

# CS 294-5: Statistical Natural Language Processing



Question Answering  
Dan Klein, UC Berkeley

(from Chris Manning's slides, which includes slides originally borrowed from Sanda Harabagiu, ISI, Nicholas Kushmerick)

## Assignment 3 Honors



## Project Presentations

- By popular demand: in-class presentations on the last class (Friday 12/10, unless we prefer Wednesday 12/8)
  - You've got 6-8 minutes!
  - Tell us:
    - The problem: why do we care?
    - Your concrete task: input, output, evaluation
    - A simple baseline for the task
    - Your method (half the time here)
    - Any serious surprises, challenges, etc?
    - Headline results (if any)
- Put your slides (if any) on the web before class

## Question Answering from Text

- Question Answering:
  - Give the user a (short) answer to their question, perhaps supported by evidence.
  - An idea originating from the IR community
  - With massive collections of full-text documents, simply finding *relevant documents* is of limited use: we want *answers* from textbases
- The common person's view? [From a novel]
  - "I like the Internet. Really, I do. Any time I need a piece of shareware or I want to find out the weather in Bogota ... I'm the first guy to get the modem humming. But as a source of information, it sucks. You got a billion pieces of data, struggling to be heard and seen and downloaded, and anything I want to know seems to get trampled underfoot in the crowd."
    - M. Marshall. *The Straw Men*. HarperCollins Publishers, 2002.

## People *want* to ask questions?

### Examples from AltaVista query log

who invented surf music?  
how to make stink bombs  
where are the snowdens of yesteryear?  
which english translation of the bible is used in official catholic liturgies?  
how to do clayart  
how to copy psx  
how tall is the sears tower?

### Examples from Excite query log (12/1999)

how can i find someone in texas  
where can i find information on puritan religion?  
what are the 7 wonders of the world  
how can i eliminate stress  
What vacuum cleaner does Consumers Guide recommend

Around 10-15% of query logs

## AskJeeves

- Probably the most hyped example of "question answering"
- It largely does pattern matching to match your question to their own knowledge base of questions
- If that works, you get the human-curated answers to that known question
- If that fails, it falls back to regular web search
- A potentially interested middle ground, but a fairly weak shadow of real QA

## A Brief (Academic) History

- Question answering is not a new research area
- Question answering systems can be found in many areas of NLP research, including:
  - Natural language database systems
    - A lot of early NLP work on these
  - Spoken dialog systems
    - Currently very active and commercially relevant
- The focus on open-domain QA is new
  - MURAX (Kupiec 1993): Encyclopedia answers
  - Hirschman: Reading comprehension tests
  - TREC QA competition: 1999-

## 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?".
- For the first three years systems were allowed to return 5 ranked answer snippets (50/250 bytes) to each question.
  - IR think
  - Mean Reciprocal Rank (MRR) scoring:
    - 1, 0.5, 0.33, 0.25, 0.2, 0 for 1, 2, 3, 4, 5, 6+ doc
  - Mainly Named Entity answers (person, place, date, ...)
- From 2002 the systems are only allowed to return a single *exact* answer and the notion of confidence has been introduced.

## The TREC Document Collection

- The current collection uses news articles from the following sources:
  - AP newswire, 1998-2000
  - New York Times newswire, 1998-2000
  - Xinhua News Agency newswire, 1996-2000
- In total there are 1,033,461 documents in the collection. 3GB of text
- Clearly this is too much text to process entirely using advanced NLP techniques so the systems usually consist of an initial information retrieval phase followed by more advanced processing.
- Many supplement this text with use of the web, and other knowledge bases

## 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 debts did Quintex group leave?
8. What is the name of the rare neurological disease with symptoms such as: involuntary movements (tics), swearing, and incoherent vocalizations (grunts, shouts, etc.)?

