

# COMET: Un Marco Neural para la Evaluación de MT

Ricardo Rei

Craig Stewart

Ana C Farinha

Alon Lavie

Unbabel AI

{ ricardo.rei, craig.stewart, catarina.farinha, alon.lavie

} @unbabel.com

## Resumen

Presentamos COMET, un marco neural para entrenar modelos de evaluación de traducción automática multilingüe que obtiene nuevos niveles de estado del arte. Nuestro marco aprovecha los recientes avances en el modelado de lenguaje preentrenado multilingüe resultando en modelos de evaluación de MT altamente multilingües y adaptables que explotan información tanto de la entrada de origen como de una traducción de referencia en el idioma objetivo con el fin de predecir con mayor precisión la calidad de MT. Para mostrar nuestro marco, entrenamos tres modelos con diferentes tipos de juicios humanos: Evaluaciones Directas, Tasa de Edición de Traducción Mediada por Humanos y Métricas de Calidad Multidimensional. En la tarea compartida de traducción automática (MT) se han estado abor- dando los desafíos principales para la evaluación de MT que la traducción automática (MT) se han estado abor- dando los desafíos principales para la evaluación de MT que the-art performance on the WMT 2019 Metrics shared task and demonstrate robustness to high-performing systems.

## 1 Introduction

Historically, metrics for evaluating the quality of machine translation (MT) have relied on assessing the similarity between a hypothesis generated by MT and a reference translation generated by a human in the target language. The traditional metrics have focused on basic characteristics at the lexical level, such as the number of n-grams coincident between the MT hypothesis and the reference translation. Metrics such as BLEU (Papineni et al,2002) and METEOR (Graham et al,2013), Tasa de Edición de Traducción Mediada por Humanos (Lavie y Denkowski ,2009) siguen siendo populares debido a sus métricas compatibles con el marco de Métricas de Calidad Multidimensional (Lommel et al,2014).

Los enfoques neuronales modernos para MT resultan en una calidad de traducción mucho más alta que a menudo se logra con enfoques basados en reglas. Por esta razón, se ha vuelto cada vez más evidente que ya no podemos confiar en métricas como BLEU (Papineni et al,2002) y METEOR (Graham et al,2013), Tasa de Edición de Traducción Mediada por Humanos (Lavie y Denkowski ,2009) para proporcionar una estimación precisa de la calidad de MT (Barrault et al,2019).

Mientras que un creciente interés de investigación en métodos para entrenar modelos y sistemas de MT ha resultado en una reciente y dramática mejora en la calidad de la evaluación de MT, se ha quedado atrás. La comunidad de investigación de MT todavía depende en gran medida de métricas que no ha surgido ningún nuevo estándar ampliamente adoptado. En 2019, la Tarea Compartida de Traducción de Noticias recibió un total de 153 envíos de sistemas de MT (Barrault et al,2019). La Tarea Compartida de Noticias del mismo año vio solo 24 sub-misiones, casi la mitad de las cuales fueron participantes en la Tarea Compartida de Estimación de Calidad, adaptada como la Tarea Compartida de Estimación de Calidad (Ma et al,2019).

Los hallazgos de la tarea mencionada anteriormente desafiaron los supuestos de las métricas tradicionales para la evaluación de MT que se han estado abor- dando los desafíos principales para la evaluación de MT que Es decir, que las métricas tradicionales por correlacionarse con precisión con el juicio humano a nivel de segmento y no logran diferenciar adecuadamente los sistemas de MT de mayor rendimiento.

En este trabajo, presentamos COMET, un marco basado en PyTorch para entrenar modelos de evaluación de MT altamente multilingües y adaptables que pueden funcionar como métricas. Nuestro marco se beneficia de los recientes avances en el modelado de lenguaje cruzado (Artetxe y Schwenk ,2019; Devlin et al,2019; Conneau y Lample ,2019; Conneau et al,2019) para generar estimaciones de predicción de juicios humanos como Evaluaciones Directas (Lavie y Denkowski ,2009) y Tasa de Edición de Traducción Mediada por Humanos (Lavie y Denkowski ,2009) y un medio para evaluar los sistemas de MT de manera compatible con el marco de Métricas de Calidad Multidimensional (Lommel et al,2014).

Inspirados por trabajos recientes sobre Estimación de Calidad de Traducción (QEST) que demostraron que es posible lograr altos niveles de correlación con juicios humanos incluso sin una traducción de referencia (Fenici et al. ,2019), proponemos un enfoque novedoso para incorporar

Optimizado para la Evaluación de Traducción.

ing la entrada en idioma fuente en nuestros modelos de evaluación MT. Tradicionalmente, solo los modelos QE han hecho uso de la entrada de la fuente, mientras que las métricas CDA y WMT se basan en cambio en la traducción de referencia. Como [Terahashi et al \(2020\)](#), mostramos que el uso de un espacio de incrustación multilingüe para aprovechar la información de las tres entradas y demostrar el valor añadido por la fuente como entrada a nuestros modelos de evaluación MT.

Para ilustrar la eficacia y flexibilidad de el COMETmarco, entrenamos tres modelos que estiman diferentes tipos de juicios humanos y muestran un progreso prometedor hacia una mejor correlación a nivel de segmento y robustez a alta-calidad MT.

Vamos a liberar tanto el COMETmarco y los modelos de evaluación MT entrenados descritos en este documento a la comunidad de investigación tras su publicación.

## 2 Arquitecturas de Modelos

Los juicios humanos sobre la calidad de MT suelen venir en forma de puntuaciones a nivel de segmento, MQM y HTER. Para DA, es práctica común convertir las puntuaciones en clasificaciones relativas, el número de anotaciones por segmento es limitado ([Bojar et al, 2017b; Ma et al., 2018, 2019](#)). Esto significa que, para dos hipótesis  $h_i$  y  $h_j$  de la misma fuente  $s$ , si la puntuación DA asignada a  $h_i$  es mayor que la puntuación asignada a  $h_j$ , se considera como una hipótesis "mejor". Para abarcar estas diferencias, nuestro marco admite dos distintas arquitecturas: El modelo Estimator y el modelo de Clasificación de Traducción. La diferencia entre ellos es el objetivo de entrenamiento. Mientras que el Estimator se entrena para regresar directamente en una puntuación de calidad, el modelo de Clasificación de Traducción se entrena para minimizar la distancia entre una hipótesis "mejor" y tanto su referencia correspondiente como su fuente original. Ambos modelos están compuestos por un codificador cruzado de lenguaje y una capa de agrupación.

### 2.1 Codificador Cruzado de Lenguaje

El bloque de construcción principal de todos los modelos en nuestro marco es un modelo cruzado de lenguaje preentrenado, como BERT multilingüe ([Devlin et al., 2019](#)), XLM ([Conneau y Lample, 2019](#)) o XLM-RoBERTa ([Conneau et al, 2019](#)). Estos modelos contienen varias capas de codificador transformado.

<sup>2</sup>En la Tarea de Métricas Compartidas WMT, si la diferencia es entre las puntuaciones DA no es mayor a 25 puntos, esos segmentos son excluidos de los Datos RR.

Truuir tokens enmascarados al descubrir la relación entre esos tokens y los tokens en el idioma de destino. Cuando se entrenan con datos de múltiples idiomas, este objetivo preentrenado ha demostrado ser altamente efectivo en tareas interlingüística como la clasificación de documentos y la inferencia de lenguaje natural ([Conneau et al, 2019](#)), generando bien a idiomas y guiones no vistos ([Bres et al., 2019](#)). Para los experimentos en este documento, confiamos en XLM-RoBERTa (base) como nuestro codificador.

Dada una secuencia de entrada  $x_1, \dots, x_n$ , el codificador produce una incrustación para cada token  $x_j$  y cada capa  $k \in \{0, 1, \dots, K\}$ . En nuestro marco, aplicamos este proceso a la fuente, hipótesis MT, y referencia para mapearlos a un espacio de características compartido.

