



Entreculturas 14 (2024) pp. 80-99 — ISSN: 1989-5097

Análisis cuantitativo de traducciones del inglés al español poseídas por traductores institucionales

Quantitative analysis of English-Spanish translations post-edited by institutional translators

 Miguel Vega Expósito
Universidad de Granada (España)

Recibido: 29 de septiembre de 2023

Aceptado: 13 de enero de 2024

Publicado: 27 de febrero de 2024

ABSTRACT

Human post-editing of machine translation keeps growing in the professional translation sector and is getting increasingly prominent in the academic world. However, there is limited research focused on professional practice in specific language combinations which may be used to create training materials for prospective professionals. In this article, we present our findings from quantitative analyses of texts that were translated from a variety of English-language sources into Spanish by the neural machine translation system eTranslation and post-edited by experienced translators from the Directorate-General for Translation of the European Commission. We discuss our findings regarding the use of automated suggestions, text length, and lexical richness.

KEYWORDS: post-editing, machine translation, text length, lexical richness.

RESUMEN

La práctica de la posesición humana de traducciones automáticas no deja de crecer en la traducción profesional y cobra cada día más protagonismo en el ámbito académico. Sin embargo, no existen suficientes estudios centrados en la práctica profesional en combinaciones lingüísticas concretas que permitan desarrollar materiales pedagógicos para futuros profesionales. En este artículo damos cuenta de las observaciones que hemos hecho tras diversos análisis cuantitativos de textos traducidos a la lengua española por el sistema de traducción automática neuronal eTranslation y poseídos por traductores profesionales de la Dirección General de Traducción de la Comisión Europea a partir de textos de diversa tipología en lengua inglesa. Presentamos observaciones hechas relativas al aprovechamiento de las propuestas automáticas, la extensión de los textos y la riqueza léxica.

PALABRAS CLAVE: posesición, traducción automática, extensión del texto, riqueza léxica.

1. Introducción

La práctica de la posesición humana de traducciones automáticas crece año tras año (ELIS, 2020, 2021, 2022, 2023). Frente a otras formas más tradicionales de ejercer la actividad traductora, no cabe duda de que el uso combinado de la traducción automática (TA) y la posesición en diversos contextos profesionales acelera el proceso de traducción y, por lo tanto, aumenta la productividad (Jia y Lai, 2022: 34-35). Los sistemas de TA no han dejado de prosperar, especialmente desde que se comenzaron a popularizar los sistemas neuronales a partir del año 2016.

Sin embargo, estos sistemas no llegan a alcanzar los mismos niveles de calidad y fiabilidad que los traductores profesionales (Koehn, 2020: 20). Los programas de TA neuronal, al igual que ChatGPT y otros sistemas de aprendizaje profundo, son herramientas basadas en estadísticas y probabilidades. Reproducen lo que encuentran en los datos, pero no tienen una comprensión en el sentido humano de la palabra. Mientras sigan basándose, como lo hacen hoy, en el entrenamiento con ingentes cantidades de datos de buena calidad y no en el desarrollo de capacidades humanas como la comprensión y el razonamiento, los sistemas de TA seguirán cometiendo errores que un traductor profesional no comete. Es cierto que esa tasa de errores se reduce conforme pasan los años, se usan memorias de traducción y corpus bilingües de mayor calidad para entrenar los sistemas, y se mejoran los algoritmos, pero incluso en el caso de que en algún momento la tasa de errores se aproximara a cero, la posesición humana seguiría siendo imprescindible en multitud de contextos, porque los poseedores también deben validar las propuestas automáticas y consecuentemente asumir la responsabilidad de la calidad del texto final (Trojszczak, 2022: 188-189).

La posesición es una forma de traducción diferente de los otros modos más tradicionales de practicar esta actividad y requiere algunas destrezas especiales (Gaspari et al., 2015: 334-336; Cadwell et al., 2018: 309-310; Ginovart Cid et al., 2020; Do Carmo y Moorkens, 2020; Pym y Torres-Simón, 2021: 14-15; O'Brien, 2022: 118). Por ello, el

uso de herramientas de TA y la práctica profesional de la posesición deberían ocupar un lugar más prominente en la formación de traductores profesionales en nuestro país. Esto exige que quienes nos dedicamos a la enseñanza conozcamos mejor tales tecnologías y rutinas profesionales, sepamos distinguir la posesición de otras actividades muy emparentadas como las de revisar o traducir con memorias de traducción, sigamos venciendo las resistencias que aún quedan frente a la TA (Vieira, 2020; Ragni y Vieira, 2022), y estemos convencidos de su utilidad. Todo esto es un proceso muy lento, pero afortunadamente la producción científica en estos ámbitos comienza a ser abundante y confiamos en que pronto las herramientas de TA y la posesición ocupen el lugar que merecen en la formación de futuros traductores.

Los aspectos más estudiados por ahora en relación con la posesición son la calidad, el esfuerzo, la productividad y el proceso. También existen estudios comparativos de directrices de trabajo, propuestas formativas, algunos manuales de carácter general muy útiles, etc. Sin embargo, no conocemos aún investigaciones que se hayan centrado en analizar las diferencias entre las traducciones automáticas y los correspondientes textos poseídos por profesionales con el objeto de definir y recopilar buenas prácticas. Estudios de esta índole, especialmente si están centrados en pares de lenguas concretas, permitirían dar pasos muy importantes en el terreno de la enseñanza.

La pequeña aportación que queremos hacer aquí es exponer algunos de los resultados de una investigación en la que hemos analizado posesiciones completas en lengua española elaboradas por profesionales de la Dirección General de Traducción de la Comisión Europea a partir de textos traducidos por un sistema de TA neuronal desde el inglés. Los objetivos de tal proyecto han consistido en recabar toda la información que nos pudiera aportar sobre la práctica profesional de la posesición un análisis cuantitativo y cualitativo de los textos poseídos publicados, así como de los textos intermedios involucrados en el proceso, y elaborar a partir de los resultados algunas propuestas pedagógicas. En esta publicación presentamos solamente

los resultados de los análisis cuantitativos. Queda fuera de los límites de este artículo la presentación del análisis cualitativo de los datos, la recopilación de prácticas recomendadas para la posesición completa inglés-español y los desarrollos pedagógicos.

2. Materiales y métodos

La entidad que hemos escogido para conocer mejor la práctica de la posesición es el Departamento de Lengua Española de la Dirección General de Traducción (DGT) de la Comisión Europea. En él, las traducciones automáticas se realizan desde el año 2019 con el sistema neuronal eTranslation, que es entrenado con las memorias de traducción de las propias instituciones europeas, recogidas en el conjunto de bases de datos Euramis (Oravec et al., 2019: 320), por lo que está perfectamente adaptado para la traducción de nuevos documentos de la misma institución.

Mientras los profesionales de la DGT trabajan en su traducción, lo hacen en un entorno adaptado de Trados Studio en el que pueden identificar a simple vista la procedencia de las propuestas de traducción que reciben. Si lo desean, no hay celdas vacías y nunca tienen que partir de cero, a no ser que tomen esa decisión voluntariamente. También pueden decidir libremente qué propuestas de traducción quieren usar y cómo quieren hacerlo. Si las propuestas proceden de las memorias de traducción, aparecen en color verde o naranja y los traductores tienen que hacer las adaptaciones necesarias. Si proceden de eTranslation, el color azul y una penalización de 25 puntos porcentuales les sirven de aviso. En ese caso su tarea es la de hacer una posesición completa para que el producto final alcance cotas de calidad similares a las de las traducciones humanas y sea apto para ser publicado.

Para este estudio hemos confeccionado una base de datos y diversos corpus de textos a partir de 406 archivos XLIFF y TMX cedidos por la DGT. Estos ficheros corresponden, por un lado, a 203 documentos distintos en formato XLIFF con un total de 219 643 palabras redactadas originalmente en lengua inglesa y las correspondientes 260 384 palabras de las versiones finales en lengua español

la aprobadas por los traductores de la DGT. Por otro lado, los 203 archivos TMX corresponden a los mismos documentos originales que los XLIFF y contienen las mismas 219 643 palabras en inglés, así como las correspondientes traducciones automáticas a la lengua española hechas por eTranslation, que suman un total de 258 757 palabras. Los metadatos de los archivos TMX revelan que los traductores solo utilizaron una parte de estas traducciones automáticas y desearon o utilizaron solo de forma parcial el resto de propuestas. El número total de segmentos asciende a 16 826 unidades.

