# WebLicht – A Service Oriented Architecture for Language Resources and Tools

Erhard Hinrichs, Thomas Zastrow

University of Tübingen

Munich, October 17 2009

## Current Situation

- 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 ….

- …on the user side:

  - Every potential user has to download and install them on his own machine: this may cause problems with operating systems, compiler versions, missing libraries, …

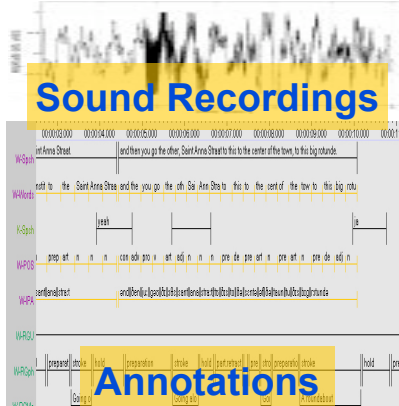  - Keep an eye on updates, (security) patches, new versions etc.

## Current Situation

- … on the developers side:

  - How to publish LRT?

  - Question of license, user permissions, …

  - Combination and comparison with other tools/resources
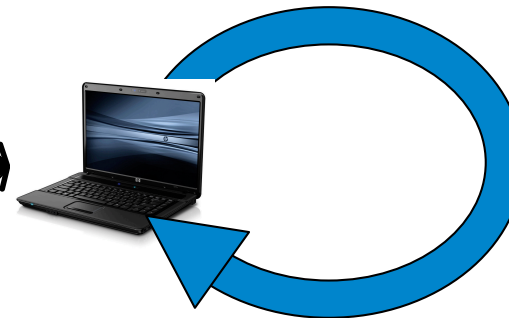
  - Sustainability, long term support

# Current Download-First Paradigm

**Video Recordings**

**Sound Recordings**

**Annotations**

**Texts**

- installing
- downloading
- adapting
- converting
- scripting
- etc etc

- some like it very much as I liked the VW Käfer
- not very efficient, rights problems, etc
- many are cut off since one needs IT skills
- cyberspace needs to overcome this scenario

Munich, October 17 2009

## One Possible Solution

### - Make LRT available on the web! –

- 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)

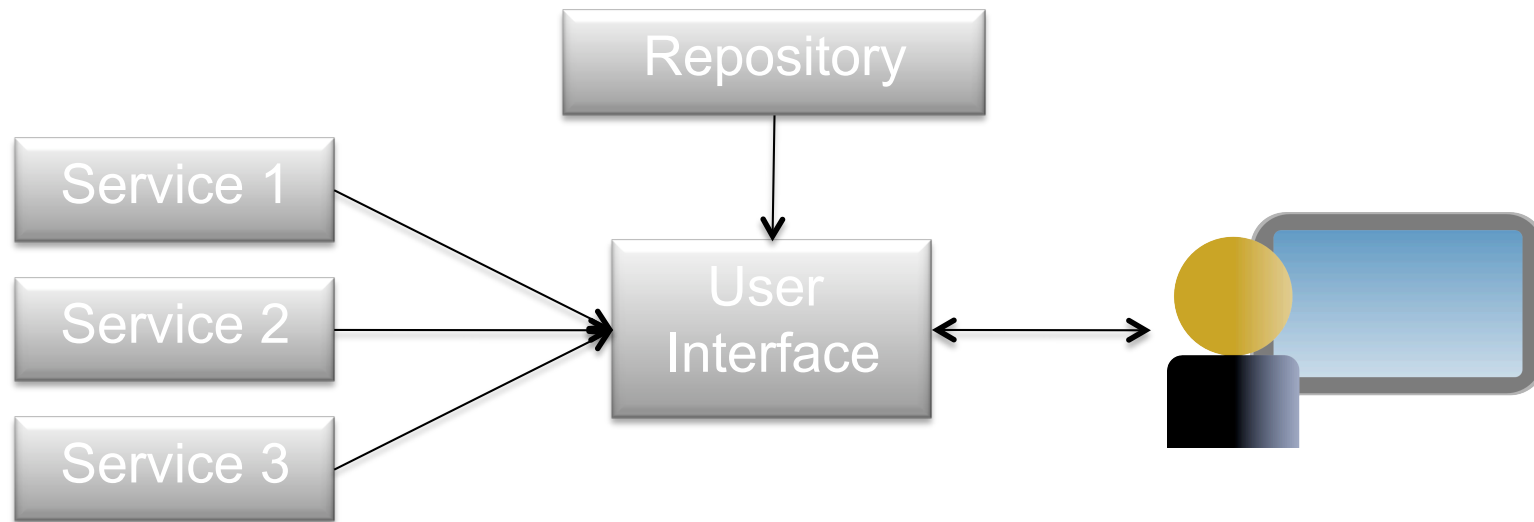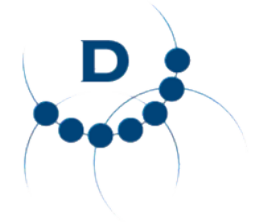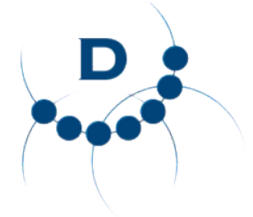   ➔ Solution: a *Service Oriented Architecure* **(SOA)**