### 2.2 Capa de Agrupación

Los incrustaciones generadas por la última capa de los codificadores preentrenados suelen utilizarse para el ajuste de modelos a nuevas tareas. Sin embargo, [Tembargo et al, 2019](#), mostró que diferentes capas dentro de la red pueden capturar información lingüística que es relevante para diferentes tareas posteriores. En el caso de la evaluación [Mañg et al, 2020](#)) demostró que diferentes capas pueden lograr diferentes niveles de correlación y que utilizar solo la última capa a menudo resulta en un rendimiento inferior. En este trabajo, nosotros utilizamos el enfoque descrito en [Peters et al. \(2018\)](#) y agrupamos información de las capas de codificación más una única incrustación para cada token  $x_j$ , utilizando un mecanismo de atención por capas. Esta incrustación se calcula entonces como:

$$e_{x_j} = \mu \sum_{k=0}^K \alpha_k x_j^{(k)} \quad (1)$$

donde  $\mu$  es un coeficiente de peso entre 0 y 1,  $x_j^{(k)}$  corresponde al vector de incrustaciones de capa  $k$  para el token  $x_j$ ,  $\alpha = \text{softmax}([\alpha^{(1)}, \alpha^{(2)}, \dots, \alpha^{(K)}])$  es un vector correspondiente a los pesos de cada capa para evitar el sobreajuste a la información contenida en la última capa.  $\alpha^{(i)}$  es un vector de probabilidad con una suma de 1. Finalmente, como [Reimers y Gurevych, 2019](#)), aplicamos un promedio de agrupación a la palabra resultante para derivar una incrustación de oración para cada segmento.

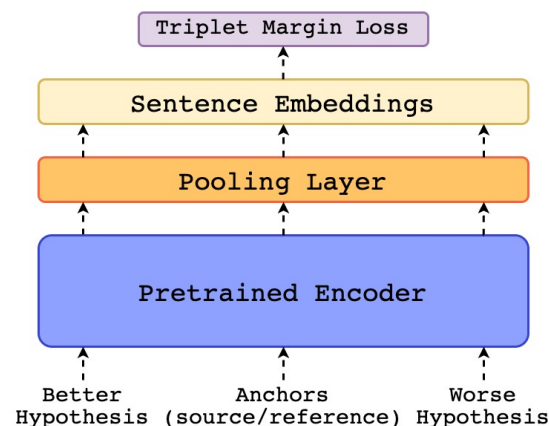
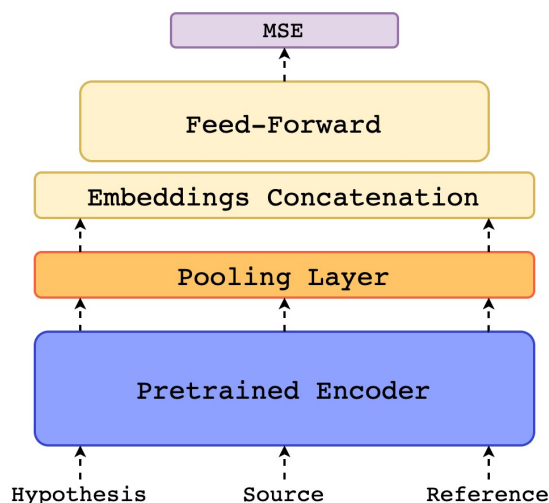


Figura 1: Arquitectura del modelo estimador. La fuente, la hipótesis y la referencia se codifican independientemente usando un codificador cruzado preentrenado. El resultado son incrustaciones de palabras que luego se pasan a través de segmentos de codificación independiente para crear una incrustación de oración para cada segmento. Finalmente, las incrustaciones de oraciones resultantes se combinan y concatenan en un solo vector que se pasa a un regresor de avance. Todo el modelo es entrenado minimizando el Error Cuadrático Medio (MSE).

Figura 2: Arquitectura del modelo de Clasificación de Traducción. El modelo recibe 4 segmentos: la fuente, la referencia, una hipótesis "mejor" y una "peor". Los segmentos se codifican independientemente usando un codificador preentrenado y una capa de agrupación en el nivel superior. Finalmente, utilizando el margen de pérdida de margen de triplet (Schroff et al., 2015) optimizamos el espacio de incrustación resultante para minimizar la distancia entre la hipótesis "mejor" y las "anclas" (fuente y referencia).

## 2.3 Modelo Estimador

Dado un incrustación de oración dimensional para la fuente, la hipótesis y la referencia, adoptamos el enfoque propuesto en RUSKIMANAKA et al. (2018) y extraemos las siguientes características combinadas:

- Producto fuente elemento por elemento:
- Producto de referencia elemento por elemento;
- Diferencia absoluta de la fuente elemento por elemento:  $|h \cdot s|$
- Diferencia absoluta de referencia elemento por elemento:  $|h \cdot r|$

Luego se concatenan estas características combinadas a la incrustación de referencia hipótesis embebida  $h$  en un solo vector  $x = [h; r; h \cdot s; h \cdot r; |h \cdot s|; |h \cdot r|]$  que sirve como entrada para un regresor de avance rápido. La fuerza de estas características radica en destacar las diferencias entre las incrustaciones en el espacio de características semánticas.

Luego se entrena el modelo para minimizar la media del error cuadrado entre las puntuaciones predichas y las evaluaciones de calidad (DA, HTER o MQM). Figura 1 ilustra la arquitectura propuesta.

Nota que elegimos no incluir la fuente cruda incrustación en nuestra entrada concatenada. Temprano en la experimentación reveló que el valor añadido por la incrustación de la fuente como características de entrada al regresor fue insignificante en el mejor de los casos. Una vez que nuestro modelo de estimador HTER entrenado con el vector  $x = [h; s; r; h \cdot s; h \cdot r; |h \cdot s|; |h \cdot r|]$  como entrada para el avance solo logra impulsar el rendimiento a nivel de segmento en 8 de los 18 pares de idiomas descritos en la sección 6.1. En la continuación y el la mejora promedio en el Tau de Kendall en esos conjuntos fue de +0.0009. Como se señaló en (Zhang et al., 2020), mientras que los modelos preentrenados multilingües son a menudo útiles para varios idiomas, el espacio de características entre los idiomas está mal alineado. Sobre esta base decidimos a favor de excluir la incrustación de origen sobre la intuición de que la información más importante proviene de la incrustación de referencia y reducir el espacio de características permitiría al modelo enfocarse más en la información relevante. Esto no sin embargo, niega el valor general de la fuente para nuestro modelo; donde incluimos características de combinación tales como  $|h \cdot s|$  notamos ganancias en correlación como se explora más adelante en la continuación.

## 2.4 Modelo de Clasificación de Traducción

Nuestro modelo de Clasificación de Traducción (Figura 2) recibe como entrada una tupla  $(s, h^+, h^-, r)$  donde  $h^+$  denota una hipótesis que fue clasificada más alta que otra hipótesis  $h^-$ . Luego pasamos a través de nuestro codificador y capa de agrupación cruzada para tener una incrustación de oración para cada segmento de oración. Finalmente, usando las incrustaciones  $(s, h^+, h^-, r)$ , calculamos la pérdida de margen  $L(s, h^+, h^-, r)$  en relación con la fuente y la referencia:

$$L(\cdot) = L(s, h^+, h^-, r) + L(r, h^+, h^-, r) \quad (2)$$

donde:

$$L(s, h^+, h^-, r) = \max\{0, d(s, h^+) - d(s, h^-) + \cdot\} \quad (3)$$

$$L(r, h^+, h^-, r) = \max\{0, d(r, h^+) - d(r, h^-) + \cdot\} \quad (4)$$

$d(u, v)$  denota la distancia euclidiana entre  $u$  y  $v$  y  $\cdot$  es un margen. Por lo tanto, durante el entrenamiento el modelo optimiza el espacio de incrustación para que la distancia entre los anclajes  $(s, r)$  y el "peor" hipótesis  $h^-$  es mayor al menos que la distancia entre los anclajes y la hipótesis "mejor"  $h^+$ .

Durante la inferencia, el modelo descrito recibe un trío  $(s, h, r)$  con solo una hipótesis. La puntuación de calidad asignada es la armónica entre la distancia a la fuente  $d(s, h)$  y la distancia a la referencia  $d(r, h)$ :

$$f(s, h, r) = \frac{2 \times d(r, h) \times d(s, h)}{d(r, h) + d(s, h)} \quad (5)$$

Finalmente, convertimos la distancia resultante en un puntaje asignado  $h_i$  a  $r$  estrictamente mayor que el puntaje de similitud limitado entre 0 y 1 como si-  
 gue:

$$f(s, h, r) = \frac{1}{1 + f(s, h, r)} \quad (6)$$

## 3 Corpora

Para demostrar la efectividad de nuestras descritas arquitecturas de modelo (3), entrenamos tres MT modelos de evaluación donde cada modelo apunta a un diferente tipo de juicio humano. Para entrenar estos modelos, utilizamos datos de tres corpora diferentes: el corpus QT21, el DA RR del WMT Metrics tarea compartida (2017 a 2019) y un propietario Corpus anotado MQM.

