# Building an infrastructure for accessing and analysing figures in scholarly publications

Piotr Praczyk – CERN, 10.03.2011
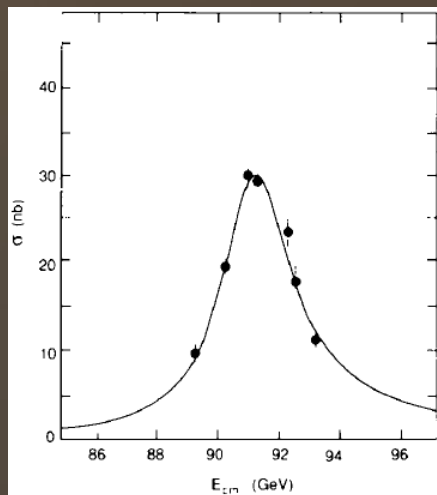
# *Usage of graphics in scholarly communication*

- Describe experiments

- Summarise large amounts of data

- Illustrate relations between results
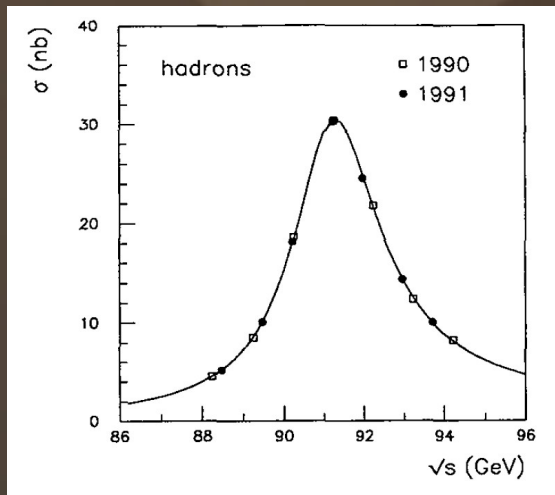
- Present ideas in a schematic manner
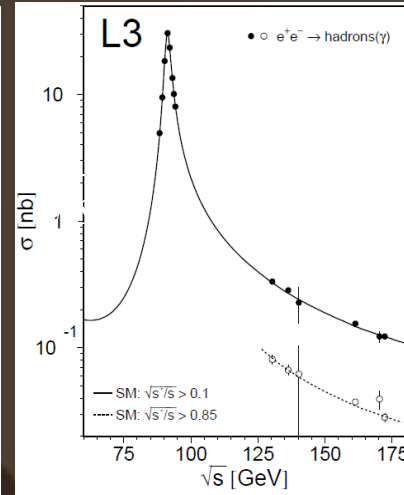
# *Different measurements of the same quantity*

