



Language Technology for Multilingual Europe

Georg Rehm

Network Manager META-NET
DFKI (German Research Center for Artificial Intelligence), Berlin, Germany

georg.rehm@dfki.de

Human Language Technology Day – Belgrade, Serbia
October 29, 2012

The CESAR logo features the word "CESAR" in a large, light grey sans-serif font. The letters are partially transparent, allowing a colorful circular graphic to be seen behind them. This graphic consists of several colored segments: blue, green, red, and orange, arranged in a circular pattern that suggests a sun or a gear.

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).

Outline



- Introduction
- META-SHARE
- Language White Paper Series
- Strategic Research Agenda
- Recent Developments – Next Steps

- **Challenge:** Providing each language community with the most advanced technologies for communication and information so that maintaining their mother tongue does not turn into a disadvantage.
- While research has made considerable progress in recent years, the pace of progress is not fast enough to meet the challenge within the next 10-20 years.
- All stakeholders – researchers, LT user and provider industries, language communities, funding programmes, policy makers – should **team up for a major dedicated push.**



Objectives



META-NET is a network of excellence dedicated to fostering the technological foundations of the European multilingual information society.

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

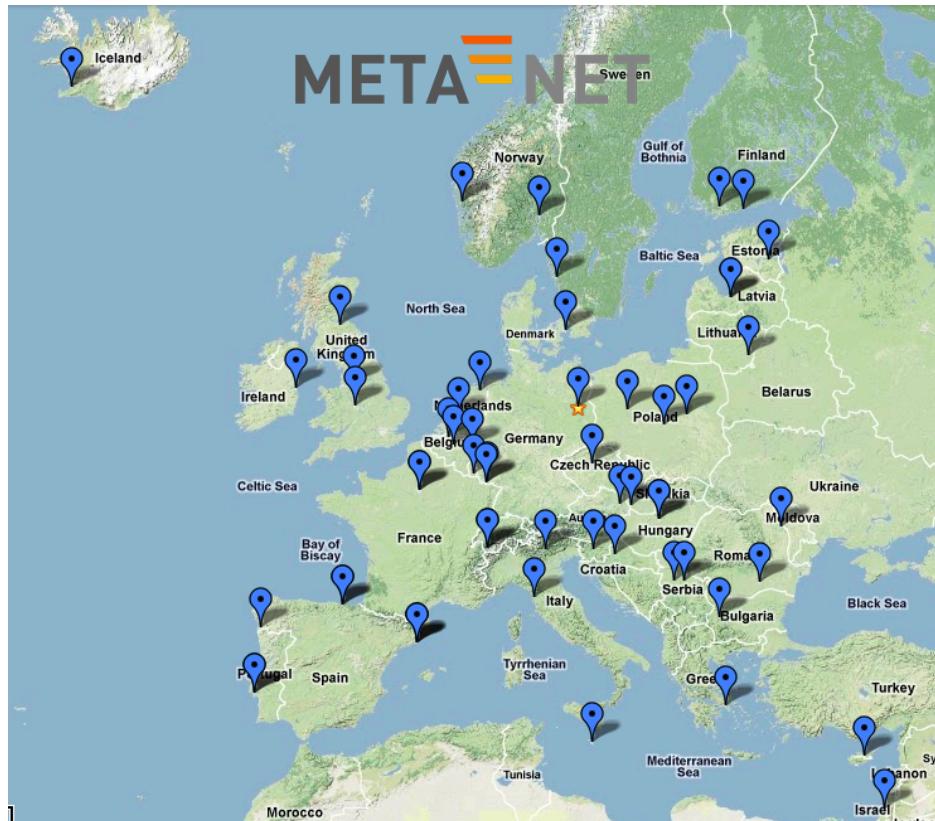
META-SHARE: Building an open resource exchange infrastructure

META-RESEARCH: Building bridges to neighbouring technology fields

Four EU-Funded Projects



- ❑ Initial project: T4ME (FP7;
13 partners, 10 countries)
- ❑ Three ICT-PSP consortia
since Feb. 2011: CESAR,
METANET4U, META-NORD
- ❑ All EU member states and
several non-member states
covered.
- ❑ META-NET in Oct. 2012:
60 members in **34** countries.



