SEMANTIC
TECHNOLOGY &
BUSINESS
CONFERENCE
10th ANNUAL
August 19-21, 2014 | San Jose, California USA

YAHOO!
LABS

# Yahoo Knowledge Graph
## Making Knowledge Reusable at Yahoo

PRESENTED BY **Nicolas Torzec**| August 20, 2014

# Background & Context

# Google Knowledge Graph

**Introducing the Knowledge Graph: Things, Not Strings.**

May 16th, 2012

## Brad Pitt

Actor

William Bradley "Brad" Pitt is an American actor and film producer. He has received a Golden Globe Award, a Screen Actors Guild Award, and three Academy Award nominations in acting categories, and ...
Wikipedia

**Born:** December 18, 1963 (age 50), Shawnee, OK

**Height:** 5' 11" (1.80 m)

**Partner:** Angelina Jolie (2005–)

**Spouse:** Jennifer Aniston (m. 2000–2005)

**Children:** Shiloh Nouvel Jolie-Pitt, Vivienne Marcheline Jolie-Pitt, More

### Movies
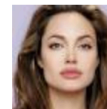View 45+ more

World War Z
2013

The Curious Case of B...
2008

Mr. & Mrs. Smith
2005

Seven
1995

Moneyball
2011

### People also search for
View 15+ more

Angelina Jolie
Partner

Jennifer Aniston
Former spouse

Tom Cruise

George Clooney

Johnny Depp

# Bing Knowledge Graph

**Understand Your World with Bing.**

[March 21st, 2013](March 21st, 2013)

# Yahoo Knowledge Graph



**Yahoo Entity Search.**

[Soft launch in Nov. 2013](#)

# Other "Knowledge Graphs"

| | |
|---|---|
| **Rich, domain-specific, graphs** | Wolfram Alpha, BBC, Rovi, TMS, Baseline, Gracenote, Amazon, Walmart Labs |
| **Interest graphs** | Gravity Adchemy |
| **Social graphs** | Facebook LinkedIn |
| **Reference knowledge graphs** | Freebase + Yago, Wikidata, DBpedia, and other Wikipedia-based projects… |
| **Common-sense knowledge graphs** | Cyc |

YAHOO!

# Scope

# Vision

- **A unified knowledge graph for Yahoo**
  - › All entities and topics relevant to Yahoo (users)
  - › Rich information about entities: facts, relationships, features
  - › Identifiers, interlinking across data sources, and links to relevant services

- **To power knowledge-based services at Yahoo**
  - › **Search:** display, and search for, information about entities
  - › **Discovery:** relate entities, interconnect data sources, link to relevant services
  - › **Understanding:** recognize entities in queries and text

- **Managed and served by a central knowledge team / platform**

YAHOO!

# Value Proposition

- **Data breadth, depth, and accuracy**
  - › Combine information from multiple complementary/overlapping data sources

- **Centralized expertise**
- **Common technologies**     **leveraged across the company**
- **Same knowledge graph**

- **Speed and agility at launching new and richer experiences**

YAHOO!

# In a Nutshell

**Knowledge Acquisition**

**Knowledge Integration**

**Knowledge Consumption**

Ongoing information extraction, from complementary sources.

Reconciliation into a unified knowledge repository.

Enrichment and serving…

YAHOO!

# Making knowledge reusable at Yahoo

# Key Tasks

**Knowledge Repository**
(common ontology)

| Knowledge Acquisition | Knowledge Integration | Knowledge Consumption |
|---|---|---|
| Data Acquisition | Blending | Serving / Export |
| Information Extraction | Entity Reconciliation | Editorial Curation |
| Schema Mapping | | Enrichment |

Data Quality Monitoring

# Data Acquisition

| Knowledge Acquisition | Knowledge Integration | Knowledge Consumption |
|---|---|---|

**Knowledge Acquisition**
- Data Acquisition
- Information Extraction
- Schema Mapping

**Knowledge Integration**
- Blending
- Entity Reconciliation

**Knowledge Consumption**
- Serving
- Export
- Editorial Curation
- Enrichment

**Knowledge Repository**
(common ontology)

Data Quality Monitoring

YAHOO!