## 3.1 El corpus QT21

El corpus QT21 es públicamente disponible y contiene oraciones generadas por la industria de ya sea que dominio de tecnología de la información o ciencias de la vida (Specia et al, 2017). Este corpus contiene un total de 1786 pares de oración fuente, respectiva referencia generada por humanos, hipótesis MT (ya sea de un MT estadístico basado en frases o de un neural MT), y MT post-editado (PE). Los pares de idiomas representados en este corpus son: inglés a alemán (en-de), letón (en-lt) y checo (en-cs), y alemán a inglés (de-en).

La puntuación HTER se obtiene calculando la tasa de edición de traducción (TER) entre la hipótesis MT y el PE correspondiente. Finalmente, después de calcular el HTER para cada MT, construimos un conjunto de datos de entrenamiento donde  $s_i$  denota el texto fuente,  $h_i$  denota la MT hipótesis,  $r_i$  la traducción de referencia, y la puntuación HTER para la hipótesis  $h_i$  de esta manera buscamos aprender una regresión  $r_i$  que predice el esfuerzo humano requerido para corregir la hipótesis mirando la fuente, hipótesis, y referencia (pero no la hipótesis post-editada).

## 3.2 El WMT DA RR corpus

Desde 2017, los organizadores de la Tarea Compartida de Traducción (Tarea Compartida de Traducción) han recogido juicios humanos en forma de ad-ecuaciones (DARR) (Draheim et al, 2013, 2014, 2017). Estos DAs se mapean entonces en rangos relativos-ings (DARR) (Ma et al., 2019). El resultado de datos para cada año (2017-19) forman un conjunto de datos  $\{s_i, h_i^+, h_i^-, r_i\}_{i=1}^N$  donde  $h_i^+$  denota una hipótesis "mejor" y  $h_i^-$  denota una "peor". Aquí buscamos aprender una función de tal manera que el puntaje asignado  $h_i$  a  $r_i$  es estrictamente mayor que el puntaje asignado a  $h_i^-$  ( $r(s_i, h_i^+, r_i) > r(s_i, h_i^-, r_i)$ ).

Estos datos contienen un total de 24 pares de idiomas de alto recursos como el chino al inglés (zh-en) y el inglés al gujarati (en-gu).

## 3.3 El corpus MQM

El corpus MQM es una base de datos interna propietaria de traducciones generadas por MT de soporte al cliente

<sup>3</sup>Datos QT21: <https://lindat.mff.cuni.cz/repository/xmlui/handle/11372/LRT-2390>

<sup>4</sup>Los datos brutos para cada año de la tarea compartida de métricas son disponibles públicamente en la página de resultados (ejemplo de ejemplo <http://www.statmt.org/wmt19/results.html>). Sin embargo, ten en cuenta que los archivos se resalta que estos ocasionalmente requieren utilidades personalizadas que no están disponibles.



mensajes de chat que fueron anotados de acuerdo con las pautas establecidas por [Biderman y Lommel \(2014\)](#). Estos datos contienen un total de 12K tuplas, cubriendo 12 pares de idiomas del inglés al: alemán (en-de), español (en-es), español latinoamericano (en-es-latam), francés (en-fr), italiano (en-it), japonés (en-ja), holandés (en-nl), portugués (en-pt), portugués brasileño (en-pt-br), ruso (en-ru), sueco (en-sv), y turco (en-tr). Nota que en este corpus, el inglés siempre se ve como el idioma fuente pero nunca como el idioma objetivo. Cada tupla consta de una oración fuente, una referencia generada por humanos, una hipótesis de MT, y su puntuación MQM, derivada de las anotaciones de error por uno (o más) anotadores entrenados. La métrica MQM a lo largo de este documento es una métrica interna de acuerdo con el marco MQM ([Lommel et al,2014](#)) (MQM). Los errores se anotan bajo una tipología interna definida bajo tres tipos de errores: 'Estilo', 'Fluidez' y 'Precisión'. Nuestros puntajes MQM varían de 0 a 100 y se definen como:

$$MQM = 100 \cdot \left( \frac{\text{Menor}}{L} \times \frac{\text{Mayor}}{100} + \frac{\text{Crítico}}{100} \right) \quad (7)$$

Longitud de la oración

donde  $Y_{\text{Menor}}$  denota el número de errores menores,  $Y_{\text{Mayor}}$  el número de errores mayores, y  $Y_{\text{Crítico}}$  el número de errores críticos.

Nuestra métrica MQM tiene en cuenta la severidad de los errores identificados en la hipótesis MT, lo que lleva a una métrica más detallada que HTER o DA. Cuando se usaron en nuestros experimentos, estos valores se dividieron por 100 y se truncaron a 0. Como se describe en la sección 2.4, construimos un conjunto de datos  $D = \{s_i, h_i, r_i, y_i\}_{i=1}^N$ , donde  $s_i$  denota el texto fuente,  $h_i$  denota la hipótesis MT,  $r_i$  la traducción de referencia, y  $y_i$  la puntuación MQM para la hipótesis  $h_i$ .

## 4 Experimentos

Entrenamos dos versiones del modelo Estimador descrito en la sección 2.3: uno que regresa en HTER (COMET-HTER) entrenado con el corpus QT21, y otro que regresa en nuestra implementación propia de MQM (COMET-MQM) entrenado con nuestro corpus MQM interno. Para el modelo de Clasificación de Traducción, descrito en la sección 2.4, entrenamos con el WMTDA RR corpus de 2017 y 2018 (COMET-RANK). En esta sección, introducimos la configuración

de entrenamiento para estos modelos y la correspondiente configuración de la configuración.

### 4.1 Configuración de Entrenamiento

Las dos versiones de los Estimadores (COMET-HTER/MQM) comparten la misma configuración de entrenamiento: hiper-parámetros (los detalles están incluidos en los Apéndices). Para el entrenamiento, cargamos el codificador pre-entrenado y inicializamos tanto la capa de agrupación como la capa de avance. Mientras que los escalares de capa de agrupación se establecen inicialmente a 0, los escalares de avance se establecen a 1. Durante el entrenamiento, dividimos los parámetros del modelo en dos grupos: los parámetros del codificador, que incluyen el modelo del codificador pre-entrenado, y los parámetros del regresor, que incluyen los parámetros de la red de avance directo superior. Aplicamos descongelamiento gradual y discriminación de tasas de aprendizaje ([Ruder, 2018](#)), lo que significa que el modelo de codificador está congelado mientras que el avance directo se optimiza con una tasa de aprendizaje de 5. Después de la primera época, el modelo completo se ajusta finamente pero la tasa de aprendizaje de los parámetros del codificador se establece en 0.5 para evitar el olvido catastrófico. En contraste con los dos Estimadores, para el COMET-RANK, dado que este modelo no añade ningún nuevo parámetro encima de XLM-RoBERTa (base) aparte de los escalares de capa de agrupación, usamos una tasa de aprendizaje de 5 para todo el modelo.

Configuramos la tasa de aprendizaje y la configuración de la tarea de prueba y la configuración de la tarea de entrenamiento. Como se describe en [Ma et al., 2019](#), compararemos el COMET-RANK con los modelos con el mejor rendimiento en las presentaciones compartidas de la tarea y otras métricas de vanguardia como BERTSCORE y BLEURT.<sup>5</sup> El método de evaluación utilizado es la fórmula oficial de Kendall Tau,  $\tau$ , del WMT 2019 Metrics Shared Task ([Ma et al,2019](#)) definido como:

$$\tau = \frac{\text{Concordante} - \text{Discordante}}{\text{Concordante} + \text{Discordante}} \quad (8)$$

donde  $C$  es el número de veces que una métrica asigna una puntuación más alta a la hipótesis "mejor" y  $D$  es el número de veces que una métrica asigna una puntuación más alta a la hipótesis "peor".

<sup>5</sup>Para facilitar la investigación futura, también proporcionaremos, dentro del marco, instrucciones detalladas y scripts para ejecutar otras métricas como CHR F, BLEU, BERTSCORE, y BLEURT.

