


## Language Technology

underlying technologies for  
research infrastructure in the  
social sciences and humanities

**Tamás Váradi**  
Hungarian Academy of Sciences  
varadi@nytud.hu

7th December 2005 ECRIUK Conference Nottingham 1



## The Role of Language

- Dominant medium now and in future
- Congenial for humans, challenge for machines
- not merely streams of characters, words but rather concepts, properties, relations, arguments, facts etc.
- accessible to humans, partially so to machines
- information explosion of electronic texts
- tapping the huge resources of the web


7th December 2005 ECRIUK Conference Nottingham 2



## Language Technology

- Intelligent handling of language through automated procedures
  - information extraction/data mining
  - text understanding
  - speech recognition
  - semantic web
  - multilingual information extraction
  - machine translation

7th December 2005 ECRIUK Conference Nottingham 3



## Language Resources

- Data
  - Corpora – linguistically annotated textual databases
    - monolingual and parallel aligned corpora
  - lexicons
  - terminological databases
  - ontologies
- Annotation technology
- Tools


7th December 2005 ECRIUK Conference Nottingham 4



## Tools

- Tokenizers
- Morphological analyzers
- Taggers
- Syntactic parsers
- Semantic taggers
- Parallel corpus aligners
- Automatic term extractors
- Intelligent data mining

7th December 2005 ECRIUK Conference Nottingham 5



## The TELRI Story

- 1996 – 2002 EU founded networking project
- A truly Pan-European exercise
- Consortium of 26 nations involving 24 languages
- Tried and tested network of research centres of excellences
- Representing state of the art language technology expertise

7th December 2005 ECRIUK Conference Nottingham 6



## TRACTOR

TELRI RESEARCH ARCHIVE OF  
COMPUTATIONAL TOOLS AND  
RESOURCES

- Massive collection of language resources in 24 languages
- Accessible on-line and available free for research
- Currently hosted by Oxford University
- Not (just) for linguists and not necessarily by linguists!

7th December 2005

ECRIUK Conference Nottingham

7



## Summary

- Modern research methodology in SSH requires effective means to deal with the bottleneck presented by the explosive growth of textual data in electronic form
- Language technology is a fundamental enabling infrastructure that can be deployed in many areas to address the need to process texts intelligently
- SSH research infrastructure should greatly benefit from the utilization of state of the art language technologies

7th December 2005

ECRIUK Conference Nottingham

8