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



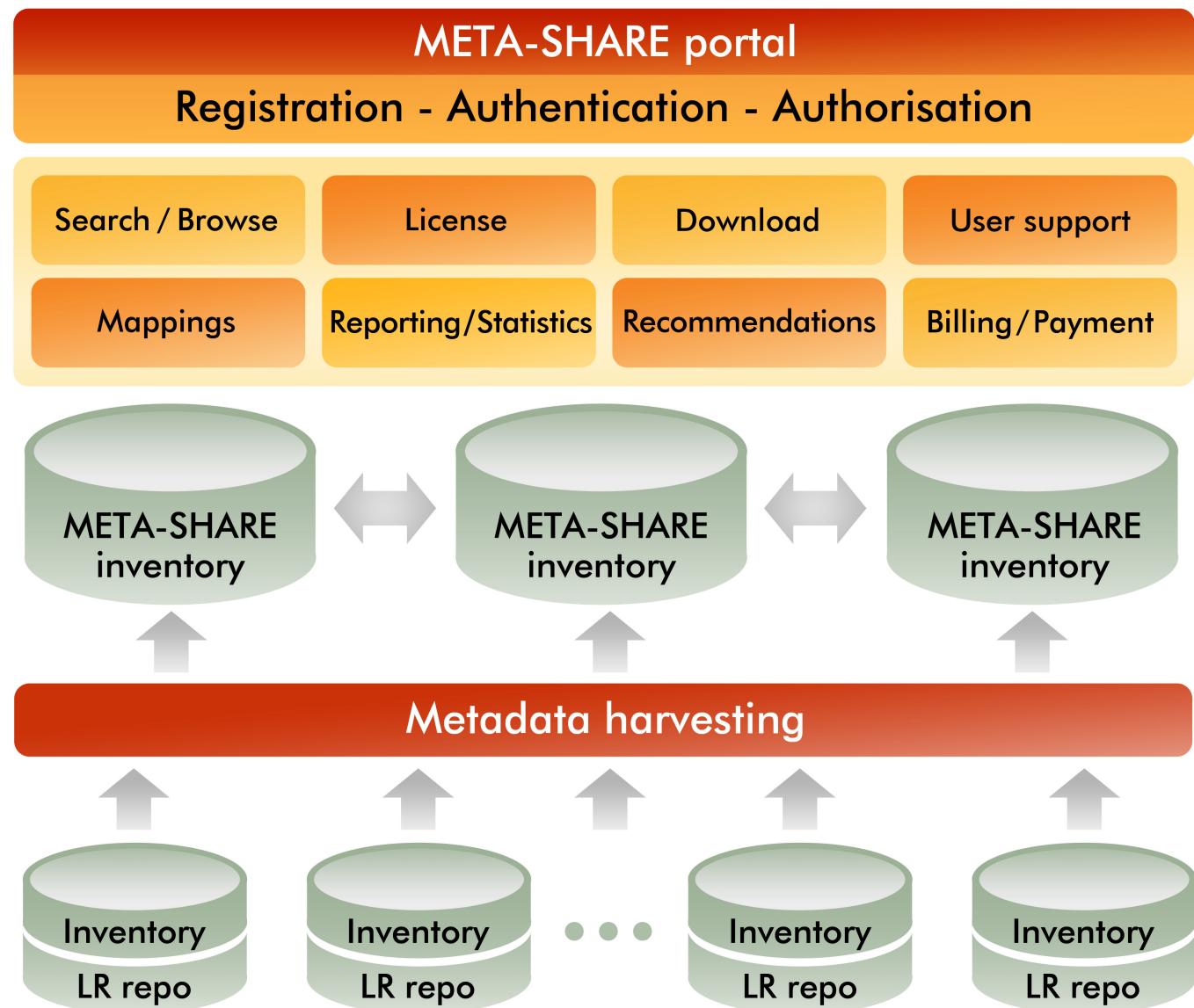
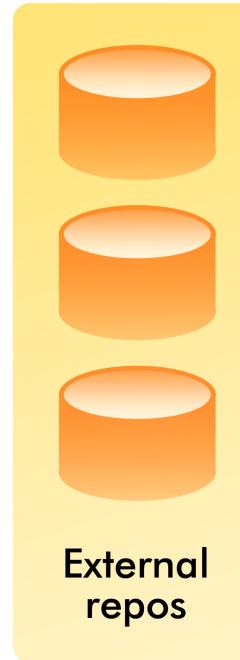
META-NET

META-SHARE

META-SHARE at a Glance



- ❑ Open exchange infrastructure for language resources and tools.
- ❑ Language resources and tools are documented, uploaded, stored in repositories, catalogued, can be downloaded, shared, discussed.
- ❑ Improve their visibility, documentation, identification, availability, preservation, interoperability.
- ❑ Long-term goal: boost research, technology and innovation through wide availability, pooling, openness and sharing of resources.
- ❑ Repositories store and maintain inventories of resources and tools.
- ❑ Metadata inventories are exported and harvested in the network.
- ❑ Currently 13 repositories up and running; ~1.300 LRs available.





?

Username: t4me-admin

META=SHARE

1,248 language resources at your disposal...

Type in your keywords, please...

Search

Browse

Statistics

What is it? - About the project

META-NET aims at creating META-SHARE, a sustainable network of repositories of language data, tools and related web services documented with high-quality metadata, aggregated in central inventories allowing for uniform search and access to resources. Data and tools can be both open and with restricted access rights, free and for a fee. META-SHARE targets existing but also new and emerging language





Browse Resources

Community

Documentation

Statistics

Filter by:

Language

Serbian (14)

English (5)

Bulgarian (4)

Albanian (3)

French (3)

more

Resource Type

Media Type

Availability

Licence

Restrictions of Use

Validated

Foreseen Use

Use Is NLP Specific

Linguality Type

Multilinguality Type

Modality Type

MIME Type

14 Language Resources

Order by: Resource Name A-Z

- [Corpus of Contemporary Serbian](#) [\[details\]](#) [↓ 0](#) [👁 1](#)
[Serbian](#)
- [Corpus of Contemporary Serbian Newspapers and Magazines](#) [\[details\]](#) [↓ 0](#) [👁 1](#)
[Serbian](#)
- [ECI/MCI \(European Corpus Initiative/Multilingual Corpus I\)](#) [\[details\]](#) [↓ 0](#) [👁 1](#)
[Albanian](#) [Bulgarian](#) [Chinese](#) [Czech](#) [Danish](#) [Dutch](#) [English](#) [Estonian](#) [French](#) [Gaelic](#) [German](#) [Greek](#)
[Italian](#) [Japanese](#) [Latin](#) [Lithuanian](#) [Malay](#) [Norwegian](#) [Portuguese](#) [Russian](#) [Serbian](#) [Spanish](#) [Swedish](#)
[Turkish](#) [Uzbek](#)
- [English-Serbian Aligned Corpus](#) [\[details\]](#) [↓ 0](#) [👁 1](#)
[English](#) [Serbian](#)
- [French-Serbian Aligned Corpus](#) [\[details\]](#) [↓ 0](#) [👁 1](#)
[French](#) [Serbian](#)
- [INTERA Corpus](#) [\[details\]](#) [↓ 0](#) [👁 0](#)
[Bulgarian](#) [English](#) [Modern Greek \(1453-\)](#) [Serbian](#) [Slovenian](#)
- [KORLEX – Serbian Lexicon](#) [\[details\]](#) [↓ 0](#) [👁 0](#)

META-VISION

Language White Paper Series

Language White Paper Series



- ❑ Reports on the state of our languages in the digital age and the level of support through language technology.
- ❑ Series covers 30 languages.
- ❑ Key communication instruments to address decision makers and journalists.
- ❑ Inform about societal and technological problems and challenges as well as economic opportunities.
- ❑ >2 years in the making.
- ❑ >200 national experts as contributors.
- ❑ >8.000 copies printed and distributed to politicians and journalists.



