**CLARIN**

Common Language Resources and Technology Infrastructure

# CLARIN - a European Research Infrastructure

**Peter Wittenburg**

**Max-Planck Institut für**

**Psycholinguistik, Nijmegen**

J. Taylor

"eScience is about global collaboration in key areas of science and the next generation of infrastructures that will enable it"

Requires new persistent platforms

- to enable researchers to combine resources and tools to solve the big challenges of today (global migration, crisis of cultures and minds)
- to increase the efficiency of researchers in the many small tasks
    - 40 % of the time of  "knowledge workers" is spent, to find useful material (Forrester Research)

# CLARIN Goal

**What:**

- Offer a distributed Research Infrastructure of integrated and interoperable Language Resources and Tools that serves researchers and students in the SSH

**How:**

- allow the combination of existing and web-accessible digital centers hosting resources in a common federation
- offer language tools and services as distributed services with a common web interface
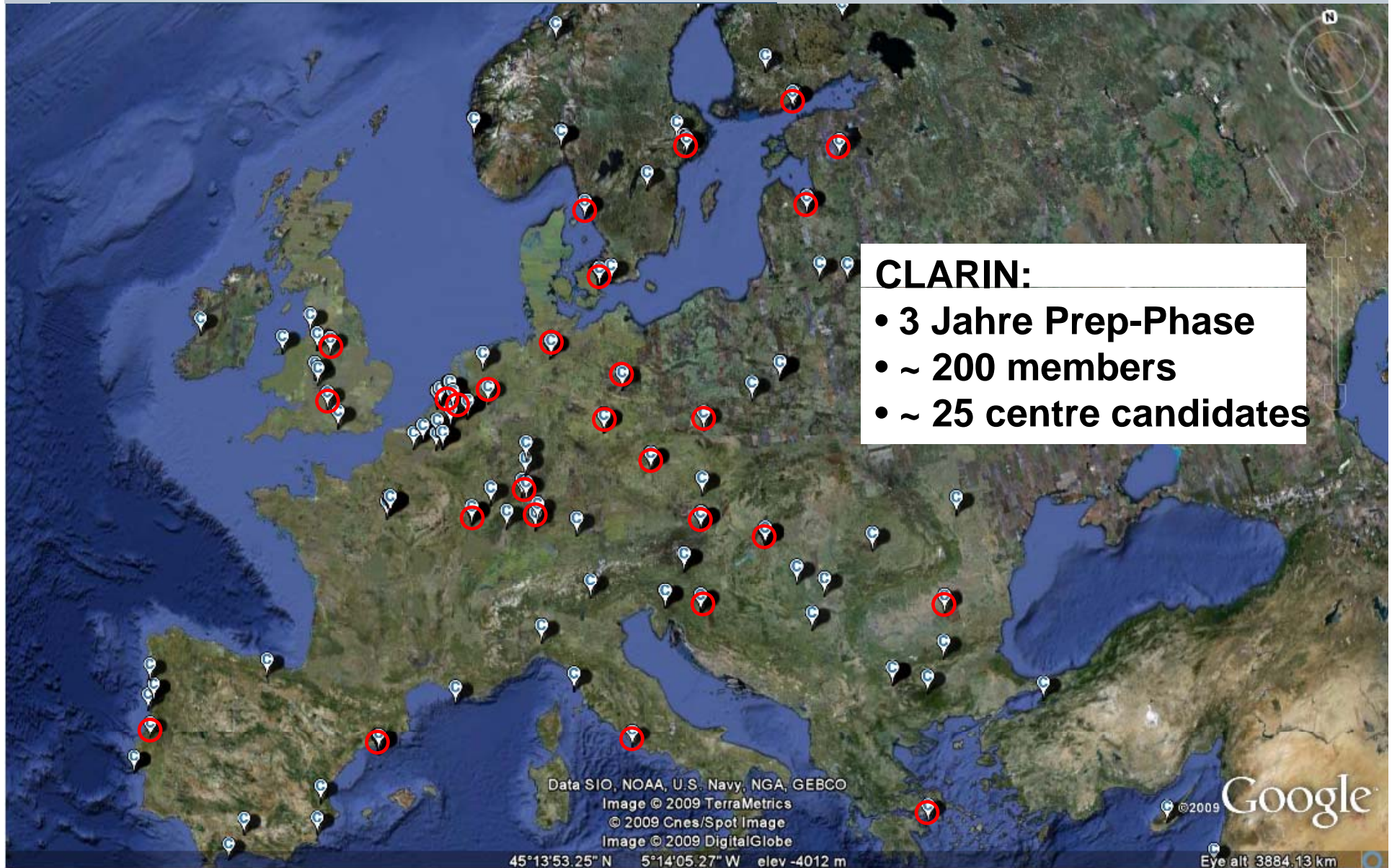
A researcher authenticates at his own organization and creates a *virtual collection* of resources from different repositories and executing a *virtual pipeline of processes* on them.
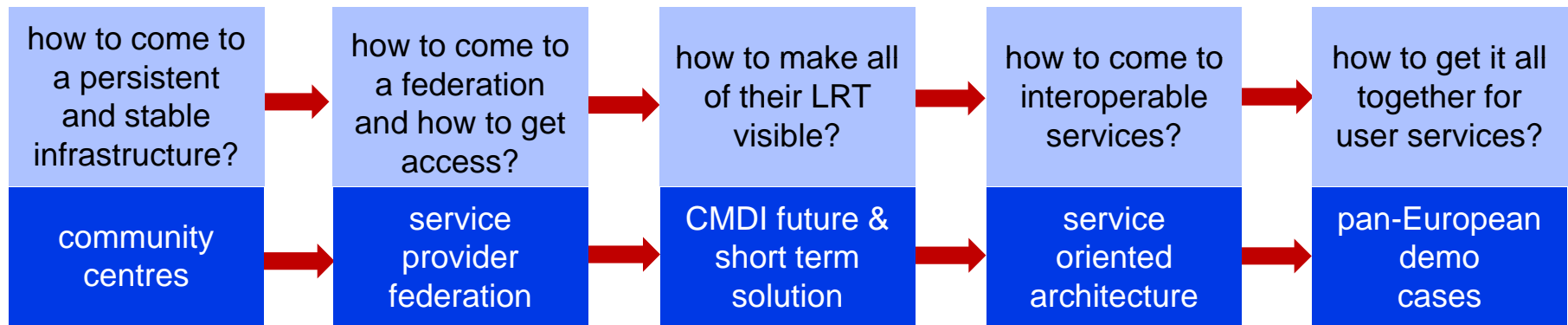


