# Access Management at
# Kielipankki – The Language Bank of Finland

Martin Matthiesen

DELAD Workshop Cork, 16.11.2017

*CSC – Suomalainen tutkimuksen, koulutuksen, kulttuurin ja julkishallinnon ICT-osaamiskeskus*

# Contents

- CSC?

- A use case example: HS.fi

- Registration of new resources

- Application and Approval Process

- Time for questions

# CSC – IT Center for Science

- Non-profit, (mostly) publicly funded company

- Owned by the Finnish ministry of education

- Supporting research & education institutions
  - "We don't do research, we support research."
  - Network (FUNET)
  - High-performance super computing
  - Data services

- Home to E-Infrastructures like Elixir, Kielipankki

# The Use Case HS.fi

- **Helsingin Sanomat newspaper corpus** 9/2011-9/2012

- 94 000 articles

- 600,000 sentences

- 8 Million words

- Licensed with the **CLARIN ACA +NC license**
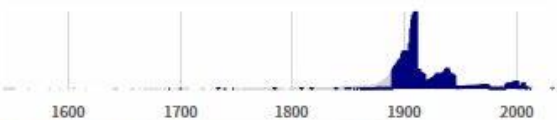
- Available at https://korp.csc.fi/

# Ingredients

We need: The data, a contract and

| lbr.csc.fi (LBR) | "Language Bank Rights" Access Management |
|---|---|
| korp.csc.fi (KORP) | The corpus viewer |
| metashare.csc.fi (METASHARE) | Metadata management |
| pid.csc.fi (PID) | PID management (URN/Handle) |
| www.kielipankki.fi (PORTAL) | License information |

# Prepare registration of corpus

- Create metadata in METASHARE

- Assign Persistent Identifier, like a URN/Handle (here: URN:HS.fi) in PID

- License text available in PORTAL

- Register corpus in LBR

# Dependencies

Create metadata in METASHARE

Register a URN in PID ("URN:HS.fi") to point to this page

stration  Community  Documentation  Statistics  Your Profile, Kielipankin

CSC

👁97 ✔15

# The HS.fi News and Comments Corpus

▷ View resource name in all available languages

*HS.fi*

http://urn.fi/urn:nbn:fi:lb-2014052717

ID: http://urn.fi/urn:nbn:fi:lb-2014052718    Here: URN:HS.fi

The HS.fi News and Comments Corpus contains the domestic news of the Helsingin Sanomat website and their comments from 5.9.2011 to 4.9.2012. The corpus starts with the first news of 5.9.2011 and ends with a news published in the morning on 3.9.2012 and the comments published on the website by 5.9.20 12.

The corpus has been published at https://korp.csc.fi.

Important: pseudonyms should be anonymized in publications referring to the corpus.

For detailed information on the license of the resource see https://kitwiki.csc.fi/twiki/bin/view/FinCLARIN/ClarinEulaEngACANc. Read Less

▷ View resource description in all available languages

« Back     Download     Edit Resource

text

| Distribution | Monolingual text corpus | Metadata |
|---|---|---|
| **Availability** | **Languages** | **Created:** 01/16/2014 |
| Available - Restricted Use | Finnish | **Last Updated:** 01/16/2014 |
| **Licence** | **Linguality** | **Metadata Language:** English, Finnish (en, fi) |
| *CLARIN ACA - NC* | **Linguality type:** Monolingual | **Metadata Creator** |
| **Restrictions:** Academic - Non Commercial Use, Attribution, No Redistribution, Other | **Size** | Imre Bartis |
| **Distribution Access/Medium:** Accessible Through Interface | 8,054,894 Tokens | |
| **Execution location:** *hidden* | 594,209 Sentences | |
| **Attribution Details:** HS.fi- | 93,602 Texts | |
| | **Modalities** | |

# Make the license text available in PORTAL using templates



**License type**

Choose the appropriate CLARIN license category. Note that RES implies: ID, PLAN, BY, NORED; ACA implies NORED. (to change edit the license template)

○ CLARIN PUB
◉ CLARIN ACA
○ CLARIN RES

**Resource name (EN) ***

The english name of the resource. Must be identical to the name used in Metashare and the Portal's corpus list.