Tabla 1: Tau de Kendall (correlaciones en pares de idiomas con inglés como fuente para las Métricas MT19 corpus. Para BERTSCORE informamos los resultados con el modelo de codificador predeterminado para una comparación con XLM-RoBERTa (base) para ser justos con nuestros modelos. Los valores reportados para YiSi-1 se toman directamente del documento de la tarea (Ma et al, 2019).

Métrica	en-cs	en-de	en-es	en-gu	en-kk	en-lt	en-ru	en-zh
BLEU	0.364	0.248	0.395	0.463	0.363	0.333	0.469	0.235
CHRF	0.444	0.321	0.518	0.548	0.510	0.438	0.548	0.241
YiSi-1	0.475	0.351	0.537	0.551	0.546	0.470	0.585	0.355
BERTSCORE(predeterminado)	0.500	0.363	0.527	0.568	0.540	0.464	0.585	0.356
BERTSCORE(xlmr-base)	0.503	0.369	0.553	0.584	0.536	0.514	0.599	0.317
COMET-HTER	0.524	0.383	0.560	0.552	0.508	0.577	0.539	0.380
COMET-MQM	0.537	0.398	0.567	0.564	0.534	0.574	0.615	0.378
COMET-RANK	0.603	0.427	0.664	0.611	0.693	0.665	0.580	0.449

hí o las puntuaciones asignadas a ambas hipótesis son las métricas recientemente publicadas B y B LEURT. Como en la Tabla El modelo RR muestra

Como se mencionó en los hallazgos (Mazzeo et al, 2019), fuertes correlaciones con los juicios humanos, superando las correlaciones a nivel de segmento de todas las métricas presentadas para este específico para inglés fueron frustrantemente bajas. Además, todas las métricas en cinco de siete pares de idiomas. Las presentadas mostraron una dramática falta de capacidad. El Estimator MQM muestra sorprendentemente fú clasificar correctamente los sistemas de MT fuertes. Podemos explicar de que este modelo fue entrenado si nuestros nuevos modelos de evaluación de MT también se entrenaron con inglés como objetivo. Como resultado, seguimos la configuración de evaluación de código utilizado en nuestros modelos entrenados utilizada en el análisis presentado en (Ma et al, 2019), donde se examinan los niveles de correlación "zero-shot" se debe a la inclusión de para partes de los RR datos que incluyen solo los datos de la fuente en nuestros modelos.

10, 8, 6 y 4 sistemas MT superiores.

## 5 Resultados

### 5.1 Del inglés al X

Tabla 1 muestra resultados para todos los ocho pares de idiomas con inglés como fuente. Contrastamos nuestros tres modelos contra métricas de referencia como BLEU y CHRF, la métrica ganadora de la tarea 2019. Nuestros resultados son consistentes con las observaciones de YiSi-1, así como la más reciente BERTSCORE. Observamos que en general nuestros tres modelos entrenados con datos de alta calidad superan a menudo por márgenes significativos, todas las métricas. Nuestro DARR Ranker supera a los dos Estimadores en siete de los ocho pares de idiomas. Además, los sistemas MT de mejor rendimiento para aunque el Estimator MQM se entrena solo en 12 segmentos anotados, se desempeña aproximadamente como el Estimator HTER para la mayoría de los pares de idiomas y supera todas las demás métricas en en-ru.

### 5.2 De X al inglés

Tabla 2 muestra resultados para los siete idiomas que los tres pares donde el inglés es el objetivo, nuestros tres modelos son mejores o competitivos con todos los demás modelos contra métricas de referencia. Como B LEURT, la métrica ganadora de la tarea 2019,

### 5.3 Pares de idiomas que no involucran inglés

Nuestros tres modelos fueron entrenados con datos que involucran inglés (ya sea como fuente o como objetivo). Sin embargo, para demostrar que nuestras métricas generalizan bien, las probamos en los tres pares de idiomas que no incluyen inglés en ninguna fuente o objetivo. Como se puede ver en la Tabla 3, nuestros resultados son consistentes con las observaciones de las Tablas 2.

### 5.4 Robustez ante MT de Alta Calidad

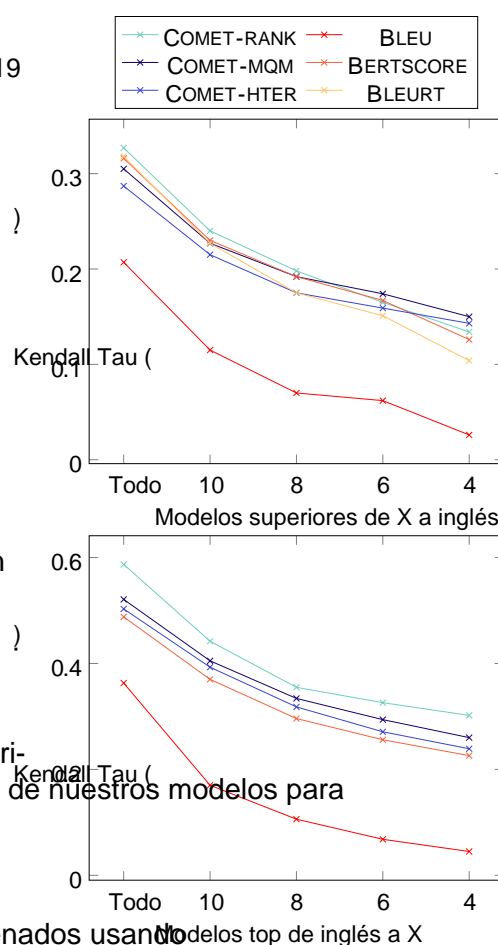
Para el análisis, utilizamos el DARR Ranker del Tarea Compartida 2019 y evaluamos en el subconjunto de los sistemas MT de mejor rendimiento para cada par de idiomas. Incluimos pares de idiomas para los cuales pudimos recuperar datos de al menos diez pares de idiomas MT (es decir, todos menos kk-en y gu-en). Contrastamos contra el fuerte recientemente propuesto BERTSCORE y B LEURT, con BLEU como base-línea. Los resultados se presentan en la Figura 3. En general, nuestras métricas superan el rendimiento de otros-

Tabla 2: Tau de Kendall ( $\tau$ ) correlaciones en pares de idiomas con inglés como objetivo para las Métricas WMT19 DACorpus RR. En cuanto a BLEURT, para BLEURT informamos los resultados de dos modelos: el modelo base, que es comparable en tamaño con el codificador que usamos y el modelo grande que es el doble de tamaño.

Métrica	de-en	fi-en	gu-en	kk-en	lt-en	ru-en	zh-en
BLEU	0.053	0.236	0.194	0.276	0.249	0.177	0.321
CHRF	0.123	0.292	0.240	0.323	0.304	0.115	0.371
YISI-1	0.164	0.347	0.312	0.440	0.376	0.217	0.426
BERTSCORE(predeterminado)	0.190	0.354	0.292	0.351	0.381	0.221	0.432
BERTSCORE(xlmr-base)	0.171	0.335	0.295	0.354	0.356	0.202	0.412
BLEURT (base-128)	0.171	0.372	0.302	0.383	0.387	0.218	0.417
BLEURT (large-512)	0.174	0.374	0.313	0.372	0.388	0.220	0.436
COMET-HTER	0.185	0.333	0.274	0.297	0.364	0.163	0.391
COMET-MQM	0.207	0.343	0.282	0.339	0.368	0.187	0.422
COMET-RANK	0.202	0.399	0.341	0.358	0.407	0.180	0.445

Tabla 3: Tau de Kendall ( $\tau$ ) correlaciones en idioma pares no involucrados en inglés para las Métricas WMT19 DACorpus RR.

Métrica	de-cs	de-fr	fr-de
BLEU	0.222	0.226	0.173
CHRF	0.341	0.287	0.274
YISI-1	0.376	0.349	0.310
BERTSCORE(predeterminado)	0.358	0.329	0.300
BERTSCORE(xlmr-base)	0.386	0.336	0.309
COMET-HTER	0.358	0.397	0.315
COMET-MQM	0.386	0.367	0.296
COMET-RANK	0.389	0.444	0.331



ers. Incluso el Estimador MQM, entrenado solo con 12K segmentos, es competitivo, lo que destaca el poder de nuestro marco propuesto.

