

COMET: Un marco neuronal para la evaluación de MT

Ricardo Rei Craig Stewart Ana C Farinha Alon Lavie Unbabel IA

{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 correlación de vanguardia con los juicios humanos. Nuestro marco aprovecha los recientes avances en el modelado de lenguaje preentrenado cruzado, lo que resulta en modelos de evaluación de TA altamente multilingües y adaptables que explotan información tanto de la entrada de origen como de una traducción de referencia en el idioma de destino para predecir con mayor precisión la calidad de la TA. Para mostrar nuestro marco, entrenamos tres modelos con diferentes tipos de juicios humanos:

Direct Assessments, *Human-mediated Translation Edit Rate* y *Multidimensional Quality Metrics*. Nuestros modelos logran un rendimiento de vanguardia en la tarea compartida de métricas WMT 2019 y demuestran robustez frente a sistemas de alto rendimiento.

1 Introducción

Históricamente, las métricas para evaluar la calidad de la traducción automática (TA) se han basado en evaluar la similitud entre una hipótesis generada por TA y una traducción de referencia generada por humanos en el idioma de destino. Las métricas tradicionales se han centrado en características básicas a nivel léxico, como contar el número deigramas coincidentes entre la hipótesis de TA y la traducción de referencia. Métricas como BLEU (Papineni et al., 2002) y METEOR (Lavie y Denkowski, 2009) siguen siendo populares como medio para evaluar sistemas de TA debido a su ligereza y rápida computación.

Los enfoques neuronales modernos para la traducción automática resultan en una calidad de traducción mucho más alta que a menudo se desvía de la transferencia léxica monótona entre idiomas. Por esta razón, se ha vuelto cada vez más evidente que ya no podemos confiar en métricas como BLEU para proporcionar una estimación precisa de la calidad de la traducción automática (Barrault et al., 2019).

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

Los hallazgos de la tarea mencionada anteriormente destacan dos desafíos importantes para la evaluación de MT que buscamos abordar aquí (Ma et al., 2019). A saber, que las métricas actuales tienen dificultades para correlacionarse con precisión con el juicio humano a nivel de segmento y no logran diferenciar adecuadamente los sistemas de MT de mejor rendimiento.

En este artículo, presentamos COMET¹, un marco basado en PyTorch para entrenar modelos de evaluación de MT altamente multilingües y adaptables que pueden funcionar como métricas. Nuestro marco aprovecha los recientes avances en 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 *Direct Assessments* (DA) (Graham et al., 2013), *Human-mediated Translation Edit Rate* (HTER) (Snover et al., 2006) y métricas compatibles con el marco *Multidimensional Quality Metric* (Lommel et al., 2014).

Inspirados en trabajos recientes sobre Estimación de Calidad (QE) que demostraron que es posible alcanzar altos niveles de correlación con los juicios humanos incluso sin una traducción de referencia (Fonseca et al., 2019), proponemos un enfoque novedoso para incorporar-

¹Crosslingual Optimized Metric for Evaluation of Translation.

ingresar la entrada en el idioma fuente en nuestros modelos de evaluación de MT. Tradicionalmente, solo los modelos de QE han utilizado la entrada fuente, mientras que las métricas de evaluación de MT se basan en la traducción de referencia. Como en (Takahashi et al., 2020), mostramos que el uso de un espacio de incrustación multilingüe nos permite aprovechar la información de las tres entradas y demostrar el valor añadido por la fuente como entrada a nuestros modelos de evaluación de MT.

Para ilustrar la efectividad y flexibilidad del marco COMET, entrenamos tres modelos que estiman diferentes tipos de juicios humanos y mostramos un progreso prometedor hacia una mejor correlación a nivel de segmento y robustez ante traducciones automáticas de alta calidad.

Publicaremos tanto el marco COMET como los modelos de evaluación de MT entrenados descritos en este documento a la comunidad de investigación tras la publicación.

2 Arquitecturas de Modelos

Los juicios humanos sobre la calidad de la MT suelen presentarse en forma de puntuaciones a nivel de segmento, como DA, MQM y HTER. Para DA, es una práctica común convertir las puntuaciones en clasificaciones relativas (DARR) cuando el número de anotaciones por segmento es limitado (Borjar et al., 2017b; Ma et al., 2018, 2019). Esto significa que, para dos hipótesis de MT 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 , h_i se considera una hipótesis "mejor". Para abarcar estas diferencias, nuestro marco admite dos arquitecturas distintas: el modelo Estimador y el modelo de Clasificación de Traducción. La diferencia fundamental entre ellos es el objetivo de entrenamiento. Mientras que el Estimador se entrena para hacer una regresión directa sobre 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 multilingüe y una capa de agrupamiento.

2.1 Codificador multilingüe

El bloque de construcción principal de todos los modelos en nuestro marco es un modelo multilingüe 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 codificadores transformar que son

²In the WMT Metrics Shared Task, if the difference between the DA scores is not higher than 25 points, those segments are excluded from the DARR data.

entrenados para reconstruir tokens enmascarados al descubrir la relación entre esos tokens y los circundantes. Cuando se entrena con datos de múltiples idiomas, se ha encontrado que este objetivo preentrenado es altamente efectivo en tareas multilingües como la clasificación de documentos y la inferencia en lenguaje natural (Conneau et al., 2019), generalizando bien a idiomas y escrituras no vistas (Pires et al., 2019). Para los experimentos en este documento, nos basamos en XLM-RoBERTa (base) como nuestro modelo de codificador.

Dada una secuencia de entrada $x = [x_0, x_1, \dots, x_n]$, el codificador produce un embedding $e_j^{(\ell)}$ para cada token x_j y cada capa $\ell \in \{0, 1, \dots, k\}$. En nuestro marco, aplicamos este proceso a la fuente, la hipótesis de MT y la referencia para mapearlas en un espacio de características comparativo.

2.2 Capa de Agrupamiento

Las incrustaciones generadas por la última capa de los codificadores preentrenados se utilizan generalmente para ajustar modelos a nuevas tareas. Sin embargo, (Tenney et al., 2019) mostraron que diferentes capas dentro de la red pueden capturar información lingüística que es relevante para diferentes tareas posteriores. En el caso de la evaluación de MT, (Zhang et al., 2020) mostraron 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, utilizamos el enfoque descrito en Peters et al. (2018) y agrupamos información de las capas de codificadores más importantes en una sola incrustación para cada token, e_j , utilizando un mecanismo de atención por capas. Esta incrustación se calcula como:

$$e_{x_j} = \mu E_{x_j}^\top \alpha \quad (1)$$

donde μ es un coeficiente de peso entrenable, $E_j = [e_j^{(0)}, e_j^{(1)}, \dots, e_j^{(k)}]$ corresponde al vector de incrustaciones de capa para el token x_j , y $\alpha = \text{softmax}([\alpha^{(1)}, \alpha^{(2)}, \dots, \alpha^{(k)}])$ es un vector correspondiente a los pesos entrenables por capa. Para evitar el sobreajuste a la información contenida en cualquier capa individual, utilizamos la deserción de capa (Kondratyuk y Straka, 2019), en la que con una probabilidad p el peso $\alpha^{(i)}$ se establece en $-\infty$.

Finalmente, como en (Reimers y Gurevych, 2019), aplicamos agrupamiento promedio a las incrustaciones de palabras resultantes para derivar una incrustación de oración para cada segmento.

