

Gestión terminológica, corpus especializados y extracción automática de terminología en español



Sergio Rodríguez-Tapia

EDITORIAL COMARES



Sergio Rodríguez-Tapia

Gestión terminológica, corpus especializados y extracción automática de terminología en español

Granada, 2024

Colección indexada en la MLA International Bibliography desde 2005

EDITORIAL COMARES

INTERLINGUA

371

Colección fundada por:

EMILIO ORTEGA ARJONILLA y PEDRO SAN GINÉS AGUILAR

Comité Científico (Asesor):

ESPERANZA ALARCÓN NAVÍO Universidad de Granada	ÓSCAR JIMÉNEZ SERRANO Universidad de Granada
JESÚS BAIGORRI JALÓN Universidad de Salamanca	ÁNGELA LARREA ESPIRAL Universidad de Córdoba
CHRISTIAN BALLIU ISTI, Bruxelles	HELENA LOZANO Università di Trieste
LORENZO BLINI LUSPIO, Roma	MARIA JOAO MARÇALO Universidade de Évora
ANABEL BORJA ALBÍ Universitat Jaume I de Castellón	FRANCISCO MATTE BON LUSPIO, Roma
NICOLÁS A. CAMPOS PLAZA Universidad de Murcia	JAVIER MARTÍN PÁRRAGA Universidad de Córdoba
MIGUEL Á. CANDEL-MORA Universidad Politécnica de Valencia	ANTONIO RAIGÓN RODRÍGUEZ Universidad de Córdoba
ÁNGELA COLLADOS AÍS Universidad de Granada	CHELO VARGAS-SIERRA Universidad de Alicante
MIGUEL DURO MORENO Universidad de Málaga	MERCEDES VELLA RAMÍREZ Universidad de Córdoba
FRANCISCO J. GARCÍA MARCOS Universidad de Almería	ÁFRICA VIDAL CLARAMONTE Universidad de Salamanca
GLORIA GUERRERO RAMOS Universidad de Málaga	GERD WOTJAK Universidad de Leipzig
CATALINA JIMÉNEZ HURTADO Universidad de Granada	

ENVÍO DE PROPUESTAS DE PUBLICACIÓN:

Las propuestas de publicación han de ser remitidas (en archivo adjunto, con formato PDF) a alguna de las siguientes direcciones electrónicas: anabelen.martinez@uco.es, psgines@ugr.es

Antes de aceptar una obra para su publicación en la colección INTERLINGUA, ésta habrá de ser sometida a una revisión anónima por pares. Para llevarla a cabo se contará, inicialmente, con los miembros del comité científico asesor. En casos justificados, se acudiría a otros especialistas de reconocido prestigio en la materia objeto de consideración.

Los autores conocerán el resultado de la evaluación previa en un plazo no superior a 60 días. Una vez aceptada la obra para su publicación en INTERLINGUA (o integradas las modificaciones que se hiciesen constar en el resultado de la evaluación), habrán de dirigirse a la Editorial Comares para iniciar el proceso de edición.

Este volumen ha sido financiado por el grupo de investigación HUM-060 «Pensamiento, lenguas y textos: estudios teóricos, aplicados y didácticos» (IDEA-TEXT) de la Universidad de Córdoba.

Maquetación: Miriam L. Puerta

© Sergio Rodríguez-Tapia

© Editorial Comares, 2024

Polígono Juncaril • C/ Baza, parcela 208 • 18220 Albolote (Granada) • Tlf.: 958 465 382

