



Grupo de Pesquisa em  
Inteligência e Imagens

# Sistemas de *Question-Answering* (QA)



Enter what you want to calculate or know about:



 Examples  Random

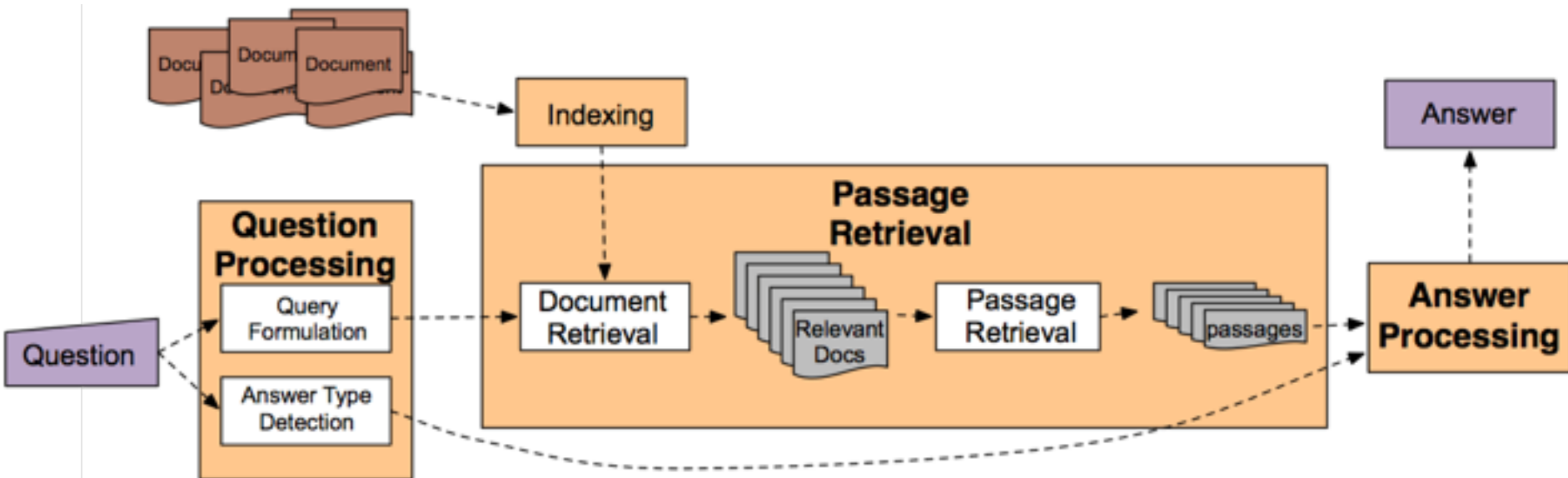


# Abordagens para QA

- Baseadas em Recuperação de Informação (IR)
  - TREC; Google
- Baseadas em Conhecimento (KB)
  - Apple Siri; Wolfram Alpha
- Híbridas
  - IBM Watson



# QA via IR





## QA via KB

- Constrói representação semântica da query
  - Datas, localidades, entidades, quantidades, etc.
- Representação semântica —> query em dados estruturados
  - Bases de dados geoespaciais
  - Ontologias (Wikipedia, dbPedia, WordNet, etc)
  - Bases de dados científicas



