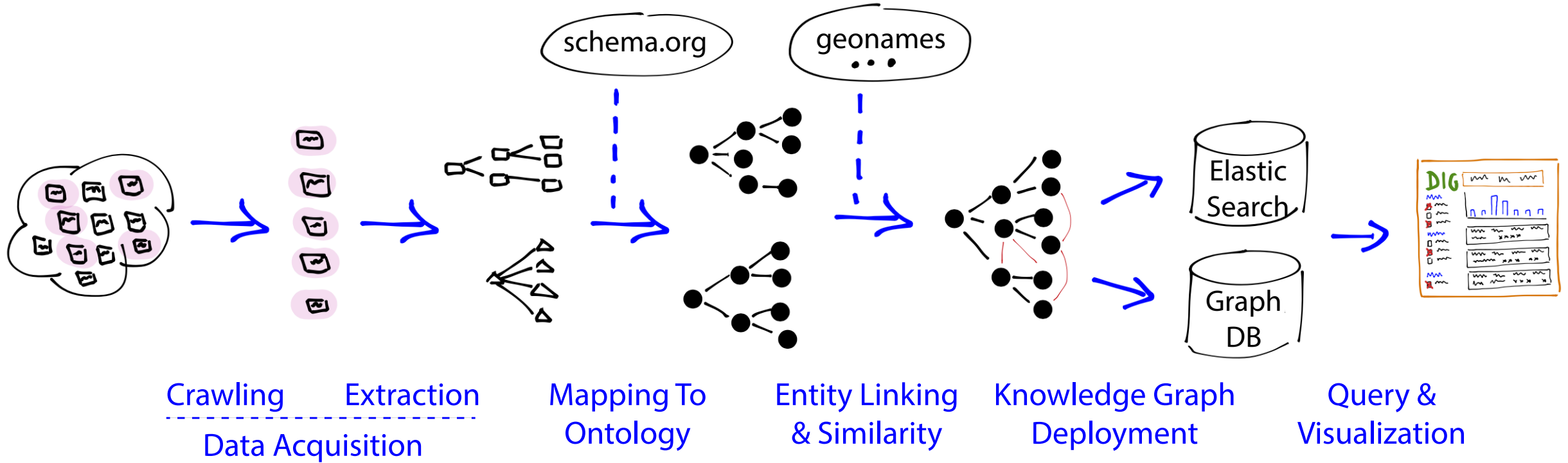# Information integration on the web

Mayank Kejriwal

Let's go deeper into the architecture

# Typical Web information integration architecture



Crawling     Extraction     Mapping To Ontology     Entity Linking & Similarity     Knowledge Graph Deployment     Query & Visualization

Data Acquisition

schema.org     geonames

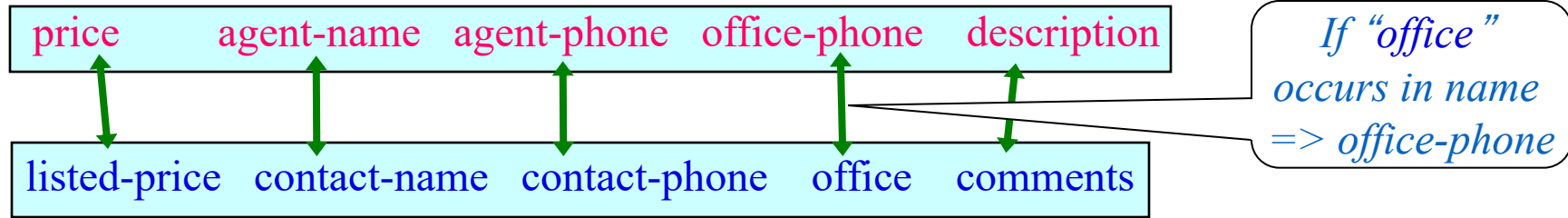Elastic Search     Graph DB     DIG

# Schema Mapping

- Given two different sources with different schemas, how do we automatically align the information

- Research Topics
  - Automatic schema alignment based on structure and naming
  - Automatic alignment based on the source contents

# Schema Mapping

**Mediated schema**

| price | agent-name | agent-phone | office-phone | description |
|---|---|---|---|---|

| listed-price | contact-name | contact-phone | office | comments |
|---|---|---|---|---|