## 5.5 La Importancia de la Fuente

Para arrojar algo de luz sobre el valor real y la contribución del input del lenguaje fuente en la capacidad de aprender predicciones precisas, entrenamos dos versiones de nuestro modelo RR Ranker: uno que usa solo la referencia, y otro que usa tanto la referencia como la fuente. Ambos modelos fueron entrenados usando el corpus WMT 2017 que solo incluye pares de idiomas

del inglés (en-de, en-cs, en-es, en-tr). En otras palabras, mientras que el inglés nunca fue observado como

un idioma objetivo durante el entrenamiento para ambas variantes

del modelo, el entrenamiento de la segunda variante incluye

incrustaciones de fuente en inglés. Luego probamos claramente muestra que para la arquitectura de clasificación

estas dos variantes del modelo en el corpus WMT2018 la fuente mejora la correlación general

para estos pares de idiomas y para las direcciones inversas humanos. Además,

ciones (con la excepción de en-cs porque cs-en la inclusión de la fuente expuso la segunda vari-

no existe para WMT 2018). Los resultados en la Tabla 4 muestran el efecto de la inclusión de la fuente en los resultados de los modelos.

Tabla 4: Comparación entre COMET-RANK (sección 2.4) y una versión solo de referencia de la misma en datos WMT18. Ambos modelos fueron entrenados con WMT17 lo que significa que el modelo solo de referencia nunca se expone al idioma durante el entrenamiento.

Métrica	en-cs	en-de	en-es	en-tr	cs-en	de-en	es-en	tr-en
COMET-RANK (solo ref.)	0.660	0.764	0.630	0.539	0.249	0.390	0.159	0.128
COMET-RANK	0.711	0.799	0.671	0.563	0.356	0.542	0.278	0.260
· ·	0.051	0.035	0.041	0.024	0.107	0.155	0.119	0.132

reflejado en un más alta para los pares de idiomas con un espacio de incrustación y luego calcular una puntuación que refleja la similitud semántica entre esos segmentos. Sin embargo, los juicios humanos como DA y MQM, capturan mucho más que solo semántica similar, resultando en una correlación superior entre los juicios humanos y las puntuaciones producidas por tales métricas.

## 6 Reproducibilidad

Vamos a liberar tanto la base de código COMET-C como el marco y los modelos de evaluación de MT entrenados descritos en este documento a la comunidad de investigación tras la publicación, junto con los scripts detallados necesarios para ejecutar todas las líneas de código reportadas. Los modelos reportados en este documento fueron entrenados en un único Tesla T4 (16GB) GPU. Además, nuestro marco de trabajo se basa en PyTorch Lightning (Falcon 2019), un envoltorio ligero de PyTorch, que fue creado para la máxima flexibilidad y reproducibilidad.

## 7 Trabajo relacionado

Las métricas clásicas de evaluación de MT se caracterizan comúnmente como métricas de coincidencia por programas utilizando características hechas a mano, estiman la calidad de MT calculando el número y la fracción de palabras que aparecen simultáneamente en una hipótesis de traducción y una o más humanas referencias. Métricas como BLEU (Papineni et al. 2002), METEOR (Lavie y Denkowski, 2009), y CHRF (Popović, 2015) han sido ampliamente estudiados y mejorados (Koehn et al., 2007; Popović, 2017; Denkowski y Lavie, 2011; Guo y Hu, 2019), pero, por diseño, generalmente fallan en reconocer y capturar la similitud semántica más allá del nivel léxico.

En años recientes, los incrustaciones de palabras (Mikolajev et al., 2013; Pennington et al., 2014; Peters et al., 2018; Devlin et al., 2019) han surgido como una común alternativa utilizada para la coincidencia para capturar la similitud semántica de las palabras. Las métricas basadas en incrustaciones como METEOR-VECTOR (Servan et al., 2016), BLEU2VEC (Tättar y Fishel, 2017), YISI-1 (Lo, 2019), MOVERSCORE (Zhao et al., 2019), y BLEURTSCORE (Zhang et al., 2020) crean alineaciones suaves entre referencia e hipótesis

<sup>6</sup>Estos serán alojados en <https://github.com/Unbabel/COMET>

que refleja la similitud semántica entre esos segmentos. Sin embargo, los juicios humanos como DA y MQM, capturan mucho más que solo semántica similar, resultando en una correlación superior entre los juicios humanos y las puntuaciones producidas por tales métricas. Métricas aprendibles (Shimanaka et al., 2018; Mathur et al., 2019; Shimanaka et al., 2019) intentan optimizar directamente la correlación con juicios humanos, y han mostrado recientemente promising results. BLEURT (Sellam et al., 2020), una métrica aprendible basada en BERT (Devlin et al., 2019), reclama un rendimiento de vanguardia durante los últimos años de la tarea compartida de métricas WMT. Porque BLEURT se construye sobre English-BERT (Devlin et al., 2019), solo puede ser utilizado cuando el inglés es el idioma objetivo que limita su aplicabilidad. Además, hasta donde sabemos, todos los anteriormente propuestos métricas aprendibles se han centrado en optimizar DA que, debido a la escasez de anotadores, puede resultar inherentemente sesgado.

Evaluación de MT sin referencias, también conocida como Estimación de Calidad (QE), históricamente a menudo regresado en HTER para la evaluación a nivel de segmento (Barrault et al., 2013, 2014, 2015, 2016, 2017a). Más recientemente, MQM ha sido utilizado para la evaluación a nivel de segmento (Specia et al., 2018; Fonseca et al., 2019). Al aprovechar codificadores pre-entrenados altamente multilingües como BERT multilingüe (Devlin et al., 2019) y XLM (Conneau y Lample, 2019), los sistemas QE han estado mostrando auspiciosas correlaciones con los juicios humanos (Kociský et al., 2019a). Concurrentemente, el OpenKiwi marco (Kepler et al., 2019b) ha facilitado para los investigadores avanzar en el campo y construir modelos QE más fuertes.

## 8 Conclusiones y Trabajo Futuro

En este documento presentamos COMET-C, un nuevo marco neural para entrenar modelos de evaluación MT que pueden servir como métricas automáticas y ser fácilmente



adaptado y optimizado para diferentes tipos de humil-  
juicios de calidad de MT.

Para mostrar la efectividad de nuestro marco, buscamos abordar los desafíos reportados en el 2019 WMT Metrics Shared Task (Ma et al, 2019).

Entrenamos tres modelos distintos que logran nuevos resultados de vanguardia para la correlación a nivel de segmento con juicios humanos, y muestran una capacidad prometedora para diferenciar mejor los sistemas de alto rendimiento.

Uno de los desafíos de aprovechar el poder de los modelos preentrenados es el peso oneroso de parámetros y tiempo de inferencia. Un camino principal para el trabajo futuro COMET se centrará en el impacto de soluciones más compactas como DistilBERT (Sanh et al, 2019).

Además, mientras esbozamos la posible importancia del texto fuente anterior, notamos que nuestro COMET-RANK modelo pondera la fuente y la referencia de manera diferente durante la inferencia pero igual en su función de pérdida de entrenamiento. El trabajo futuro optimizará esta formulación y examinará más a fondo la interdependencia de las diferentes entradas.

Agradecimientos

Estamos agradecidos a Andr

Edoardo Azzurro, Diarmuid O'Gartaigh, Miguel Varona, los revisores, por sus valiosos comentarios y discusiones. Este trabajo fue apoyado en parte por el Programa P2020 a través de los proyectos MAIA y Unbabel4EU, supervisados por ANI bajo los números de contrato 045909 y 042671, respectivamente.

Referencias

Mikel Artetxe y Holger Schwenk. 2019.

Mas-incrustaciones de oraciones multilingües para cero-transferencia cruzada lingüística y más Transacciones de la Asociación para la Lingüística Computacional, 7:597-610.