- WebLicht: a Service Oriented Architecture for incremental automatic annotation of text corpora

- Work started in October 2008

- Participants (September 2009):

    - BBAW Berlin

    - ASV Leipzig

    - IDS Mannheim

    - IMS Stuttgart

    - SfS Tübingen

## Service Oriented Architectures

- Components of a SOA

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

Munich, October 17 2009

## 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

Munich, October 17 2009

## 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?

## 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

Munich, October 17 2009

# WebLicht

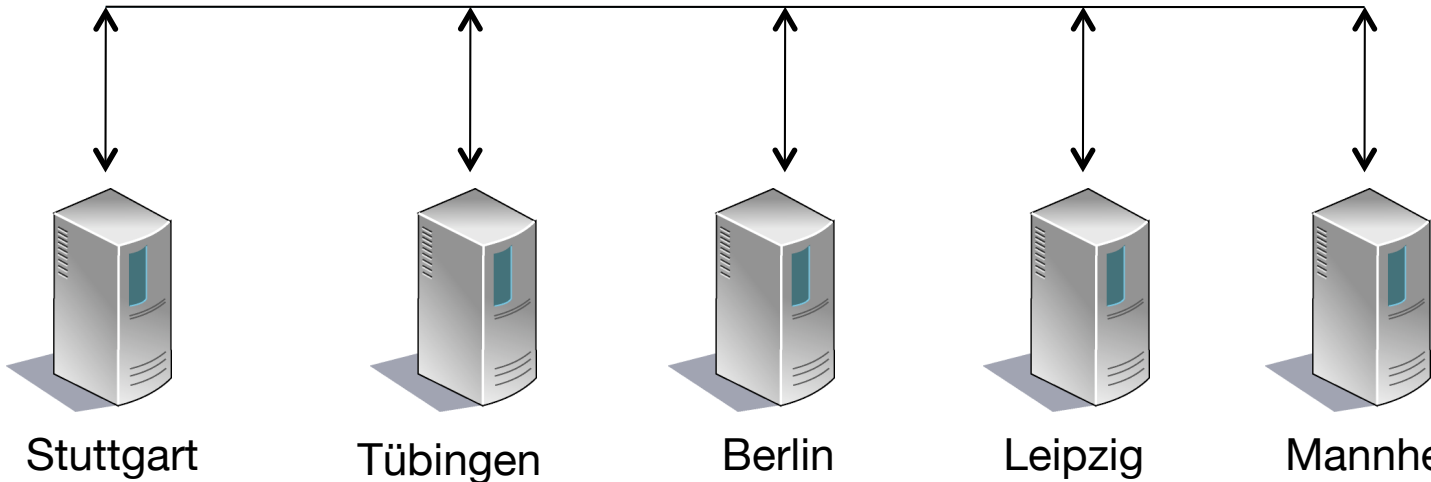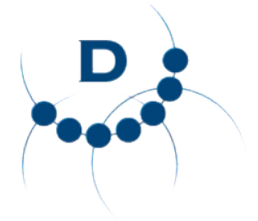*Stuttgart*

*Tübingen*

*Leipzig*



Standard-conformant
Text Corpus Encoding

Web 2.0 Application for
Tool Chaining
and Execution

Repository

Stuttgart

Tübingen

Berlin

Leipzig

Mannheim

Munich, October 17 2009

# WebLicht

## D-SPIN

CLARIN

**Live Presentation ....**

Munich, October 17 2009