# D-Spin & WebLicht

Thomas Zastrow

University Tübingen

thomas.zastrow@uni-tuebingen.de

# Outline

- The D-Spin Project

- Clarin Federation

- WebLicht

- Further Work

# The D-Spin Project

# The D-Spin Project

- *eScience is about global collaboration in key areas of science and the next generation of infrastructure that will enable it.* (J. Taylor)

- The aim of D-SPIN along with CLARIN is to establish a *virtual research infrastructure* based on available language resources and tools

- D-Spin stands for *Deutsche Sprach Resourcen Infrastructure* (German Language Resource Infrastructure)

- D-SPIN is the German contribution to the European CLARIN-Projekt

# The D-Spin Project intends:

- 7 to 9 centres are created

- These centres will work together within a resource-provider-federation and are embedded in the **Clarin Federation infrastructure** (AAI = Authentification and Authorization Infrastructure)

- German resources, data, and tools are gradually made available via **state-of-the-art-registries and web services**

- The regulatory framework is designed in such a way that researchers from DFN AAI institutions are able to merge and exchange data

- Simple workflow models and tools can be defined on these interoperable resources

- Various projects are accomplished along with humanists, in order to develop specific solutions and basic services

- Training sessions are carried out

# The D-Spin Project

- Partners are:
  - Max Planck Institute for Psycholinguistics, Nijmegen
  - Department of Linguistics, Computational Linguistics, Tübingen (Coordinator)
  - IDS - Institute for the German Language, Mannheim
  - BBAW - Berlin-Brandenburgische Akademie der Wissenschaften, Berlin
  - ASV - Department of Computer Science, NLP Group, Leipzig
  - Comparative Linguistics, Frankfurt a.M.
  - DFKI - German Research Center for Artificial Intelligence, Saarbrücken
  - IMS - Institute for Natural Language Processing, Stuttgart
  - FB05 - Applied Linguistics and Computational Linguistics, Gießen

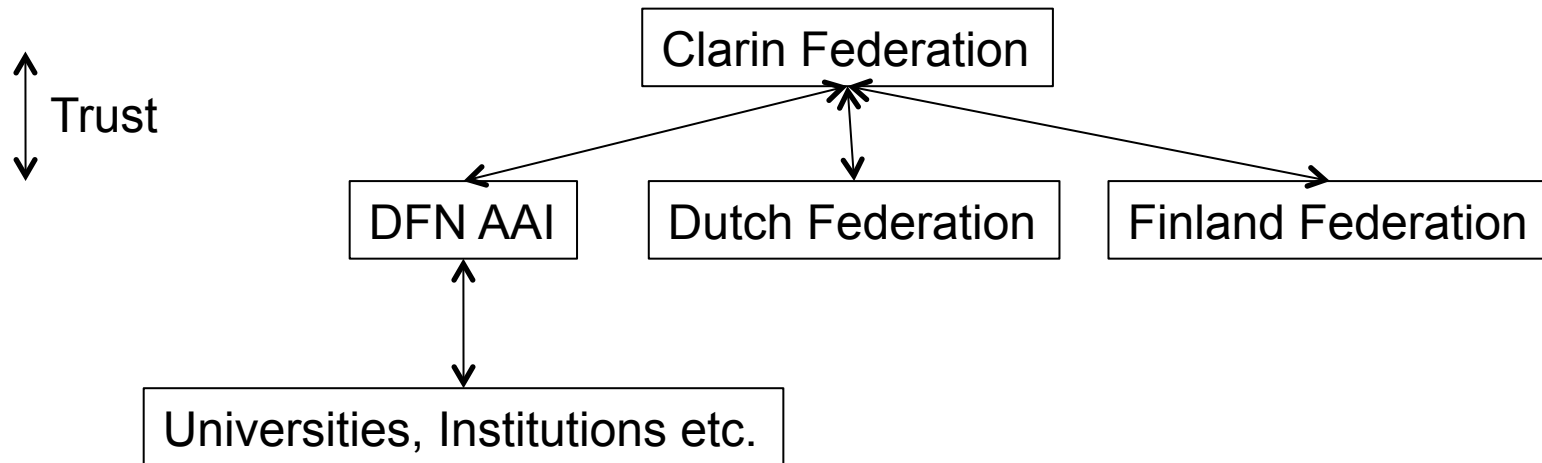# The Clarin Federation

# The Clarin Federation

- The CLARIN Service Provider Federation consists of a set of formal agreements and some software technology specifications to authenticate and identify the researchers in a secure way when they work on distributed language resources and applications

