

THE SPANISH LANGUAGE IN THE DIGITAL AGE
LA LENGUA ESPAÑOLA EN LA ERA DIGITAL

Maite Melero
Toni Badia
Asunción Moreno



White Paper Series

Serie de Libros Blancos

THE SPANISH
LANGUAGE IN
THE DIGITAL
AGE

LA LENGUA
ESPAÑOLA
EN LA ERA
DIGITAL

Maite Melero Barcelona Media Centre d'Innovació

Toni Badia Universitat Pompeu Fabra

Asunción Moreno Universitat Politècnica de Catalunya

Georg Rehm, Hans Uszkoreit
(editores, 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 978-3-642-30840-6 ISBN 978-3-642-30841-3 (eBook)
DOI 10.1007/978-3-642-30841-3
Springer Heidelberg New York Dordrecht London

Library of Congress Control Number: 2012946618

© 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)



PRÓLOGO

Este documento es parte de una serie de Libros Blancos que promueve el conocimiento sobre las tecnologías del lenguaje y su potencial, dirigida a educadores, periodistas, políticos y las propias comunidades lingüísticas.

La disponibilidad y el uso de las tecnologías lingüísticas en Euro-pa varían según el idioma. Por lo tanto, las acciones requeridas para apoyar la investigación y el desarrollo de estas tecnologías también varían para cada idioma. Las acciones necesarias también dependen de factores como la complejidad intrínseca de la lengua y el tamaño de su comunidad.

META-NET, una red de excelencia financiada por la Comisión Europea, ha llevado a cabo un análisis de los recursos lingüísticos y las tecnologías actuales para cada lengua (p. 79). Este análisis se ha centrado en las 23 lenguas oficiales europeas, así como en otros idiomas importantes a nivel nacional y regional en Europa. Los resultados sugieren que existen todavía muchas lagunas por cubrir en este área.

META-NET se compone de 54 centros de investigación de 33 países que están trabajando con representantes de empresas comerciales, agencias gubernamentales, industria, organizaciones de investigación, empresas de software, proveedores de tecnología y universidades europeas [1].

Juntos están creando una visión tecnológica común, y al mismo tiempo están desarrollando una agenda estratégica de investigación que, a 10 años vista, permita a las aplicaciones basadas en tecnología lingüística abordar las deficiencias detectadas.

PREFACE

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. 79). 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 [1]. 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.

Los autores de este documento agradecen a los autores del Libro Blanco para el alemán el permiso para reutilizar material seleccionado de su texto original [2].

La elaboración de este Libro Blanco ha sido financiada por el Séptimo Programa Marco y el Programa de Políticas de Apoyo a las TIC de la Comisión Europea mediante los contratos T4ME (GA 249 119), CESAR (GA 271 022), META-NET4U (GA 270 893) i META-NORD (GA 270 899).

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 [2].

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 249 119), CESAR (Grant Agreement 271 022), METANET4U (Grant Agreement 270 893) and META-NORD (Grant Agreement 270 899).



CONTENIDO CONTENTS

LA LENGUA ESPAÑOLA EN LA ERA DIGITAL

1	Resumen ejecutivo	1
2	Un riesgo para nuestras lenguas y un reto para la tecnología lingüística	4
2.1	Las fronteras lingüísticas son un obstáculo para la Sociedad de la Información Europea	5
2.2	Nuestras lenguas en peligro	5
2.3	La tecnología lingüística es clave en el desarrollo de muchas aplicaciones	6
2.4	Oportunidades para la tecnología lingüística	6
2.5	Retos de la tecnología lingüística	7
2.6	Adquisición del lenguaje por los seres humanos y por las máquinas	8
3	El español en la Sociedad de la Información Europea	10
3.1	Datos generales	10
3.2	Particularidades de la Lengua Española	10
3.3	Acontecimientos recientes	11
3.4	El cultivo del idioma en España	12
3.5	La lengua en la educación	14
3.6	Aspectos internacionales	14
3.7	El español en Internet	16
4	Recursos de tecnología lingüística para el español	18
4.1	Arquitectura de las aplicaciones	18
4.2	Áreas principales de aplicación	20
4.3	Otras áreas de aplicación	26
4.4	Programas Educativos	28
4.5	Proyectos e iniciativas nacionales	29
4.6	Disponibilidad de herramientas y recursos	30
4.7	Comparación entre lenguas	32
4.8	Conclusiones	33
5	Acerca de META-NET	36

