

Data-oriented Content Query System: Searching for Data into Text on the Web

Mianwei Zhou, Tao Cheng, Kevin Chen-Chuan Chang
WSDM 2010, New York, USA



Many Web Applications Try to Exploit the “Content” of Web Pages.

In most cases, what we really want are not pages, but the information units inside.

Web Info
Extraction

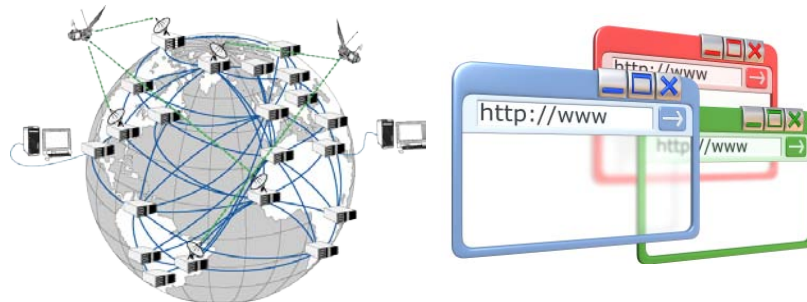
Typed Entity
Search

Web-based
Q/A



?

?



Content-Exploiting Application 1: Web Information Extraction

Web Information Extraction (WIE)

(Marius 2006, Cafarella 2005, Etzioni 2004)

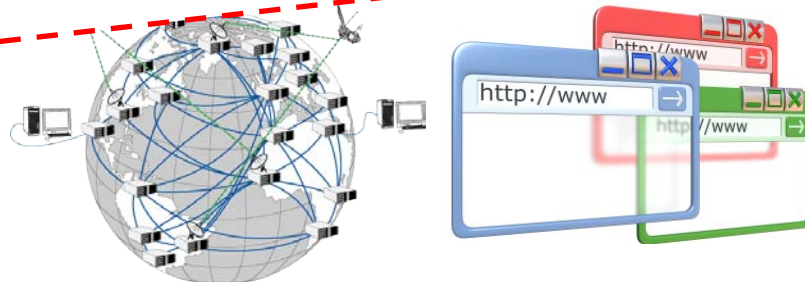
Pattern:

"X is CEO of Y"

Company	CEO
Google	Eric Schmidt

Limitation

- Focus on simple patterns.
- Lack of interactivity.



Content-Exploiting Application 2: Web-based Question Answering

Web-based Question Answering (WQA)

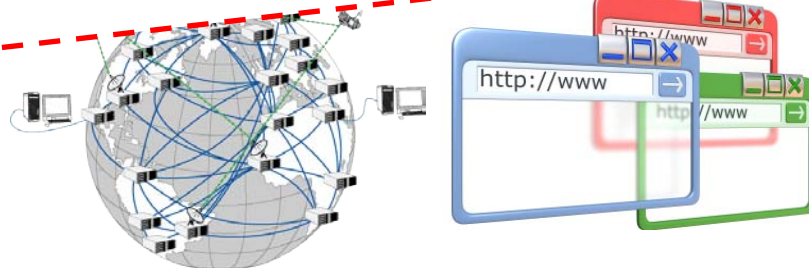
(Wu 2007, Lin 2003, Brill 2002)

Who is CEO
of Dell?

Michael Dell

Limitation

- Only rely on top-k pages to retrieve the answer.



Content-Exploiting Application 3: Typed-Entity Search

Typed-Entity Search (TES)

(Cheng 2007, Cafarella 2007, Chakrabarti 2006)

Amazon
Phone



Entity Search

Ranked Entity List

Limitation

- Limited Number of Data Type
- Lack of Flexibility



But ... Where is CEO ?



General System for Querying Text “Content”, Much Like How DBMS Supports Data Application