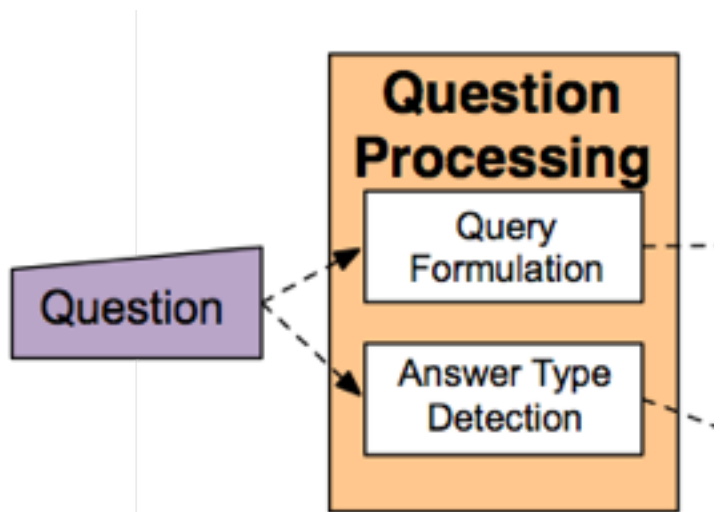
## QA híbrida

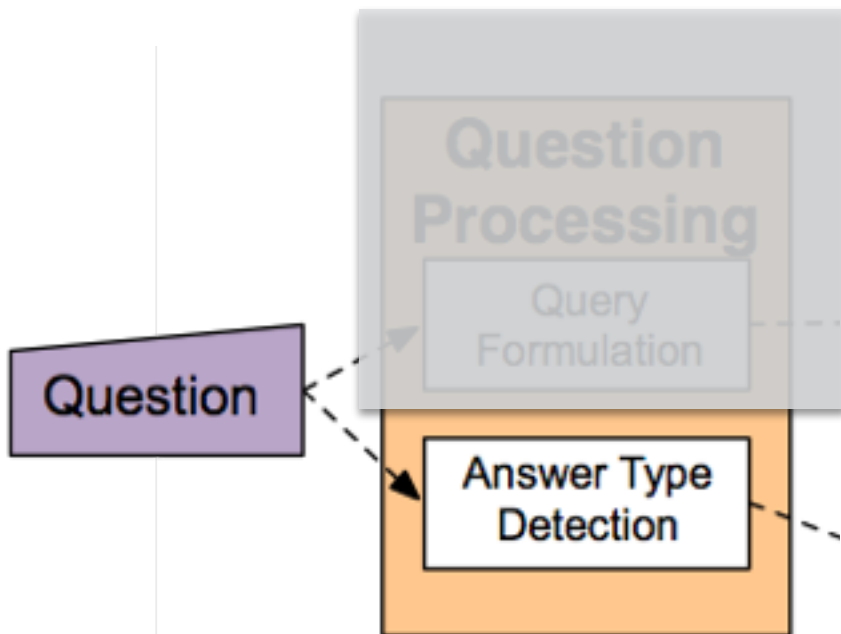
- Constrói representação semântica da query
- Gera respostas candidatas via IR
- Pontua cada candidata usando fontes de conhecimento mais ricas



