

Unsupervised Information Extraction by Text Segmentation

Eli Cortez

Advisor: Altigran Soares da Silva

Universidade Federal do Amazonas



Agenda

- ▶ **Introduction**

- ▶ Information Extraction by Text Segmentation (IETS)
- ▶ Contributions

- ▶ **Related Work**

- ▶ Web Extraction Methods and Tools
- ▶ Probabilistic Graph-Based Methods

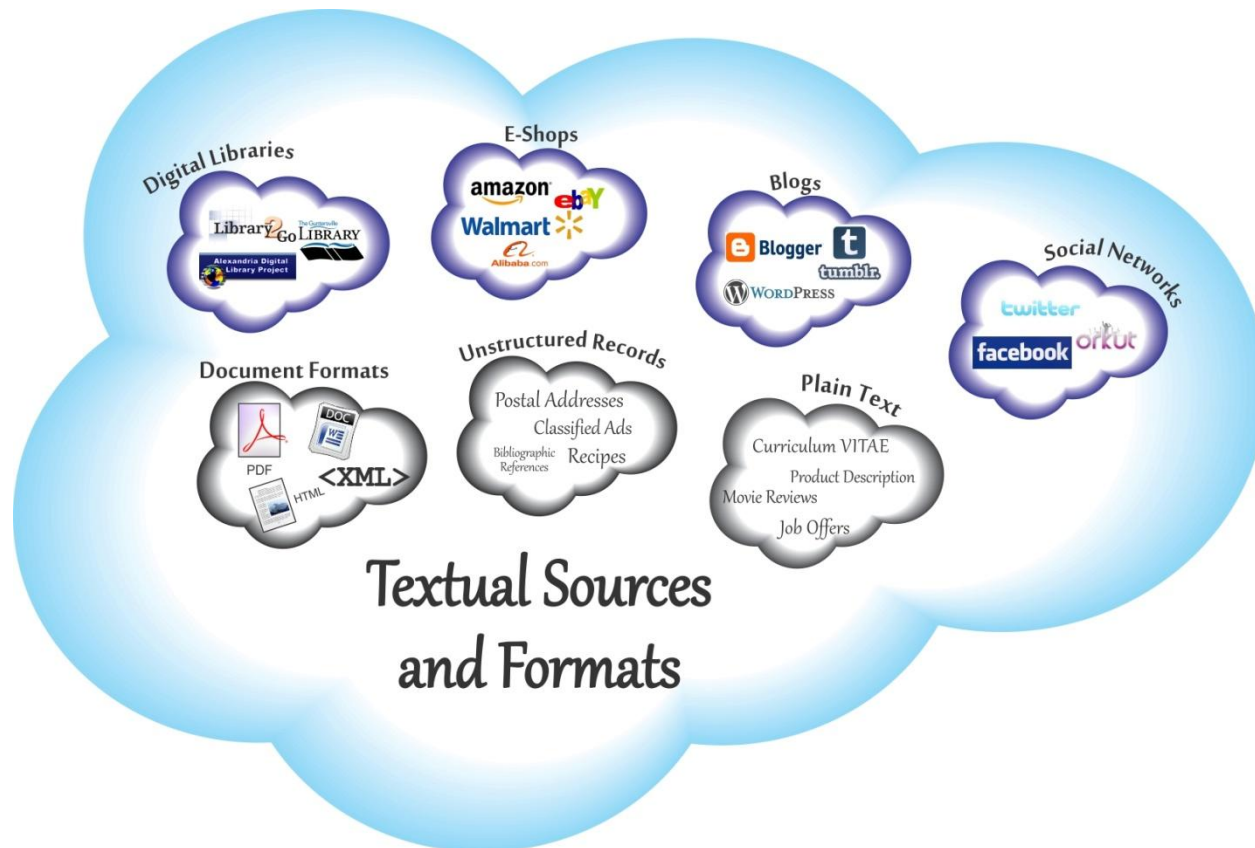
- ▶ **Our Proposed Approach for IETS**

- ▶ Ondux
- ▶ Judie
- ▶ iForm

- ▶ **Conclusions and Future Work**

Introduction

- ▶ Steady increasing in the number and the types of sources of textual information available in the World-Wide Web



Introduction

- ▶ These sources constitute large repositories of valuable data on a variety of domains.
- ▶ Data referring to different “things” such as:
 - ▶ Personal Information;
 - ▶ Products;
 - ▶ Publication;
 - ▶ Companies;
 - ▶ Cities;

Introduction

- ▶ Important restrictions on the way data they contain can be manipulated.
- ▶ Text snippets (product descriptions, movie reviews) can hardly be subject to **automated processing**.
 - ▶ Difficult to automatically identify data of interest.

The screenshot shows the Allrecipes.com website interface. At the top, there is a search bar with the text "Example: cupcakes" and a "Search" button. Below the search bar, there are navigation tabs for "new at", "recipes", "videos", "menus", "holidays", "thebuzz", and "join for FREE". The main content area features a recipe for "Chocolate Graham Nut Cake" by Carol. The recipe includes a description, a star rating, and a "Rate/Review" button. The ingredients list is highlighted with a red box and labeled "Unstructured Records". The ingredients list is as follows:

Quantity	Unit	Ingredient
6	eggs	
1/2	teaspoon	cream of tartar
1	cup	white sugar
1/2	teaspoon	vanilla extract
1/2	cup	finely ground graham cracker crumbs
1	teaspoon	baking powder
1/4	teaspoon	salt
3/4	cup	ground walnuts
3/4	cup	semisweet chocolate chips

The directions section is also visible, starting with "1. In a large bowl, beat the egg whites with the cream of tartar until stiff. Gradually beat in 1/4 cup sugar until the mixture is slightly glossy. Set aside."

Introduction

- ▶ **The Information Extraction (IE) Problem**
 - ▶ Automatic extract **structured information** such as **entities, relationships** between entities, and **attributes** describing entities from **noisy unstructured sources**.
 - ▶ Named Entity Recognition;
 - ▶ Open Information Extraction;
 - ▶ Relationship Extraction;
 - ▶ **Information Extraction by Text Segmentation (IETS)**

Introduction

- ▶ Information Extraction by Text Segmentation (IETS)
 - ▶ The problem of **extracting attribute values** occurring in implicit **semi-structured data records** in the form of continuous text.

Eli Cortez - Rua 15 n 324 Japiim 1 - 69075 - Manaus

Name	Street	Number	Neigh.	Zip	City
Eli Cortez	Rua 15	n 324	Japiim 1	69075	Manaus

- ▶ Why is it important to extract information?
 - ▶ Query structured data; Data Mining; Record Linkage.

Introduction

▶ Contributions

- ▶ In this work we tackle the Information Extraction by Text Segmentation Problem (IETS)
 - ▶ Important and practical problem frequently addressed in the recent literature.
 - Borkar@SIGMOD'01, McCallum@ICML'01, Agichtein@SIGKDD'04, Mansuri@ICDE'06, Zhao@SICDM'08, Cortez@JASIST'09
- ▶ We propose and implement an **unsupervised approach** to this problem.
 - ▶ Relies on information available on pre-existing data.
 - ▶ Learn **content-based** features (i.e., *domain knowledge*).
 - ▶ Exploit content-based features to directly learn **structure-based** features (i.e., *source knowledge*) from test data.
- ▶ Eliminate the need of a user involved in any source specific training process.

Introduction

▶ Contributions

- ▶ Based on our approach we produced a number of results.
 - ▶ ONDUX – On-Demand Unsupervised Learning for IE
 - **SIGMOD'10, IDAR'10, SBBD'11**
 - ▶ JUDIE – Joint Unsupervised Structure Discovery and IE
 - **SIGMOD'11**
 - ▶ iForm – A Probabilistic Approach for Automatically Filling Form-Based Web Interfaces
 - **WWW'09, PVLDB'10**

Related Work

- ▶ **Language for Wrapper Development.**

- ▶ Alternative to general purpose languages such as Perl and Java.
- ▶ Minerva, WEB-OQL

- ▶ **Wrapper Induction Methods**

- ▶ Machine Learning usage to semi-automatically induce wrappers.
- ▶ WEIN, StalKer

- ▶ **NLP-based Methods**

- ▶ Usage of Natural Language Processing techniques (semantic class, POS)
- ▶ WHISK, TEXTRUNNER

- ▶ **Ontology-based Methods**

- ▶ Usage of an ontology and conceptual description of the data of interest

- ▶ **HTML-aware Methods**

- ▶ Explore the HTML Structure (Tags) and their representation (DOM)
- ▶ RoadRunner, Webtables

Related Work

Methods	Disadvantages
Language for Wrapper Development.	Rely on the Regularity of the HTML format
Wrapper Induction Methods.	Rely on the Regularity of the HTML format
NLP-based Methods	Require Linguistic and Grammatical Elements
Ontology-based Methods	Require a huge human effort to manually create ontologies
HTML-aware Methods	Rely on the Regularity of the HTML format

These disadvantages precludes their usage in a large number of textual sources that are available on the Web.



Related Work

▶ Probabilistic Graph-Based Methods

- ▶ Deal with the limitations of the extraction methods that are based on the HTML structure.
- ▶ Based on probabilistic frameworks such as: Conditional Random Fields (CRF) and Hidden Markov Models (HMM)

▶ Supervised Methods

- ▶ Rely on human-created training sets to generate graphical models able to extract information
- ▶ Require training data from each source

▶ Unsupervised Methods

- ▶ Rely on pre-existing datasets for easing the training process of probabilistic methods.
 - Dictionaries, Knowledge Bases

Related Work

- ▶ Probabilistic Graph-Based Methods
 - ▶ Supervised Methods

Regent Square \$228,900 1028 Mifflin Ave.; 6 Bedrooms; 2 Bathrooms. 412-638-7273

Related Work

- ▶ Probabilistic Graph-Based Methods
 - ▶ Supervised Methods

Regent Square \$228,900 1028 Mifflin Ave.; 6 Bedrooms; 2 Bathrooms. 412-638-7273

1. <Neighborhood>**Regent Square**</Neighborhood>
2. <Price>**\$228,900**</Price>
3. <Number>**1028**</Number>
4. <Street>**Mifflin Ave.;**</Street>
5. <Bedroom>**6 Bedrooms**</ Bedroom>
6. <Bathroom>**2 Bathrooms**</Bathroom>
7. <Phone>**412-638-7273**</Phone>

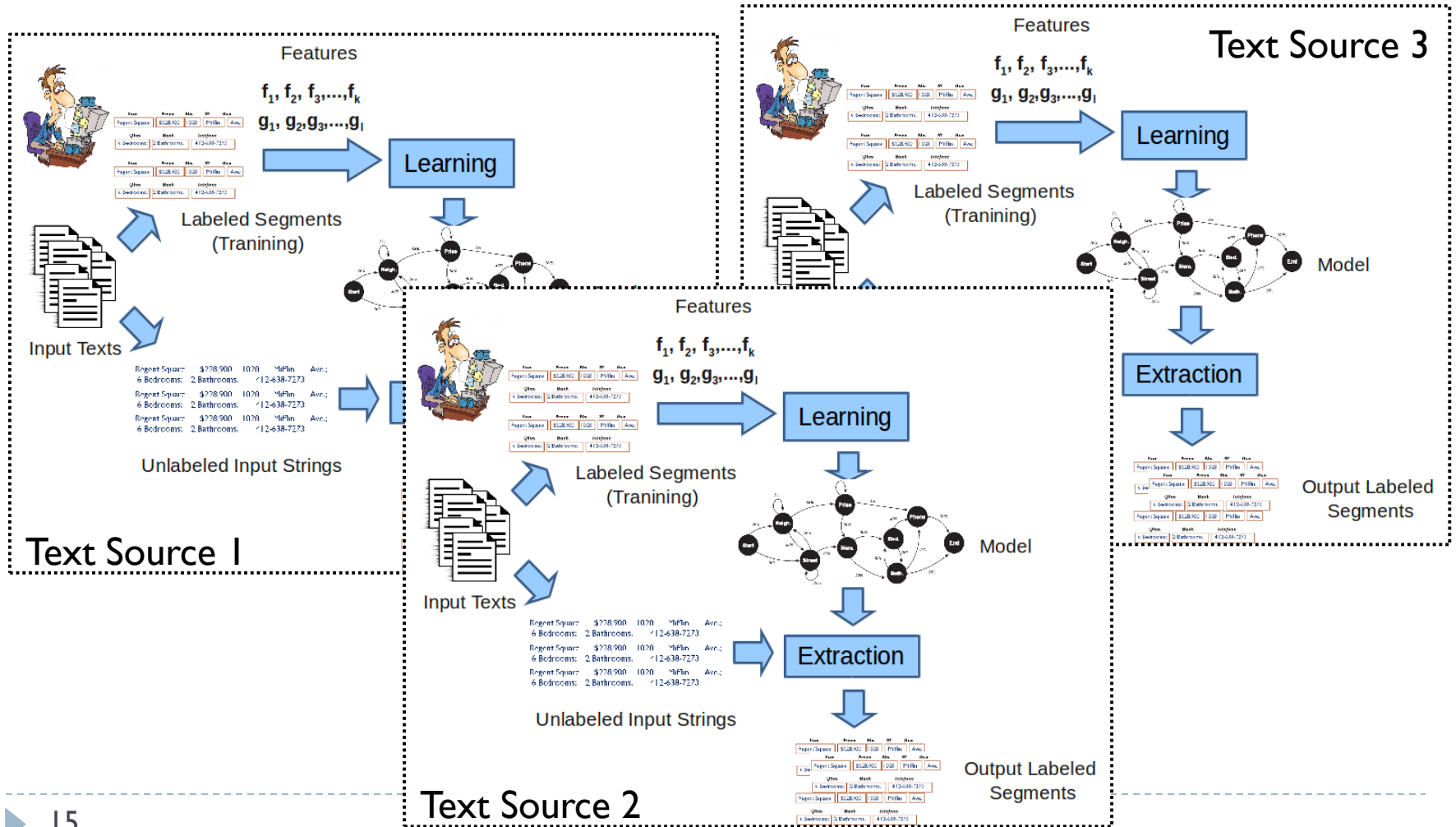
CRF and HMM
methods learn from
given examples, lexical,
style (content)
positioning and
sequencing (structure)
features

