

# Guías de anotación para “Generación de descripciones de imágenes mediante métodos de Deep Learning”

*Autor: Carlos García Hernán*

## Introducción

El principal objetivo de este documento es el de construir las guías de anotación necesarias para llevar a cabo la curación o naturalización de las traducciones automáticas de las descripciones del dataset MS-COCO.

Estas descripciones serán usadas junto las imágenes para alimentar modelos de *Image Captioning*, que serán capaces de generar una descripción a partir de una imagen dada.

La decisión de construir estas guías y su formación se apoyan en los siguientes tres pilares:

- **El criterio original de anotación del conjunto de datos MS-COCO:** Tratando de respetarse en el que será el nuevo dataset las instrucciones y conceptos que recibieron los anotadores originales [1,2]. *Realmente esto no es para el TFM todavía sino para los anotadores, que van a construir el dataset.*
- **Otros trabajos de adaptación de MS-COCO a diferentes lenguas:** Estableciendo estos trabajos como referentes a seguir en el nuestro [3,4,5].
- **Procedimientos y consejos de corrección de traducciones automáticas:** Apoyando nuestras guías en trabajos de post-edición de especialistas [6].

Al delimitar unos procedimientos comunes se busca reducir en lo posible el sesgo humano que conlleva el uso de anotadores no expertos como es el caso.

## Criterio de anotación original del dataset MS-COCO

El dataset MS-COCO es un conjunto de datos creado por Microsoft, en el cual se almacenan un gran número de imágenes cotidianas (coches, personas, comida, etc). Con este dataset, tal como se puede observar en [1], se tiene como objetivo el avanzar en el estado del arte de múltiples tareas de visión artificial, como son la clasificación de imágenes, reconocimiento de objetos, segmentación y generación de descripciones. Este último objetivo es el relevante para nosotros, ya que utilizaremos este conjunto de datos, para entrenar, validar y testear diferentes algoritmos que precisamente tienen como fin generar descripciones automáticas.

Todas las anotaciones de MS-COCO han sido generadas por anotadores noveles humanos, los cuales han recibido una serie de instrucciones que podemos encontrar en [2] y a continuación:

- Describir todas las partes importantes de una escena.
- No comenzar una frase con "there is".
- No describir detalles sin importancia.

- No describir cosas que pueden suceder en el pasado o en el futuro.
- No describir lo que una persona puede decir.
- No dar nombres propios.
- Las frases deben contener al menos 8 palabras.

Recordemos que estas anotaciones fueron originalmente creadas en inglés, la adaptación de estas normas no cambia en su análogo español, exceptuando la segunda norma, que como es lógico cambiará a:

- No comenzar una frase con un “Hay”.

Dado que las anotaciones en español que se van a generar parten de un traductor automático que recibe como inputs las anotaciones originales, estas nuevas anotaciones no deberían incumplir dichas normas, pero debido a que las originales fueron anotadas por personas inexpertas y no se realizó ningún tipo de revisión, es posible encontrar anotaciones que incumplan una o varias de las instrucciones. Es por este motivo por el que, a la hora de realizar la curación de las descripciones, se debe reparar este tipo de errores.

## Trabajos similares de conversión del dataset MS-COCO a otros idiomas

En [3] se tiene como objetivo el construir un conjunto de datos para *image captioning* en japonés, ya que existen recursos muy limitados en este idioma. En dicho trabajo se pensó como primera aproximación el usar Google Translate, pero las traducciones resultantes eran demasiado literales y por lo tanto poco naturales. Por lo que se optó por anotar de nuevo las imágenes de MS-COCO. Para ello 2100 anotadores procedieron a generar 820.310 descripciones en japonés, respetando las proporciones de 5 descripciones por imagen del dataset original.

Las instrucciones dadas a esos nuevos anotadores son las mismas que las recibidas en [2] a los anotadores originales, solo que ciertas instrucciones fueron adaptadas a las peculiaridades del idioma.

En [4] encontramos otra aproximación, en este caso, para adaptar MS-COCO al chino. Para ello se utilizaron 39 personas para anotar y donde de nuevo, se especificaron una serie de instrucciones en la línea de lo mostrado en [2]. A diferencia de [3] en este caso, se realizaba un muestreo periódico de las nuevas descripciones por parte de un tribunal especialista evaluador que revisaba que estas anotaciones mantenían cierto umbral de calidad.

