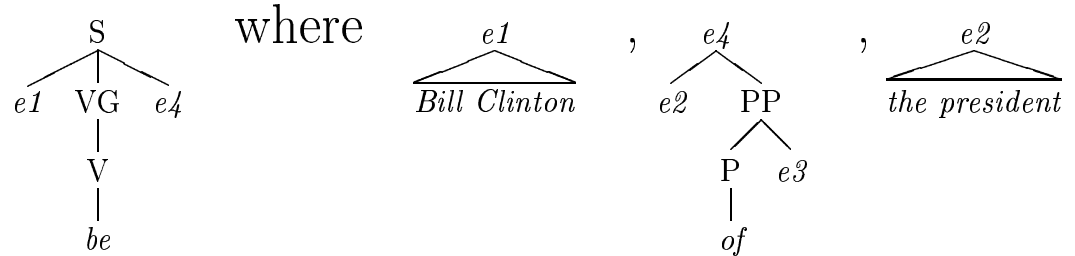
Class	Description	Examples	Re
loc	Location	in, outside	
tim	Time	in, before	
man	Manner	by	
pos	Possession	in, of	
xxx	Unknown	(anything)	

## PP Attachment



# Coreference Resolution

- Entity indexing - provide “pointers”



- Resolution of pronouns → NGs - by semantic class & number
- Resolution of proper/definite NGs & number anaphora —  
ing more info
- Resolution of time & location anaphora
- Non-deterministic

## Syntactic Simplification

- Complex sentences split into list of simple “equivalents”
- Conjunctions, punctuation
  - “*John likes Mary but Mary likes Bill*”  
→ “*John likes Mary*”, “*Mary likes Bill*”
- Subordinate, relative clauses
  - “*Mary, although John likes her, likes Bill*”  
→ “*Although John likes Mary*”, “*Mary likes Bill*”
  - “*Mary, who John likes, likes Bill*”  
→ “*John likes Mary*”, “*Mary likes Bill*”
- Queries & answers require different treatment later

## Structure Extraction

- Predicate-argument structures formed from NG-VG-NG  
s: [ like, john, mary ]  
s: [ like, who, mary ]
- State structures formed if VG is existential, or from NG  
s: [ [bill clinton], [[president], [of, [the  
s: [ who, [[president], [of, [the usa]]] ]

## Lexical Matching

- Query words matched with suitable entities by semantic type:

where	loc	PP
when	tim	PP, NG, NP
who, whom	per, org	NG, NP
what	abs, obj, org	NG, NP
which	abs, obj, org	NG, NP
whose	pos	PP
how	man	PP, AdvP
why	rea	PP

- VGs matched by predicate name
- NGs matched by content (complex NGs require many rules etc.)
- PPs matched by semantic class and NG matching (including stacking)

## Structural Matching

- Set of equivalence rules for tight matching of significant structures

- Active-Passive

$$s : [VG_{pas}, Arg1, pp : [by, Arg2]] \Leftrightarrow s : [VG_{act}, Arg1, pp : [by, Arg2]]$$

- Existential Ordering

$$s : [Arg1, Arg2] \Leftrightarrow s : [Arg2, Arg1]$$

- Verb Nominalisation

$$s : [VG_{util}, Arg1, Arg2, \dots] \Leftrightarrow s : [VG_{lemma}, Arg1', \dots]$$

$$\text{where } Arg1' = Arg1 \setminus NP_{lemma}$$

- PP-Verbs (e.g. possession)

$$s : [VG_{pos}, Arg1, Arg2] \Leftrightarrow s : [Arg2, pp : [pos] : [o]]$$



# Answer Passage Phenomena

Direct Match	<i>"Snowdon is in Wales"</i>
Singular/Plural	<i>"There are several racecourse in Newmarket"</i>
Unnecessary Modifiers	<i>"Snowdon, the highest mountain in the UK, is in Wales"</i>
Stem Matching	<i>"The biggest IT company is Microsoft"</i>
Inclusions	<i>"The Pacific islands like Hawaii produce tropical fruit like b"</i>
Expansions	<i>"In the Suffolk town of Newmarket there is a racecourse"</i>
Simple PPs	<i>"Livingstone is mayor of London"</i>
Possessive PPs	<i>"Ken Livingstone defeated all the other candidates and is Lo"</i>
Verb Form Variations	<i>"Hawaii is producing bananas"</i>
VG Modifiers	<i>"Hawaii produced bananas in 1991"</i>
Verbs of Possession	<i>"Newmarket has a racecourse"</i>
Passives	<i>"One is made very fat by bananas"</i>
Existential Ordering	<i>"In Newmarket there is a racecourse"</i>
Existential Compounds	<i>"Snowdon, in Wales, is a serious mountain"</i>
Coreference	<i>"Pat is in the kitchen and Mike is there too"</i>
Conjunctions	<i>"Pat is in the lounge and Mike is in the garden"</i>
Relative Clauses (internal)	<i>"Hawaii specializes in the production of tropical fruit, which and pineapples"</i>
Relative Clauses (external)	<i>"Pat, who is reading, is in the lounge"</i>
Relative Clauses (coindexed)	<i>"Tropical fruit production, which includes banana production [...] of Hawaii"</i>
Subordinate Clauses	<i>"There is a library, with all the texts the student needs, in Cambridge"</i>
Verb Nominalisation	<i>"In Hawaii, which is not far from California, there is large s bananas"</i>

## Evaluation

- Answers were assessed by a “notional user” with no detailing of the system
- Portions of surface text marked manually and compared to output
- Difficulty of consistency in marking → “narrow-scope” overlap allowed
- Conventional IR performance measures calculated:

$$Recall = n(\text{correct \& identified}) / n(\text{correct})$$

$$Precision = n(\text{correct \& identified}) / n(\text{identified})$$

- 86% recall, 97% precision on training data
- 50-80% recall, >90% precision in blind test (cf. TREC-8 5 answers)

## Summary

- Structural relations used successfully, especially to reject
- Structural equivalences useful in the majority of cases
- Requires complex coding of rules, so large amounts of training (robust to unseen structural phenomena)  
*“Who is the author of the book [...]?” ↔ “[...] by I”*
- What do we do about inference & negation?  
*“about 7 inches higher than 14,776 feet 2 inches” ↔ inches”*