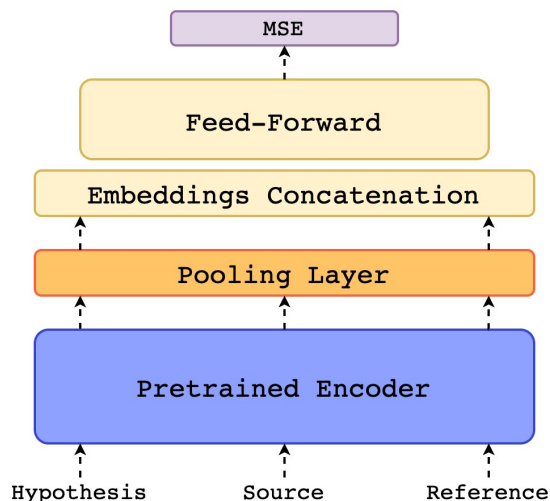


Figura 1: Arquitectura del modelo estimador. La fuente, la hipótesis y la referencia se codifican de manera independiente utilizando un codificador multilingüe preentrenado. Las incrustaciones de palabras resultantes se pasan a través de una capa de agrupamiento para crear una incrustación de oración para cada segmento. Finalmente, las incrustaciones de oración resultantes se combinan y concatenan en un solo vector que se pasa a un regresor de avance. Todo el modelo se entrena minimizando el Error Cuadrático Medio (MSE).

2.3 Modelo de Estimador

Dada una incrustación de oración de d dimensiones para la fuente, la hipótesis y la referencia, adoptamos el enfoque propuesto en RUSE (Shimanaka et al., 2018) y extraemos las siguientes características combinadas:

- Producto fuente elemento a elemento: $\mathbf{h} \odot \mathbf{s}$
- Producto de referencia elemento a elemento: $\mathbf{h} \odot \mathbf{r}$
- Diferencia absoluta elemento por elemento de la fuente: $|\mathbf{h} - \mathbf{s}|$
- Diferencia de referencia absoluta elemento por elemento: $|\mathbf{h} - \mathbf{r}|$

Estas características combinadas se concatenan luego con la incrustación de referencia \mathbf{r} y la incrustación de hipótesis \mathbf{h} en un solo vector $\mathbf{x} = [\mathbf{h}; \mathbf{r}; \mathbf{h} \odot \mathbf{s}; \mathbf{h} \odot \mathbf{r}; |\mathbf{h} - \mathbf{s}|; |\mathbf{h} - \mathbf{r}|]$ que sirve como entrada para un regresor de avance. La fuerza de estas características radica en resaltar las diferencias entre las incrustaciones en el espacio de características semánticas.

El modelo se entrena para minimizar el error cuadrático medio entre las puntuaciones predichas y las evaluaciones de calidad (DA, HTER o MQM). La Figura 1 ilustra la arquitectura propuesta.

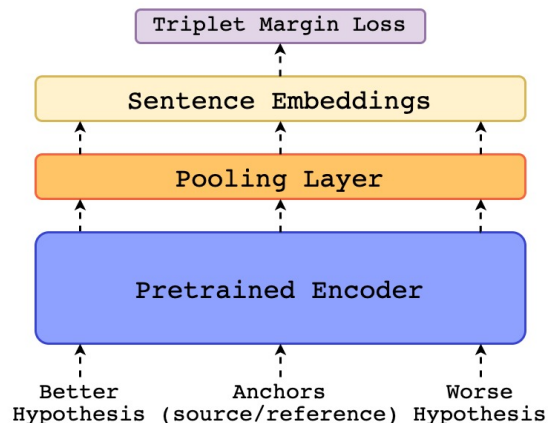


Figura 2: Arquitectura del modelo de clasificación de traducción. Esta arquitectura recibe 4 segmentos: el origen, la referencia, una hipótesis "mejor" y una "peor". Estos segmentos se codifican de manera independiente utilizando un codificador multilingüe preentrenado y una capa de agrupamiento en la parte superior. Finalmente, utilizando la pérdida de margen de tripletas (Schroff et al., 2015) optimizamos el espacio de incrustación resultante para minimizar la distancia entre la hipótesis "mejor" y los "anclajes" (origen y referencia).

Tenga en cuenta que elegimos no incluir la incrustación de origen sin procesar (\mathbf{s}) en nuestra entrada concatenada. Experimentos iniciales revelaron que el valor añadido por la incrustación de origen como características de entrada adicionales para nuestro regresor era, en el mejor de los casos, insignificante. Una variación de nuestro modelo estimador HTER entrenado con el vector $\mathbf{x} = [\mathbf{h}; \mathbf{s}; \mathbf{r}; \mathbf{h} \odot \mathbf{s}; \mathbf{h} \odot \mathbf{r}; |\mathbf{h} - \mathbf{s}|; |\mathbf{h} - \mathbf{r}|]$ como entrada para el modelo de solo avance logró mejorar el rendimiento a nivel de segmento en 8 de los 18 pares de idiomas descritos en la sección 5 a continuación y la mejora promedio en Tau de Kendall en esos entornos fue +0.0009. Como se señala en Zhao et al. (2020), si bien los modelos preentrenados multilingües son adaptativos a múltiples idiomas, el espacio de características entre idiomas está mal alineado. Con base en esto, decidimos a favor de excluir la incrustación de origen, decidimos a favor de excluir la incrustación de origen con la intuición de que la información más importante proviene de la incrustación de referencia y reducir el espacio de características permitiría al modelo centrarse más en la información relevante. Sin embargo, esto no niega el valor general de la fuente para nuestro modelo; donde incluimos características de combinación como $\mathbf{h} \odot \mathbf{s}$ y $|\mathbf{h} - \mathbf{s}|$ notamos ganancias en la correlación, como se explora más a fondo en la sección 5.5 a continuación.

2.4 Modelo de Clasificación de Traducción

Nuestro modelo de clasificación de traducción (Figura 2) recibe como entrada una tupla $\chi = (s, h^+, h^-, r)$ donde h^+ denota una hipótesis que fue clasificada más alto que otra hipótesis h^- . Luego pasamos χ a través de nuestro codificador multilingüe y capa de agrupamiento para obtener una incrustación de oración para cada segmento en el χ . Finalmente, utilizando las incrustaciones $\{s, h^+, h^-, r\}$, calculamos la pérdida de margen de tripletas (Schroff et al., 2015) en relación con la fuente y la referencia:

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

donde:

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

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

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

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

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

Finalmente, convertimos la distancia resultante en una puntuación de similitud limitada entre 0 y 1 de la siguiente manera:

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

3 Corpora

Para demostrar la efectividad de nuestras arquitecturas de modelo descritas (sección 2), entrenamos tres modelos de evaluación de MT donde cada modelo se dirige a un tipo diferente de juicio humano. Para entrenar estos modelos, utilizamos datos de tres corpora diferentes: el corpus QT21, el DARR de la tarea compartida de métricas WMT (2017 a 2019) y un corpus anotado MQM propietario.

3.1 El corpus QT21

El corpus QT21 es un conjunto de datos disponibles públicamente³ que contiene oraciones generadas por la industria de los dominios de tecnología de la información o ciencias de la vida (Specia et al., 2017). Este corpus contiene un total de 173K tuplas con la oración fuente, la referencia generada por humanos respectiva, la hipótesis de MT (ya sea de un MT estadístico basado en frases o de un MT neuronal), y el 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 al calcular la tasa de edición de traducción (TER) (Snover et al., 2006) entre la hipótesis de MT y el PE correspondiente. Finalmente, después de calcular el HTER para cada MT, construimos un conjunto de datos de entrenamiento $D = \{s_i, h_i, r_i, y_i\}_{i=1}^N$, donde s_i denota el texto fuente, h_i denota la hipótesis de MT, r_i la traducción de referencia y y_i la puntuación HTER para la hipótesis h_i . De esta manera, buscamos aprender una regresión $f(s, h, r) \rightarrow y$ que prediga el esfuerzo humano requerido para corregir la hipótesis al observar el fuente, la hipótesis y la referencia (pero no la hipótesis post-editada).