30 Languages Covered



- ❑ Basque
- ❑ Bulgarian*
- ❑ Catalan
- ❑ 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*
- ❑ Croatian

* = Official EU language

Cross-Lingual Ranking

- ❑ In four application areas, each language is assigned to one of five clusters, ranging from *excellent LT support* to *weak/no support*:
 1. Machine Translation
 2. Speech Processing
 3. Text Analysis
 4. Resources
- ❑ Results finalised at a meeting in Berlin with representatives of all 30 languages
(October 21/22, 2011).



Text Analysis

MT

excellent	good	moderate	fragmentary	weak or no support
	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

excellent	good	moderate	fragmentary	weak or no support
	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

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	Croatian, Icelandic, Latvian, Lithuanian, Maltese, Romanian

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, Norwegian, Portuguese, Romanian, Serbian , Slovak, Slovene	Icelandic, Irish, Latvian, Lithuanian, Maltese

Europe's Languages and LT

English

Dutch
French
German
Italian
Spanish

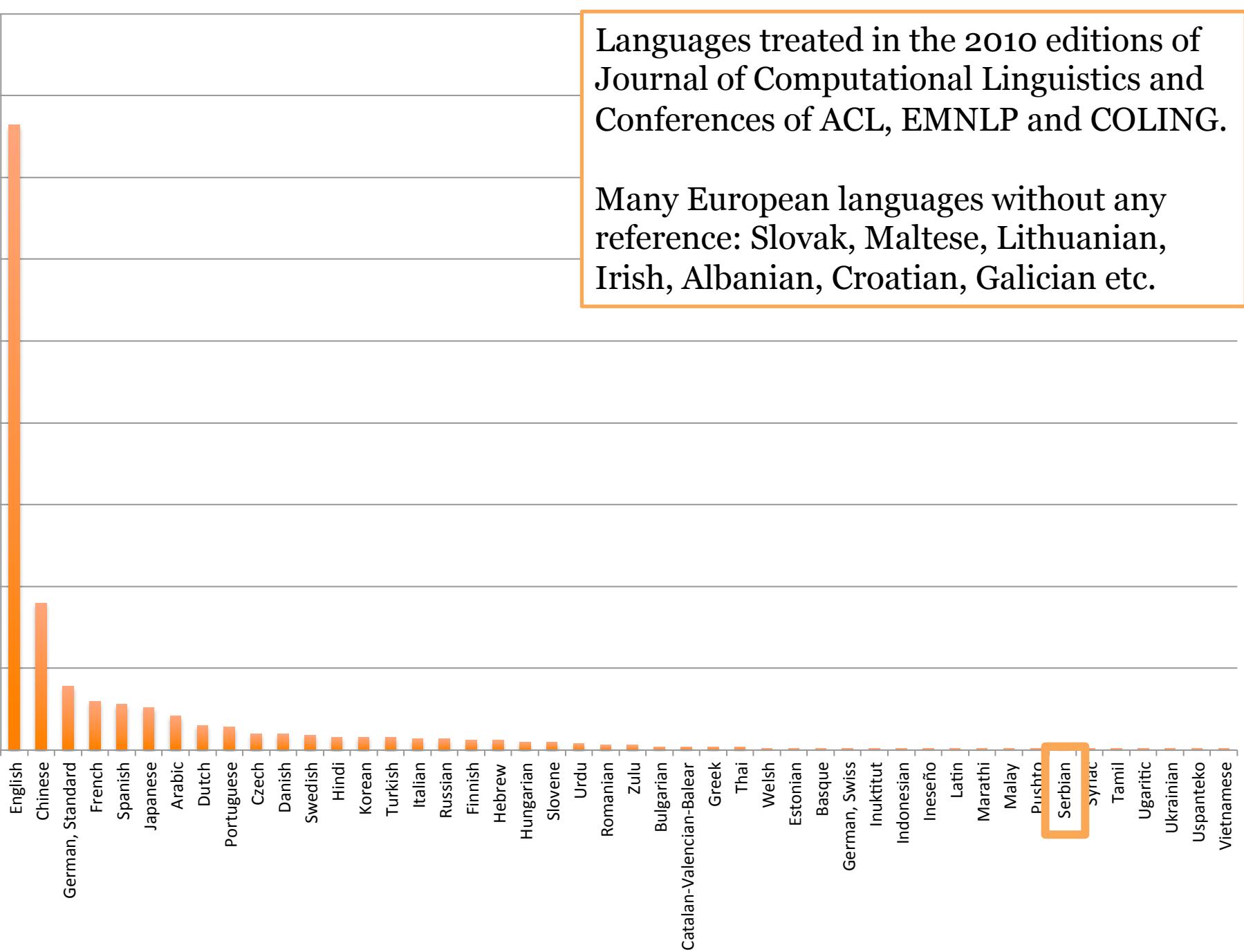
Catalan
Czech
Finnish
Hungarian
Polish
Portuguese
Swedish

Basque
Bulgarian
Danish
Galician
Greek
Norwegian
Romanian
Slovak
Slovene

Croatian
Estonian
Icelandic
Irish
Latvian
Lithuanian
Maltese
Serbian

*good support through
Language Technology*

*weak or
no support*



Languages treated in the 2010 editions of Journal of Computational Linguistics and Conferences of ACL, EMNLP and COLING.