<https://www.comares.com> • E-mail: libreriacomares@comares.com

<https://www.facebook.com/Comares> • <https://twitter.com/comareseditor>

<https://www.instagram.com/editorialcomares>

ISBN: 978-84-1369-743-7 • Depósito legal: Gr. 787/2024

Impresión y encuadernación: COMARES

Sumario

AGRADECIMIENTOS	XI
Capítulo I—LA TECNOLOGÍA EN LA LINGÜÍSTICA Y LA TERMINOLOGÍA EN LA ACTUALIDAD	1
1. LA TECNOLOGÍA EN LA GESTIÓN Y EN EL TRATAMIENTO DEL DISCURSO ESPECIALIZADO	2
2. LA GESTIÓN TERMINOLÓGICA COMO OBJETO DE ESTUDIO DE LA LINGÜÍSTICA	4
3. OBJETIVOS Y ESTRUCTURA DE ESTA OBRA	5
3.1. Criterios de selección de los recursos del catálogo	7
3.2. Criterios de selección para la descripción de los recursos del catálogo....	11
Capítulo II—LA GESTIÓN TERMINOLÓGICA Y LAS HERRAMIENTAS INFORMÁTICAS..	13
4. LAS BASES DE DATOS TERMINOLÓGICAS.....	14
4.1. Estructuración y presentación de los datos.....	18
4.1.1. <i>Organización onomasiológica</i>	18
4.1.2. <i>Organización semasiológica</i>	19
4.2. Recuperación de información: filtros y selección de información	19
4.3. Ejemplos de bases de datos terminológicas en español.....	22
4.3.1. <i>Cercaterm</i>	22
4.3.2. <i>DeCS/MeSH. Descriptores en Ciencias de la Salud</i>	24
4.3.3. <i>Dicciomed. Diccionario médico-biológico, histórico y etimológico</i>	26
4.3.4. <i>DicoAdventure</i>	28
4.3.5. <i>EuroVoc</i>	30
4.3.6. <i>IATE. Interactive Terminology for Europe</i>	32
4.3.7. <i>TERMDAT. Die Terminologiedatenbank der Bundesverwaltung</i> ...	34
4.3.8. <i>TERMIUM Plus</i>	36
5. EL PROCESO DE GESTIÓN TERMINOLÓGICA: FASES Y CRITERIOS	38
5.1. Definir el alcance y los límites del recurso	39
5.1.1. <i>Establecer la motivación, los objetivos y los destinatarios</i>	39
5.1.1.1. <i>Definir la actividad profesional en la que se utilizará prototípicamente el recurso</i>	40
5.1.1.2. <i>Definir las necesidades prototípicas de los usuarios que suelen realizar la actividad profesional</i>	41

5.1.1.3.	<i>Determinar el valor de las unidades seleccionadas para incluirse en el proyecto</i>	41
5.1.1.4.	<i>Determinar las lenguas del proyecto</i>	41
5.1.2.	<i>Establecer las características del recurso: los campos de la entrada terminológica</i>	42
5.1.2.1.	<i>Información fonético-fonológica</i>	42
5.1.2.2.	<i>Información morfológica</i>	43
5.1.2.3.	<i>Información sintáctica</i>	44
5.1.2.4.	<i>Información semántica</i>	45
5.1.2.5.	<i>Información pragmática</i>	50
5.1.2.6.	<i>Información administrativa</i>	52
5.1.2.7.	<i>Algunas reflexiones sobre los campos</i>	56
5.1.3.	<i>Determinar los recursos humanos y técnicos y planificar el proyecto</i>	56
5.2.	Elaborar el recurso.....	57
5.3.	Revisar el recurso.....	57
5.4.	Editar, actualizar y dar difusión al producto.....	58
6.	LOS SISTEMAS DE GESTIÓN TERMINOLÓGICA.....	59
6.1.	Ejemplos de sistemas de gestión terminológica.....	60
6.1.1.	<i>MultiTerm</i>	60
6.1.2.	<i>TermStar NXT</i>	62
6.1.3.	<i>WebTerm</i>	64
7.	LAS ONTOLOGÍAS COMO BASES DE CONOCIMIENTO TERMINOLÓGICO.....	65
7.1.	La ingeniería del conocimiento y la terminología.....	66
7.2.	El concepto de ontología.....	67
7.3.	La formalización en las ontologías.....	69
7.4.	Componentes de una ontología.....	72
7.5.	Aplicaciones de una ontología.....	74
Capítulo III—EL DISEÑO Y USO DE LOS CORPUS ESPECIALIZADOS.....		77
8.	MOTIVACIÓN DE LA LINGÜÍSTICA DE CORPUS.....	77
8.1.	El corpus: el objeto de la lingüística de corpus.....	78
8.2.	Consideraciones epistemológicas y ventajas instrumentales.....	79
8.3.	La relevancia de los corpus en el análisis del discurso especializado y sus aplicaciones.....	81
9.	CRITERIOS DE CLASIFICACIÓN DE LOS CORPUS.....	83
9.1.	Los corpus específicos y especializados: denominaciones y ventajas.....	89
10.	FASES Y CRITERIOS EN LA COMPILACIÓN DE CORPUS ESPECIALIZADOS.....	91
10.1.	Compilar el corpus: alcanzar la representatividad del corpus especializado.....	93
10.1.1.	<i>Criterios cualitativos de representatividad</i>	94
10.1.1.1.	<i>La fiabilidad de las fuentes</i>	94
10.1.1.2.	<i>Los límites del ámbito de especialidad</i>	96
10.1.1.3.	<i>El grado de especialización</i>	100
10.1.1.4.	<i>Diversificación de las muestras: homogeneidad y heterogeneidad</i>	101
10.1.2.	<i>Criterios cuantitativos de representatividad</i>	103
10.1.2.1.	<i>Tamaño del corpus</i>	103
10.1.2.2.	<i>Distribución de las muestras del corpus</i>	105