3.2 El corpus DARR de WMT

Desde 2017, los organizadores de la Tarea Compartida de Traducción de Noticias WMT (Barrault et al., 2019) han recopilado juicios humanos en forma de DAs de adecuación (Graham et al., 2013, 2014, 2017). Estos DAs se mapean en clasificaciones relativas (DARR) (Ma et al., 2019). Los datos resultantes para cada año (2017-19) forman un conjunto de datos $D = \{s_i, h_i^+, h_i^-, r_i\}_{i=1}^N$ donde h_i^+ denota una hipótesis "mejor" y h_i^- denota una "peor". Aquí buscamos aprender una función $r(s, h, r)$ tal que la puntuación asignada a h_i^+ sea estrictamente mayor que la puntuación asignada a h_i^- ($r(s_i, h_i^+, r_i) > r(s_i, h_i^-, r_i)$). Estos datos⁴ contienen un total de 24 pares de idiomas de alto y bajo recurso, como chino a inglés (zh-en) e inglés a 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.

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

⁴The raw data for each year of the WMT Metrics shared task is publicly available in the results page (2019 example: <http://www.statmt.org/wmt19/results.html>). Note, however, that in the README files it is highlighted that this data is not well documented and the scripts occasionally require custom utilities that are not available.

mensajes de chat que fueron anotados de acuerdo con las directrices establecidas en Burchardt y Lommel (2014). Estos datos contienen un total de 12 K tuplas, que cubren 12 pares de idiomas del inglés a: alemán (en-de), español (en-es), español latinoamericano (en-es-latam), francés (en-fr), italiano (en-it), japonés (en-ja), neerlandés (en-nl), portugués (en-pt), portugués brasileño (en-pt-br), ruso (en-ru), sueco (en-sv) y turco (en-tr). Tenga en cuenta que en este corpus el inglés siempre se considera como el idioma fuente, pero nunca como el idioma objetivo. Cada tupla consiste en una oración fuente, una referencia generada por humanos, una hipótesis de MT y su puntuación MQM, derivada de anotaciones de errores por uno (o más) anotadores entrenados. La métrica MQM a la que se hace referencia a lo largo de este documento es una métrica interna definida de acuerdo con el marco MQM (Lommel et al., 2014) (MQM). Los errores se anotan bajo una tipología interna definida en tres tipos principales de errores: 'Estilo', 'Fluidez' y 'Precisión'. Nuestras puntuaciones MQM varían de $-\infty$ a 100 y se definen como:

$$\text{MQM} = 100 - \frac{I_{\text{Minor}} + 5 \times I_{\text{Major}} + 10 \times I_{\text{Crit.}}}{\text{Sentence Length} \times 100} \quad (7)$$

donde I_{Minor} denota el número de errores menores, I_{Major} el número de errores mayores y $I_{\text{Crit.}}$ el número de errores críticos.

Nuestra métrica MQM tiene en cuenta la gravedad de los errores identificados en la hipótesis de 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 en 0. Como en la sección 3.1, construimos un conjunto de datos de entrenamiento

$D = \{s_i, h_i, r_i, y_i\}_{i=1}^N$, donde s_i denota el texto fuente, h_i denota la hipótesis de 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: una que regresa en HTER (COMET-HTER) entrenada con el corpus QT2 1, y otra que regresa en nuestra implementación propietaria de MQM (COMET-MQM) entrenada con nuestro corpus interno de MQM. Para el modelo de Clasificación de Traducción, descrito en la sección 2.4, entrenamos con el corpus WMT DARR de 2017 y 2018 (COMET-RANK). En esta sección, introducimos el entrenamiento.

configuración para estos modelos y la configuración de evaluación correspondiente.

4.1 Configuración del Entrenamiento

Las dos versiones de los Estimadores (COMET-HTER/MQM) comparten la misma configuración de entrenamiento y hiperparámetros (los detalles se incluyen en los Apéndices). Para el entrenamiento, cargamos el codificador preentrenado e inicializamos tanto la capa de agrupamiento como el regresor de avance. Mientras que los escalares por capa α de la capa de agrupamiento se establecen inicialmente en cero, los pesos del avance se inicializan aleatoriamente. Durante el entrenamiento, dividimos los parámetros del modelo en dos grupos: los parámetros del codificador, que incluyen el modelo del codificador y los escalares de α ; y los parámetros del regresor, que incluyen los parámetros de la red de avance superior. Aplicamos descongelamiento gradual y tasas de aprendizaje discriminativas (Howard y Ruder, 2018), lo que significa que el modelo del codificador está congelado durante una época mientras que el avance se optimiza con una tasa de aprendizaje de $3e-5$. Después de la primera época, todo el modelo se ajusta finamente, pero la tasa de aprendizaje para los parámetros del codificador se establece en $1e-5$ para evitar el olvido catastrófico.

En contraste con los dos Estimadores, para el modelo COMET-RANK afinamos desde el principio. Además, dado que este modelo no añade ningún nuevo parámetro además de los escalares de capa α sobre XLM-RoBERTa (base), utilizamos una única tasa de aprendizaje de $1e-5$ para todo el modelo.

4.2 Configuración de Evaluación

Utilizamos los datos de prueba y la configuración de la Tarea Compartida de Métricas WMT 2019 (Ma et al., 2019) para comparar los modelos COMET con las mejores presentaciones de la tarea compartida y otras métricas recientes de vanguardia como BERTSCORE y BLEURT.⁵ El método de evaluación utilizado es la formulación oficial similar a Tau de Kendall, τ , de la Tarea Compartida de Métricas WMT 2019 (Ma et al., 2019) definida como:

$$\tau = \frac{\text{Concordant} - \text{Discordant}}{\text{Concordant} + \text{Discordant}} \quad (8)$$

donde *Concordant* es el número de veces que una métrica asigna una puntuación más alta a la hipótesis "mejor" h^+ y *Discordant* es el número de veces que una métrica asigna una puntuación más alta a la hipótesis "peor"

⁵To ease future research we will also provide, within our framework, detailed instructions and scripts to run other metrics such as CHRF, BLEU, BERTSCORE, and BLEURT

Tabla 1: Correlaciones de Tau de Kendall (τ) en pares de idiomas con inglés como fuente para el corpus de métricas DARR WMT19. Para BERTSCORE, informamos los resultados con el modelo de codificador predeterminado para una comparación completa, pero también con XLM-RoBERTa (base) para la equidad con nuestros modelos. Los valores reportados para YiSi-1 se toman directamente del artículo de la tarea compartida (Ma et al., 2019).

Metric	en-cs	en-de	en-fi	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 (default)	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

h^- o las puntuaciones asignadas a ambas hipótesis son las mismas.

Como se mencionó en los hallazgos de (Ma et al., 2019), las correlaciones a nivel de segmento de todas las métricas enviadas fueron frustrantemente bajas. Además, todas las métricas enviadas mostraron una falta dramática de capacidad para clasificar correctamente los sistemas de MT fuertes. Para evaluar si nuestros nuevos modelos de evaluación de MT abordan mejor este problema, seguimos el conjunto de evaluación descrito utilizado en el análisis presentado en (Ma et al., 2019), donde se examinan los niveles de correlación para porciones de los datos de DARR que incluyen solo los 10, 8, 6 y 4 sistemas de MT principales.