*Schema of realestate.com*

*If "office" occurs in name => office-phone*

**realestate.com**

| listed-price | contact-name | contact-phone | office | comments |
|---|---|---|---|---|
| $250K | James Smith | (305) 729 0831 | (305) 616 1822 | Fantastic house |
| $320K | Mike Doan | (617) 253 1429 | (617) 112 2315 | Great location |
| ....... | ....... | ....... | ....... | ....... |

**homes.com**

| sold-at | contact-agent | extra-info |
|---|---|---|
| $350K | (206) 634 9435 | Beautiful yard |
| $230K | (617) 335 4243 | Close to Seattle |

*If "fantastic" & "great" occur frequently in data instances => description*
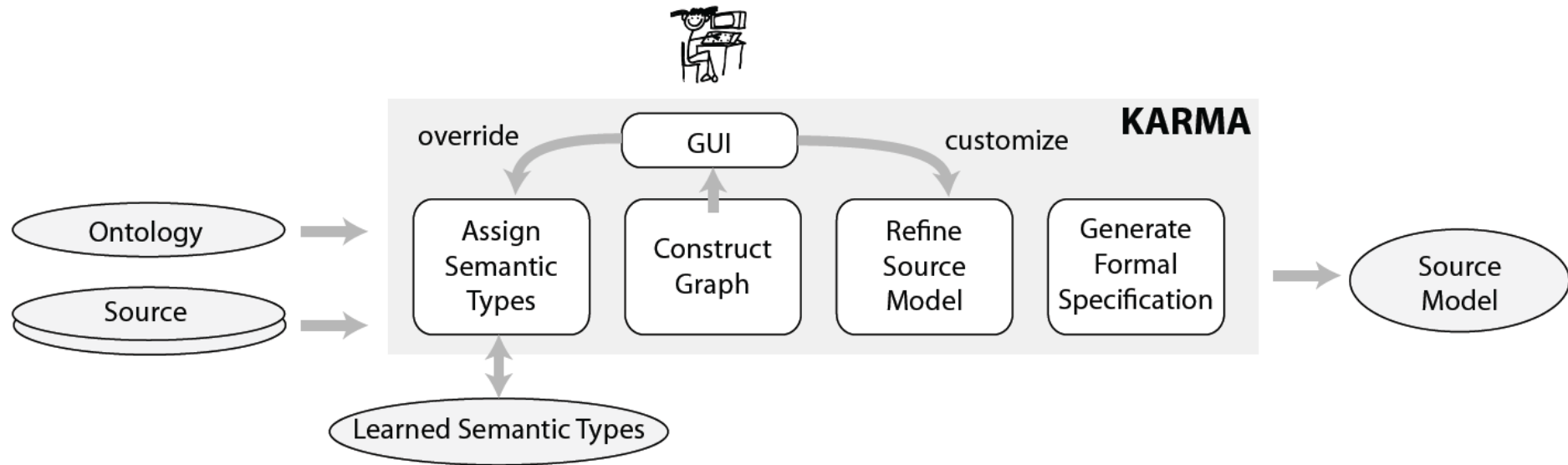
# Source Modeling

- Semantic typing
- Source discovery
- Automatic source modeling
- Interactive source modeling

# Automatic Source Modeling

- How to learn semantic descriptions of sources and services

# Semi-Automatic Source Modeling



Relational database to RDF mapping

Anything to RDF mapping

# String Similarity:
# Why Strings Don't Match Perfectly?

typos      "Joh" vs "John"

OCR errors      "J0hn" vs "John"

formatting conventions      "03/17" vs "March 17"

abbreviations      "J. S. Sargent" vs "John Singer Sargent"

nick names      "John" vs "Jock"

word order      "Sargent, John S." vs "John S. Sargent"

# String Similarity Problem Definition

Given X and Y sets of strings

Find pairs (x, y)
such that both x and y
refer to the same real world entity



"John S. Sargent"

"John Singer Sargent"

# Record Linkage



How can the same objects be identified
when they are stored in inconsistent text formats?

# Record Linkage

- Align entities across sources

- Research Topics:
  - Blocking
  - Matching individual attributes
  - Matching records
  - Matching entities

Silk - A Link Discovery Framework for the Web of Data

Freie Universität Berlin

Robert Isele (Freie Universität Berlin)
Anja Jentzsch (Freie Universität Berlin)
Chris Bizer (Freie Universität Berlin)
Julius Volz (Google)

# Mashup Construction

# Ontology-based data access and integration

- Use ontology language as domain model
  - OWL2 profiles
- Answering queries under description logic constraints
  - Unions of conjunctive queries
  - Datalog