Grupo de Pesquisa em  
Inteligência e Imagens

## Tipos de respostas e formulação da query



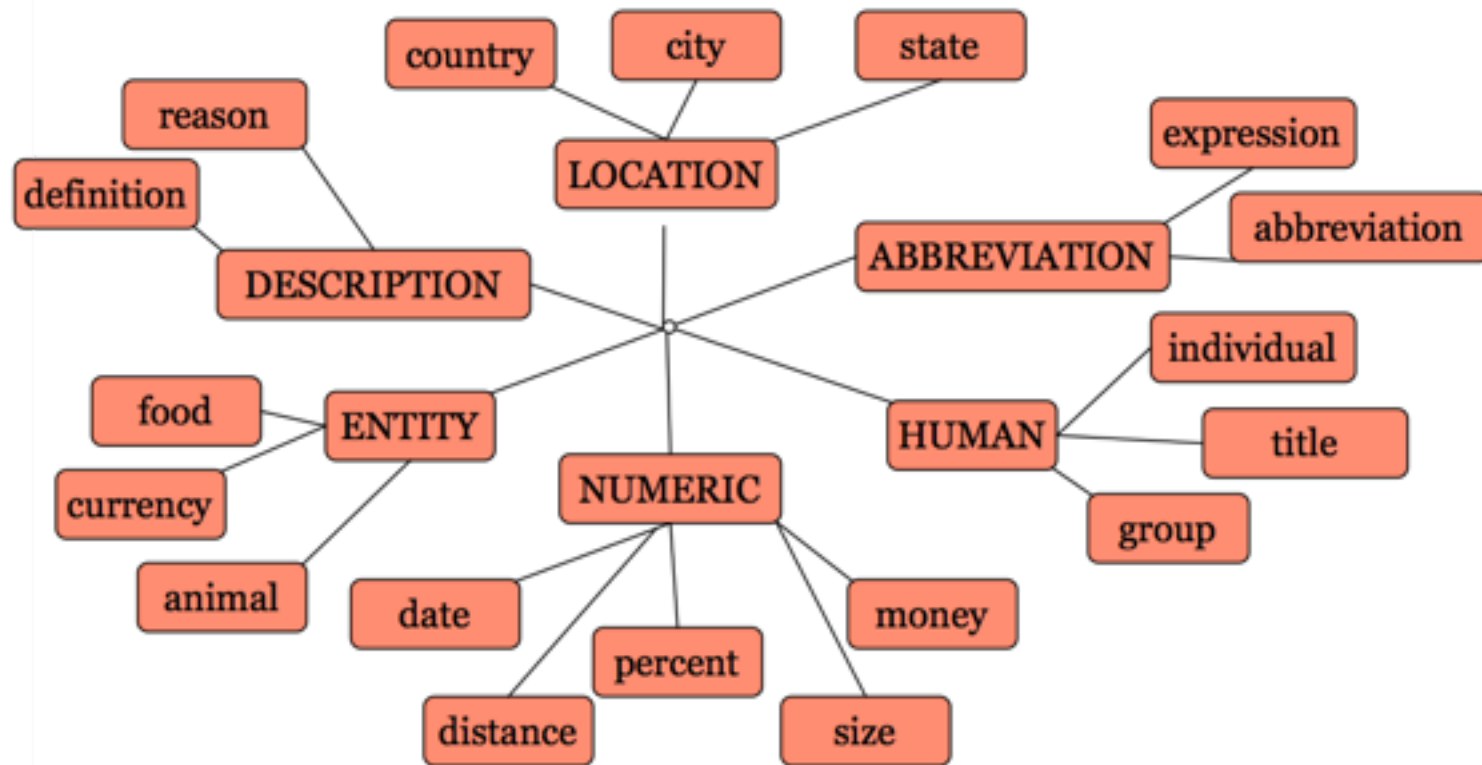






# Taxonomia do tipo de resposta

Xin Li, Dan Roth. 2002. Learning Question  
Classifiers. COLING'02





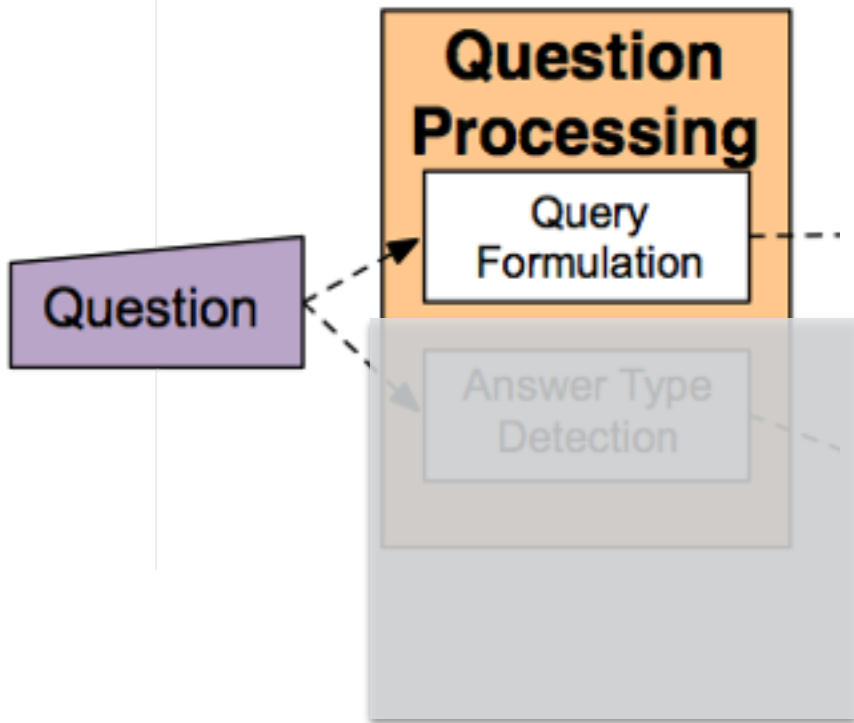
# Detecção do tipo de resposta

- Regras definidas manualmente
  - Uso de Expressões Regulares
    - Who {islwaslarelwere} PERSON
    - PERSON (YEAR – YEAR)
  - Palavra central da 1a frase nominal depois do pronome relativo (Que? Qual? Quem? Como? Quando? Onde?)
    - Quem é o **autor** de Dracula?
    - Qual é a **flor** do estado da California?



# Detecção do tipo de resposta

- Aprendizagem de Máquina
  - Definir taxonomia de tipos de questões
  - Rotular dados de treinamento para cada tipo de questão
  - Treinar classificador para cada tipo de questão usando features:
    - Palavras e frases de questionamento
    - POS tags
    - Palavras centrais (headwords)
    - Entidades nomeadas
    - Palavras relacionadas semanticamente





# Algoritmo para seleção de palavras-chave

Dan Moldovan, Sanda Harabagiu, Marius Păcă, Rada Mihalcea, Richard Goodrum, Roxana Girju and Vasile Rus. 1999. Proceedings of TREC-8.

1. *Select all non-stop words in quotations*
2. *Select all NNP words in recognized named entities*
3. *Select all complex nominals with their adjectival modifiers*
4. *Select all other complex nominals*
5. *Select all nouns with their adjectival modifiers*
6. *Select all other nouns*
7. *Select all verbs*
8. *Select all adverbs*
9. *Select the QFW word (skipped in all previous steps)*
10. *Select all other words*



1. *Select all non-stop words in quotations*

...

4. *Select all other complex nominals*

...

7. *Select all verbs*

...

~~Who~~ ~~coined~~ ~~the~~ ~~term~~ ~~“cyberspace”~~ ~~in~~ ~~his~~ ~~novel~~ ~~“Neuromancer”~~?