THE SPANISH LANGUAGE IN THE DIGITAL AGE

1	Executive Summary	37
2	Languages at Risk: a Challenge for Language Technology	39
2.1	Language Borders Hold back the European Information Society	40
2.2	Our Languages at Risk	40
2.3	Language Technology is a Key Enabling Technology	40
2.4	Opportunities for Language Technology	41
2.5	Challenges Facing Language Technology	42
2.6	Language Acquisition in Humans and Machines	42
3	The Spanish Language in the European Information Society	44
3.1	General Facts	44
3.2	Particularities of the Spanish Language	44
3.3	Recent Developments	45
3.4	Language Cultivation in Spain	46
3.5	Language in Education	47
3.6	International Aspects	48
3.7	Spanish on the Internet	50
4	Language Technology Support for Spanish	52
4.1	Application Architectures	52
4.2	Core Application Areas	53
4.3	Other Application Areas	60
4.4	Educational Programmes	62
4.5	National Projects and Initiatives	63
4.6	Availability of Tools and Resources	63
4.7	Cross-language comparison	65
4.8	Conclusions	66
5	About META-NET	69
A	Referencias – References	71
B	META-NET Miembros – META-NET Members	75
C	La serie de Libros Blancos de META-NET – The META-NET White Paper Series	79

RESUMEN EJECUTIVO

Durante los últimos 60 años, Europa ha adquirido una estructura política y económica singular y, sin embargo, cultural y lingüísticamente todavía es muy diversa. Esto implica que, del portugués al polaco y del italiano al islandés, la comunicación cotidiana entre los ciudadanos europeos, así como la comunicación en el ámbito de los negocios y de la política se enfrenta inevitablemente con las barreras del idioma. Las instituciones de la UE gastan casi mil millones de euros al año en mantener su política de multilingüismo, es decir, traduciendo documentos e interpretando intervenciones orales. Sin embargo, ¿es necesario que estas tareas sigan siendo una carga pesada y costosa? Las modernas tecnologías de la lengua y la investigación lingüística pueden contribuir significativamente a derribar las barreras lingüísticas. En el futuro, las tecnologías lingüísticas, combinadas con dispositivos y aplicaciones inteligentes, serán capaces de ayudar a los europeos a comunicarse fácilmente entre sí y a hacer negocios, incluso si no hablan una lengua común.

La tecnología lingüística construye puentes para el futuro de Europa.

El mercado único europeo resulta beneficioso para los países que lo integran. Sin embargo, las barreras del idioma pueden suponer un freno para los negocios, especialmente para las pequeñas y medianas empresas, que no tienen los medios económicos para afrontar la situación. La única (e impensable) alternativa a esta Europa multilingüe sería permitir que una sola lengua tomara una posición dominante y terminara reemplazando al resto de lenguas.

El aprendizaje de lenguas extranjeras siempre ha sido una forma habitual de superar las barreras lingüísticas. Sin embargo, sin apoyo tecnológico, llegar a dominar las 23 lenguas oficiales de los Estados miembros de la Unión Europea, sumadas a 60 lenguas más no oficiales, supone un obstáculo insuperable para los ciudadanos de Europa así como para su economía, para el debate político y para el progreso científico.

La solución consiste en invertir en las tecnologías clave para la superación de estas barreras, es decir, las tecnologías lingüísticas. Estas tecnologías ofrecen ventajas enormes, no sólo dentro del mercado común europeo, sino también en las relaciones comerciales con terceros países, especialmente en las economías emergentes. Para lograr este objetivo, y preservar la diversidad cultural y lingüística de Europa, es conveniente considerar las particularidades lingüísticas de todos los idiomas europeos, y analizar el estado actual de las tecnologías lingüísticas para cada uno de ellos. Las tecnologías lingüísticas constituirán en el futuro un puente único entre las lenguas de Europa.