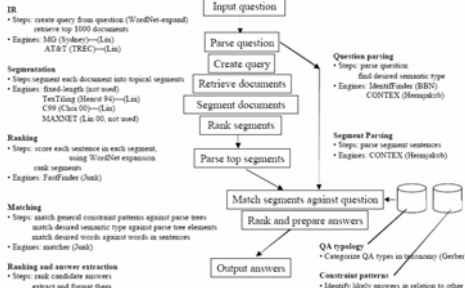
## Top Performing Systems

- Currently the best performing systems at TREC can answer approximately 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, 2001
    - Notably Harabagiu, Moldovan et al. - SMU/UTD/LCC
  - AskMSR system stressed how much could be achieved by very simple methods with enough text (and now various copycats)
  - Middle ground is to use large collection of surface matching patterns (ISI)

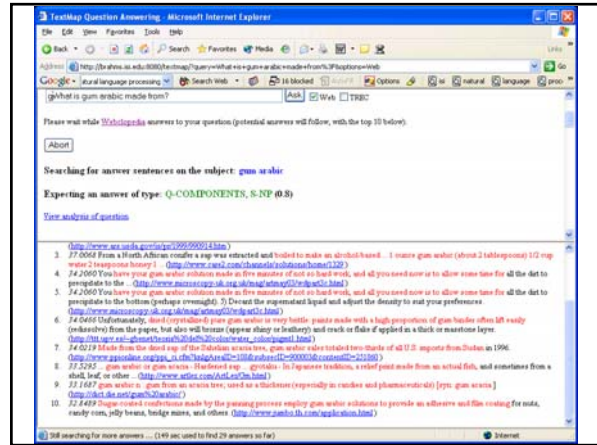
## Online QA System Examples

- Examples
  - **AnswerBus** is an open-domain question answering system: [www.answerbus.com](http://www.answerbus.com)
  - **Ionaut**: <http://www.ionaut.com:8400/>
  - **LCC**: <http://www.languagecomputer.com/>
  - **EasyAsk, AnswerLogic, AnswerFriend, Start, Quasm, Mulder, Webclopedia, etc.**
  - **ISI TextMap**  
<http://brahms.isi.edu:8080/textmap/>

# Webclopedia Architecture

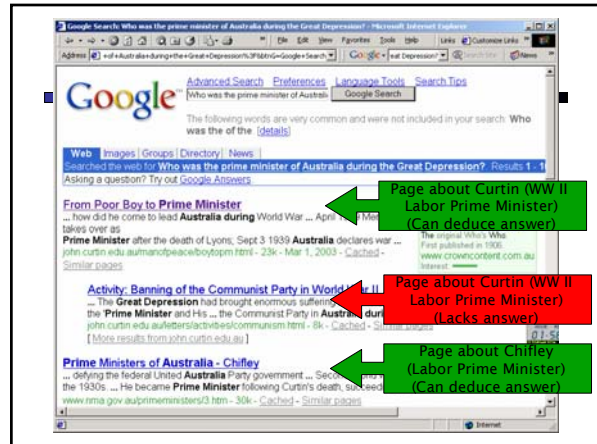


- IR**
- Steps: create query from question (WordNet-expand) return top 1000 documents
  - Engines: MG (Syntex)-(Lis) ATAT (TRAC)-(Lis)
- Segmentation**
- Steps: segment each document into topical segments
  - Engines: Rank-length (just used) TextKing (Hearst 94)-(Lis) C9 (Cass 00)-(Lis) MANNET (Lis 00, not used)
- Ranking**
- Steps: score each sentence in each segment, using WordNet expressions rank segments
  - Engines: Fastfinder (rank)
- Matching**
- Steps: search general concourse patterns against parse trees match desired semantic type against parse tree elements match desired words against words in sentences
  - Engines: searcher (rank)
- Ranking and answer extraction**
- Steps: rank candidate answers extract and format them
  - Engines: part of searcher (rank)



# The Google answer #1

- Include question words etc. in your stop-list
- Do standard IR
- Sometimes this (sort of) works:
- Question: *Who was the prime minister of Australia during the Great Depression?*
- Answer: *James Scullin (Labor) 1929-31.*