Tanto en [4] como en [5], se aprovecha el corpus cruzado resultante chino-inglés para hacer un ejercicio de *transfer learning*, donde primero se entrena al modelo en inglés, para posteriormente extraer las capas superficiales de la parte lingüística del modelo e introducir capas nuevas y realizar el mismo entrenamiento, pero con las descripciones en chino. De

esta forma se consigue aprovechar ambos corpus en distintos idiomas, encontrando en este proceso una mejora considerable a entrenar el modelo en un solo idioma.

En el estudio de trabajos similares al pretendido en este proyecto podemos establecer dos conclusiones:

- La creación de un nuevo corpus desde cero conlleva un trabajo humano enorme, que no es razonable para los objetivos de este TFM. Por lo que se va a realizar un trabajo previo de traducción automática del que partir. Para de esta forma se conseguir mantener la variabilidad de descripciones del dataset original, al mantener el contenido semántico de las descripciones, que no sería posible con un número bajo de anotadores. Otro de los motivos es el de reducir esfuerzo, al ser en un principio menor el trabajo de naturalización de descripciones a la creación de cero de estas.
- El uso del dataset tanto en inglés como en español para la potenciación del aprendizaje del modelo.

## Creación de las guías de naturalización de traducciones automáticas

Hasta ahora hemos establecido que se traducen las descripciones originales por los motivos expuestos en las dos secciones anteriores. Pero al ser realizadas de forma automática, las nuevas descripciones se encuentran limitadas por la tecnología actual, pudiendo ser demasiado literales y por lo tanto poco naturales, causando un peor rendimiento del modelo. Por esto se ha decidido someter a las traducciones a un ejercicio de naturalización por parte de revisores humanos.

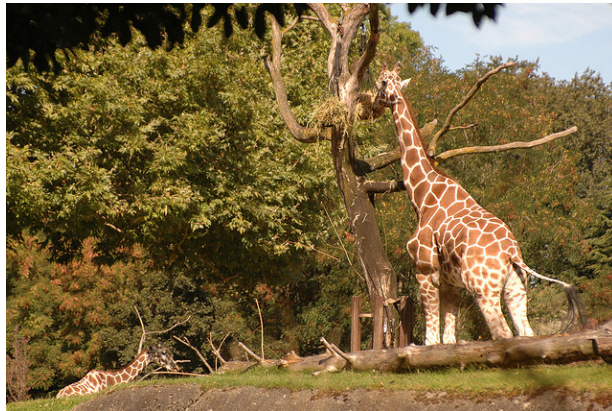
La primera aproximación a la post-edición de traducciones automáticas la recibimos desde [6], donde se recopilan diferentes consejos a la hora de realizar dicha actividad, de los cuales se han seleccionado las más convenientes para nuestro objetivo:

- **Tener guías de errores:** Se debe facilitar a los anotadores un listado con los errores comunes y como tratarlos.
- **La regla “5-10 second evaluation”:** El anotador debe ser capaz de decidir si una traducción debe ser realizada entre los primeros 5-10 segundos.

Tal como [6] recomienda, vamos a establecer una guía con errores comunes que encontraremos en nuestro dataset. Para ello se han revisado las primeras 100 descripciones traducidas automáticamente y recopilados errores encontrados en dicho proceso y la forma en la que deben ser tratados:

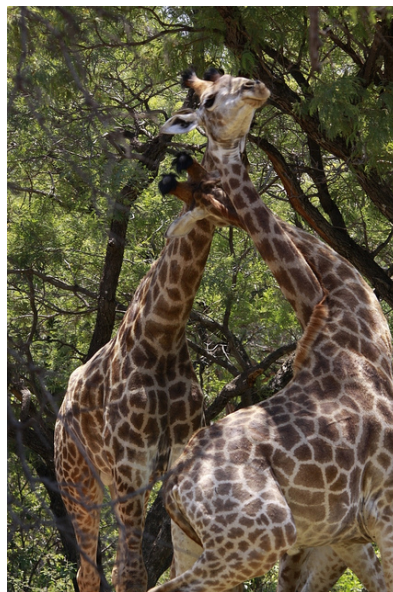
- 1) **Tratar de mantener el sentido de la descripción original:** Uno de los principales motivos de traducir las descripciones en lugar de crearlas de nuevo es para mantener la variabilidad [con variabilidad me refiero a la diversidad de descripciones que aportan los diferentes

puntos de vista de las personas, donde una reanotación con menos anotadores eliminaría dicha diversidad] del conjunto de datos original, por lo que es contraproducente cambiar el significado original. Podemos ver un ejemplo válido y otro que no:



**TA:** Una jirafa comiendo comida de la copa de un árbol  
**Curación correcta:** Una jirafa comiendo de la copa de un árbol  
**Curación incorrecta:** Una jirafa de pie cerca de un árbol

- 2) **No se deben mantener errores originales:** Hay casos en los que está claro que la descripción original es errónea (tanto en inglés como en la traducción). En estos casos es lógico reparar dicho error, aunque se pierda el sentido original, ya que el beneficio es mayor.