Web Info
Extraction

Typed Entity
Search

Web-based
QA



Data-oriented Content Query System

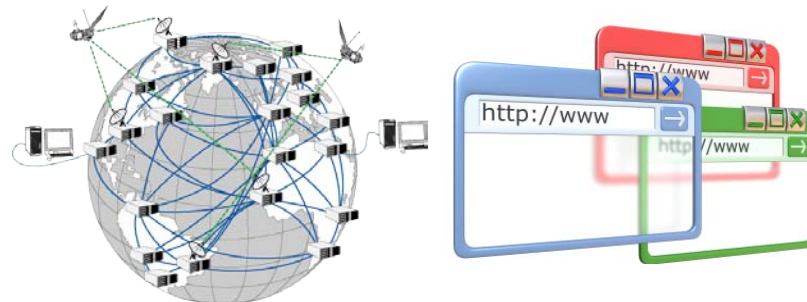
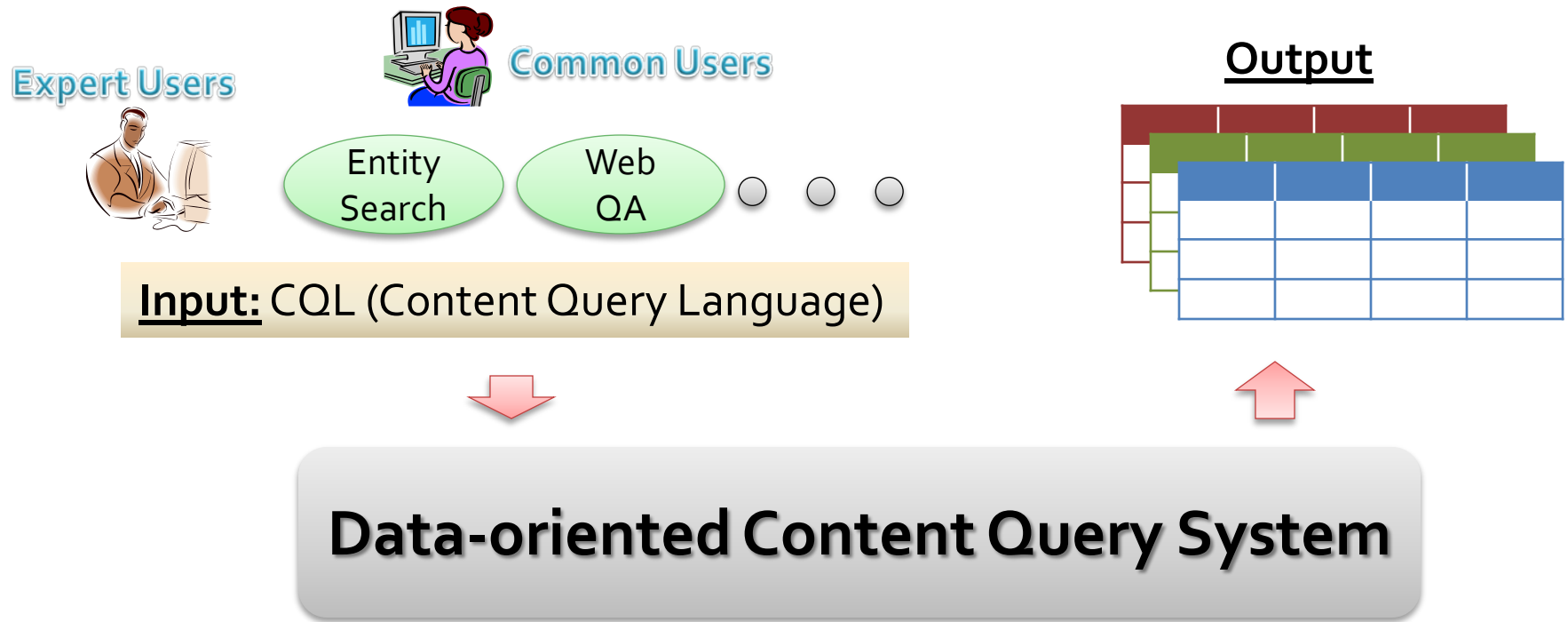
Requirements