Examples are source-dependent



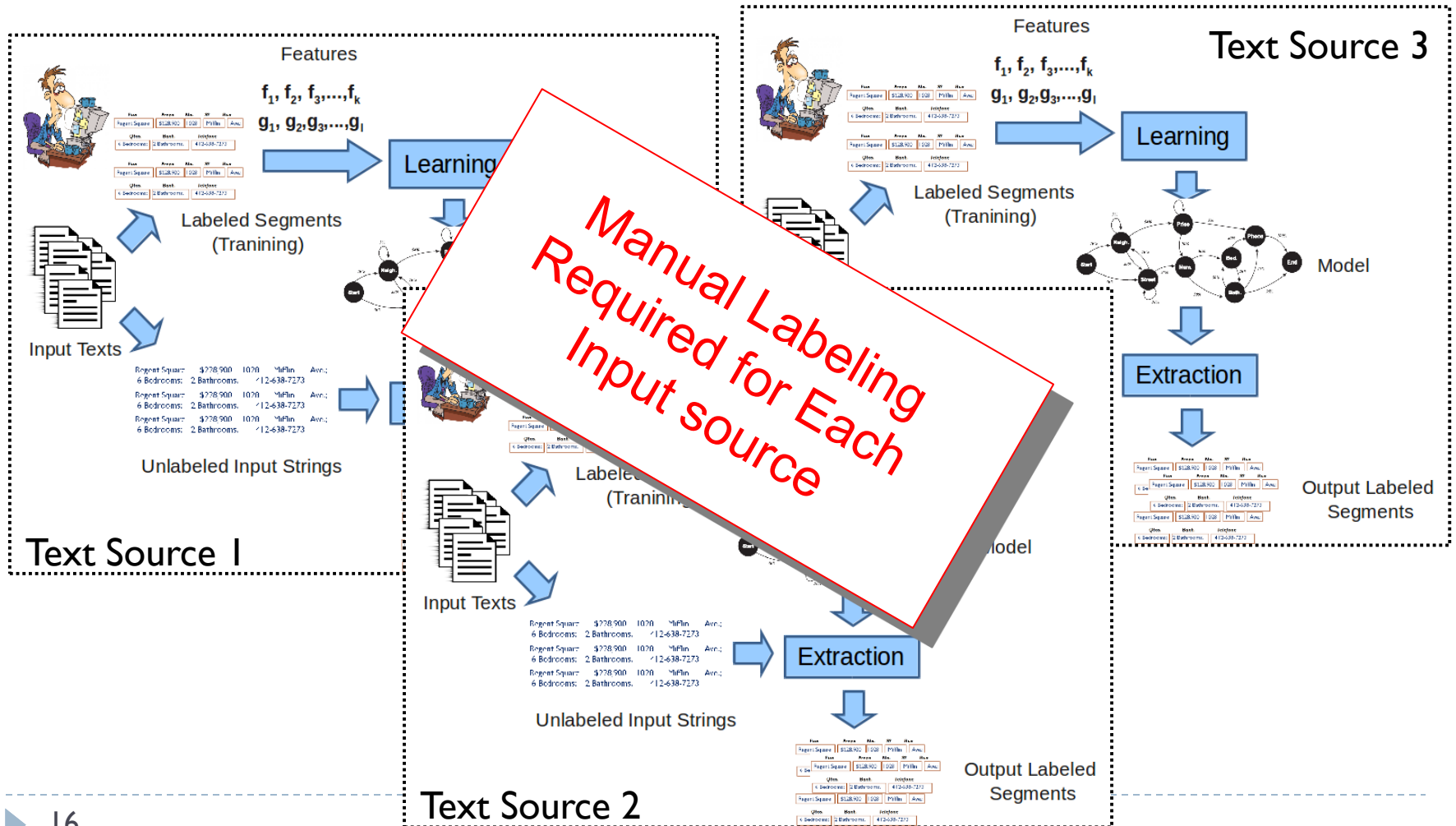
Related Work

Supervised Methods



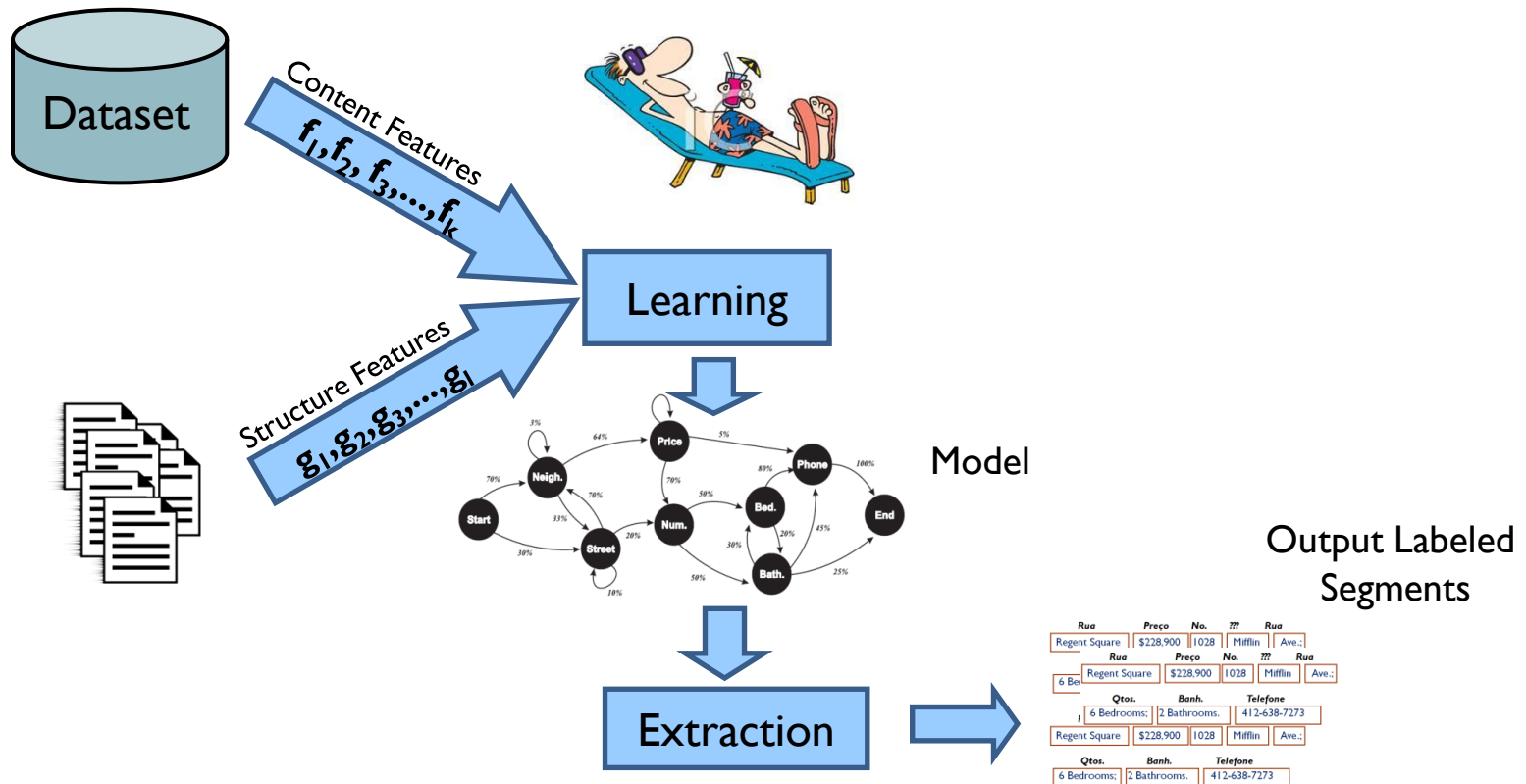
Related Work

Supervised Methods



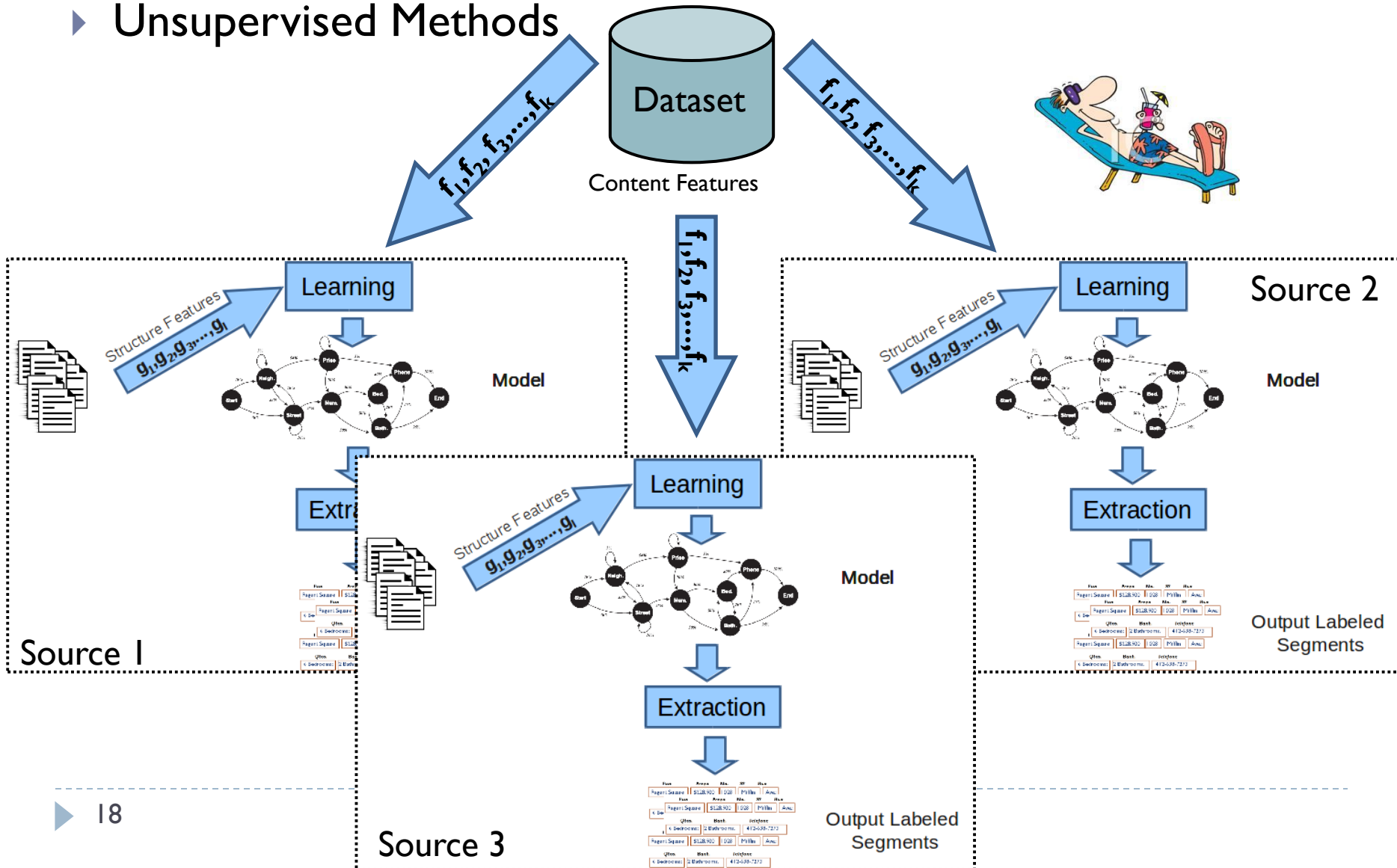
Related Work

► Unsupervised Methods



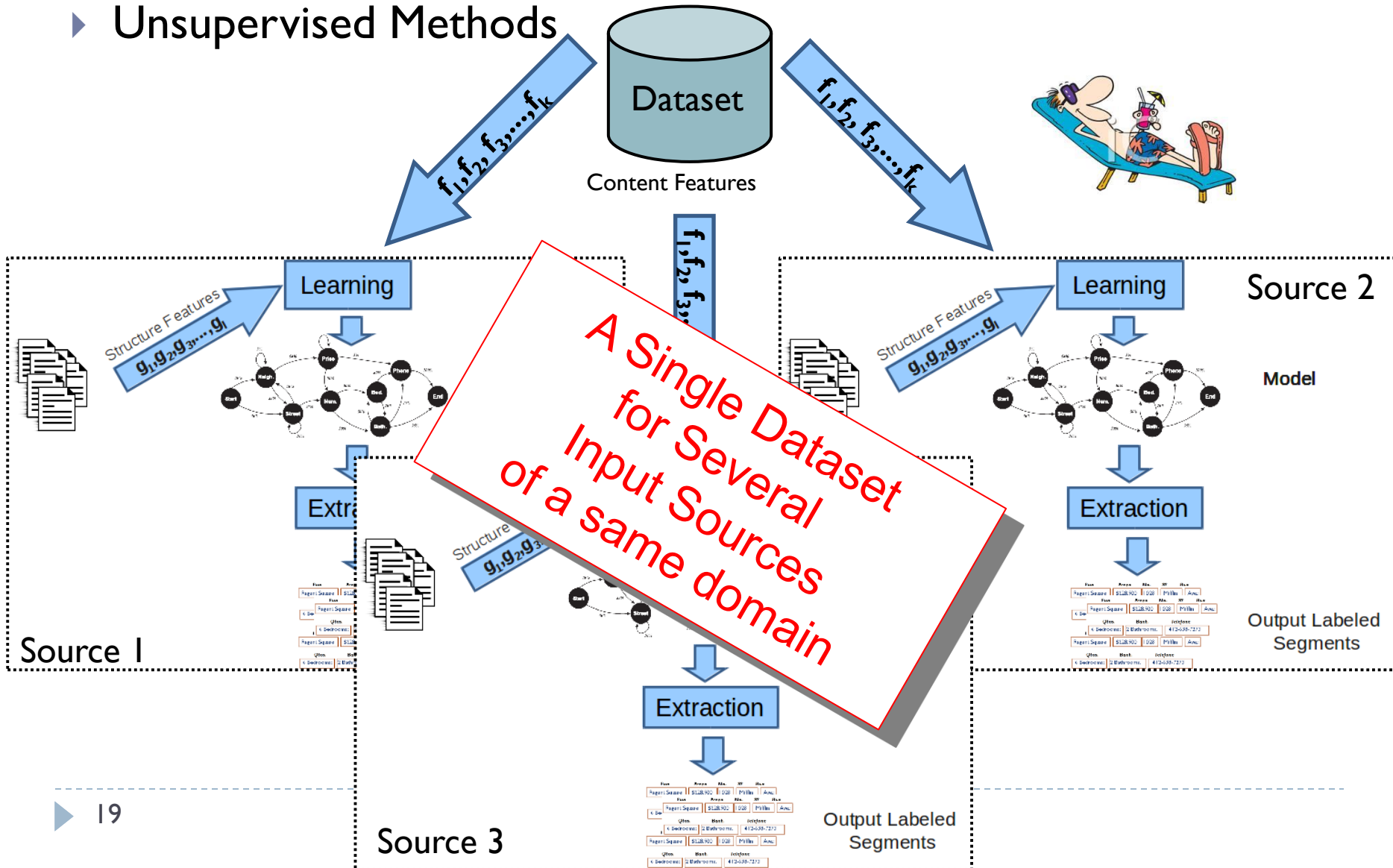
Related Work

► Unsupervised Methods



Related Work

► Unsupervised Methods



Related Work

▶ Probabilistic Graph-Based Methods

Supervised

X

UNsupervised

**Hand-labeled
examples**

Source Dependent

Scalability Problem

Reusability



**Pre-existing
information**

**Source
Independent**

Easily adaptable

Related Work

- ▶ Probabilistic Graph-Based Methods
 - ▶ Unsupervised Methods
 - ▶ [Agichtein et al @ SIGKDD 2004]
 - ▶ Usage of Reference Tables to create an unsupervised model using Hidden Markov Models (HMM)
 - ▶ [Zhao et al. @ SIAM ICDM 2008]
 - ▶ Usage of reference tables to create unsupervised CRF models - (U-CRF)

- ▶ [Sarawagi et al. @ ICDE 2006]
 - ▶ Usage of pre-existing data and hand labeled training sets to create an semi-supervised model using CRF

Related Work

▶ Probabilistic Graph-Based Methods

▶ Unsupervised Methods

▶ [Agichtein et al @ SIGKDD 2004]

- ▶ Usage of Reference Tables to create an unsupervised model using Hidden Markov Models (HMM)

▶ [Zhao et al. @ SIAM ICDM 2008]

- ▶ Usage of reference tables to create unsupervised CRF models - (U-CRF)

Both models assume single positioning and ordering of attributes in all test instances.