Many European languages without any reference: Slovak, Maltese, Lithuanian, Irish, Albanian, Croatian, Galician etc.

Serbian

White Paper Website



www.meta-net.eu/whitepapers

META-NET White Paper Series

A Network of Excellence forging the Multilingual Europe Technology Alliance

META META-NET META-VISION META-SHARE META-RESEARCH News & Events LT-World Contact Search Site English

Accessibility | Contact | Log in

META-NET White Paper Series

- > Overview and List of Volumes
- > Quotes and Testimonials
- > Authors and Contributors
- > Key Results
- > Press Release
- > Press Coverage
- > The Team Behind the Series

Europe's Languages in the Digital Age

31 Volumes cover 30 European Languages

Basque, Bulgarian, Catalan, Croatian, Czech, Danish, Dutch, English, Estonian, Finnish, French, Galician, German, Greek, Hungarian, Icelandic, Irish, Italian, Latvian, Lithuanian, Maltese, Norwegian (bokmål), Norwegian (nynorsk), Polish, Portuguese, Romanian, Serbian, Slovak, Slovene, Spanish, Swedish.

At a Glance

- Key Results and Cross-Language Comparison

Aims and Scope

META-NET, a Network of Excellence consisting of 60 research centres from 34 countries, is dedicated to building the technological foundations of a multilingual European information society.

META-NET is forging META, the Multilingual Europe Technology Alliance. The benefits offered by Language Technology differ from language to language. So do the actions that need to be taken within META-NET, depending on the factors such as the complexity of the respective language, the size of its community, and the existence of active research centres in this area.

The META-NET Language White Paper series "Languages in the European Information Society" reports on the state of each European language with respect to Language Technology and explains the most urgent risks and chances. The series will cover all official European languages and several other languages spoken in geographical Europe. While there have been a number of valuable and comprehensive scientific studies on certain aspects of languages and technology, there exists no generally understandable compendium that takes a stand by presenting the main findings and challenges for each language. The META-NET white paper series will fill this gap.

Quotes and Testimonials

- Andrius Kubilius (Prime Minister of the Republic of Lithuania): "Having

White Paper Press Campaign



- Headline of press release:

**At Least 21 European Languages in Danger of Digital Extinction.
Good News and Bad News on the European Day of Languages.**

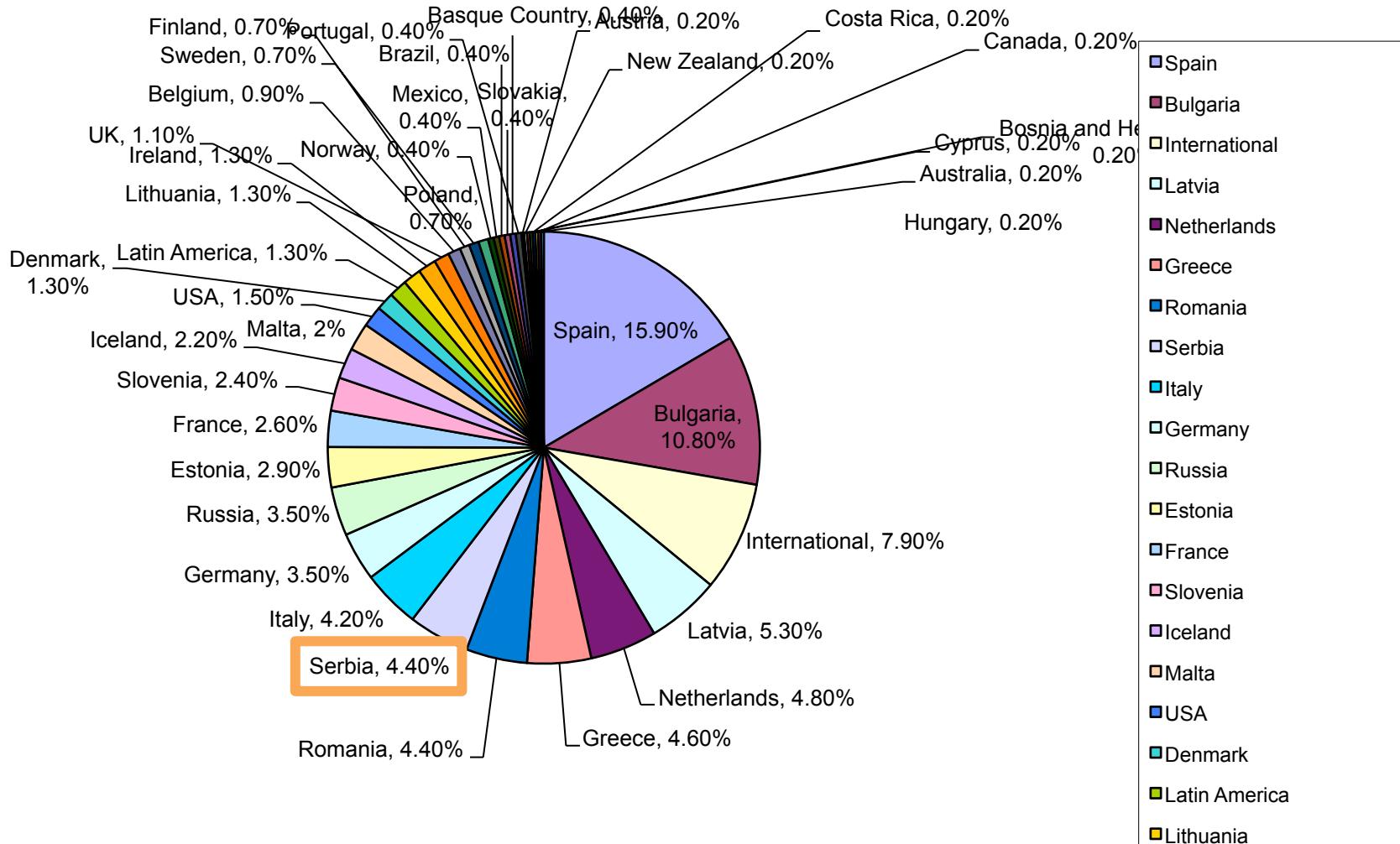
- Sent out to journalists, politicians and other stakeholder groups before the European Day of Languages (September 26).
- Overwhelmed by the huge interest in the topic and our key findings!
- 470+ mentions in the online and traditional press.
- 40+ interviews with META-NET representatives (television, radio).
- News came in from 41 countries in 35 different languages.

