# Domain-Independent Data Extraction: Person Names

Carl Christensen and Deryle Lonsdale
*Brigham Young University*
cvchristensen@gmail
lonz@byu.edu

# Challenge

- Extraction software and techniques yield good results with domain specific data extraction
- Person names and information rarely domain specific
- Identification and extraction difficult because of noisy data, lack of formatting

# WePS task

- Web People Search
- 18 attribute values on person names
  - Date of birth, Birth place, Other name, Occupation, Affiliation, Work, Award, School, Major, Degree, Mentor, Location, Nationality, Relatives, Phone, FAX, Email, Web site
- Training corpus – 17 names, approx. 100 web pages per
- Script given to evaluate performance
- Test corpus of comparable size
- New ground in information extraction

# Ontos

- Software developed by BYU data extraction group
- Ontology based method leveraged to organize data
- Off the shelf performance
- Similar uses for obituaries and car ads information extraction
  - Good performance on these tasks

# Ontos obituary results

Beloved son, brother and uncle, Robert Gene Larkins

(obituary,DeceasedName,2,0)

, 46, passed away Wednesday, October 6, 2004, at his home in Goshen, Utah. Bobby was born December 24, 1957, in Riverside, California. He was the fifth of six children; third son of four, born to Robert A. and Je... Ramona High School in Arlington, California. Bobby's passion was anything that would go FAST. When he was 14 years old, the local police pulled him over and issued him a citation for going above the speed limit on ... people, especially relatives, respectfully referred to Bobby as, "McGuiver." He could fix EVERYTHING with almost ANYTHING; his specialty being cars and motorcycles. Many times he would receive telephone calls ... very adept at diagnosing the problem and would tell the caller just what they needed to do to fix it. Bobby was an incredible human being. He was the type of person who would give you the shirt off his back if he thoug... "the going rate" for all the vehicle repairs he made for people. But instead, he chose to serve those who couldn't afford to pay. He will be greatly missed by so many. His absence has already been felt, and will continue to ... by many loved ones who have gone on before him; including his older brother Timothy Fred, who passed away two years ago. Also, his beloved dog, Patches. He is survived by his parents, and siblings, Judy Tatum, of ... Peterson, of Alpine, Utah; John Larkins, of Albany, Oregon; also, many nephews, nieces and friends, who will miss his humor, advice and presence. We appreciate the Warenski Funeral Home for their assistance. May ... HEAL." Dear Bobby, we celebrate your life and the great love you shared with us. We know that we shall see you again in time. Funeral services for Robert will be held Monday, October 11, 2004 at 11:00 a.m., at 17... blocks north of the Mt. Timpanogos Temple). Family and friends are invited to call, from 10:00 to 11:00 a.m. Interment for Robert will take place at the Alpine City Cemetery, located at 350 North Grove Street, Alpine...

