# The PAISA' Project

Marco Baroni

HLT for Italian Workshop

# PAISA'

- **P**iattaforma per l'**A**pprendimento dell'**I**taliano **S**u Corpora **A**nnotati
- MIUR FIRB project 2009-2012
- Partners:
  - University of Bologna: <u>Sergio Scalise</u>, Claudia Borghetti, Emiliano Guevara
  - University of Trento: <u>Marco Baroni</u>, Marco Brunello, Sara Castagnoli, Egon Stemle
  - ILC, Pisa: <u>Vito Pirrelli</u>, Alessandro Lenci, Felice Dell'Orletta
  - EURAC, Bolzano: <u>Andrea Abel</u>, Verena Lyding, Christopher Culy

# Goals

- Build and annotate Web-derived "reference" corpus of Italian language
- Innovative corpus visualization and exploration tools
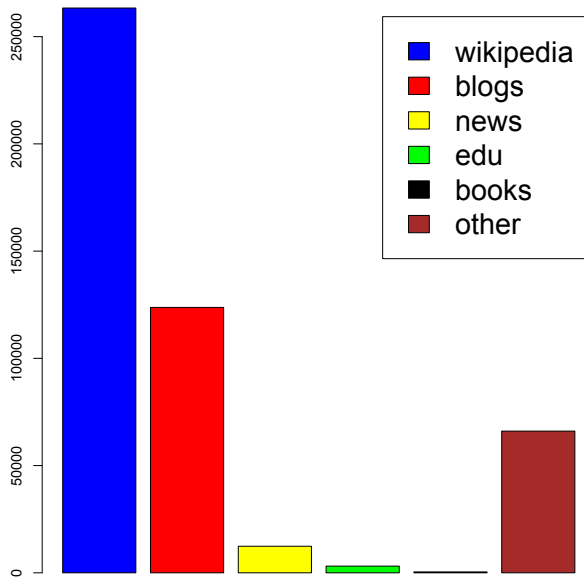- Applications in Italian language and culture teaching

# The WaCky initiative

- ▶ Very large Web-crawled, automatically cleaned and annotated corpora for English, French, German and Italian (2005-2009)
- ▶ Lessons learnt:
  - ▶ Copyright
  - ▶ Corpus composition
  - ▶ Cleaning
  - ▶ Annotating Web text

# Copyright

- PAISA' corpus constructed by crawling Web pages published under CreativeCommons License only:
    - Attribution
    - Derived works OK
    - Non-commercial derived works only
- Same licensing propagates to PAISA' corpus

# Corpus composition

470K documents

# Web cleaning:
## The boilerplate problem

# Web cleaning on steroids: `krdwrd`
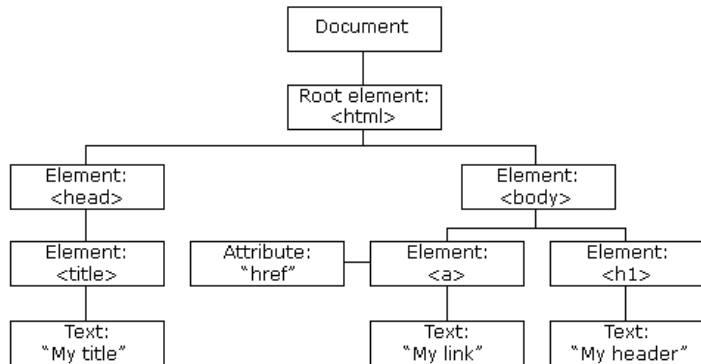
Main Index Page

General Ratings Page

One true sign of a truly great band is when said band
ardently defies categorisation, that is, when for every
"well, they sound like this reggae-influenced heavy metal
band playing avantgarde bebop" remark you can have
yourself a "funny, I thought they were this raw punk
outfit doing acoustic folk" counterproposal. And I don't
simply mean "being diverse" here, I mean "being different".
Blazing off every colour of the spectrum. Baring one's
soul in all of its existing aspects. That sort of thing.

READER COMMENTS SECTION

Return to the Index page! NOW!

# Web cleaning on steroids: `krdwrd`

Cues from the DOM tree (figure from `http://www.w3schools.com/`)

# Web cleaning on steroids: `krdwrd`

## Cues from visual rendering

# Web cleaning on steroids

Cross-validation on gold standard

|                  | features | precision | recall | F   |
|------------------|----------|-----------|--------|-----|
| wacky heuristics | NA       | 80%       | 99%    | 88% |
| text cues only   | 21       | 92%       | 93%    | 92% |
| DOM cues only    | 13       | 89%       | 91%    | 90% |
| visual cues only | 8        | 90%       | 93%    | 91% |
| full krdwrd      | 42       | 93%       | 92%    | 92% |

# Corpus comparison

- ITWAC: 1.9B tokens
- PAISA': 700M tokens
- La Repubblica: 380M tokens

# PAISA' vs. La Repubblica

| nouns | | verbs | | adjs | |
|---|---|---|---|---|---|
| *PAISA'* | *Rep.* | *PAISA'* | *Rep.* | *PAISA'* | *Rep.* |
| punto | miliardo | registrare | dire | informatico | scorso |
| blog | anno | pubblicare | chiedere | simile | politico |
| commento | presidente | leggere | andare | esterno | pubblico |
| settembre | governo | segnare | spiegare | hi-tech | primo |
| guida | ministro | inserire | arrivare | opzionale | stesso |
| articolo | lira | riservare | restare | famoso | grande |
| galleria | paese | inviare | parlare | successivo | generale |
| video | giorno | usare | sembrare | gay | finanziario |
| set | partito | contattare | cominciare | digitale | americano |
| pubblicità | stato | mostrare | mettere | bello | lungo |

# PAISA' vs. ITWAC

| nouns | | verbs | | adjs | |
|---|---|---|---|---|---|
| *PAISA'* | *ITWAC* | *PAISA'* | *ITWAC* | *PAISA'* | *ITWAC* |
| punto | numero | registrare | svolgere | ultimo | pubblico |
| blog | articolo | pubblicare | chiedere | pubblicitario | regionale |
| commento | legge | segnare | ritenere | hi-tech | presente |
| guida | comma | leggere | presentare | simile | sociale |
| settembre | servizio | riservare | dire | opzionale | previsto |
| galleria | lavoro | commentare | tenere | esterno | nazionale |
| set | presidente | mostrare | riguardare | informatico | necessario |
| video | decreto | contattare | andare | famoso | generale |
| pubblicità | commissione | continuare | consentire | nuovo | seguente |
| e-mail | caso | deselezionare | prevedere | grande | legislativo |

# Ongoing work

- ▶ Annotation (POS tagging, dependency parsing. . . )
- ▶ Classification by genre and topic
- ▶ Visualization, user-friendly interface
- ▶ E-mail me (`marco.baroni@unitn.it`) if you want a copy of the raw corpus!