Loïc Barrault, Ondrej Bojar, Marta R. Costa-jussà, 2017b. Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, Christof Monz, Mathias Müller, Santanu Pal, Matt Post, y Marcos Zampieri. 2019. Hallazgos de la conferencia 2019 sobre traducción automática (WMT19) En Actas de la Cuarta Conferencia sobre Traducción Automática (Volumen 2: Documentos de Trabajo Compartidos, Día 1), páginas 1-61, Florencia, Italia. Asociación para la Lingüística Computacional. Ondrej Bojar, Christian Buck, Chris Callison-Burch, Christian

Federman Barry Haddow Philipp Koehn, Christof Monz, Matt Post, Radu Soricut, y Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco

- Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, y Veselin Stoyanov. 2019. [Aprendizaje no supervisado de representación cruzada de idiomas a gran escala](#). preimpresión arXiv:1911.02116
- Alexis Conneau y Guillaume Lample. 2019. [Pretraining with masked word prediction](#). En [ICML](#). Larochelle, H., Beygelzimer, A., H. Larochelle, A. Beygelzimer, M. de Lencq, para la clasificación de texto. Actas de la 56ª Reunión Anual de la Asociación para la Lingüística Computacional (Volumen 1: 7069. Curran Associates, Inc. Michael Denkowski y Do Almondo, páginas 328-339, Melbourne, Australia. Asociación para la Lingüística Computacional.
- Métrica automática para la optimización de sistemas de traducción automática. En [Actas del Sexto Taller sobre Traducción Automática](#), páginas 1-5. En [Participación de Unbabel en la tarea compartida de estimación de calidad de traducción](#). Jacob Devlin, Ming-Wei Chang, Kenton Lee, y Kristina Toutanova. 2019. [Ensamblado de transformadores bidireccionales profundos para la comparación de documentos](#). En [Actas de la Conferencia 2019 del Capítulo Norteamericano de la Asociación para la Lingüística Computacional](#), páginas 116-126. En [PyTorch Lightning: El envoltorio ligero de PyTorch para la investigación de IA de alto rendimiento](#). GitHub
- WA Falcon. 2019. [PyTorch Lightning: The lightweight PyTorch wrapper for high-performance AI research](#). GitHub.
- Erick Fonseca, Lisa Yankovskaya, André F. T. Martins, Mark Fishel y Christian Federmann. 2019. [Hallazgos de las tareas de aprendizaje de máquina de código abierto para la traducción automática](#). En [Actas de la Cuarta Conferencia sobre Traducción Automática \(Volumen 3: Documentos de Tarea Compartida, Día 2\)](#), páginas 1-10, Florencia, Italia. Asociación para la Lingüística Computacional.
- Yvette Graham, Timothy Baldwin, Alistair Moffat y Justin Zobel. 2018. [Escalas de medición continua en la evaluación humana de la traducción automática e Interoperabilidad con el Discurso](#). , páginas 33-41, 21. 22. Dan Kondratyuk y Milan Straka. 2019. [75 lecciones de la 7ª Conferencia sobre Traducción Automática](#). En [Actas de la Conferencia sobre Traducción Automática](#), páginas 177-180, Praga, República Checa. Asociación para la Lingüística Computacional.
- Yvette Graham, Timothy Baldwin, Alistair Moffat y Justin Zobel. 2014. [¿Está mejorando la traducción automática con el tiempo? En](#) [Actas de la 14ª Conferencia del Capítulo Europeo de la Asociación para la Lingüística Computacional](#), páginas 443-451, Gotemburgo, Suecia. Asociación para la Lingüística Computacional. 33. [El meteor métrico para la evaluación de la traducción automática](#). En [Actas de la Cuarta Conferencia sobre Traducción Automática](#), páginas 507-513, Florencia, Italia. Asociación para la Lingüística Computacional.
- Yvette Graham, Timothy Baldwin, Alistair Moffat y Justin Zobel. 2017. [¿Pueden los sistemas de traducción automática ser evaluados solo por la multitud?](#) Ingeniería del Lenguaje Natural, 23(1):330. 41. [Métricas de calidad de traducción automática](#). En [Actas de la Cuarta Conferencia sobre Traducción Automática](#), páginas 507-513, Florencia, Italia. Asociación para la Lingüística Computacional.
- Yinuo Guo y Junfeng Hu. 2019. [Meteor++ 2.0: Adoptar conocimientos de paráfrasis a nivel sintáctico](#). En [Actas de la Cuarta Conferencia sobre Traducción Automática](#), páginas 45. 45. [Métricas de calidad de traducción automática](#). En [Actas de la Cuarta Conferencia sobre Traducción Automática](#), páginas 507-513, Florencia, Italia. Asociación para la Lingüística Computacional.

- marco para declarar y describir la traducción métricas de calidad. *Traducción automática: tecnologías de la traducción*, páginas 455-463. 5001, Florencia, Italia. Asociación para la Computación Lingüística.
- Qingsong Ma, Ondrej Bojar, y Yvette Graham. 2018. **Resultados de la tarea compartida de métricas WMT18: Ambos caracteres y incrustaciones logran un buen rendimiento**. En *Actas de la Tercera Conferencia sobre Traducción Automática: Documentos de Tarea Compartida*, páginas 671-688, Bélgica, Bruselas. Asociación para la Computación Lingüística.
- Qingsong Ma, Johnny Wei, Ondrej Bojar, y Yvette Graham. 2019. **Resultados de las métricas WMT19 tarea compartida: sistemas de MT a nivel de segmento presentan grandes desafíos**. En *Actas de la Cuarta Conferencia sobre Traducción Automática (Volumen 2: Documentos de Tarea Compartida)*, páginas 372-390, Florencia, Italia. Asociación para la Computación Lingüística.
- Nitika Mathur, Timothy Baldwin, y Trevor Cohn. 2019. **Poniendo la evaluación en contexto: Los incrustaciones contextuales mejoran la evaluación de la traducción**. En *Actas de la 57ª Reunión Anual de la Asociación para la Computación Lingüística*, páginas 2799-2808, Florencia, Italia. Asociación para la Computación Lingüística.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, y Jeff Dean. 2013. **Representaciones distribucionales de palabras y frases y sus analogías**. En *Sistemas de Procesamiento de Información Neurales*, páginas 3111-3119. Curran Associates, Inc.
- Kishore Papineni, Salim Roukos, Todd Ward, y Wei-Jing Zhu. 2002. **Bleu: un método para la evaluación automática de la traducción**. En *Actas de la 40ª Reunión Anual de la Asociación para la Computación Lingüística*, páginas 311-318, Filadelfia, Pensilvania, EE. UU. Asociación para la Computación Lingüística.
- Jeffrey Pennington, Richard Socher, y Christopher Manning. 2014. **Glove: Vectores globales para la representación de palabras**. En *Actas de la Conferencia sobre Métodos Empíricos en el Procesamiento del Lenguaje Natural (EMNLP)*, páginas 1532-1543, Doha, Qatar. Asociación para la Computación Lingüística.
- Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, y Luke Zettlemoyer. 2018. **Representaciones de palabras contextualizadas profundamente**. En *Actas de la Conferencia 2018 del Capítulo Norteamericano de la Asociación para la Computación Lingüística: Tecnologías del Lenguaje Humano, Volumen 1 (Documentos largos)*, páginas 2227-2237, Nueva Orleans, Luisiana. Asociación para la Computación Lingüística.
- Telmo Pires, Eva Schlinger, y Dan Garrette. 2019. **¿Qué tan multilingüe es BERT multilingüe?** En *Procedimientos de la 57ª Reunión Anual de la Asociación para la Computación Lingüística*, páginas 4996-5001, Florencia, Italia. Asociación para la Computación Lingüística.
- Maja Popović. 2015. **chrF: puntuación f de n-gramas de caracteres para la evaluación automática de MT**. En *Decimo Taller de Traducción Automática Estadística*, páginas 392-395, Lisboa, Portugal. Asociación para la Computación Lingüística.
- Maja Popović. 2017. **chrF++: palabras ayudando a los caracteres n-gramas**. En *Actas de la Segunda Conferencia de Traducción Automática*, páginas 612-618, Copenhague, Dinamarca. Asociación para la Computación Lingüística.
- Nils Reimers e Iryna Gurevych. 2019. **Sentence-BERT: Incrustaciones de frases usando redes Siamese BERT-networks**. En *Actas de la Conferencia 2019 sobre Métodos Empíricos en Procesamiento del Lenguaje Natural y la 9na Conferencia Internacional Conjunta sobre Procesamiento del Lenguaje Natural (EMNLP-IJCNLP)*, páginas 3982-3992, Hong Kong, China. Asociación para la Computación Lingüística.
- Victor Sanh, Lysandre Debut, Julien Chaumond, y Thomas Wolf. 2019. **Distilbert, una versión destilada de BERT: más pequeña, más rápida, más barata y más ligera**. preimpresión arXiv:1910.01108
- F. Scrocco, D. Kaler, y J. Philbin. 2015. **Facenet: Una incrustación unificada para el reconocimiento facial**. En *Conferencia IEEE 2015 sobre Visión por Computadora y Reconocimiento de Patrones (CVPR)*, páginas 815-823.
- Thibault Sellam, Dipanjan Das, y Ankur Parikh. 2020. **BLEURT: Aprendiendo métricas robustas para la generación de texto**. En *Actas de la 58va Reunión Anual de la Asociación para la Computación Lingüística*, páginas 7881-7892, Online. Asociación para la Computación Lingüística.
- Christophe Servan, Alexandre B'érard, Zied Elloumi, Hervé Blanchon, y Laurent Besacier. 2016. **Word2Vec vs DBnary: ¿Aumentando METEOR usando representaciones vectoriales o recursos léxicos?** En *Actas de COLING 2016, la 26ª Conferencia Internacional sobre Lingüística Computacional: Documentos Técnicos*, páginas 1159-1168, Osaka, Japón. El Comité Organizador de COLING 2016.
- Hiroki Shimanaka, Tomoyuki Kajiwar, y Mamoru Komachi. 2018. **RUSE: Regresor utilizando oraciones incrustaciones para la evaluación automática de la traducción**. En *Actas de la Tercera Conferencia sobre Traducción Automática: Documentos de Tarea Compartida*, páginas 751-758, Bélgica, Bruselas. Asociación para la Computación Lingüística.
- Hiroki Shimanaka, Tomoyuki Kajiwar, y Mamoru Komachi. 2019. **Evaluación de la Traducción de Máquinas con BERT Regresor**. preimpresión de arXiv arXiv:1907.12679