## But often it doesn't...

- Question: *How much money did IBM spend on advertising in 2002?*
- Answer: *I dunno, but I'd like to ...* ☹

The screenshot shows a Google search result page. Annotations with arrows point to specific search results:

- A green box labeled "Lot of ads on Google these days!" points to the top of the search results.
- A red box labeled "No relevant info (Marketing firm page)" points to a result from "Advertising metrics".
- A red box labeled "No relevant info (Mag page on ad exec)" points to a result from "Business 2.0 - Magazine Article - Shelly Lazarus".
- A red box labeled "No relevant info (Mag page on MS-IBM)" points to a result from "B2B requires simple marketing materials - 2002-07-20".

## The Google answer #2

- Take the question and try to find it as a string on the web
- Return the next sentence on that web page as the answer
- Works brilliantly if this exact question appears as a FAQ question, etc.
- Works lousily most of the time
- Reminiscent of the line about monkeys and typewriters producing Shakespeare
- But a slightly more sophisticated version of this approach has been revived in recent years with considerable success...

## AskMSR

- Web Question Answering: Is More Always Better?**
  - Dumais, Banko, Brill, Lin, Ng (Microsoft, MIT, Berkeley)

- Q: "Where is the Louvre located?"
- Want "Paris" or "France" or "75058 Paris Cedex 01" or a map
- Don't just want URLs

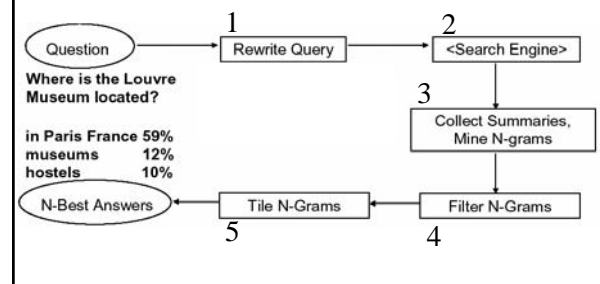
The screenshot shows a Google search for "Where is the Louvre located?". The results include a map, a link to the Louvre Museum website, and several news articles. The text "Paris Cedex 01" is visible in one of the results.

## AskMSR: Shallow approach

- In what year did Abraham Lincoln die?*
- Ignore hard documents and find easy ones

The screenshot shows a search result for "Abraham Lincoln, 1809-1865". It includes a portrait of Lincoln, a snippet of text from a document, and a green box with the text "Sixteenth President of the United States Born in 1809 - Died in 1865".

## AskMSR: Details



## Step 1: Rewrite queries

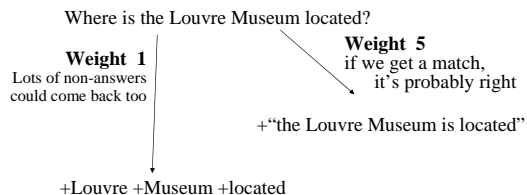
- Intuition: The user's question is often syntactically quite close to sentences that contain the answer
  - Where is the Louvre Museum located?
  - The Louvre Museum is located in *Paris*
  - Who created the character of Scrooge?
  - Charles Dickens* created the character of Scrooge.

## Query Rewriting: Variations

- Classify question into seven categories
    - Who** is/was/are/were...?
    - When** is/did/will/are/were ...?
    - Where** is/are/were ...?
  - a. Category-specific transformation rules
    - eg "For Where questions, move 'is' to all possible locations"
    - "Where is the Louvre Museum located"
    - "is the Louvre Museum located"
    - "the is Louvre Museum located"
    - "the Louvre is Museum located"
    - "the Louvre Museum is located"
    - "the Louvre Museum located is"
  - b. Expected answer "Datatype" (eg, Date, Person, Location, ...)
    - When** was the French Revolution? → DATE
- Nonsense, but who cares? It's only a few more queries to Google.
- Hand-crafted classification/rewrite/datatype rules (Could they be automatically learned?)

