

# TextGrid: der GWES im Einsatz

Grid Workflow Workshop 2009  
TU Berlin, 5./6. März 2009

Martin Haase  
DAASI International GmbH  
Martin.Haase@DAASI.de



## Überblick

- **GWES? Grid Workflow Execution Service, die Workflow Engine von Fraunhofer FIRST**
- **Hintergrund: TextGrid und Anforderungen an Workflow**
- **Prototyp im Einsatz**
- **Ausblick: Erweiterungen für Benutzerfreundlichkeit**



## TextGrid: Projekthintergrund

- D-Grid-Projekt, Laufzeit 10/2005 – 05/2009
- „Modulare Plattform für verteilte und kooperative wissenschaftliche Textdatenverarbeitung - ein Community-Grid für die Geisteswissenschaften“
- Partner: TU Darmstadt, SUB Göttingen, IDS Mannheim, Universität Trier, FH Worms, Universität Würzburg, DAASI International GmbH, Saphor GmbH



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT



INSTITUT FÜR  
DEUTSCHE SPRACHE

**DAASI**  
International

Directory Applications  
for Advanced Security  
and Information Management



BAYERISCHE JULIUS-MAXIMILIANS  
UNIVERSITÄT  
WÜRZBURG

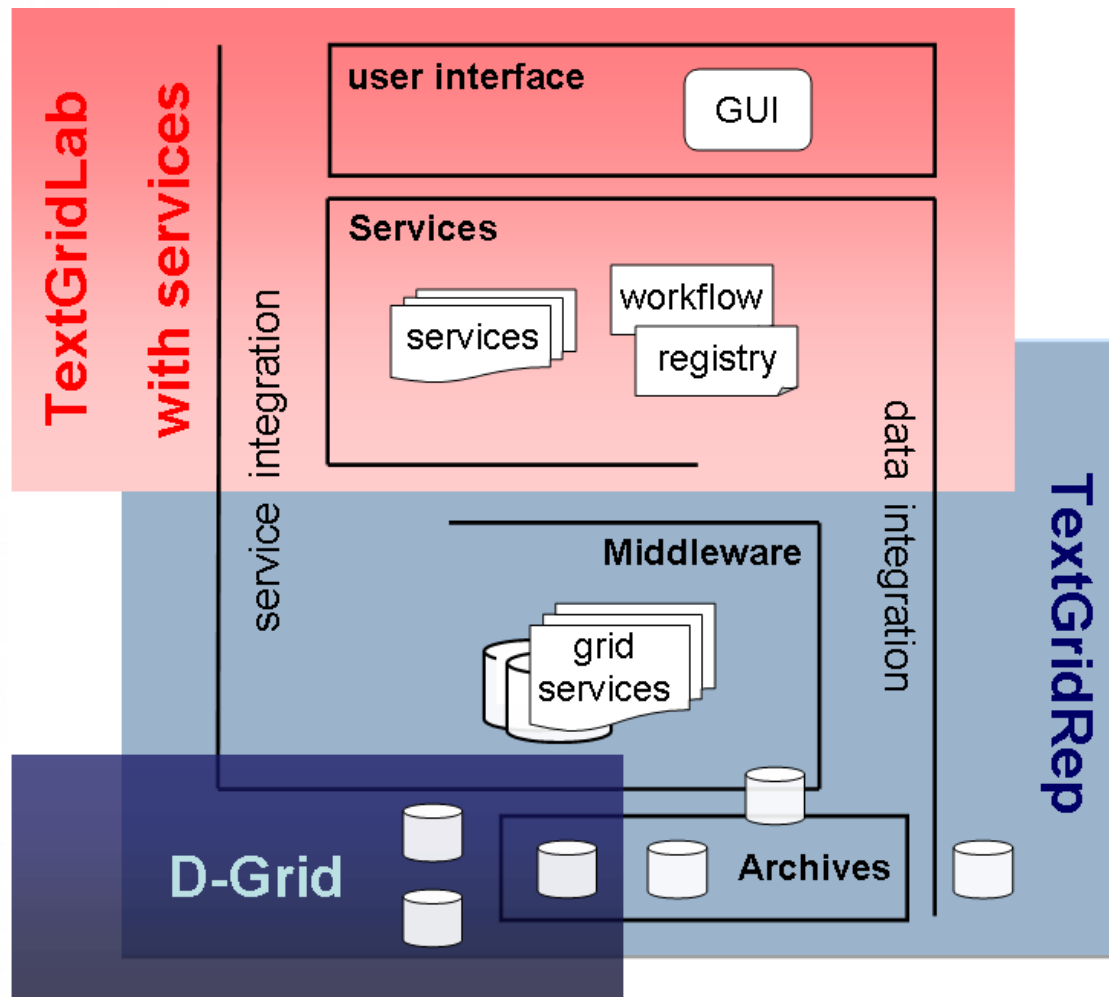
  
**saphor**

**DAASI**  
International

Directory Applications  
for Advanced Security  
and Information Management



# TextGrid: Architektur



- **TextGridRep**
  - **Storage Grid (GT4/GAT) und angebundene Archive**
  - **Middleware**
    - **TG-auth\***: Authentifizierung (Shibboleth) und Autorisierung (RBAC)
    - **TG-crud**: Operationen Create, Read, Update, Delete für Datenobjekte mit Registrierung in sämtlichen TG-Services
