

Ricardo Rei

Craig Stewart

Ana C Farinha

Alon Lavie

Unbabel IA

{craig.stewart, catarina.farinha, alon.lavie}@unbabel.com

@unbabel.com

## Résumé

Nous présentons COMET, un cadre neuronal pour la formation d'évaluation de la traduction automatique multilingue. Les modèles d'évaluation que nous obtenons atteignent un nouvel état de l'art de corrélation avec le jugement humain. Notre cadre tire parti des récentes percées dans le langage pré-entraîné multilingue et des modèles d'évaluation MT adaptables qui exploitent les informations à la fois de la source d'entrée et d'une traduction de référence en langue cible dans l'ordre pour prédire plus précisément la qualité de la MT. À la fois pour présenter notre cadre, nous formons trois modèles avec différents types de jugements humains : Évaluations Directes, Transmissions et Qualité d'Édition. Nos modèles atteignent un nouvel état de la performance de pointe sur le WMT 2019 Meta-task partagée et démontrent une robustesse à des systèmes haute performance.

## 1 Introduction

Historiquement, les mesures pour évaluer la qualité de la traduction automatique (TA) s'est appuyée sur l'évaluation de la similarité entre une hypothèse générée par MT et une traduction de référence générée par l'homme dans la langue cible. Les mesures traditionnelles se sont concentrées sur des caractéristiques de base, au niveau lexical, telles que le comptage du nombre de n-grammes correspondants entre le MT et la traduction de référence. Métriques telles que BLEU (Papineni et al., 2002) et METEOR (Sutskever et al., 2006) restent populaires en tant que moyen d'évaluer les systèmes MT en raison de leur calcul léger et rapide.

Les approches neuronales modernes à la traduction automatique donnent des résultats bien meilleurs. Une qualité de traduction plus élevée qui dévie souvent de transfert lexical monotone entre les langues. Pour cette raison, il est devenu de plus en plus évident que nous ne pouvons plus nous fier à des mesures telles que BLEU pour fournir une estimation précise de la qualité de la MT (Barrault et al., 2019).

Alors qu'un intérêt de recherche accru dans les réseaux neuronaux des méthodes pour former des modèles et des systèmes de MT ont conduit à une amélioration récente et spectaculaire en MT qualité, l'évaluation MT a pris du retard. Le MT la communauté de recherche repose toujours largement sur des méthodes obsolètes des mesures et aucune nouvelle norme largement adoptée n'a émergé. En 2019, la Traduction de Nouvelles WMT La Tâche Partagée a reçu un total de 153 systèmes MT soumissions (Barrault et al., 2019). Les Mesures La Tâche Partagée de la même année n'a vu que 24 soumissions, dont presque la moitié étaient des nouveaux venus dans le Tâche Partagée d'Estimation de Qualité, adaptée en tant que métriques (Ma et al., 2019).

Les résultats de la tâche susmentionnée sont élevés-éclaircissent deux défis majeurs à l'évaluation MT que nous cherchons à aborder ici (Ma et al., 2019). À savoir, que les mesures actuelles peinent à accuser-corrélent précisément avec le jugement humain à segment-niveau de ment et échouent à différencier adéquatement les systèmes MT les plus performants.

Dans cet article, nous présentons COMET<sup>1</sup>, un PyTorch-cadre basé pour la formation hautement multilingue et des modèles d'évaluation MT adaptables qui peuvent fonctionner comme des métriques. Notre cadre utilise cet avantage des avancées récentes dans le langage interlinguistique de modélisation (Artetxe et Schwenk, 2019; Devlin et al., 2019; Conneau et Lample, 2019; Conneau et al., 2019) pour générer des estimations de prédictions de jugements humains tels que les Évaluations Directes (ED) (Conneau et al., 2013), Taux d'Édition humaine (HTER) (Snover et al., 2006) et le taux d'édition humaine (HTER) (Snover et al., 2006) et des métriques conformes à la Qualité Multidimensionnelle (Lommel et al., 2014).