Los 203 documentos de partida son de muy diversa índole y temática. Todos proceden de las instituciones europeas y se pusieron a disposición del público general a lo largo del año 2019 en diferentes medios. Se redactaron en lengua inglesa y comprenden textos legales, documentos administrativos, comunicados de prensa, textos de páginas web, etc. Los temas tratados en los textos son también muy variados: agricultura, transporte, cultura, justicia, medio ambiente, educación, tecnología, salud pública, migración y energía son solo unos pocos de ellos. La selección de los textos fue llevada a cabo por personal de la DGT sin más criterio que el de que ya estuvieran publicados o fueran de carácter no sensible y, por lo tanto, pudieran usarse con fines de investigación. Las personas que hicieron la selección sabían que nuestro estudio se centraría en la posesición, pero el hecho de que se hubiera usado o no la traducción automática al verter los textos a la lengua española no fue uno de los criterios de selección usados, pues de hecho en 67 de los documentos no hay registro de su uso.

Las versiones finales de los documentos traducidos al español recogidas en los 203 archivos XLIFF fueron elaboradas por personal propio y externo de la DGT en una proporción que desconocemos. Sí sabemos que en todos los casos se trata de versiones finales, por lo que los textos fueron validados por personal de la DGT conforme a sus exigentes directrices de calidad (Comisión Europea, 2010; European Commission, 2015, 2020). Gracias a los metadatos, también sabemos que la traducción automática se usó al menos en un 30,7 % del volumen total de palabras. Entendemos que en estos casos siempre se persiguió una

posedición completa. Destacamos en este sentido una de las mencionadas directrices:

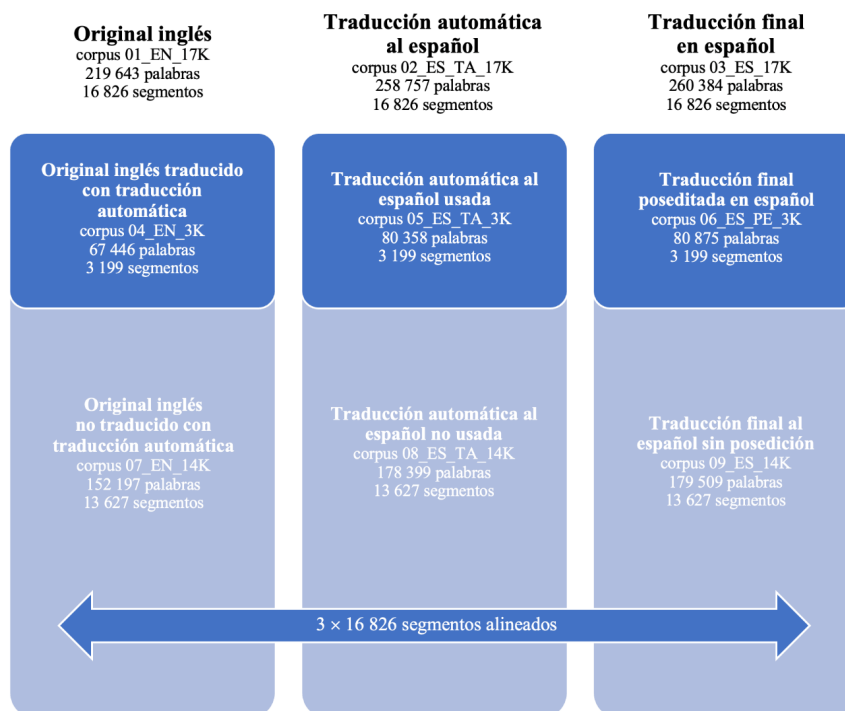
[...] la calidad en traducción depende en cierta medida de las circunstancias en las que se desenvuelve el trabajo. En muchos casos, el peticionario antepone la adecuación a la finalidad de la traducción (el principio de *fit for purpose*) a su propia calidad lingüística. De cualquier forma, el mandato de la DGT exige la excelencia en los textos que produce, especialmente los destinados al público o los de carácter legislativo (Comisión Europea, 2010: 14).

Los archivos XLIFF permiten identificar los segmentos en los que los traductores tomaron como punto de partida las traducciones automáticas, pero no registran los casos en que, después de tomar esa decisión, optaron finalmente por borrar la traducción automática completamente para traducir desde cero. Asimismo, los archivos XLIFF

tampoco permiten identificar otros usos de la traducción automática a través de funciones como el autocompletado inteligente (*autosuggest*) o la búsqueda de concordancias (Arnejšek y Unk, 2020). Por razones de protección de datos de carácter personal y por problemas técnicos durante la seudonimización de las traducciones, los archivos XLIFF no nos han permitido distinguir entre las tareas de diversos profesionales.

A partir de todos estos documentos, hemos generado una base de datos que aglutina toda la información en formato CSV. Asimismo, hemos creado diversos grupos de corpus alineados en Sketch Engine (figura 1). El primer grupo lo componen las 16 826 unidades de traducción con datos válidos extraídas de los 406 archivos procesados. Es decir, el primer corpus consta de 16 826 segmentos y 219 643 palabras originales en inglés (lo hemos denominado 01_EN_17K). El segundo consta de 16 826 segmentos alineados con los del primero y contiene 258 757 palabras traducidas automáticamente al español (02_ES_TA_17K).

Figura 1. Descripción de los nueve corpus alineados creados para este estudio.



El tercero comprende la misma cantidad de segmentos alineados y contiene las 260 384 palabras de las traducciones finales en español (03_ES_17K).

En segundo lugar, hemos creado otros tres corpus alineados usando subconjuntos de los tres primeros. Así, hemos creado uno con los 3 199 segmentos originales en los que los traductores decidieron partir de la traducción automática (04_EN_3K), otro para los mismos segmentos traducidos automáticamente al español que realmente se utilizaron (05_ES_TA_3K) y un tercero con los mismos segmentos poseídos en español (06_ES_PE_3K). Respectivamente, constan de 67 446 palabras en inglés, 80 358 palabras traducidas automáticamente al español y 80 875 palabras poseídas.

Finalmente, hemos creado un tercer grupo de corpus alineados usando los subconjuntos complementarios. Es

decir, un corpus con los 13 627 segmentos originales en inglés con un total de 152 197 palabras para las que los traductores decidieron no usar directamente la traducción automática (07_EN_14K), otro con la misma cantidad de segmentos y 178 399 palabras traducidas automáticamente al español (08_ES_TA_14K) y un último corpus con las correspondientes 179 509 palabras de las traducciones finales en español en las que no queda registro de que se haya usado la posesión (09_ES_14K).

Reproducimos en la tabla 1 un pequeño extracto del segundo grupo de corpus alineados, aquel que contiene las traducciones poseídas. Ocultamos aquí los metadatos irrelevantes para la ocasión, aunque mostramos, eso sí, una cuarta columna con los cambios que nos interesan resaltados. Hemos analizado los datos descritos hasta aquí haciendo

Tabla 1. Extracto de los corpus alineados.

Segmento origen 04_EN_3K	Traducción automática 05_ES_TA_3K	Traducción poseída 06_ES_PE_3K	Cambios resaltados
...
<i>The EU worked on a Strategy on connecting Europe and Asia, with the fight against climate change and regional cooperation at its heart.</i>	<i>La UE trabajó en una estrategia para conectar Europa y Asia, con la lucha contra el cambio climático y la cooperación regional en su centro.</i>	<i>La UE trabajó en una estrategia para conectar Europa y Asia, entre cuyos pilares se encuentran la lucha contra el cambio climático y la cooperación regional.</i>	<i>La UE trabajó en una estrategia para conectar Europa y Asia, <u>en</u>entre cuyos pilares se encuentran la lucha contra el cambio climático y la cooperación regional en su centro.</i>
<i>To develop relations to their full potential, the EU adopted a new Strategy on India and ratified the EU-Philippines Partnership and Cooperation Agreement.</i>	<i>Para desarrollar todo su potencial, la UE adoptó una nueva estrategia sobre la India y ratificó el Acuerdo de Colaboración y Cooperación entre la UE y Filipinas.</i>	<i>Para desarrollar las relaciones en todo su potencial, la UE adoptó una nueva estrategia sobre la India y ratificó el Acuerdo de Colaboración y Cooperación entre la UE y Filipinas.</i>	<i>Para desarrollar <u>las relaciones en</u> todo su potencial, la UE adoptó una nueva estrategia sobre la India y ratificó el Acuerdo de Colaboración y Cooperación entre la UE y Filipinas.</i>
<i>The EU and Thailand launched a labour dialogue to promote decent labour conditions.</i>	<i>La UE y Tailandia han iniciado un diálogo laboral para promover unas condiciones laborales dignas.</i>	<i>En el ámbito laboral, la UE y Tailandia iniciaron un diálogo para promover unas condiciones de trabajo dignas.</i>	<i>LaEn el ámbito laboral, <u>la</u> UE y Tailandia han iniciado<u>iniciaron</u> un diálogo laboral para promover unas condiciones laborales<u>de trabajo</u> dignas.</i>
...

uso de diversos métodos y herramientas de observación. En primer lugar, hemos utilizado una simple hoja de cálculo, concretamente Microsoft Excel, para cuantificar las unidades de traducción, computar los recursos usados por los profesionales y medir la extensión de los segmentos en número de palabras, así como la distancia de edición y el grado de similitud porcentual entre las traducciones automáticas y sus posesiciones. La distancia de edición o de Levenshtein calcula el número mínimo de inserciones, eliminaciones o sustituciones requeridas para pasar de un texto a otro, en nuestro caso del segmento creado por el sistema de TA al segmento final poseditado. Este dato comparativo ofrece una idea sobre el alcance de la intervención humana en las traducciones automáticas.