    - **TG-search**: Suche auf Original- und normalisiertem Format (TEI-XML, TextGrid Baseline Encoding)
    - **TG-Log**: zentrales Logging
- **Services: Tokenizer, Lemmatizer, Workflow-Engine, Streaming Editor, Sortierer, ...**
- **TextGridLab: GUI (Eclipse Rich Client), Frontend zu Services bzw. TextGridRep**



- Operationen auf TextGridObjekten, bestehend aus
  - Nutzdaten (Text, XML, TEI-XML, JPEG, ...)
  - Metadaten
    - Agent (author, contributor, translator, editor, ...)
    - Title, subtitle, date, language, page, volume, ISBN, ...
    - Administratives: project, URI, format
    - Relationen zu anderen Objekten (isVersionOf, hasSchema, ...)
    - Anwendungsspezifisches (custom element)
- Infrastruktur baut auf Web Services auf (SOAP und REST)
- Services akzeptieren URIs oder Nutzdaten direkt
- Services bekommen Session-Token und Log-Token
- CRUD akzeptiert Base64-kodierte Nutzdaten oder MTOM





## Anforderungen an Workflow

- **Batch Processing: viele TextGridObjekte verarbeiten**
- **Front-End im TextGridLab zu den Services**
- **TextGridObjekte lesen, verarbeiten, neu erstellen → regelgesteuerte Transformation der Metadaten für zu erstellende Objekte**
- **Workflow Editor soll den Endbenutzern des TextGridLab abnehmen:**
  - **Routinearbeit (z.B. kopieren des Session-Tokens, Gerüst zum Einbinden des CRUD, Base64-Kodierung von Konfigurationsdaten)**
  - **Komplexität von Workflow-Beschreibungssprachen**