## Query Rewriting: Weights

- One wrinkle: Some query rewrites are more reliable than others



## Step 2: Query search engine

- Send all rewrites to a Web search engine
- Retrieve top N answers (100?)
- For speed, rely just on search engine's "snippets", not the full text of the actual document

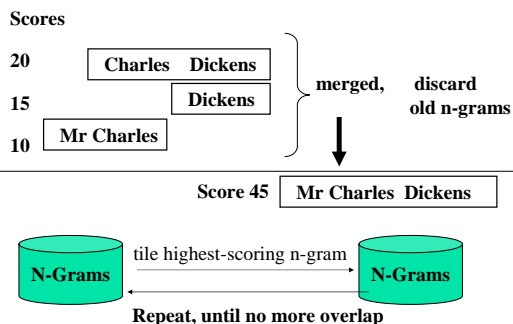
## Step 3: Mining N-Grams

- Simple: Enumerate all N-grams (N=1,2,3 say) in all retrieved snippets
  - Use hash table and other fancy footwork to make this efficient
- Weight of an n-gram: occurrence count, each weighted by "reliability" (weight) of rewrite that fetched the document
- Example: "Who created the character of Scrooge?"
  - Dickens - 117
  - Christmas Carol - 78
  - Charles Dickens - 75
  - Disney - 72
  - Carl Banks - 54
  - A Christmas - 41
  - Christmas Carol - 45
  - Uncle - 31

## Step 4: Filtering N-Grams

- Each question type is associated with one or more "data-type filters" = regular expression
  - When... → Date
  - Where... → Location
  - What ... → Person
  - Who ... → Person
- Boost score of n-grams that do match regexp
- Lower score of n-grams that don't match regexp
- Details omitted from paper....

## Step 5: Tiling the Answers



## Results

- Standard TREC contest test-bed:
  - ~1M documents; 900 questions
- Technique doesn't do too well (though would have placed in top 9 of ~30 participants!)
  - MRR = 0.262 (ie, right answered ranked about #4-#5 on average)
  - Why? Because it relies on the enormity of the Web!
- Using the Web as a whole, not just TREC's 1M documents... MRR = 0.42 (ie, on average, right answer is ranked about #2-#3)

## Issues

- In many scenarios (e.g., monitoring an individuals email...) we only have a small set of documents
- Works best/only for "Trivial Pursuit"-style fact-based questions
- Limited/brittle repertoire of
  - question categories
  - answer data types/filters
  - query rewriting rules

## Ravichandran and Hovy 2002: Learning Surface Patterns

- Use of Characteristic Phrases
- "When was <person> born"
  - Typical answers
    - "Mozart was born in 1756."
    - "Gandhi (1869-1948)..."
  - Suggests phrases like
    - "<NAME> was born in <BIRTHDATE>"
    - "<NAME> ( <BIRTHDATE> -"
  - as Regular Expressions can help locate correct answer

## Use Pattern Learning

- Example: Start with "Mozart 1756"
  - Results:
    - "The great composer Mozart (1756-1791) achieved fame at a young age"
    - "Mozart (1756-1791) was a genius"
    - "The whole world would always be indebted to the great music of Mozart (1756-1791)"
  - Longest matching substring for all 3 sentences is "Mozart (1756-1791)"
  - Suffix tree would extract "Mozart (1756-1791)" as an output, with score of 3
  - Reminiscent of IE pattern learning

## Pattern Learning (cont.)

- Repeat with different examples of same question type
  - "Gandhi 1869", "Newton 1642", etc.
- Some patterns learned for BIRTHDATE
  - a. born in <ANSWER>, <NAME>
  - b. <NAME> was born on <ANSWER> ,
  - c. <NAME> ( <ANSWER> -
  - d. <NAME> ( <ANSWER> - )

## Experiments: (R+H, 2002)

- 6 different Question types
  - from Webclopedia QA Typology (Hovy et al., 2002a)
    - BIRTHDATE
    - LOCATION
    - INVENTOR
    - DISCOVERER
    - DEFINITION
    - WHY-FAMOUS