Response: Examples

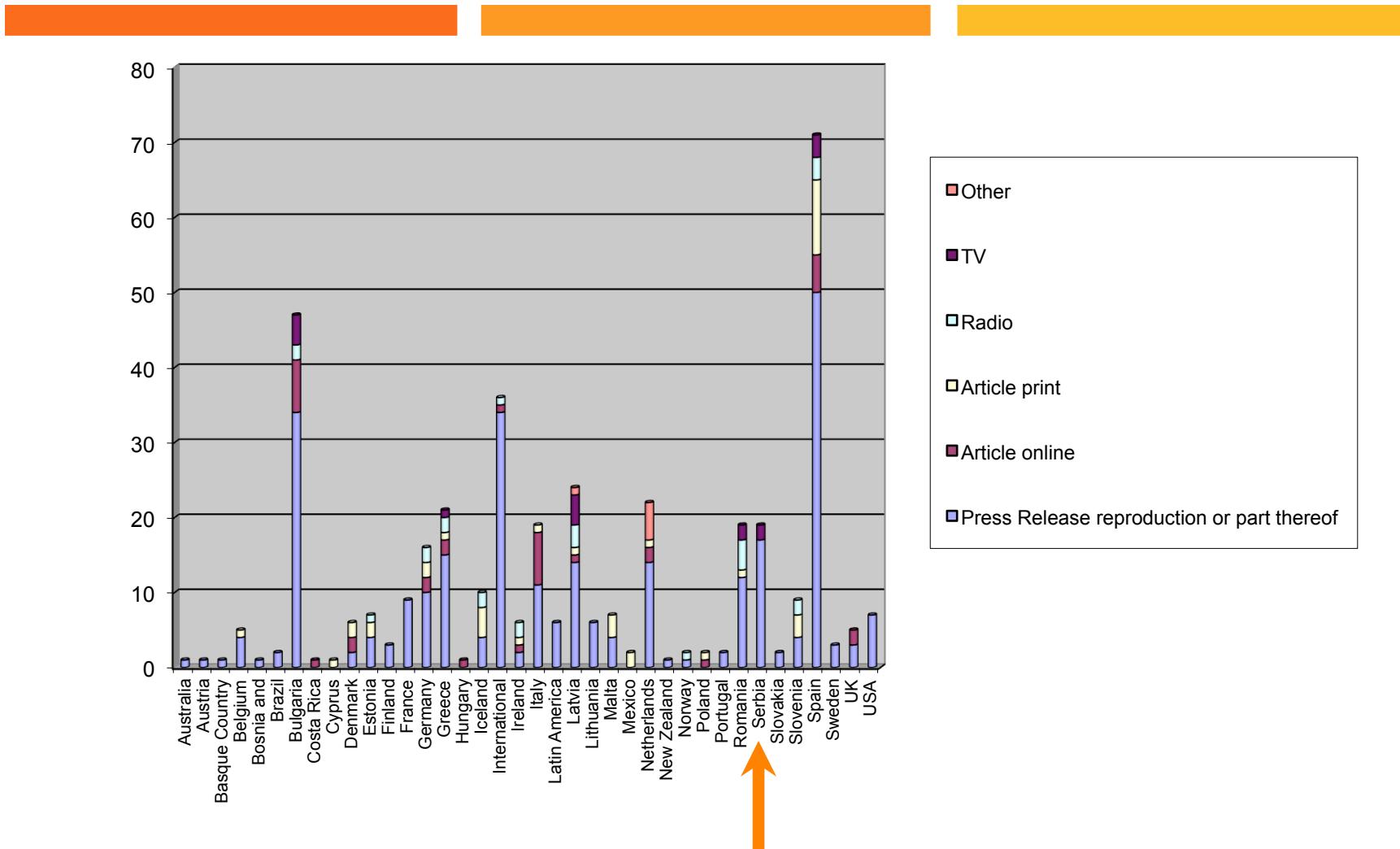
- ❑ **Austria:** Der Standard.
- ❑ **Denmark:** Politiken, Berlingske Tidende.
- ❑ **Finland:** Tiede.
- ❑ **Germany:** Heise Newsticker, Süddeutsche Zeitung.
- ❑ **Greece:** in.gr, Πρώτο Θέμα, Prosilipsis.
- ❑ **Iceland:** Fréttablaðið, Morgunblaðið.
- ❑ **Italy:** Wired.
- ❑ **Norway:** Computerworld.
- ❑ **Slovenia:** Delo, Dnevnik, Demokracija.
- ❑ **Serbia:** Politika, PTC1.
- ❑ **Spain:** El Mundo.
- ❑ **UK:** Huffington Post.
- ❑ **USA:** Mashable, NBC News, Reddit.



Coverage by Country



Coverage by Category



Press Campaign: Highlights



Ord. Forskere arbejder på at forbedre danske oversættelser på internettet.

Dårlig sprogteknologi truer dansk på nettet

Af Jens Ejding
// ejs@berlingske.dk

Det danske sprog har det svært i den digitale verden.

Det konstaterer danske sprogforskere- og eksperter i forbindelse med den nye internationale undersøgelse META-NET, der ser nærmere på, hvordan en lang række mindre, europæiske sprog som dansk klarer sig i den digitale verden.