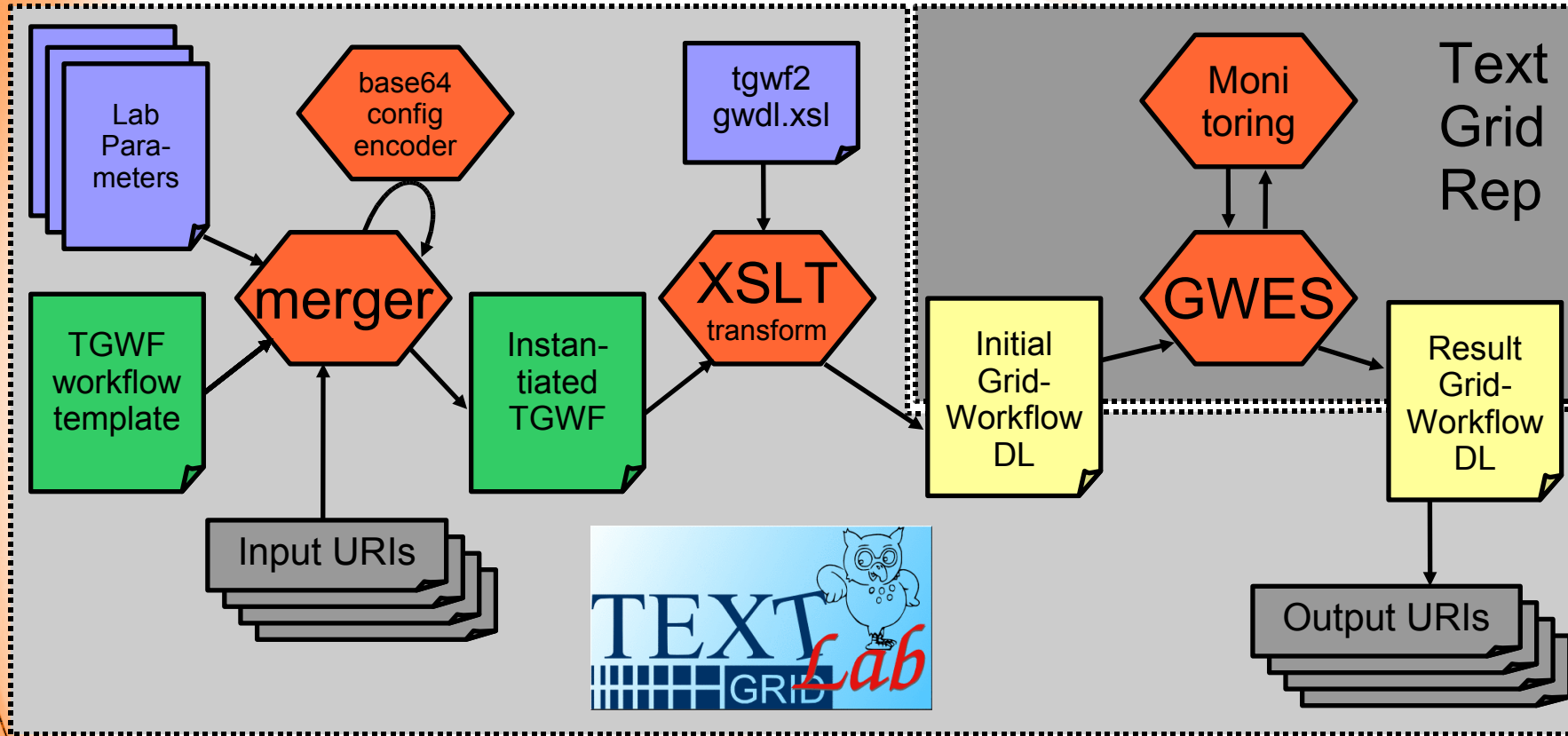
## Prototyp im Einsatz

- **Workflow Engine: Grid Workflow Execution Service (GWES) von Fraunhofer FIRST**
- **Erstellung einer auf die Bedürfnisse von TextGrid abgestimmten Workflow-Beschreibungssprache (TGWF)**
- **Bearbeitung der TGWF-Dokumente im TextGrid-XML-Editor**
- **Auswahl von Eingabeobjekten über die vorhandenen Mechanismen im TextGridLab**
- **Konfigurationsdaten der Services eingebettet in TGWF**
- **Metadaten transformation mit dem TextGrid-Streaming-Editor, XSLT-Stylesheet eingebettet in TGWF**
- **Auf Knopfdruck Erzeugung eines Workflow-Dokuments für den GWES mit Einfügung aller Variablen und XSL-Transformation in GridWorkflowDL**