King Arthur failed by the way will CLARIN fail as well?

# CLARIN is pan-European



CLARIN:
- 3 Jahre Prep-Phase
- ~ 200 members
- ~ 25 centre candidates

... at least IT oriented aspects

| how to come to a persistent and stable infrastructure? | how to come to a federation and how to get access? | how to make all of their LRT visible? | how to come to interoperable services? | how to get it all together for user services? |
|---|---|---|---|---|
| community centres | service provider federation | CMDI future & short term solution | service oriented architecture | pan-European demo cases |

CLARIN has other very important aspects:
• Relation with SSH disciplines - mainly driven by national funds
• Education/Training, Help/Support/Advice, Dissemination
• Harmonization of licencing and Code of Conducts
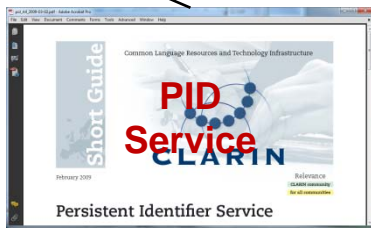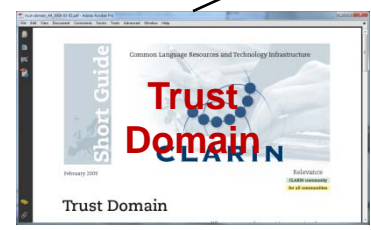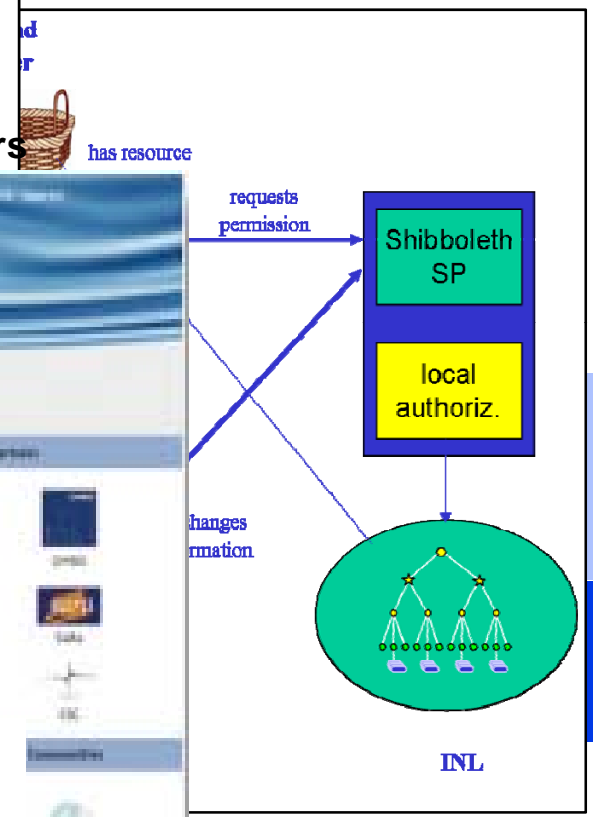• Specification of the ERIC legal framework to ensure persistency

