The MIT Information Quality Industry Symposium, 2007

# Demonstrations of linguistic data matching, consolidation, and cleansing

**Jeff Fried**
**VP Advanced Solutions**
**Jeff.Fried@fastsearch.com**

**:::fast**

The MIT Information Quality Industry Symposium, 2007

# Data Quality

- Still people entering the data
  - ## We do typos
  - ## We duplicate
  - ## We mess up

# Data Consolidation and ETL Studio

**Direct access to RDBMs for info from some Telco's**

**XML feed from other Telco's**

**XML**

**Flat files (CSV or fixed) from the 'laggards'**
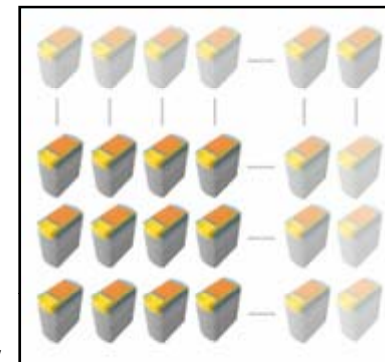
**Logic for matching and cleansing**

**Logic for ETL**

**Lookup to ESP for Matching**

**Ordered hits (by quality)**

**Cleansed data Back to ESP for next round Cleansing or for 'search'**

**clean data**

**Master database for persistant storage**
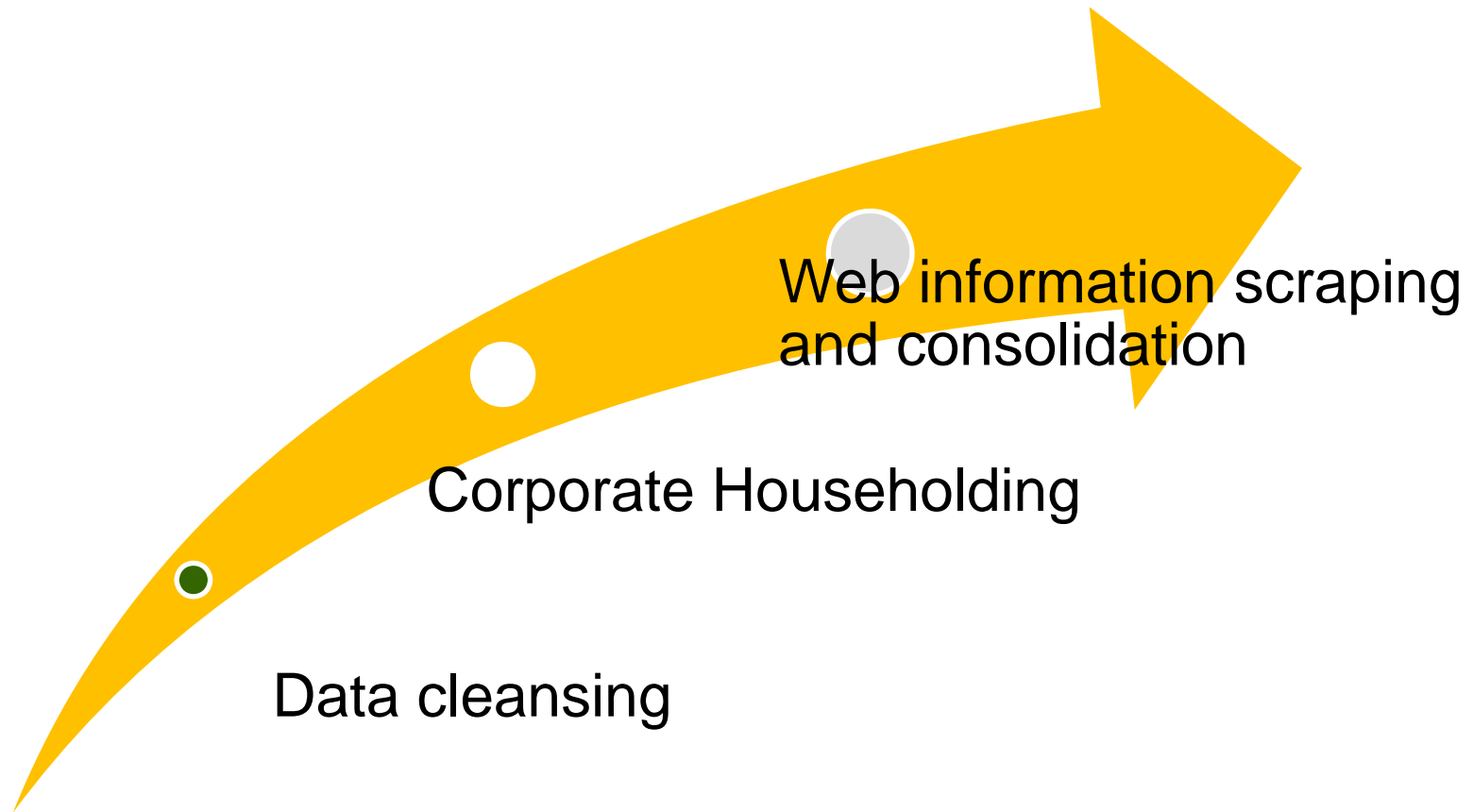
**Ambigous data (close hits or unidentified)**