Inspiré par des travaux récents sur l'Estimation de Qualité (QE) qui a démontré qu'il est possible d'atteindre des niveaux élevés de corrélation avec les jugements humains même sans une traduction de référence (Fonseca et al., 2019), nous proposons une nouvelle approche pour incorporer-

<sup>1</sup> Translinguistique Métrique pour l'Évaluation de Traduction. Evaluation of

ation en traduisant l'entrée de la langue source dans notre évaluation MT des modèles d'ation. Traditionnellement, seuls les modèles QE ont a utilisé l'entrée de la source, tandis que l'évaluation MT les métriques d'évaluation se basent plutôt sur la traduction de référence-tion. Comme dans (Takahashi et al., 2020), nous montrons que, we show that l'utilisation d'un espace d'incorporation multilingue nous permet pour exploiter les informations provenant des trois entrées et démontrer la valeur ajoutée par la source en tant qu'entrée à nos modèles d'évaluation MT.

Pour illustrer l'efficacité et la flexibilité de le cadre COMET, nous formons trois modèles qui estimer différents types de jugements humains et montre des progrès prometteurs vers une meilleure cor-relation au niveau du segment et robustesse face à des niveaux élevés de qualité MT.

Nous publierons à la fois le cadre COMET et les modèles d'évaluation MT formés décrits dans ce papier à la communauté de recherche lors de la publication.

## 2 Architectures de Modèles

Les jugements humains sur la qualité de la TA proviennent généralement sous la forme de scores au niveau du segment, tels que DA, MQM et HTER. Pour DA, il est courant de convertir les scores en classements relatifs (DARR) lorsque le nombre d'annotations par segment est limité (Devlin et al., 2017b; Ma et al., 2018, 2019). Ceci signifie que, pour deux hypothèses MT  $h_i$  et  $h_j$  de la même source, si le score DA attribué à  $h_i$  est supérieur au score attribué à  $h_j$ ,  $h_i$  est considéré comme une "meilleure" hypothèse.<sup>2</sup> Pour encoder ces différences, notre cadre prend en charge deux architectures distinctes : Le modèle Estimator et le Modèle de classement de traduction. Le fondamental la différence entre eux est l'objectif de formation. Alors que l'Estimateur est formé pour régresser directement sur un score de qualité, le modèle de classement de traduction est formé pour minimiser la distance entre un "meilleur" hypothèse et ses deux références correspondantes et sa source originale. Les deux modèles sont composés d'un encodeur inter-langues et d'une couche de regroupement.

### 2.1 Encodeur Cross-lingual

Le bloc de construction principal de tous les modèles dans notre cadre est un pré-entraîné, cross-linguistique modèle tel que BERT multilingue (Devlin et al., 2019), XLM (Conneau et Lample, 2019) ou XLM-RoBERTa (Conneau et al., 2019). Ces modèles contiennent plusieurs couches de transformateur encodeur qui sont

<sup>2</sup>Dans la Tâche Partagée des Métriques WMT, si la différence entre les scores DA n'est pas supérieure à 25 points, ces segments sont exclus des données DARR.

formé à reconstruire des jetons masqués en les découvrant-établissant la relation entre ces jetons et le ceux environnants. Lorsqu'entraîné avec des données provenant de plusieurs langues cet objectif pré-entraîné a a été trouvé très efficace dans le contexte interlinguistique tâches telles que la classification de documents et le traitement naturel inférence de langue (Conneau et al., 2019), Gennep 2019), généralisant bien aux langues et scripts inconnus (Pires et al., 2019). Pour les expériences dans ce document, nous nous appuyons sur XLM-RoBERTa (base) comme notre encodeur modèle.

Étant donné une séquence d'entrée  $x = [x_0, x_1, \dots, x_n]$ , l'encodeur produit un encastrement  $e(\cdot)$  pour chaque jeton  $x_j$  et couche  $k$ . Dans notre cadre, nous appliquons ce processus à la source, Hypothèse MT, et référence afin de les cartographier dans un espace de caractéristiques partagé.