▶ [Sarawagi et al. @ ICDE 2006]

- ▶ Usage of pre-existing data and **hand labeled training sets** to create an semi-supervised model using CRF



Our Proposed Approach for IETS



Our Proposed Approach for IETS

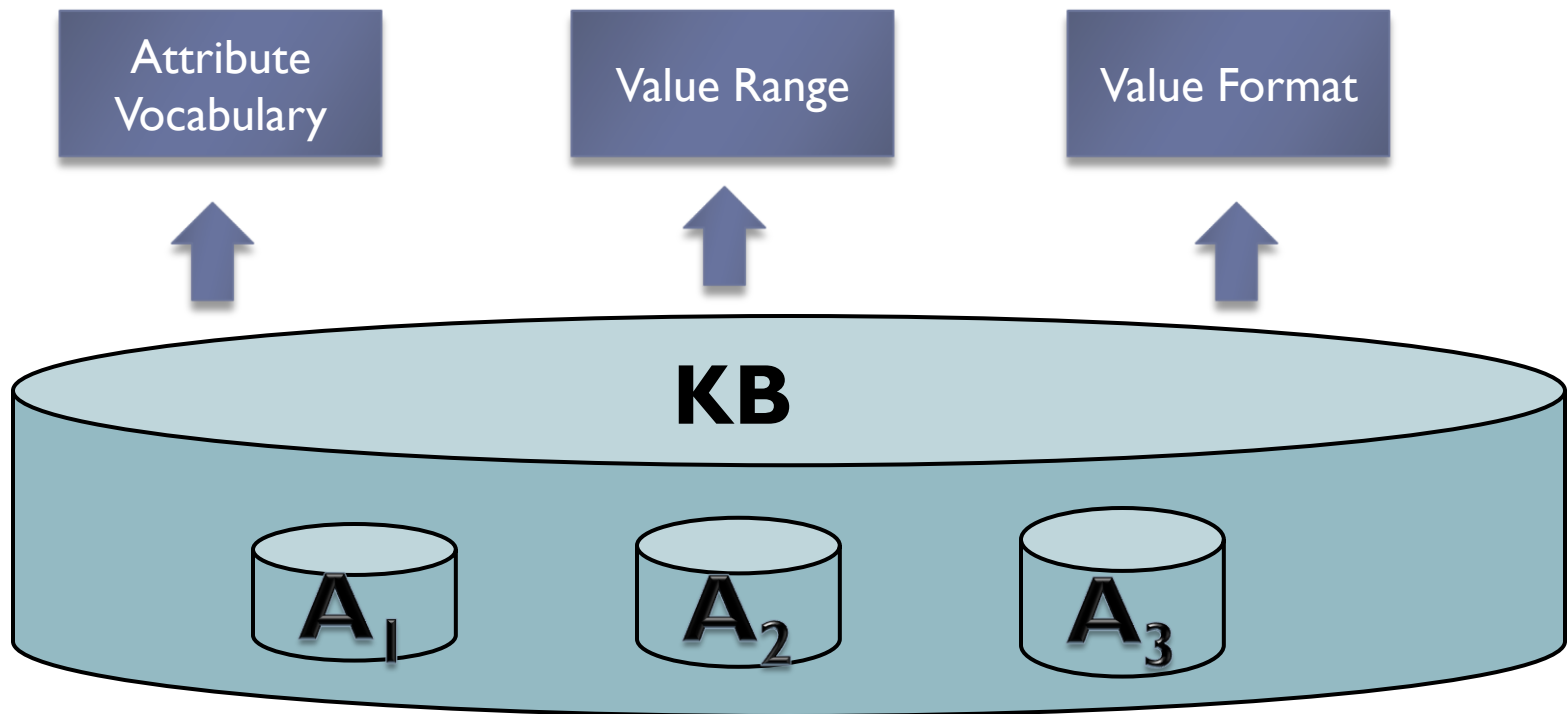
► Knowledge Bases

- Set of pairs $KB = \{(m_1, O_1), \dots, (m_n, O_n)\}$
- Easily built from pre-existing sources
 - Bibliographic DBs, Freebase, **Wikipedia**, etc.

$$\begin{aligned} K &= \{ \langle Author, O_{Author} \rangle, \langle Title, O_{Title} \rangle \} \\ O_{Author} &= \{ "J. K. Rowling", "Galadriel Waters", "Beatrix Potter" \} \\ O_{Title} &= \{ "Harry Potter and the Half-Blood Prince", \\ &\quad "A Guide to Harry Potter", "The Rabbit's Halloween" \} \end{aligned}$$

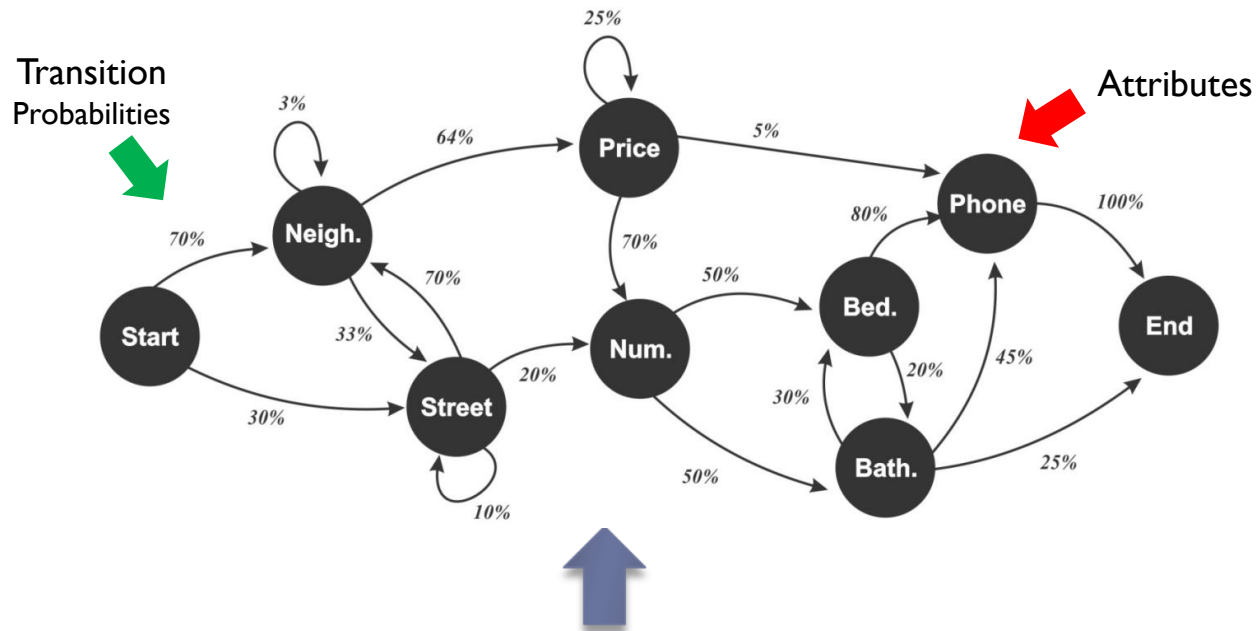
Our Proposed Approach for IETS

- ▶ Our approach relies on two types of features:
 - ▶ State or content-based features;



Our Proposed Approach for IETS

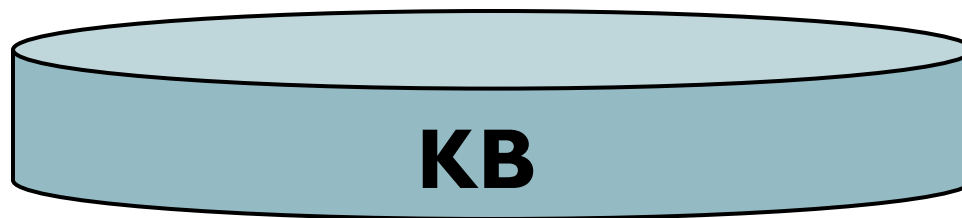
- ▶ Our approach relies on two types of features:
 - ▶ Transition or Structure-based Features;



Our Proposed Approach for IETS

- ▶ Knowledge bases implicitly encode domain knowledge.
 - ▶ Very suitable source for learning content-based features
- ▶ Attribute Vocabulary
 - ▶ Exploit the common vocabulary often shared by values of textual attributes

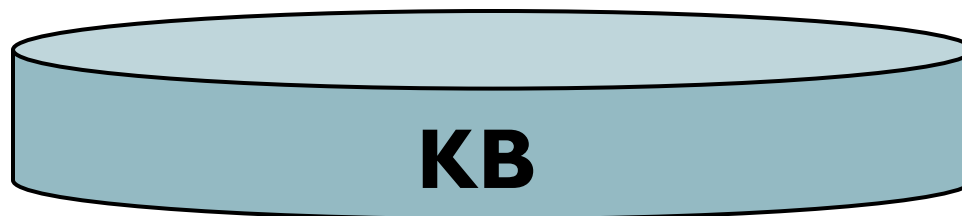
$$AF(s, A) = \frac{\sum_{t \in T(A) \cap T(s)} fitness(t, A)}{|T(s)|} \quad fitness(t, A) = \frac{f(t, A)}{N(t)} \times \frac{f(t, A)}{f_{max}(A)}$$



Our Proposed Approach for IETS

- ▶ Knowledge bases implicitly encode domain knowledge.
 - ▶ Very suitable source for learning content-based features
- ▶ Attribute Value Range
 - ▶ For the case of numeric candidate values, it measures the similarity between a numeric value and the set of values of a numeric attribute

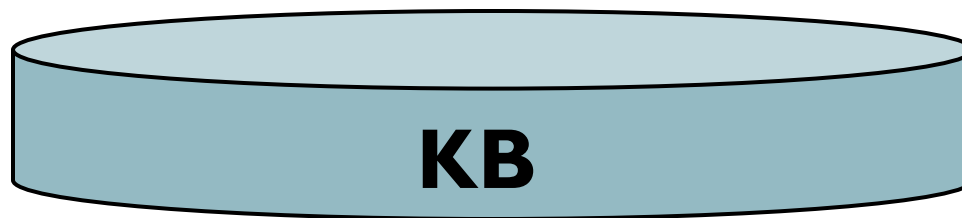
$$NM(s, A) = e^{-\frac{v_s - \mu}{2\sigma^2}}$$



Our Proposed Approach for IETS

- ▶ Knowledge bases implicitly encode domain knowledge.
 - ▶ Very suitable source for learning content-based features
- ▶ Attribute Value Format
 - ▶ Exploits the common format often used to represent values of some attributes.

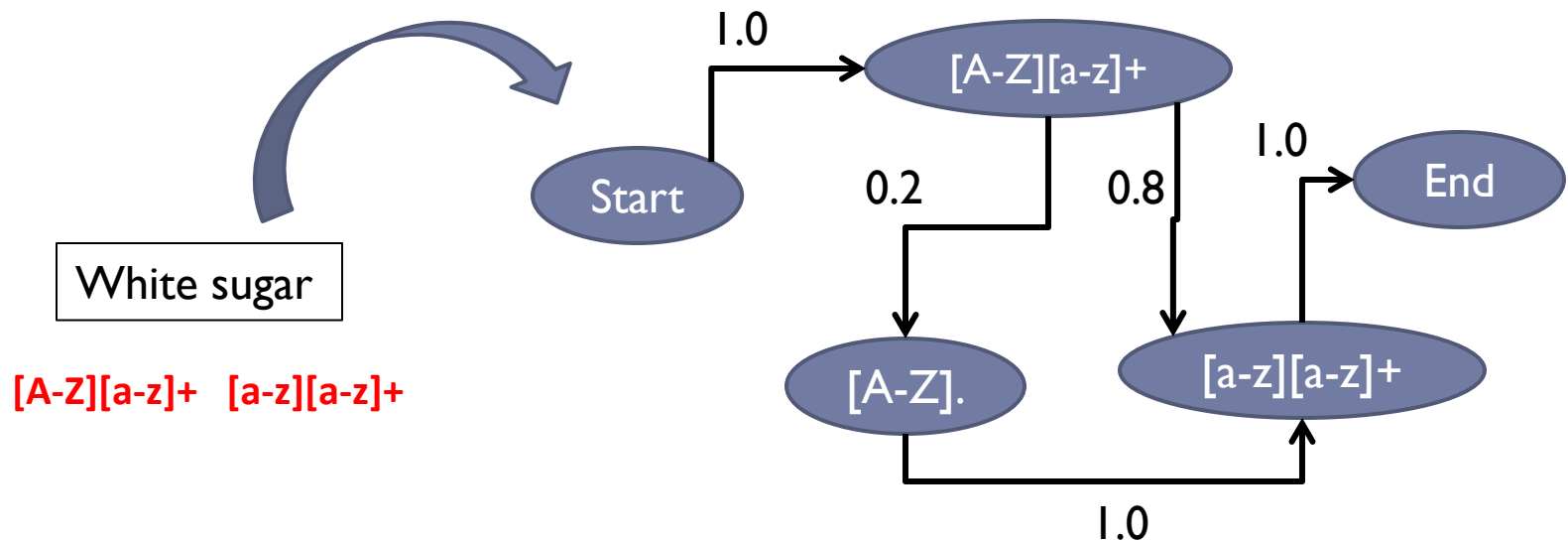
$$format(s, A) = \frac{\sum_{\langle n_x, n_y \rangle \in path(s)} w(n_x, n_y)}{|path(s)|}$$



Our Proposed Approach for IETS

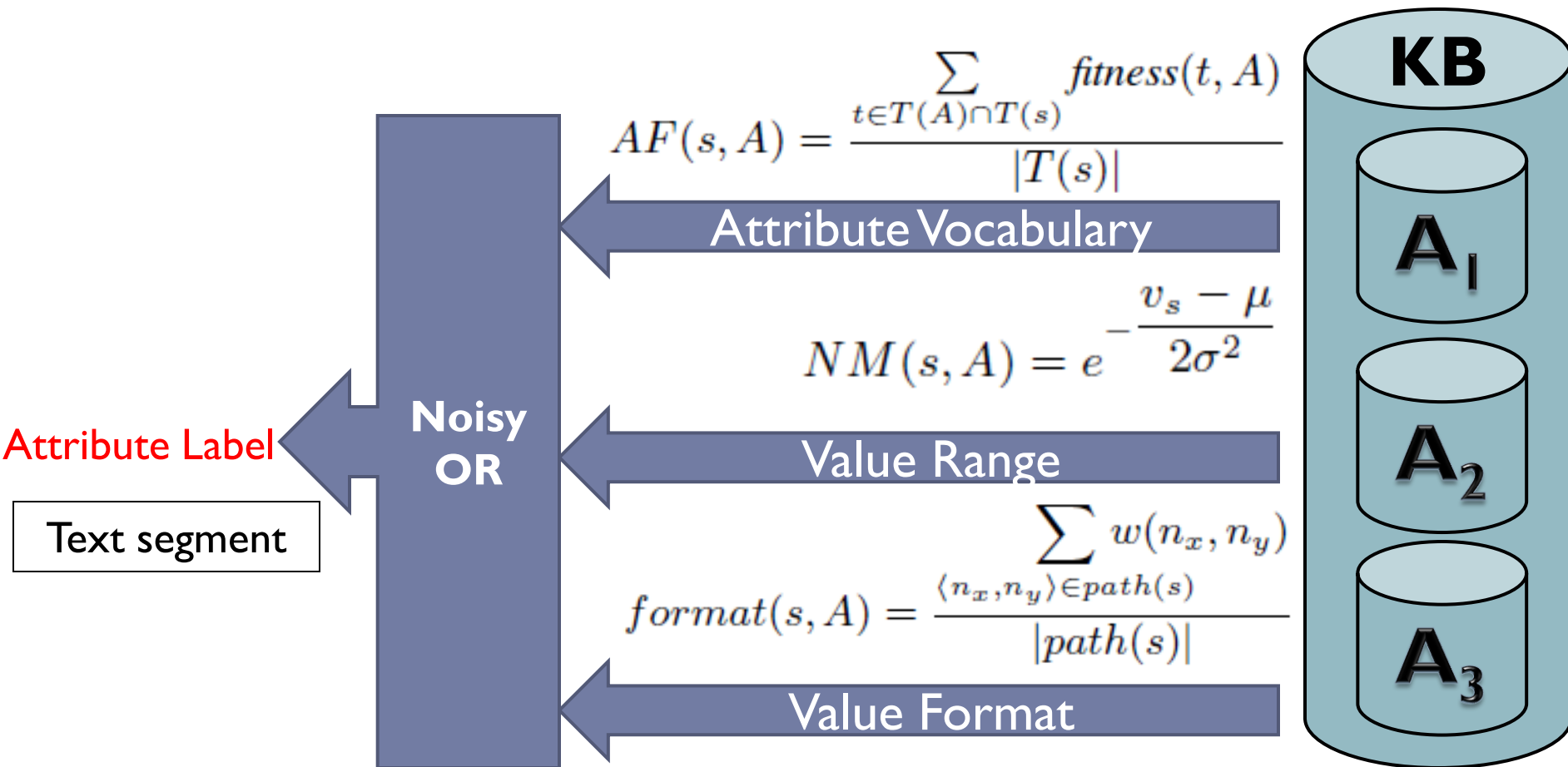
▶ Attribute Value Format (Style)

- ▶ First a Markov model is generated for each attribute.
- ▶ Computes the probability of the input mask sequence represents a path in each Markov model of each attribute.