- The initial CLARIN Service Provider Federation already gives more than one million researchers and students from 3 countries (Netherlands, Finland, Germany) access to resources and applications of the participating centers

# The Clarin Federation

- The Clarin federation is based on the Single-Sign-On system *Shibboleth*

- Shibboleth is already implemented as a nationwide service in many countries

# The Clarin Federation

# WebLicht

# WebLicht - Motivation

- Many linguistic resources (corpora, dictionaries, …) and tools (tokenizer, tagger, parser, …) are available

- Most of them are implemented to run on local machines. This can be inconvenient and error-prone

- One possible solution: *Make them available on the web!*

# WebLicht - Motivation

- For some kinds of LRT, its easy to put them online (make resources downloadable, offer search engines etc.)

- For other kinds, more effort is necessary (limiting access to resources, how to make tools online usable)

- Solution: a *Service Oriented Architecure (SOA)*

  - ➔ The end user needs just a browser: no more installation, configuration etc. of software is necessary

# WebLicht - Architecture

- WebLicht is a *SOA* for building annotated text corpora

- Development started in October 2008

- WebLicht consists of the following components:

  - **Distributed Services:** offering functionality (resources & tools) over the (inter-)net. Implemented as webservices

  - **Repository:** stores metadata and technical information about the services

  - **User interface:** interacts with the user and combines services and information from the repository. Access still possible via scripts / programming code

# WebLicht - Architecture

*Stuttgart*

*Tübingen*

*Leipzig*



Repository

Standard-conformant
Text Corpus Encoding

Web 2.0 Application for
Tool Chaining
and Execution

Stuttgart    Tübingen    Berlin    Leipzig    Finland    Romania

# WebLicht – The Services

- Services are implemented as REST style webservices

- HTTPs POST method is used to send data from the UI to the services

- As client, anything which is able to use the HTTP protocol, can be used:

    - Browser

    - Commandline tools (wget, curl)

    - Programming Languages

- Anyone can implement his/her own interface to WebLicht

# WebLicht – The Repository

- Implemented at the ASV Leipzig

- It offers information and a query engine for the services:

  - Which services are available?

  - How can I combine them?

  - Which input/output format does a service accept/produce?

- Example: a tokenizer is already applied to a plain text, which services can be used next?

# WebLicht – The User Interface

- Web 2.0 Application for Tool Chaining and Execution
- Implemented at the SfS Tübingen
- Java application, deployed in Apache Tomcat
- Allows the user to:
  - upload a text (plain text, MS Word, RTF or PDF files)
  - construct a text from corpora in Leipzig
  - use some hardwired example texts
  - Build a chain of linguistic tools
  - Executes the tool chain with the uploaded text and presents the results
- During the chaining process, it queries the repository for available services

# WebLicht – Integrating new Services

- Building a webservice for WebLicht consists of the following steps:

  - Create a RESTstyle webservice around the tool as wrapper

  - Make in- and output compatible with WebLicht's TCF format

  - Register the service in the repository

Input in
TCF format
→

**Wrapper (RESTstyle webservice)**

**Tool or Resource**

Output in
TCF format
→

- *You can find further information and a tutorial online:*

  - *http://weblicht.sfs.uni-tuebingen.de/englisch/weblicht.shtml*

# WebLicht - Combinations
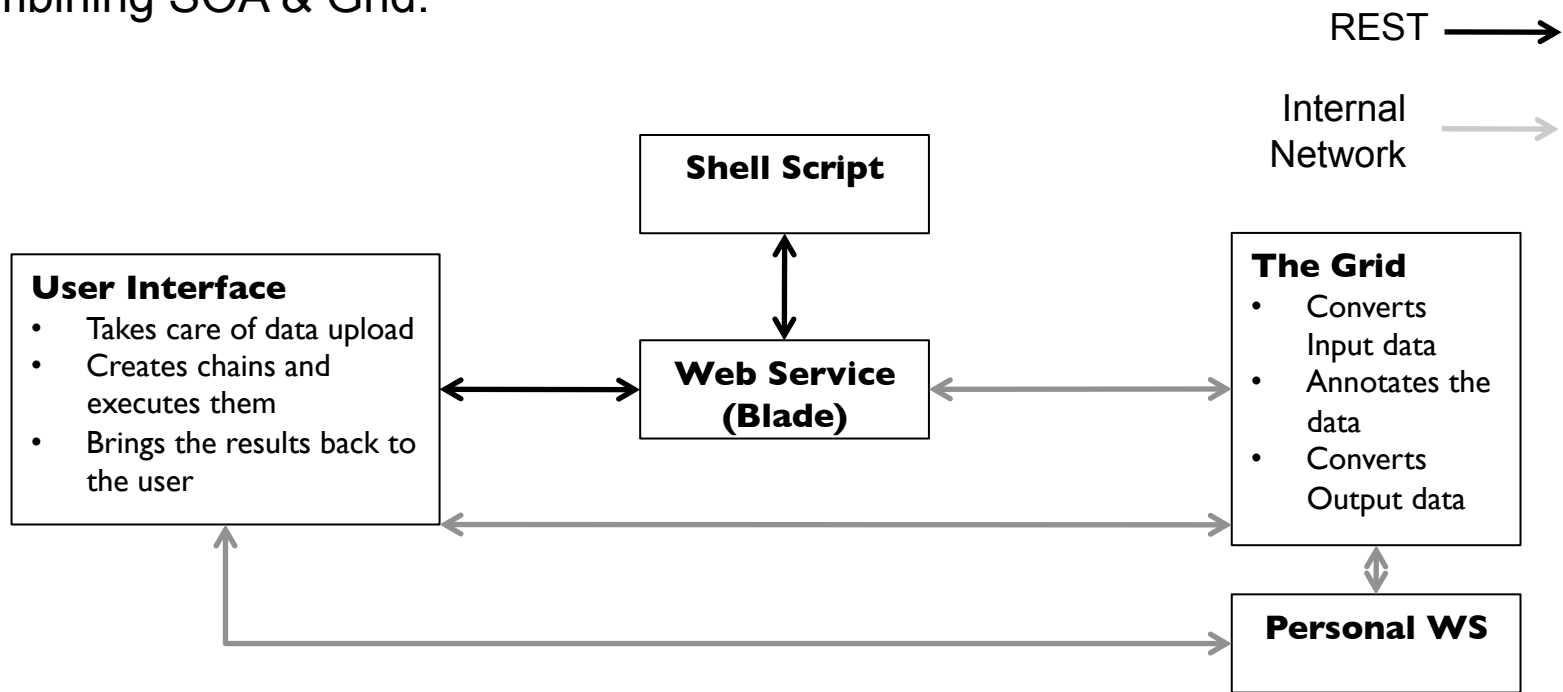
# WebLicht – Further Work

- Add more services and functionality
- Usability tests
- Make WebLicht more scaleable
- Creating personal workspaces
- Implement an asynchronous workflow, based on Grid technology

# WebLicht goes Grid

Combining SOA & Grid:

REST →

Internal
Network →

**Shell Script**

**User Interface**
- Takes care of data upload
- Creates chains and executes them
- Brings the results back to the user

**Web Service (Blade)**

**The Grid**
- Converts Input data
- Annotates the data
- Converts Output data

**Personal WS**

# Links etc.

- The D-Spin homepage: http://www.d-spin.org
- DFN AAI-Federation based login to WebLicht and some other webapplications: https://weblicht.sfs.uni-tuebingen.de/

Thomas Zastrow
Seminar für Sprachwissenschaft
Universität Tübingen

Wilhelmstr. 19
D-72074 Tübingen

thomas.zastrow@uni-tuebingen.de
http://www.thomas-zastrow.de

Tel.: 07071/29-73968
Fax: 07071/29-5214