En segundo lugar, hemos utilizado una herramienta de análisis de textos, concretamente Sketch Engine, para tareas sencillas de recuento de palabras por casos y tipos, así como por categorías gramaticales, con el objetivo primordial de identificar diferencias numéricas manifiestas entre los textos traducidos automáticamente y los textos poseditados, sobre todo en lo referente a su riqueza léxica.

3. Resultados

3.1. Grado de utilización de la traducción automática

Un vistazo rápido al mero número de palabras y segmentos que contienen cada uno de los nueve corpus utilizados para este estudio ya nos permiten hacer algunas observaciones significativas. Las proporciones en la tabla 3 ya nos indican que, en el contexto profesional que aquí analizamos, se hace uso de la traducción automática en un porcentaje muy significativo, pero no mayoritario, del volumen de los textos. Sobre un total de 16 826 segmentos, los traductores utilizaron la propuesta de eTranslation en 3 199 de ellos, es decir, en el 19 % de los segmentos. Convertida esta cuenta en palabras, sobre un total de 219 643 palabras en lengua inglesa, la traducción automática se aprovechó en 67 446 de ellas, es decir, en un 30,7 %.

Tabla 2. Número de palabras y segmentos de todos los corpus usados.

Todos los corpus de textos				
	Palabras	%	Segmentos	%
Textos originales inglés 01_EN_17K	219 643	100 %		
Traducción automática español 02_ES_TA_17K	258 757	100 %	16 826	100 %
Versión final español 03_ES_17K	260 384	100 %		

Tabla 3. Número de palabras y segmentos de los corpus de textos en los que se usó traducción automática y posesición.

Corpus de textos poseditados				
	Palabras	%	Segmentos	%
Textos originales inglés 04_EN_3K	67 446	30,7 %		
Traducción automática español 05_ES_TA_3K	80 358	31,1 %	3 199	19,0 %
Versión poseditada español 06_ES_PE_3K	80 875	31,1 %		

Tabla 4. Número de palabras y segmentos de los corpus de textos en los que no se usó traducción automática.

Corpus de textos no poseditados				
	Palabras	%	Segmentos	%
Textos originales inglés 07_EN_14K	152 197	69,3 %		
Traducción automática español 08_ES_TA_14K	178 399	68,9 %	13 627	81,0 %
Versión final español 09_ES_14K	179 509	68,9 %		

Tabla 5. Origen de las propuestas de traducción.

	Origen de las propuestas de traducción usadas para textos originales en inglés 01_EN_17K	
	Palabras	Porcentaje
eTranslation	67 446	30,7 %
Memorias de traducción	108 562	49,4 %
Traducción humana nueva	42 346	19,3 %
Copia del original	1 289	0,6 %
Total	219 643	100,0 %

Como término de comparación, en la tabla 5 puede verse que las memorias de traducción se usaron en un 49,4 % del volumen de los textos. Se tradujeron desde cero un 19,3 % de las palabras y se transfirieron directamente desde el segmento original el 0,6 % de palabras restantes.

En otras palabras, las memorias de traducción se siguen priorizando, aunque no exclusivamente por decisión del traductor, pues recordemos que el Departamento de Español de la DGT establece que en Trados Studio las traducciones automáticas se deben penalizar con 25 puntos.

También es reseñable el hecho de que la traducción automática no parece haberse usado, al menos no de forma directa, en 67 de los 203 documentos analizados, como se puede desprender de la tabla 6. De tales documentos, en 8 casos la razón para no usar la traducción automática fue que para todos los segmentos existían propuestas en las memorias de traducción. Sin embargo, en los otros 59 documentos había un total de 1 923 segmentos y casi 40 000 palabras para los que no existían propuestas en las memorias de traducción. En esos casos los profesionales o bien tomaron la decisión consciente de prescindir por completo de la traducción automática, lo cual no es habitual en la DGT (Caddwell et al., 2016: 234), o bien la usaron de forma parcial a través de recursos como el autocompletado inteligente. Representan un 29,1 % de los documentos de nuestro estudio y un 18,1 % de las palabras.

Tabla 6. Combinaciones de recursos de traducción utilizados por documentos originales.

Recursos utilizados	Documentos	Memorias de traducción		Traducción automática poseditada		Traducción humana nueva		Copia del original		Total	
		Segm.	Palab.	Segm.	Palab.	Segm.	Palab.	Segm.	Palab.	Segm.	Palab.
Todos los recursos	64	5 458	44 887	2 155	45 624	268	2 559	511	732	8 392	93 802
Memorias y traducciones humanas nuevas	53	2 659	22 260	0	0	1 866	38 173	232	281	4 757	60 714
Memorias y traducciones automáticas poseditadas	56	1 699	27 333	955	19 667	0	0	133	166	2 787	47 166
Solo memorias de traducción	8	725	14 076	0	0	0	0	1	1	726	14 077
Solo traducciones humanas nuevas	6	0	0	0	0	57	1 439	9	9	66	1 448
Solo traducciones automáticas poseditadas	9	0	0	59	1 273	0	0	0	0	59	1 273
Traducciones automáticas poseditadas y humanas nuevas	7	0	0	30	787	9	78	0	0	39	865
Total	203	10 541	108 556	3 199	67 351	2 200	42 249	886	1 189	16 826	219 345

3.2. Extensión del texto traducido

Otra observación interesante que podemos hacer a partir de los datos generales de las tablas 3 y 4 es que, cuando el profesional hace uso de la traducción automática, la diferencia de extensión entre el original inglés y la traducción española es mayor que cuando no hace uso de esta. En la tabla 4 podemos comprobar que los textos originales en inglés tienen 152 197 palabras, mientras la versión final española tiene 179 509 palabras, un 17,9 % más. Sin embargo, en los textos en los que intervino activamente la traducción automática, las palabras pasan de 67 446 a 80 875 entre un idioma y otro, un 19,9 % más. Se trata de una diferencia de 2 puntos porcentuales, que en un volumen de textos como el que hemos usado no creemos que sea fruto del azar: los traductores usaron más palabras en español cuando se dejaron ayudar por la traducción automática. Esto se debe posiblemente a que las máquinas suelen traducir hasta los últimos pormenores del original y llegan hasta un nivel de detalle al que los traductores en otras circunstancias no le habrían dado importancia.

Ante algunas traducciones poseditadas, cabe preguntarse, por lo tanto, si el profesional lo habría trasladado de forma tan detallada al español de no haber sido por la traducción automática. El segmento «*These proposals are now under final consideration by Member States and the European Parliament*», por ejemplo, fue traducido con la ayuda de eTranslation y después de introducir varios cambios por «*Estas propuestas están siendo actualmente examinadas por los Estados miembros y el Parlamento Europeo*», cuando en otras circunstancias probablemente se habría omitido el adverbio *actualmente*, puesto que no aporta información que no esté expresa en el tiempo verbal y resta algo de fluidez a la lectura. De hecho, hemos encontrado ejemplos de omisiones similares a la que proponemos en el resto del corpus.

Abundando en esta misma idea, es interesante comprobar que, cuando los traductores tomaron como punto de partida traducciones automáticas, no solo usaron más palabras que cuando no lo hicieron, sino que además usaron algunas palabras más que el propio sistema automático al corregir sus textos. Así se puede comprobar en la tabla 3,

donde vemos que la máquina produjo 80 358 palabras y el traductor las engrosó ligeramente hasta llegar a 80 875. Una vez más, parece que el poseedor, ante la profusión de detalles que recibe en las propuestas, abunda algo más en tales detalles en lugar de aligerar el texto.