Our Proposed Approach for IETS

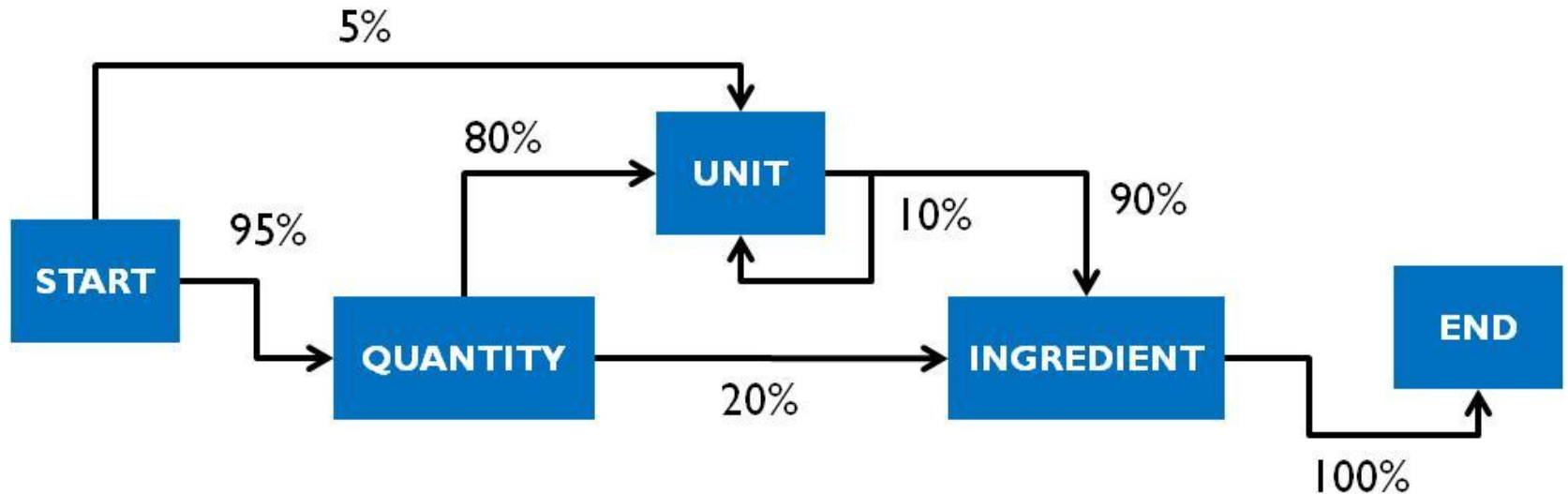
► Content-based Features



Our Proposed Approach for IETS

- ▶ Structure-based features are automatically induced from content-based Features
 - ▶ HMM-like graph called Positioning and Sequencing Model (PSM)
- ▶ Positioning and Sequencing Model
 - ▶ Automatically learned *On-Demand* from test instances
 - ▶ No *a priori* training required
- ▶ Structure-based features
 - ▶ Dependent of the placement of attributes values on the input
 - ▶ Thus, they are **input-dependent**

Our Proposed Approach for IETS

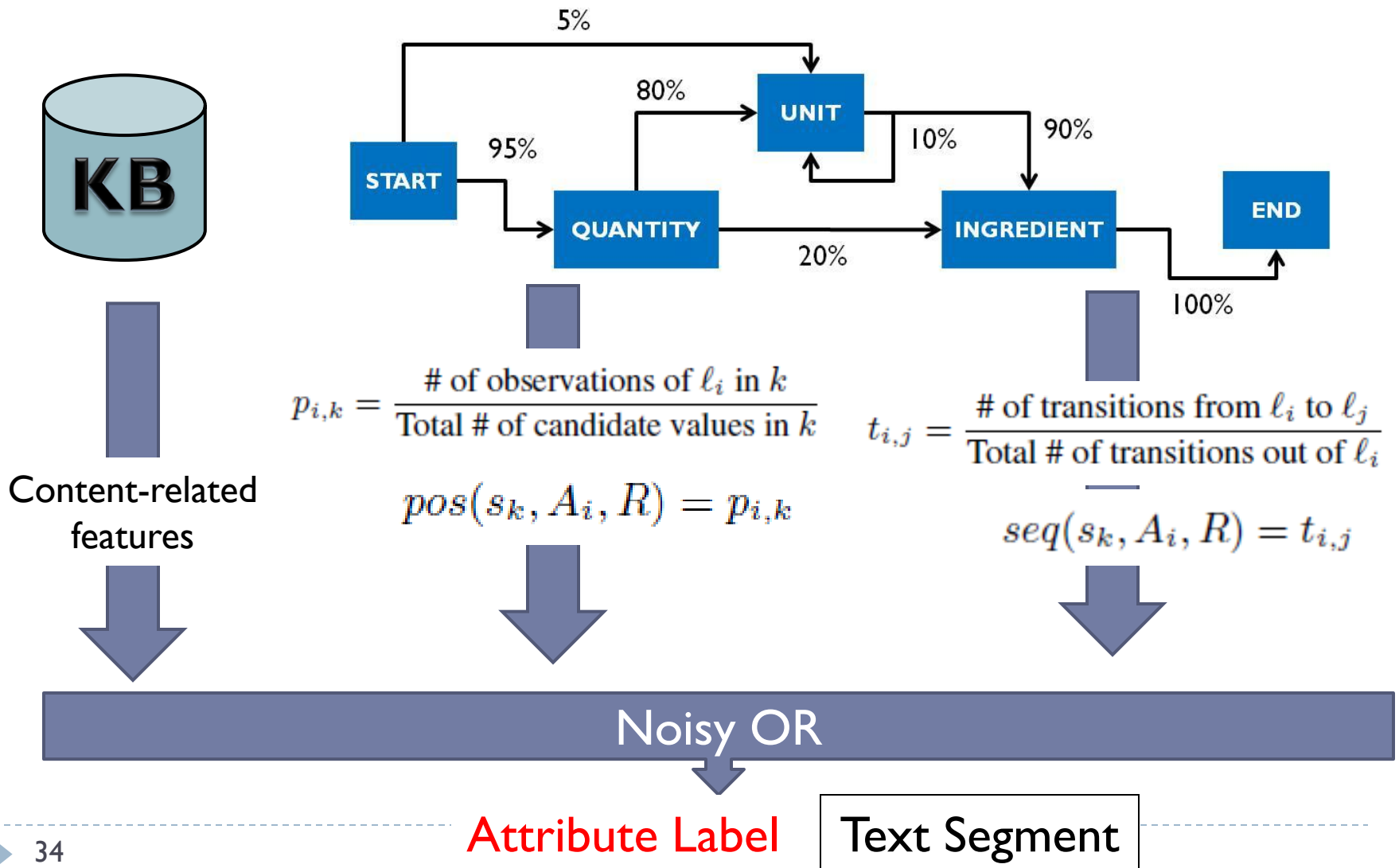


$$p_{i,k} = \frac{\text{\# of observations of } l_i \text{ in } k}{\text{Total \# of candidate values in } k} \quad t_{i,j} = \frac{\text{\# of transitions from } l_i \text{ to } l_j}{\text{Total \# of transitions out of } l_i}$$

$$pos(s_k, A_i, R) = p_{i,k}$$

$$seq(s_k, A_i, R) = t_{i,j}$$

Our Proposed Approach for IETS



Our Proposed Approach for IETS

- ▶ **Combination Strategy**

- ▶ *Bayesian Noise-OR-Gate*

$$or(p_1, \dots, p_n) = 1 - ((1 - p_1) \times \dots \times (1 - p_n))$$

- ▶ We assume that the features we use exploit different properties of the attributes of the KB, i.e., they are **independent**.
 - ▶ Probabilistic methods such as CRF and HMM deploy optimization process to combine their features.
 - ▶ Not using optimization can, in theory, lead to sub-optimal results, our experiments demonstrates that our combination works very well in practice.

Our Proposed Approach for IETS

- ▶ Based on our approach
 - ▶ We developed unsupervised information extraction by text segmentation methods
 - ▶ ONDUX
 - ▶ **On Demand Unsupervised Information Extraction**
 - ▶ JUDIE
 - ▶ **Joint Unsupervised Structure Discovery and Information Extraction**
 - ▶ iForm
 - ▶ **A Probabilistic Approach for Automatically Filling Form-Based Web Interfaces**

ONDUX

On-Demand Unsupervised Learning for Information Extraction

Cortez et al. - SIGMOD 2010, Cortez and Silva – IDAR 2010

ONDUX

- ▶ Deals with text documents containing implicit semi-structured data records
 - ▶ Addresses
 - ▶ Bibliographic References
 - ▶ Classified Ads
 - ▶ Product Descriptions

Postal Address

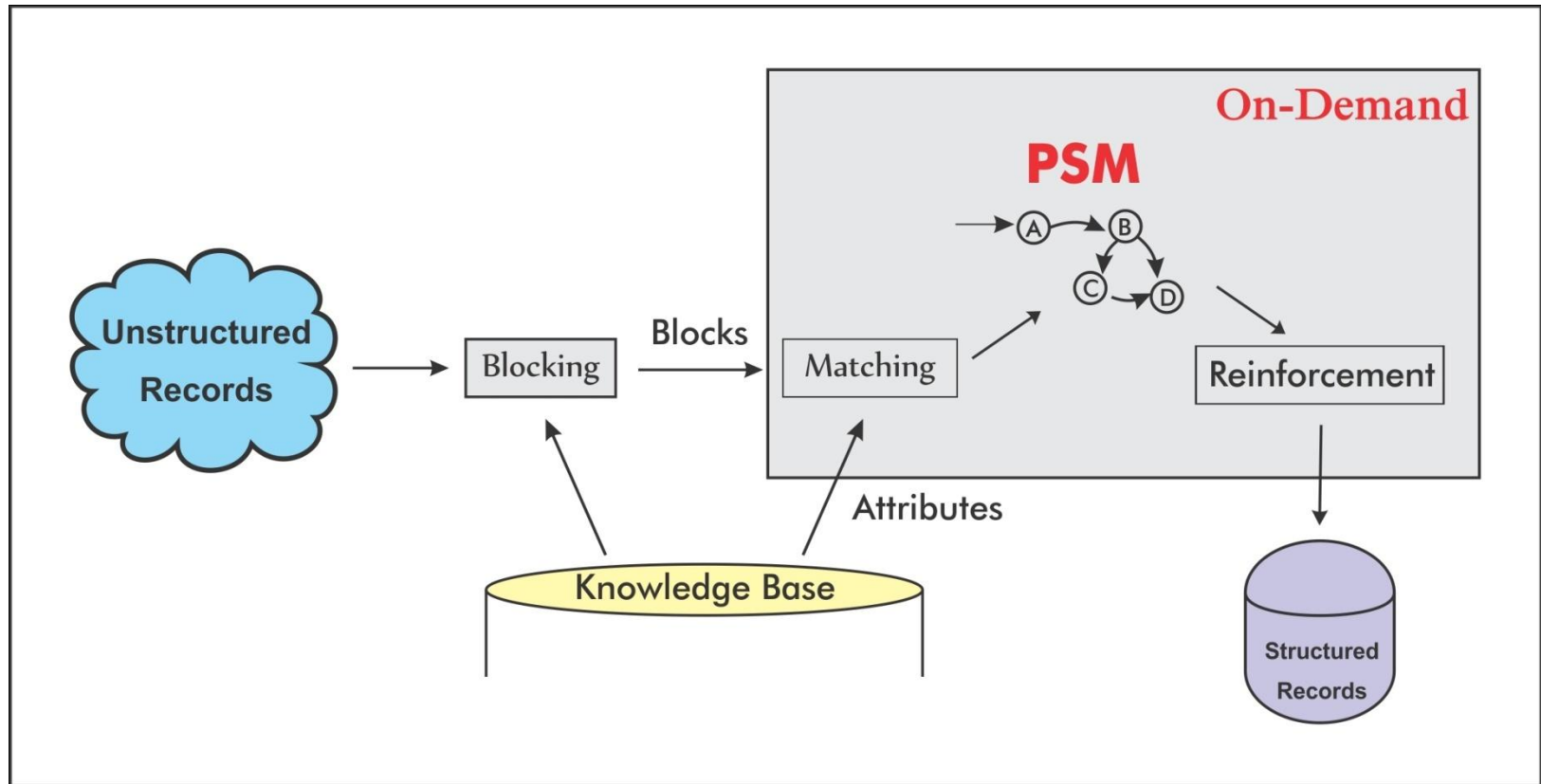
Dr. Robert A. Jacobson, 8109 Harford Road, Baltimore,
MD 21214

Bibliographic Reference

Pável Calado, Marco Cristo, Marcos André Gonçalves,
Edleno S. de Moura, Berthier Ribeiro-Neto, Nivio Ziviani.
Link-based similarity measures for the classification of Web
documents. JASIST, v. 57 n.2, p. 208-221, January 2006

ONDUX

► General View



ONDUX

▶ Blocking

- ▶ Split the input text in *substrings* called *blocks*;
- ▶ Consider the co-occurrence of consecutive terms based on the KB

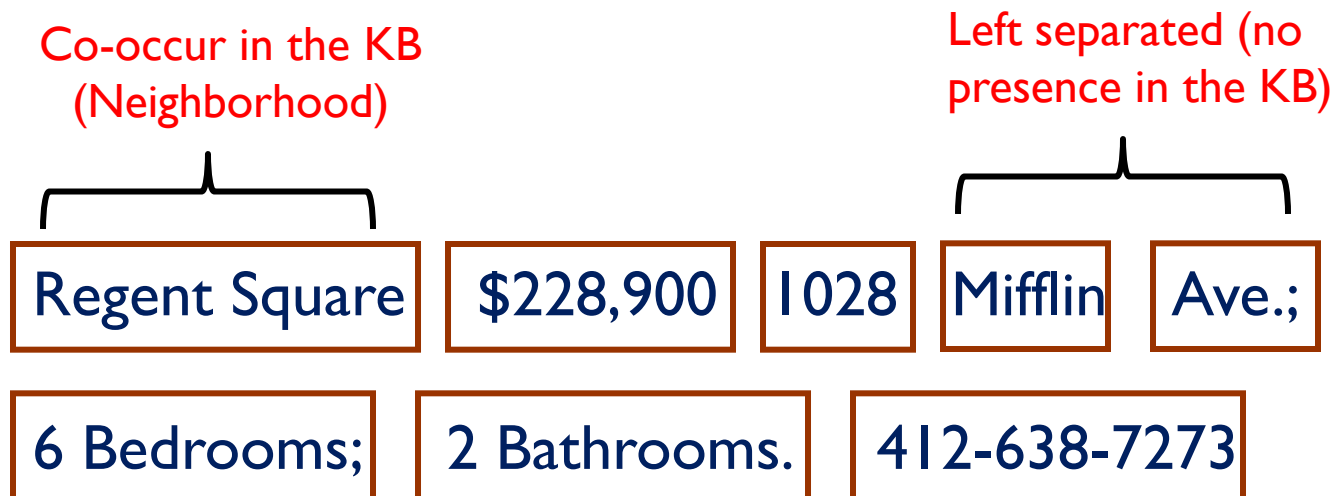
Regent Square \$228,900 1028 Mifflin Ave.;

