

LINDAT-Clarín

Language Research Infrastructure

Jan Hajič

Charles University in Prague

Czech Republic

LINDAT Project principal investigator



ZÁPADOČESKÁ
UNIVERZITA
V PLZNI



LM2010013

Supported by



Research in the area of Natural Language



- Language research paradigm:
 - Empirical research
 - Theoretical linguistics: evidence, counterexamples, ...
 - Quantitative and formal linguistics: preferences, patterns, ...
 - Computational linguistics / Natural Language Processing
 - Statistical (machine) learning methods prevalent
- Need for language-related data
 - Language Corpora
 - Raw & analyzed (manually, automatically)
 - Monolingual, parallel
 - Text, audio, video, multimedia
 - As much as possible, easy (free) access



ZÁPADOČESKÁ
UNIVERZITA
V PLZNI



LINDAT-Clarín at
CESAR Roadshow
June 7, 2012

LM2010013
Supported by



Research Infrastructure



- Provides support for research
 - Measurements, data, expertise, ...
- Language infrastructure
 - Language and language-related data
 - Easy access and licensing
 - Possibility to include researcher's data
 - Tools for language processing
 - Web services



ZÁPADOČESKÁ
UNIVERZITA
V PLZNI



LINDAT-Clarín at
CESAR Roadshow
June 7, 2012

LM2010013
Supported by



LINDAT-Clarín



- Infrastructure for Language Research in...
 - Humanities (linguistics and related)
 - Computational Linguistics
- Provides
 - Data repository, linked to European networks
 - Clarin ERIC, META-SHARE
 - National Center for language data
 - Acquisition, specification, licensing
 - Linguistic and language-related annotation
 - Archiving, preservation
 - Distribution



ZÁPADOČESKÁ
UNIVERZITA
V PLZNI



LINDAT-Clarín at
CESAR Roadshow
June 7, 2012

LM2010013
Supported by



LINDAT-Clarín Partners



- Charles University in Prague (coordinator)
 - Institute of Formal and Applied Linguistics
 - Faculty of Mathematics and Physics, CS School
- Masaryk University (Brno)
 - Natural Language Lab
 - Faculty of Informatics
- University of West Bohemia (Pilsen)
 - Dept. of Cybernetics (Speech and Language group)
 - Faculty of Applied Sciences
- Czech Language Institute
 - Part of the Academy of Sciences of the Czech Rep.



ZÁPADOČESKÁ
UNIVERZITA
V PLZNI



LINDAT-Clarín at
CESAR Roadshow
June 7, 2012

LM2010013
Supported by



LINDAT-Clarín



- Project Timeframe
 - 2010-2015: funding secured (EUR 4.7M)
 - Started in October 2010
 - 2010: building the technical backbone, hiring
 - New server room
 - Ready for expansion (up to 400 CPU, 100TB)
 - 2011-2013
 - Construction phase: node of Clarín ERIC (Feb. 2012)
 - » Compatibility with other networks (META-SHARE)
 - 2014-
 - Operational phase



ZÁPADOČESKÁ
UNIVERZITA
V PLZNI



LINDAT-Clarín at
CESAR Roadshow
June 7, 2012

LM2010013
Supported by



LINDAT-Clarín

Project structure



- Clarín ERIC coordination (Netherlands, MPI)
 - National Coordinator
- Full Clarín “Centre”
 - Node of the Clarín ERIC network
 - 24/24 service, data repository, full compatibility
 - Software and services development, contracts and agreements with service providers, licensing
- National Centre for Czech Language Resources
 - acquisition, conversion, description (metadata)
 - annotation, preservation
 - evaluation data, shared tasks preparation
 - methodology development (e.g. crowdsourcing)



ZÁPADOČESKÁ
UNIVERZITA
V PLZNI



LINDAT-Clarín at
CESAR Roadshow
June 7, 2012

LM2010013
Supported by



LINDAT-Clarín Clarín Centre



- Data-related services
 - Repository (dSpace)
 - Universal access and authentication
 - Single sign-on (Shibboleth) within Clarín ERIC
 - » 2011: problematic in several respects; local registration added
 - Permissions, licensing (active, passive)
 - Persistency, archiving, metadata API (Clarín ERIC, META)
 - Persistent ID assignment (EPIC handles)
- Language services (plans only, not started yet)
 - Clarín ERIC compatible API (Weblicht?), workflows
 - (Mostly) Czech-language related
 - morphology, tagging, parsing, anaphora resolution, ...



