# A Scalable Architecture for Extracting, Aligning, Linking, and Visualizing Multi-Int Data

Craig Knoblock & Pedro Szekely

University of Southern California

# Introduction

- Massive quantities of data available for analysis
  - OSING, HUMINT, SIGINT, MASINT, GEOINT, …
- Data is spread across multiple sources, multiple sites and multiple formats
  - Databases, text, web sites, XML, JSON, etc…
- If an analyst could exploit all of this data, it could transform analysis
  - Disruptive technology for analysis

# Solution:
# Domain-specific Insight Graphs

- Innovative architecture
  - Extracting, aligning, linking, and visualizing massive amounts of data
  - Domain-specific content from structured and unstructured sources
- State-of-the art open source software
  - Open architecture with flexible APIs
  - Cloud-based infrastructure (HDFS, Hadoop, ElasticSearch, etc.)

# Example Scenario

- Want to determine the nuclear know-how of a given country from open source data

- Analyze the universities, academics, publications, reports, articles within the country

# Scenario Results

- Exploit the data available from
  - Web pages, publications, articles, etc.

- Produce a knowledge graph
  - Key people and connections
  - Technical capabilities and how they have changed over time

# DIG Pipeline

- Crawling
- Extracting
- Cleaning
- Integration
- Computing simlarity
- Entity resolution
- Graph construction
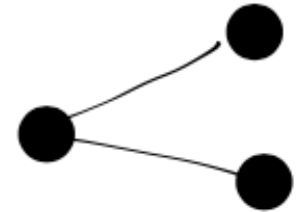- Query, analysis, and visualization

# Crawling

- Challenge: how to crawl just the relevant pages
- Approach:
  - Uses the Apache Nutch framework for Web pages
  - Uses Karma to extract pages from the deep Web

# Extracting

- Need to produce a structured representation for indexing and linking
- Highly configurable architecture for extractors
  - Learning of landmark extractors for structured data
  - Trainable CRF-based extractors for unstructured data
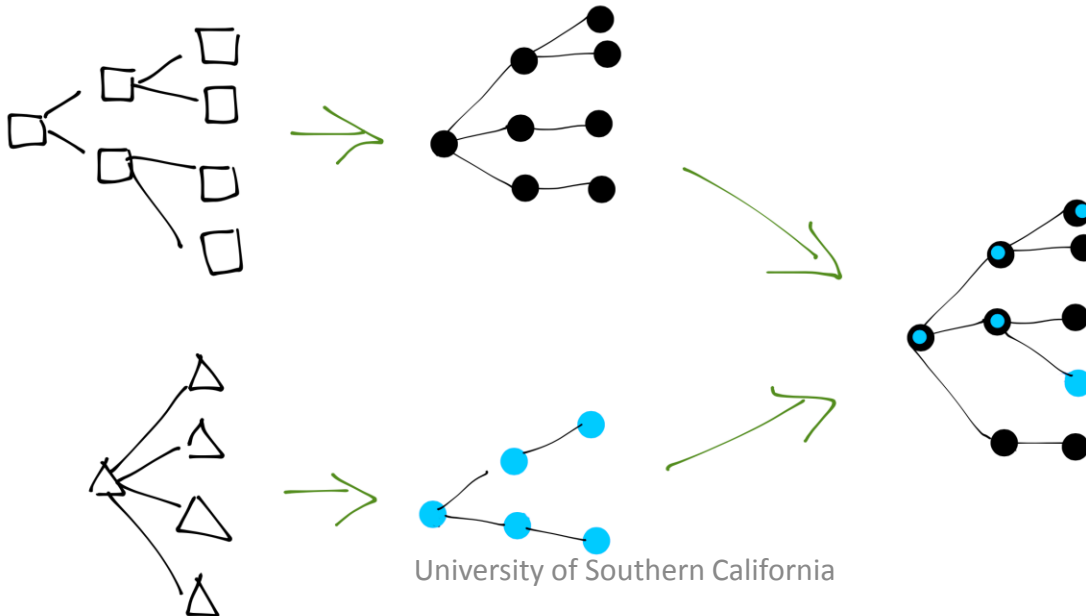  - Uses Mechanical Turk to crowd source training data