- Matthew Snover, Bonnie Dorr, Richard Schwartz, Linnea Micciulla y John Makhoul. 2006. [Un estudio de la tasa de edición de traducción con anotación humana](#). En las Actas de la Asociación para la Traducción Automática en línea, páginas 223-231.
- Lucia Specia, Frédéric Blain, Varvara Logacheva, Ramón Astudillo y André F. T. Martins. 2018. [Hallazgos de la tarea compartida WMT 2018 sobre la calidad estimación](#). En Actas de la Tercera Conferencia sobre Traducción Automática: Documentos compartidos, páginas 689-709, Bélgica, Bruselas. Asociación para la Computación Lingüística.
- Lucia Specia, Kim Harris, Frédéric Blain, Aljoscha Burchardt, Vivien Macketanz, Inguna Skadina, Matteo Negri, y Marco Turchi. 2017. [Traducción de calidad y productividad: Un estudio sobre idiomas ricos en morfología](#). En Cumbre de Traducción Automática XVI, páginas 55-71, Nagoya, Japón.
- Kosuke Takahashi, Katsuhito Sudoh y Satoshi Nakamura. 2020. [Evaluación automática de la traducción utilizando entradas en el idioma fuente y modelo de idioma cruzado lingüístico](#). En Actas de la 58ª Reunión Anual de la Asociación para la Computación Lingüística, páginas 3553-3558, En línea. Asociación para la Computación Lingüística.
- Andre T'attar y Mark Fishel. 2017. [bleu2vec: la métrica dolorosamente familiar en el espacio vectorial continuo estereotipado](#). En Actas de la Segunda Conferencia sobre Traducción Automática, páginas 619-622, Copenhague, Dinamarca. Asociación para la Computación Lingüística.
- Ian Tenney, Dipanjan Das y Ellie Pavlick. 2019. [BERT redescubre el clásico pipeline de NLP](#). En Actas de la 57ª Reunión Anual de la Asociación para la Lingüística Computacional, páginas 4593-4601, Florencia, Italia. Asociación para la Lingüística Computacional.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, y Yoav Artzi. 2020. [Bertscore: Evaluando la generación de texto con BERT](#). En Conferencia Internacional sobre Representaciones de Aprendizaje, Wei Zhao, Goran Glavač, Maxime Peyrard, Yang Gao, .
- Robert West, y Steffen Eger. 2020. [Sobre las limitaciones de los codificadores cruzados como se expone por la evaluación de la traducción automática sin referencia](#). En Actas de la 58ª Reunión Anual de la Asociación para la Lingüística Computacional, páginas 1656-1671, En línea. Asociación para la Lingüística Computacional.
- Wei Zhao, Maxime Peyrard, Fei Liu, Yang Gao, Christian M. Meyer, y Steffen Eger. 2019. [MoverScore: Evaluación de la generación de texto con incrustaciones contextualizadas y la distancia de Levenshtein](#). En Conferencia Internacional Conjunta sobre Procesamiento del Lenguaje Natural (EMNLP-IJCNLP), páginas 563-578, Hong Kong, China. Asociación para la Lingüística Computacional.



## A Apéndices

En la [tabla 5.1](#) numeramos los hiperparámetros utilizados para entrenar nuestros modelos. Antes de inicializar estos modelos, se estableció una semilla aleatoria a 3 en todas las bibliotecas que realizan operaciones "aleatorias" (`torch`, `numpy`, `random` y `cuda` ).

Tabla 5: Hiper-parámetros utilizados en nuestro Q para entrenar los modelos presentados.

Hiper-parámetro	COMET (Est-HTER/MQM)	COMET-RANK
Modelo del codificador	XLM-RoBERTa (base)	XLM-RoBERTa (base)
Optimizador	Adam (parámetros predeterminados)	Adam (parámetros predeterminados)
n épocas congeladas	1	0
Tasa de aprendizaje	3e-05 y 1e-05	1e-05
Tamaño del lote	16	16
Función de pérdida	MSE	Margen de triada (1.0)
Abandono por capas	0.1	0.1
Precisión FP	32	32
Unidades ocultas de avance rápido	2304, 1152	.
Activaciones de avance rápido	Tanh	.
Abandono de avance rápido	0.1	.

Tabla 6: Estadísticas para el corpus QT21.

	en-de	en-cs	en-lv	de-en
Total de tuplas	54000	42000	35474	41998
Prom. de t <sub>q</sub> (referencia)	17.80	15.56	16.42	17.71
Prom. de t <sub>q</sub> (k=1)	16.70	17.37	18.39	17.18
Prom. de t <sub>q</sub> (MT)	17.65	15.64	16.42	17.78

Tabla 7: Estadísticas para WMT 2017 Corpus RR.

	en-cs	en-de	en--	en-lv	en-tr
Total de tuplas	32810	6454	3270	3456	247
Prom. de t <sub>q</sub> (referencia)	19.70	22.15	15.59	21.42	17.57
Prom. de t <sub>q</sub> (k=1)	22.37	23.41	21.73	26.08	22.51
Prom. de t <sub>q</sub> (MT)	19.45	22.58	16.06	22.18	17.25

	zh-es	914203.539.70
	ru-es	3925274.20.80
	lt-en	2186235.25.25
	kk-en	9728.36.32.68
Pares de idiomas RR a inglés.		
DA	gu-en	2011704.92.02
	--en	321855.42.76
	de-es	8536529.22.22
	(ref. (MT))	
Tabla 8: Estadísticas para el WMT 2019		
	Total Pares de idiomas	

	de-fr	486723226.68
	de-es	23224722.89
	fr-de	36922823.36
	en-fr	186524.6.83
	en-es	242447423.37
	en-ll	174012026.97
	en-kl	18182319.92
Pares de idiomas RR desde inglés y sin inglés.		
DA	en-gu	3532432.97
	en-31	2205223.69
	en-99	2402527.98
	en-25	2722422.60
	(ref. (MT))	
Tabla 9: Estadísticas para el WMT 2019		
	Total Pares de idiomas	

en-es	10.13.33	17
en-es	10.13.33	17
en-es	12.13.42	21
en-es	17.950.399	
en-es	12.13.42	19
en-es	12.13.64	22.02
estadísticas	17.473.75.85	59
en-es	2590.90.23	88
en-es	1043.37.93	19
en-es	2756.73.73	41
Tabla 10: Corpus MQM (sección		
en-es	1520.32.62	84
en-es	3704.23.33	91
en-es	2447.14.23	66
(referred)		
Total de palabras		

et-es	5623.18.23	52
en-es	1358.24.19	61
en-es	2228.25.24	86
en-es	19805.22.82	15
en-es	2216.23.48	37
en-es	192815.24.23	74
en-es	5419.20.19	73
Pares de idiomas RR.		
DA	6723.20.25	64
tr-es	6523.25.22	80
ru-es	1044.27.25	25
fi-es	15248.15.23	46
Tabla 11: Estadísticas para el WMT 2018		
en-es	2824.28.24	94
zh-es	3325.26.25	45
(referred)		
Total de palabras		