Forskerne fra bl.a. Københavns Universitet og Dansk Sprognævn når frem til, at dansk i fremtiden kan få det endnu sværere i den digitale verden, fordi Google Translate, GPSer, applikationer til smartphones og andre sprogtækologiske programmer ikke i tilstrækkelig grad formår at behandle de mange nuancer i det danske sprog.

Professor i sprogtækologi på Københavns Universitet, Bolette Sandford Pedersen, mener, at der er brug for en slags digital dansk sprogbank fyldt med data, så bl.a. oversættelser bliver så præcise og gode som muligt. Med

hjælp fra sprogbanken kan forskere ifølge professoren hjælpe virksomheder med at forbedre programmer, der skal håndtere sproglig viden om bl.a. maskinoversættelse, talegenkendelse og informationsøsgning.

Dermed vil der blive længere mellem fejlagtige oversættelser, som når »hæld olie på panden« med Google Translate bliver til »pour oil on the forehead« på engelsk. Oversættelser, der er i værste fald er så uprecise, at danskere ender med at fravælge deres eget sprog i den digitale verden.

Sproghjælp til virksomheder

Hun anerkender dog, at »teknologien til automatiske oversættelser på mange måder er fantastisk«.

»Den er bare ikke god nok, når det gælder dansk,« siger hun:

»Det er som om, at vi i et vist omfang lægger det i hænderne på Google eller andre virksomheder at afgøre, om dansk skal behandles godt nok eller ej. Men det danske marked er ikke stort for dem. Spørgsmålet er derfor,

fakta ■

Sprog i Europa

■ Der er omkring 80 sprog i EU. For 21 af dem – også dansk – gælder det, at der er store sprogtækologiske mangler, når det gælder bl.a. maskinoversættelse, talegenkendelse og informationsøsgning.

■ Ifølge en EU-undersøgelse køber et stigende antal europæiske internetbrugere varer eller tjenester på nettet, hvor det sprog, der bliver anvendt, ikke er deres eget. Det gælder over halvdelen af brugerne.

■ Over hver tredje anvender et fremmedsprog til at skrive mail eller indlæg på nettet.

om vi ikke i højere grad selv skal gøre noget for at sikre, at det fornødne datamateriale er til rådighed, så vi får gode oversættelser og anden god sprogtækologi. Det kunne f.eks. være ved, at vi gjorde en indsats for at få oprettet en sprogbank med en masse beriget materiale om dansk.«

»Hvis vi hele tiden oplever, at oversættelser er behæftede med fejl, tør vi ikke stole på dem,« siger hun og understreger, at »fejlagtige oversættelser kan føre til store misforståelser.«

I følge Dansk Sprognævns direktør, Sabine Kirchmeier-Andersen, kan dårlig sprogtækologi have konsekvenser for mange danskere, der ikke er så gode til engelsk.

»Hvis vi har ambitioner om at bruge det danske sprog i fremtidens teknologiske univers, skal der gøres en indsats nu for at fastholde ekspertise og udbygge den viden, vi har,« mener hun:

»Ellers risikerer vi, at kun folk, der taler flydende engelsk, vil få glæde af de nye generatører af web-, tele- og robotteknologi, der er på vej.« ■

Press Campaign: Highlights

METANET

Latviešu valoda apdraudēta

Alma ORUPE, tālr. 67886751

Eiropas Savienībā (ES) ir aptuveni 60 runāto valodu, taču – cik ilgi tās kopiena varēs lepoties ar tādu kultūras un valodu bagātību? Pētījumi liecina, ka digitalizācijas laikmetā vismaz 21 no tām draud izmiršana. Arī latviešu valodai izdzīvot nebūs viegli, un valodu tehnoloģiju jomā tai nav daudz, ar ko lepoties – vairākumā sadaļu tā ir pēdējās vietās.

Starptautiskajā konferenčē *Valoda, tehnoloģijas un Eiropas naktne*, kur tika vērtēti valodas tehnoloģiju sasniegumi un problēmas, vairākkārt tika uzsvērts, ka ES valstu dažādās valodas ir gan bagātība, gan arī barjeira, kura traucē saziņu un sadarbišķi, kas sniedzas pār valstu robežām, kavē preču un pakalpojumu eksportu, uzņēmumu spēju strādāt ārējos tirgos. «Lai Eiropa saglabātu savu vadošo lomu pasaulei, tā būs vajadzīgas visām Eiropas valodām pie lägotas valodu tehnoloģijas, kas būs viegli pieejamas un efektivas,» norādīja Vācijas Mākslīgā intelekta pētījuma centra zinātniskais direktors Hanss Uškoreits, piebilstot, ka tas nozīmē, ka šajā jomā ES nāksies ieguldīt lielus līdzekļus. Ja 2008. gadā tulkošanai, programmatūru lokalizācijai un timekļa vietnē globa-

Atbalsts valodu tehnoloģijām (30 Eiropas valodās)				
Sektors	labs	viduvējs	fragmentārs	vājš vai nav
Runas apstrāde	1 (angļu)	9 (vācu, spānu u.c.)	14 (igauņu, polu u.c.)	6 (latviešu, lietuviešu u.c.)
Mašīntulksošana	1 (angļu)	2 (franču, spānu)	7 (itāliešu, polu)	19 (latviešu, lietuviešu u.c.)
Teksta analīze	1 (angļu)	5 (franču, itāliešu u.c.)	14 (sorbu, basku u.c.)	8 (latviešu, lietuviešu u.c.)
Runas un teksta resursi	1 (angļu)	9 (čehu, franču u.c.)	15 (igauņu, dāņu u.c.)	5 (latviešu, lietuviešu u.c.)