### 2.2 Couche de Pooling

Les incorporations générées par la dernière couche de la les encodeurs pré-entraînés sont généralement utilisés pour le réglage fin modèles à de nouvelles tâches. Cependant, (Tenney et al., 2019) a montré que différentes couches à l'intérieur du net-le travail peut capturer des informations linguistiques qui sont rel-évants pour différentes tâches en aval. Dans le cas de l'évaluation MT, (Zhang et al., 2020) a montré que showed that différentes couches peuvent atteindre différents niveaux de cor-relation et que l'utilisation de la dernière couche seulement souvent donne des résultats inférieurs. Dans ce travail, nous a utilisé l'approche décrite dans Peters et al. (2018) (2018) et rassembler des informations des sources les plus importantes en-coder les couches en une seule intégration pour chaque ken, ej, en utilisant un mécanisme d'attention couche par couche. Cet encastrement est ensuite calculé comme suit :

$$ex_j = \mu_{x_j} \quad (1)$$

où  $\mu$  est un coefficient de poids entraînable,  $E_j = [e_j^{(0)}, e_j^{(1)}, \dots, e_j^{(k)}]$  correspond au vecteur de incrustations de couche pour le jeton  $x_j$  et  $\alpha$  est un vecteur correspondant aux poids entraînables couche par couche. Dans afin d'éviter le surajustement à l'information con-tenu dans une seule couche, nous avons utilisé l'abandon de couche (Kendall et al., 2019), dans lequel, in which with a probabilité  $p$  est fixé à  $-\infty$ .

Enfin, comme dans (Reimers et Gurevych, 2019) (2019), nous appliquons une moyenne de regroupement au mot résultant des embeddings pour dériver un embedding de phrase pour chaque segment.

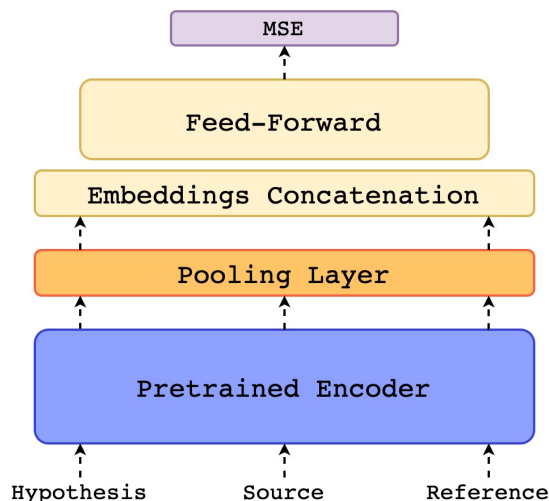


Figure 1 : Architecture du modèle d'estimation. La source, l'hypothèse et la référence sont codées indépendamment utilisant un encodeur pré-entraîné inter-langues. Le résultat des incorporations de mots sont ensuite transmises à travers une mise en commun couche pour créer une incorporation de phrase pour chaque segment. Enfin, les plongements de phrases résultants sont combiné et concaténé en un seul vecteur qui est transmis à un régresseur feed-forward. Le modèle entier est formé en minimisant l'Erreur Quadratique Moyenne (EQM).

## 2.3 Modèle Estimateur

Étant donné une incorporation de phrase  $d$ -dimensionnelle pour le source, l'hypothèse, et la référence, nous adoptons l'approche proposée dans RUSE (Shimanaka et al., 2018) et extrayez les caractéristiques combinées suivantes :

- Produit source élément par élément  $h \odot s$
- Produit de référence élément par élément  $h \odot r$
- Différence absolue élément par élément de la source :  $|h - s|$
- Différence de référence absolue élément par élément:  $|h - r|$

Ces caractéristiques combinées sont ensuite concaténées à l'incorporation de référence et à l'hypothèse en un seul vecteur  $x = [h; r; h \odot s; h \odot r; |h - s|; |h - r|]$  qui sert d'entrée à un régresseur feed-forward. La force de ces caractéristiques sont dans la mise en évidence des différences entre incorporations dans l'espace des caractéristiques sémantiques.

Le modèle est ensuite entraîné pour minimiser la moyenne l'erreur quadratique entre les scores prédits et évaluations de qualité (DA, HTER ou MQM). La figure 1 illustre l'architecture proposée.