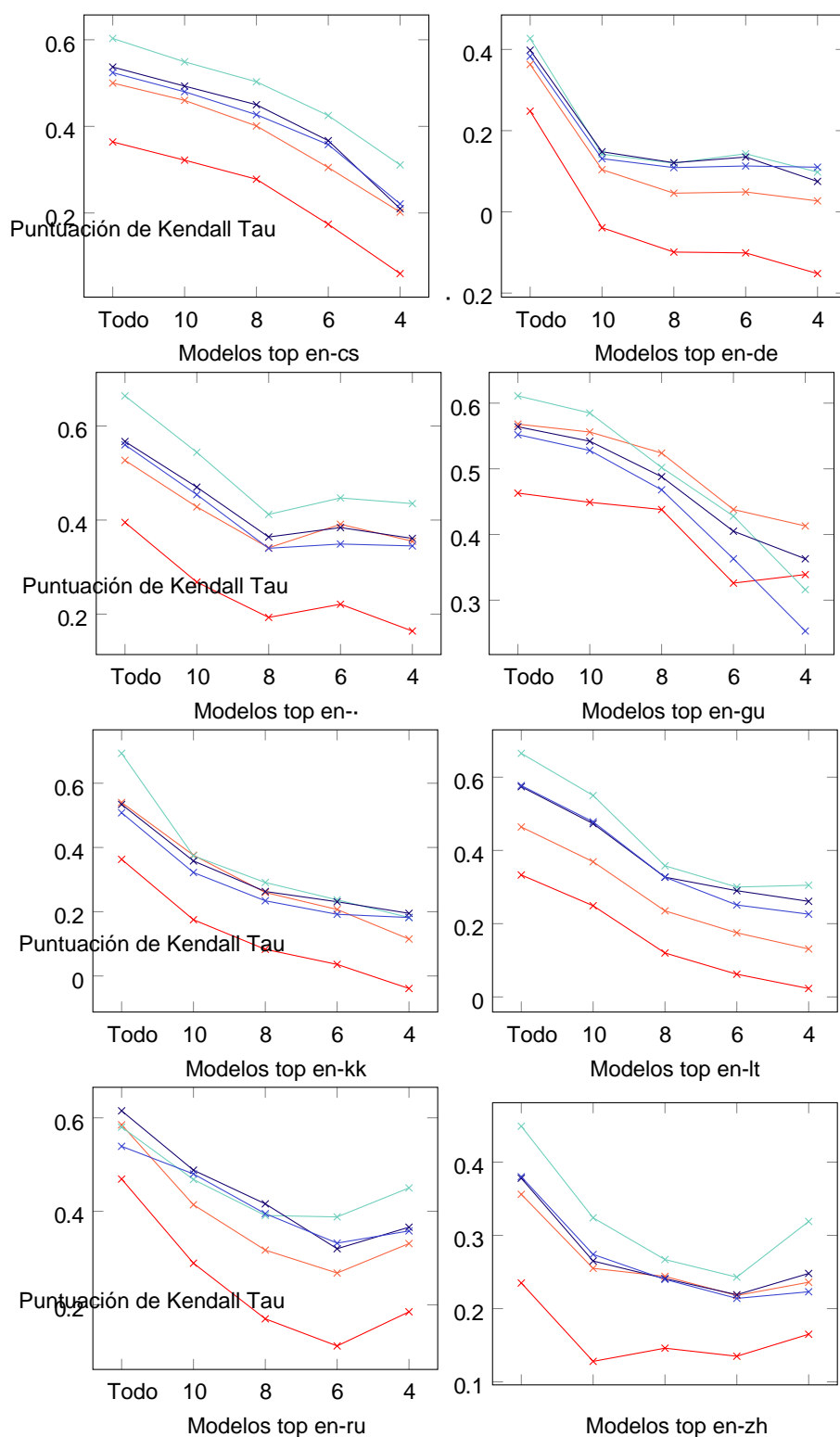


Tabla 12: Rendimiento de métricas en todos y los mejores (10,8, 6 y 4) sistemas MT para todos los pares de idiomas d. El esquema de color es el siguiente: COMET-RANK , COMET-HTER , COMET-MQM , BLEU , BERTSCORE

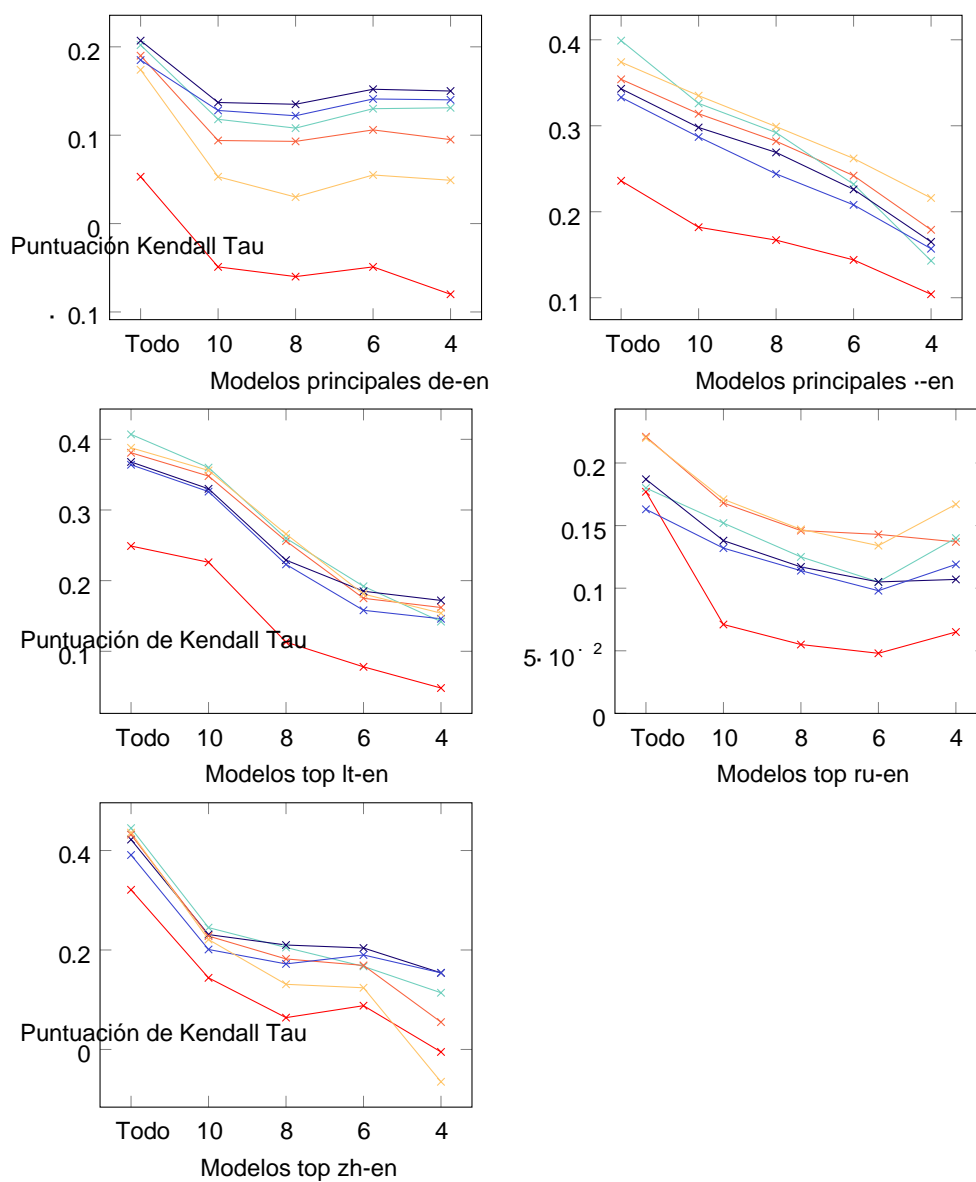


Tabla 13: Rendimiento de las métricas en todos y los mejores (10,8, 6 y 4) sistemas de MT para todos los idiomas en pares. El esquema de color es el siguiente: COMET-RANK, COMET-HTER, COMET-MQM, BLEU, BERTSCORE, BLEURT