# Community Centres

**25 Centre Candidates**

**all are busy with restructuring plans**

**2 already give long-term preservation service**

| how to come to a persistent and stable infrastructure? | how to come to a federation and how to get access? | how to make all of their LRT visible? | how to come to interoperable services? | how to get it all together for user services? |
|---|---|---|---|---|
| community centres | service provider federation | CMDI future & short term solution | service oriented architecture | pan-European demo cases |

**CLARIN Centres**

CLARIN Centres

**Centres Criteria**

Criteria for CLARIN Centres

**Long-term Preservation**

Long-term Preservation and Access

**REPLIX Replication**

REPLIX

# Service Provider Federation



- **Service Provider Federation**
  - **Agreement 1**
  - **n centers members**

http://www.pidconsortium.eu

Shibboleth SP

local authoriz.

has resource

requests permission

changes information

INL

Trust Domain

Initial Federation

PID Service

Trust Domain

Initial CLARIN
Service Provider Federation

Persistent Identifier Service

# Metadata Domain



about 270.000 resources/corpora included

ISOc...

IMDI Domain

- CGN (12.000)
- End.Lang. (35.000)
- MPI (33.000)
- BAS (7.400)
- AILLA (1.800)
- OLAC (40.000)
- LRT Inventory (800/137)
- DEKL...
- ELRA (60)
- others

OAI PMH transformation

GIS overlay

Facet...

Catalogue

hard problem:
- mapping
- granularity
- curation

**this is where the ILSP team played a central role**

Component Metadata — Component Metadata

Metadata now — Create CLARIN Metadata Now

Virtual Collection — Virtual Collections

ISOcat Registry — Concept Registry Service

VLO Observatory — Virtual Language Observatory

# Service Oriented Architecture

# Demo Cases (just started)

The Language Archive

| C4/WebLicht Corpus Case | EU Identity Index Case | Multimedia/multi modal Case | Folkstory Case |
|---|---|---|---|

| how to come to a persistent and stable infrastructure? | how to come to a federation and how to get access? | how to make all of their LRT visible? | how to come to interoperable services? | how to get it all together for user services? |
|---|---|---|---|---|
| community centres | service provider federation | CMDI future & short term solution | service oriented architecture | pan-European demo cases |

# not alone ...

# need to take care of data ...



**Trust**

**Data Curation**

**Data generators**

**Users**

User functionalities
Data capture & transfer
Virtual Research
Environments

**Community Support Services**

CLARIN, DARIAH etc

Data discovery & navigation
Workflow generation
Annotation,
Interpretability

**Common Data Services**

Daten e-Infrastructure

Safe & persistent storage
Identifiers, Authenticity,
Workflow execution,
Mining

**Architecture created by EC High Level Expert Group
will be a guideline for coming decades**

- live in a multilingual Europe with a joint historical tradition and need to exploit this strength
- many research questions are cross-national
- required standards cannot be national