**TA:** Dos cebras parecen estar abrazándose en la naturaleza  
**Curación correcta:** Dos jirafas parecen estar abrazándose en la naturaleza

No solo pueden ser errores semánticos, esta norma se extiende a errores de forma como son descripciones todas en mayúsculas, errores gramaticales, léxicos o incumplir algunas de las instrucciones originales como empezar la descripción con “Hay”.

No se debe abusar de esta norma, ya que hay que estar totalmente seguro de que la descripción original está equivocada para no aumentar en la revisión los errores del dataset.

- 3) **Es posible mantener palabras originales:** Un error común es que los traductores automáticos traduzcan palabras que normalmente son usadas en su versión inglesa. Ejemplos son parking, skate, manager, etc. Si el uso de la palabra en inglés está muy extendido en nuestro idioma es lícito mantenerla en la traducción.
- 4) **Reparar errores de concordancia:** Cuando el género y/o número no concuerdan entre el sujeto y el verbo o entre el adjetivo y sustantivo.
- 5) **Evitar las repeticiones de palabras:** O con la misma raíz en una misma descripción.

**TA:** Motocicleta estacionada en el estacionamiento de asfalto.

**Curación correcta:** Motocicleta estacionada en un parking.

- 6) **Orden de palabras artificial:** Un elemento muy común en los traductores automáticos.

**TA:** Una cebra pastando en la hierba en un campo abierto y verde.

**Curación correcta:** Una cebra pastando hierba verde en campo abierto

- 7) **Uso incorrecto de las preposiciones:**

**TA:** Un poste con un reloj que marca las 6:10 por un coche blanco.

**Curación correcta:** Un poste con un reloj que marca las 6:10 cerca de un coche blanco.

- 8) **Errores en los tiempos verbales.**

- 9) **Exceso de palabras:** En ocasiones la descripción tiene un exceso de información que hace perder naturalidad a la descripción. Es posible eliminar aquellos elementos no determinantes para la descripción.

**TA:** Una cebrá pastando en una exuberante hierba verde en un campo

**Curación correcta:** Una cebrá pastando hierba en un campo

- 10) **Elementos inadecuados:** Es lícito eliminar elementos no apropiados como opiniones subjetivas del anotador original.

**TA:** Son valientes por cabalgar en la selva con esos elefantes.

**Curación correcta:** Gente cabalgando en elefantes por la selva.

Estos son los errores más comunes encontrados al corregir las primeras 100 traducciones automáticas producidas con la API de DeepL, pero es muy posible que aparezcan nuevos tipos de errores. La norma más importante que el anotador debe tener en mente es la siguiente:

**Se debe tratar de crear la descripción más fiel a la traducción original que un humano pudiese escribir.**

## Referencias

- [1] Lin, T. Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., ... & Zitnick, C. L. (2014, September). Microsoft coco: Common objects in context. In *European conference on computer vision* (pp. 740-755). Springer, Cham.
- [2] Chen, X., Fang, H., Lin, T. Y., Vedantam, R., Gupta, S., Dollár, P., & Zitnick, C. L. (2015). Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*.
- [3] Yoshikawa, Y., Shigeto, Y., & Takeuchi, A. (2017). Stair captions: Constructing a large-scale japanese image caption dataset. *arXiv preprint arXiv:1705.00823*.
- [4] Li, X., Xu, C., Wang, X., Lan, W., Jia, Z., Yang, G., & Xu, J. (2019). COCO-CN for Cross-Lingual Image Tagging, Captioning, and Retrieval. *IEEE Transactions on Multimedia*, 21(9), 2347-2360.
- [5] Miyazaki, T., & Shimizu, N. (2016, August). Cross-lingual image caption generation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (pp. 1780-1790).
- [6] Garrido, Y. F. (2016). Posedición: entre la productividad y la calidad. *redit: Revista electrónica de didáctica de la traducción y la interpretación*, (10), 22-42.