

Informatie Extractie met Reguliere Expressies

Gosse Bouma

Information Science
University of Groningen

Natuurlijke Taalverwerking II, 2006/2007

- ▶ Web Search, Information Extraction, Question Answering,
- ▶ Reguliere Expressies en Informatie Extractie,
- ▶ Question Answering als Informatie Extractie.



Web Search (Information Retrieval)

- ▶ Zoek relevante documenten over een bepaald onderwerp.
- ▶ **Vraag:** keywords (en boolese connectieven),
- ▶ Armstrong **AND** Cycling **AND NOT** Tour
- ▶ Zoeken in (indexen van) **Documenten**,
- ▶ **Antwoord:** lijst van relevante documenten.



Information Retrieval, Web Search

- ▶ Robuust (Gigabytes aan documenten, 1000-en queries),
- ▶ Technieken goed onderzocht,
- ▶ Nadeel:
 - ▶ Niet *direct* antwoord op een vraag,
 - ▶ Documenten bevatten relevante en irrelevante informatie.



Information Extraction

- ▶ Extraheer relevante informatie uit ongestructureerde tekst,
- ▶ Vul hiermee een **database**.
- ▶ **Vraag:** database query,
- ▶ **Antwoord:** matching entries uit de database.
- ▶ **Nadeel:** Systeem bepaalt wat relevant is.



Academic Transfer

Twee Promovendi Waarnemen en Bewegen (100% elk)
Groningen, 38 uur per week
Faculteit der Medische Wetenschappen
Functiebeschrijving
.....
Functie-eisen Opleidingsniveau/vaardigheden:
Universitair - afgeronde universitaire opleiding
bewegingswetenschappen, psychologie, (bio)fysica
.....
Aanbod De RUG biedt een salaris van minimaal ? 1.702,-
bruto per maand in het eerste jaar...

Contractduur: 4 jaar Maximaal aantal uur per week: 38



Een vacature-site

- ▶ Haalt vacatures van het web, uit de krant,
- ▶ **Ongestructureerde** informatie,
- ▶ Wordt doorzocht op relevante velden,
- ▶ b.v. functie, opleiding, bedrijf, plaats, salaris,...
- ▶ Informatie wordt opgeslagen in **database**.



citeseer.nj.nec.com

- ▶ Haalt papers van het web (ps, pdf),
- ▶ Wordt doorzocht op auteur, titel, abstract, verwijzingen,..
- ▶ Database, met bovendien:
 - ▶ Active bibliography,
 - ▶ Similar documents,
 - ▶ Users who viewed this document also viewed,
 - ▶ Citatie-index,



- ▶ **Open Domain:** CLEF evaluation exercise
 - ▶ Find answers for **factoid questions**
 - ▶ in Algemeen Dagblad, NRC 1994, 1995
 - ▶ 80 mln words
 - ▶ **CLEF 2007:** use `wikipedia.nl` as well (50 mln words)
- ▶ **Closed Domain:** Medical QA
 - ▶ Find answers to **medical questions**
 - ▶ in encyclopedia and reference material, web documents
 - ▶ 3 mln words

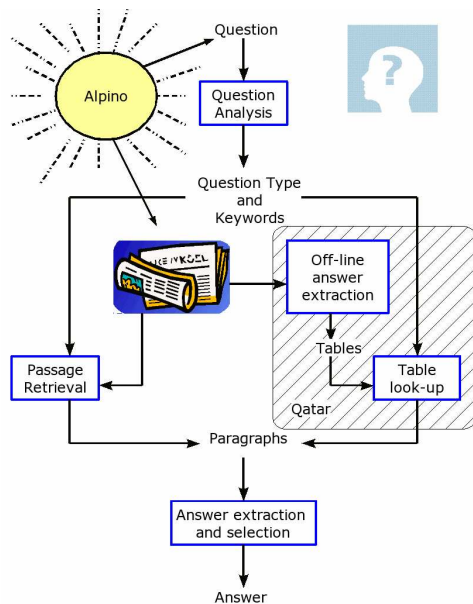
High-quality syntactic information is useful for QA

- ▶ question analysis
 - ▶ answer extraction
 - ▶ on-line and off-line
 - ▶ acquiring ontological information
- ▶ answer ranking
- ▶ information retrieval



Joost: a QA system for Dutch

Question Answering



[Frequently Asked Questions](#)
[Best Questions](#)
[Recent Questions](#)
[Random Sample Questions](#)

Sample Questions

Wanneer werd het embargo tegen Irak ingesteld?

Hoeveel hunebedden zijn er in Nederland?

Wie won de Nobelprijs voor economie

Wat is de hoofdstad van Iran

Wanneer was de slag bij heiligerlee

Waarvoor staat de afkorting CIA

Wat betekent KI?