7

4

1

4

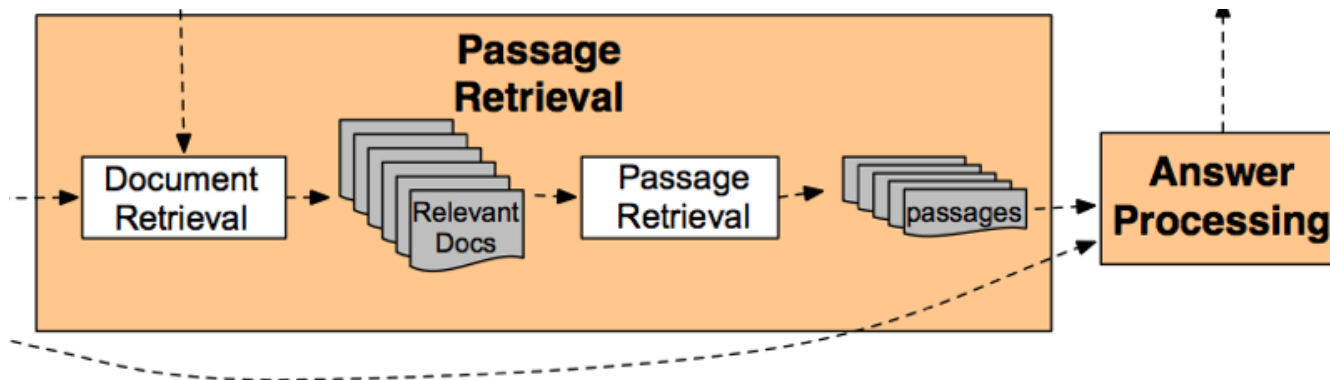
1

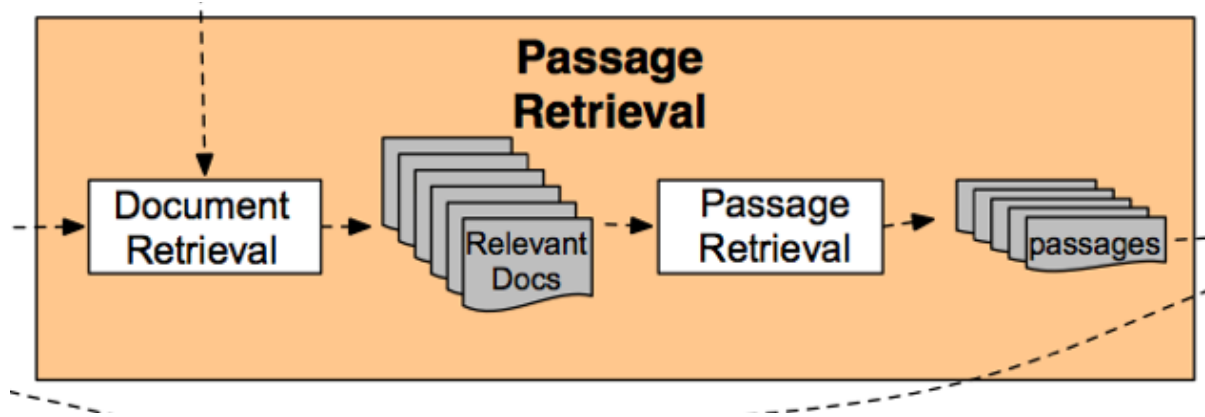
cyberspace/1 Neuromancer/1 term/4 novel/4 coined/7



Grupo de Pesquisa em  
Inteligência e Imagens

## Recuperação do trecho e Extração da resposta





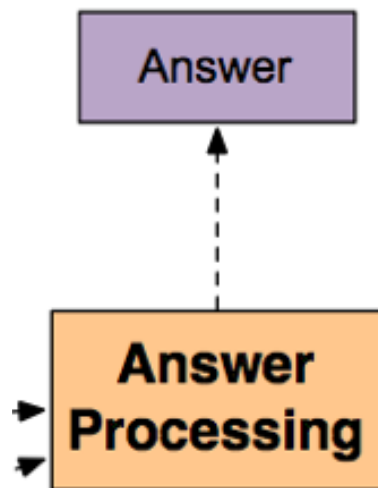
- Passo 1: IR recupera documentos a partir da query
- Passo 2: Segmenta documentos em trechos (ex. parágrafos)
- Passo 3: Ranqueamento do trecho
  - utiliza tipo da resposta para re-ranquear o trecho





# Features para ranqueamento do trecho

- # de Entidades Nomeadas do tipo certo no trecho
- # de palavras da query existentes no trecho
- Proximidade das palavras-chave da query para qualquer outra no trecho
- Ranqueamento do documento contendo o trecho





# Extração da resposta

- Executar um rotulador de entidades nomeadas para tipos de resposta nos trechos
  - Cada tipo de resposta requer um rotulador que a detecte
  - Se o tipo é CITY, o rotulador deve rotular CITY
    - pode ser via RegEx
- Retorna a string com o tipo correto
  - Quem é o Primeiro Ministro da Índia? **(PERSON)**  
**Manmohan Singh**, Primeiro Ministro da Índia, havia contado a líderes da esquerda que o acordo não seria renegociado.
  - Quanto alto é o Monte Everest? **(LENGTH)**  
A altura oficial do monte Everest é **29035 feet**



# Ranquear respostas candidatas...

- ... caso existam mais de uma.

Q: Who was Queen Victoria's second son?