## Experiments: pattern precision

- BIRTHDATE table:
  - 1.0 <NAME> ( <ANSWER> - )
  - 0.85 <NAME> was born on <ANSWER> ,
  - 0.6 <NAME> was born in <ANSWER>
  - 0.59 <NAME> was born <ANSWER>
  - 0.53 <ANSWER> <NAME> was born
  - 0.50 - <NAME> ( <ANSWER>
  - 0.36 <NAME> ( <ANSWER> -
- INVENTOR
  - 1.0 <ANSWER> invents <NAME>
  - 1.0 the <NAME> was invented by <ANSWER>
  - 1.0 <ANSWER> invented the <NAME> in

## Experiments (cont.)

- WHY-FAMOUS
  - 1.0 <ANSWER> <NAME> called
  - 1.0 laureate <ANSWER> <NAME>
  - 0.71 <NAME> is the <ANSWER> of
- LOCATION
  - 1.0 <ANSWER>'s <NAME>
  - 1.0 regional : <ANSWER> : <NAME>
  - 0.92 near <NAME> in <ANSWER>
- Depending on question type, get high MRR (0.6-0.9), with higher results from use of Web than TREC QA collection

## Shortcomings & Extensions

- Need for POS &/or semantic types
  - "Where are the Rocky Mountains?"
  - "Denver's new airport, topped with white fiberglass cones in imitation of the Rocky Mountains in the background , continues to lie empty"
  - <NAME> in <ANSWER>
- NE tagger &/or ontology could enable system to determine "background" is not a location

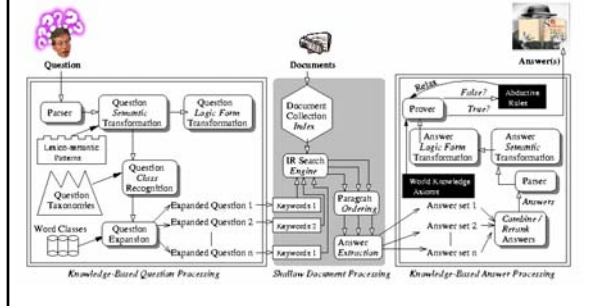
## Shortcomings... (cont.)

- Long distance dependencies
  - "Where is London?"
  - "London, which has one of the most busiest airports in the world, lies on the banks of the river Thames"
  - would require pattern like:  
<QUESTION>, (<any\_word>)\*, lies on <ANSWER>
- But: abundance & variety of Web data helps system to find an instance of patterns w/o losing answers to long distance dependencies

## Shortcomings... (cont.)

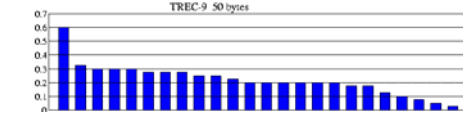
- Their system uses only one anchor word
  - Doesn't work for Q types requiring multiple words from question to be in answer
    - "In which county does the city of Long Beach lie?"
    - "Long Beach is situated in Los Angeles County"
    - required pattern:  
<Q\_TERM\_1> is situated in <ANSWER> <Q\_TERM\_2>
- Does not use case
  - "What is a micron?"
  - "...a spokesman for Micron, a maker of semiconductors, said SIMMs are..."

## LCC: Harabagiu, Moldovan et al.

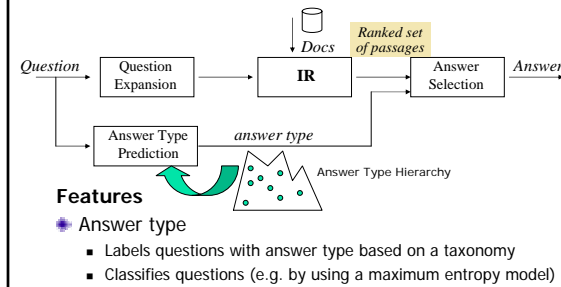


## Value from Sophisticated NLP – Pasca and Harabagiu (2001)

- Good IR is needed: SMART paragraph retrieval
- Large taxonomy of question types and expected answer types is crucial
- Statistical parser used to parse questions and relevant text for answers, and to build KB
- Query expansion loops (morphological, lexical synonyms, and semantic relations) important
- Answer ranking by simple ML method