3.3. Extensión de los segmentos traducidos

En la tabla 3 llama la atención el hecho de que la proporción de los segmentos poseditados con respecto al total de segmentos sea del 19 %, mientras la proporción de palabras poseditadas supera el 30 %. Esto quiere decir que los poseedores aprovecharon más las propuestas de eTranslation cuando los segmentos eran considerablemente más largos que la media. De hecho, como se puede observar en la tabla 7, la longitud media de todos los segmentos originales del corpus bruto es de 13,0 palabras, mientras la longitud media de los segmentos originales que se tradujeron automáticamente y poseditaron es de 21,1 palabras. Sin duda, esto está relacionado con dos hechos evidentes: por un lado, la probabilidad de contar con propuestas de las memorias de traducción es mayor cuando los segmentos son cortos; por otro lado, la probabilidad de copiar originales directamente a las celdas de la traducción es mayor cuando los segmentos son muy cortos.

Tabla 7. Longitud media de segmentos según la forma de traducir.

	Longitud media de segmentos originales	
	Palabras	σ
1. Todos los segmentos originales 01_EN_17K	13,0	14,5
2. Segmentos en los que se usó TA 04_EN_3K	21,1	13,9
3. Segmentos en los que no se usó TA 07_EN_14K	11,2	14,0
4. Segmentos en los que se usaron las memorias	10,3	13,8
5. Segmentos en que se hizo una traducción humana nueva	19,2	14,1
6. Segmentos que se copiaron del original	1,3	1,1

Nada de lo anteriormente dicho explica, sin embargo, por qué la longitud media de los segmentos que se tradujeron sin ayuda es de 19,2 palabras y la de los segmentos que se poseitaron es de 21,1 palabras. La justificación de esta diferencia debe buscarse exclusivamente en aquellos documentos donde se usan tanto la traducción automática como las traducciones humanas nuevas, puesto que no tiene sentido buscarla en los documentos en los que los profesionales decidieron prescindir expresamente de una u otra forma de traducir de principio a fin. Cuando prescinden de una de ellas, está claro que tuvieron algún motivo para decidirse por una u otra forma de trabajar y lo que encontramos son segmentos de longitudes medias prácticamente idénticas: 20,6 palabras (con una desviación típica de 13,5) de promedio en los segmentos de los documentos que los profesionales decidieron traducir sin usar en ningún caso la ayuda directa de eTranslation; 20,7 palabras de promedio (con una desviación típica de 13,3) en los segmentos traducidos con ayuda de eTranslation cuando los profesionales decidieron usarlo en todos los casos de un mismo documento en los que Trados Studio no les ofreció nada mejor.

Sin embargo, si nos limitamos a observar aquellos documentos en los que los profesionales optaron por usar las dos formas de traducir –poseitar traducciones automáticas en unos casos y traducir desde cero en otros cuando no había propuestas procedentes de las memorias–, la diferencia de extensión de los segmentos se acentúa notablemente: la longitud media de los segmentos que tradujeron sin ayuda es de 9,5 palabras (13,8), mientras la longitud media de los segmentos en los que optaron por partir de la propuesta de eTranslation es de 21,2 palabras (14,1). Estos datos concuerdan perfectamente con la diferencia observada entre las líneas 2 y 5 de la tabla 7. A pesar de la alta dispersión de los datos, es evidente que esto no es fruto del azar y debe de haber alguna explicación. Es posible que esté relacionada con la expectativa de esfuerzo y calidad por parte del profesional antes de acometer la traducción de un segmento: si este consta de unas pocas palabras, será muy fácil y rápido traducirlo sin ayuda, y la calidad estará garantizada; si el segmento es largo merecerá la pena aceptar la oferta de ayuda del sistema. Esta conjetura, sin embargo, no concuerda

del todo con el grado de aprovechamiento de las propuestas automáticas en segmentos cortos, cercana al 90 %.

3.4. Aprovechamiento de las propuestas automáticas

Centrándonos en los segmentos en los que sí se hizo uso de eTranslation para después poseitar la propuesta automática, lo primero que hemos querido saber es hasta qué punto se aprovecharon tales propuestas. Para ello, hemos calculado la distancia de Levenshtein o distancia de edición entre el segmento creado por el sistema de TA y el segmento final poseitado. En otras palabras, hemos calculado para cada segmento el número mínimo de inserciones, eliminaciones o sustituciones de caracteres requeridas para pasar de la propuesta automática a la versión final y lo hemos convertido para cada caso en el grado de similitud porcentual. Los datos observados dan una idea de la elevada calidad de las propuestas de eTranslation, pues el grado medio de similitud es del 85,6 %. Dicho de otro modo, los profesionales no cometieron el error de hacer cambios excesivos (Nunziatini y Marg, 2020; Nitzke y Gros, 2021) y en promedio tuvieron que cambiar menos de un 15 % de cada uno de los segmentos traducidos automáticamente para considerar que podían ser publicados. No olvidemos que eTranslation fue entrenado, entre otras cosas, con memorias de traducción de la propia Comisión Europea.

Como muestra de lo que es una unidad de traducción en la que el grado de similitud entre la propuesta de eTranslation y el segmento final poseitado ronda el promedio hallado, reproducimos en la tabla 8 (ver página siguiente) un ejemplo extraído de nuestros corpus.

También hemos querido averiguar si el grado de aprovechamiento es notoriamente mayor o menor en los segmentos más cortos o más largos. Como se puede observar en la tabla 9, efectivamente, en los cien segmentos más cortos en los que se usó la propuesta de eTranslation, la afinidad se eleva unos puntos hasta alcanzar el 88,7 %. En otras palabras, en los segmentos más cortos los poseidores hicieron menos correcciones. No obstante, en el caso de los cien segmentos más largos, la similitud se mantiene estable en un 85,3 %.

Tabla 8. Ejemplo de una unidad de traducción con un grado de similitud del 85,6 %.

Original	Directive 2008/56/EC does not prohibit the construction of wastewater treatment plants as long as these works do not impair the achievement of good environmental status in the region.
eTranslation	La Directiva 2008/56/CE no prohíbe la construcción de plantas de tratamiento de aguas residuales en la medida en que estas obras no vayan en detrimento de la consecución de un buen estado medioambiental en la región.
Posedición	La Directiva 2008/56/CE no prohíbe la construcción de depuradoras de aguas residuales siempre que estas obras no vayan en detrimento de la consecución de un buen estado medioambiental en la región.
Cambios	<i>La Directiva 2008/56/CE no prohíbe la construcción de <u>plantas de tratamiento</u> depuradoras de aguas residuales en la medida en <u>siempre</u> que estas obras no vayan en detrimento de la consecución de un buen estado medioambiental en la región.</i>

Tabla 9. Similitud media entre traducciones automáticas y traducciones poseídas.

	Similitud media entre TA y PE	
	%	σ
Los 100 segmentos más cortos (TA de 1 a 20 caracteres)	88,7 %	24,9 %
Los restantes 2999 segmentos (TA de 21 a 369 caracteres)	85,5 %	15,6 %
Los 100 segmentos más largos (TA de 370 a 759 caracteres)	85,3 %	12,2 %
Todos los segmentos	85,6 %	15,9 %

Sabemos que los segmentos más largos son los más complicados y los que más esfuerzo cognitivo exigen por parte del traductor. Sin embargo, no observamos que ese hecho se traduzca en nuestros corpus, con respecto a segmentos de menor extensión, en un mayor o menor grado de similitud entre la versión traducida en bruto y la versión poseída. Mayor similitud habría podido interpretarse como una señal de menor atención al detalle o mayor confianza en el sistema cuando los segmentos son largos. Menor similitud habría podido interpretarse como una renuncia a tomar las propuestas del sistema de traducción automática como punto de partida o como señal de menor confianza en el sistema cuando los segmentos son largos. Sin embargo, ninguna de las dos cosas se ha observado. No olvidemos que estamos analizando el trabajo de profesionales experimentados que ejercen su labor con altos niveles de exigencia de calidad.

Como término de comparación de los datos que acabamos de interpretar, conviene resaltar que, mientras la similitud entre las propuestas de eTranslation y las poseídas finales es del 85,6 %, el parecido entre las ofertas de eTranslation y las traducciones humanas que no tomaron el sistema automático como punto de partida es del 75,7 %. Como es lógico, este último dato refleja una mayor distancia entre las traducciones automáticas y las humanas, pero el hecho de que sea tan alto sugiere que la traducción automática se aprovechó a través del autocompletado inteligente o *autosuggest*, del que no queda registro en los datos. De hecho, sabemos que es una práctica común descartar las traducciones automáticas completas y aprovecharlas parcialmente con *autosuggest*, especialmente al trabajar con segmentos largos, puesto que es una forma de trabajar que a algunas personas les resulta más cómoda por exigir menor concentración.