La traducción automática y las herramientas de procesamiento del habla que están actualmente disponibles en el mercado, aún están lejos de alcanzar este ambicioso objetivo. Los agentes dominantes en estos ámbitos son principalmente empresas de propiedad privada radicadas en América del Norte. Ya en la década de 1970, la UE se dio cuenta de la enorme relevancia de la tecnología lingüística como conductor de la unidad europea, y comenzó a financiar sus primeros grandes proyectos de investigación, como EUROTRA. Al mismo tiempo, se establecieron proyectos nacionales que generaron resul-

tados valiosos, pero que nunca llevaron a una acción europea concertada. En contraste con este esfuerzo financiero altamente selectivo, otras sociedades multilingües, como la India (22 lenguas oficiales) y Sudáfrica (11 lenguas oficiales) han creado recientemente programas nacionales a largo plazo para la investigación lingüística y el desarrollo tecnológico.

Las tecnologías de la lengua, clave para el futuro.

Muchas de las aplicaciones de tecnología lingüística utilizan actualmente métodos estadísticos, que se basan en grandes cantidades de datos e ignoran las propiedades intrínsecas de la lengua. Por ejemplo, algunos de los sistemas de traducción más populares traducen automáticamente mediante la comparación de la frase a traducir con centenares de miles de frases previamente traducidas por humanos. La calidad de la producción depende en gran medida de la cantidad y la calidad de la muestra de corpus disponible. Así, mientras que la traducción automática de oraciones sencillas, en idiomas con una cantidad suficiente de texto disponible, puede obtener resultados útiles, los métodos estadísticos puros están condenados al fracaso en el caso de idiomas con cantidades mucho menores de datos, o en el caso de las construcciones sintácticas complejas. Dada esta situación, la Unión Europea ha decidido financiar proyectos tales como EuroMatrix y EuroMatrixPlus (desde 2006) y iTranslate4 (desde 2010) que llevan a cabo investigación básica y aplicada y generan recursos lingüísticos para todos los idiomas europeos. El análisis de las propiedades estructurales más profundas de la lengua es el único camino posible si queremos crear aplicaciones que funcionen bien para toda la gama de las lenguas de Europa.

La investigación europea en este ámbito ya ha logrado varios éxitos. Por ejemplo, los servicios de traducción de la Unión Europea actualmente están utilizando MO-

SES, una aplicación de traducción automática de código abierto, que se ha desarrollado principalmente a través de proyectos de investigación europeos. Muchos de los laboratorios de investigación y desarrollo (por ejemplo, IBM y Philips) han cerrado o se han trasladado a otro lugar. En lugar de construir sobre los resultados de sus proyectos de investigación, Europa ha tendido a realizar actividades de investigación aisladas, con menor impacto en el mercado. Incluso el valor económico de estos primeros esfuerzos, puede verse en el número de spin-offs. Una compañía como Trados, que fue fundada en 1984, fue vendida a SDL, con sede en el Reino Unido, en 2005.

Basándose en los conocimientos acumulados hasta el momento, parece claro que la actual tendencia a la tecnología híbrida, mezcla de procesamiento lingüístico de la lengua con métodos estadísticos, será capaz de reducir la brecha entre todas las lenguas europeas e ir más allá. Como esta serie de “libros blancos” muestra, existen grandes diferencias de preparación entre los diferentes idiomas y estados europeos con respecto a las tecnologías de la lengua. Sin embargo, incluso los idiomas “grandes” como el español, todavía necesitan dedicar más recursos a la investigación con objeto de que las soluciones tecnológicas estén realmente listas para su uso cotidiano.

El objetivo a largo plazo de META-NET es introducir tecnología lingüística de alta calidad para todas las lenguas a fin de lograr la unidad política y económica a través de la diversidad cultural. La tecnología ayudará a derribar las barreras existentes y a construir puentes entre las lenguas de Europa. Esto requiere que todas las partes interesadas – del mundo de la política, de la investigación, la industria y la sociedad – unan sus esfuerzos de cara al futuro.

Las tecnologías lingüísticas
contribuyen a unificar Europa.