- Tipo da resposta: **Person**
- Trecho:

The Marie biscuit is named after **Marie Alexandrovna**, the daughter of **Czar Alexander II of Russia** and wife of **Alfred**, the second son of **Queen Victoria** and **Prince Albert**



# Aprendizagem de máquina

## Features para ranquear respostas candidatas

Answer type match: candidata contém uma frase com o tipo de resposta correto

Pattern match: a identidade da expressão regular que casa com a candidata

Question keywords: # de palavras-chave da questão existentes na candidata

Keyword distance: # médio de palavras comuns as palavras-chave da query e da candidata

Novelty factor: pelo menos uma palavra da candidata não está na query

Apposition features: a candidata é um aposto para uma frase contendo termos da questão

Punctuation location: a candidata é imediatamente seguida por uma vírgula, ponto, aspas, ponto e vírgula ou ponto de exclamação.

Sequences of question terms: o comprimento da maior sequência de termos da questão que ocorre na resposta candidata.

**IBM Watson: >50 componentes**



# Métricas de avaliação

1. Accuracy (a resposta casa com a desejada/correta?)
2. *Mean Reciprocal Rank (MRR)*
  - Para cada query, retorna uma lista ranqueada de M respostas candidatas
  - O score da query é 1/Rank da primeira resposta correta
    - se a 1a é correta: 1
    - senão, se a segunda é a correta:  $\frac{1}{2}$
    - senão, se a terceira é a correta:  $\frac{1}{3}$ , etc.
    - se nenhuma correta  $\Rightarrow$  Score = 0
  - Considere a média sobre todas as N queries

$$MRR = \frac{\sum_{i=1}^N \frac{1}{rank_i}}{N}$$



Grupo de Pesquisa em  
Inteligência e Imagens

# Utilizando conhecimento em QA



# Extração de relação

- Respostas: bases de dados de relações
  - born-in(“Emma Goldman”, “June 27 1869”)
  - author-of(“Cao Xue Qin”, “Dream of the Red Chamber”)
  - retirado de Wikipedia infoboxes, DBpedia, FreeBase, etc.
- Questões: extração de relações em Questões

Whose granddaughter starred in E.T.?

(acted-in ?x “E.T.”)

(granddaughter-of ?x ?y)





# Raciocínio temporal

- Bases de dados de relações
  - (e obituários, dicionários biográficos, etc.)
- IBM Watson

"In 1594 he took a job as a tax collector in Andalusia"

Candidatos:

  - **Thoreau** é uma resposta ruim (nasceu em 1817)
  - **Cervantes** é possível (estava vivo em 1594)



# Conhecimento geospacial

- Beijing é uma boa resposta para "cidade asiática"
- Califórnia é "sudoeste de Montana"
- geonames.org:

The screenshot shows a web browser window with the URL [www.geonames.org/search.html?q=palo+alto&country=](http://www.geonames.org/search.html?q=palo+alto&country=). The page header includes links for [GeoNames Home](#), [Postal Codes](#), [Download](#), [Webservice](#), and [About](#), along with a [login](#) link. The search bar contains the text "palo alto" and a dropdown menu is set to "all countries". Below the search bar are buttons for "search", "show on map", and a link for "advanced search". The results section indicates "459 records found for 'palo alto'".

	Name	Country	Feature class	Latitude	Longitude
1	<a href="#">Palo Alto</a> Palo Al'to, Palo Alto, pa luo ao duo, paroeruto, Пало Алто, Пало Альто, פאלו אלטו, パロアルト, 帕羅奧多	<a href="#">United States</a> , California Santa Clara County	populated place population 64,403, elevation 9m	N 37° 26' 30"	W 122° 8' 34"
2	<a href="#">Palo Alto Township</a> Palo Alto Township	<a href="#">United States</a> , Iowa Jasper County	administrative division elevation 256m	N 41° 38' 15"	W 93° 2' 57"
3	<a href="#">Borough of Palo Alto</a>	<a href="#">United States</a> , Pennsylvania Schuylkill County	administrative division population 1,032, elevation 210m	N 40° 41' 21"	W 76° 10' 2"



Grupo de Pesquisa em  
Inteligência e Imagens

# Sistemas de *Question-Answering* (QA)