



Digital Language Extinction as a Challenge for the Multilingual Web

Georg Rehm

Network Manager META-NET
DFKI, Berlin, Germany

georg.rehm@dfki.de

Multilingual Web Workshop 2014: New Horizons for the Multilingual Web
Madrid, Spain – May 8, 2014



Co-funded by the 7th Framework Programme and the ICT Policy Support Programme of the European Commission through the contracts T4ME, CESAR, METANET4U, META-NORD (grant agreements no. 249119, 271022, 270893, 270899).

Digital Language Extinction



- Many smaller languages are experiencing problems digitally:
 - **Loss of function** – other languages take over entire functional areas such as, e.g., texting, email, search, e-commerce etc.
 - **Loss of prestige** – if it's not on the web, the language doesn't exist
 - **Loss of competence** – can you raise a digital native in your language?
- Andras Kornai's classification – corresponds to the amount of digital communication in that language:
 1. **digitally thriving languages** (comfort zone languages)
 2. vital languages
 3. heritage languages
 4. still/moribund/dead languages

potentially facing digital extinction ...
- Implications for the European/global multilingual web?

- ❑ Network of Excellence dedicated to fostering the technological foundations of the European multilingual information society.
- ❑ Projects: T4ME, CESAR, METANET4U, META-NORD.
- ❑ First funded phase ended on Jan. 31, 2013; new projects such as, e.g., QTLaunchPad and QTLeap are contributing.
- ❑ All EU member states and several non-member states covered.
- ❑ META-NET: **60** research centres in **34** European countries.



<http://www.meta-net.eu/members>

META-VISION: Building a community with a shared vision and strategic research agenda

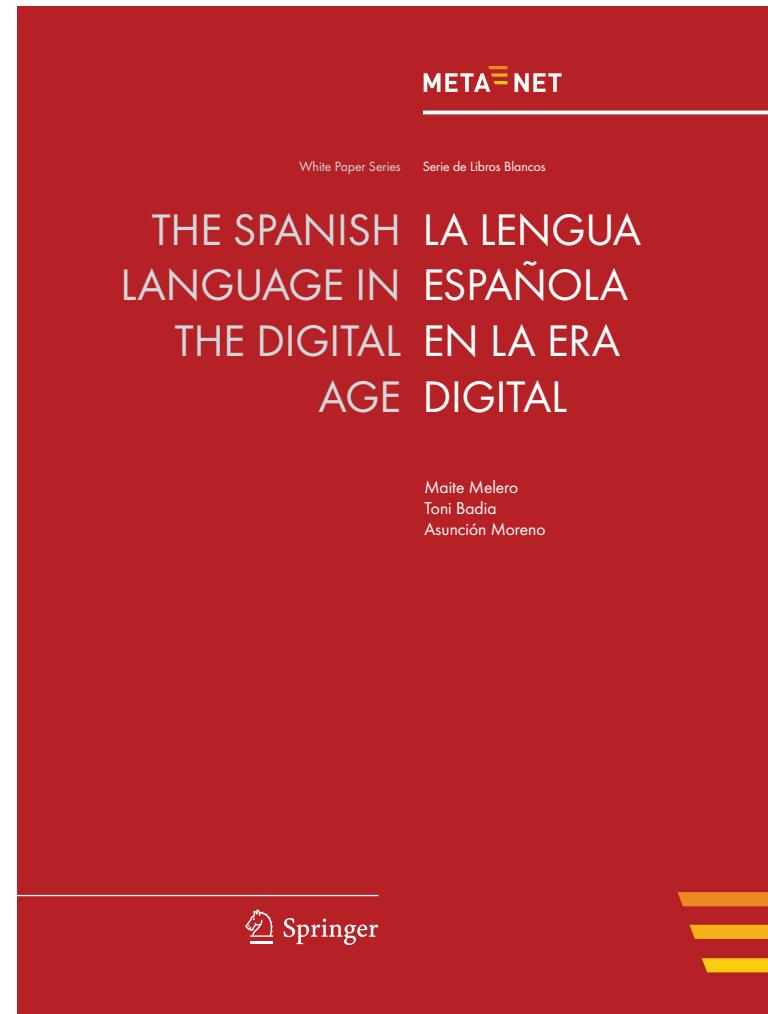
META-SHARE: Building an open resource exchange infrastructure

META-RESEARCH: Building bridges to neighbouring technology fields

Language White Paper Series



- “Europe’s Languages in the Digital Age”
- Series covers 31 languages in 31 volumes.
- Reports on the state of our languages in the digital age and the level of support through language technology.
- >2 years in the making.
- >215 experts as contributors.
- >8.000 copies distributed to politicians and journalists.



