



**„Al. I. Cuza” University, Iasi, Romania**



**Faculty of Computer Science**

# **Grammar-based Automatic Extraction of Definitions. Applications for Romanian**

**Adrian Iftene, Diana Trandabăț, Ionuț Pistol  
{adiftene, dtrandabat, ipistol}@info.uaic.ro**

September, 26, 2007, Borovets, Bulgaria

# Overview

- Introduction
- Romanian grammar
  - Categorization of Definitions
  - Distribution of the definitions into types
  - Rules
  - Evaluation
- Applications
  - Question Answering
  - Textual Entailment

## Introduction

- Under the framework of the project LT4eL was created an environment for collecting and (semi)automatic exploiting language resources (Monachesi et al, 2006)
- 9 languages involved (bul, cze, dut, eng, ger, mal, pol, por and rom)
- Manually annotation of keywords, definitions of various terms and semantic concepts
- A grammar was created for the automatic identification of definitions in texts

## Categorization of Definitions

- **“is\_def”** – *“HTML **este** tot un protocol folosit de World Wide Web.”* (HTML **is** also a protocol used by World Wide Web).
- **“verb\_def”** – *“Poșta electronică **reprezintă** transmisia mesajelor prin intermediul unor rețele electronice.”* (Electronic mail represents sending messages through electronic networks).
- **“punct\_def”** – *“Bit – prescurtarea pentru binary digit”* (Bit – shortcut for binary digit)

## Categorization of Definitions (cont...)

- “**layout\_def**”

<b>Ro:</b>	
<i>Organizarea datelor</i>	<i>Cel mai simplu mod de organizare este cel secvențial.</i>
<b>En:</b>	
<b>Data organizing</b>	<b>The simplest method is the sequential one.</b>

- “**pron\_def**” – “...definirii conceptului de baze de date. *Acesta descrie metode de ...*” (...defining the database concept. It describes methods of ....)
- “**other\_def**” – “*triunghi echilateral, adică cu toate laturile egale*” (equilateral triangle **i.e.** having all sides equal).

## Distribution of the definitions

Type	Manual	%	Automatic	%
is_def	70	33.8	204	32.8
<b>verb_def</b>	116	<b>56.0</b>	272	<b>43.8</b>
<b>punct_def</b>	15	<b>7.2</b>	124	<b>20.0</b>
layout_def	2	1.0	21	3.4
pron_def	4	2.0	0	0.0
<b>Total</b>	<b>207</b>		<b>621</b>	

# Rules

- Simple grammar rules
- Composed grammar rules

- "is\_def" grammar rule:

```

<rule name="may_be_term">
  <seq>
    <query match="tok[@base='fi' and
                  substring(@ctag,1,5)='vmip3']"/>

    <first>
      <ref name="UndefNominal" />
      <ref name="DefNominal" />
    </first>
  </seq>
</rule>