William Keith Romney, 75, died Wednesday, October 6, 2004 after a lingering illness. He was born on January 17, 1929 in Salt Lake City, Utah, to Lois and Carl Romney. Bill will be remembered for his dedication to ... his children his love for skiing, golf, and his appreciation of music and art. He is survived by his soul mate, Joan L. Romney, two sons: William Bruce Bach, Coeur d'Alene, ID; Michael (Keri) Keith Romney, Ogden; and ... Jordan; Heather (Peter) Romney Redgrave, Berlin, Germany; and Laurel Romney, Alta; 12 grandchildren, and three great-grandchildren. Bill is also survived by his brother, Carl (Barbara) Romney, Alexandria, VA; sist... and nephews. He is preceded in death by his parents, sister Marie, brother Bob, daughter Valerie, and beloved companion Putter. The family would like to express their appreciation to the V.A. hospital, and his health c... Bill's life at his home on Saturday, October 9, 2004 from 1:00-3:00 p.m. In lieu of flowers, please send a donation to the Red Cross in Bill's name. Funeral directors, Myers Mortuary of Ogden. Send condolences to fam...
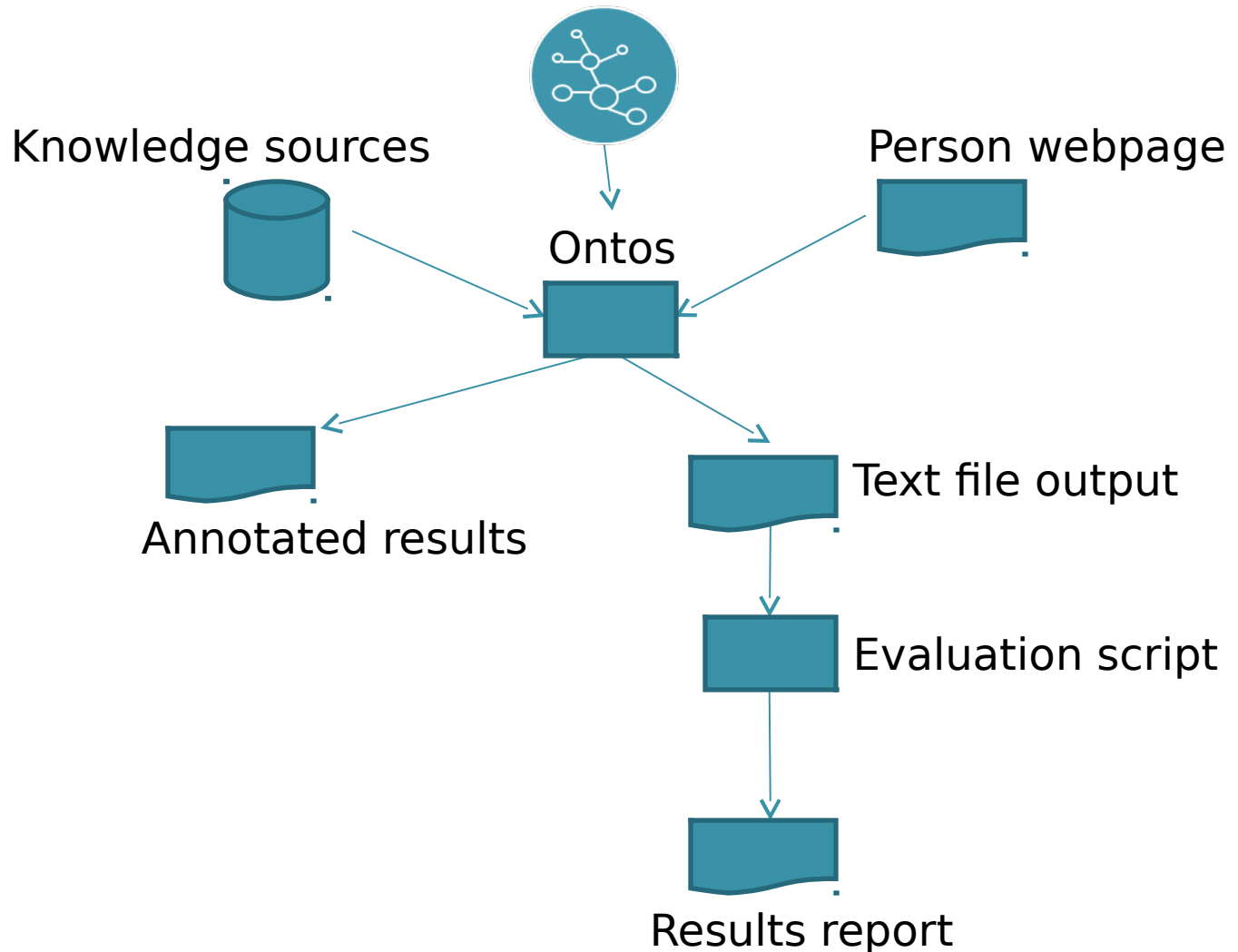
Ernesto J. Sandoval, 80, passed away October 7, 2004 at his home surrounded by his family. Born November 9, 1923 to Adelaido and Onofre Aguilar Sandoval in La Puente, New Mexico. Survived by his daughters ...

## Annotation Results

| DeceasedName | Age | DeathDate | BirthDate | Relationship | IntermentDate | IntermentAddress | FuneralAddress | FuneralTime | FuneralDate | V |
|---|---|---|---|---|---|---|---|---|---|---|
| Robert Gene Larkins | 46 | October 6, 2004 | December 24, 1957 | Show | | 350 North Grove Street, Alpine, Utah | 1776 North 900 East, American Fork, Utah | 11:00 | October 11, 2004 | |

| Relationship | RelativeName |
|---|---|
| sibling | Gloria Larkins |
| parent | Jean Larkins |
| neices | CANNOT HEAL |