# Cleaning

- Cleaning and normalization to support analysis and linking
    - Visualization showing data distribution
    - Learned transformations from examples
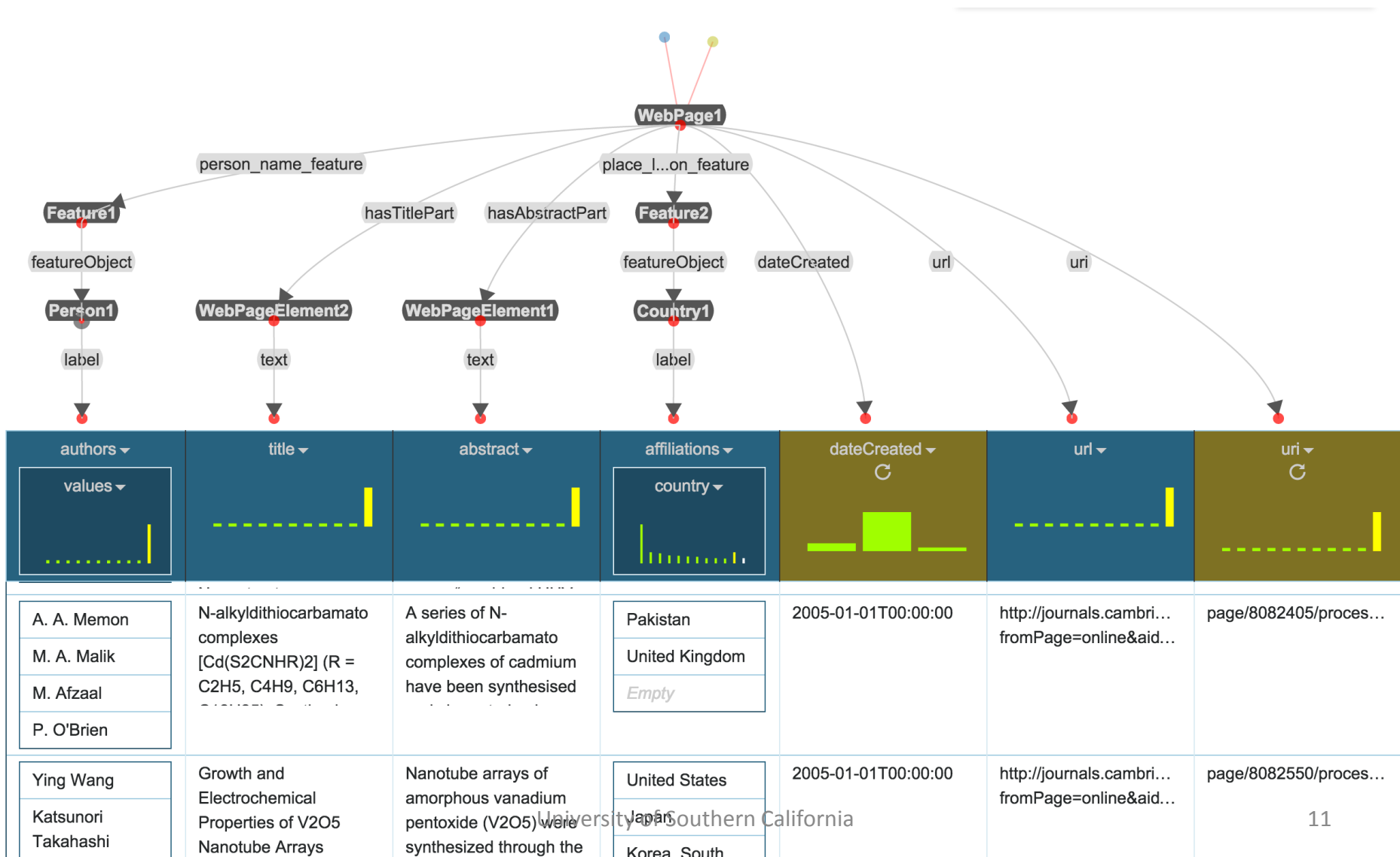    - Cleaning programs written in Python