- Basque
- Bulgarian*
- Catalan
- Croatian*
- Czech*
- Danish*
- Dutch*
- English*
- Estonian*
- Finnish*
- French*
- Galician
- German*
- Greek*
- Hungarian*
- Icelandic
- Irish*
- Italian*
- Latvian*
- Lithuanian*
- Maltese*
- Norwegian
- Polish*
- Portuguese*
- Romanian*
- Serbian
- Slovak*
- Slovene*
- Spanish*
- Swedish*
- Welsh

* Official EU language

Cross-Lingual Comparison

META[≡]NET

- 1. Machine Translation
- 3. Speech Processing/Synthesis
- 2. Text Analytics
- 4. Language Resources
- Ranking: from *excellent LT support* to *weak/no support*.
- Cross-lingual comparison discussed and finalised at a network meeting with representatives of all languages (Oct., 2011).



	excellent	good	moderate	fragmentary	weak or no support through LT
MT		English	French, Spanish	Catalan, Dutch, German, Hungarian, Italian, Polish, Romanian	Basque, Bulgarian, Croatian, Czech, Danish, Estonian, Finnish, Galician, Greek, Icelandic, Irish, Latvian, Lithuanian, Maltese, Norwegian, Portuguese, Serbian, Slovak, Slovene, Swedish, Welsh
Text Analytics	excellent	good	moderate	fragmentary	weak or no support through LT
		English	Dutch, French, German, Italian, Spanish	Basque, Bulgarian, Catalan, Czech, Danish, Finnish, Galician, Greek, Hungarian, Norwegian, Polish, Portuguese, Romanian, Slovak, Slovene, Swedish	Croatian, Estonian, Icelandic, Irish, Latvian, Lithuanian, Maltese, Serbian, Welsh
Speech	excellent	good	moderate	fragmentary	weak or no support through LT
		English	Czech, Dutch, Finnish, French, German, Italian, Portuguese, Spanish	Basque, Bulgarian, Catalan, Danish, Estonian, Galician, Greek, Hungarian, Irish, Norwegian, Polish, Serbian, Slovak, Slovene, Swedish	Croatian, Icelandic, Latvian, Lithuanian, Maltese, Romanian, Welsh
Resources	excellent	good	moderate	fragmentary	weak or no support through LT
		English	Czech, Dutch, French, German, Hungarian, Italian, Polish, Spanish, Swedish	Basque, Bulgarian, Catalan, Croatian, Danish, Estonian, Finnish, Galician, Greek, Norwegian, Portuguese, Romanian, Serbian, Slovak, Slovene	Icelandic, Irish, Latvian, Lithuanian, Maltese, Welsh

Observations and Results

- ❑ When it comes to technology support, there are *massive differences* between Europe's languages and technology areas.
- ❑ Support for English is ahead of any other language.
- ❑ But: even support for English is *far* from being perfect.
- ❑ Several languages get the weakest score in *all four areas* (e.g., Icelandic, Latvian, Lithuanian, Maltese)!



Digital Language Extinction!



- ❑ “At Least 21 European Languages in Danger of Digital Extinction!”
- ❑ Press release on European Day of Languages (Sept. 26, 2012).
- ❑ Huge global interest in the topic and our key findings!
- ❑ 600+ mentions in the press.
- ❑ News from 40+ countries in 35+ different languages.
- ❑ 20+ television reports and 30+ broadcast interviews (radio, tv) with META-NET representatives.
- ❑ Two Parliamentary Questions in the EP on the “digital extinction of languages” topic.