5 Resultados

5.1 Del inglés al X

La Tabla 1 muestra los resultados para los ocho pares de idiomas con inglés como fuente. Contrastamos nuestros tres modelos COMET contra métricas de referencia como BLEU y CHRF, la métrica ganadora de la tarea de 2019 YISI-1, así como la más reciente BERTSCORE. Observamos que en general nuestros tres modelos entrenados con el marco COMET superan, a menudo por márgenes significativos, todas las demás métricas. Nuestro modelo DARR Ranker supera a los dos Estimadores en siete de los ocho pares de idiomas. Además, aunque el Estimador MQM está entrenado solo con 12K segmentos anotados, rinde aproximadamente al mismo nivel que el Estimador HTER para la mayoría de los pares de idiomas, y supera a todas las demás métricas en en-ru.

5.2 De X al inglés

La Tabla 2 muestra los resultados para los siete pares de idiomas hacia el inglés. Nuevamente, contrastamos nuestros tres modelos COMET contra métricas de referencia como BLEU y CHRF, la métrica ganadora de la tarea de 2019 YISI-1, como

así como las métricas BERTSCORE y BLEURT publicadas recientemente. Como se muestra en la Tabla 1, el modelo DARR muestra fuertes correlaciones con los juicios humanos, superando la métrica BLEURT específica para inglés propuesta recientemente en cinco de siete pares de idiomas. Nuevamente, el Estimador MQM muestra resultados sorprendentemente fuertes a pesar de que este modelo fue entrenado con datos que no incluían inglés como objetivo. Aunque el codificador utilizado en nuestros modelos entrenados es altamente multilingüe, hipotetizamos que este poderoso resultado "zero-shot" se debe a la inclusión de la fuente en nuestros modelos.

5.3 Pares de idiomas que no involucran inglés

Los tres modelos COMET fueron entrenados con datos que involucran inglés (ya sea como fuente o como objetivo). Sin embargo, para demostrar que nuestras métricas se generalizan bien, las probamos en los tres pares de idiomas de WMT 2019 que no incluyen inglés ni en la fuente ni en el objetivo. Como se puede ver en la Tabla 3, nuestros resultados son consistentes con las observaciones en las Tablas 1 y 2.

5.4 Robustez ante MT de Alta Calidad

Para el análisis, utilizamos el corpus DARR de la tarea compartida de 2019 y evaluamos en el subconjunto de los datos de los sistemas de MT de mejor rendimiento para cada par de idiomas. Incluimos pares de idiomas para los cuales pudimos recuperar datos de al menos diez sistemas de MT diferentes (es decir, todos menos kk-en y gu-en). Contrastamos con el fuerte BERTSCORE y BLEURT propuestos recientemente, con BLEU como línea base. Los resultados se presentan en la Figura 3. Para los pares de idiomas donde el inglés es el objetivo, nuestros tres modelos son mejores o competitivos con todos los demás; donde el inglés es la fuente, notamos que en general nuestras métricas superan el rendimiento de otros.

Tabla 2: Correlaciones de Tau de Kendall (τ) en pares de idiomas con inglés como objetivo para el corpus de métricas DARR de WMT19. En cuanto a BERTSCORE, para BLEURT informamos resultados para dos modelos: el modelo base, que es comparable en tamaño con el codificador que utilizamos, y el modelo grande que es el doble de tamaño.

Metric	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 (default)	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: Correlaciones de Tau de Kendall (τ) en pares de idiomas que no involucran inglés para el corpus de métricas DARR WMT19.

Metric	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 (default)	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 con solo 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 de la entrada del idioma fuente en la capacidad de nuestros modelos para aprender predicciones precisas, entrenamos dos versiones de nuestro modelo DARR Ranker: una que utiliza solo la referencia y otra que utiliza tanto la referencia como la fuente. Ambos modelos fueron entrenados utilizando el corpus WMT 2017 que solo incluye pares de idiomas del inglés (en-de, en-cs, en-fi, en-tr). En otras palabras, aunque el inglés nunca se observó como un idioma objetivo durante el entrenamiento de ambas variantes del modelo, el entrenamiento de la segunda variante incluye incrustaciones de fuente en inglés. Luego probamos estas dos variantes del modelo en el corpus WMT 2018 para estos pares de idiomas y para las direcciones inversas (con la excepción de en-cs porque cs-en no existe para WMT 2018). Los resultados en la Tabla

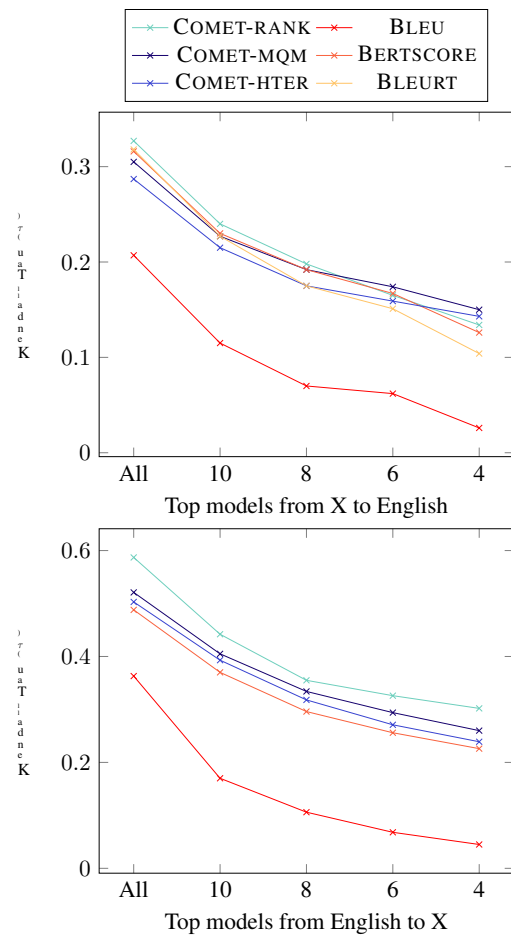


Figura 3: Rendimiento de métricas sobre todos y los mejores (10, 8, 6 y 4) sistemas de MT.

4 muestra claramente que para la arquitectura de clasificación de traducción, incluir la fuente mejora la correlación general con los juicios humanos. Además, la inclusión de la fuente expuso la segunda variante del modelo a los embeddings en inglés, que es

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

Metric	en-cs	en-de	en-fi	en-tr	cs-en	de-en	fi-en	tr-en
COMET-RANK (ref. only)	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
$\Delta\tau$	0.051	0.035	0.041	0.024	0.107	0.155	0.119	0.132

reflejado en un $\Delta\tau$ más alto para los pares de idiomas con inglés como objetivo.

6 Reproducibilidad

Publicaremos tanto la base de código del marco COMET como 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 base reportadas.⁶ Todos los modelos reportados en este documento fueron entrenados en una sola GPU Tesla T4 (16GB). Además, nuestro marco se basa en PyTorch Lightning (Falcon, 2019), un envoltorio ligero de PyTorch, que fue creado para una flexibilidad y reproducibilidad máximas.

7 Trabajo Relacionado

Las métricas clásicas de evaluación de MT se caracterizan comúnmente como métricas de coincidencia de n -gramas porque, utilizando características elaboradas a mano, estiman la calidad de MT contando el número y la fracción de n -gramas que aparecen simultáneamente en una hipótesis de traducción candidata y una o más referencias humanas. Métricas como BLEU (Papineni et al., 2002), METEOR (Lavie y Denkowski, 2009) y CHRF (Popović, 2015) han sido ampliamente estudiadas y mejoradas (Koehn et al., 2007; Popović, 2017; Denkowski y Lavie, 2011; Guo y Hu, 2019), pero, por diseño, generalmente no logran reconocer y capturar la similitud semántica más allá del nivel léxico.