**'Error' database for manual inspection, correction, storage/learning**

**⁞⁞⁞fast**

The MIT Information Quality Industry Symposium, 2007

Web information scraping and consolidation

Corporate Householding

Data cleansing

:::fast

# The MIT Information Quality Industry Symposium, 2007

## Viewing data, we notice issues

# The MIT Information Quality Industry Symposium, 2007

## We click on the gauge

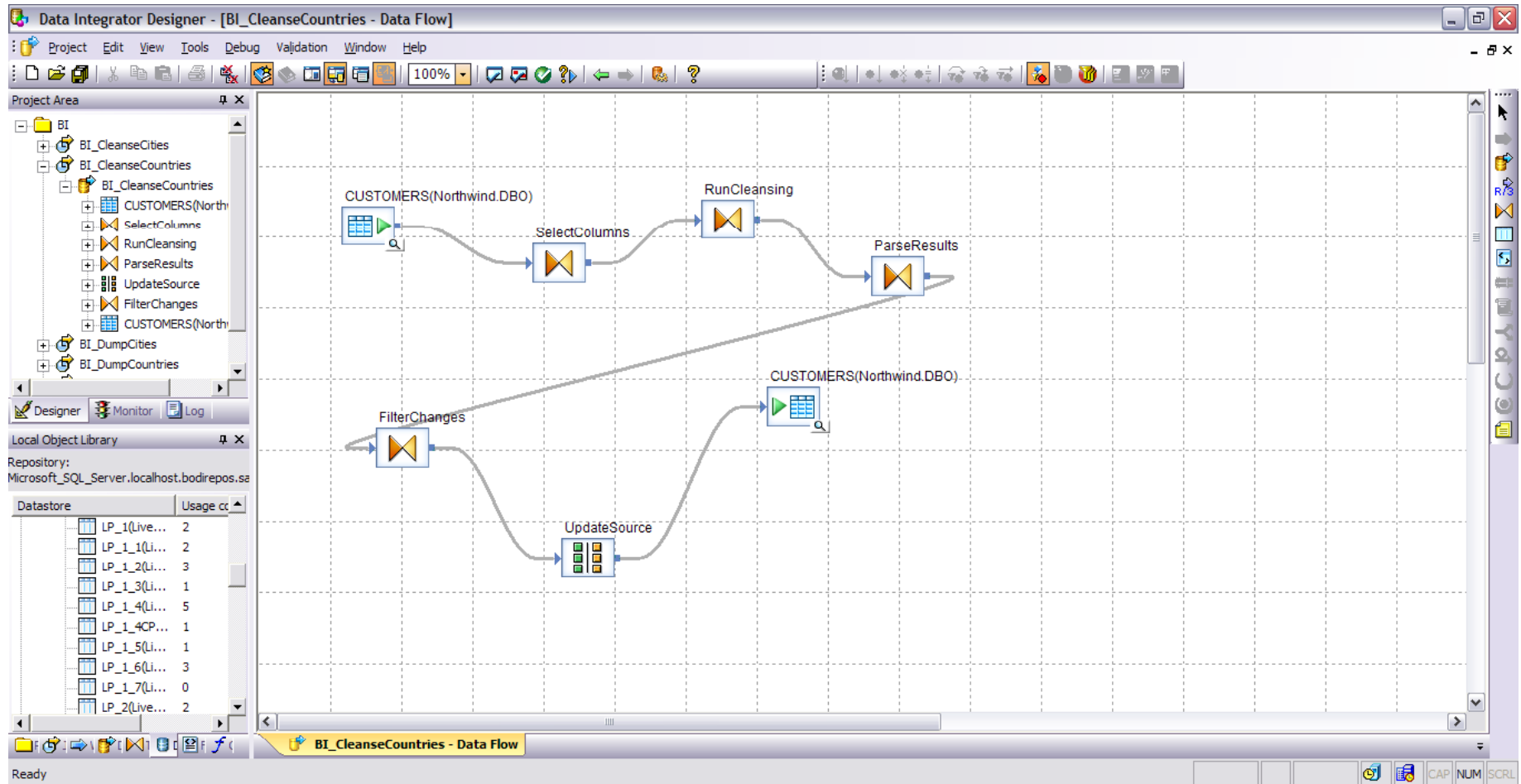The MIT Information Quality Industry Symposium, 2007