Update of the Study (2014)

- ❑ Study comprised 31 volumes/languages.
- ❑ Many languages missing! Need for extension – at least of the comparison.
- ❑ We invited three language community bodies to participate in the update:

European Federation of National Institutions for Language (EFNIL)

Network to Promote Linguistic Diversity (NPLD)

Experts Committee of the European Language Charter (Council of Europe)



<http://www.meta-net.eu>



An Update and Extension of the META-NET Study “Europe’s Languages in the Digital Age”

Georg Rehm¹, Hans Uszkoreit¹, Ido Dogan², Vartkes Goetcherian³,
Mehmet Ugur Dogan⁴, Coskun Mermer⁴, Tamás Varadi⁵, Sabine Kirchmeier-Andersen⁶,
Gerhard Stöckl⁷, Meirion Prys Jones⁸, Stefan Oeter⁹, Sigríð Gramstad¹⁰

META-NET
DFKI GmbH
Berlin, Germany¹

META-NET
Bar-Ilan University
Tel Aviv, Israel²

META-NET
Hungarian Academy of Sciences
Budapest, Hungary³

META-NET
EFLNIL, METANET
Hungarian Academy of Sciences
Budapest, Hungary⁴

META-NET
Council of Europe, Com. of Experts
Hamburg, Germany⁵

META-NET
EFNIL
Dutch Language Council
Copenhagen, Denmark⁶

META-NET
Council of Europe, Com. of Experts
Bergen, Norway^{7,10}

META-NET
EFNIL
Institut für Deutsche Sprache
Mannheim, Germany⁸

META-NET
Tubitak Bilgeş
Gebze, Turkey⁹

NPLD
Network to Promote Ling. Diversity
Cardiff, Wales¹⁰

Abstract

This paper extends and updates the cross-language comparison of LT support for 30 European languages as published in the META-NET Language White Paper Series. The updated comparison confirms the original results and paints an alarming picture: it demonstrates that there are even more dramatic differences in LT support between the European languages.

Keywords: LR National/International Projects, Infrastructure/Policy Issues, Multilinguality, Machine Translation

1. Introduction and Overview

The multilingual setup of our European society imposes societal challenges on political, economic and social integration and inclusion, especially in the creation of the single digital market and unified information space targeted by the Digital Agenda (EC, 2010). Language technology is the missing piece of the puzzle, it is the key enabler and solution to boosting growth and strengthening Europe’s competitiveness.

Recognising Europe’s exceptional demand and opportunities, 60 leading research centres in 34 European countries joined forces in META-NET, a Network of Excellence dedicated to the technological foundations of a multilingual European information society. META-NET was partially supported through four projects funded by the EC: T4IME, CESAR, METANETEU and META-NORD. META-NET is forming the Multilingual Europe Technology Alliance (META) with more than 760 organisations and experts representing multiple stakeholders and signed collaboration agreements with more than 40 other projects and initiatives. META-NET’s goal is monolingual, crosslingual and multilingual technology support for all European languages (Rehm and Uszkoreit, 2013). We recommend focusing on three priority research themes connected to application scenarios that will provide European R&D with the ability to compete with other markets and achieve benefits for European society and citizens as well as opportunities for our economy and future growth.

This paper extends and updates one important result of the work carried out within the META-VISION pillar of the initiative, the cross-language comparison of LT support for 30 European languages as published in the META-NET Language White Paper Series (Rehm and Uszkoreit, 2012).

2. The Language White Paper Series

Answering the question on the current state of a whole R&D field is difficult and complex. For LT nobody had collected these indicators and provided comparable reports for a substantial number of European languages yet. To arrive at a first rough overview, the META-NET prepared the Language White Paper Series “Europe’s Languages in the Digital Age” (Rehm and Uszkoreit, 2012) that describes the current state of LT support for 30 European languages (including all 24 official EU languages). This undertaking had been in preparation with more than 200 experts since mid 2010 and was published in the summer of 2012. The study included a comparison of the support all languages receive in four areas: MT, speech, text analytics, language resources. The differences in technology support between the various languages and areas are dramatic and alarming. In the four areas, English is ahead of the other languages but even support for English is far from being perfect. While there are good quality software and resources available for a few larger languages and application areas, others, usually smaller languages, have substantial gaps. Many languages lack basic technologies for text

