

THE БЪЛГАРСКИЯТ  
BULGARIAN ЕЗИК В  
LANGUAGE ДИГИТАЛНАТА  
IN THE ЕПОХА  
DIGITAL AGE

Diana Blagoeva  
Svetla Koeva  
Vladko Murdarov



---

White Paper Series

Серия Бели книги

THE BULGARIAN  
LANGUAGE  
IN THE  
DIGITAL AGE

БЪЛГАРСКИЯТ  
ЕЗИК В  
ДИГИТАЛНАТА  
ЕПОХА

Diana Blagoeva Bulgarian Academy of Sciences

Svetla Koeva Bulgarian Academy of Sciences

Vladko Murdarov Bulgarian Academy of Sciences

---

Georg Rehm, Hans Uszkoreit

(редактори, editors)

*Editors*

Georg Rehm  
DFKI  
Alt-Moabit 91c  
Berlin 10559  
Germany  
e-mail: georg.rehm@dfki.de

Hans Uszkoreit  
DFKI  
Alt-Moabit 91c  
Berlin 10559  
Germany  
e-mail: hans.uszkoreit@dfki.de

ISSN 2194-1416                      ISSN 2194-1424 (electronic)  
ISBN 978-3-642-30167-4            ISBN 978-3-642-30168-1 (eBook)  
DOI 10.1007/978-3-642-30168-1  
Springer Heidelberg New York Dordrecht London

Library of Congress Control Number: 2012940341

© Springer-Verlag Berlin Heidelberg 2012

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publisher's location, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)



## ПРЕДГОВОР

## PREFACE

Бялата книга е част от серия документи, представящи развитието в областта на езиковите технологии и техния потенциал. Документите са предназначени за преподаватели, журналисти, политици, различни езикови общности и т. н. Достъпът до езикови технологии за езиците, които се говорят в Европа, е много различен. Затова и необходимите действия за подкрепа на изследванията и развитието на езиковите технологии също са различни. Те зависят от много фактори, например сложността на даден език и броя на неговите носители.

META-NET, мрежа за високи постижения, изградена с подкрепата на Европейската комисия, предлага анализ на съществуващите езикови ресурси и технологии в серията Бели книги (р. 81). Анализът е съсредоточен върху 23-те официални европейски езика, наред с други по-важни национални и регионални езици. Резултатите показват значителен недостиг за всеки език. Задълбоченият експертен анализ и оценка на актуалната ситуация ще помогнат за увеличаване на ефекта от изследванията и намаляване на риска от пропуски.

В META-NET участват 54 изследователски центъра от 33 страни (р. 77), работещи съвместно с търговски и правителствени организации, научни институции, софтуерни компании, фирми за информационни технологии и европейски университети. Те разработват заедно визия за технологично развитие и стратегическа програма за научни изследвания, които показват как езиковите технологии могат да отговорят на научните предизвикателства до 2020 г.

This white paper is part of a series that promotes knowledge about language technology and its potential. It addresses journalists, politicians, language communities, educators and others. The availability and use of language technology in Europe varies between languages. Consequently, the actions that are required to further support research and development of language technologies also differ. The required actions depend on many factors, such as the complexity of a given language and the size of its community.

META-NET, a Network of Excellence funded by the European Commission, has conducted an analysis of current language resources and technologies in this white paper series (p. 81). The analysis focused on the 23 official European languages as well as other important national and regional languages in Europe. The results of this analysis suggest that there are tremendous deficits in technology support and significant research gaps for each language. The given detailed expert analysis and assessment of the current situation will help maximise the impact of additional research.

As of November 2011, META-NET consists of 54 research centres from 33 European countries (p. 77). META-NET is working with stakeholders from economy (software companies, technology providers and users), government agencies, research organisations, non-governmental organisations, language communities and European universities. Together with these communities, META-NET is creating a common technology vision and strategic research agenda for multilingual Europe 2020.

Авторите на този документ благодарят сърдечно на авторите на Бялата книга за немски за предоставената възможност да използват избрани езиково независими части [1].

Разработката на настоящата Бяла книга е финансирана по Седма рамкова програма и Програма ICT Policy Support Programme на Европейската комисия, договори T4ME (договор за финансиране 249119), CESAR (договор за финансиране 271022), METANET4U (договор за финансиране 270893) и META-NORD (договор за финансиране 270899).

---

The authors of this document are grateful to the authors of the White Paper on German for permission to re-use selected language-independent materials from their document [1].

The development of this white paper has been funded by the Seventh Framework Programme and the ICT Policy Support Programme of the European Commission under the contracts T4ME (Grant Agreement 249119), CESAR (Grant Agreement 271022), METANET4U (Grant Agreement 270893) and META-NORD (Grant Agreement 270899).