No deben sorprendernos otros parecidos elevados que se pueden observar en la tabla 10, como los que existen entre las proposiciones automáticas y los segmentos copiados del original (87,6 %), puesto que habitualmente son direcciones URL, códigos y nombres propios que el sistema también transfiere tal cual. Asimismo, es razonable la semejanza entre eTranslation y las memorias de traducción (85,1 %), puesto que aquel fue entrenado con estas.

Tabla 10. Similitud media entre traducciones automáticas y versiones finales.

	Similitud media entre TA y versión final	
	%	σ
1. Todos los segmentos 02_ES_TA_17K vs 03_ES_17K	84,1 %	20,6 %
2. Segmentos en los que se usó TA 05_ES_TA_3K vs 06_ES_PE_3K	85,6 %	15,9 %
3. Segmentos en los que no se usó TA 08_ES_TA_14K vs 09_ES_14K	83,7 %	21,5 %
4. Segmentos en los que se usaron las memorias	85,1 %	21,1 %
5. Segmentos en que se hizo una traducción humana nueva	75,7 %	19,0 %
6. Segmentos que se copiaron del original	87,6 %	27,6 %

Si nos volvemos a concentrar en los segmentos en los que se utilizó eTranslation para después poseer las propuestas automáticas, hemos querido observar los datos

segmentándolos según los grados de similitud. Así lo hacemos en la tabla 11.

Estos datos nos indican algunas cosas sobre cómo los profesionales aprovecharon la traducción automática en nuestro estudio. Para empezar, conviene resaltar que el número medio de palabras por segmento ronda las 24-25 cuando hablamos de afinidades de entre el 70 % y el 99 %. Sin embargo, el número medio de palabras por segmento en las traducciones automáticas que los poseedores dieron por válidas sin alteración alguna es de aproximadamente 12. Es decir, cuando los segmentos son cortos los poseedores están más dispuestos a dejar la traducción automática inalterada.

Por otro lado, el hecho de que el grueso de los segmentos poseídos se concentre en la parte alta de la tabla es muy positivo y significa que eTranslation les es de mucha utilidad a los traductores: más de la mitad de los segmentos poseídos requirieron muy pocos cambios. La información que nos ofrece la tabla 11 si la analizamos junto con los datos generales es muy interesante también en lo que nos puede revelar sobre el acierto o desacierto de los traductores cuando decidieron usar o no las propuestas de eTranslation. Le dedicamos un nuevo apartado.

Tabla 11. Número y extensión de segmentos por rangos de similitud entre traducciones automáticas y traducciones poseídas.

Similitud entre TA y PE	N.º de segmentos	% segmentos	Palabras inglés	% palabras inglés	Promedio de palabras por segmento inglés	
					Palabras	σ
100 %	681	21,3 %	8 269	12,3 %	12,1	10,5
90 % - 99 %	953	29,8 %	22 342	33,2 %	23,4	13,8
80 % - 89 %	723	22,6 %	18 115	26,9 %	25,1	12,8
70 % - 79 %	398	12,4 %	9 945	14,8 %	25,0	13,4
60 % - 69 %	221	6,9 %	4 947	7,3 %	22,4	14,4
50 % - 59 %	94	2,9 %	1 853	2,8 %	19,7	13,2
0 % - 49 %	129	4,0 %	1 880	2,8 %	14,6	13,0
Total	3 199	100,0 %	67 351	100,0 %	21,1	13,9

3.5. Decisión de usar o no usar las propuestas de traducción automática

Si volvemos la vista atrás y nos fijamos brevemente en la tabla 6, concretamente en los documentos en los que los profesionales decidieron usar la traducción automática, vemos que eso ocurrió en 136 documentos de los 203 analizados. Esos 136 documentos suman 11 275 segmentos y 143 105 palabras. En ellos usaron las memorias de traducción o copiaron los segmentos originales en 7 799 segmentos (69,2 %) y 73 117 palabras (51,1 %). Eso nos deja con 3 476 segmentos y 69 988 palabras en los que solo quedaban dos opciones: o bien usar eTranslation, o bien no usarlo y traducir desde cero.

Todo indica que los traductores decidieron usarlo en 3 199 segmentos (92 %) y 67 351 palabras (96,2 %). El hecho de que en 2 755 (86,1 %) de esos segmentos haya más de un 70 % de similitud entre la traducción automática y la final, tal y como se desprende de la tabla 11, puede interpretarse como un indicador de que los poseedores acertaron al tomar la decisión de partir de la TA como base porque les pareció suficientemente buena. La decisión de partir de eTranslation no fue tan acertada con los otros 444 segmentos (13,9 %) en que se usó el sistema. En otras palabras, si establecemos el 70 % de similitud como el límite entre los casos en que mereció y no mereció la pena usar la traducción automática como punto de partida, estos datos nos llevan a pensar que, por lo general, los poseedores profesionales saben bien cuándo deben usar la traducción automática. Las tablas 12 y 13 muestran ejemplos con similitudes del 70 % y del 55 %, respectivamente.

Debemos tratar estos datos con cautela puesto que Trados Studio registra como usos de eTranslation aquellos segmentos en los que el sistema inserta la propuesta automática en la celda correspondiente del editor incluso cuando el profesional decide posteriormente borrar a mano todas y cada una de las palabras de la propuesta (no así cuando escoge la combinación de teclas Alt + Supr o la opción de borrar el segmento de destino). Es probable que en una parte nada despreciable de los 444 supuestos errores haya ocurrido lo que acabamos de describir. Esta peculiaridad de

Tabla 12. Ejemplo de una unidad de traducción con un grado de similitud del 70 %.

Original	Public authorities will be able to access data for scrutiny and supervisory control wherever it is stored or processed in the EU.
eTranslation	Las autoridades públicas podrán acceder a los datos con fines de control y supervisión dondequiera que se almacenen o procesen en la UE.
Posedición	Las autoridades públicas podrán acceder a los datos para realizar supervisiones y controles de vigilancia dondequiera que se almacenen o traten los datos en la UE.
Cambios	<i>Las autoridades públicas podrán acceder a los datos con fines para realizar supervisiones y controles de control y supervisión <u>vigilancia</u> dondequiera que se almacenen o procesen <u>traten</u> los datos en la UE.</i>

Tabla 13. Ejemplo de una unidad de traducción con un grado de similitud del 55 %.

Original	Within three months, the Africa-Europe Alliance was already delivering on its first projects.
eTranslation	En el plazo de tres meses, la Alianza África-Europa ya estaba cumpliendo sus primeros proyectos.
Posedición	En el plazo de tres meses, los primeros proyectos de la Alianza África-Europa ya estaban dando sus frutos.
Cambios	<i>En el plazo de tres meses, los <u>primeros</u> proyectos de la Alianza África-Europa ya estaba cumpliendo <u>estaban dando sus primeros proyectos</u> frutos.</i>

Trados Studio complica la interpretación de nuestros datos en relación con la decisión de usar o no usar las propuestas de eTranslation una vez vistas por el traductor.

Veamos lo que ha ocurrido cuando la decisión de los traductores fue la de descartar las propuestas de eTranslation.

Tabla 14. Ejemplo de una unidad de traducción en la que no se usó la propuesta automática.

Original	<i>ECODISTRICT-ICT - Integrated decision support tool for retrofit and renewal towards sustainable districts</i>
Propuesta de eTranslation desechada	<i>ECODISDISTIC-TIC: herramienta de apoyo a la toma de decisiones integrada para modernizar y renovar los distritos sostenibles</i>
Traducción humana	<i>ECODISTRICT-ICT: instrumento de apoyo a la toma de decisiones integrada para modernizar y renovar los distritos sostenibles</i>
Diferencias	<i>ECODISDISTIC-TIC: herramienta</i> <i>ECODISTRICT-ICT: instrumento de apoyo a la toma de decisiones integrada para modernizar y renovar los distritos sostenibles</i>

Pues bien, en 277 (8 %) de los 3 476 segmentos optaron por la traducción humana en lugar de la automática poseída. El dato llamativo aquí es que en solamente 87 de tales segmentos la similitud resultó siendo menor del 70 %. En otras palabras, tradujeron desde cero 190 segmentos que se acabaron pareciendo mucho a las propuestas de eTranslation. Estos datos nos hacen pensar que los profesionales acertaron en su decisión de no tomar eTranslation como punto de partida solamente en 87 segmentos (31,4 %) y erraron en 190 segmentos (68,6 %). La tabla 14 muestra uno de esos 190 casos. No olvidemos, sin embargo, que existe la posibilidad de que hicieran uso de las propuestas de eTranslation a través del autocompletado inteligente y procedieran de esta forma por comodidad.