## Ablaufschema



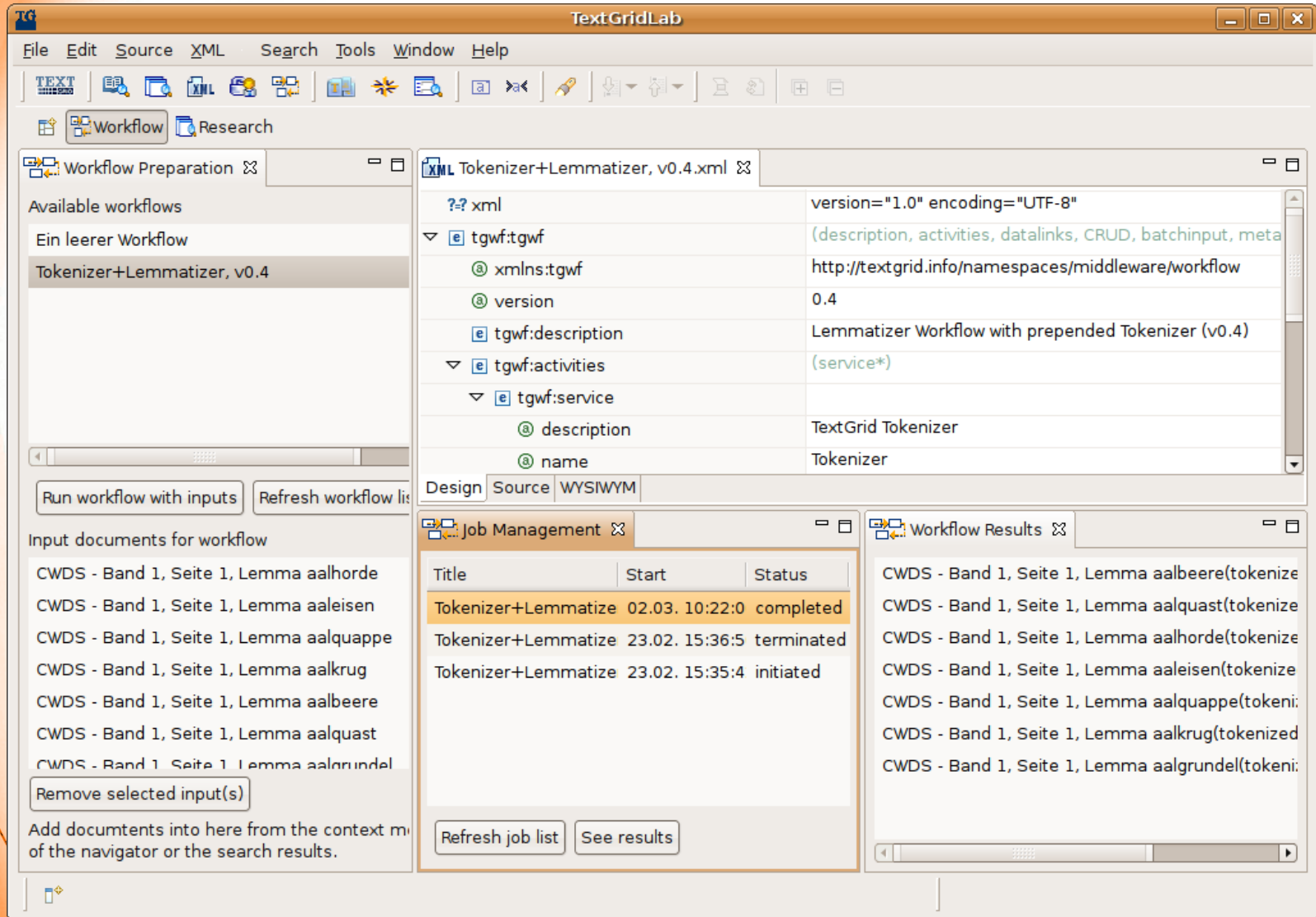
## TGWF (vereinfacht)

```

<tgwf>
  <activities>
    <service serviceID="tok" wsdl="..tok.wsdl"/>
    <service serviceID="lem" wsdl="..lem.wsdl"/>
  </activities>
  <datalinks>
    <link from="crudRead" to="tok"/>
    <link from="tok" to="lem"/>
    <link from="lem" to="crudCreate"/>
  </datalinks>
  <CRUD instance="" sessionID="" log=""/>
  <batchinput/>
  <metadatatransformation>
    <xsl:transform>...</xsl:transform>
  </metadatatransformation>
  <inputconstants>
    <activity serviceID="tok">
      <const name="config">...</const>
    </activity>
  </inputconstants>
</tgwf>

```





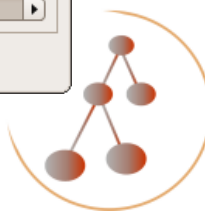
The screenshot displays the TextGridLab application window with the following components:

- Workflow Preparation:** Shows available workflows, including "Tokenizer+Lemmatizer, v0.4".
- XML Editor:** Displays the XML structure for the workflow, including elements like `xmlns:tgwf`, `version`, `tgwf:activities`, and `tgwf:service`.
- Job Management:** A table showing the status of workflow jobs.
- Workflow Results:** A list of processed documents and their corresponding lemmas.