# СЪДЪРЖАНИЕ CONTENTS

## ЕЗИЦИТЕ В ЕВРОПЕЙСКОТО ЕЗИКОВО ИНФОРМАЦИОННО ОБЩЕСТВО

<b>1</b>	<b>Резюме</b>	<b>1</b>
<b>2</b>	<b>Заплаха за езиците и предизвикателство пред езиковите технологии</b>	<b>6</b>
2.1	Езиковите граници – пречка пред европейското информационно общество . . . . .	7
2.2	Рискът за нашите езици . . . . .	7
2.3	Езиковите технологии предоставят възможности . . . . .	8
2.4	Перспективи пред езиковите технологии . . . . .	8
2.5	Предизвикателства пред езиковите технологии . . . . .	10
2.6	Как хората и машините учат език? . . . . .	10
<b>3</b>	<b>Българският език в европейското информационно общество</b>	<b>12</b>
3.1	Общи данни . . . . .	12
3.2	Особености на българския език . . . . .	12
3.3	Актуално . . . . .	14
3.4	Езикова политика в България . . . . .	15
3.5	Езикът в образованието . . . . .	17
3.6	Международен статут на българския език . . . . .	18
3.7	Българският език в интернет . . . . .	19
<b>4</b>	<b>Приложение на езиковите технологии за български</b>	<b>21</b>
4.1	Архитектура на стандартна система за езикова обработка . . . . .	21
4.2	Основни сфери на приложение . . . . .	23
4.3	Други сфери на приложение . . . . .	30
4.4	Образователни програми за езикови технологии . . . . .	31
4.5	Национални проекти и инициативи . . . . .	32
4.6	Налични програми и ресурси . . . . .	33
4.7	Сравнение между езиковите технологии за отделните езици . . . . .	33
4.8	Заклучение . . . . .	35
<b>5</b>	<b>За META-NET</b>	<b>39</b>

# THE BULGARIAN LANGUAGE IN THE DIGITAL AGE

<b>1</b>	<b>Executive Summary</b>	<b>41</b>
<b>2</b>	<b>Languages at Risk: a Challenge for Language Technology</b>	<b>45</b>
2.1	Language Borders Hold back the European Information Society . . . . .	46
2.2	Our Languages at Risk . . . . .	46
2.3	Language Technology is a Key Enabling Technology . . . . .	46
2.4	Opportunities for Language Technology . . . . .	47
2.5	Challenges Facing Language Technology . . . . .	48
2.6	Language Acquisition in Humans and Machines . . . . .	48
<b>3</b>	<b>The Bulgarian Language in the European Information Society</b>	<b>50</b>
3.1	General Facts . . . . .	50
3.2	Particularities of the Bulgarian Language . . . . .	50
3.3	Recent Developments . . . . .	52
3.4	Official Language Protection in Bulgaria . . . . .	53
3.5	Language in Education . . . . .	54
3.6	International Aspects . . . . .	55
3.7	Bulgarian on the Internet . . . . .	56
<b>4</b>	<b>Language Technology Support for Bulgarian</b>	<b>58</b>
4.1	Application Architectures . . . . .	58
4.2	Core Application Areas . . . . .	59
4.3	Other Application Areas . . . . .	66
4.4	Educational Programmes . . . . .	67
4.5	National Projects and Initiatives . . . . .	67
4.6	Availability of Tools and Resources . . . . .	68
4.7	Cross-language comparison . . . . .	69
4.8	Conclusions . . . . .	70
<b>5</b>	<b>About META-NET</b>	<b>74</b>
<b>A</b>	<b>Цитирани източници – References</b>	<b>75</b>
<b>B</b>	<b>Организации членки на META-NET – META-NET Members</b>	<b>77</b>
<b>C</b>	<b>Серия Бели книги на META-NET – The META-NET White Paper Series</b>	<b>81</b>

## РЕЗЮМЕ

Информационните технологии променят живота ни. Ежедневно използваме компютри, за да пишем, редактираме, изчисляваме, търсим информация, и все по-често за да четем, слушаме музика, да разглеждаме снимки или да гледаме филми. Носим малки компютри в джобовете си и ги използваме, за да позвъним, напишем имейл, да получим необходима информация или просто за да се позабавяваме, където и да се намираме. Как широкообхватната дигитализация на информация, знание и ежедневна комуникация влияе върху езика ни? Дали езикът ни ще се промени и дали е възможно дори да изчезне?