MT

Text Analytics

Speech

Resources

excellent	good	moderate	fragmentary	weak or no support
	English	French, Spanish	Catalan, Dutch, German, Hungarian, Italian, Polish, Romanian	Albanian, Asturian, Basque, Bosnian, Breton, Bulgarian, Croatian, Czech, Danish, Estonian, Finnish, Frisian, Friulian, Galician, Greek, Hebrew, Icelandic, Irish, Latvian, Limburgish, Lithuanian, Luxembourgish, Macedonian, Maltese, Norwegian, Occitan, Portuguese, Romany, Scots, Serbian, Slovak, Slovene, Swedish, Turkish, Vlax Romani, Welsh, Yiddish
excellent	good	moderate	fragmentary	weak or no support
	English	Dutch, French, German, Hebrew , Italian, Spanish	Basque, Bulgarian, Catalan, Czech, Danish, Finnish, Galician, Greek, Hungarian, Norwegian, Polish, Portuguese, Romanian, Slovak, Slovene, Swedish	Albanian, Asturian, Bosnian, Breton, Croatian, Estonian, Frisian, Friulian, Icelandic, Irish, Latvian, Limburgish, Lithuanian, Luxembourgish, Macedonian, Maltese, Occitan, Romany, Scots, Serbian, Turkish, Vlax Romani, Welsh, Yiddish
excellent	good	moderate	fragmentary	weak or no support
	English	Czech, Dutch, Finnish, French, German, Italian, Portuguese, Spanish	Basque, Bulgarian, Catalan, Danish, Estonian, Galician, Greek, Hungarian, Irish, Norwegian, Polish, Serbian, Slovak, Slovene, Swedish, Turkish	Albanian, Asturian, Bosnian, Breton, Croatian, Frisian, Friulian, Hebrew, Icelandic, Latvian, Limburgish, Lithuanian, Luxembourgish, Macedonian, Maltese, Occitan, Romanian, Romany, Scots, Vlax Romani, Welsh, Yiddish
excellent	good	moderate	fragmentary	weak/no support
	English	Czech, Dutch, French, German, Hungarian, Italian, Polish, Spanish, Swedish	Basque, Bulgarian, Catalan, Croatian, Danish, Estonian, Finnish, Galician, Greek, Hebrew , Norwegian, Portuguese, Romanian, Serbian, Slovak, Slovene	Albanian, Asturian, Bosnian, Breton, Frisian, Friulian, Icelandic, Irish, Latvian, Limburgish, Lithuanian, Luxembourgish, Macedonian, Maltese, Occitan, Romany, Scots, Turkish, Vlax Romani, Welsh, Yiddish