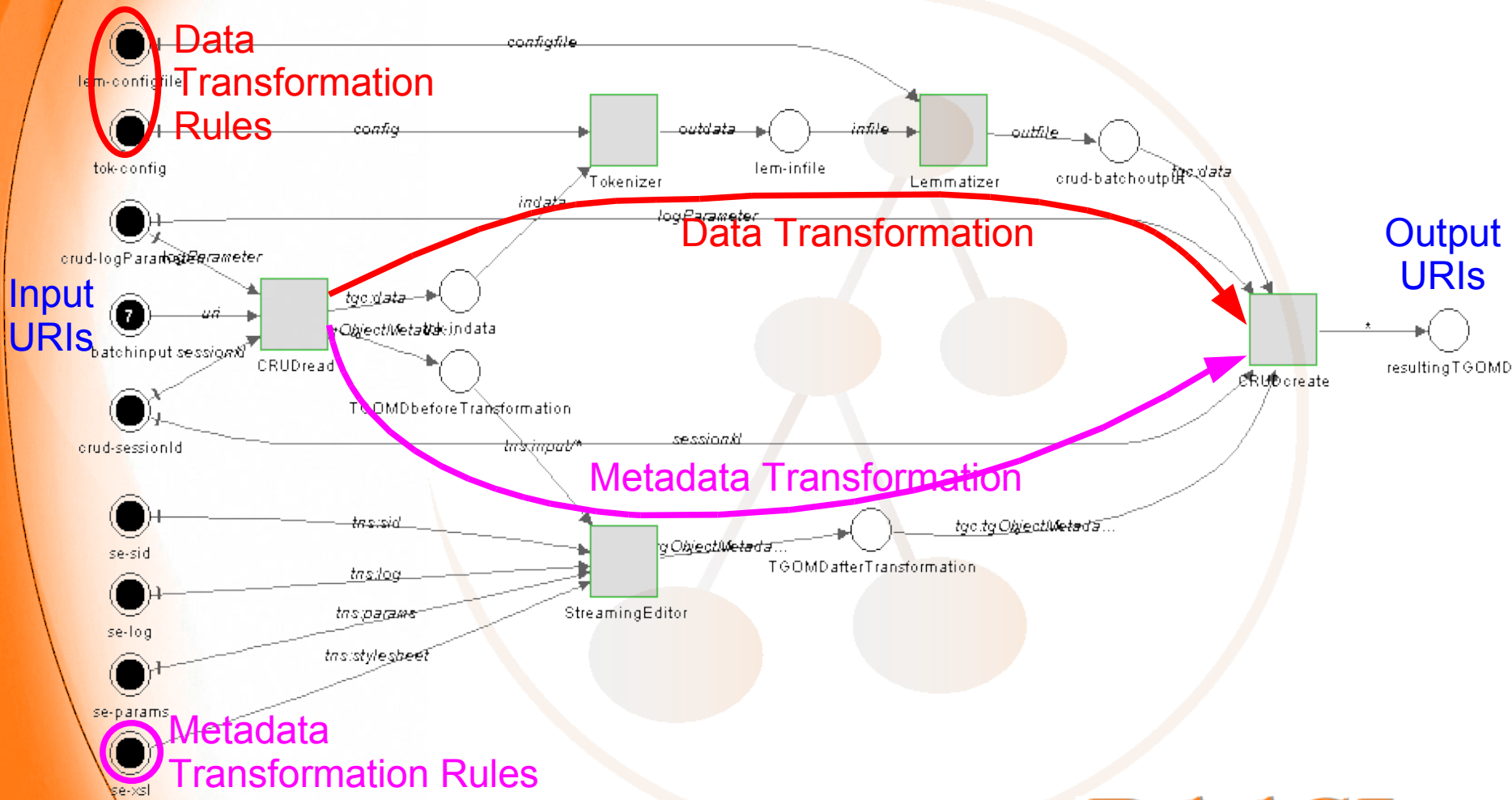
Title	Start	Status
Tokenizer+Lemmatize	02.03. 10:22:0	completed
Tokenizer+Lemmatize	23.02. 15:36:5	terminated
Tokenizer+Lemmatize	23.02. 15:35:4	initiated