# Data Acquisition

- **Multiple complementary data sources**
  - › Combine and cross-validate data from *authoritative\** sources
  - › Reference data sources such as Wikipedia and Freebase form our backbone
  - › Specialized data sources such as TMS and Music Brainz adds breadth/depth
  - › Optimize for relevance, comprehensiveness, correctness, freshness, consistency

- **Ongoing acquisition of raw data**
  - › Feed acquisition from open data sources and paid providers
  - › Web/Targeted crawling, online fetching, ad hoc acquisition (e.g. Wikipedia monitoring)
  - › Deal w/ operational complexity: data size, bandwidth, update frequency, license, ©

YAHOO!

# Information Extraction

| Knowledge Acquisition | Knowledge Integration | Knowledge Consumption |
|---|---|---|
| Data Acquisition | | Serving / Export |
| **Information Extraction** | Blending | Editorial Curation |
| Schema Mapping | Entity Reconciliation | Enrichment |

**Knowledge Repository**
(common ontology)

Data Quality Monitoring

**Knowledge Acquisition**    **Knowledge Integration**    **Knowledge Consumption**

# Information Extraction

- **Extraction of entities, attributes, relationships, features**
  - › Deal w/ scale, volatility, heterogeneity, inconsistency, schema complexity, breakage
  - › Expensive to build and <u>maintain</u> (i.e. declarative rules, expert's knowledge, ML…)
  - › Being able to measure and monitor data quality is key

- **Mixed approach**
  1. Parsing of large data feeds and online data APIs
  2. Structured data extraction on the Web: markup, Web scraping, Wrapper induction,
  3. Wikipedia mining, Web mining, News mining, open information extraction

YAHOO!

# Schema Mapping



**Knowledge Acquisition**
- Data Acquisition
- Information Extraction
- Schema Mapping

**Knowledge Integration**
- Blending
- Entity Reconciliation

**Knowledge Consumption**
- Serving
- Export
- Editorial Curation
- Enrichment

**Knowledge Repository**
(common ontology)

Data Quality Monitoring

# Schema Mapping

- **Normalization to common ontology, schemas, and data types/units**
  - › Upfront normalization: uniform data facilitate downstream usage
  - › Validation against the ontology to ensure well-formedness, validity, and consistency

| | | |
|---|---|---|
| **Ontology alignment** | `<Mad_Men, isA, TVSeries>` | Classifiers: heuristics + ML |
| **Schema mapping** | `<Jon_Hamm, birthplace, St._Louis>` | Template-driven ; mostly declarative |
| **Data normalization** | `<Jon_Hamm, birthdate, "1971-03-10">` | Common plugins |

- **Challenges**
  - › Noisy information extraction: e.g. **strong types vs. inferred types**
  - › Discrepancies between source/target ontologies: e.g. can Pal_(dog) be an actor?
  - › Schema complexity and schema evolutions…

YAHOO!

# Knowledge Representation

| Knowledge Acquisition | Knowledge Integration | Knowledge Consumption |

Data Acquisition

Information Extraction

Schema Mapping

Blending

Entity Reconciliation

Serving    Export

Editorial Curation

Enrichment

**Knowledge Repository**
(common ontology)

Data Quality Monitoring

**Knowledge Acquisition**    **Knowledge Integration**    **Knowledge Consumption**

# Knowledge Representation

- **Property Graph data model**
  - › JSON-LD serialization when needed

- **Common ontology**
  - › OWL ontology. Focuses on representation & validation, not reasoning
  - › Covers domains relevant to Yahoo: 300 classes, 500 object properties, 800 data prop.

**Challenges**
  - › Modeling/managing temporality, provenance, license, localization
  - › Soundness, expressiveness and comprehensiveness … vs. practicality
  - › Collaborative development, conflicting modeling, schema evolution over time

YAHOO!