ZÁPADOČESKÁ
UNIVERZITA
V PLZNI



LINDAT-Clarín at
CESAR Roadshow
June 7, 2012

LM2010013
Supported by



LINDAT-Clarín

Building Czech Resources



- All partners... history of resource creation
 - Charles University in Prague
 - Prague Dependency Treebank, parallel corpora for Machine Translation, spoken language corpora, NLP tools
 - Digital Libraries: The Malach project (Visual History, USC)
 - University of West Bohemia
 - Speech-related data (ASR, TTS), speech tools
 - Masaryk University
 - Lexical resources and tools
 - Institute of Czech Language
 - Lexical and historic language resources
- LINDAT-Clarín: continued resource building



ZÁPADOČESKÁ
UNIVERZITA
V PLZNI



LINDAT-Clarín at
CESAR Roadshow
June 7, 2012

LM2010013
Supported by



LINDAT-Clarín

Education and Training



- University partners
 - Master's and Ph.D. programmes
 - Linguistics, Computer Science, NLP, ...
- LINDAT opportunities for students
 - NLP and speech-related experiments
 - new data specification annotation
 - linguistic research
- Training
 - Winter school(s), tutorials for researchers
 - Czech researchers & international



ZÁPADOČESKÁ
UNIVERZITA
V PLZNI



LINDAT-Clarín at
CESAR Roadshow
June 7, 2012

LM2010013
Supported by



LINDAT-Clarín

International Cooperation



- Language Resources (world-wide)
 - Same needs, similar datasets
 - Since early 1990s (Statistical methods in NLP)
- History of cooperation
 - LDC (since 2001), ELDA (since 2005)
 - MU, UWB and CUNI has already provided some datasets
 - PDT, U.S. Malach project, Czech Wordnet, parallel corpora
 - Evaluation datasets (CLEF, CoNLL, MT – several years)
- Cooperation with research organizations
 - Assessment of needs, research plans
 - Future / emerging technologies: data first
 - Tools (latest results – high quality tools, including for Czech)



ZÁPADOČESKÁ
UNIVERZITA
V PLZNI



LINDAT-Clarín at
CESAR Roadshow
June 7, 2012

LM2010013
Supported by



Access to Resources



- Repository
 - Technological use – clear (download)
 - Accompanied with software for technological use
 - Problem: use in humanities, some licenses
- Use in humanities (Clarín-like)
 - Plan: web applications based on web services
 - Current technology:
 - Web-based access
 - Example:
 - LINDAT lexical resources
 - » VALLEX 2.6: <http://ufal.mff.cuni.cz/vallex/2.6>
 - » PDT-Vallex 2.5+PDT: <http://ufal.mff.cuni.cz/lindat/PDT-Vallex.html>
 - Dialogs: <http://ujc.dialogy.cz>, <http://ufal.mff.cuni.cz/~toman/pdtsc>, <http://ufal.mff.cuni.cz/~toman/pdtse>
 - Treebanks
 - » PEDT <http://ufal.mff.cuni.cz/pedt2.0>
 - » PCEDT <http://ufal.mff.cuni.cz/pcedt2.0>



ZÁPADOČESKÁ
UNIVERZITA
V PLZNI



LINDAT-Clarín at
CESAR Roadshow
June 7, 2012

LM2010013
Supported by



Plans for 2012 and beyond: repository and coordination



- Clarin ERIC – finally!
 - Start collaboration with Clarin-D, NL etc.
 - Authentication – top priority
 - Present licensing solution developed at UFAL
 - Central site? (search / metadata sharing)
 - Learn from META-SHARE experience
 - Implementation of Clarin standards
 - Compatibility (conversions) with META-SHARE
 - Legal issues (cooperation with Clarin-D)



ZÁPADOČESKÁ
UNIVERZITA
V PLZNI



LINDAT-Clarin at
CESAR Roadshow
June 7, 2012

LM2010013
Supported by



Plans for 2012 – Czech-related data



- Acquisition of data
 - Goal:
 - 2 GW in repository
 - 5 GW total (searchable)
- Data annotation
 - Towards PDT 3.0:
 - Bridging anaphora
 - Discourse
 - Corrections
 - Spoken annotation
 - PDT-like annotation on (reconstructed) spoken corpora



ZÁPADOČESKÁ
UNIVERZITA
V PLZNI



LINDAT-Clarín at
CESAR Roadshow
June 7, 2012

LM2010013
Supported by



Plans for 2012 - dissemination



- LINDAT-Clarín
 - Lectures at various events
 - CLARA EU winter school – Feb. (tutorial, treebanks)
 - Workshop at LREC (META-NET organized)
 - Advanced treebank annotation
- Offering services
 - Web-based, access to large corpora
 - Annotation services
 - such as the tasks done for ELRA on small scale in 2010/11



ZÁPADOČESKÁ
UNIVERZITA
V PLZNI



LINDAT-Clarín at
CESAR Roadshow
June 7, 2012

LM2010013
Supported by



LINDAT-Clarín



<http://www.lindat.cz>

Thank you for your attention!



ZÁPADOČESKÁ
UNIVERZITA
V PLZNI



LINDAT-Clarín at
CESAR Roadshow
June 7, 2012

LM2010013
Supported by