3.6. Riqueza léxica y las diferencias entre categorías gramaticales

Lo siguiente que hemos querido observar en nuestros corpus ha sido la riqueza léxica (palabras diferentes divididas entre palabras totales) de los diferentes conjuntos de textos. Queremos saber si las traducciones automáticas usan más o menos variedad de palabras que las traducciones poseídas y, en caso de que así sea, nos interesa conocer el

comportamiento de las diferentes categorías gramaticales y algunos ejemplos de las expresiones que se usan en unos casos, pero no en otros. También queremos saber si las traducciones poseídas tienen mayor o menor variedad léxica que aquellas en las que eTranslation no ha intervenido.

Para todo ello, hemos utilizado Sketch Engine, concretamente las funciones de recuento de palabras por ocurrencias o casos, por palabras diferentes o tipos léxicos, distinguiendo entre categorías gramaticales, y hemos creado listas paralelas de frecuencias de palabras. Puesto que la riqueza léxica solamente tiene sentido compararla en corpus de extensión similar, aparte de utilizar los tres corpus correspondientes a los segmentos en lengua inglesa para los que realmente se usó la traducción automática (04_EN_3K, véase la figura 1), a los segmentos traducidos al español por eTranslation (05_ES_TA_3K) y a los segmentos finales poseídos (06_ES_PE_3K), también hemos creado tres nuevos corpus de extensión similar a estos tres a partir de los textos en los que la traducción automática no intervino de forma directa. Concretamente, hemos hecho una selección aleatoria de segmentos alineados del corpus de originales ingleses con los que los traductores decidieron no usar la traducción automática (07_EN_14K), del corpus de segmentos traducidos automáticamente al español que no se tomaron como punto de partida (08_ES_TA_14K) y del corpus de traducciones finales en español en las que no se hizo uso de la posesión (09_ES_14K). Estos corpus reducidos los hemos bautizado con los nombres 10_EN_3K, 11_ES_TA_3K y 12_ES_3K, respectivamente.

El análisis de los corpus con Sketch Engine confirma lo que ya observamos con otras herramientas en el apartado 3.2.: los traductores humanos usaron al traducir del inglés al español más palabras que los sistemas automáticos y más aún cuando poseitaron. Lo interesante es que, tal y como se puede desprender de las tablas 15 y 16, algo parecido ocurre también con la riqueza léxica: los humanos usaron en nuestros corpus más palabras distintas que el sistema de TA, lo cual concuerda también con las observaciones hechas por Toral (2019: 276). Posiblemente esto se deba a que, a diferencia de los sistemas automáticos, entre los humanos en general, y por supuesto entre los traductores de la

Tabla 15. Número de palabras totales y palabras diferentes en los corpus de textos en los que se usó traducción automática y posesición (Sketch Engine).

	CORPUS DE TEXTOS CON POSEDICIÓN		
	Palabras totales	Palabras diferentes	Riqueza léxica
Textos originales inglés 04_EN_3K	65 590	7 756	11,82 %
Traducción automática español 05_ES_TA_3K	78 273	8 416	10,75 %
Versión poseeditada español 06_ES_PE_3K	78 768	8 701	11,05 %

Tabla 16. Número de palabras totales y palabras diferentes en corpus de textos similares en los que no se usó traducción automática y posesición (Sketch Engine).

	CORPUS DE TEXTOS SIN POSEDICIÓN		
	Palabras totales	Palabras diferentes	Riqueza léxica
Textos originales inglés 10_EN_3K	64 347	8 064	12,53 %
Traducción automática español 11_ES_TA_3K	75 863	8 821	11,63 %
Versión final español 12_ES_3K	76 188	9 018	11,84 %

Tabla 17. Ejemplo en el que el poseedor introduce mayor variedad léxica.

Original	<i>Artificial intelligence is already part of our everyday lives – from helping doctors make faster and more accurate medical diagnoses to assisting farmers in using fewer pesticides on crops.</i>
Posedición	<i>La inteligencia artificial ya forma parte de nuestra vida cotidiana: <u>de ayudar</u> <u>facilita</u> a los médicos <u>a hacer diagnósticos</u> <u>médicos diagnosticar con más rápidos</u> <u>rapidez</u> y <u>precisos para ayudar</u> <u>exactitud</u>, <u>ayuda</u> a los agricultores a utilizar menos plaguicidas en los cultivos, <u>y un largo</u> etcétera.</i>

lengua española en particular, se intenta evitar la repetición de vocablos dentro de una misma oración o en oraciones cercanas si con ello no se compromete la claridad del mensaje, además de que está muy bien valorada la capacidad de utilizar un vocabulario variado. Los sistemas de TA actuales buscan siempre la traducción más plausible sin importar si ello implica que dos palabras distintas de una sola oración original o de dos segmentos consecutivos se viertan a la lengua meta con la misma palabra traducida (como ocurre en el ejemplo de la tabla 17), mientras los traductores humanos hacemos un esfuerzo añadido no solo para evitar repeticiones de palabras a escasa distancia sino también para introducir variedad cuando la ocasión lo permite.

Como se desprende de las tablas 18, 19, 20 y 21, este mismo fenómeno se puede observar en varias categorías gramaticales, aunque en algunos casos las diferencias entre la riqueza léxica de los textos humanos y de los textos automáticos son más notables que en otros. Es especialmente notorio el caso de los verbos, donde vemos que tanto traductores como poseedores usaron casi un 10 % más de vocablos distintos que los sistemas automáticos. No es infrecuente encontrar en nuestro corpus que los traductores humanos sustituyeran estructuras como «el requisito de sostenibilidad de las actividades pesqueras» por «el requisito de que las actividades pesqueras sean sostenibles» o estructuras como «exista falseamiento de la competencia»

por «se falsee la competencia». Esto nos hace pensar que posiblemente los traductores y poseedores hacen de esta forma un esfuerzo por contrarrestar el estilo nominal de los textos administrativos y lo sustituyen por un estilo verbal más eficiente, más claro y menos abstracto. La tabla 22 muestra otro claro ejemplo de esta práctica.

Tras advertir esta destacable diferencia entre los números de verbos distintos que usan humanos y máquinas al traducir los mismos textos de partida, quisimos averiguar cuáles son los verbos que habían introducido los poseedores y que eTranslation no había usado ni una sola vez. Para ello, elaboramos listas paralelas de frecuencias de palabras y casamos los lemas idénticos de los corpus 05_ES_TA_3K y 06_ES_PE_3K. Algunos de esos verbos que no usó en ningún caso el sistema automático, pero sí los poseedores, son los siguientes: *atajar, aunar, atañer, escolarizar, colegir, perpetrar*. La máquina usó en su lugar expresiones más comunes como *poner fin, unir, respetar, a la escuela* o *aparecer*.

Evidentemente, el *software* usa el vocabulario que aparece en los textos con los que ha sido entrenado y escoge siempre la manera más probable (es decir, frecuente) de traducir algo de acuerdo con esos textos. Precisamente por esa razón, no es de extrañar que la lista de verbos de los que prescindió eTranslation esté compuesta de voces mucho menos comunes que sus alternativas. En otras palabras, hay vocablos y expresiones poco comunes que pueden desaparecer de los textos traducidos y, por lo tanto, existe el peligro de que tales textos se empobrezcan si los poseedores no consiguen remediarlo. Por lo pronto, es mala señal que, como veremos más abajo, la riqueza léxica del corpus de textos poseídos 06_ES_PE_3K sea menor que la del corpus de tamaño y características similares 12_ES_3K, en el que la traducción automática o bien no intervino o lo hizo solo de forma parcial.

Con los adjetivos las observaciones son similares a las de los verbos, aunque con cifras mucho más moderadas: los poseedores y traductores usaron respectivamente un 2,5 % y un 5,3 % más de adjetivos distintos que los sistemas automáticos. De nuevo, es interesante comprobar que aquellos adjetivos que no usó en ningún caso el sistema eTranslation, pero sí los poseedores, como *extensivo, ci-*

trícola, probatorio, septentrional e inveterado, son vocablos mucho menos cotidianos en la lengua española actual que las expresiones que usó la máquina en su lugar: *ampliarse, de los cítricos, aportar pruebas, norte, firme*.

En el caso de los sustantivos, también es algo mayor la riqueza léxica en los textos poseídos y traducidos directamente por humanos. El vocabulario introducido por los poseedores y pasado por alto por el sistema en este caso también confirma lo observado ya con verbos y adjetivos: *portavocía, memorando, escrutinio, banderola, psicopedagogía, cibercapacidades, agarradera* y *cualidad* son algunos de los sustantivos inexistentes en las propuestas de eTranslation y que los poseedores introdujeron en sustitución de otros más habituales, como *portavoz, memorándum, controlado, banner, psicólogo educativo, competencias digitales, empuñadura* y *calidad*.