# Integration

- Need to align the data across extracted data and structured sources

- Performed using a data integration tool called Karma
    - Karma maps arbitrary sources into a shared domain vocabulary (schema alignment)
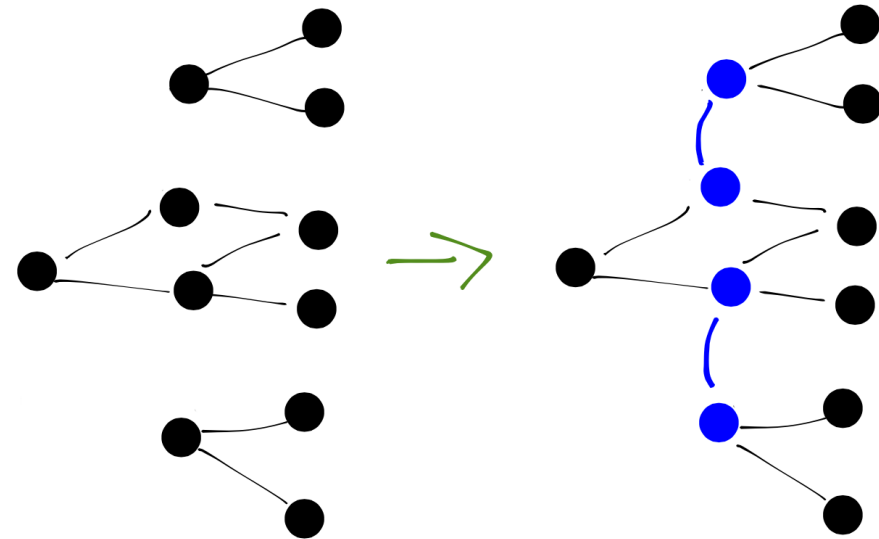    - Uses machine learning to minimize user effort
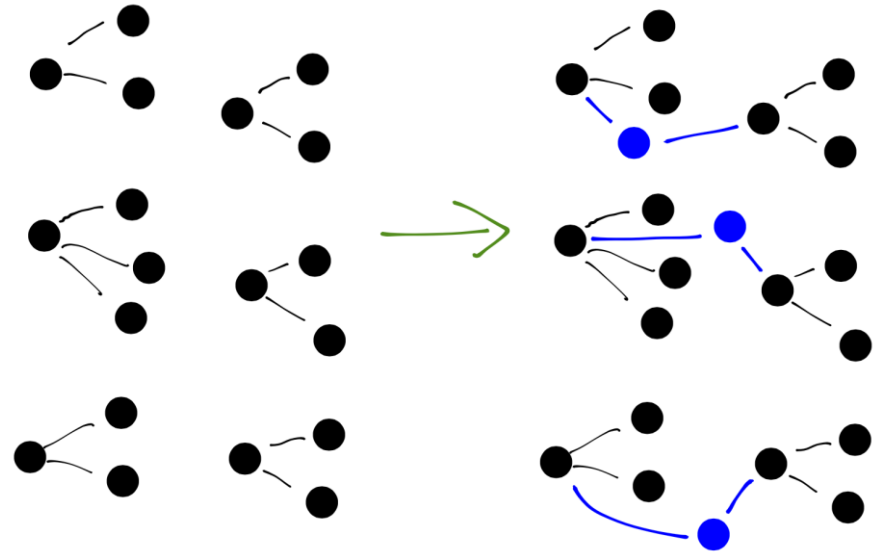
# Integration Using Karma