```

# Evaluation

- Lxtransduce (Tobin 2005) is used to match the grammar in files

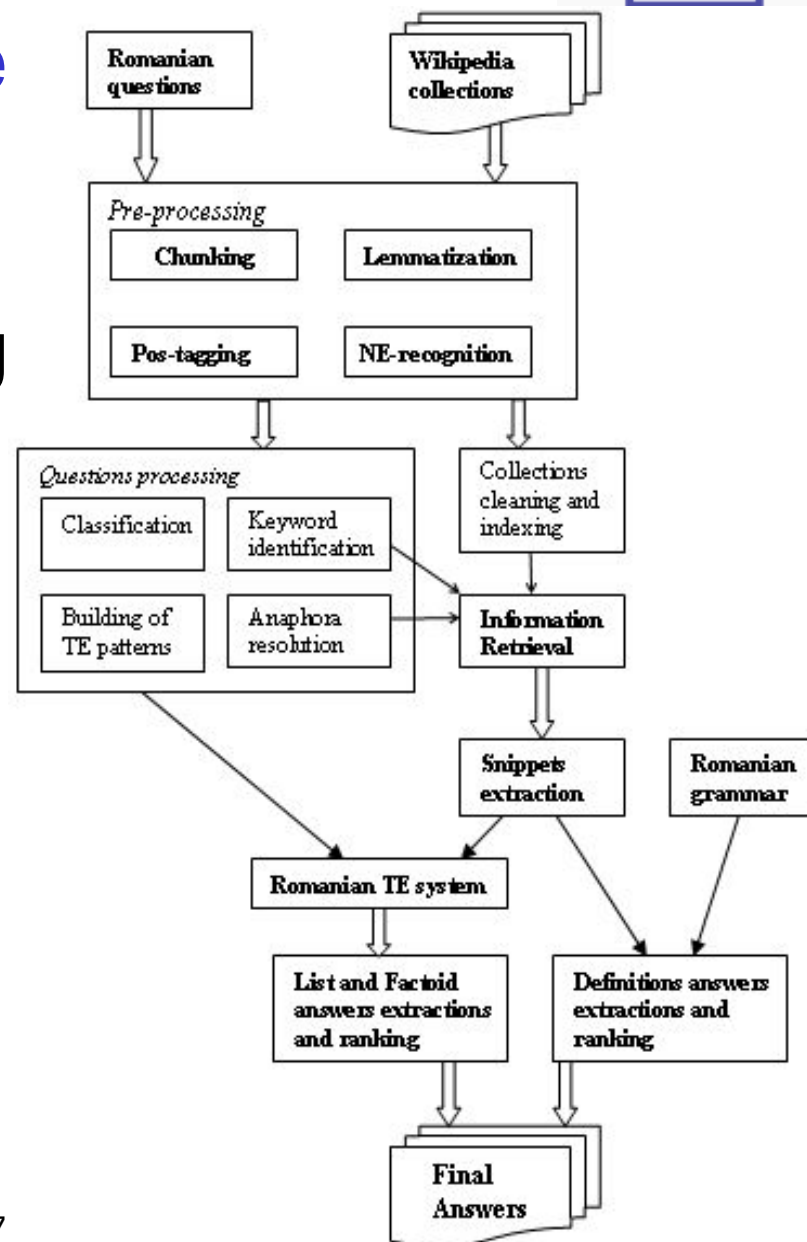
Definition Type	Result
is_def	<b>Sentence-level matching:</b> P: 0.5366, R: <b>1.0</b> , F2: 0.7765 <b>Token-level matching:</b> P: 0.0648, R: 0.3328, F2: 0.14
<b>verb_def</b>	<b>Sentence-level matching</b> P: 0.7561, R: <b>1.0</b> , F2: 0.9029 <b>Token-level matching</b> P: 0.0471, R: 0.1422, F2: 0.085
punct_def	<b>Sentence-level matching</b> P: 0.1463, R: <b>1.0</b> , F2: 0.3396 <b>Token-level matching</b> P: 0.0025, R: 0.1163, F2: 0.0072
layout_def	<b>Sentence-level matching</b> P: 0.0488, R: <b>1.0</b> , F2: 0.1333 <b>Token-level matching</b> P: 0.0007, R: 0.1020, F2: 0.0022

## Question Answering

- Accordingly to the answer type, we have the following type of questions (Harabagiu, Moldovan 2007):
  - **Factoid** – *“Who discovered the oxygen?” or “When did Hawaii become a state?” or “What football team won the World Cup in 1992?”.*
  - **List** – *“What countries export oil?” or “What are the regions preferred by the Americans for holidays?”.*
  - **Definition** – *“What is quasar?” or “What is a question-answering system?”*

# QA: System architecture

- Pre-processing
- Questions processing
- Collection indexing
- Information retrieval
- Snippet extraction
- Answer extraction



## QA – Example

- Question: Cine este Zeus?