6 Bedrooms; 2 Bathrooms. 412-638-7273

ONDUX

▶ Blocking

- ▶ Split the input text in *substrings* called *blocks*;
- ▶ Consider the co-occurrence of consecutive terms based on the KB



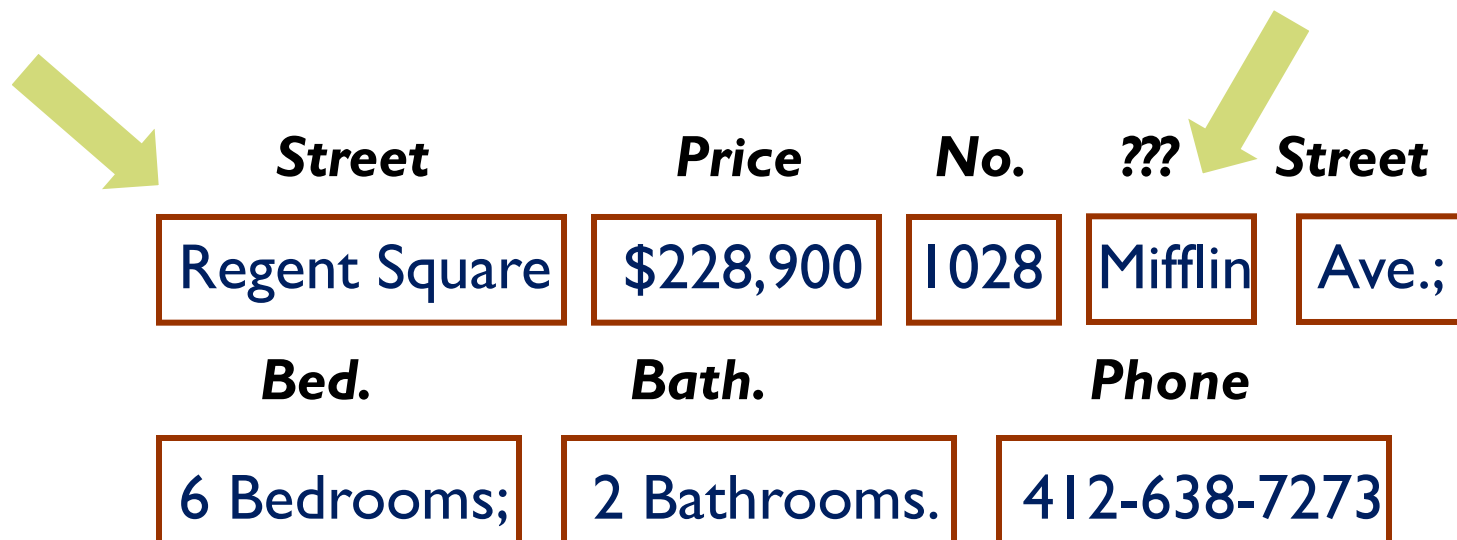
ONDUX

▶ Matching

- ▶ Associate each blocks with attributes according to content-based features.

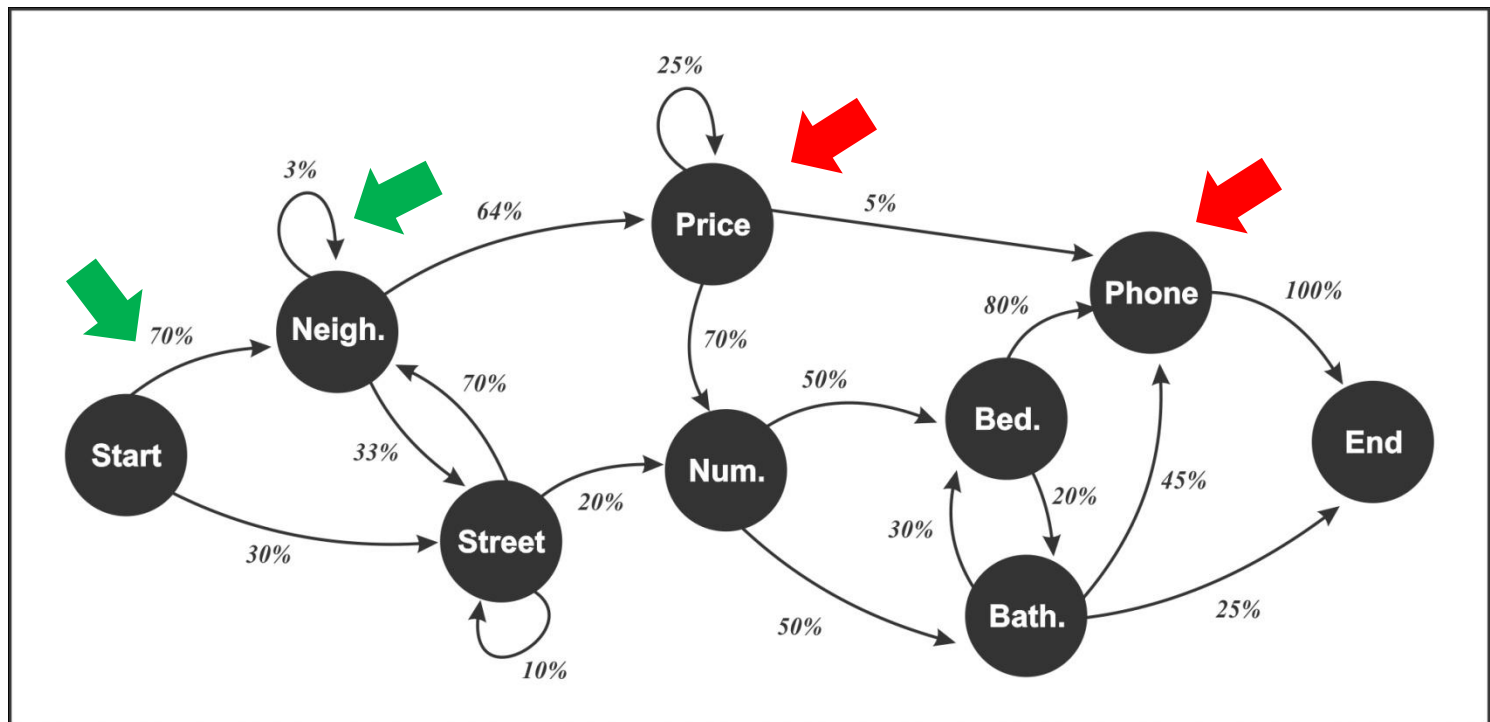
- ▶ Attribute Vocabulary $AF(s, A) = \frac{\sum_{t \in T(A) \cap T(s)} fitness(t, A)}{|T(s)|}$

- ▶ Value Range $NM(s, A) = e^{-\frac{v_s - \mu}{2\sigma^2}}$



ONDUX

► Reinforcement – PSM



Ordering and Positioning Features are learned *On-Demand* based on the test instances trough the Matching Phase

ONDUX

► Reinforcement

- Once the PSM is built, we combine the content-based and the structure-based features using the Bayesian operator *OR*.

Street	Price	No.	???	Street
Regent Square	\$228,900	1028	Mifflin	Ave.;
Bed.	Bath.	Phone		
6 Bedrooms;	2 Bathrooms.	412-638-7273		

ONDUX

► Reinforcement

- Once the PSM is built, we combine the content-based and the structure-based features using the Bayesian operator *OR*.

<i>Neighborhood</i>	<i>Price</i>	<i>No.</i>	<i>Street</i>
Regent Square	\$228,900	1028	Mifflin Ave.;
<i>Bed.</i>	<i>Bath.</i>	<i>Phone</i>	
6 Bedrooms;	2 Bathrooms.	412-638-7273	

ONDUX - Experiments

▶ Setup

- ▶ We tested our proposed method with several sources from 3 distinct domains:

- ▶ Addresses
- ▶ Bibliographic Data
- ▶ Classified Ads

▶ Metrics

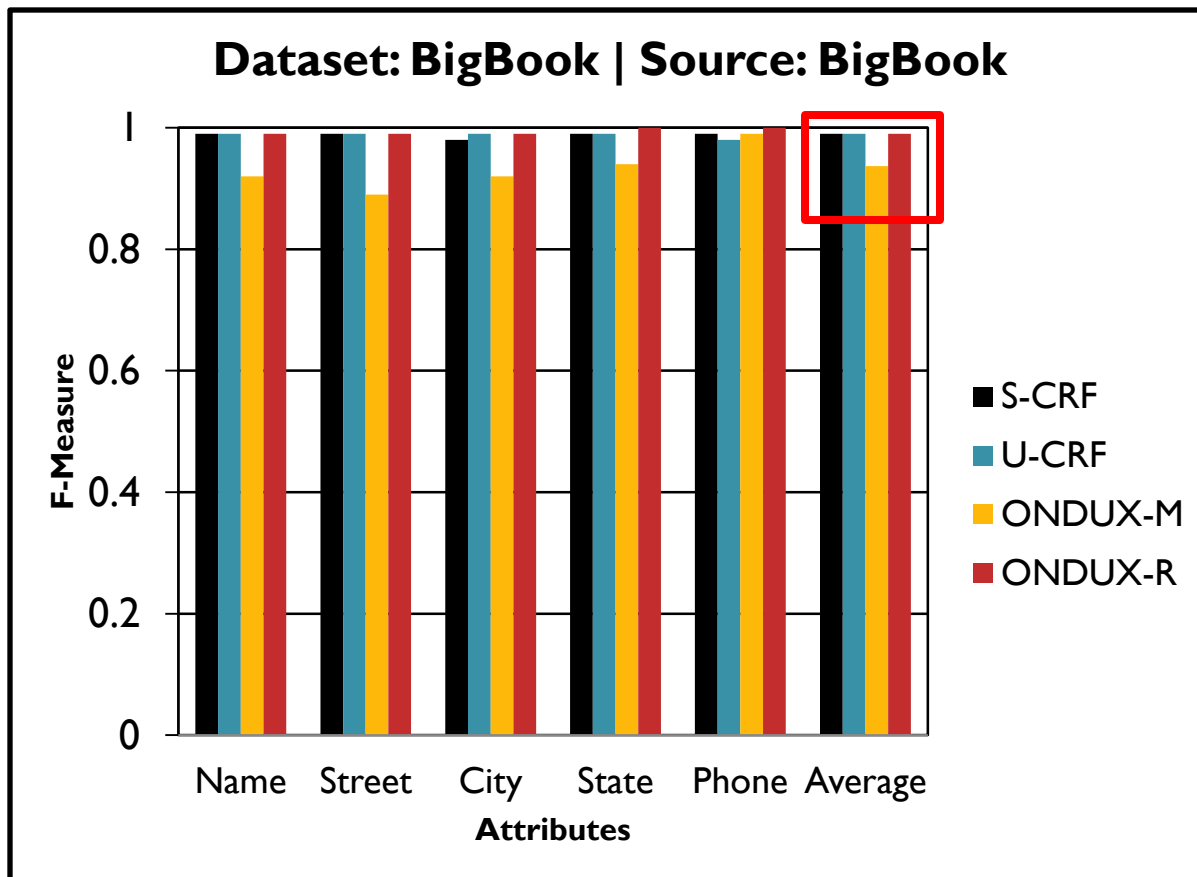
- ▶ Precision, Recall and F-Measure
 - T-Test for the statistical validation of the results

▶ Baselines

- ▶ U-CRF and S-CRF

ONDUX - Experiments

▶ Extraction Quality



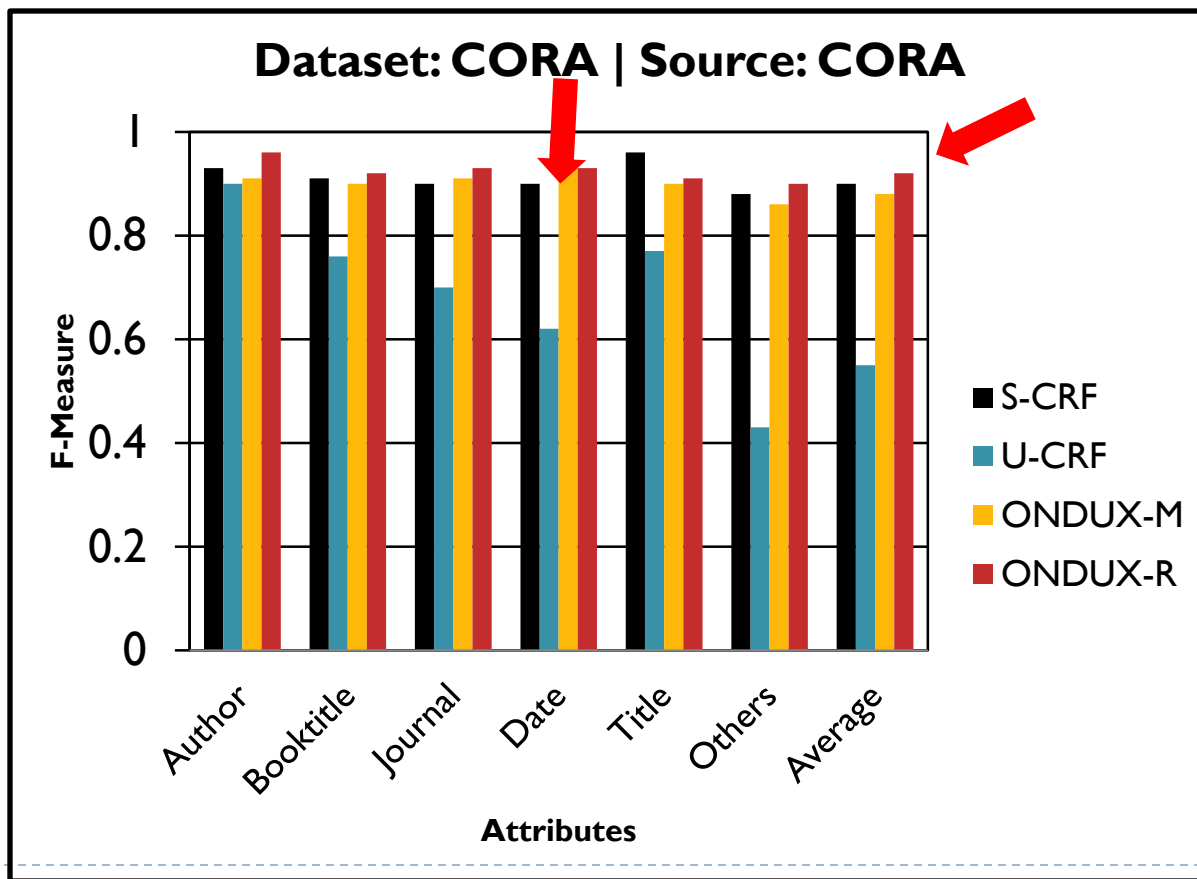
U-CRF results similar to Zhao@SICDM (validation)