Measured cross section for $e^+e^- \to$ hadrons as a function of $\sqrt{s}$



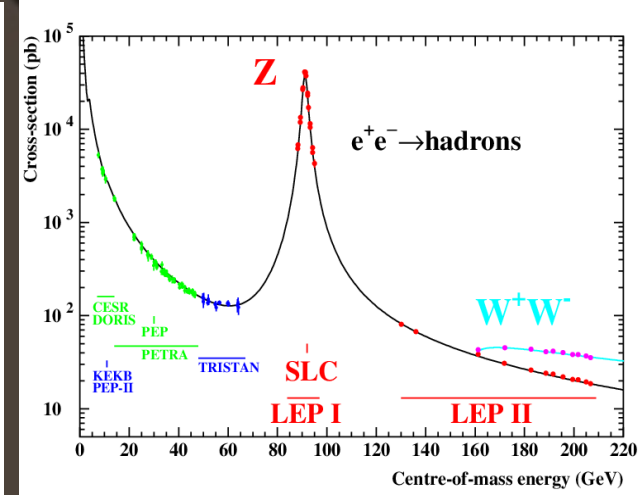1989        1991        1997        2005

# Current understanding of the work
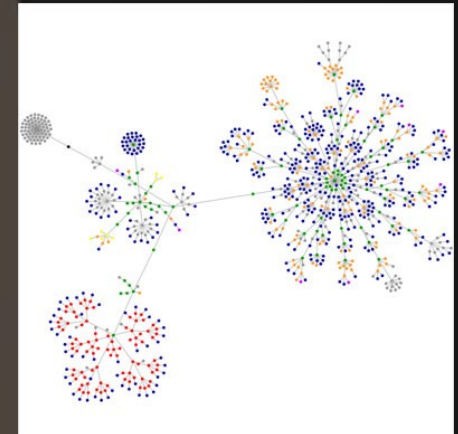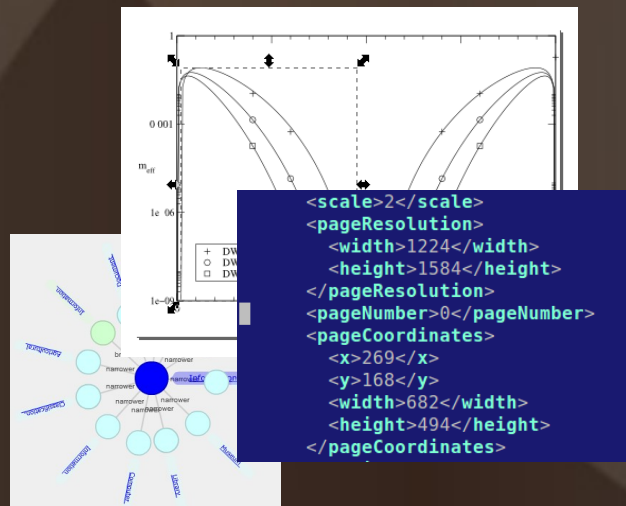
Extraction

Indexing

Scholarly publications
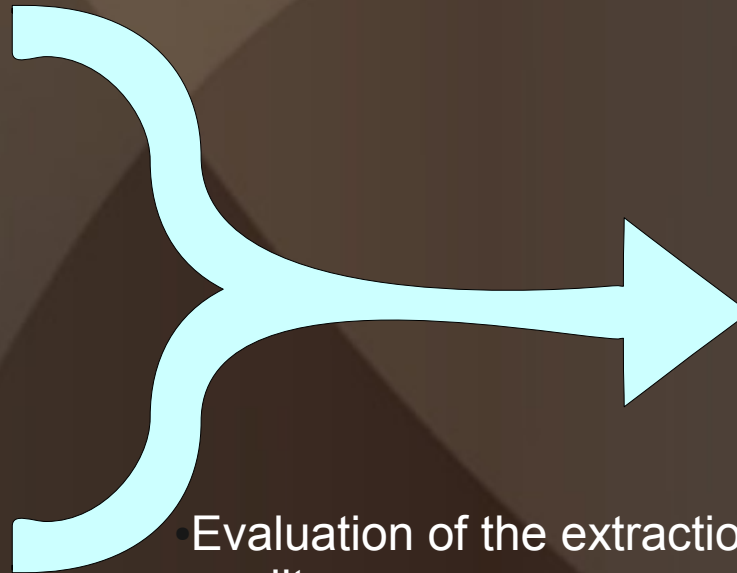
Description of figures as separate entities

Collective description Of figures

```xml
<scale>2</scale>
<pageResolution>
    <width>1224</width>
    <height>1584</height>
</pageResolution>
<pageNumber>0</pageNumber>
<pageCoordinates>
    <x>269</x>
    <y>168</y>
    <width>682</width>
    <height>494</height>
</pageCoordinates>
```
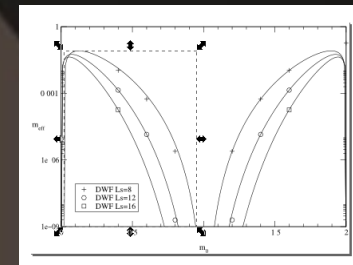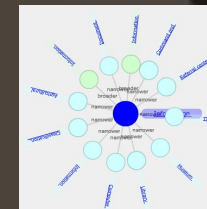
# *Automatic extraction of figures*