Strategic Research Agenda

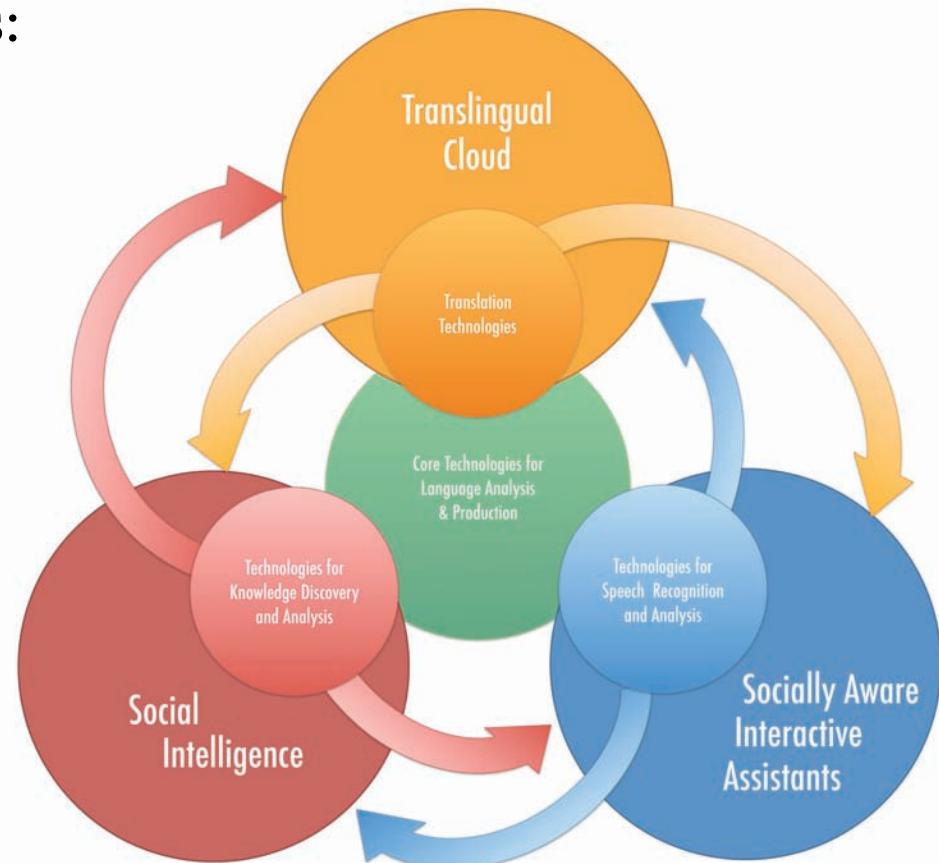


- ❑ Addresses the problems we identified when preparing the white papers.
- ❑ Can put Europe ahead of its competitors in this technology area.
- ❑ 200 contributors; >2 years.
54% industry; 46% research;
4% (inter)national institutions.
- ❑ Presented and discussed at 90+ conferences and major workshops.
- ❑ Published & presented in early 2013.
- ❑ <http://www.meta-net.eu/sra>



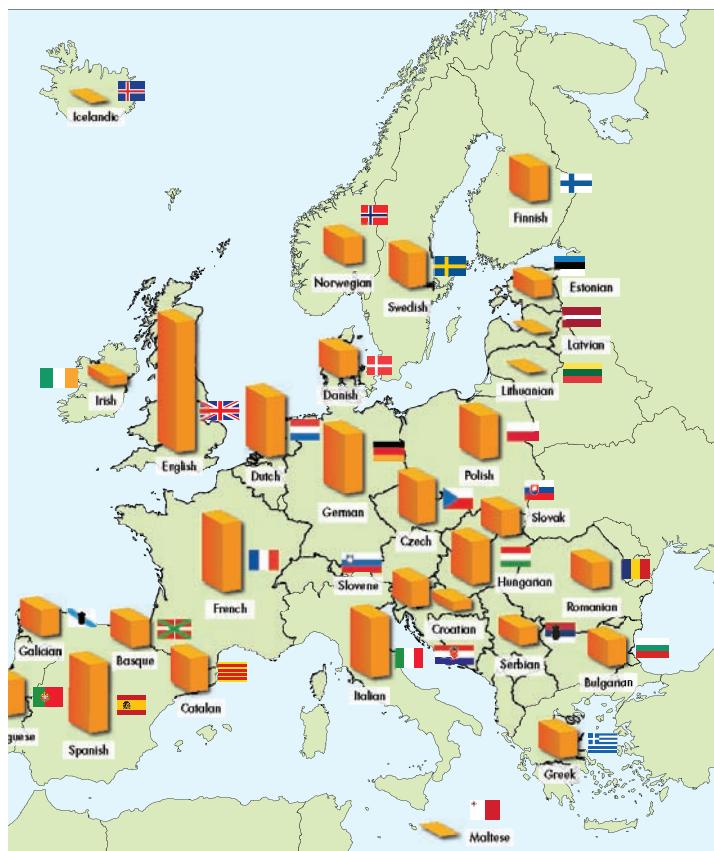
Priority Research Themes

- Three priority research themes:
 - Translingual Cloud
 - Social Intelligence and e-Participation
 - Socially-Aware Interactive Assistants
- Two additional themes:
 - European Service Platform for Language Technologies
 - Core Technologies for Language Analysis and Production



Resources, Tools, Technologies

METANET



2013



2020

Summary and Conclusions



- ❑ Many languages suffer from insufficient technology support.
- ❑ These languages are threatened by digital extinction.
- ❑ Multilingual technologies are a key enabler to overcoming language barriers and a good means to preserving our languages.
- ❑ Policy implications! We need different support for Kornai's classes:
 - Only **digitally thriving languages** can take care of themselves!
 - **Vital, heritage, still/moribund languages** need our help!
 - Focus on fostering R&D for smaller, less-resourced languages as well as research and technology transfer between languages.
- ❑ We suggest setting up an interdisciplinary research effort so that Europe can overcome language barriers, benefit from language diversity and fight digital language extinction (see META-NET SRA).