## Answer types in SOA QA systems



## QA Typology (from ISI USC)

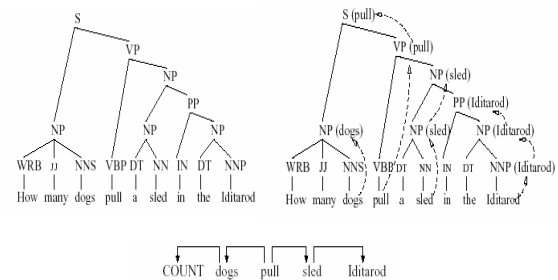
- Typology of typical Q forms—94 nodes (47 leaf nodes)
- Analyzed 17,384 questions (from answers.com)

(STEM)	(SPATIAL-QUANTITY
((AGENT	(VOLUME-QUANTITY AREA-QUANTITY DISTANCE-QUANTITY))
((NAME (PERSON-FIRST-NAME (EYE MARY ...))	(PERCENTAGE))
((NAME-FIRST-NAME (LAWRENCE SAM ...)))	(UNIT
(COMPANY-NAME (BOJING AMERICAN-EXPRESS))	((INFORMATION-UNIT (BIT BYTE ... EXABYTE))
(JESUS ROMANOFF ...)	(MASS-UNIT (OUNCE ...)) (ENERGY-UNIT (BTU ...))
(ANIMAL-NOUN (ANIMAL (WOODCHUCK YAK ...))	(CURRENT-UNIT (KILOTT PICO ...))
(PERSON	(TEMPORAL-UNIT (ATTORSECOND ... MILLISECOND))
(ORGANIZATION (SQUADRON DICTATORSHIP ...))	(TEMPERATURE-UNIT (FAHRENHEIT CELSIUS))
(GROUP-OF-PEOPLE (FISHES CRUIS ...))	(ILLUMINATION-UNIT (LUX CANDLEA))
(STATE-DISTRICT (VIRGO MISSISSIPPE ...))	(SPATIAL-UNIT
(CITY (OLAN-BATHY VIRGINIA ...))	((VOLUME-UNIT (DECILITER ...))
(COUNTRY (SUDANESE EINMANNE ...)))	(AREA-UNIT (ACRE)) ... PERCENT))
(PLACE	(TEMPERER-QUANTITY
(STATE-DISTRICT (CITY COUNTY ...))	((FOOD (HUMAN-FOOD (FISH CHEESE ...))
(OROLOGICAL-FORMATION (STAR CANTON ...))	(SUBSTANCE
AIRPORT COLLEGE CAPITOL ...)	((LIQUID (LIMONADE MAGLIUM BLEAD ...))
(ABSTRACT	(GAS-FORM-SUBSTANCE (GAS AIR)) ...))
(LANGUAGE (LETTER-CHARACTER (A B ...)))	(INSTUMENT (DRUM DRILL GUNDRON (ARM GUN) ...))
(QUANTITY	(BODY-PART (ARM HEART ...))
(NUMERICAL-QUANTITY INFORMATION-QUANTITY	(MEDICAL-INSTRUMENT (PIANO))
(MASS-QUANTITY MORTALITY-QUANTITY	... (GOVERNMENT (PLANT (DISEASE))
(TEMPORAL-QUANTITY ENERGY-QUANTITY	
(TEMPERATURE-QUANTITY ILLUMINATION-QUANTITY	

## Named Entity Recognition for QA

- The results of the past 5 TREC evaluations of QA systems indicate that current state-of-the-art QA is determined by the recognition of Named Entities:
  - Precision of recognition
  - Coverage of name classes
  - Mapping into concept hierarchies
  - Participation into semantic relations (e.g. predicate-argument structures or frame semantics)