También queremos llamar la atención sobre lo que ocurre con el número total de sustantivos. A pesar de que, tal y como se puede comprobar en las tablas 15 y 16, las traducciones al español tienen un 18-20 % más de palabras que los textos originales ingleses, se pueden constatar dos hechos interesantes en las tablas 18 y 19: por un lado, las traducciones tienen menos sustantivos en total que los originales, sean automáticas o humanas; por otro, tienen menos sustantivos todavía si pasan por las manos de poseedores o traductores humanos. Lo primero es una diferencia habitual entre los dos idiomas. Lo segundo es compatible con la idea expuesta más arriba de que los traductores y poseedores hacen un esfuerzo por contrarrestar el estilo nominal.

Algo parecido ocurre con las preposiciones, que se usan en menor número en los textos traducidos o poseídos por humanos que en las traducciones automáticas en bruto. Junto con los sustantivos, son las únicas categorías verbales en las que se da este fenómeno, pues verbos, adjetivos, adverbios, conjunciones y pronombres se usan en mayor número cuando los textos pasan por las manos de poseedores o traductores.

De hecho, resulta muy llamativo lo que ocurre con los pronombres, evidentemente no tanto por la variedad de los que se usan como por la cantidad. No sorprende que en los textos poseídos se use el pronombre *cuanto* (por

ejemplo, «habida cuenta de cuanto precede») y en los automáticos no, o que los poseedores eliminaran las propuestas automáticas de usar el pronombre *usted* (por ejemplo, «puede usted contribuir»). Sin embargo, es muy interesante observar que los poseedores introdujeron en los textos un 10 % más de pronombres cuando hicieron la revisión de las traducciones automáticas. Así, cambiaron por ejemplo «el acceso a esas aguas» por «el acceso a las mismas», «podrían ahorrar a las empresas» por «podría suponer para estas un ahorro» o sencillamente introdujeron el pronombre

allí donde la TA optó por la omisión: «Habida cuenta del interés [...] para la Unión, esta debe adoptar [...]». Este fenómeno lógicamente está relacionado con la dificultad que tienen los sistemas de traducción automática para tratar correctamente las referencias anafóricas, derivado a su vez, como sugiere Koehn (2020: 7-9), de su incapacidad para razonar, aplicar el sentido común o el conocimiento del mundo. Las inferencias de información que conlleva cualquier anáfora siguen siendo complejas para una máquina y sencillas para un redactor humano.

Tabla 18. Número de sustantivos, verbos, adjetivos y adverbios en los corpus en los que se usó traducción automática y posesición (Sketch Engine).

CORPUS DE TEXTOS CON POSEDICIÓN												
	Sustantivos			Verbos			Adjetivos			Adverbios		
	Totales	Diferentes	Riqueza	Totales	Diferentes	Riqueza	Totales	Diferentes	Riqueza	Totales	Diferentes	Riqueza
Textos originales inglés 04_EN_3K	25 304	3 185	12,6 %	8 174	786	9,6 %	5 689	943	16,6 %	1 635	261	16,0 %
TA español 05_ES_TA_3K	25 257	2 845	11,3 %	8 210	684	8,3 %	6 736	794	11,8 %	1 565	169	10,8 %
Versión poseeditada español 06_ES_PE_3K	25 151	2 867	11,4 %	8 615	750	8,7 %	6 813	814	11,9 %	1 655	177	10,7 %

Tabla 19. Número de sustantivos, verbos, adjetivos y adverbios en corpus similares en los que no se usó traducción automática y posesición (Sketch Engine).

CORPUS DE TEXTOS SIN POSEDICIÓN												
	Sustantivos			Verbos			Adjetivos			Adverbios		
	Totales	Diferentes	Riqueza	Totales	Diferentes	Riqueza	Totales	Diferentes	Riqueza	Totales	Diferentes	Riqueza
Textos originales inglés 10_EN_3K	24 856	3 076	12,4 %	7 712	809	10,5 %	5 307	946	17,8 %	1 466	264	18,0 %
TA español 11_ES_TA_3K	25 010	2 887	11,5 %	7 422	730	9,8 %	6 526	780	12,0 %	1 392	160	11,5 %
Versión final español 12_ES_3K	24 813	2 892	11,7 %	7 643	799	10,5 %	6 692	821	12,3 %	1 440	158	11,0 %

Tabla 20. Número de preposiciones, conjunciones y pronombres en los corpus en los que se usó traducción automática y posesición (Sketch Engine).

CORPUS DE TEXTOS CON POSEDICIÓN									
	Preposiciones			Conjunciones			Pronombres		
	Totales	Diferentes	Riqueza	Totales	Diferentes	Riqueza	Totales	Diferentes	Riqueza
Textos originales inglés 04_EN_3K	10 349	64	0,6 %	2 837	9	0,3 %	906	19	2,1 %
Traducción automática español 05_ES_TA_3K	17 169	21	0,1 %	3 747	14	0,4 %	1 621	38	2,3 %
Versión posesitada español 06_ES_PE_3K	17 139	21	0,1 %	3 829	13	0,3 %	1 804	42	2,3 %

Tabla 21. Número de preposiciones, conjunciones y pronombres en corpus similares en los que no se usó traducción automática y posesición (Sketch Engine).

CORPUS DE TEXTOS SIN POSEDICIÓN									
	Preposiciones			Conjunciones			Pronombres		
	Totales	Diferentes	Riqueza	Totales	Diferentes	Riqueza	Totales	Diferentes	Riqueza
Textos originales inglés 10_EN_3K	10 952	63	0,6 %	2 891	8	0,3 %	805	18	2,2 %
Traducción automática español 11_ES_TA_3K	17 183	21	0,1 %	3 936	14	0,4 %	1 876	38	2,0 %
Versión final español 12_ES_3K	17 161	21	0,1 %	4 055	14	0,4 %	2 048	38	1,9 %

Tabla 22. Ejemplo en el que el posesitor sustituye sustantivos por verbos.

Original	<i>The rules on electronic identification for electronic transactions in the internal market have entered into force offering Europeans a new freedom: to rely on a notified electronic identification means to securely access digital services provided everywhere across Europe ranging from enrolling in a foreign university, accessing electronic health records, registering a company and filing tax returns online to digitally opening a bank account.</i>
Posedición	<i>Las <u>Entraron en vigor las</u> normas sobre <u>la</u> identificación electrónica para las transacciones electrónicas en el mercado interior <u>han entrado en vigor ofreciendoque ofrecen</u> a los europeos una nueva libertad: utilizar medios de identificación electrónica notificados para acceder de forma segura a servicios digitales prestados en <u>toda-cualquier otro lugar de</u> Europa, <u>que van desde la matriculaciónpor ejemplo matricularse</u> en una universidad extranjera, <u>el accesoacceder</u> a los historiales médicos electrónicos, <u>el registro de registrar</u> una empresa y <u>la presentación de, cumplimentar</u> declaraciones de <u>impuestosrenta</u> en línea <u>parao</u> abrir una cuenta bancaria en línea.</i>

En lo que se refiere a la diferencia en riqueza léxica entre los textos poseditados, por un lado, y los textos traducidos desde cero, por otro, observamos en nuestros corpus una ligera diferencia a favor de los segundos. Es decir, cuando comparamos la densidad de palabras totales, sustantivos, verbos, adjetivos y adverbios del corpus 06_ES_PE_3K con la del corpus 12_ES_3K en las tablas 15, 16, 18 y 19, vemos que esta es en todos los casos más elevada en las traducciones humanas, aunque no de forma muy acusada. En otras palabras, parece desprenderse de los datos que, cuando el profesional no toma las propuestas automáticas como punto de partida, usa más palabras distintas que cuando sí lo hace. Ya concluimos más arriba que los traductores humanos utilizan un vocabulario más variado que los sistemas automáticos. Probablemente estos datos nos sugieren que además las propuestas automáticas tienen un efecto de facilitación (*priming*) en el comportamiento del poseedor, tal y como indica Toral (2019: 276). Es decir, condicionan su decisión de traducir de una forma u otra y provocan en última instancia que la variedad léxica de las posesiciones sea menor que la de las traducciones humanas. En otras palabras, estos datos refuerzan la idea de que el uso de la traducción automática deja una huella en la traducción final (Toral, 2019; Castilho y Resende, 2022).

4. Conclusiones

El estudio que aquí hemos presentado se enmarca en una investigación más amplia cuyos objetivos son entender en qué consiste la posesición completa a través de ejemplos reales de textos poseditados por traductores profesionales y hallar en ellos los principios generales de una buena praxis profesional. En este artículo damos cuenta de las observaciones que hemos hecho tras diversos análisis eminentemente cuantitativos de textos traducidos automáticamente a la lengua española y poseditados por traductores de la Dirección General de Traducción de la Comisión Europea a partir de textos de diversa tipología en lengua inglesa.