Current Funding Opportunities



- ❑ Horizon 2020-ICT-17 (“Cracking the language barrier”): Focus on languages with *fragmentary* or *weak/no* support through LT.
 - Official EU Member State languages only!  
 - Funding is limited: only 15M€ (out of 658M€ for ICT in 2015/2016).  
- ❑ CEF: MT is considered an obligatory component, core building block, key component of Europe’s future digital infrastructure!
 - Budget of only 4M€ in CEF-WP-2014 for procurement.  
 - Prep phase only! MT not considered ripe for production use just yet.  

ICT-17c Proposal: CRACKER

METANET

- CRACKER – Cracking the Language Barrier: Coordination, Evaluation and Resources for European MT Research

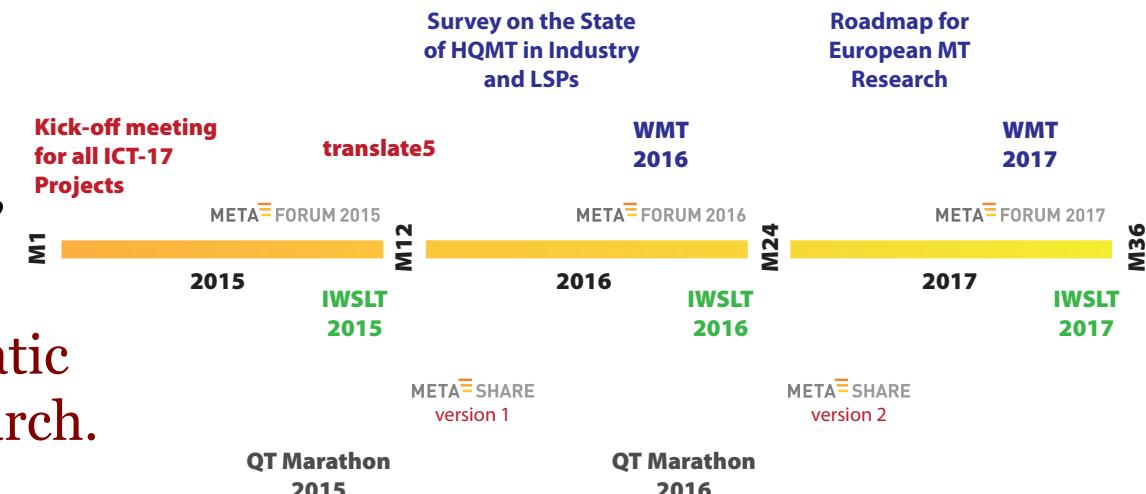
- CSA (budget: 1M€),
2015 to 2017.

- DFKI, CUNI, ELDA, FBK,
ATH, UEDIN, USFD

- Pushes towards a systematic improvement of MT research.

- Nucleus of strategy towards HQ-MT is the future ICT-17 group of projects.

- Activities: evaluation campaigns (WMT, IWSLT), META-FORUM, META-SHARE, training (MT Marathons), coordination.



¡Muchas gracias!

- ❑ Digital language extinction is a serious threat for many languages.
- ❑ The Multilingual Web community can help by providing resources, tools, technologies for these languages.
- ❑ Together we can fight digital language extinction!

<http://www.meta-net.eu>

<http://www.facebook.com/META.Alliance>

Acknowledgements: This work would not have been possible without the dedication and commitment of our colleagues Aljoscha Burchardt, Kathrin Eichler, Tina Klüwer, Arle Lommel, Felix Sasaki and Hans Uszkoreit (all DFKI), the 60 member organisations of the META-NET network of excellence, the ca. 70 members of the Vision Groups, the ca. 30 members of the META Technology Council, the more than 200 authors of and contributors to the META-NET Language White Paper Series and the ca. 200 representatives from industry and research who contributed to the META-NET Strategic Research Agenda.