Wie hebben de Nobelprijs voor economie gewonnen

Wanneer verdween de sjah uit Iran

Waar ligt de ijsbaan van Groningen?

[more ...](#)

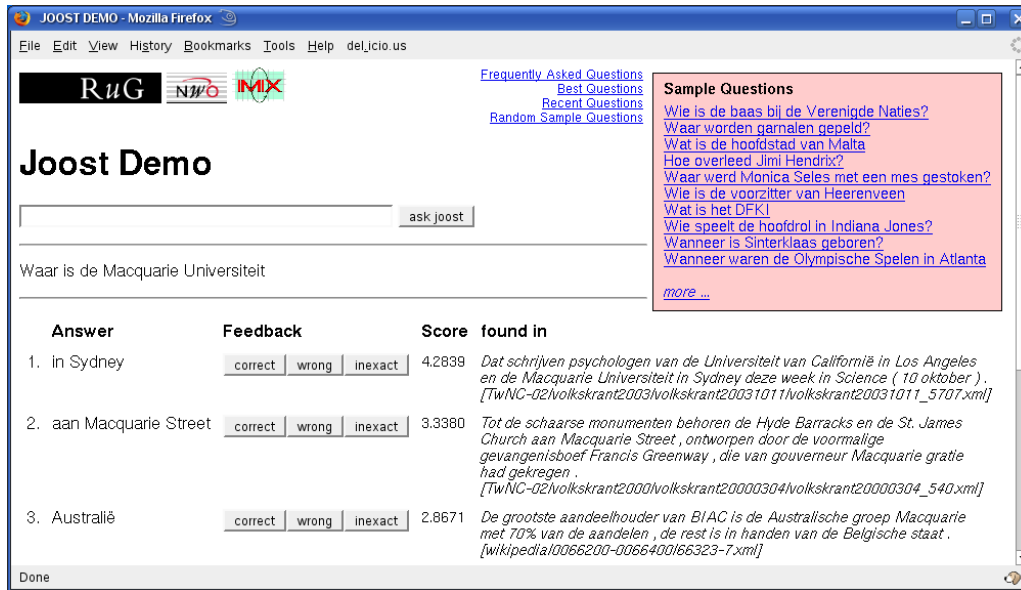
Joost Demo

ask joost

Wat is het DFKI

Answer	Feedback	Score	found in
1. universitair onderzoekscentrum in Saarbrucken	<input checked="" type="radio"/> correct (1) <input type="radio"/> wrong <input type="radio"/> inexact	7.6556	<i>Verbmobil is bedacht door DFKI, een universitair onderzoekscentrum in Saarbrucken.</i> [TwnC-02/volkskrant2001/volkskrant20010830/volkskrant20010830_3131.xml]





- ▶ For frequently asked question types, answers are searched off-line
 - ▶ How many **inhabitants** does **Location** have?
 - ▶ When was **Person** born?
 - ▶ Who won the **Nobelprize** for literature in 1990?
 - ▶ What does the **abbreviation** ADHD mean?
 - ▶ What **causes** Frei syndrome?
 - ▶ What are the **symptoms** of poisoning by mushrooms?

Information Extraction

Information Extraction Strategie

- ▶ Extraheer de geboortedata van personen
- ▶ Extraheer de geboorteplaats (of het geboorteland) van personen

1. Zoek naar voorbeelden,
2. Ontdek patronen,
3. Test precisie van deze patronen,
4. (Pas patronen aan, en itereer naar 1)
5. Doe extractie.

Eigenlijke Extractie

- ▶ Vind de relevante delen:
- ▶ **Naam** + geboortedatum,
- ▶ Print de naam en de geboortedatum

Andere patronen?

```
> zgrep 'geboren in 19[0-9][0-9]-' *.gz
Mac Rebennack , geboren in 1941
Ze is geboren in 1957 in Revnice
Carl Barks ( geboren in 1901
Oppenheim ( Amerikaan , geboren in 1938 )
Kanan werd geboren in 1949
De Amerikaan John Baldessari bijvoorbeeld ( geboren
in 1931 )
```



Precision en Recall

- ▶ Precision : $\frac{\text{Goede geëxtraheerde relaties}}{\text{Geëxtraheerde relaties}}$
- ▶ Recall : $\frac{\text{Goede geëxtraheerde relaties}}{\text{Relaties in de tekst}}$

Informatie Extractie Methodologie

- ▶ Handmatig
 - ▶ Moeilijk (m.n. bedenken van nieuwe patronen),
 - ▶ Tijdrovend.
- ▶ Machine Learning
 - ▶ Aantal voorbeelden handmatig coderen,
 - ▶ Programma ontdekt patronen om vergelijkbare gevallen te vinden