# We cleanse the data

The MIT Information Quality Industry Symposium, 2007
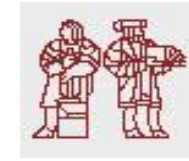
# Now look at the result

The MIT Information Quality Industry Symposium, 2007

# Details

The MIT Information Quality Industry Symposium, 2007

# Linguistic Fundamentals

## Lexicon Base

- Special terminology lexica
- Geographical and people's names
- Spellcheck dictionaries
- Subject-specific ontologies
- Synonymy Dictionaries
- Part-of-speech Dictionaries
- Inflection Dictionaries
- Language-specific Common Words

## Basic Linguistic Algorithms

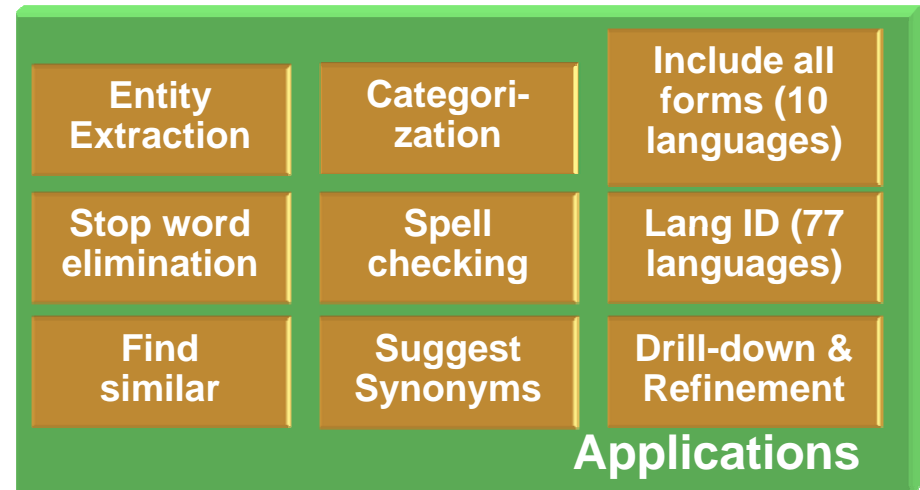| Vectorization | Part-of-speech Tagging |
|---|---|
| Language normalization | Pattern extraction |
| | Stemming / Lemmatization |