Konferencē prezentētajā pētījumā *Latviešu valoda digitālā laikmetā*, kurā piedalījusies 200 eksperti, vilktas paraleles starp diviem laikmetiem: kad parādījās grāmatu spiedi un pašreizējo – digitālo. Tās valodas, kurā netika drukātas grāmatas, ar laiku izmira, un tas varot notikti arī ar valodām, kas netiek liecotās globālajā tīmeklī. Patlaban privilēģētā stāvoklī esot plašāk liecotās Eiropas valodas, jo īpaši – angļu va-

loda. Salīdzinošā analīze rāda, ka tikai tai sniegtais atbalsts vērtējams kā labs. Holandiešu, franču, vācu, itāliešu un spānu valodai tas ir viduvējs, bet skaitliski nelielās Eiropas valodas, to skaitā lietuviešu un latviešu, vērtējas kā visapdraudētākās, liecina pētījums. Latviešu valoda četrās sektoros (runas un teksta resursi, tehnoloģijas runas apstrādei, mašīntulksošanai un teksta analīzei) ierakstīta sadaļa: vājš atbalsts vai nav

LATVIEŠU VALODAS ILGTSPĒJAS

POZITĪVIE ASPEKTI:

- ES oficiāla valoda
- Diegzān liels lietotāju skaits
- Valodai ir liela lingvistiskā kvalitāte
- ES fondu pieejamība projektos

RISKI:

- Maze tirgus
- Latviešiem saglabājties minoritātes kompleks
- Demogrāfiskā situācija
- Ekonomiskā atpalicība
- Emigrācija
- Minoritāsu pašpieliekamība
- Krievijas spiediens



PAREZAMI LIELI TĒRINI. Vācijas Mākslīgā intelekta pētījuma centra zinātniskais direktors Hanss Uškoreits norādīja, ka Eiropai, tā saglabātu savu vadošo lomu pasaulei, būs vajadzīgas visām Eiropas valodām pie lägotas valodu tehnoloģijas

atbalsta. Jāteic, ka igaunji Latviju kārtējo reizi ir apsteigusi. Kadrija Videre no Igaunijas valodas resursu centra skaidroja, ka tas noticis, jo Igaunijā valodu tehnoloģijas atbalstam ir ilgstošs un pastāvīgs valsts finansējums.

Latvijas pārstāvē datorzinātnieks doktore Inguna Skadipa bija spiesta atlīdzīt, ka kopumā neizskatāmies pārāk labi. «Līdzšinējais valsts atbalsts bijis nepietiekams, un tas novēdis pie

diezgan lielas atpalicības valodas ilgtspējai nepieciešamo resursu un rīku jomā. Arī pārreizējais pētniecības darbs bijis fragmentārs. Tāpat problēmas rada datorlingvistikas mācību trūkums. Latvijas augstskolās, un šobrīd tādu kursu var apgūt tikai Liepājas Universitātē un Rēzeknes augstskolā,» rezumēja I. Skadina, norādot, ka tikai ilgtspējīga politika un valsts atbalsts var glābt latviešu valodu. ■

Press Campaign: Highlights



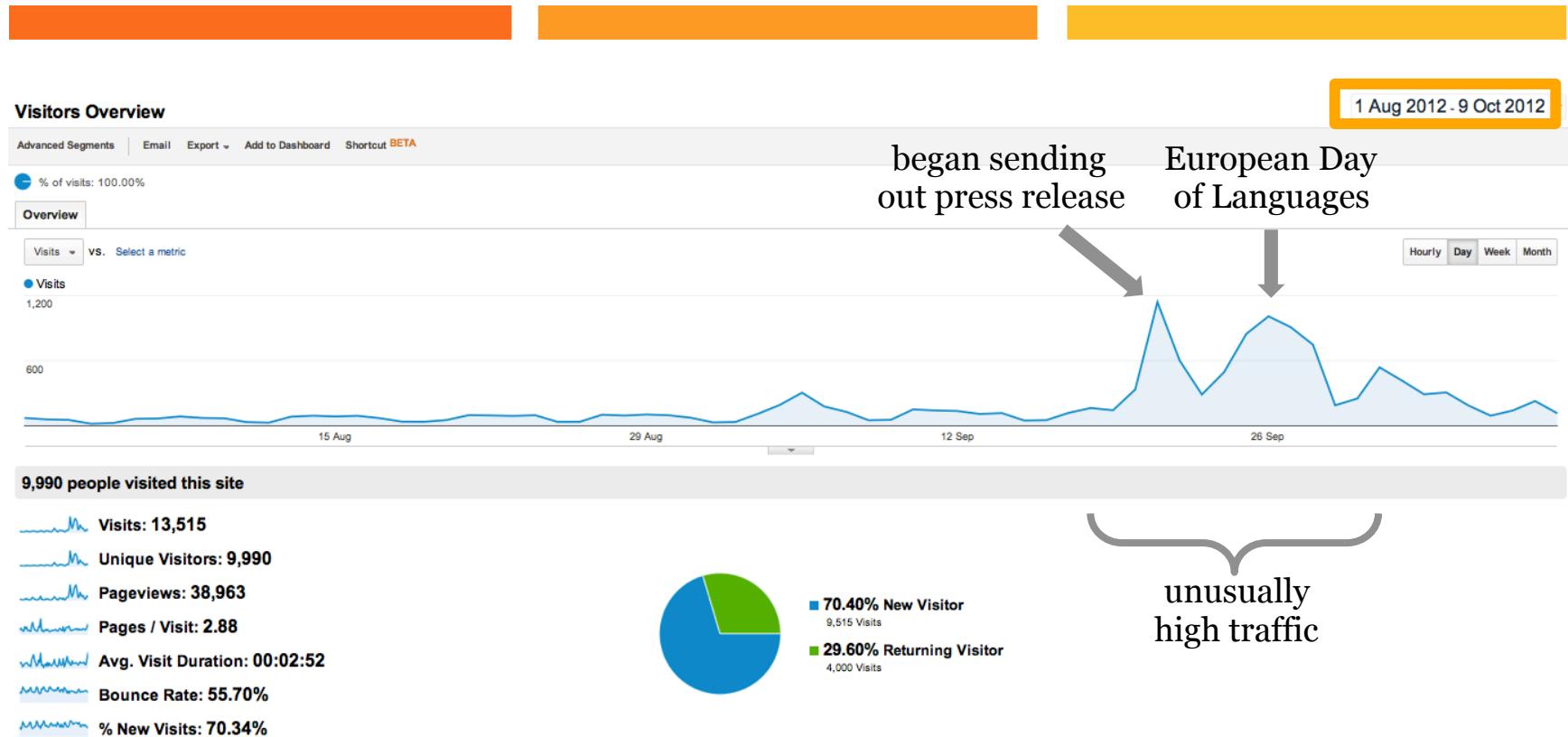
Published on Sep 28, 2012 by [Ivan Obradovic](#)

