

META-NET

The META-NET Whitepaper Series on European Languages

Andrejs Vasiljevs

Tilde, Latvia

andrejs@tilde.com

META-NET TEAM

META-NET FORUM 2011 Solutions for Multilingual Europe
Budapest, Hungary, June 27/28

eu 2011.hu



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

- ❑ “Yet another study”?
- ❑ General Approach
- ❑ Preliminary Results
- ❑ Conclusions

Why Yet Another Survey?

- ❑ META-NET aims at a concerted effort to improve monolingual and multilingual LT support for all European languages.
- ❑ The degree to which LT is used in Europe varies from language to language depending on the commercial relevance of the language, the problems the language poses for automatic processing, and the research already devoted to it.
- ❑ So far, no one has ever evaluated the state of European languages in regards to LT support or their state in the digital information age.

The Language White Papers

- ❑ Survey of the state of the respective language in the digital society.
- ❑ Provide expert estimations about the current status and availability of language resources and technologies.
- ❑ Are meant to inform politicians, journalists and the public at large about societal and technological problems, challenges, and economic opportunities.



29 Languages Covered so far

- 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*

* = Official EU language

Structure of the White Papers

- ❑ Executive Summary
- ❑ Part 1: Introduction – A Risk for Our Languages and a Challenge for Language Technology
- ❑ Part 2: *Language* in the European Information Society
- ❑ Part 3: Language Technology Support for *Language*
- ❑ Part 4: About META-NET
- ❑ References

Assessing LT Support

- ❑ How to assess Language Technology support for a certain language?
- ❑ How to arrive at a result that can be communicated?
 - Count all existing tools and resources? => Does not result in a message.
 - Define quality criteria and perform a comparative evaluation? => Complicated, complex, time-consuming process, would take too long.
- ❑ For the White Papers, experts provided estimations condensed in a one table assessing core technology areas and resources such as:
 - Parsing, Information Retrieval, Machine Translation, Speech Recognition, Speech Synthesis, Reference Corpora, Language Models, Thesauri, etc.
- ❑ Assessment is done along criteria such as availability, quality, or coverage, maturity, sustainability and adaptability (see the example White Paper in your conference bag).

Preliminary Results

Preliminary Results

- Overall ranking of the quality and coverage of tools and resources is plausible (see the top and bottom tools and technologies on the right).
- Scale:
0 = non-existent
6 = perfect

3.8	Tokenization, Morphology
3.3	Speech Synthesis
2.7	Parsing
2.4	Machine Translation
2.4	Speech Recognition
	...
1.2	Text Semantics
0.9	Language Generation
0.8	Advanced Discourse Processing

Preliminary Results

- If one takes a value of at least 4 as threshold for practical usability of a technology/resource, then, e.g.,:
 - 13 of the 30 languages lack sufficient support in speech synthesis
 - 18 of the 30 languages lack sufficient support in parsing
 - 26 of the 30 languages lack sufficient support in machine translation

Example: MT Results

- For an example ranking and comparison,
 - a few values were manually adjusted and
 - normalization was carried out for each language, all technologies: $(x-m)/d$ with m mean and d standard deviation.

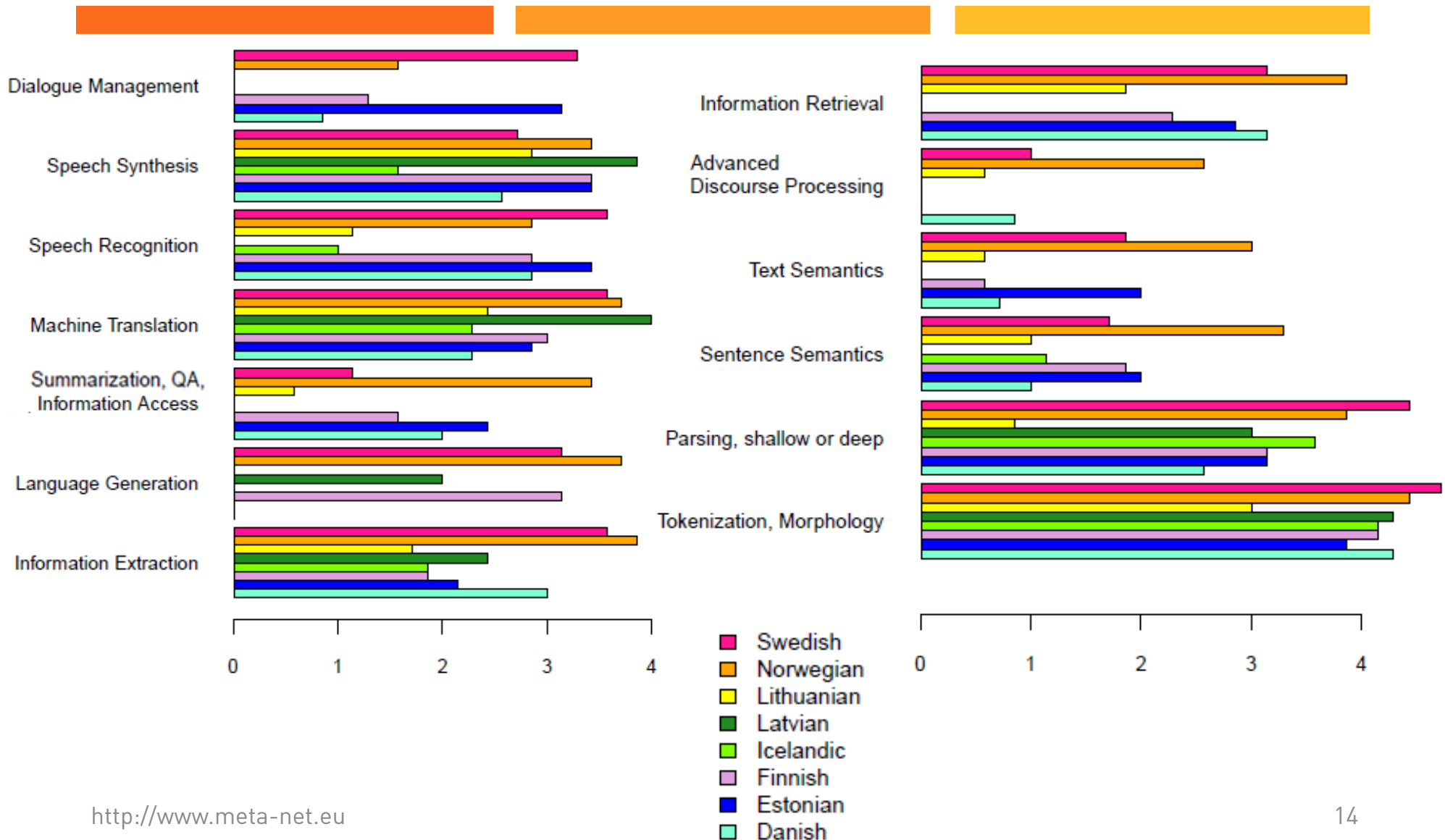
Cluster1 Higher ranks	Cluster2 Medium ranks	Cluster 3 Lower ranks
Spanish, English, Latvian, Lithuanian, Maltese, Galician, Slovene, Hungarian, Catalan	Polish, Icelandic, French, Irish, Basque, Italian, Romanian, Bulgarian, Dutch, Portuguese, Estonian, Finnish	Croatian, German, Swedish, Czech, Norwegian, Danish, Serbian, Greek
<i>Includes romanian languages with active research.</i>		<i>Difficult languages and languages with little research.</i>