# Similarity

- Computes similarity across text fields and images
  - Image similarity done using DeepSentiBank
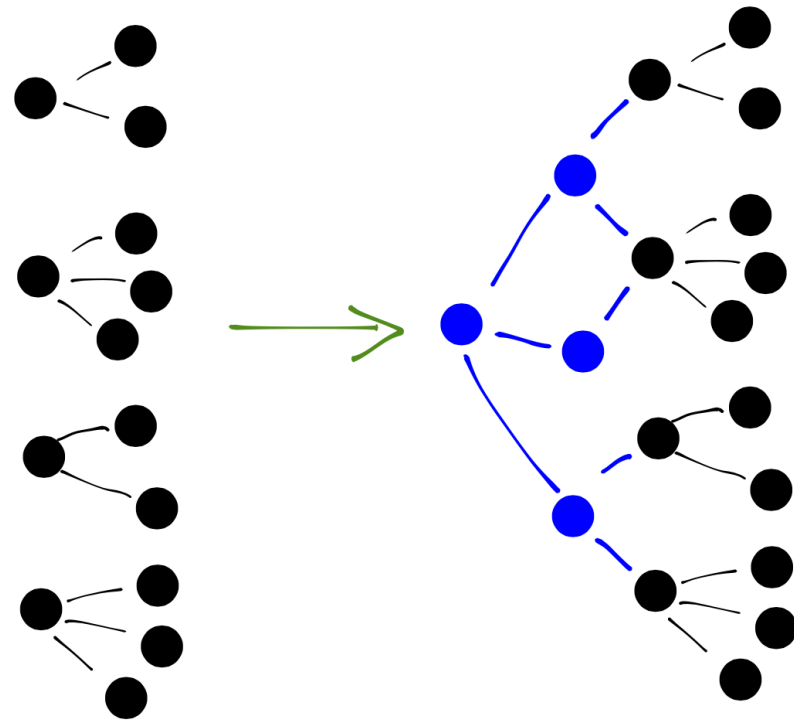  - Text similarity done using Minhash/LSH

# Entity Resolution

- Finds matching entities
- Reference source
  - Match against source to disambiguate entities
  - E.g., geonames for locations
- No reference source
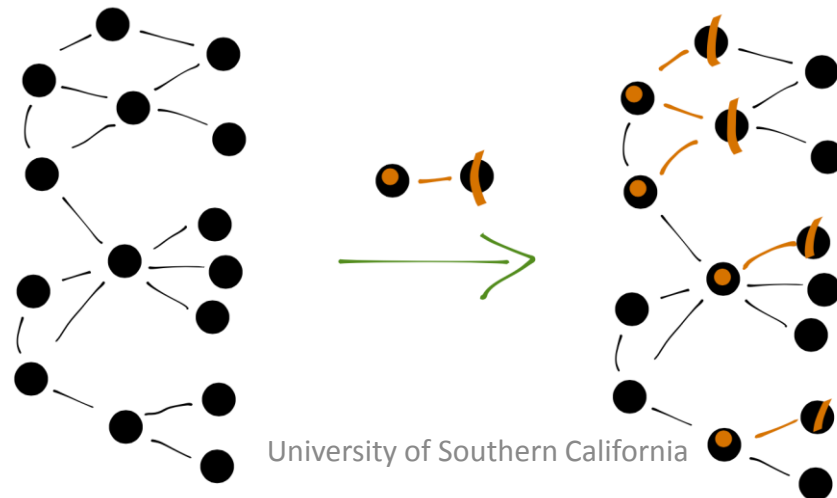  - Combine entities by considering the similarity across multiple fields

# Graph Construction

- Data is integrated into a graph that can be queries and analyzed
  - Data stored in HDFS
  - Data represented in a common language JSON-LD
  - Represented using a common terminology

# Query, Analysis and Visualization

- Challenge: support efficient querying against the graph

- Employ ElasticSearch to provide keyword querying, faceted browsing, and aggregation queries

# Query, Analysis & Visualization

- Visualization interface that provides faceted queries, timeslines, maps, etc.

# Discussion

- Technology that can provide dramatic new insights from data that is already available
- Applies to a wide range of problems
  - Determining the nuclear know-how of a given country
    - Technologies, key scientists, relevant organizations
  - Combating human trafficking
  - Understanding trends in technical areas
    - E.g., Material Science
  - Analyzing the competitive landscape of companies
  - and many other domains with massive quantities of data

# USC DIG Team

# Acknowledgements

- Collaborators



- Sponsor
  - DARPA
    - AFRL contract number FA8750-14-C-0240

# Thanks!

- More information:
  - Homepage
    - isi.edu/~knoblock
  - DIG
    - usc-isi-i2.github.io/dig
  - Karma
    - usc-isi-i2.github.io/karma