En los últimos años, las incrustaciones de palabras (Mikolov et al., 2013; Pennington et al., 2014; Peters et al., 2018; Devlin et al., 2019) han surgido como una alternativa comúnmente utilizada al emparejamiento de n -gramas para capturar la similitud semántica de las palabras. 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 BERTSCORE (Zhang et al., 2020) crean alineaciones suaves entre la referencia y la hipótesis.

⁶These will be hosted at: <https://github.com/Unbabel/COMET>

en un espacio de incrustación y luego calcular una puntuación que refleje la similitud semántica entre esos segmentos. Sin embargo, los juicios humanos como DA y MQM capturan mucho más que solo la similitud semántica, lo que resulta en un límite superior de correlación entre los juicios humanos y las puntuaciones producidas por tales métricas.

Las métricas aprendibles (Shimanaka et al., 2018; Mathur et al., 2019; Shimanaka et al., 2019) intentan optimizar directamente la correlación con los juicios humanos y han mostrado resultados prometedores recientemente. BLEURT (Sellam et al., 2020), una métrica aprendible basada en BERT (Devlin et al., 2019), afirma tener un rendimiento de vanguardia durante los últimos 3 años de la tarea compartida de métricas WMT. Debido a que BLEURT se basa en English-BERT (Devlin et al., 2019), solo se puede utilizar cuando el inglés es el idioma objetivo, lo que limita su aplicabilidad. Además, hasta donde sabemos, todas las métricas aprendibles propuestas anteriormente se han centrado en optimizar DA, que, debido a la escasez de anotados, puede resultar inherentemente ruidoso (Ma et al., 2019).

La evaluación de MT sin referencia, también conocida como Estimación de Calidad (QE), históricamente ha retrocedido a menudo en HTER para la evaluación a nivel de segmento (Bojar et al., 2013, 2014, 2015, 2016, 2017a). Más recientemente, MQM se ha utilizado para la evaluación a nivel de documento (Specia et al., 2018; Fonseca et al., 2019). Al aprovechar codificadores multilingües preentrenados altamente efectivos como multilingual BERT (Devlin et al., 2019) y XLM (Conneau y Lample, 2019), los sistemas de QE han mostrado correlaciones auspiciosas con los juicios humanos (Kepler et al., 2019a). Concurrentemente, el marco OpenKiwi (Kepler et al., 2019b) ha facilitado a los investigadores avanzar en el campo y construir modelos de QE más robustos.

8 Conclusiones y Trabajo Futuro

En este artículo presentamos COMET, un nuevo marco neuronal para entrenar modelos de evaluación de MT que pueden servir como métricas automáticas y ser fácilmente

adaptado y optimizado para diferentes tipos de juicios humanos sobre la calidad de la MT.

Para demostrar la efectividad de nuestro marco, buscamos abordar los desafíos reportados en la Tarea Compartida de Métricas WMT 2019 (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 los parámetros y el tiempo de inferencia. Una vía principal para el trabajo futuro en COMET examinará el impacto de soluciones más compactas como DistilBERT (Sanh et al., 2019).

Además, aunque esbozamos la posible importancia del texto fuente arriba, señalamos que nuestro modelo COMET-RANK pondera la fuente y la referencia de manera diferente durante la inferencia, pero de manera equitativa en su función de pérdida de entrenamiento. Trabajos futuros investigarán la optimalidad de esta formulación y examinarán más a fondo la interdependencia de las diferentes entradas.

Agradecimientos

Agradecemos a André Martins, Austin Matthews, Fabio Kepler, Daan Van Stigt, Miguel Vera y a 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 Umbra-4EU, supervisados por ANI bajo los números de contrato 045909 y 042671, respectivamente.

Referencias

Mikel Artetxe y Holger Schwenk. 2019. Embeddings de oraciones masivamente multilingües para transferencia cruzada sin entrenamiento y más allá. *Transactions of the Association for Computational Linguistics*, 7:597–610.

Loïc Barrault, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malasi, Christof Monz, Mathias Müller, Santanu Pal, Matt Post, y Marcos Zampieri. 2019. Hallazgos de la conferencia de 2019 sobre traducción automática (WMT19). En *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, páginas 1–61, Florencia, Italia. Asociación de Lingüística Computacional.

Ondřej Bojar, Christian Buck, Chris Callison-Burch, Christian Federmann, Barry Haddow, Philipp Koehn, Christof Monz, Matt Post, Radu Soricut, y

Lucia Specia. 2013. Resultados del Taller de Traducción Automática Estadística 2013. En *Proceedings of the Eighth Workshop on Statistical Machine Translation*, páginas 1–44, Sofía, Bulgaria. Asociación de Lingüística Computacional.

Ondřej Bojar, Christian Buck, Christian Federmann, Barry Haddow, Philipp Koehn, Johannes Leveling, Christof Monz, Pavel Pecina, Matt Post, Herve Saint-Amand, Radu Soricut, Lucia Specia y Aleš Tamchyna. 2014. Resultados del taller de 2014 sobre traducción automática estadística. En *Proceedings of the Ninth Workshop on Statistical Machine Translation*, páginas 12–58, Baltimore, Maryland, EE. UU. Asociación de Lingüística Computacional.

Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Shujian Huang, Matthias Huck, Philipp Koehn, Qun Liu, Varvara Logacheva, Christof Monz, Matteo Negri, Matt Post, Raphael Rubino, Lucia Specia y Marco Turchi. 2017a. Hallazgos de la conferencia de 2017 sobre traducción automática (WMT17). En *Proceedings of the Second Conference on Machine Translation*, páginas 169–214, Copenhague, Dinamarca. Asociación de Lingüística Computacional.

Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Aurélie Névoul, Mariana Neves, Martin Popel, Matt Post, Raphael Rubino, Carolina Scarton, Lucia Specia, Marco Turchi, Karin Verspoor, y Marcos Zampieri. 2016. Hallazgos de la conferencia de 2016 sobre traducción automática. En *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, páginas 131–198, Berlín, Alemania. Asociación de Lingüística Computacional.

Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Barry Haddow, Matthias Huck, Chris Hokamp, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Matt Post, Carolina Scarton, Lucia Specia y Marco Turchi. 2015. Resultados del taller de 2015 sobre traducción automática estadística. En *Proceedings of the Tenth Workshop on Statistical Machine Translation*, páginas 1–46, Lisboa, Portugal. Asociación de Lingüística Computacional.

Ondřej Bojar, Yvette Graham y Amir Kamran. 2017b. Resultados de la tarea compartida de métricas WMT17. En *Proceedings of the Second Conference on Machine Translation*, páginas 489–513, Copenhague, Dinamarca. Asociación de Lingüística Computacional.

Aljoscha Burchardt y Arle Lommel. 2014. Directrices prácticas para el uso de MQM en la investigación científica sobre la calidad de la traducción. (fecha de acceso: 2020-05-26).