# Knowledge Repository

| Knowledge Acquisition | Knowledge Integration | Knowledge Consumption |
|---|---|---|

**Knowledge Acquisition**
- Data Acquisition
- Information Extraction
- Schema Mapping

**Knowledge Integration**
- Blending
- Entity Reconciliation

**Knowledge Consumption**
- Serving
- Export
- Editorial Curation
- Enrichment

**Knowledge Repository**
(common ontology)

Data Quality Monitoring

**Knowledge Acquisition**   **Knowledge Integration**   **Knowledge Consumption**

# Knowledge Repository

- **Present knowledge repository backed by a column-oriented store :-(**
  - › Store de-normalized graph persistently and provide some random access via $2^{ndary}$ indices
  - › Scale out nicely and smooth integration with Hadoop workflows
  - › But simplistic data model and limited API make working with graph data tedious

- **Moving to a graph-oriented repository and workflow engine :-)**
  - › Scale to 100s of millions of nodes and billions of facts? (processing, storage, retrieval)
  - › Mix large record-oriented ETL workflows <u>and</u> distributed graph processing?
  - › Efficient graph traversal and query? Built-in inference mechanism?
  - › Schema-less? Data versioning?

**Challenges**

YAHOO!

# Entity Reconciliation & Blending

| Knowledge Acquisition | Knowledge Integration | Knowledge Consumption |
|---|---|---|
| Data Acquisition | **Blending** | Serving / Export |
| Information Extraction | **Entity Reconciliation** | Editorial Curation |
| Schema Mapping | | Enrichment |

**Knowledge Repository**
(common ontology)

Data Quality Monitoring

**Unified Graph**

normalized graph

normalized graph

normalized graph

normalized graph

raw data

raw data

raw data

raw data

Source 4

Source 3

Source 2

Wikipedia

Brad Pitt according to YK

Brad Pitt according to Wikipedia

YAHOO!

# Entity Reconciliation & Blending

- **Disambiguate and merge entities across/within data sources**

| | | |
|---|---|---|
| **Blocking** | Select candidates most likely to refer to the same real world entity | Fast approximate similarity search Hashing techniques |
| **Scoring** | Compute similarity score between all pair of candidates | ML classifier or heuristics |
| **Clustering** | Decide which candidates refer to the same entity and interlink them | ML clustering or heuristics |
| **Merging** | Build a unified object for each cluster. Populate with *best* properties | ML selection or heuristics |

- **Challenges**

  - › Hard Science and Tech problems !
  - › Scale and adapt to new entity types, data sources, data sizes, update frequencies…
  - › Ongoing reconciliation/blending/evaluation. Need for consistent entity IDs. Provenance.

YAHOO!

# Enrichment



**Knowledge Acquisition**
- Data Acquisition
- Information Extraction
- Schema Mapping

**Knowledge Integration**
- Blending
- Entity Reconciliation

**Knowledge Consumption**
- Serving
- Export
- Editorial Curation
- Enrichment

**Knowledge Repository**
(common ontology)

Data Quality Monitoring

YAHOO!