Workflow Results:

- CWDS - Band 1, Seite 1, Lemma aalbeere(tokenize
- CWDS - Band 1, Seite 1, Lemma aalquast(tokenize
- CWDS - Band 1, Seite 1, Lemma aalhorde(tokenize
- CWDS - Band 1, Seite 1, Lemma aaleisen(tokenize
- CWDS - Band 1, Seite 1, Lemma aalquappe(tokeni:
- CWDS - Band 1, Seite 1, Lemma aalkrug(tokenized
- CWDS - Band 1, Seite 1, Lemma aalgrundel(tokeni:



## Petrinetz im GWES

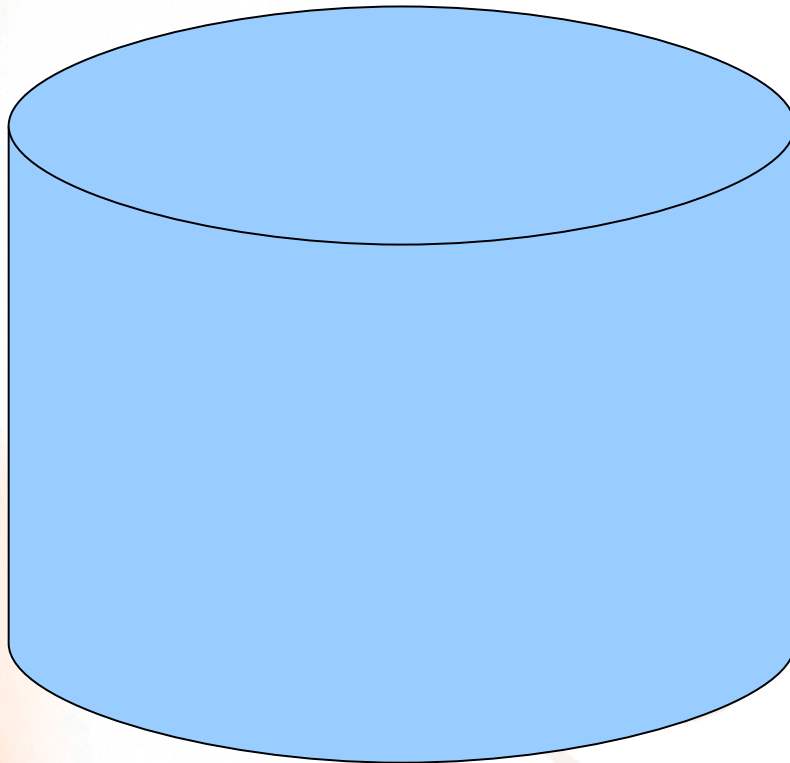


## Mehr Benutzerfreundlichkeit

- **Momentan: Verwendung des TextGridLab-eigenen XMLEditors für das TGWF-Dokument.**
- **Für unerfahrene Benutzer zu komplex:**
  - **Informationen über vorhandene Services sind evtl. nicht bekannt bzw. WSDLs nicht vertraut → Service Registry**
  - **Eingebettete Konfigurationsparameter der Services überfrachten das Dokument → spezieller Editor sowie vorgefertigte Konfigurationsbeispiele in Registry verlinkt**
  - **XSLT zur Metadatentransformation evtl. nicht vertraut → GUI für die wichtigsten Operationen mit rudimentärem Stylesheet-Generator (abh. v. Metadaten!)**
  - **Links unübersichtlich im XML → langfristig grafisches Drag&Drop für Services und Links**



## Service Registry



Config:Tokenizer: **config** URI1  
Tokenizer: config: URI2

Service : Tokenizer

WSDL: ...