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco

- Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer y Veselin Stoyanov. 2019. Aprendizaje de representación cruzada no supervisado a gran escala. *arXiv preprint arXiv:1911.02116*.
- Alexis Conneau y Guillaume Lample. 2019. Preentrenamiento de modelos de lenguaje multilingües. En H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox y R. Garnett, editores, *Advances in Neural Information Processing Systems 32*, páginas 7059–7069. Curran Associates, Inc.
- Michael Denkowski y Alon Lavie. 2011. Meteor 1.3: Métrica automática para la optimización y evaluación confiables de sistemas de traducción automática. En *Proceedings of the Sixth Workshop on Statistical Machine Translation*, páginas 85–91, Edimburgo, Escocia. Asociación de Lingüística Computacional.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee y Kristina Toutanova. 2019. BERT: Pre-entrenamiento de transformadores bidireccionales profundos para la comprensión del lenguaje. En *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, páginas 4171–4186, Minneapolis, Minnesota. Asociación de Lingüística Computacional.
- WA Falcon. 2019. PyTorch Lightning: El envoltorio ligero de PyTorch para la investigación de IA de alto rendimiento. *GitHub*.
- Erick Fonseca, Lisa Yankovskaya, André F. T. Martins, Mark Fishel y Christian Federmann. 2019. Resultados de las tareas compartidas de WMT 2019 sobre estimación de calidad. En *Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2)*, páginas 1–10, Florencia, Italia. Asociación de Lingüística Computacional.
- Yvette Graham, Timothy Baldwin, Alistair Moffat y Justin Zobel. 2013. Escalas de medición continuas en la evaluación humana de la traducción automática. En *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, páginas 33–41, Sofía, Bulgaria. Asociación de Lingüística Computacional.
- Yvette Graham, Timothy Baldwin, Alistair Moffat y Justin Zobel. 2014. ¿Está mejorando la traducción automática con el tiempo? En *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, páginas 443–451, Gotemburgo, Suecia. Asociación de Lingüística Computacional.
- Yvette Graham, Timothy Baldwin, Alistair Moffat y Justin Zobel. 2017. ¿Pueden los sistemas de traducción automática ser evaluados solo por la multitud? *Natural Language Engineering*, 23(1):330.
- Yinuo Guo y Junfeng Hu. 2019. Meteor++ 2.0: Adoptar el conocimiento de paráfrasis a nivel sintáctico en la evaluación de traducción automática. En *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, páginas 501–506, Florencia, Italia. Asociación de Lingüística Computacional.
- Jeremy Howard y Sebastian Ruder. 2018. Ajuste fino de un modelo de lenguaje universal para la clasificación de texto. En *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, páginas 328–339, Melbourne, Australia. Asociación de Lingüística Computacional.
- Fabio Kepler, Jonay Trénous, Marcos Treviso, Miguel Vera, António Góis, M. Amin Farajian, António V. Lopes, y André F. T. Martins. 2019a. La participación de Unbabel en la tarea compartida de estimación de calidad de traducción WMT19. En *Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2)*, páginas 78–84, Florencia, Italia. Asociación de Lingüística Computacional.
- Fabio Kepler, Jonay Trénous, Marcos Treviso, Miguel Vera y André F. T. Martins. 2019b. OpenKiwi: Un marco de código abierto para la estimación de calidad. En *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, páginas 117–122, Florencia, Italia. Asociación de Lingüística Computacional.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zenith, Chris Dyer, Ondřej Bojar, Alexandra Constantin y Evan Herbst. 2007. Moses: herramienta de código abierto para la traducción automática estadística. En *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, páginas 177–180, Praga, República Checa. Asociación de Lingüística Computacional.
- Dan Kondratyuk y Milan Straka. 2019. 75 idiomas, 1 modelo: Analizando dependencias universales de manera universal. En *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, páginas 2779–2795, Hong Kong, China. Asociación de Lingüística Computacional.
- Alon Lavie y Michael Denkowski. 2009. La métrica meteor para la evaluación automática de la traducción automática. *Machine Translation*, 23:105–115.
- Chi-kiu Lo. 2019. YiSi - una métrica unificada de evaluación y estimación de calidad de MT semántica para lenguas con diferentes niveles de recursos disponibles. En *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, páginas 507–513, Florencia, Italia. Asociación de Lingüística Computacional.
- Arle Lommel, Aljoscha Burchardt y Hans Uszkoreit. 2014. Métricas de calidad multidimensional (MQM): A

- marco para declarar y describir métricas de calidad de traducción. *Tradumtica: technologies de la traducci*, 0:455–463.
- Qingsong Ma, Ondřej Bojar y Yvette Graham. 2018. Resultados de la tarea compartida de métricas WMT18: Tanto los caracteres como los embeddings logran un buen rendimiento. En *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, páginas 671–688, Bélgica, Bruselas. Asociación de Lingüística Computacional.
- Qingsong Ma, Johnny Wei, Ondřej Bojar y Yvette Graham. 2019. Resultados de la tarea compartida de métricas WMT19: Los sistemas de MT a nivel de segmento y fuertes plantean grandes desafíos. En *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, páginas 62–90, Florencia, Italia. Asociación de Lingüística Computacional.
- Nitika Mathur, Timothy Baldwin y Trevor Cohn. 2019. Poniendo la evaluación en contexto: Las incrustaciones contextuales mejoran la evaluación de la traducción automática. En *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, páginas 2799–2808, Florencia, Italia. Asociación de Lingüística Computacional.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado y Jeff Dean. 2013. Representaciones distribuidas de palabras y frases y su composicionalidad. En *Advances in Neural Information Processing Systems 26*, 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 automática. En *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, páginas 311–318, Filadelfia, Pennsylvania, EE. UU. Asociación de Lingüística Computacional.
- Jeffrey Pennington, Richard Socher y Christopher Manning. 2014. Glove: vectores globales para la representación de palabras. En *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, páginas 1532–1543, Doha, Catar. Asociación de Lingüística Computacional.
- Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee y Luke Zettlemoyer. 2018. Representaciones de palabras contextualizadas profundas. En *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, páginas 2227–2237, Nueva Orleans, Luisiana. Asociación de Lingüística Computacional.
- Telmo Pires, Eva Schlinger y Dan Garrette. 2019. ¿Qué tan multilingüe es BERT multilingüe? En *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, páginas 4996–5001, Florencia, Italia. Asociación de Lingüística Computacional.
- Maja Popović. 2015. chrF: puntuación f de n-gramas de caracteres para la evaluación automática de MT. En *Proceedings of the Tenth Workshop on Statistical Machine Translation*, páginas 392–395, Lisboa, Portugal. Asociación de Lingüística Computacional.
- Maja Popović. 2017. chrF++: palabras que ayudan a los n-grams de caracteres. En *Proceedings of the Second Conference on Machine Translation*, páginas 612–618, Copenhague, Dinamarca. Asociación de Lingüística Computacional.
- Nils Reimers e Iryna Gurevych. 2019. Sentence-BERT: incrustaciones de oraciones utilizando redes Siamese BERT. En *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, páginas 3982–3992, Hong Kong, China. Asociación de Lingüística Computacional.
- 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. *arXiv preprint arXiv:1910.01108*.
- F. Schroff, D. Kalenichenko y J. Philbin. 2015. Facenet: Un embedding unificado para el reconocimiento y agrupamiento de rostros. En *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, páginas 815–823.
- Thibault Sellam, Dipanjan Das y Ankur Parikh. 2020. BLEURT: Aprendiendo métricas robustas para la generación de texto. En *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, páginas 7881–7892, En línea. Asociación de Lingüística Computacional.
- Christophe Servan, Alexandre Bérard, Zied Elloumi, Hervé Blanchon y Laurent Besacier. 2016. Word2Vec vs DBnary: ¿Aumentando METEOR utilizando representaciones vectoriales o recursos léxicos? En *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, páginas 1159–1168, Osaka, Japón. El Comité Organizador de COLING 2016.
- Hiroki Shimanaka, Tomoyuki Kajiwar y Mamoru Ko machi. 2018. RUSE: Regresor utilizando incrustaciones de oraciones para la evaluación automática de traducción automática. En *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, páginas 751–758, Bélgica, Bruselas. Asociación de Lingüística Computacional.
- Hiroki Shimanaka, Tomoyuki Kajiwar y Mamoru Ko machi. 2019. Evaluación de Traducción Automática con Regresor BERT. *arXiv preprint arXiv:1907.12679*.