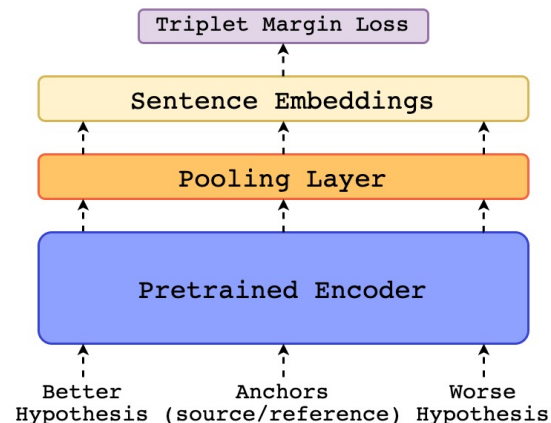


Figure 2: Architecture du modèle de classement de traduction. Cette architecture reçoit 4 segments : la source, le référence, une hypothèse "meilleure" et une "pire". Ces segments sont codés indépendamment en utilisant un pré-encodeur interlinguistique formé et une couche de regroupement sur haut. Enfin, en utilisant la perte de marge triplet (Schroff et al., 2015) nous optimisons l'espace d'incrustation résultant pour minimiser la distance entre l'hypothèse "meilleure" et les "ancres" (source et référence).

Notez que nous avons choisi de ne pas inclure la source brute incorporation  $s$  dans l'entrée concaténée. Tôt l'expérimentation a révélé que la valeur ajoutée par l'incorporation de la source en tant que caractéristiques d'entrée supplémentaires à notre régresseur était négligeable au mieux. Une variation sur notre modèle estimateur HTER formé avec le vecteur  $x = [h; s; r; h \odot s; h \odot r; |h - s|; |h - r|]$  as l'entrée dans le seul feed-forward réussit à stimuler-évaluant les performances au niveau du segment dans 8 des 18 langues paires de langues décrites dans la section 5 ci-dessous et le amélioration moyenne du Tau de Kendall dans ces ensembles-tings était +0.0009. Comme noté dans Zhao et al. (2020), (2020), tandis que les modèles pré-entraînés interlinguistiques sont adaptatifs à plusieurs langues, l'espace de caractéristiques entre les langues sont mal alignées. Sur cette base, nous décidés en faveur de l'exclusion de l'incorporation de la source sur l'intuition que l'information la plus importante provient de l'incorporation de référence et de la réduction la caractérisation de l'espace des fonctionnalités permettrait au modèle de se concentrer davantage sur les informations pertinentes. Cela ne fait pas cependant nie la valeur générale de la source à notre modèle; où nous incluons des caractéristiques de combinaison tel que  $h \odot s$  nous notons des gains en corrélation comme exploré plus loin dans la section 5.5 ci-dessous.

## 2.4 Modèle de Classement de Traduction

Notre modèle de classement de traduction (Figure 2) reçoit comme entrée un tuple  $(s, h^+, r, h^-)$  where  $h^+$  désigne une hypothèse qui a été classée plus haut que une autre hypothèse  $h^-$ . Nous faisons ensuite passer  $(s, h^+, r, h^-)$  à travers notre encodeur interlinguistique et notre couche de regroupement pour obtenir une incorporation de phrase pour chaque segment dans le tuple. Enfin, en utilisant les embeddings  $\{s, h^+, h^-, r\}$ , nous calculons la perte de marge triplet (Schroff et al., 2015) en relation avec la source et la référence :

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

où :

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

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

où  $d(\cdot)$  désigne la distance euclidienne entre un vecteur et un autre.  $\epsilon$  est une marge. Ainsi, pendant l'entraînement, le modèle optimise l'espace d'incorporation afin que la distance entre les ancres  $(s, r)$  et la "pire" hypothèse  $h^-$  soit plus grande que la distance entre les ancres et l'hypothèse "meilleure"  $h^+$ .

Pendant l'inférence, le modèle décrit reçoit un triplet  $(s, h, r)$  et une hypothèse. Le score de qualité attribué à  $h$  est la moyenne harmonique entre la distance à la source  $d(s, h)$  et la distance à la référence  $d(r, h)$  :

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

Enfin, nous convertissons la distance résultante en un score de similarité limité entre 0 et 1 comme suit :

suit comme :

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

## 3 Corpus

Pour démontrer l'efficacité de notre méthode décrite dans les architectures de modèles (section 2), nous avons utilisé trois modèles d'évaluation où chaque modèle cible un différent type de jugement humain. Pour former ces modèles, nous utilisons des données de trois corpus différents : le corpus QT21, le DARR du WMT Met-tâche partagée (2017 à 2019) et une propriétaire Corpus annoté MQM.