The HS.fi News and Comments Corpus

**Kielivaran nimi (FI) ***

The finnish name of the resource. Again ensure it is the same as in Metashare and /aineistot.

HS.fi-uutiskommenttiaineisto

**URN of resource ***

Fill in the URN (without "http://urn.fi/", start with "urn:nbn:") to the Metashare article that describes all corpus variants covered by this license.

urn:nbn:fi:lb-2014052718

**Copyright holder ***

The copyright holder according to Metashare.

Helsingin Sanomat, Jutta Salminen

**Tags: ID & Access**

Select restrictions for "Identification and Access Conditions"

☑ ID: The user needs to be authenticated or identified.
☐ AFFIL=EDU: The user needs to be affiliated with a community of researchers through a university or research institution.
☐ AFFIL=META: The user needs to be affiliated with the general community of language research and technology researchers.

## The license in PORTAL



### CLARIN ACA end-user license +NC 1.0

Resource: The HS.fi News and Comments Corpus (URN: urn:nbn:fi:lb-2014052718)

Copyright holder: Helsingin Sanomat, Jutta Salminen

The Copyright holder grants the End-User a free, non-exclusive and perpetual (for the duration of the copyrig such, as modified, or as part of a compilation or derived work. The permission applies to all known or future mo Resource on other devices and in other formats. Distribution of copies is not allowed.

### Additional license terms as defined in the Terms of Service Agreement:

**Identification and Access Conditions**

- ID: The user needs to be authenticated or identified.

**General Use conditions**

- BY: Attribution, i.e. acknowledgement of authorship, is required.
- NC: The content is available for non-commercial purposes only.

**Distribution conditions**

- NORED: The user is not permitted to redistribute the resource.

This license has been made in compliance with copyright agreements by WIPO – the World Intellectual Prop intellectual property laws grant rights not mentioned in this license, they are also regarded as part of the rights different legal systems. Additional rights to the Resource may be agreed separately in writing.

# Language Bank Rights at Kielipankki

- Manages applications, access rights

- Instance of CSC's Resource Entitlement Management System REMS

- REMS was originally developed for Bioscience (Elixir)

- Used at Kielipankki since April 2015 as Language Bank Rights (LBR)

- Controls 47 resources (10/2015: 35 resources)

# Register the corpus in **LBR**



- Register the license(s).

- Register HS.fi as a "resource" (using "URN:HS.fi").

- Create forms to be filled in by applicants.

- Create a workflow that is followed during an application.

- Connect all of the above to a "Catalog item".

# Registration of corpus completed

- Created metadata in METASHARE

- Assigned URN/Handle (here: URN:HS.fi) in PID

- License text now available in PORTAL

- Registered corpus in LBR

# Language Bank Rights (LBR) Application and Approval

## Applicant's view

• Name is link to URN (to metadata)
• Applicant can add items to basket (webshop metaphor)

## The application

- Fill out form
- Accept licenses
- "Submit" sends application to approvers defined in workflow

**The approver view (1/4)**

• The application itself

## The approver view (2/4)

- Applied items (could be more)



Application #64

Application | Applied items | Applicant and members | History | Publications

| ID | Catalog item |
|----|--------------|
| 62 | The HS.fi News and Comments Corpus |

# The approver view (3/4)

- Applicant
- More info available behind name.

## The approver view (4/4)

- History of application
- Who what when
- Comments can be added

(Publications: Applicant can augment application later with relevant publications, skipped here)

# The applicant's view

- Status of application(s)
- Applicant can Close or Update

# Recap: Language Bank Rights at Kielipankki

- Based on Resource Entitlement Management System (REMS)
- Connects
    - resource users
    - approvers / licensors
    - licenses
    - resources
- Keeps electronic **paper trail**
- Heavily relies on **SAML2 Single Sign On** (eduGAIN, local federations, CLARIN IdP, Finnish Eduuni)

# Thank you!

- The Language Bank:

- https://www.kielipankki.fi/language-bank/

- Language Bank Rights:
  https://lbr.csc.fi/

- Demo:
  https://remsdemo.csc.fi/

- Time for comments, questions.