1. Extensible Data Types
2. Flexible Contextual Patterns
3. Customizable Scoring

System Framework

Input: CQL

Output: Table



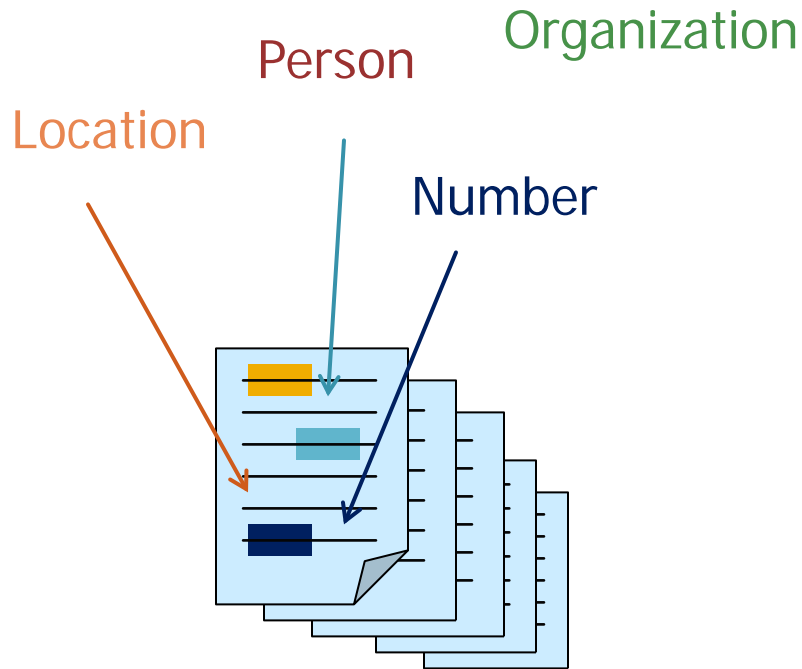
Demo

Online Demo

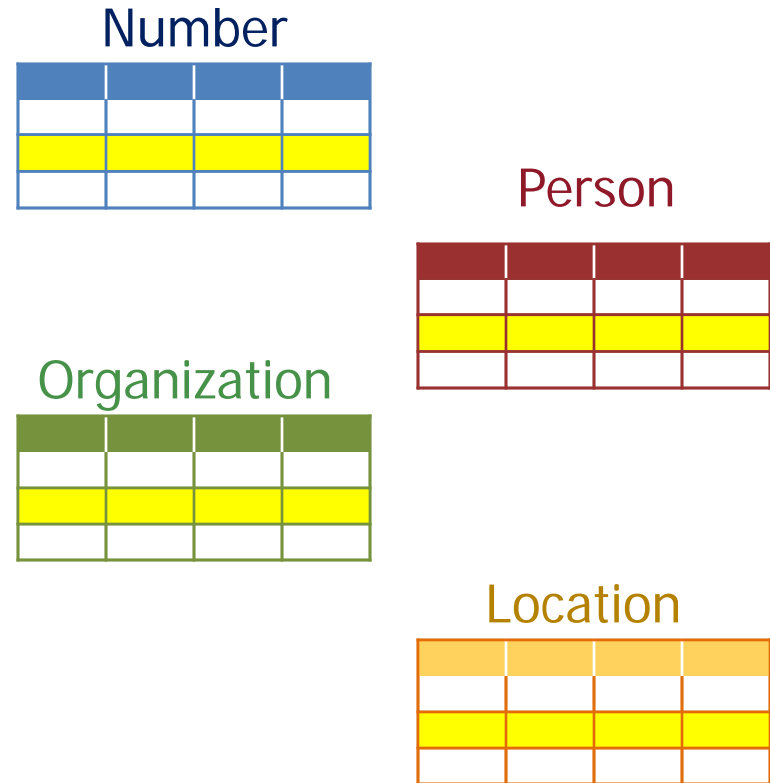
<http://wisdm.cs.uiuc.edu/demos/docqs>

Why Relational Model: Occurrence->Tuple

What we need



Relational Model



Why Relational Model: Content Query->Relational Operation

What we need

Find the population of China



China has a population of **1.3 billion**

China with its population of **1.3 billion** people

China is established in **1949**.

Shanghai is the largest city with **15 million** inhabitants in China



1.3 bi
1. 1.3 billion
2. 15 million
...

Relational Model

FROM #number



WHERE pattern(...)



GROUP BY #number



ORDER BY conf()

Why Relational Model

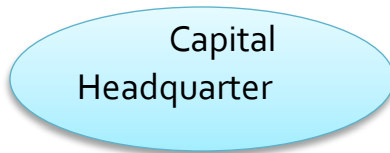
Extensible Data Type->View

What we need

Number



Location



Person



Relational Model

Table

Number			

View

population

price

phone

Main Technique: Highly Efficient and Scalable Index Structure

Data Type Definition

INPUT
SELECT ...
FROM ...
WHERE ...

OUTPUT

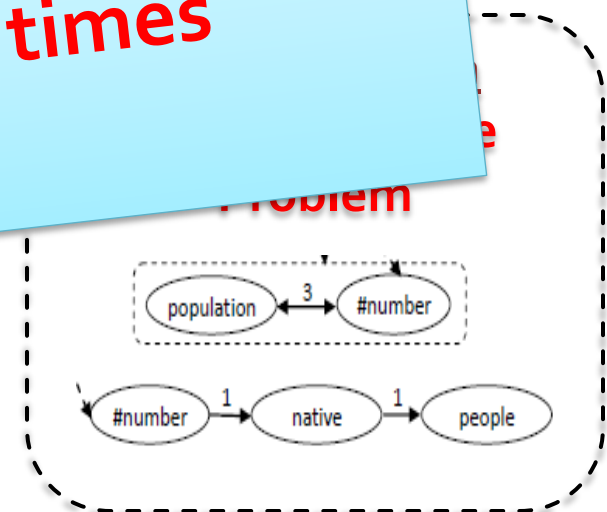
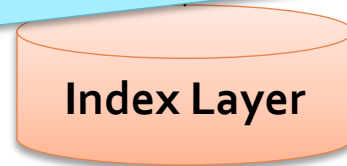
Experimental Result

- Speed improvement: **6-10 times**
- Space overhead: **Around 2 times** original corpus size.

Index Design

Special I

- Contextual
- Join Index



Data-oriented Content Query System: Supporting Ad-hoc, Interactive "Content" Search.

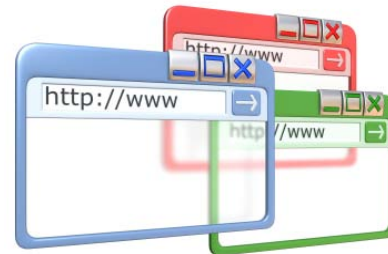
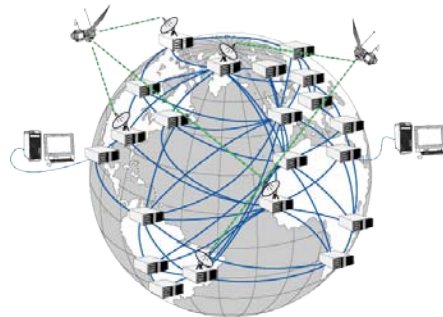
Web Info
Extraction

Typed Entity
Search

Web-based
Q/A



Data-oriented Content Query System
Interactive *Ad-hoc*



Q & A

Thank You!