Example: Comparing the Situation for Parsing

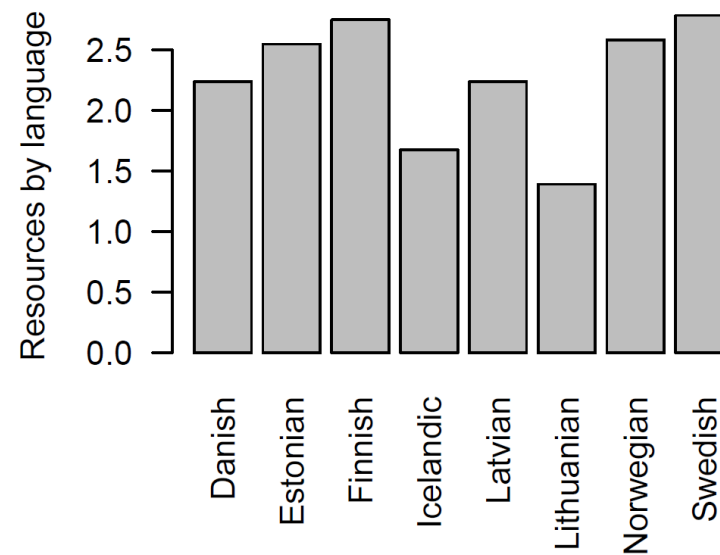
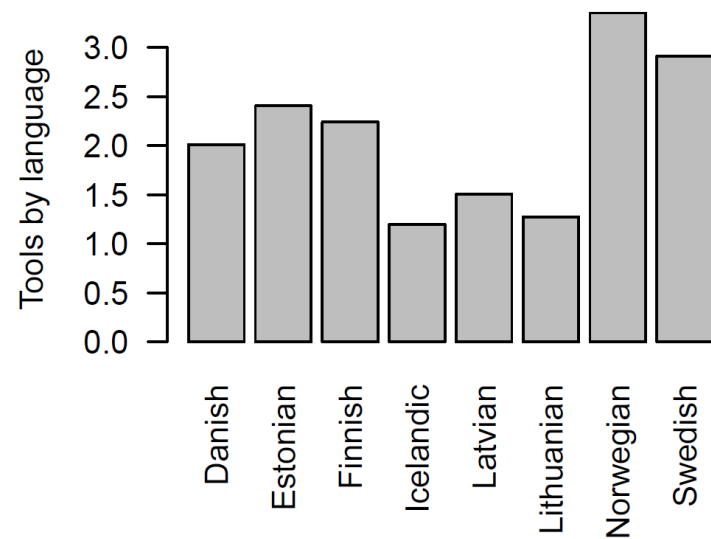
Cluster1	Cluster2	Cluster 3	Cluster 4
Czech, Dutch, English, German	Bulgarian, French, Norwegian, Polish, Portuguese, Spanish, Swedish	Basque, Catalan, Danish, Finnish, Galician, Hungarian, Irish, Italian, Romanian	Croatian, Estonian, Greek, Icelandic, Latvian, Lithuanian, Maltese, Serbian, Slovak, Slovene

For this example ranking, a few values were manually adjusted.

Example: Findings from the Nordic and Baltic White Papers



Example: Findings from the Nordic and Baltic White Papers



General Conclusions

- ❑ Speech processing and synthesis appears to be more mature than processing of written text. Advanced information access technology is in its infancy.
- ❑ Research was successful in designing particular high quality results in some areas, but many of the resources lack standardization, i.e., even if they exist, sustainability is not given; concerted programmes and initiatives are needed to standardize data and interchange formats.
- ❑ Most (very) large companies have stopped working in the area, leaving the field to SMEs, which can hardly attack an international market.

Q/A

META=NET

Thank you.

office@meta-net.eu

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

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