## Syntax to Logical Forms



- Syntactic analysis plus semantic => logical form
- Mapping of question and potential answer LFs to find the best match



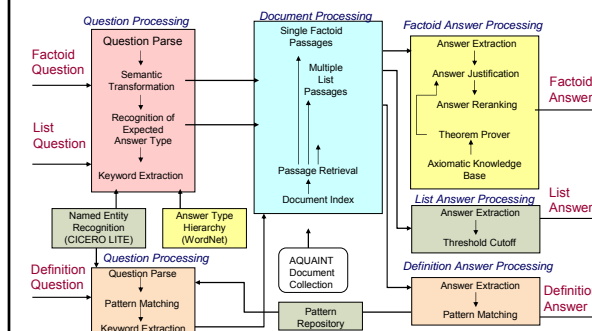
## Abductive inference

- System attempts inference to justify an answer (often following lexical chains)
- Their inference is a kind of funny middle ground between logic and pattern matching
- But quite effective: 30% improvement
- Q: *When was the internal combustion engine invented?*
- A: *The first internal-combustion engine was built in 1867.*
- invent -> create\_mentally -> create -> build

## Question Answering Example

- How hot does the inside of an active volcano get?
- get(TEMPERATURE, inside(volcano(active)))
- "lava fragments belched out of the mountain were as hot as 300 degrees Fahrenheit"
- fragments(lava, TEMPERATURE(degrees(300)), belched(out, mountain))
  - volcano ISA mountain
  - lava ISPARTOF volcano
  - lava inside volcano
  - fragments of lava HAVEPROPERTIESOF lava
- The needed semantic information is in WordNet definitions, and was successfully translated into a form that was used for rough 'proofs'

## The Architecture of LCC's QA System around 2003



## References

- AskMSR: Question Answering Using the Worldwide Web
  - Michele Banko, Eric Brill, Susan Dumais, Jimmy Lin
  - <http://www.ai.mit.edu/people/jimmylin/publications/Banko-etal-AAAI02.pdf>
- In Proceedings of 2002 AAAI SYMPOSIUM on Mining Answers from Text and Knowledge Bases, March 2002
- Web Question Answering: Is More Always Better?
  - Susan Dumais, Michele Banko, Eric Brill, Jimmy Lin, Andrew Ng
  - <http://research.microsoft.com/~sdumais/SIGIR2002-QA-Submit-Conf.pdf>
- D. Ravichandran and E.H. Hovy. 2002. Learning Surface Patterns for a Question Answering System. ACL conference, July 2002.

## References

- S. Harabagiu, D. Moldovan, M. Pasca, R. Mihalcea, M. Surdeanu, R. Bunescu, R. Girju, V. Rus and P. Morarescu. *FALCON: Boosting Knowledge for Answer Engines*. The Ninth Text Retrieval Conference (TREC 9), 2000.
- Marius Pasca and Sanda Harabagiu. High Performance Question/Answering, in *Proceedings of the 24th Annual International ACL SIGIR Conference on Research and Development in Information Retrieval (SIGIR-2001)*, September 2001, New Orleans LA, pages 366-374.
- L. Hirschman, M. Light, E. Breck and J. Burger. *Deep Read: A Reading Comprehension System*. In Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics, 1999.
- C. Kwok, O. Etzioni and D. Weld. *Scaling Question Answering to the Web*. ACM Transactions in Information Systems, Vol 19, No. 3, July 2001, pages 242-262.
- M. Light, G. Mann, E. Riloff and E. Breck. *Analyses for Elucidating Current Question Answering Technology*. Journal of Natural Language Engineering, Vol. 7, No. 4 (2001).
- M. M. Soubbotin. *Patterns of Potential Answer Expressions as Clues to the Right Answers*. Proceedings of the Tenth Text Retrieval Conference (TREC 2001).