- sharing costs in all respects is more efficient
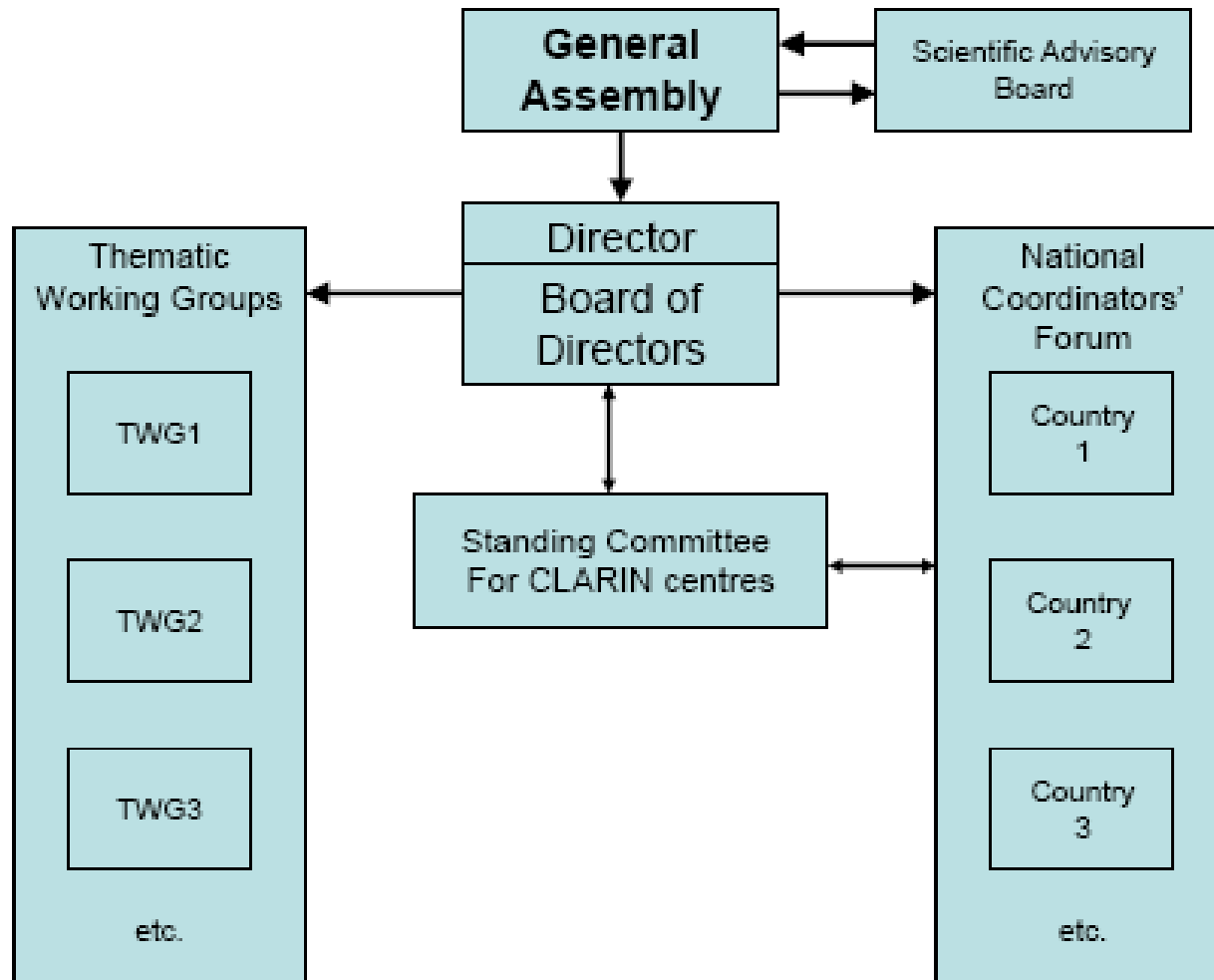- finally it's about global competition also in SSH

# Why now?

- there is the ESFRI process and all countries are synchronized which is a unique chance to build infrastructures

- in total 44 initiatives on the ESFRI roadmap and there is the potential of gain by an eco system of RI

- we need to organize our resource domain due to huge increase of data (MPI: 200 TB)

- we need to take care to not loose our cultural and scientific memory

- there is a huge uptake of RI and there will be many funding streams!!!

# who and when?

- current EU CLARIN consortium in prep phase (08-10): 32 partners from 24 countries

- CLARIN construction phase from 2011; main funds by national programs - but additional funding streams by EC connected to RI

- legal issue: foundation of a European Research Infrastructure Consortiums (ERIC) as basis for future with automatic qualification to participate in programs

# Organisation of the CLARIN ERIC



CLARIN
Utrecht

# who seems to be on board?

Belgium, Bulgaria, Germany, Denmark, Estonia, Latvia, Finland, Croatia, <u>Netherlands</u>, Norwegen, Austria, Portugal, Spain, Czech Republic, Hungary, South Tirol, ?

Some are discussing: FR, SW, GR?, etc.

