# Automating the Extraction of Genealogical Information from the Web

GeneTIQS

Troy Walker & David W. Embley
Family History Technology Conference
March 25, 2004

# Genealogical Information on the Web

- **Hundreds of thousands of sites**
  - Some professional (Ancestry.com, Familysearch.org)
  - Mostly hobbyist (203,200 indexed by Cyndislist.com)
- **Search engines**
  - "Walker genealogy" on Google: 199,000 results
  - 1 page/minute = 5 months to go through
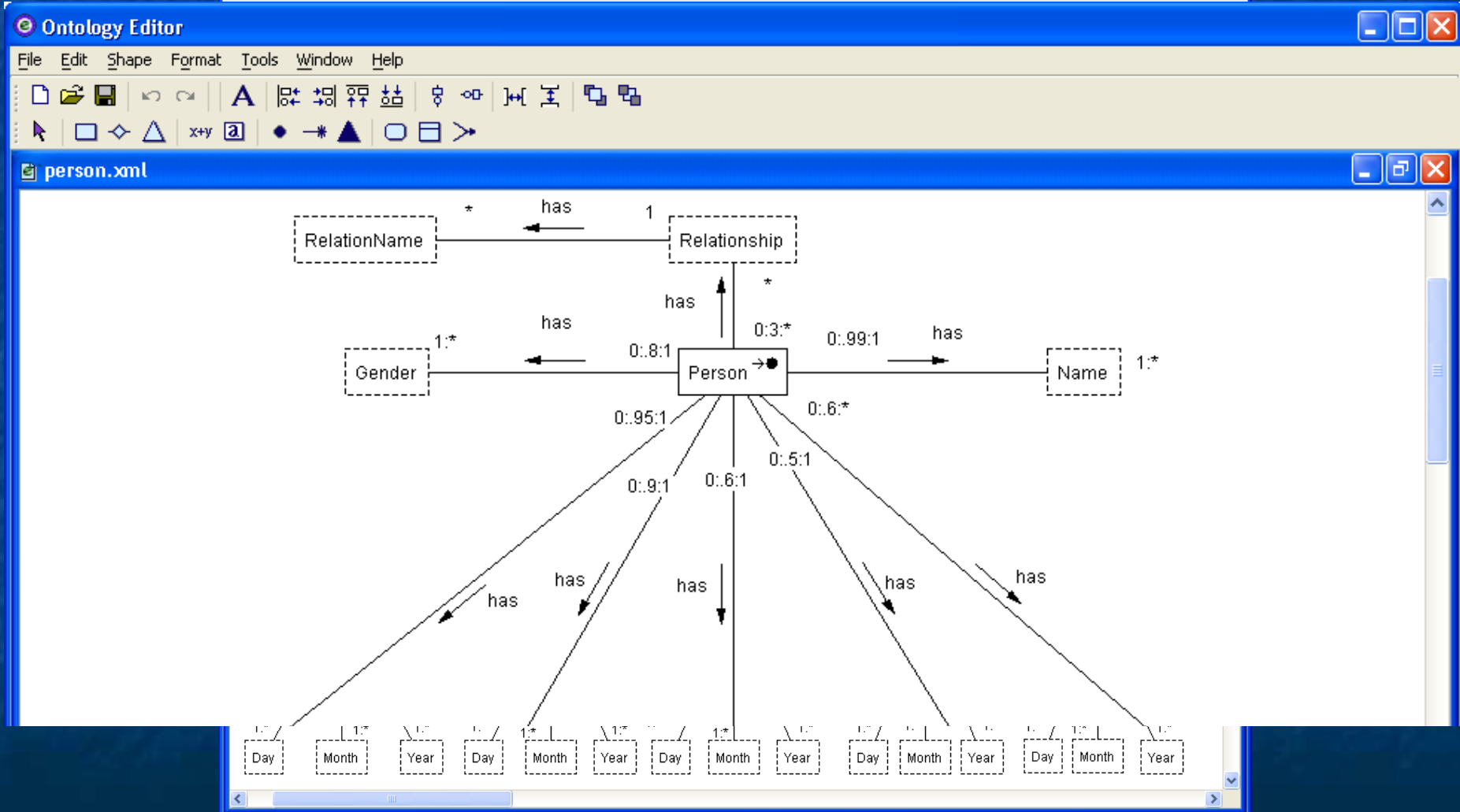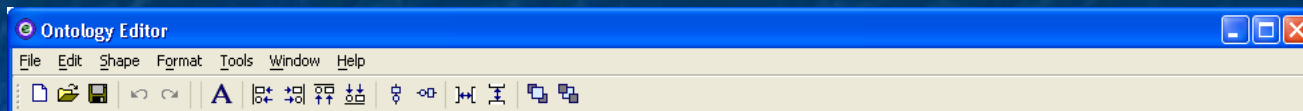- **Why not enlist the help of a computer?**