# Enrichment
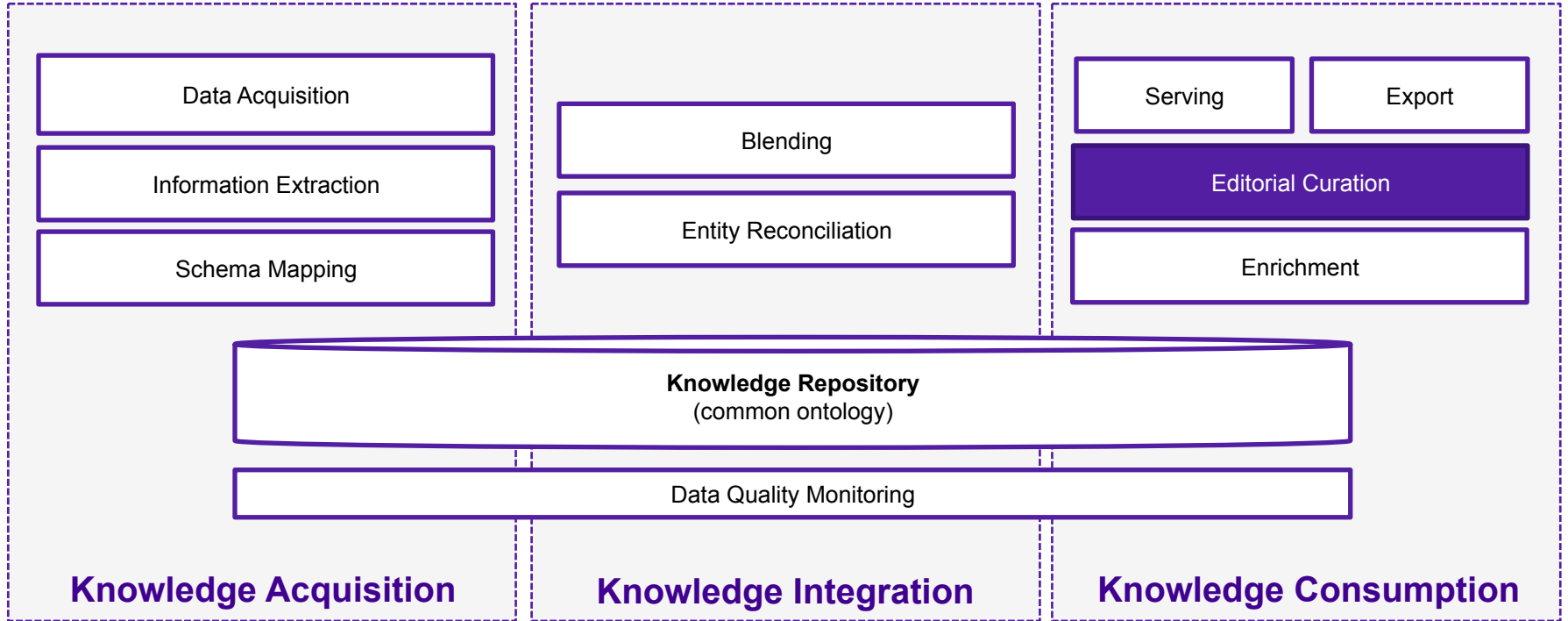
- **Enrich the graph with complementary and/or inferred information**
  - › Generic enrichments vs. context-specific and application-specific enrichments

- **Examples:**
  - › Entity description  cleanup and summarization
  - › Ranking of related entities
  - › Entity categorization

- **Challenges**
  - › Integrating, managing, and running a large number of, possibly conflicting, enrichers.

# Editorial Curation

| Knowledge Acquisition | Knowledge Integration | Knowledge Consumption |
|---|---|---|
| Data Acquisition | | Serving / Export |
| Information Extraction | Blending | **Editorial Curation** |
| Schema Mapping | Entity Reconciliation | Enrichment |