## Applications

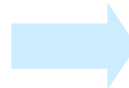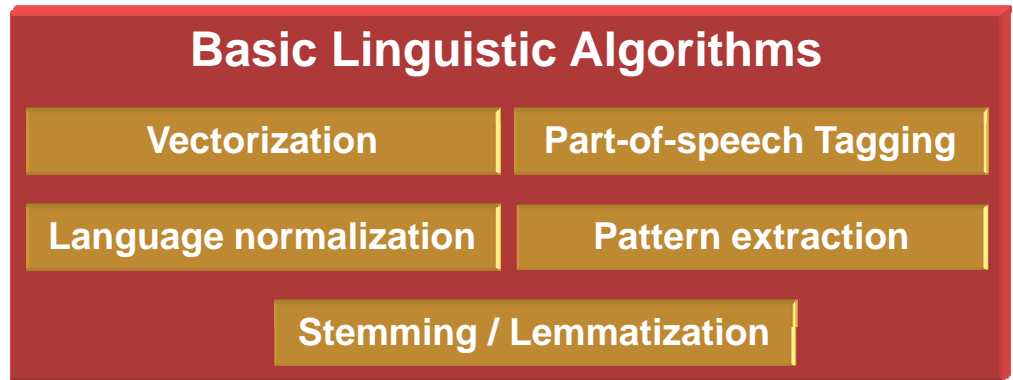| Entity Extraction | Categori-zation | Include all forms (10 languages) |
|---|---|---|
| Stop word elimination | Spell checking | Lang ID (77 languages) |
| Find similar | Suggest Synonyms | Drill-down & Refinement |

::fast

**11**

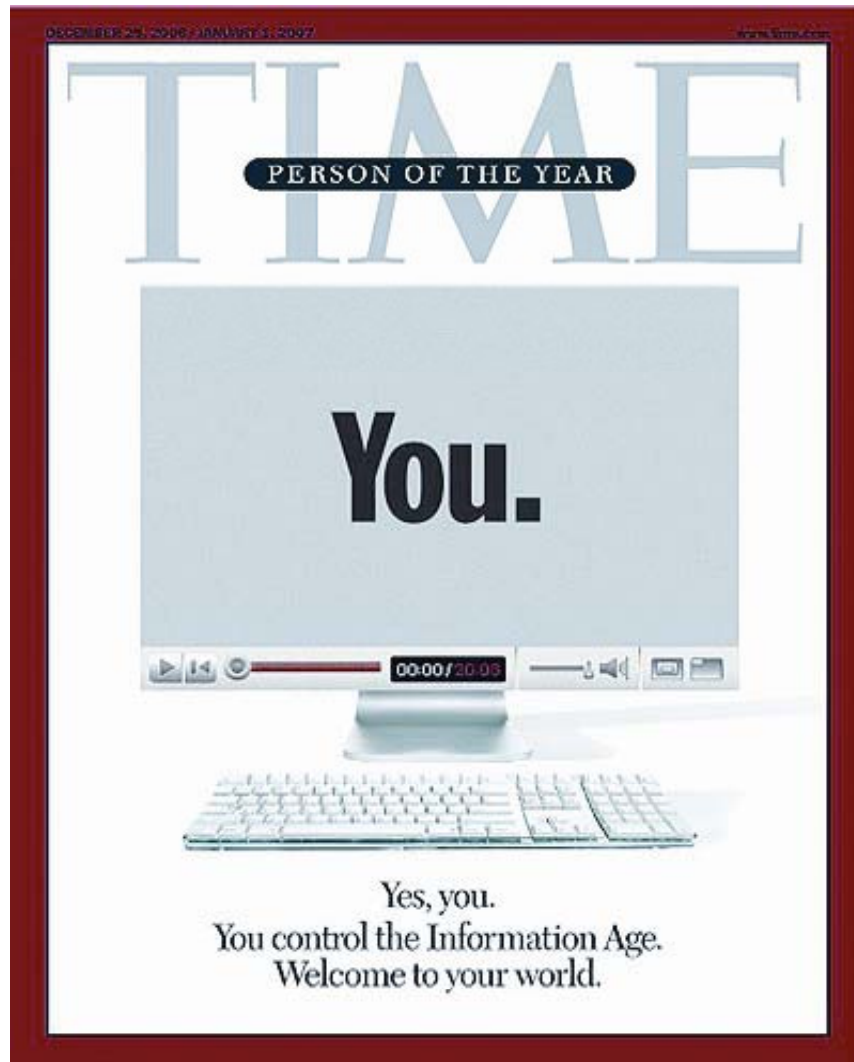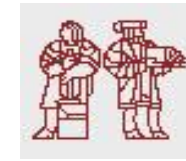The MIT Information Quality Industry Symposium, 2007

# Q&A

The MIT Information Quality Industry Symposium, 2007

# Thank you!

**find the real value of search**

**Jeff Fried**
**VP Advanced Solutions**
**Jeff.Fried@fastsearch.com**

**13**