SUMARIO

10.2. Codificación y etiquetado.....	106
10.3. Recuperación y análisis de información. Las funciones de los gestores de corpus.....	113
11. EJEMPLOS DE CORPUS ESPECIALIZADOS EN ESPAÑOL	119
11.1. <i>Corpus de las Sexualidades de México (CSMX)</i>	119
11.2. <i>Corpus Tècnic</i>	121
11.3. <i>DIACOM-es</i>	123
 Capítulo IV—LA EXTRACCIÓN TERMINOLÓGICA	 127
12. DELIMITACIÓN CONCEPTUAL: EXTRACCIÓN TERMINOLÓGICA MANUAL Y AUTOMÁTICA	127
13. VENTAJAS DE LA EXTRACCIÓN AUTOMÁTICA DE TERMINOLOGÍA.....	128
14. APLICACIONES DE LA EXTRACCIÓN AUTOMÁTICA DE TERMINOLOGÍA.....	130
15. EXTRACCIÓN MANUAL DE TERMINOLOGÍA: LA RELATIVIDAD DEL VALOR TERMINOLÓGICO	132
16. SISTEMAS DE EXTRACCIÓN AUTOMÁTICA DE TERMINOLOGÍA.....	135
16.1. Sistemas lingüísticos de extracción automática de terminología	135
16.2. Sistemas estadísticos de extracción automática de terminología	137
16.3. Sistemas híbridos de extracción automática de terminología	138
17. CRITERIOS PARA LA EXTRACCIÓN AUTOMÁTICA DE TERMINOLOGÍA.....	139
17.1. La unicidad	139
17.1.1. <i>Identificación de la unicidad mediante método lingüístico</i> .	139
17.1.2. <i>Identificación de la unicidad mediante método estadístico</i> .	141
17.2. La terminologicidad.....	142
17.2.1. <i>Identificación de la terminologicidad mediante método estadístico</i>	142
17.2.1.1. <i>Método distribucional. Relación TF-IDF</i>	142
17.2.1.2. <i>Método contextual: algoritmo C-value</i>	143
17.2.1.3. <i>Método contextual: algoritmo NC-value</i>	143
17.2.2. <i>Identificación de la terminologicidad mediante método lingüístico</i>	144
17.2.3. <i>Identificación de la terminologicidad mediante método contrastivo</i>	144
17.3. La variación terminológica	145
17.4. La evaluación y la validación	146
18. EJEMPLOS DE SISTEMAS PARA LA EXTRACCIÓN AUTOMÁTICA DE TERMINOLOGÍA	147
18.1. <i>OneClick Terms de Sketch Engine</i>	148
18.2. <i>TermoStat Web 3.0</i>	150
 Capítulo V—RECAPITULACIÓN Y REFLEXIONES FINALES.....	 155
 Capítulo VI—REFERENCIAS BIBLIOGRÁFICAS.....	 161
 Capítulo VII—CATÁLOGO DE RECURSOS TECNOLÓGICOS SOBRE GESTIÓN TERMINOLÓGICA EN ESPAÑOL.....	 175

Agradecimientos

Este proyecto nace por una motivación personal y docente. Esta monografía ha sido fruto de un proceso de introspección que ha pretendido responder a muchas de las preguntas que me he ido haciendo durante mi etapa investigadora, así como a las cuestiones planteadas por el alumnado de Terminología de la Universidad de Córdoba a lo largo de los años. Esta obra ha sido resultado de la consulta a muchos colegas y compañeros que admiro, quienes han ofrecido su tiempo y dedicación para brindarme su opinión.