Dataset follows the single order assumption

After Reinforcement ONDUX achieved similar quality

ONDUX - Experiments

▶ Extraction Quality



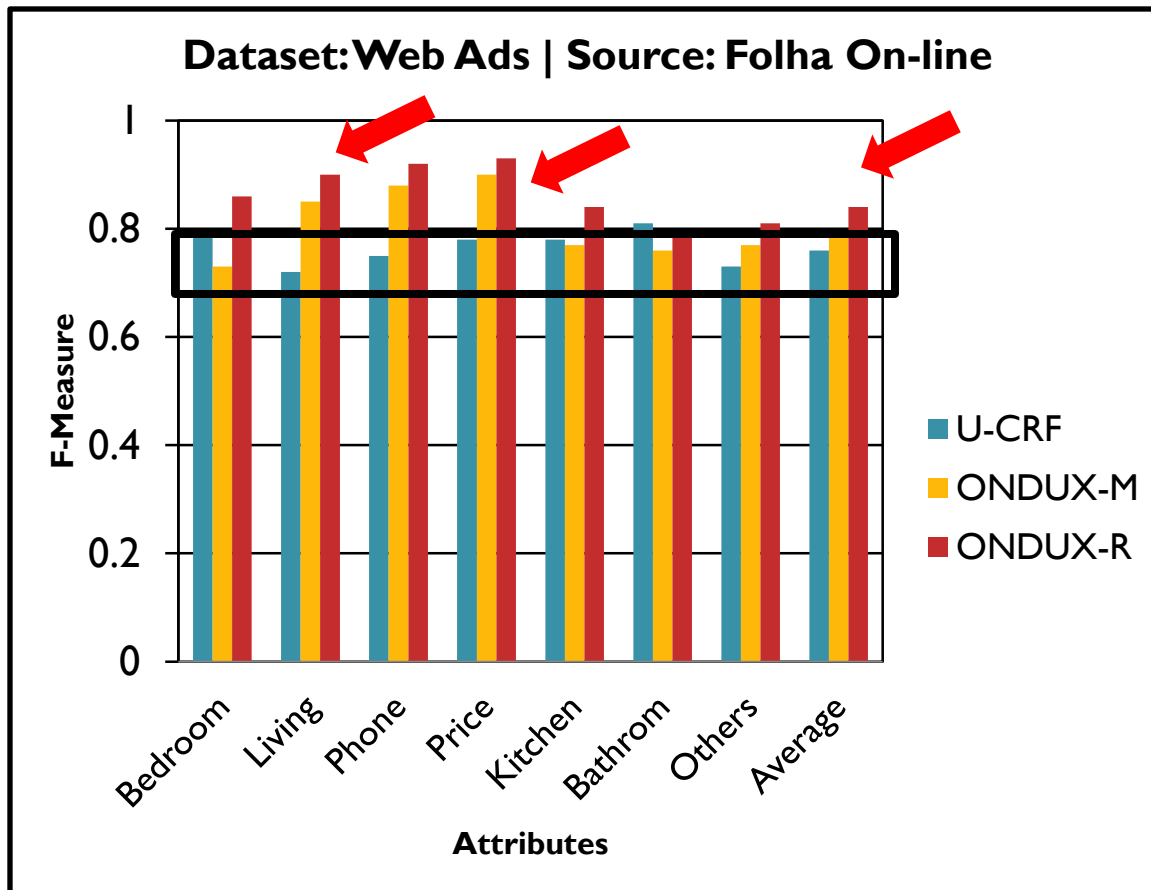
CORA includes a variety of citation styles (conference, journal, books, etc.)

S-CRF achieved results higher than U-CRF due to the hand-labeled training

In general, ONDUX outperformed CRF models

ONDUX - Experiments

► Extraction Quality



U-CRF presented a poor performance (very heterogeneous dataset)

Due to the Matching Phase and the PSM that is learned *On-Demand*, ONDUX achieve very high quality results

JUDIE

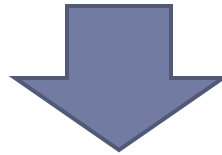
Joint Unsupervised Structure Discovery and Information Extraction

Cortez et al. - SIGMOD 2011

JUDIE

Chocolate Cake Recipe

1/2 cup butter 2 eggs 4 cups white sugar ground cinnamon 2 tablespoons dark rum 6 chopped pecans 1/2 cup milk 1 1/2 cups applesauce 2 cups all-purpose flour 1/4 cup cocoa powder 2 teaspoons baking soda 1/8 teaspoon salt 1 cup raisins 1/4 cup dark rum



Quantity	Unit	Ingredient
1/2	cup	butter
2		eggs
4	cups	white sugar
		ground cinnamon
2	tablespoons	dark rum
6		chopped pecans

JUDIE

- ▶ **Joint Unsupervised Structure Discovery and Information Extraction**
 - ▶ Detects the structure of each individual record being extracted without any user intervention
 - ▶ Looks for frequent patterns of label repetitions or **cycles**
- ▶ **Integrates this algorithm in the IE process**
 - ▶ Accomplished by successive refinement steps that alternate information extraction and structure discovery.

JUDIE

1/2 cup raising flour 2 level Tbsp Cocoa pinch Salt 1/4 cup Melted butter 1 Egg a little Vanilla

Q U I Q U ? I U I Q U I Q I I I

1/2	cup	raising	flour	2	level	Tbsp	Cocoa	pinch	Salt	1/4	cup	Melted	butter	1	Egg	a little	Vanilla
-----	-----	---------	-------	---	-------	------	-------	-------	------	-----	-----	--------	--------	---	-----	----------	---------

Q	U	I	Q	U	?	I	U	I	Q	U	I	Q	I	I	I		
1/2	cup	raising	flour	2	level	Tbsp	Cocoa	pinch	Salt	1/4	cup	Melted	butter	1	Egg	a little	Vanilla

Phase 1
Structure-free Labeling
Structure Sketching

Q	U	I	Q	U	U	I	U	I	Q	U	I	Q	I	U	I		
1/2	cup	raising	flour	2	level	Tbsp	Cocoa	pinch	Salt	1/4	cup	Melted	butter	1	Egg	a little	Vanilla

Phase 2
Structure-aware Labeling
Structure Refinement

Q	U	I	Q	U	U	I	U	I	Q	U	I	Q	I	U	I		
1/2	cup	raising	flour	2	level	Tbsp	Cocoa	pinch	Salt	1/4	cup	Melted	butter	1	Egg	a little	Vanilla

JUDIE – Structure-free Labeling

▶ What is the best label for each segment?

- ▶ No information on the structure of the data records
- ▶ Resort only to content-based features

$$AF(s, A) = \frac{\sum_{t \in T(A) \cap T(s)} \text{fitness}(t, A)}{|T(s)|}$$

Attribute Vocabulary

$$NM(s, A) = e^{-\frac{v_s - \mu}{2\sigma^2}}$$

Value Range

$$\text{format}(s, A) = \frac{\sum_{\langle n_x, n_y \rangle \in \text{path}(s)} w(n_x, n_y)}{|\text{path}(s)|}$$

Value Format

Noisy
OR

Ingredient

White sugar

KB

A₁

A₂

A₃

JUDIE – Structure-free Labeling

- ▶ Initially labels potential values with attribute names.
 - ▶ No information on the structure of the data records
 - ▶ Resort only to content-based features
 - ▶ Learned from the pre-existing KB

1/2 cup raising flour 2 level Tbsp Cocoa pinch Salt 1/4 cup Melted butter 1 Egg a little Vanilla

Q	U	I	Q	U	?	I	U	I	Q	U	I	Q	I	I	I
1/2	cup	raising flour	2	level	Tbsp	Cocoa	pinch	Salt	1/4	cup	Melted butter	1	Egg	a little	Vanilla

Limitations:

Unmatching: “Tbsp”

Mismatching: “a little”

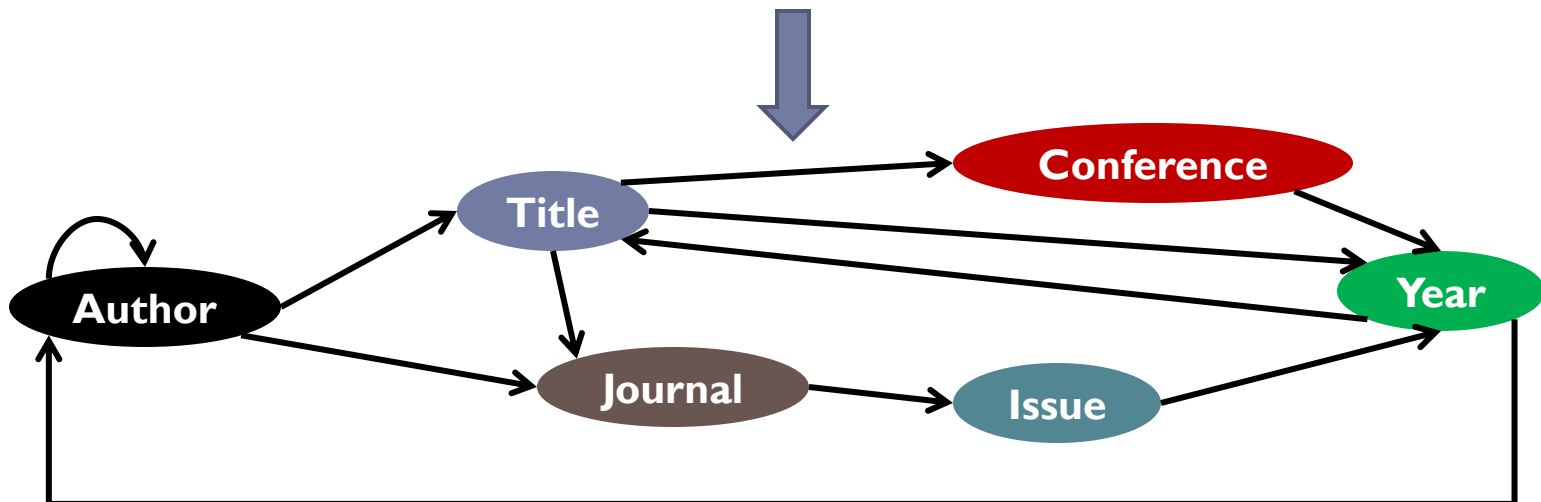
The SD Algorithm (I)

- ▶ Uncover the structure of implicit records from the input text.
 - ▶ Used in the Structure Sketching and Structure Refinement.
- ▶ Takes as input a sequence of labels and generates the structure of each record.
 - ▶ Assumption: It is possible to identify patterns of sequences by looking for cycles into a graph (Adjacency Graph) that models the ordering of labels.

The SD Algorithm (II)

- ▶ Consider a sequence of labels from a bibliographic reference input text.

Title Conference Year Author Author Title Conference Year Author Title
Conference Year ... Author Title Journal Issue Year Author Title Journal
Issue Year Author Author Journal Issue Year Title Year ... Author Title
Conference Year Author Author Author Title Journal Issue Year



The SD Algorithm (V)

Dominant Cycles

- ▶ Given the set of Coincident cycles that are also viable, the Dominant Cycle are most frequent in the input

- ▶ Finally, the algorithm works by first identifying all dominant cycles in the adjacency graph and then processing each of these cycles, the largest cycles being processed first.

- ▶ In our given examples, the dominant cycles are:
 1. [Author, Title, Journal, Issue, Year]
 2. [Author, Title, Conference, Year]
 3. [Author, Journal, Issue, Year]
 4. [Title, Conference, Year]
 5. [Title, Year]

JUDIE – Structure Sketching

- ▶ Organizes the labeled candidate values into records
 - ▶ Induces a structure on the unstructured text input.
 - ▶ Outputs labeled values grouped into records
 - ▶ Uses a novel algorithm called **Structure Discovery (SD)**

Q	U	I	Q	U	?	I	U	I	Q	U	I	Q	I	I	I
1/2	cup	raising flour	2	level	Tbsp	Cocoa	pinch	Salt	1/4	cup	Melted butter	1	Egg	a little	Vanilla

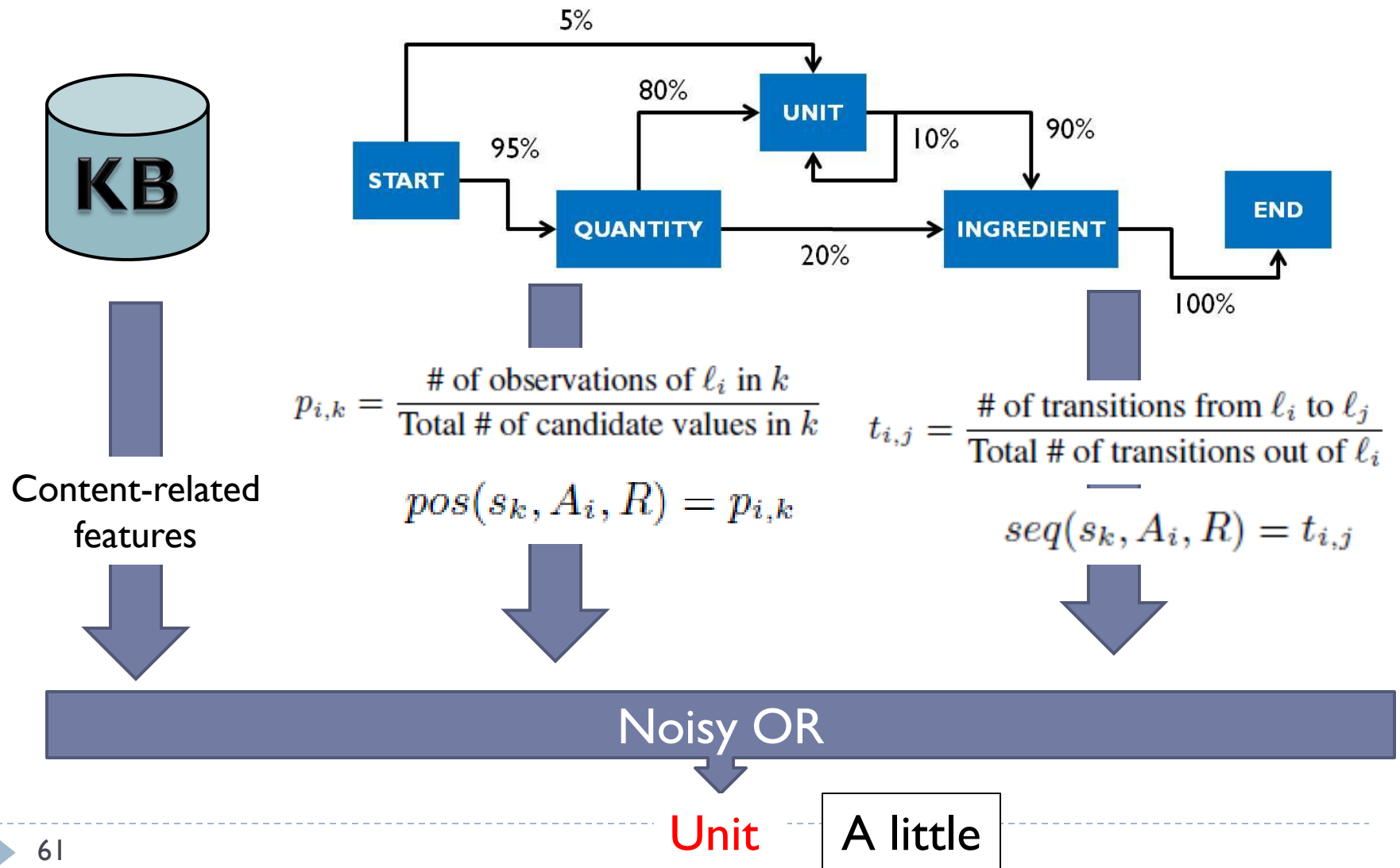
Q	U	I	Q	U	?	I	U	I	Q	U	I	Q	I	I	I
1/2	cup	raising flour	2	level	Tbsp	Cocoa	pinch	Salt	1/4	cup	Melted butter	1	Egg	a little	Vanilla