A partir de una muestra de 406 documentos, 16 826 segmentos, 219 643 palabras originales en inglés, 258 757

palabras traducidas con el sistema de TA neuronal eTranslation al español y 260 384 palabras de las versiones finales publicadas en español que nos facilitó la DGT, hemos observado que la traducción automática fue de utilidad a los traductores en más de un 30 % del volumen total de palabras y no fue el punto de partida directo en un 18 % del volumen de palabras (29 % de los documentos analizados), mientras el uso de las memorias de traducción ascendió prácticamente al 50 %.

En cuanto al aprovechamiento de las propuestas automáticas, en aquellos segmentos en los que los profesionales aceptaron la ayuda de eTranslation, tuvieron que cambiar menos de un 15 % del texto para considerar que podían ser publicados. Llama la atención que aprovecharon más las propuestas de eTranslation cuando los segmentos eran considerablemente más largos que la media. Por otro lado, en los segmentos más cortos los poseedores hicieron menos correcciones y dejaron la traducción automática inalterada con más frecuencia. También hemos concluido que, ante la decisión de usar o no usar las propuestas de traducción automática, los poseedores profesionales acertaron en más del 80 % de las ocasiones si establecemos el 70 % de similitud entre la traducción automática y la final como el límite entre los casos en que mereció y no mereció la pena usar la sugerencia de eTranslation como punto de partida en lugar de traducir desde cero.

En lo que respecta a la extensión de los textos, resulta interesante advertir que los traductores usaron más palabras en español cuando se dejaron ayudar por la traducción automática. En general, usaron un mayor número de verbos, adjetivos, adverbios, conjunciones y nada menos que un 11 % más de pronombres, aunque menos sustantivos y preposiciones que eTranslation. Asimismo, fue mayor la riqueza léxica de los traductores y poseedores, pues usaron más palabras distintas que el sistema de TA. El caso más acusado es el de los verbos: los profesionales usaron en torno a un 10 % más de verbos distintos que eTranslation. Es interesante notar que, cuando los profesionales no tomaron las propuestas automáticas como punto de partida, usaron un número ligeramente mayor de palabras distintas que cuando sí lo hicieron.

Creemos que es importante seguir investigando las diferencias entre las traducciones automáticas y los correspondientes textos poseditados por traductores expertos con el objeto de definir y recopilar buenas prácticas, y trasladar este aprendizaje al aula. Conviene para ello tratar de sortear limitaciones como las que nos hemos encontrado en este estudio, en el que no se han podido identificar el 100 % de los usos de la traducción automática (a través de funciones como el autocompletado inteligente o la búsqueda de concordancias) ni hemos podido distinguir entre las posesiones de diversos profesionales. Puede ser especialmente valioso realizar estudios de carácter cualitativo centrados en las modificaciones de las traducciones automáticas realizadas por poseedores experimentados.

Bibliografía

- Arnejšek, Mateja y Unk, Alenka (2020). Multidimensional assessment of the eTranslation output for English-Slovene. En André Martins, Helena Moniz, Sara Fumega, Bruno Martins, Fernando Batista, Luisa Coheur, Carla Parra, Isabel Trancoso, Marco Turchi, Arianna Bisazza, Joss Moorkens, Ana Guerberof, Mary Nurminen, Lena Marg y Mikel L. Forcada (Eds.), *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation* (pp. 383-392). European Association for Machine Translation. <https://aclanthology.org/2020.eamt-1.41>
- Cadwell, Patrick, Castilho, Sheila, O'Brien, Sharon y Mitchell, Linda (2016). Human factors in machine translation and post-editing among institutional translators. *Translation Spaces*, 5(2). <https://doi.org/10.1075/ts.5.2.04cad>
- Cadwell, Patrick, O'Brien, Sharon, y Teixeira, Carlos S. C. (2018). Resistance and accommodation: factors for the (non-) adoption of machine translation among professional translators. *Perspectives*, 26(3). <https://doi.org/10.1080/0907676X.2017.1337210>
- Carmo, Félix do y Moorkens, Joss (2020). Differentiating editing, post-editing and revision. En Maarit Koponen, Brian Mossop, Isabelle S. Robert y Giovanna Scocchera (Eds.), *Translation revision and post-editing* (pp. 35-49). Routledge.
- Castilho, Sheila y Resende, Natália (2022). Post-Editese in literary translations. *Information*, 13(2). <https://doi.org/10.3390/info13020066>
- Comisión Europea. Dirección General de Traducción (2010). *Guía del Departamento de Lengua Española*. (2010).
- ELIS (2020). *European Language Industry Survey 2020*.
- ELIS (2021). *European Language Industry Survey 2021*.
- ELIS (2022). *European Language Industry Survey 2022*.
- ELIS (2023). *European Language Industry Survey 2023*.
- European Commission, Directorate-General for Translation (2015). *DGT translation quality guidelines*.
- European Commission, Directorate-General for Translation (2020). *DGT guidelines for evaluation of outsourced translations (TRAD19)*.
- Gaspari, Federico, Almaghout, Hala y Doherty, Stephen (2015). A survey of machine translation competences: insights for translation technology educators and practitioners. *Perspectives*, 23(3). <https://doi.org/10.1080/0907676X.2014.979842>
- Ginovart Cid, Clara, Colominas, Carme y Oliver, Antoni (2020). Language industry views on the profile of the post-editor. *Translation Spaces*, 9(2), 283-313. <https://doi.org/10.1075/ts.19010.cid>
- Jia, Yanfang y Lai, Si (2022). Post-Editing metaphorical expressions: Productivity, quality, and strategies. *Journal of Foreign Languages & Cultures*, 6(2).
- Koehn, P. (2020). *Neural machine translation*. Cambridge University Press.
- Nitzke, Jean y Gros, Anne-Kathrin (2021). Preferential changes in revision and post-editing. En Maarit Koponen, Brian Mossop, Isabelle S. Robert y Giovanna Scocchera (Eds.), *Translation revision and post-editing* (pp. 21-34). Routledge.
- Nunziatini, Mara y Marg, Lena (2020). Machine translation post-editing levels: breaking away from the tradition and delivering a tailored service. En André Martins, Helena Moniz, Sara Fumega, Bruno Martins, Fernando Batista, Luisa Coheur, Carla Parra, Isabel Trancoso, Marco Turchi, Arianna Bisazza, Joss Moorkens, Ana Guerberof, Mary Nurminen, Lena Marg y Mikel L. Forcada (Eds.), *Proceedings of the 22nd Annual Conference of the European*

- Association for Machine Translation* (pp. 309-318). European Association for Machine Translation. <https://www.aclweb.org/anthology/2020.eamt-1.33>
- O'Brien, S. (2022). How to deal with errors in machine translation: postediting. En Dorothy Kenny (Ed.), *Machine translation for everyone: Empowering users in the age of artificial intelligence* (pp. 105-120). Language Science Press. https://library.oapen.org/bitstream/handle/20.500.12657/61713/1/external_content.pdf#page=117
- Oravecz, Csaba, Bontcheva, Katina, Lardilleux, Adrien, Tihanyi, László, y Eisele, Andreas (2019). eTranslation's submissions to the WMT 2019 news translation task. En Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, André Martins, Christof Monz, Matteo Negri, Aurélie Nével, Mariana Neves, Matt Post, Marco Turchi y Karin Verspoor (Eds.), *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)* (pp. 320-326). Association for Computational Linguistics. <https://doi.org/10.18653/v1/W19-5334>
- Pym, Anthony y Torres-Simón, Ester (2021). Is automation changing the translation profession? *International Journal of the Sociology of Language*, 2021(270), 39-57. <https://doi.org/10.1515/ijsl-2020-0015>
- Ragni, Valentina y Vieira, Lucas Nunes (2022). What has changed with neural machine translation? A critical review of human factors. *Perspectives*, 30(1). <https://doi.org/10.1080/0907676X.2021.1889005>
- Toral, Antonio (2019). Post-editeese: an exacerbated translationese. En Mikel Forcada, Andy Way, Barry Haddow y Rico Sennrich (Eds.), *Proceedings of machine translation summit XVII: research track* (pp. 273-281). European Association for Machine Translation. <https://aclanthology.org/W19-6627>
- Trojszczak, Marcin (2022). Translator training meets machine translation – Selected challenges. En Barbara Lewandowska-Tomaszczyk y Marcin Trojszczak (Eds.), *Language use, education, and professional contexts* (pp. 179-192). Springer International Publishing. https://doi.org/10.1007/978-3-030-96095-7_11
- Vieira, Lucas Nunes (2020). Automation anxiety and translators. *Translation Studies*, 13(1).