En primer lugar, agradezco a Laia Vidal e Iria Da Cunha los comentarios recibidos a la génesis de este proyecto y a una primera versión de este volumen. Gracias a sus sugerencias, este trabajo ve la luz con nueva estructura, planteamientos y preguntas. Asimismo, agradezco a Judit Freixa el hecho de facilitarme la información administrativa sobre algunos de los proyectos del IULA, que han sido integrados en este texto.

En mi universidad doy las gracias a Adela González y a María Martínez-Atienza por sus precisos comentarios acerca de la lingüística de corpus y de la fraseología respectivamente. A Eduardo J. Jacinto le agradezco especialmente la atenta lectura de partes de este volumen, pues esta obra contiene muchas de las conversaciones mantenidas desde hace siete años, porque se ha ido construyendo mediante el diálogo, el intercambio de opinión, de intereses y, especialmente, de la ilusión que compartimos por el discurso especializado. A Alfonso Zamorano le debo mi interés científico por el metalenguaje, que he pretendido plasmar en esta obra. Le agradezco todo su apoyo académico e institucional, su paciencia y sus constantes ánimos y alegrías compartidas.

A Carmen Oliva, Sergio Martínez y Yordan Apostolov les agradezco sus sugerencias y comentarios porque, con buen ojo clínico, han demostrado ser los discípulos con los que cualquier profesor querría contar. Les agradezco su atenta lectura a capítulos de este trabajo y las notas que han permitido precisar ambigüedades o confusiones, o acabar con gazapos propios de mi falta de cuidado.

En lo personal, querría agradecer las ideas de Carmen Moreno, quien me ha ayudado a reflexionar sobre mis objetos de estudio y sobre mis objetivos vitales más de lo que cree. A Rafa y Javi me gustaría agradecerles la inmensa paciencia demostrada, su incuestionable talante y su curiosidad por este mundo de «términos impronunciados y cosas difíciles que nadie entiende».

A todos ellos, doy las gracias por su tiempo e inestimable dedicación.

Sergio Rodríguez-Tapia

colección:

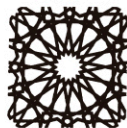
INTERLINGUA

371

Dirigida por:

Ana Belén Martínez López y Pedro San Ginés Aguilar

La gestión terminológica consiste en una serie de procedimientos de carácter técnico que involucran el empleo de la tecnología para aumentar la eficiencia en la recopilación, documentación, almacenamiento y recuperación del vocabulario especializado. Constituye, pues, un trabajo esencial para el tejido industrial, especialmente para las empresas y las instituciones. El resultado de la gestión terminológica suele ser una base de datos, donde se consignan los términos, sus usos, recomendaciones o vinculaciones a ciertos productos o empresas, aclaraciones, observaciones o definiciones, sinónimos, etc. La utilidad real de estas bases de datos terminológicas reside en que permiten consultar rápidamente información validada y contrastada, ya sea por los usuarios externos o por el personal de la empresa o institución. Este volumen presenta un estado de la cuestión donde se aborda la fundamentación teórica y metodológica en torno a la gestión terminológica, que se presenta como una de las actividades relacionadas con el tratamiento del discurso especializado. Para ello, se introduce su necesaria relación con las bases de datos terminológicas, con la lingüística de corpus y con la extracción terminológica, aplicaciones para las cuales el presente volumen recopila un catálogo de 103 recursos informáticos diferentes. El fin último de este volumen es servir como instrumento de apoyo a la investigación teórica y a la docencia universitaria de los campos relacionados con los estudios lingüísticos; en concreto, terminológicos.



COMARES
editorial

ISBN 978-84-1369-743-7



9 788413 697437