# Advantage of membership

- privilaged access to CLARIN federation

- networked with CLARIN centres (direct technology transfer)

- a word when discussing priorities, agreements, best practices

- access to EC funding streams

- access to education and training programs to make our young generation competitive

# Weitere Informationen

- CLARIN web site: http://www.clarin.eu

- CLARIN office: clarin@clarin.eu

- CLARIN Newsletter:

  http://www.clarin.eu/newsletter

- CLARIN members:

  http://www.clarin.eu/members
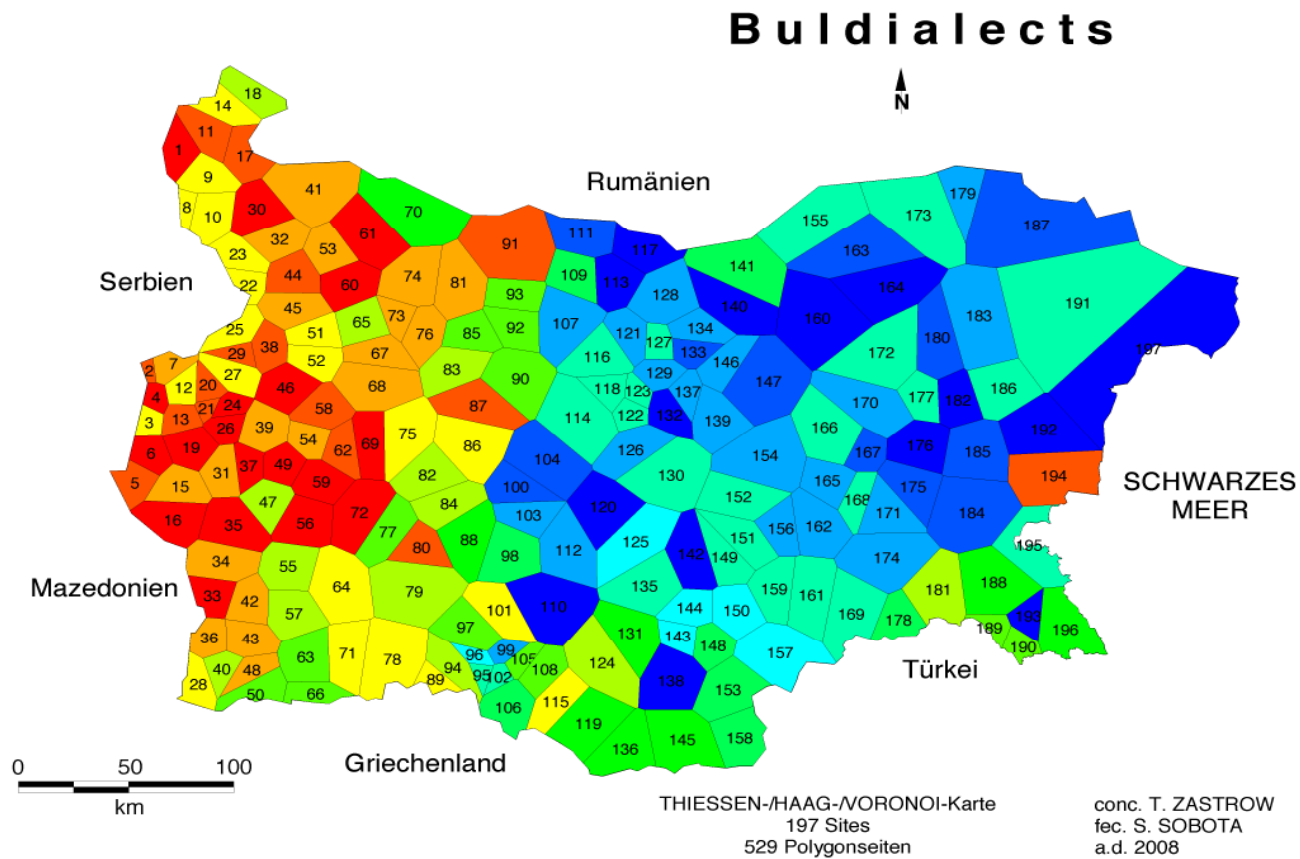
Thanks for your attention.

# CLARIN Usage Scenario

- Scenario: A Serbian and a German PhD student want to study language variation in the Balkan area

- Resource: via VLO they find all relevant language variation data for that area

- Tools/Services: Modern clustering methods available via the web allow to quickly build dialect continua on top of a geographic map; visualization services allow to pipeline this to get a nice output

## Visualization of Dialect Data: Clustering

# CLARIN Usage Scenario

- Scenario: Linguists, sociologists and ethnologists want to study the cultural and linguistic differences of parliament debates in SE, DE and GR about the swine flue and compare how such global problems are dealt with

- Resource: building a virtual collections of all debates (Audio, Video, Transkription)

- Tools/Services: allowing researchers to analyse and annotate gestures, intonation, word choices, timing etc where partly powerful computers need being used

- Vision: in 2011/12 such computational services will be made available in CLARIN 2011