D: 0026: (Cine, zeus, PERSON)

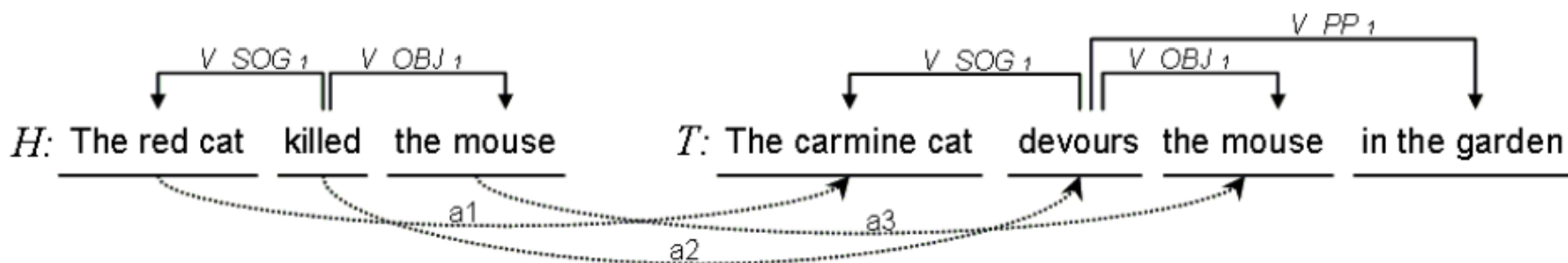
- Snippet:

0026#10014#1.0#Zeus#Zeus\zeus\ASN\P  
 este\fi\V3\ cel\cel\TSR\ mai\mai\R\  
 puternic\puternic\ASN\ dintre\dintre\S\  
 olimpieni\olimpieni\NPN\ ,\,\COMMA\  
 socotit\socoti\VP\ drept\drept\S\  
 stăpânul\stăpân\NSRY\ suprem\suprem\ASN\  
 al\al\TS\ oamenilor\om\NPOY\ și\și\CR\  
 al\al\TS\ zeilor\zeu\NPOY\ .\.\PERIOD\

- Our pattern for "is\_def" (0026.\*\zeus\.\*\P  
 .\*\fi\V3\ (.\*)) match the snippet

# Textual Entailment

- TE is defined (Dagan et al., 2006) as a directional relation between two text fragments, termed *T* (text) - the entailing text, and *H* (hypothesis) - the entailed text.
- It is then said that *T* entails *H* if, typically, a human reading *T* would infer that *H* is most likely true.
- Example:
  - T: The carmine cat devours the mouse in the garden.
  - H: The red cat killed the mouse.



## TE - Background Knowledge

- Snippet extraction for NEs from H without corresponding value in T. For each such context:
  - a) “core” identification
  - b) left NEs extraction
  - c) right NEs extraction
  - d) calculate LNEs X RNEs

## BK - Example

<pair id="748" entailment="YES">

<T>Argentina President Carlos Menem has ordered an 'immediate' investigation into war crimes allegedly committed by British troops during the 1982 Falklands War.</T>

<H>Argentine demanded an investigation of alleged war crimes during the Falklands War.</H>

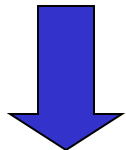
</pair>

# BK - Example

“Argentine”: Extracted Snippets from Wikipedia:

```

ar |calling_code = 54 |footnotes = Argentina also has a
territorial dispute
Argentina', , NaciÃ³n Argentina (Argentine Nation) for many
legal purposes), is
in the world. Argentina occupies a continental surface area of
Argentina national football team
  
```



Argentine [is] Argentina

Netherlands [is] Holland

2 [is] two

Los Angeles [in] California

Chinese [in] China

Netherlands [is] Dutch

Netherlands [is] Nederlandse

Netherlands [is] Antillen

Netherlands [in] Europe

Netherlands [is] Holland

Antilles [in] Netherlands

# Conclusions

- We presented the Romanian grammar used in the European LT4eL project
- The definitions were devised in 6 types
- Applications: QA and TE
- Apply the grammar to a new corpus

# Acknowledgments

- Special thanks goes to the other members of the Romanian team in the LT4eL project, Dan Cristea and Corina Forăscu
- We also acknowledge the help provided by Claudia Borg

**THANK YOU!**