JUDIE – Structure-aware Labeling

- ▶ *Now, what is the best label for each segment?*
- ▶ We already know some structural information
- ▶ Re-labels segments considering **content-based features** and **structure-based features**
- ▶ Structure-based features learned using a graphical model (PSM)

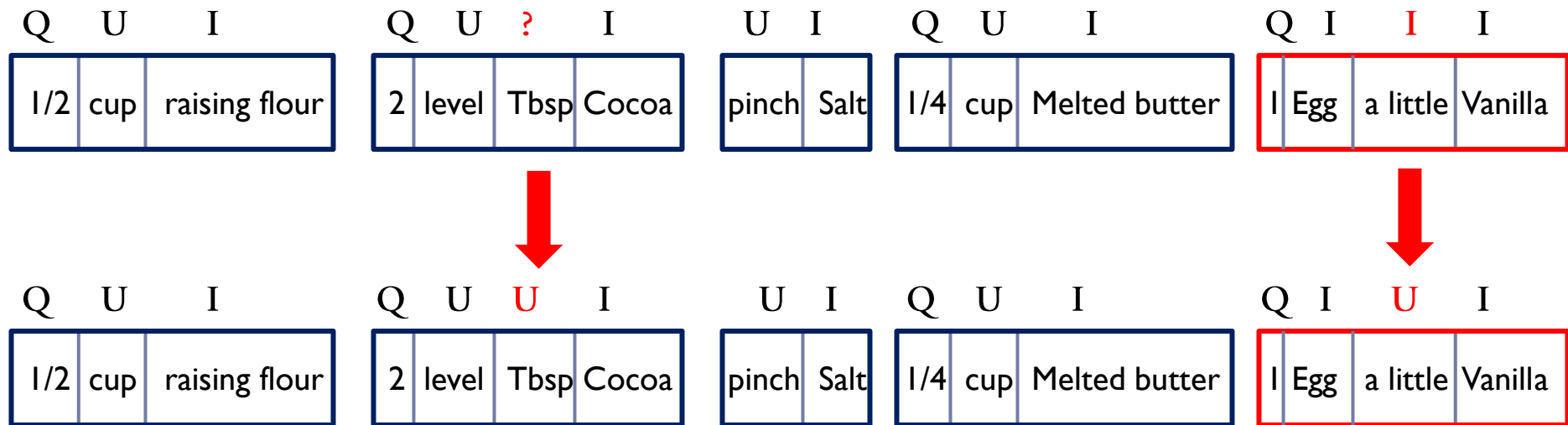
Q	U	I	Q	U	?	I	U	I	Q	U	I	Q	I	I	I
1/2	cup	raising flour	2	level	Tbsp	Cocoa	pinch	Salt	1/4	cup	Melted butter	1	Egg	a little	Vanilla

JUDIE – Structure-aware Labeling



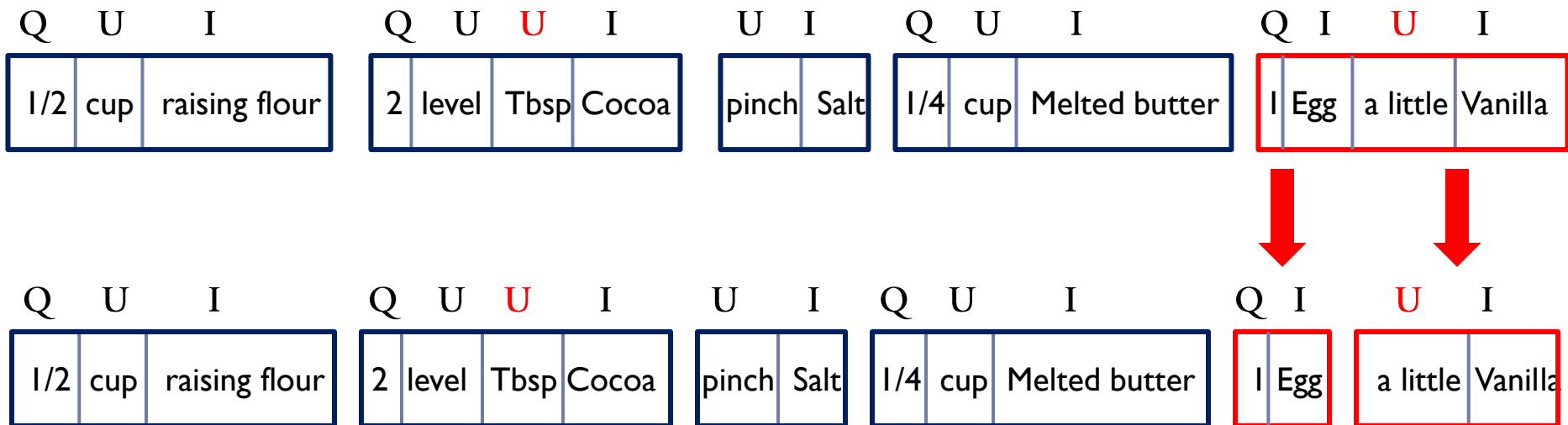
JUDIE – Structure-aware Labeling

- ▶ Labels textual values considering:
 - ▶ Uses a graphic model representing the likelihood of attribute transitions within the input text
 - ▶ Content-related features and structure-based features



JUDIE – Structure Refinement

- ▶ Applies again the SD algorithm
 - ▶ Considers the output of the structure-aware labeling
 - ▶ Fixes structural problems
 - Structure-aware labeling produces more precise results



Experiments

Domain	Dataset	Text Inputs	Attributes	Source	Attributes	Records
<i>Cooking Recipes</i>	<i>Recipes</i>	500	3	<i>FreeBase.com</i>	3	100
<i>Product Offers</i>	<i>Products</i>	10000	3	<i>Nhemu.com</i>	3	5000
<i>Postal Addresses</i>	<i>BigBook</i>	2000	5	<i>BigBook</i>	5	2000
<i>Bibliography</i>	<i>CORA</i>	500	3 to 7	<i>PersonalBib</i>	7	395
<i>Classified Ads</i>	<i>WebAds</i>	500	5 to 18	<i>Folha On-line</i>	18	125

▶ Metrics

- ▶ Precision, Recall and F-Measure
 - ▶ T-Test for the statistical validation of the results

▶ Baselines

- ▶ ONDUX and U-CRF

Evaluation – Record Level

Dataset	Phase 1	Phase 2	Gain (%)
Recipes	0.79	0.90	13.2
CORA	0.69	0.83	19.3
Web Ads	0.70	0.77	9.7

- ▶ Phase 1: acceptable. $F \approx 0.7$
- ▶ Phase 2: positive impact. Gains $> 9\%$
- ▶ In CORA, gains higher than 19%
 - ▶ Structural information led to significant improvements.

Comparison with baselines – Attribute Level

Attribute	JUDIE	ONDUX	U-CRF
Author	0.88	0.922	0.87
Title	0.70	0.79	0.69
Booktitle	0.86	0.89	0.56
Journal	0.84	0.90	0.55
Volume	0.90	0.96	0.43
Pages	0.86	0.84	0.50
Date	0.87	0.89	0.49
Average	0.86	0.88	0.58

CORA

Attribute	JUDIE	ONDUX	U-CRF
Bedroom	0.82	0.86	0.79
Living	0.89	0.90	0.72
Phone	0.87	0.92	0.75
Price	0.92	0.93	0.78
Kitchen	0.83	0.84	0.78
Bathroom	0.77	0.79	0.81
Others	0.73	0.79	0.71
Average	0.84	0.85	0.76

Web Ads

- ▶ Results very close to ONDUX and even better than U-CRF
- ▶ Recall: JUDIE faces a harder task.

iForm

A Probabilistic Approach for Automatically Filling
Form-Based Web Interfaces

Toda et al. – WWW 2009, Toda et Al. – PVLDB 2010

The Form Filling Problem

▶ **Goal:**

- ▶ To automatically fill out the fields of a given **form-based** interface with **values extracted** from a **data-rich free text document**.
 1. Extracting values from the input text;
 2. Filling out the fields of the target form using them.

Example

► Form-based interface

Vehicle Info

Type

Year

Make

Model

VIN

Mileage

Transmission

Engine

Drivetrain

Body style

Color

Int color

Int material Cloth Leather

Seating

Wheels

Tires

Roof

Truck bed

Stereo

Dealer code

Stock code

MSRP

NADA

KBB

Warranty

Text Box

Selection List

Features

Check-box

- | | | |
|--|--|---|
| <input type="checkbox"/> Power Steering | <input type="checkbox"/> Air Cond. (Rear) | <input type="checkbox"/> Roof Rack |
| <input type="checkbox"/> Power Brakes | <input type="checkbox"/> Cruise Control | <input type="checkbox"/> Fog Lamps |
| <input type="checkbox"/> Power Windows | <input type="checkbox"/> Air Bags (Driver) | <input type="checkbox"/> Sliding Rear Win |
| <input type="checkbox"/> Power Locks | <input type="checkbox"/> Air Bags (Passgr) | <input type="checkbox"/> Running Boards |
| <input type="checkbox"/> Power Mirrors | <input type="checkbox"/> Security System | <input type="checkbox"/> Bed Liner |
| <input type="checkbox"/> Power Seat (Driver) | <input type="checkbox"/> Rear Defroster | <input type="checkbox"/> Custom Bumper |
| <input type="checkbox"/> Power Seat (Passgr) | <input type="checkbox"/> Tilt Wheel | <input type="checkbox"/> Grill Guard |
| <input type="checkbox"/> Antilock Brakes | <input type="checkbox"/> Rear Wipers | <input type="checkbox"/> Winch |
| <input type="checkbox"/> Air Conditioning | <input type="checkbox"/> Tinted Windows | <input type="checkbox"/> Opt. Fuel Tank |
| <input type="checkbox"/> Towing Package | <input type="checkbox"/> Cup Holder | |
| <input type="checkbox"/> Utility | <input type="checkbox"/> Toolbox | |
| <input type="checkbox"/> Underbody Hoist | <input type="checkbox"/> Trailer Hitch | |
| <input type="checkbox"/> Hydraulic Lift | <input type="checkbox"/> Dual Rear Wheels | |
| <input type="checkbox"/> Rear Spoiler | <input type="checkbox"/> AM/FM | |
| <input type="checkbox"/> Pickup Shell | <input type="checkbox"/> CD Player | |
| <input type="checkbox"/> Tachometer | <input type="checkbox"/> D.A.B | |
| <input type="checkbox"/> Keyless Entry | | |
| <input type="checkbox"/> Digital Clock | | |

Example

▶ Data-rich free text document

2005 Honda new **Accord** Ex, Extra Clean, very **low Mileage**, Maintained By Dealer!
Vehicle Located in Stockton, Ca. Ad Id# 28147

This is a brand new car with **automatic transmission!**

Car with Air Conditioning, clock, **Cruise Control**, Digital Info Center, Dual Zone Climate Control, Heated Seats, Leather Steering Wheel, Memory Seat Position, Power Driver's Seat, **Power Steering**, **Power Brakes**, Power Passenger Seat, **Power Windows**, **Cup Holder**, **Rear Air Conditioning**, **Sunroof**, Tilt Steering Wheel, Original Owner, **Alloy Wheels.**

Am/fm, **Cd Changer**, Mp3, Satellite

Contact Us At XXX-XXXX-XXXX For More Information

Visit xxx xxx Motors

Example

► Form Filling

Vehicle Info

Type	<input type="text" value="- Please Select -"/>
Year	<input type="text" value="2005"/>
Make	<input type="text" value="Honda"/>
Model	<input type="text" value="Accord"/>
VIN	<input type="text"/>
Mileage	<input type="text" value="low"/>
Transmission	<input type="text" value="Automatic"/>
Engine	<input type="text"/>
Drivetrain	<input type="text" value="- Please Select -"/>
Body style	<input type="text" value="- Please Select -"/>
Color	<input type="text"/>
Int color	<input type="text"/>
Int material	<input type="checkbox"/> Cloth <input type="checkbox"/> Leather
Seating	<input type="text"/>
Wheels	<input type="text" value="Alloy Wheels"/>
Tires	<input type="text" value="- Please Select -"/>
Roof	<input type="text" value="- Please Select -"/>
Truck bed	<input type="text" value="- Please Select -"/>
Stereo	<input type="text" value="- Please Select -"/>
Dealer code	<input type="text"/>
Stock code	<input type="text"/>
MSRP	<input type="text"/>
NADA	<input type="text"/>
KBB	<input type="text"/>
Warranty	<input type="text" value="- Please Select -"/>

Features

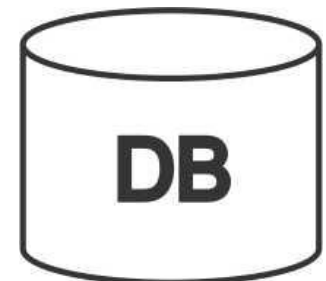
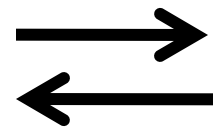
<input checked="" type="checkbox"/> Power Steering	<input checked="" type="checkbox"/> Air Cond. (Rear)	<input type="checkbox"/> Roof Rack
<input checked="" type="checkbox"/> Power Brakes	<input checked="" type="checkbox"/> Cruise Control	<input type="checkbox"/> Fog Lamps
<input checked="" type="checkbox"/> Power Windows	<input type="checkbox"/> Air Bags (Driver)	<input type="checkbox"/> Sliding Rear Win
<input type="checkbox"/> Power Locks	<input type="checkbox"/> Air Bags (Passgr)	<input type="checkbox"/> Running Boards
<input type="checkbox"/> Power Mirrors	<input type="checkbox"/> Security System	<input type="checkbox"/> Bed Liner
<input type="checkbox"/> Power Seat (Driver)	<input type="checkbox"/> Rear Defroster	<input type="checkbox"/> Custom Bumper
<input type="checkbox"/> Power Seat (Passgr)	<input type="checkbox"/> Tilt Wheel	<input type="checkbox"/> Grill Guard
<input type="checkbox"/> Antilock Brakes	<input type="checkbox"/> Rear Wipers	<input type="checkbox"/> Winch
<input type="checkbox"/> Air Conditioning	<input type="checkbox"/> Tinted Windows	<input type="checkbox"/> Opt. Fuel Tank
<input type="checkbox"/> Towing Package	<input checked="" type="checkbox"/> Cup Holder	
<input type="checkbox"/> Utility	<input type="checkbox"/> Toolbox	
<input type="checkbox"/> Underbody Hoist	<input type="checkbox"/> Trailer Hitch	
<input type="checkbox"/> Hydraulic Lift	<input type="checkbox"/> Dual Rear Wheels	
<input type="checkbox"/> Rear Spoiler	<input checked="" type="checkbox"/> AM/FM	
<input type="checkbox"/> Pickup Shell	<input type="checkbox"/> CD Player	
<input type="checkbox"/> Tachometer	<input type="checkbox"/> D.A.B	
<input type="checkbox"/> Keyless Entry		
<input type="checkbox"/> Digital Clock		