META-NET White paper press release interview, National Public Service TV Serbia

37 views

0 likes, 0 dislikes

Website: Visitors Overview



Website: Visitors' Cities





META-VISION

Strategic Research Agenda

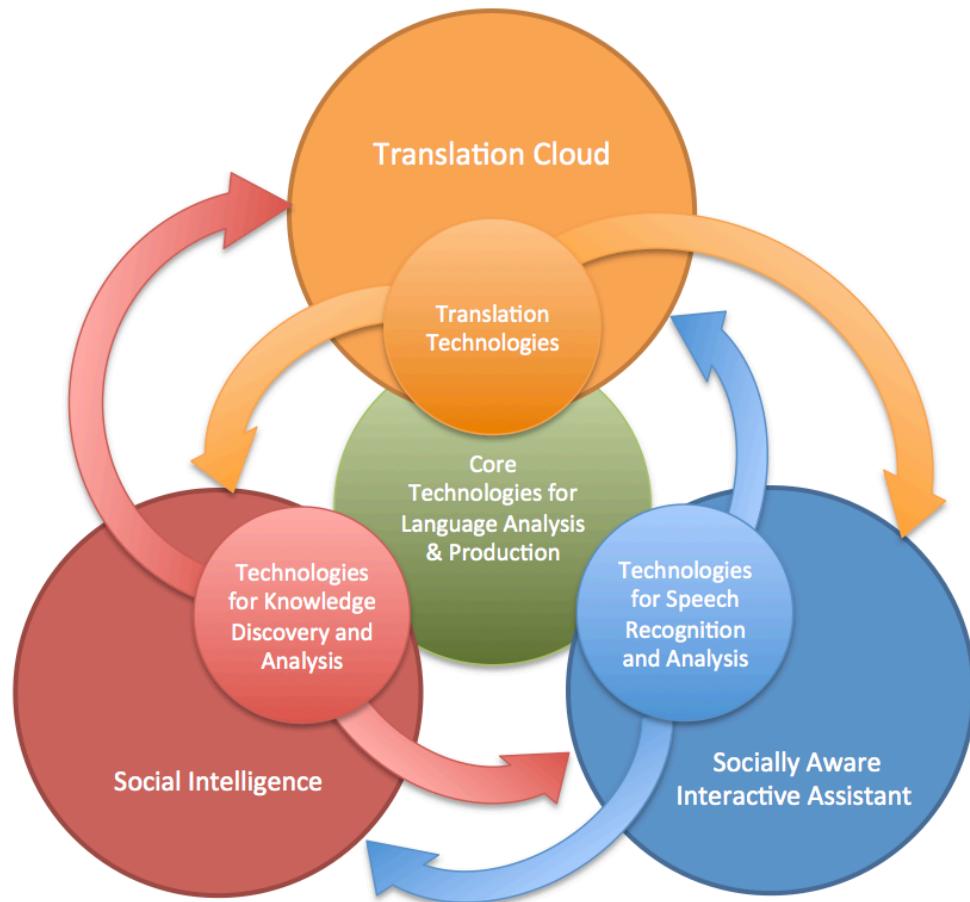
Strategic Research Agenda



- ❑ *META-NET Strategic Research Agenda for Multilingual Europe 2020.*
- ❑ Addresses the problems we found when preparing the white papers.
- ❑ Three priority research themes and application/innovation scenarios.
- ❑ Can put Europe ahead of its competitors in this technology area.
- ❑ 180+ contributors.
- ❑ Final version to be ready in Nov. 2012.
- ❑ SRA will be presented to the EC and national bodies.



Strategic Research Agenda



META

Multilingual Europe Technology Alliance



- **META-NET** is a network of excellence.
- **META** is an open and growing strategic technology alliance:
Multilingual Europe Technology Alliance.
 - 640+ members, including W3C, Google, Microsoft, GALA, research centres, LT companies etc.
 - Main goal: to support our Strategic Research Agenda.
 - **Join us!** <http://www.meta-net.eu/join>



META-NET

Recent News – Next Steps

Recent News – Next Steps



- ❑ **META-SHARE** Version 3.0 released in September 2012.
- ❑ Incoming resources from many EU-funded projects (40+ CAs), especially CESAR, METANET4U and META-NORD.
- ❑ META-SHARE launch event in January 2013.
- ❑ **META-TRUST AISBL** is an international non-profit organisation, founded on September 12, 2012: <http://www.meta-trust.eu>.
- ❑ Legal person of META-NET.
- ❑ **Strategic Research Agenda** press campaign, focus on social media (videos, infographics etc.).
- ❑ Meet with national research planners, funders, policy makers and inform them about the Strategic Research Agenda.
- ❑ **Data Liberation campaign.**

Q/A

Upcoming opportunities can provide sufficient resources to make our visions for Europe's citizens and economy, as described in the SRA, a reality (Horizon 2020; Connecting Europe Facility, CEF).

Thank you very much!

office@meta-net.eu

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

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

META^{NET}

