# Technology in language documentation

Jacquelijn Ringersma

Max Planck Institute for Psycholinguistics

# Language documentation:

## Aim:

Maintain, consolidate or revitalize endangered languages

## Through:

Creation of a representative, multipurpose and long-lasting record of languages

## By:

Recording language events, speech and gesture in natural context etc.

Storing the resources in an organized, accessible and persistent archive

# In this presentation:

Why archiving?

What is an organized, accessible and persistent archive?

Which technology is required (and offered by the MPI)

      Metadata tool

      Archive upload and access management tools

      Browsing, searching and accessing the resources

      Enrichments with ELAN, LEXUS and ViCoS

## Misconceptions about archiving

1. Your stuff is buried here and gone forever



(from Andrew Garret
Berkeley Archives)

## Misconceptions about archiving

1. Your stuff is buried here and gone forever

2. Other linguists will take advantage of your hard work and take away your good ideas



(from Andrew Garret
Berkeley Archives)

## But the actual truth is that:

1. Other linguists do not really care about your work

2. The people who do care are the members of the speech communities – and they care about it in a different way than you do

"The coolest thing to do with your data will be thought of by someone else"

(attributed to Rufus Pollack)

# Is there a danger that we loose digital data?

YES,

UNESCO: 80% of our recordings is endangered

How much of your data and files on the notebook is organized, backed-up?

How long can media and formats be accessed?

# Is there a danger that we loose the data?

YES, a few messages

> Archive data into a trusted archive (long term preservation and accessibility)
> Create high quality metadata so that you can find the way back to the data
> Use open standards

**MPI-PL archive:**
Trusted archive,

> ORGANIZED information
> Continuous extension
> Collaboration and interaction
> Commentary and relation drawing (enrichment)
> Supporting centre for cross-corpus and language work

## Correct conceptions about archiving

1. It requires discipline

2. It creates a bit of techno noise

# Task of the 'archiving instance'

Organization of corpora or data following clear principles

Creation of a coherent and consistent archive
Store data in an accessible and persistent form (long-term)

Give access to data to different users, but protect data against unauthorized access

Adhere to code of conduct and adhere to ethical and legal issues
Provide tools to researchers

## Clear principles: Data organization and access infrastructure

**IMDI**    Metadata editor

Browser and search

## Coherent, consistent and persistent: Data management

**LAMUS** Checks the content of the files, and file type check

Assigns a persistent identifier to the uploaded file

Allows the creation of corpus structures

Web based, easy to use

## Safe access: Data access rights and protection

**AMS**    All metadata in the archive is open

All resource access can be controlled by AMS (web based)

Users remain the owners and stay in control of the access

Setting of licenses and code of conducts

LAT – Language archiving technology
www.mpi-lat.eu

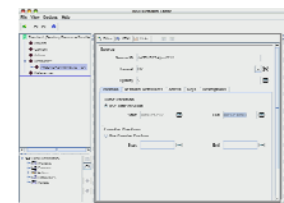# Workflow and tool requirements:

Recording



Capturing
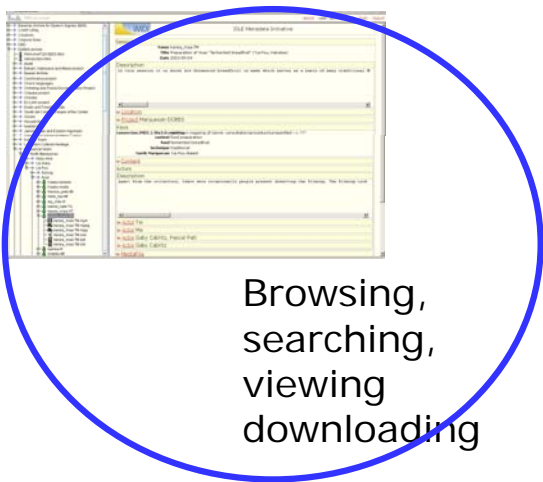(DV)



Transcoding
(MPEG1 & 2, WAV)



Enrichments



Describe
resources with
metadata



Uploading
metadata &
resources

Setting access
rights to the
resources

Browsing,
searching,
viewing
downloading

# Workflow and tool requirements:

Recording

Capturing
(DV)

Transcoding
(MPEG1 & 2, WAV)

Describe
resources with
metadata

Browsing,
searching,
viewing
downloading

Setting access
rights to the
resources

Uploading
metadata &
resources

# Describe resources with metadata

## IMDI metadata

Metadata is data about data

Structured and machine readable

Elements describe the content of the resource files

**IMDI set:**

General data: project, location

Content data: Genre, Interactivity, Modality, Language

Actor data: Age, gender, languages

Resource data: format, size

**IMDI editor – downloadable from the LAT page**

LAT – Language archiving technology

www.mpi-lat.eu

# Workflow and tool requirements:

Recording



Capturing
(DV)



Transcoding
(MPEG1 & 2, WAV)



Describe
resources with
metadata



Browsing,
searching,
viewing
downloading

Setting access
rights to the
resources

Uploading
metadata &
resources

# Upload resources and metadata files with LAMUS

## LAMUS

Checks the content of the files, and file type check

Assigns a persistent identifier to the uploaded file

Allows the creation of corpus structures

Web based, easy to use

LAT – Language archiving technology

www.mpi-lat.eu

# Workflow and tool requirements:

Recording



Capturing
(DV)



Transcoding
(MPEG1 & 2, WAV)



Describe
resources with
metadata



Browsing,
searching,
viewing
downloading

Setting access
rights to the
resources

Uploading
metadata &
resources

# Setting access rights: AMS

Access on resources: IPR, privacy, copyright etc.
(Metadata is always open!)



Default setting after
upload

# Workflow and tool requirements:

Recording

Capturing
(DV)

Transcoding
(MPEG1 & 2, WAV)

Describe
resources with
metadata

Browsing,
searching,
viewing
downloading

Setting access
rights to the
resources

Uploading
metadata &
resources

# Browsing the data: IMDI browser

# Viewing the data: ANNEX viewer

# Searching the metadata and content data:
# Search engines: IMDI and TROVA

# Different ways of accessing:

## Community portals



## Google Earth Layers

# Workflow and tool requirements:

Recording

Capturing
(DV)

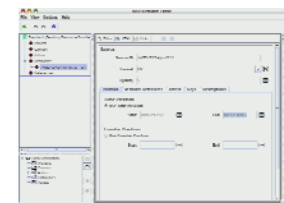Transcoding
(MPEG1 & 2, WAV)



Enrichments

Describe
resources with
metadata

Browsing,
searching,
viewing
downloading

Setting access
rights to the
resources

Uploading
metadata &
resources

- ELAN: annotating video and audio resources

LEXUS & ViCoS: Web based lexicon tool, multimedia encyclopedia

# Archiving instance,
# Max Planck Institute for Psycholinguistics

Archive managers: 3

Archive developers: 2

System manager: 1

Archiving software development: 4

Enrichment software development: 4

Archive for language data:

40 Terabyte of data

400.000 archived objects

# Training sessions:

We regularly organize training sessions on:

Audio and video handling

Archiving technology

Enrichment of data

We do welcome participants from other than DoBeS or MPI projects

Its 4-5 days in one week

Interested: please check our MPI website/Events section (www.mpi.nl/events/)

Contact:

jacquelijn.ringersma@mpi.nl

paul.trilsbeek@mpi.nl