Matthew Snover, Bonnie Dorr, Richard Schwartz, Lina Micciulla y John Makhoul. 2006. Un estudio de la tasa de edición de traducción con anotación humana rígida. En *In Proceedings of Association for Machine Translation in the Americas*, páginas 223–231.

Lucia Specia, Frédéric Blain, Varvara Logacheva, Ramón Astudillo y André F. T. Martins. 2018. Resultados de la tarea compartida de WMT 2018 sobre estimación de calidad. En *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, páginas 689–709, Bélgica, Bruselas. Asociación de Lingüística Computacional.

Lucia Specia, Kim Harris, Frédéric Blain, Aljoscha Burchardt, Vivien Macketanz, Inguna Skadina, Matteo Negri, y Marco Turchi. 2017. Calidad de traducción y productividad: Un estudio sobre lenguas con morfología rica. En *Machine Translation Summit XVI*, páginas 5–71, Nagoya, Japón.

Kosuke Takahashi, Katsuhito Sudoh y Satoshi Nakamura. 2020. Evaluación automática de traducción automática utilizando entradas en el idioma fuente y un modelo de lenguaje multilingüe. En *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, páginas 3553–3558, En línea. Asociación de Lingüística Computacional.

Andre Tättar y Mark Fishel. 2017. bleu2vec: la métrica dolorosamente familiar en esteroides de espacio vectorial continuo. En *Proceedings of the Second Conference on Machine Translation*, páginas 619–622, Copenhague, Dinamarca. Asociación de Lingüística Computacional.

Ian Tenney, Dipanjan Das y Ellie Pavlick. 2019. BERT redescubre la pipeline clásica de PLN. En *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, páginas 4593–4601, Florencia, Italia. Asociación de Lingüística Computacional.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger y Yoav Artzi. 2020. Bertscore: Evaluando la generación de texto con bert. En *International Conference on Learning Representations*.

Wei Zhao, Goran Glavaš, Maxime Peyrard, Yang Gao, Robert West y Steffen Eger. 2020. Sobre las limitaciones de los codificadores multilingües según lo expuesto por la evaluación de traducción automática sin referencia. En *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, páginas 1656–1671, En línea. Asociación de Lingüística Computacional.

Wei Zhao, Maxime Peyrard, Fei Liu, Yang Gao, Christian M. Meyer y Steffen Eger. 2019. MoverScore: Evaluación de generación de texto con incrustaciones contextualizadas y distancia de movimiento de tierra. En *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, páginas 563–578, Hong

Kong, China. Asociación de Lingüística Computacional.

A Apéndices

En la Tabla 5 enumeramos los hiperparámetros utilizados para entrenar nuestros modelos. Antes de inicializar estos modelos, se estableció una semilla aleatoria en 3 en todas las bibliotecas que realizan operaciones "aleatorias" (torch, numpy, random y cuda).

Tabla 5: Hiperparámetros utilizados en nuestro marco COMET para entrenar los modelos presentados.

Hyper-parameter	COMET(Est- HTER /MQM)	COMET-RANK
Encoder Model	XLM-RoBERTa (base)	XLM-RoBERTa (base)
Optimizer	Adam (default parameters)	Adam (default parameters)
n frozen epochs	1	0
Learning rate	3e-05 and 1e-05	1e-05
Batch size	16	16
Loss function	MSE	Triplet Margin ($\epsilon = 1.0$)
Layer-wise dropout	0.1	0.1
FP precision	32	32
Feed-Forward hidden units	2304,1152	–
Feed-Forward activations	Tanh	–
Feed-Forward dropout	0.1	–

Tabla 6: Estadísticas para el corpus QT21.

	en-de	en-cs	en-lv	de-en
Total tuples	54000	42000	35474	41998
Avg. tokens (reference)	17.80	15.56	16.42	17.71
Avg. tokens (source)	16.70	17.37	18.39	17.18
Avg. tokens (MT)	17.65	15.64	16.42	17.78

Tabla 7: Estadísticas para el corpus DARR WMT 2017.

	en-cs	en-de	en-fi	en-lv	en-tr
Total tuples	32810	6454	3270	3456	247
Avg. tokens (reference)	19.70	22.15	15.59	21.42	17.57
Avg. tokens (source)	22.37	23.41	21.73	26.08	22.51
Avg. tokens (MT)	19.45	22.58	16.06	22.18	17.25

) e c n e r e r (s n e k o o t . g v A	n e h z	0 7 0 1 3	9 8 2 4		0 7 9 3
	n e u r	2 5 8 9 3	4 7 1 2	0 0 8 1 2	
	n e a	2 5 8 1 2	5 5 6 2	2 3 0 2	5 2 5 2
	n e k k	8 7 9	6 3 0 2	2 3 6 1	8 6 9 1
	n e u g	0 1 1 0 2	4 6 7 1	2 9 1 2	2 0 7 1
	n e h	9 7 1 2 3	5 5 8 1	9 4 2 1	6 7 7 1
	n e e d	5 6 3 5 8	9 2 0 2	4 4 8 1	2 2 0 2
) e c n e r e r (s n e k o o t . g v A					
) e c r u o s (s n e k o o t . g v A					
) T M (s n e k o o t . g v A					

. s t r a p e g a u g n a i h s i g n E . o i n i
 R R A D 9 1 0 2 T M W
 e h i r o r s c i s i a S . 8 e l b a T

) e c n e r e r (s n e k o o t . g v A	r e d	2 3 7 2	6 3 1 2	8 6 5 2
	s e e d	4 9 8 2	7 2 2 2	9 8 1 2
	e d r a	9 5 1	8 6 2 2	0 6 3 2
	h z n e	8 5 8 1	5 2 4 2	3 8 6
	u r n e	4 3 5 2	9 7 4 2	7 3 3 2
	t i n e	1 0 4 7 1	0 0 4 2	7 9 0 2
	k k n e	2 7 1 8 1	9 8 7 3	2 9 9 1
	u g n e	5 6 3 1	2 3 4 2	7 9 2 3
	f i n e	0 2 8 1 3	2 1 0 2	3 5 9 1
	e d t n e	0 4 8 9 9	5 6 9 2	8 9 4 2
) e c n e r e r (s n e k o o t . g v A	s e n e	8 7 1 7 2	2 9 4 2	0 6 2 2
) e c r u o s (s n e k o o t . g v A				
) T M (s n e k o o t . g v A				

. s t r a p e g a u g n a i h s i g n E . o i n i
 R R A D 9 1 0 2 T M W
 e h i r o r s c i s i a S . 9 e l b a T