# Problems

- No standard way of presenting data
- Sites have differing schemas
- Web pages change
- New pages continuously come on line

# GeneTIQS

- **Based on work done by BYU DEG**
- **Able to extract from:**
  - Single-record documents
  - Simple multiple-record documents
  - Complex multiple-record documents
- **Robust to changes in pages**
- **Immediately works for new pages**

# Person Ontology

# Value Matchers

# Record Separation

- **Separating data related to each person**
- **Previous technique**
  - Combines many heuristics
  - Has problems
    - Assumes multiple records
    - Must be simple separation

# Single-Record Document

# Simple Multiple-Record Document



9

# Complex Multiple-Record Document

# Vector Space Modeling

- **Ontology Vector**
- **Compare to candidate records**
  - Cosine measure

$$v_1 \bullet v_2$$

  - Magnitude measure

$$\left| v_1 \right|$$

# Ontology Vector



$$\{\ 0.8,\ 0.99,\ 0.95,\ 0.9,\ 0.6,\ 0.5,\ 0.6,\ 3.0\}$$

# Vector Space Modeling

```
<!DOCTYPE…>          {0, 0, 0, 0, 0, 0, 0, 0}
<html>               {0, 141, 89, 76, 0, 0, 48, 23}
  <head>             {0, 1, 0, 0, 0, 0, 0, 0}
    …
  </head>
  <body>             {0, 140, 89, 76, 0, 0, 48, 23}
   <div>             {0, 0, 0, 0, 0, 0, 0, 0}
    …header…
   </div>
   <div>             {0, 138, 88, 76, 0, 0, 48, 23}
    …                …
```

Gender

Name

Birth

Death

Christening

Burial

Marriage

Relation
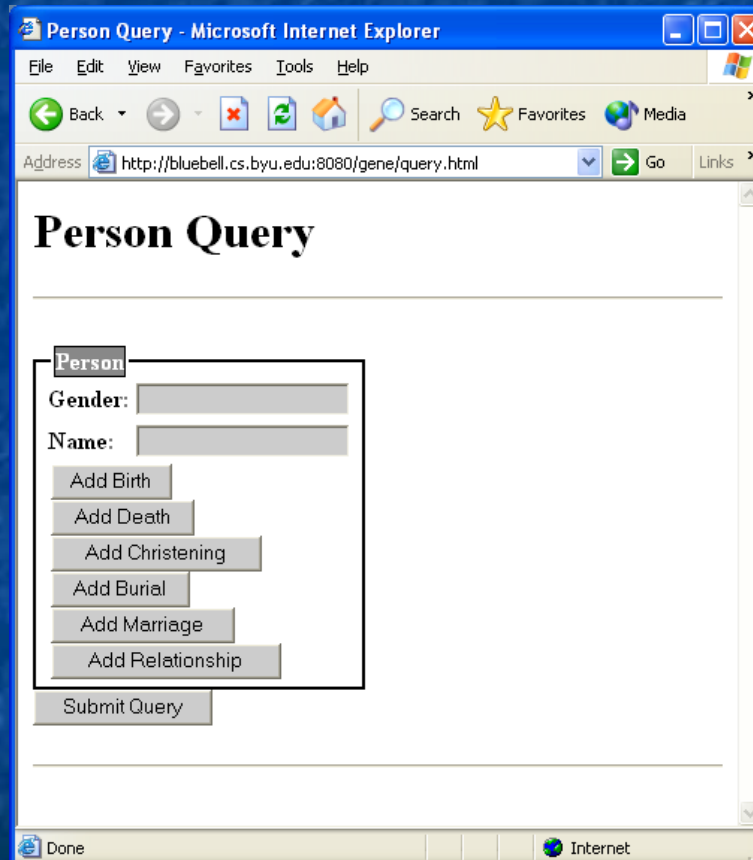
# Improvements

- **Differing schemas**
  - Low cosine measures
  - Discarded data
  - Prune dimensions

    {0.8,    0.99, 0.95, 0.9, 0.6, 0.5,   0.6,  3.0}

    {0.0, 141.0, 89.0, 76.0, 0.0, 0.0, 48.0, 23.0}

- **Richness of data in single-record documents**
  - High magnitude measure
  - Higher magnitude to split documents

# Demonstration

# Presenting Results

# Preliminary Results

- **Semi-structured Text**

  - 10 single-record documents
  - 3 simple documents containing 268 records
  - 3 complex documents containing 266 records

- **Precision and recall for record separation**

# Record Separation

|         | Recall | Precision |
|---------|--------|-----------|
| Single  | 100%   | 94.1%     |
| Simple  | 94.7%  | 97.3%     |
| Complex | 88.3%  | 93.6%     |

# Conclusion

- **Integrate, build on previous DEG work**
- **Accurate record separation**
  - Average recall: 94.3%
  - Average precision: 95.0%
- **Ontology based**
  - Robust to changes in pages
  - Immediately works with new pages