**Knowledge Repository**
(common ontology)

Data Quality Monitoring

**Knowledge Acquisition**    **Knowledge Integration**    **Knowledge Consumption**

YAHOO!
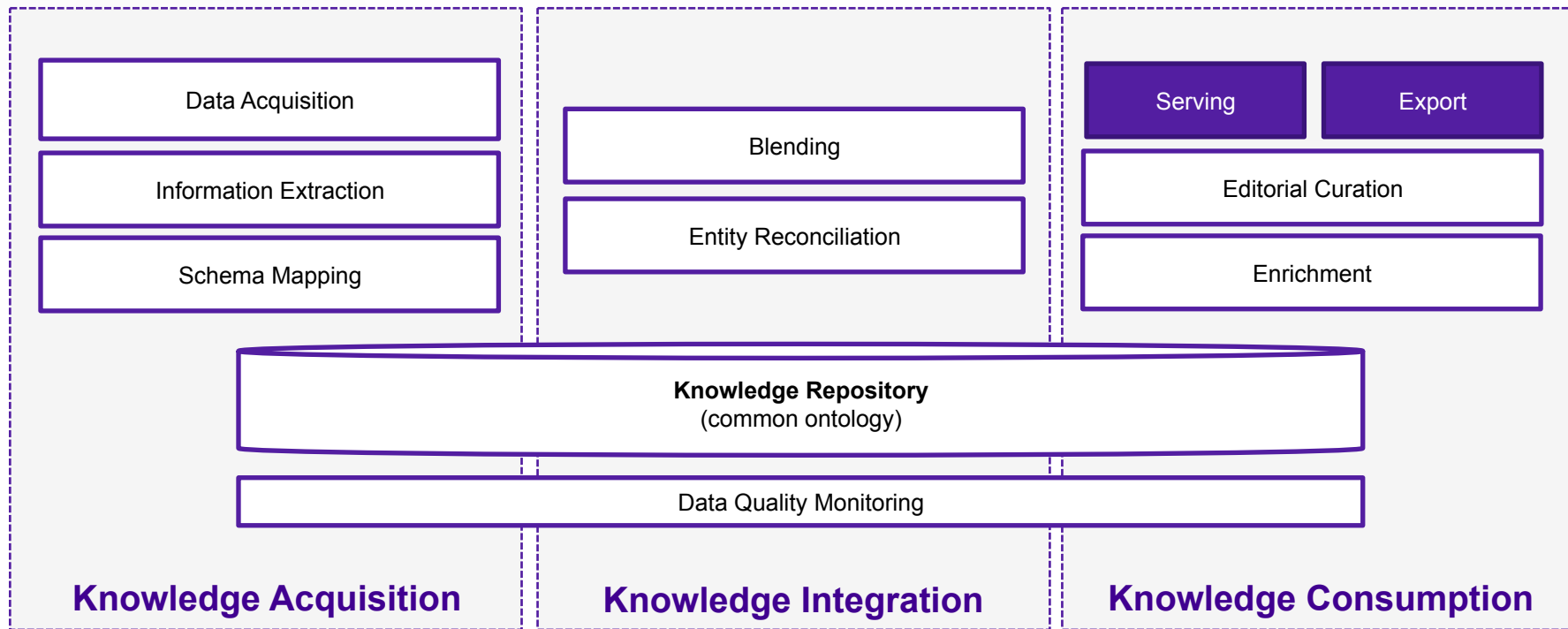
# Editorial Curation

- **Enable editors to perform hot fixes**
  - › Interactive (and safe) GUI for updating entities and associated information

- **Internal Wall of Shame**
  - › Typical issues: incorrect/outdated facts, images, categorization (examples below)
  - › Occasionally some reconciliation issues: Frankenstein objects!

- **Challenges**
  - › Instantly reflect editorial updates in knowledge graph and consuming systems
  - › Re-evaluate and manage editorial updates over time since they typically blindly overwrite
  - › Manage multiple concurrent, and possibly conflicting, editorial updates.

# Serving & Publishing

| | | |
|---|---|---|
| **Data Acquisition** | | **Serving**  **Export** |
| **Information Extraction** | **Blending** | **Editorial Curation** |
| **Schema Mapping** | **Entity Reconciliation** | **Enrichment** |

**Knowledge Repository**
(common ontology)

Data Quality Monitoring

**Knowledge Acquisition**       **Knowledge Integration**       **Knowledge Consumption**

YAHOO!

# Serving & Publishing

- **Online serving**
  - › Dedicated serving infrastructure powering various online data APIs
  - › Search layer provides efficient random access to the graph (and limited traversal)
  - › Federation layer integrates transient info from connected services at query time
  - › Customization layer provides attribute-level filtering, transformation, formatting

- **Datapack generation**
  - › Regular datapack generation for offline batch consumption
  - › Typically one single generic datapack with all the data

YAHOO!

# Knowledge-based services at Yahoo

- **Search:**
  - › display, and search for, information about entities

- **Discovery:**
  - › relate entities, interconnect data sources, link to relevant services

- **Understanding:**
  - › recognize entities in queries and text

**YAHOO!**

# Yahoo Knowledge Graph

## Making Knowledge Reusable

# Thank you.

torzecn@yahoo-inc.com
Twitter: nicolastorzec