Operation: tokenize

targetNamespace: ...

inputparams: **config** infile

outputparams: outfile

Service: Lemmatizer

...

...

...

...

...

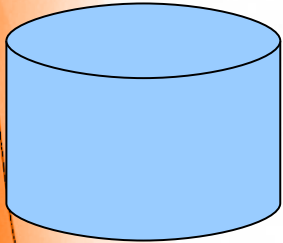




# Erweiterung: Service-Auswahl

## 1 Activities

Service Registry

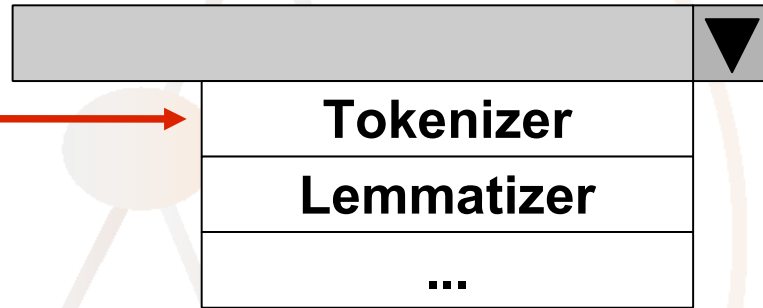


Service: **Tokenizer**  
WSDL:  
Operation: ...  
targetNamespace: ...  
inputparams: config infile  
outputparams: ...

Service: Lemmatizer

...  
...  
...  
...

Config: Tokenizer: config: URI1  
Config: Tokenizer: config: URI2



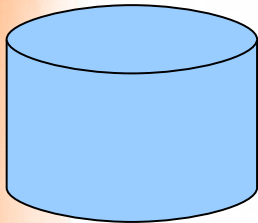
Action 1: Tokenizer  
Action 2: Lemmatizer  
Action 3: ...



# Erweiterung: Links einfügen

## 2 Data links

Service Registry



Service: Tokenizer  
WSDL: ...  
Operation: ...  
targetNamespace: ...  
inputparams: config: **infile**  
outputparams: ...

Service: Lemmatizer

Config: Tokenizer: config: URI1 ...  
Config: Tokenizer: config: URI2 ...  
...  
...

FROM		TO	
Service	Param.	Service	Param.
CRUD	Input	Tokenizer	<b>Infile</b>
Tokenizer	Outfile	Lemmatizer	Indata
Lemmatizer	Outdata	CRUD	Output

FROM  ▼  ▼

TO  ▼  ▼

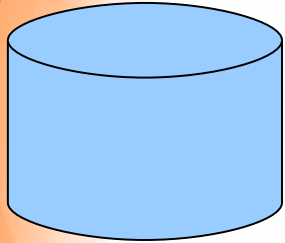


## Erweiterung: Konfigurationsauswahl

3

### Configs

#### Service Registry



Service : Tokenizer  
WSDL: ...  
Operation: ...  
targetNamespace: ...  
inputparams: config infile  
outputparams: ...

Service: Lemmatizer

...

Config: Tokenizer: config: URI1

...

Config: Tokenizer: config: **URI2**

Service ▼

Param. ▼

TextGridObject1
TextGridObject2

select

edit

new



## 4 Metadata transformation

add agent

author ▼	name	add
editor		
translator		
...		

transform title

\$title\$ (tokenisiert)	add
-------------------------	-----

(use \$title\$ for orig. title)

select project

project 1 ▼
project 2
project 3
...

edit stylesheet directly



## Ausblick

- **Stand: Prototyp, Community-Tests laufen**
- **Nutzer brauchen XML- und minimale XSLT-Kenntnisse**
- **Umsetzung für mehr Benutzerfreundlichkeit bis Projektende angestrebt – auch grafisches Drag&Drop?**
- **Weitere Services beispielhaft einbinden → Registry**



# Vielen Dank für Ihre Aufmerksamkeit!

- Fragen?
  
- **DAASI International GmbH**
  - [www.daasi.de](http://www.daasi.de)
  - [Info@daasi.de](mailto:Info@daasi.de)