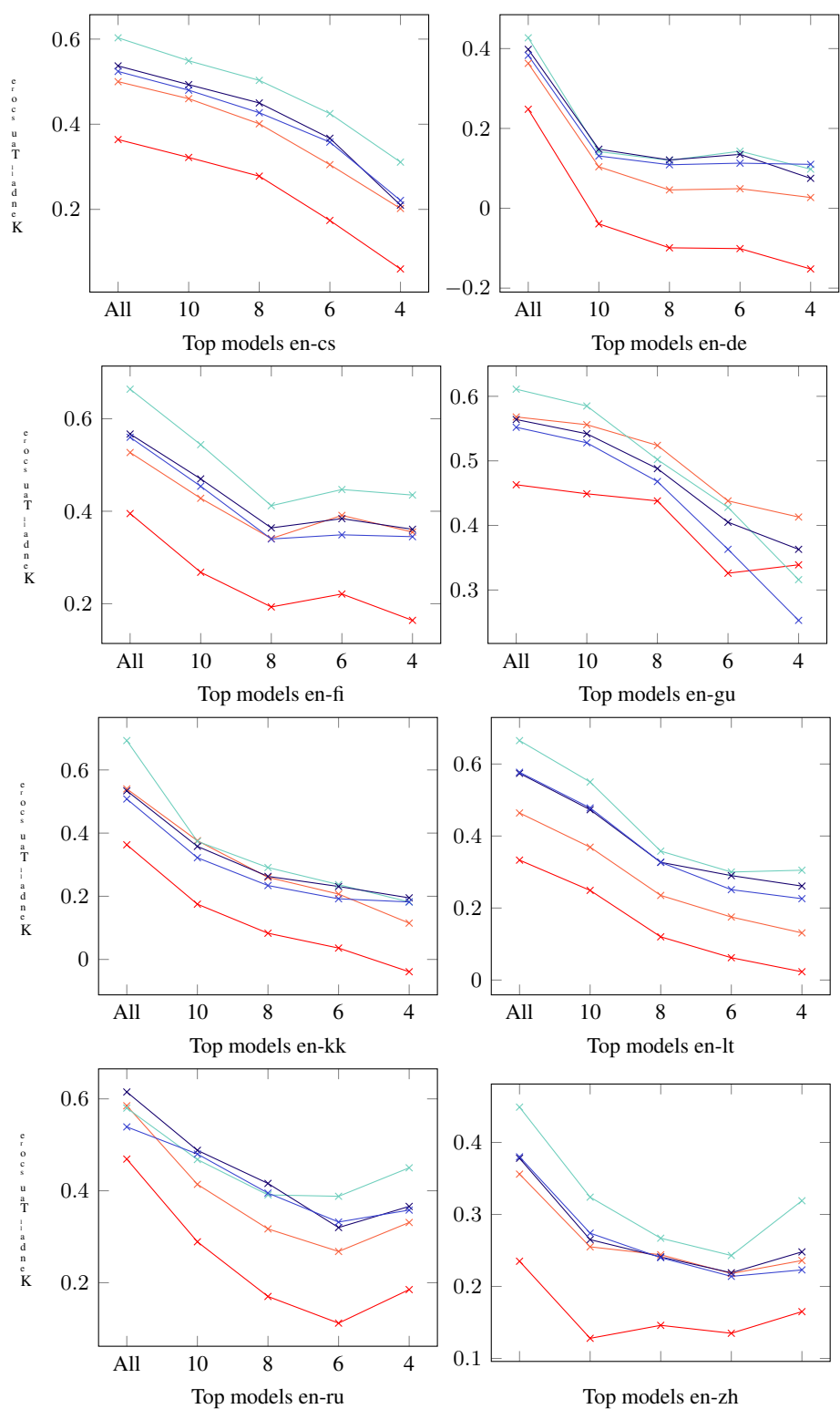


Tabla 12: Rendimiento de métricas sobre todos y los mejores (10, 8, 6 y 4) sistemas de MT para todos los pares de idiomas de inglés. El esquema de colores es el siguiente: COMET-RANK, COMET-HTER, COMET-MQM, BLEU, BERTSCORE

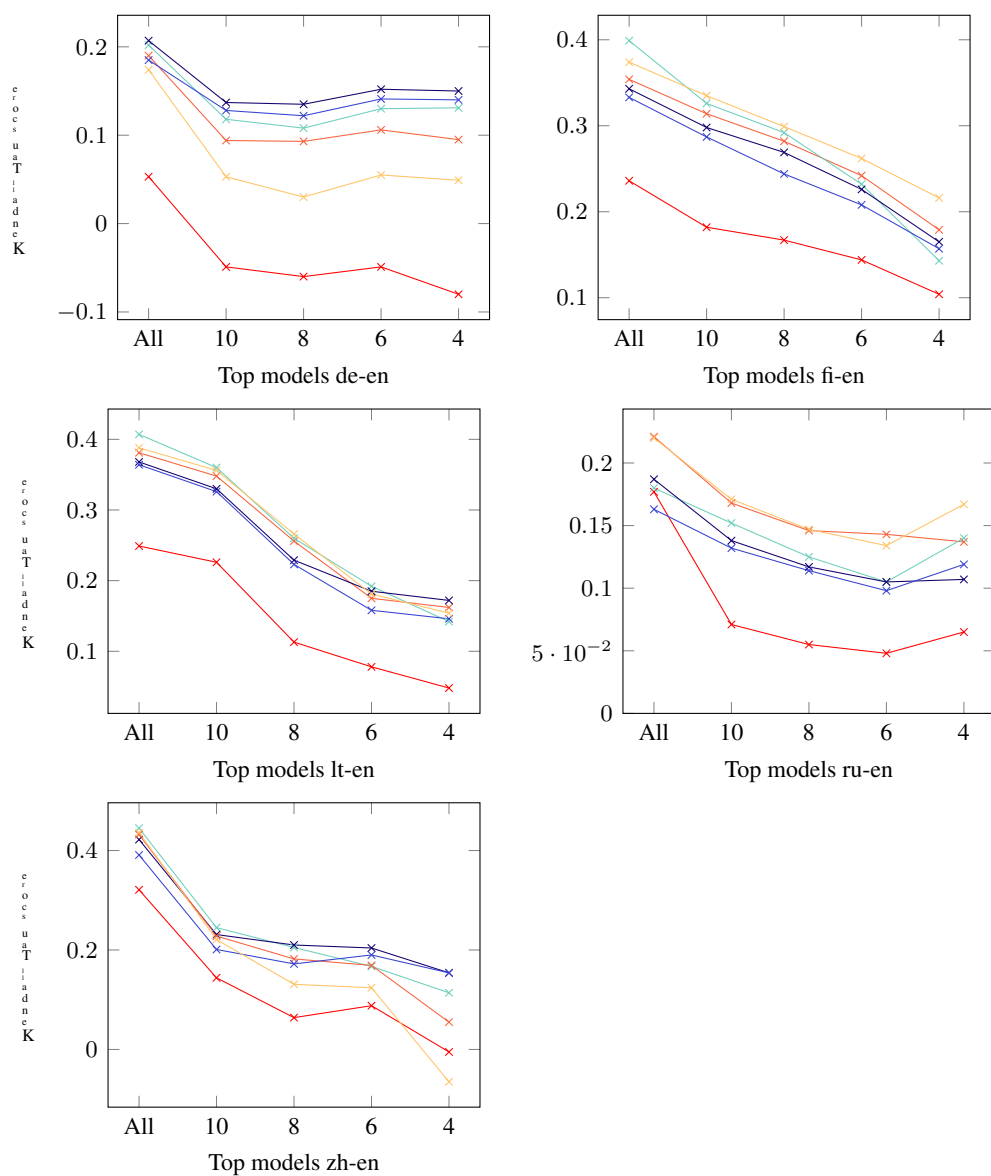


Tabla 13: Rendimiento de métricas sobre todos y los mejores (10, 8, 6 y 4) sistemas de MT para todos los pares de idiomas hacia el inglés. El esquema de colores es el siguiente: COMET-RANK, COMET-HTER, COMET-MQM, BLEU, BERTSCORE, BLEURT