| ViewingDate | ViewingAddress | BeginningTime | EndingTime |
|---|---|---|---|
| | Goshen, Utah | 10:00 | 11:00 a.m. |

| DeceasedName | Age | DeathDate | BirthDate | Relationship | IntermentDate | IntermentAddress | FuneralAddress | FuneralTime | FuneralDate |
|---|---|---|---|---|---|---|---|---|---|
| William Keith Romney | 75 | October 6, 2004 | January 17, 1929 | Show | | Alene, ID | Salt Lake City, Utah | 1:00 | October 9, 2004 |
| Ernesto J. Sandoval | 80 | October 7, 2004 | November 9, 1923 | Show | | La Puente, New Mexico | 174 East 900 South | 10:00 | |

# WePS extraction process

WePS ontology

Knowledge sources

Person webpage

Ontos

Annotated results

Text file output

Evaluation script

Results report

# Data frames

```xml
<ObjectSet order="5" x="179" y="14" lexical="true" name="BirthDate" id="osmx1" type="String">
    <DataFrame>
        <InternalRepresentation>
            <DataType typeName="java.lang.String"/>
        </InternalRepresentation>
        <ValuePhrase hint="date_9-DDMonth" caseSensitive="false" confidenceTag="">
            <ValueExpression>
                <ExpressionText>(1\d|2\d|30|31|\d)\s*{Month}\.?</ExpressionText>
            </ValueExpression>
            <RequiredContextExpression> <ExpressionText>\b{BirthTime}\b</ExpressionText></RequiredContextExpression>
        </ValuePhrase>
        <ValuePhrase hint="date_6-DDMonthYYYY" caseSensitive="false" confidenceTag="">
            <ValueExpression>
                <ExpressionText>(1\d|2\d|30|31|\d)\s*{Month}\.?\s*(\d\d\d\d)</ExpressionText>
            </ValueExpression>
            <RequiredContextExpression> <ExpressionText>\b{BirthTime}\b</ExpressionText></RequiredContextExpression>
        </ValuePhrase>
        <ValuePhrase hint="date_10-MonthDD" caseSensitive="false" confidenceTag="">
            <ValueExpression>
                <ExpressionText>{Month}\s+[0-3]?\d</ExpressionText>
            </ValueExpression>
            <LeftContextExpression>
                <ExpressionText>\b</ExpressionText>
            </LeftContextExpression>
            <RightContextExpression>
                <ExpressionText>(\s|\.)</ExpressionText>
```

- XML description of extraction ontology components

# Knowledge files

- Names, cities, countries, hypocoristics, occupations, d

- Knowledge gathered from extracting and formatting o
databases
  - - Live Journal
  - - Wikipedia
  - - Bureau of Labor Statistics
  - - etc.

- Approximately 80,000 school names
and 30,000 occupations
  - - 66% of total schools in U.S. in 2003

- All possible options for some files, small
subset for others
  - - e.g. Occupations, hypocoristics

```
James Rutter Junior High School
James W. Marshall Elementary School
Jefferson Elementary School
Jesuit High School
John Bidwell Elementary School
John Cabrillo Elementary School
John F. Kennedy High School
John H. Still High School
John H. Still Junior High School
John Reith Elementary School
John Sloat Elementary School
Jonas Salk Middle School
Keema High School
Kit Carson Middle School
La Entrada Continuation High School
Las Flores High School
Lawrence Junior High School
Leonardo Da Vinci K-8 Magnet School
Leroy F. Greene Middle School
Lincoln Continuation High School
Lincoln Law School
Lisbon Elementary School
Little John Elementary School
Loretto High School
Luther Burbank High School
Madison Elementary School
Maric College - Sacramento Campus
Maric College at Sacramento Campus
Mariemont Elementary School
Mark Hopkins Elementary School
Martin Luther King Junior High School
Maryukyokoto Elementary School
Maryland Elementary School
McGeorge School of Law
Merryhill School
Met Sacramento Charter High School
Michael J. Castori Elementary School
Mira Loma High School
Mission Avenue Open School
Natomas Charter School - Leading Edge
Natomas Charter School at Leading Edge
Natomas Charter School - Performing and Fine Arts Academy
Natomas Charter School at Performing and Fine Arts Academy
Natomas High School
Natomas Middle School
Nicholas Elementary School
Norte Del Rio High School
Northwestern California University, School of Law
O.W. Erlewine Elementary School
Oak Avenue Elementary School
Oak Ridge Elementary School
Orville Wright Elementary School
Our Savior Lutheran School
Peter Burnett Elementary School
Phoebe Hearst Elementary School
Pony Express Elementary School
```