Компютрите ни са свързани един с друг в непрекъснато развиваща се компактна и мощна глобална мрежа. Момиче от Ипанема, ученик от Жеравна и инженер от Катманду могат да общуват с приятелите си във Фейсбук, но е малко вероятно да се срещнат един с друг в онлайн общности или форуми. Ако ги интересува как да излекуват главоболието си, те могат да потърсят повече информация в Уикипедия, но дори и тогава няма да прочетат една и съща статия. Когато европейските потребители на интернет обсъждат във форуми и чатове последиците от катастрофата във Фукушима за европейската енергийна политика, те го правят в ясно разграничени езикови общности. Всичко, което интернет съдържа и предлага, е все още разделено от езика на потребителите. Винаги ли ще бъде така? В научнофантастичните филми всички говорят един език – английски, китайски или български – в зависимост от това, къде се излъчва филмът. Възможно ли е ези-

кът на космонавтите да бъде български, въпреки че те рядко биха употребявали български думи толкова естествено, колкото английски? Много от съществуващите в момента 6 000 езика едва ли ще оцелеят в глобализираното дигитално информационно общество. Предполага се, че поне 2 000 езика са обречени на изчезване в идните десетилетия. Други ще продължат да играят роля на семейно и регионално ниво, но не и в по-широки делови или академични кръгове. Какви са шансовете на българския език да оцелее? Българският език се говори от близо 9 милиона души предимно в България, но също и в Гърция, Македония, Румъния, Турция (европейската част), Украйна, Австралия, Канада, САЩ, Германия и Испания. За малка страна като България съществува относително голямо количество телевизионни канали на български език – седем национални телевизии, 16 кабелни и сателитни телевизии с многорегионално покритие и 46 – с регионално покритие. Повечето чужди филми са дублирани на български език. Книгите се връщат на мода, въпреки констатациите, че през последните години българинът е спрял да се интересува от литература. Българският е първият славянски език, който разполага със своя собствена писмена система, датираща от 9-ти век. На 1 януари 2007 г., когато България е приета за пълноправен член на Европейския съюз, кирилицата става третата официална азбука на Европейския съюз след латинската и гръцката. Някои среди изразяват недоволство от нарастващата употреба на чужди думи, особено английски, и дори съществуват



страхове, че българският език ще се „прояде“ от множество английски думи и изрази. През вековете българският език е устоял на влиянието на думи и термини от гръцки и латински – езиците на познанието, както и на навлизането на френски думи през 18-ти и 19-ти век. Добро противодействие срещу изчезването на обичаните от нас български думи е наистина да ги използваме – често и съзнателно. Главното ни притеснение не трябва да е нарастващото английско влияние върху езика, а пълното му изчезване от някои основни области на личния ни живот. Нито науката, нито авиацията или глобалният финансов пазар се нуждаят от език, разпространен по целия свят – *lingua franca*. В много области на живота е по-важно общуването с гражданите на страната, отколкото с международните партньори – вътрешната политика, например административните процедури, правото, културата и търговията. Статутът на езика зависи не само от броя на неговите носители, създадените книги и филми, телевизионните канали, които го използват, но и от присъствието на езика в дигиталното информационно пространство и софтуерните приложения. В това отношение българският език е относително добре представен: всички важни международни софтуерни продукти са локализирани за български, българската Уикипедия е на 34-а позиция сред 270 в света. Потребителите на интернет в България през 2009 г. са се увеличили с 31% спрямо 2007 г. и вече са 46% от цялото население. В областта на езиковите технологии за български също съществуват редица продукти, технологии и ресурси. Има приложения за възпроизвеждане на реч, проверка на правописа и граматиката. Съществуват и програми за автоматичен превод, макар че не винаги се предлагат лингвистично коректни преводи, особено когато преводът е от друг език на български. Това се дължи основно на специфичните езикови характеристики на българския език. Информа-

ционните и комуникационните технологии се подготвят за следваща революционна стъпка. След персоналните компютри, мрежите, миниатюризацията на техниката, мултимедията, мобилните устройства и паралелната обработка на информация, идва епохата на технологии, които ще разбират не просто букви или звукове, но и словосъчетания и изречения. Така те ще подпомагат в много по-голяма степен потребителите, тъй като ще говорят, знаят и разбират техния език. Пионери в тази сфера са например Гугъл преводачът, който предлага безплатен онлайн автоматичен превод между 57 езика, супер компютърът на IBM Watson, който победи шампиона на САЩ в играта „Jeopardy“ или мобилният асистент Siri на iPhone, който реагира на гласови команди и отговаря на въпроси на английски, немски, френски и японски.

Следващото поколение информационни технологии ще се усъвършенства в употребата на естествения език до такава степен, че потребителите ще общуват, използвайки технологиите на собствения си език. Устройствата ще могат автоматично да намерят най-важните новини и информация в световното дигитално изобилие от познание само с помощта на гласови команди. Езиковите технологии ще предлагат автоматичен превод или ще подпомагат превода, ще осигуряват резюмиране на диалог или на различни документи, а компютърно подпомогнатото обучение ще съдейства за по-лесното интегриране на малцинствени групи и чужденци. Следващото поколение информационни и комуникационни технологии ще създаде индустриални и обслужващи роботи (в момента все още в научните лаборатории), които точно ще разбират какво искат техните потребители и ще рапортуват за изпълнението на задачите си. Такова равнище на работа надхвърля простите множества от символи и речници, програми за проверка на правописа и правила за