Meta-data



Vector + Raster images

- Evaluation of the extraction quality
- Merging of results
- Acquisition of additional data



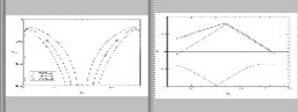(in the future)
Semantic description

# *Types of extracted meta-data*

- Boundaries of figures

- Boundaries of captions

- Text of captions

- Graphics in PNG and SVG formats

- Places, where figure is referenced

- Name of the figure inside a document

- Text present inside the figure

# *Select Your Figure*

# *Extracting data from PDF*

PDF:

- Stream of instructions
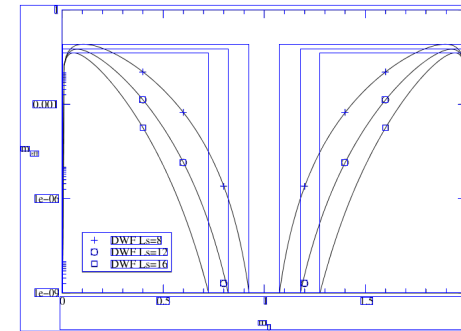- Embeded objects
  - Fonts
  - External objects
- Meta-description

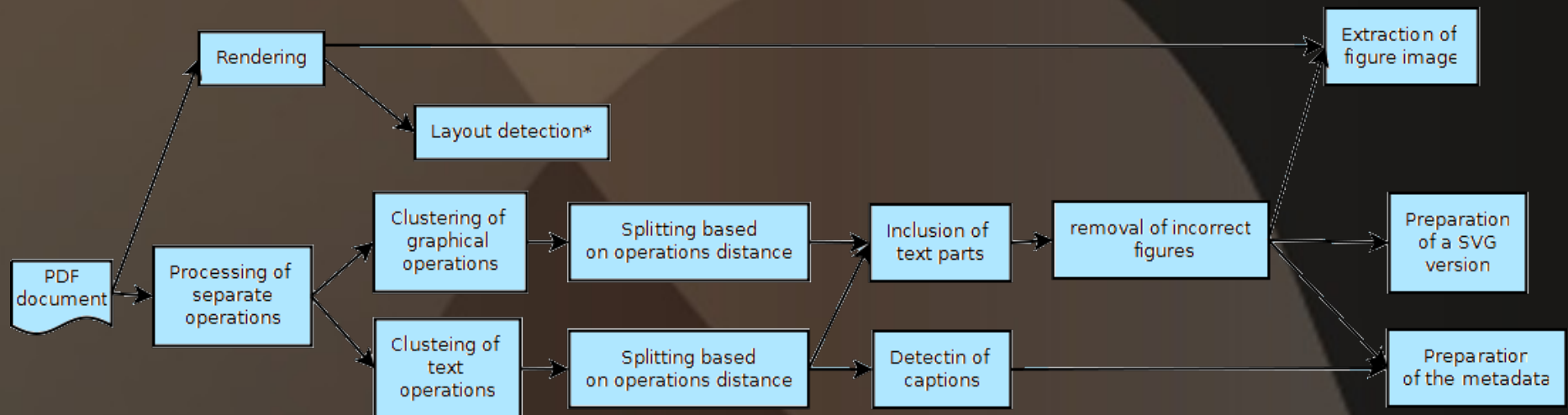# *Schema of the PDF extraction process*

# *Intermediate steps of the algorithm*



- Regions of graphics (blue)
  - Clustered graphic operations
- Regions of text (green)
  - Clustered text operations
- Elements of page layout (red)
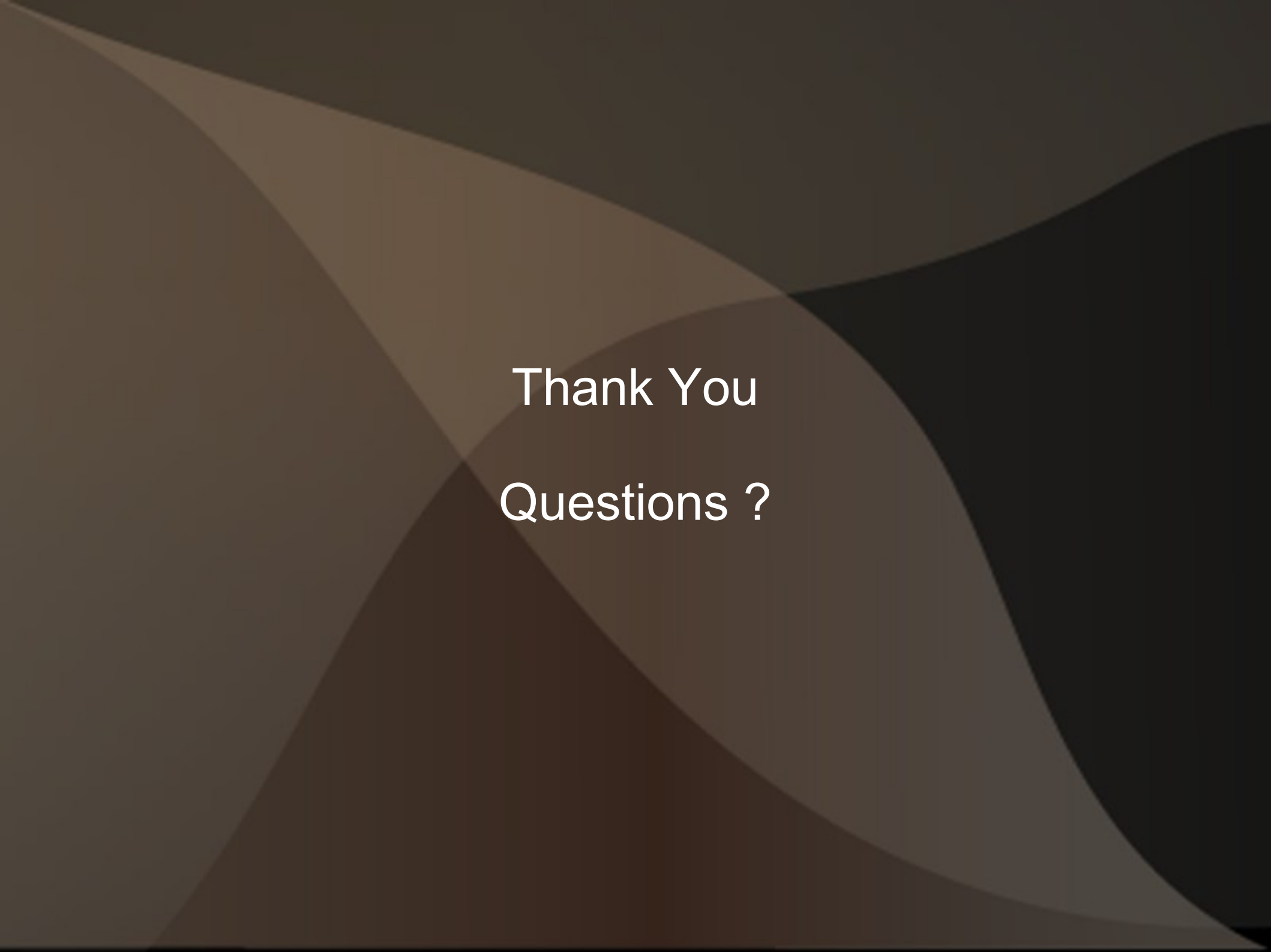
# *Future work*

- Finishing work on the PDF extractor and selction interface

  Extracting + tagging semantics of images

- Similarity measure based on semantic and graphical properties of figures

- Extraction of data described by figures

- Improvements in extractor algorithms
  - usage of different algorithms
  - usage of different types of data that are produced)

# Thank You

# Questions ?