# Constraints

- Required context expressions

```
<RequiredContextExpression>
<ExpressionText>\b{BirthTime}\b</Expre
ssionText></RequiredContextExpression>
```

- Cardinality

```
Search Person [0:*] has Occupation
[1:*];
Search Person [0:*] has Affiliation
[1:*];
Search Person [0:*] has Email [1];
```

- Regular expressions

```
Email : [\w]+[\d]*@[\w]+[.]{1}[\w]+
```

# Sample webpage

**Grand Lodge**
OF BRITISH COLUMBIA AND YUKON

# King David Lodge No. 93



**Work:**
Canadian

**Installation of officers:**
December

**District:**
District No. 17

**Lodgehall:**
1763 Bellevue Ave., West Vancouver

**Meeting day:**
1st Thursday

**Website:**
www.kingdavidlodge93.org

# Family Photo Album Dec02/Jan03

**CLICK ON THE PICTURES BELOW TO SEE THE LARGER PHOTO**

At Home.............hanging around.........



On Holidays, December/January 2003........horsing around............

# Sample annotated webpage

# Precision/Recall

$$Recall = \frac{|\{relevant\ information\} \cap \{retrieved\ information\}|}{|\{relevant\ information\}|}$$

$$Precision = \frac{|\{relevant\ information\} \cap \{retrieved\ information\}|}{|\{retrieved\ information\}|}$$

$$F_{alpha} = \frac{(1 + alpha) * precision * recall}{((alpha * precision) + recall)}$$

alpha = 0.5 for WePS

# WePS2-AE matching Result for

------------------precision=9.68320382546324 recall=37.7622377622378----------

| MATCH | MISS1 | MISS2 | Precision | Recall | F-measure |
|---|---|---|---|---|---|
| 162 | 1511 | 267 | 9.683 | 37.762 | 15.414 |

| Attribute | MATCH | MISS1 | MISS2 |
|---|---|---|---|
| affiliation | ○ Royal Society of Literature<br>○ University of Birmingham | ○ Birmingham<br>○ History<br>○ Home | |
| award | ○ \|CBE\|/\| CBE\|<br>○ \|Hawthornden Prize\|/\| Hawthornden Prize\|<br>○ \|Royal Television Society's Award for best Drama serial in the year 1989\|/\| Royal Television Society's Award for best Drama serial in the year 1989\|<br>○ \|Yorkshire Post Fiction Prize\|/\| Yorkshire Post Fiction Prize\| | ○ Booker Prize<br>○ Commonwealth Writers Prize<br>○ Criticism<br>○ Nice Work<br>○ Romance<br>○ Silver Nymph<br>○ Sunday Express Book of the Year | ○ Silver Nymph at the International Television Festival<br>○ Sunday Express Book of the Year Award<br>○ Whitbread Book of the Year |
| birthplace | | ○ MA | ○ London, England |
| dateofbirth | | | ○ January 28, 1935 |
| degree | ○ BA<br>○ PhD | | ○ MA |
| major | | ○ Modern<br>○ humanities | |
| nationality | ○ British | | |
| occupation | ○ Fellow<br>○ Honorary Professor<br>○ author<br>○ full-time writer | ○ Author<br>○ Booker<br>○ Novelist<br>○ Reader<br>○ Writers<br>○ literary critics<br>○ page<br>○ satirists | |
| school | ○ \|University of Birmingham\|/\| University of Birmingham\| | | ○ University College London |
| website | | ○ http://en.wikipedia.org/wiki/David_Lodge_%28author%29" | |

# Challenges

- Smaller/larger match preference
    - Preference for "Arizona" as place over "University of Arizona" as school
- DOM parser
    - Unofficial HTML tags cause system to fail
- Text formatting
    - Record detection for individuals intractable
- System functionality
    - Cardinality bounds, system output file

# Performance

- Very low initial precision/recall: < 1%
- Increased drastically with knowledge engineering and system constraints

  - 27 % recall with some person results approaching 40% recall

  - Approaching 10% precision

- Nothing to measure against

  - Official WePS results will be released            in April

# Future work

- Ontos robustness
- Machine learning for constraints/knowledge files
- Person name disambiguation
- Keyword probability values