Common usage of Web Forms

- ▶ A user manually fills each form field
 - ▶ Text-box, selection list, check-box and radio button
- ▶ Tedious, error prone and repetitive process

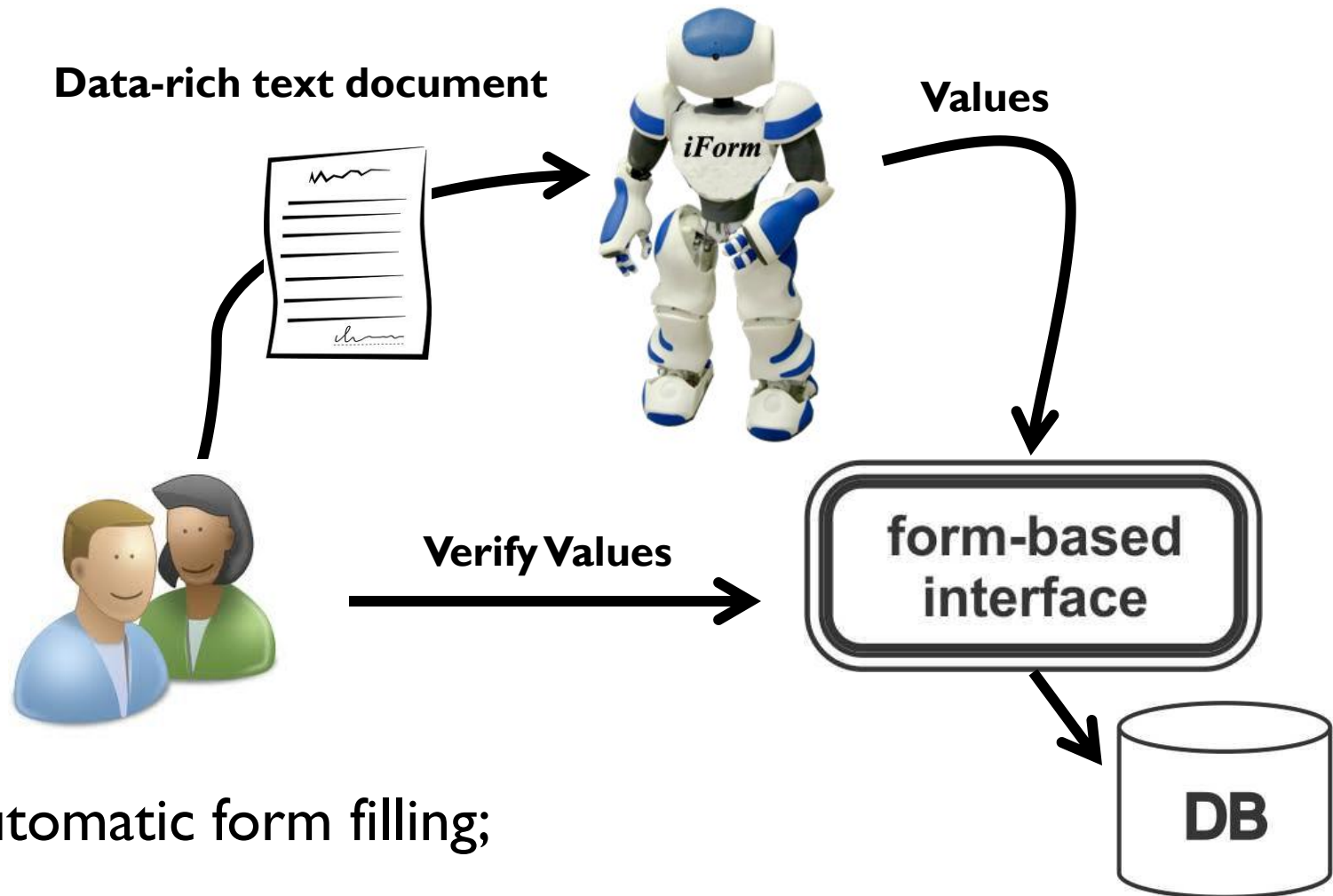


values →



iForm

- ▶ Information Extraction + Form Filling



- ▶ Automatic form filling;

iForm - Scenario



Shutter Island is a 2010 American psychological thriller film directed by Martin Scorsese. The film is based on Dennis Lehane's 2003 novel of the same name. Starring Leonardo DiCaprio, Mark Ruffalo and Ben Kingsley.

Movie Review - Data-rich text



Web Form

Web Form

Movie TV Show

Title:

Director:

Actors:

Gender:

iForm – Selecting plausible segments

- ▶ *Is this text segment a suitable value of a given field of the form?*

Shutter Island is a 2010 American psychological thriller film directed by Martin Scorsese. The film is based on Dennis Lehane's 2003 novel of the same name . Starring Leonardo DiCaprio, Mark Ruffalo and Ben Kingsley.

Shutter

Shutter Island

Shutter Island is

Shutter Island is a

...

Leonardo

Leonardo DiCaprio

Kingsley.

Redundant computation of several features can be avoided by using dynamic programming.

iForm - Features

► Features Considered:

$$TAF(F_j, S_{ab}) = \eta \sum_{\tau \in \text{tokens}(S_{ab})} \frac{\text{freq}(\tau, F_j)}{\sum_{F_i \in \mathcal{F}} \text{freq}(\tau, F_i)}$$

$$\eta = \frac{1}{k + |\text{avg}(F_j) - k|}$$

Attribute Vocabulary

$$VAF(F_j, S_{ab}) = \frac{\text{freq}(S_{ab}, F_j)}{\sum_{F_i \in \mathcal{F}} \text{freq}(S_{ab}, F_i)}$$

Attribute Value

$$\frac{\sum_{\langle n_x, n_y \rangle \in \text{path}(\mathbf{p})} w(SM(F_j), n_x, n_y)}{|\text{path}(\mathbf{p})|}$$

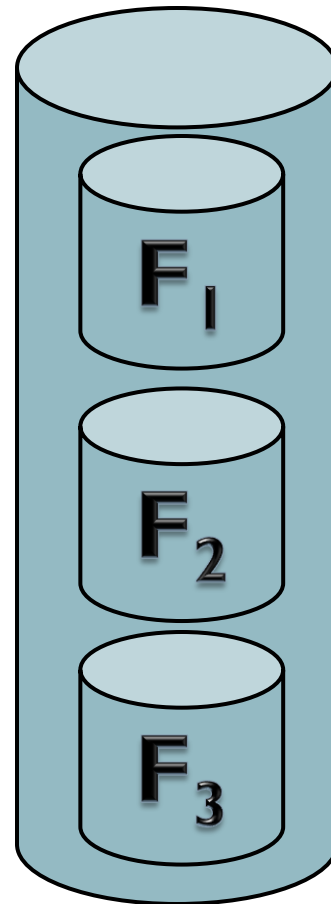
Value Format

Noisy
OR

Title

Shutter Island

Previous
Submissions

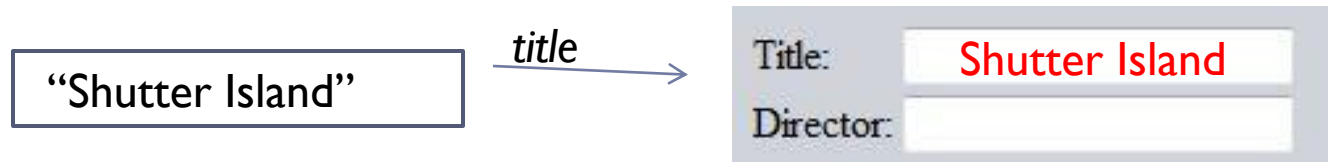


iForm – Mapping Segments to Fields

- ▶ Given the set of text segments such that their scores are above a threshold ϵ
 - ▶ iForm aims at finding a **mapping** between candidate values and form fields with a **maximum aggregate score**
 - ▶ Select non-overlapping segments.
- ▶ Accomplished by means of a two-phase procedure

iForm – Filling Form-based interfaces

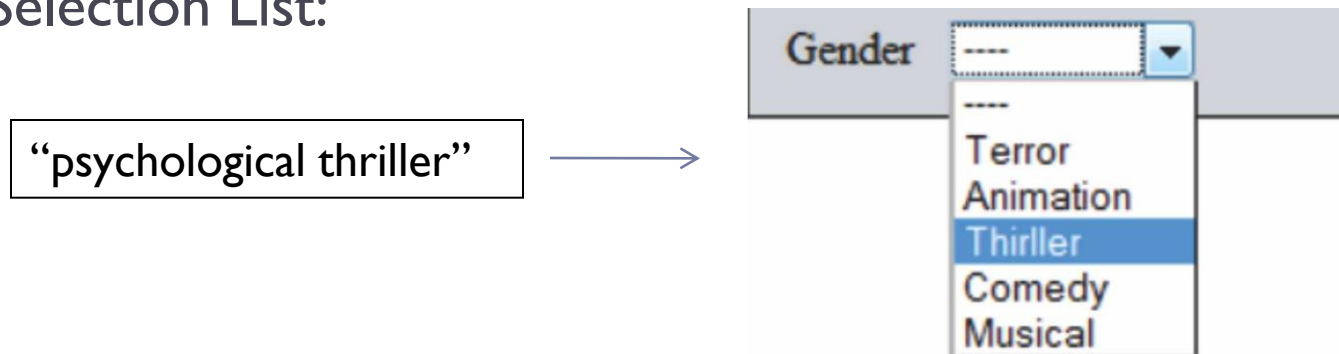
- ▶ Uses the final mapping to fill out the form fields
 - ▶ Text Boxes: Mapped text segments as a field values.



- ▶ Check boxes: *Set true for mapped fields.*



- ▶ Selection List:



iForm - Overview



Shutter Island is a 2010 American psychological thriller film directed by Martin Scorsese. The film is based on Dennis Lehane's 2003 novel of the same name. Starring Leonardo DiCaprio, Mark Ruffalo and Ben Kingsley.

Web Form

Web Form	
<input checked="" type="checkbox"/> Movie	<input type="checkbox"/> TV Show
Title:	Shutter Island
Director:	Martin Scorsese
Actors:	Leonardo DiCaprio Mark Ruffalo Ben Kingsley
Gender	Thriller

Experiments

Dataset	Test Data	Previous Data	# Fields	S - Test Data	S – Previous Data
Jobs	50	100	13	RISE	RISE
Movies	50	10000	4	IMDb	FreeBase / Wikipedia
Cars	50	10000	35	TodaOferta.com	TodaOferta.com
Cellphones	50	10000	37	TodaOferta.com	TodaOferta.com
Books 1	50	10000	5	Submarino.com	TodaOferta.com
Books 2	50	10000	4	Submarino.com	Ingenta
Books 3	50	10000	2	Submarino.com	Ourpress.com
Books 4	50	10000	3	Submarino.com	NetLibrary

▶ Baseline

- ▶ iCRF - a method for interactive form filling based on CRF
- ▶ The Jobs dataset was used for an experimental comparison between iForm and iCRF.

Evaluation – Multi-typed web forms

Movies

Type of Field	# Fields	P	R	F
Text Box	4	0.74	0.69	0.71
Submission-Level		0.73	0.67	0.69

iForm achieved high quality results in all datasets

Cellphones

Type of Field	# Fields	P	R	F
Text Box	2	0.89	0.69	0.78
Check Box	35	0.94	0.94	0.94
Average		0.94	0.93	0.93
Submission-Level		0.96	0.94	0.95

Filling quality above 0.90. In fact, more than 90% of each submission was correctly entered in the web form interface.

Evaluation – Comparison with iCRF

Jobs

Field	iForm	iCRF
Application	0.82	0.37
Area	0.18	0.23
City	0.70	0.65
Company	0.41	0.17
Country	0.77	0.87
Desired Degree	0.57	0.37
Language	0.84	0.69
Platform	0.47	0.38
Recruiter	0.44	0.22
Req. Degree	0.31	0.59
Salary	0.22	0.25
State	0.85	0.81
Title	0.72	0.49

iForm had superior F-measure levels in nine fields.

The lower quality obtained by iCRF is explained by the fact that segments to be extracted from typical free text inputs, such as jobs postings, may not appear in a regular context.

iForm was designed to conveniently exploit these field-related features from previous submissions

Conclusions

- ▶ This work proposes an unsupervised approach to the IETS problem.
 - ▶ Relies on information available on pre-existing data.
 - ▶ Exploit content-based features to directly learn from test data structure-based features.
 - ▶ Show that pre-existing datasets allow for the unsupervised learning of both content-based and structure-based features.
 - ▶ Eliminate the need of a user involved in any source specific training process.
- ▶ **Information Extraction Methods:**
 - ▶ ONDUX, JUDIE and iForm

Publications

► Thesis Core

1. Joint Unsupervised Structure Discovery and Information Extraction. **SIGMOD Conference** – 2011
2. Unsupervised Information Extraction with the ONDUX Tool. **Brazilian Symposium on Databases (SBBD)** – 2011
3. On Using Wikipedia to Build Knowledge Bases for Information Extraction by Text Segmentation. **Journal of Information and Data Management (JDIM)** – 2011
4. ONDUX: on-demand unsupervised learning for information extraction. **SIGMOD Conference**. - 2010
5. Unsupervised strategies for information extraction by text segmentation. **SIGMOD PhD Workshop on innovative Database Research (IDAR)** – 2010
6. A Probabilistic Approach for Automatically Filling Form-Based Web Interfaces. **Proceedings of the VLDB Endowment (PVLDB)** – 2010
7. Automatically filling form-based web interfaces with free text inputs. **International Conference on World Wide Web (WWW)** – 2009

Publications

▶ Related to the Information Extraction Problem

8. Building a research social network from individual perspective. **Joint Conference on Digital Libraries (JCDL) – 2011**
9. CiênciaBrasil – The Brazilian Portal of Science and Technology. **Integrated Seminar of Software and Hardware (Semish)– 2011**
10. A flexible approach for extracting metadata from bibliographic citations. **Journal of the American Society for Information Science and Technology (JASIST) – 2009**

Publications

▶ Other Publications

11. Lightweight methods for large-scale product categorization. **Journal of the American Society for Information Science and Technology (JASIST)** – 2011
12. Adaptive and Flexible blocking for record linkage tasks. **Journal of Information and Data Management (JDIM)** – 2010
13. Blocagem adaptativa e flexível para o pareamento aproximado de registros. **Brazilian Symposium on Databases (SBBD)** – 2009

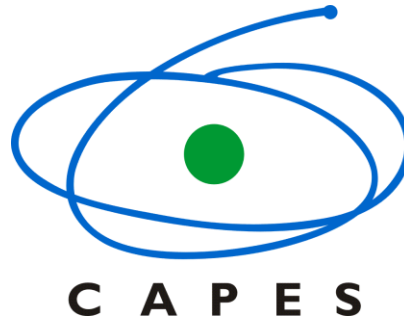
▶ Tutorials

14. Methods and techniques for information extraction by text segmentation. **Alberto Mendelzon International Workshop on Foundations of Data Management (AMW)** - 2012
15. Methods and techniques for information extraction by text segmentation. **Brazilian Symposium on Databases (SBBD)** - 2011

Future Work

- ▶ Generating transductive methods using domain knowledge
- ▶ Use our approach to extract information from HTML
- ▶ Query Extraction using our unsupervised approach
- ▶ Extraction Improvement Through User Feedback

Acknowledgments



UFAM



FAPEAM
Fundação de Amparo à Pesquisa
do Estado do Amazonas



Unsupervised Information Extraction by Text Segmentation



Eli Cortez

Advisor: Altigran Soares da Silva

Universidade Federal do Amazonas