### 3.1 Le corpus QT21

Le corpus QT21 est un ensemble de données publiquement disponibles contenant des phrases générées par l'industrie provenant soit d'un domaine de la technologie de l'information ou des sciences de la vie (Spectral, 2017). Ce corpus contient un total de 173K tuples avec phrase source, respective référence générée par l'homme, hypothèse MT (soit d'une traduction statistique basée sur des phrases ou d'une neural MT), et MT post-édité (PE). La langue des paires représentées dans ce corpus sont : de l'anglais vers l'allemand (en-de), letton (en-lt) et tchèque (en-cs), et Allemand vers Anglais (de-en).

Le score HTER est obtenu en calculant le taux d'édit de traduction (TER) (Snover et al., 2006) entre l'hypothèse MT et le PE correspondant. Enfin, après avoir calculé le HTER pour chaque MT, nous avons construit un ensemble de données d'entraînement  $\{(s_i, r_i, y_i)\}_{i=1}^n$ , où  $s_i$  désigne le texte source,  $r_i$  désigne la référence, et  $y_i$  le score HTER pour l'hypothèse  $h_i$ . Nous cherchons à apprendre une régression  $f(s, h, r)$  qui prédit l'effort humain nécessaire pour corriger le hypothèse en regardant la source, hypothèse, et référence (mais pas l'hypothèse post-éditée).

### 3.2 Le corpus WMT DARR

Depuis 2017, les organisateurs du WMT News La Tâche Partagée de Traduction (Barrault et al., 2019) ont collecté des jugements humains sous la forme de publicités-équacy DAs (Graham et al., 2013, 2014, 2017). Ces DA sont ensuite cartographiés en rang relatif (DARR) (Ma et al., 2019).

Les données pour chaque année (2017-19) forment un ensemble de données  $D = \{(s_i, h_i, r_i)\}_{i=1}^N$  where  $s_i$  note un "meilleur" hypothèse et  $h_i$  signe un "pire". Ici nous cherchons à apprendre une fonction  $f(s, h, r)$  qui attribue à  $h_i$  un score strictement supérieur au score attribué à  $s_i$  ( $f(s_i, h_i, r_i) > f(s_i, s_i, r_i)$ ). Ces données forment un total de 24 valeurs hautes et basses-paires de langues ressources telles que le chinois vers l'anglais (zh-en) et anglais vers gujarati (en-gu).

### 3.3 Le corpus MQM

Le corpus MQM est une base de données interne propriétaire des traductions générées par MT du support client

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

Les données brutes pour chaque année du partage des métriques WMT la tâche est publiquement disponible sur la page des résultats (ex 2019-Exemple : <http://www.statmt.org/wmt19/results.html>). Notez cependant que dans les fichiers README, il est hautement souligné que ces données ne sont pas bien documentées et les scripts nécessitent occasionnellement des utilitaires personnalisés qui ne sont pas disponibles.



messages de chat qui ont été annotés selon le directives établies dans Burchardt et Lommel (2014). 014).

Ces données contiennent un total de 12K tuples, couvrant traduisant 12 paires de langues de l'anglais vers : l'allemand (En-de), espagnol (en-es), espagnol d'Amérique latine espagnol (en-es-latam), français (en-fr), italien (en-it), Japonais (en-ja), Néerlandais (en-nl), Portugais (en-pt), Portugais brésilien (en-pt-br), Russe (en-ru), Suédois (en-sv), et Turc (en-tr). Notez que dans ce corpus anglais est toujours considéré comme la langue source-guage, mais jamais comme langue cible. Chaque tuple consiste en une phrase source, une génération humaine référence, une hypothèse MT, et son score MQM, dérivé des annotations d'erreur par un (ou plusieurs) annotateurs formés. La métrique MQM à laquelle on fait référence tout au long de ce document est une métrique interne définie conformément au cadre MQM (Lommel et al., 201) (MQM). Les erreurs sont annotées sous une typologie interne définie sous trois principales er-types d'erreurs: 'Style', 'Fluidité' et 'Précision'. Notre Les scores MQM varient de 0 à 100 et sont de-

défini comme :

$$MQM = 100 - IM_{minor} + \frac{5 \times IM_{major} + 10 \times ICrit}{Longueur\ de\ la\ Phrase \times 100} + \frac{10 \times I_{Crit}}{100} \quad (7)$$

où IM<sub>minor</sub> désigne le nombre d'erreurs mineures, IM<sub>major</sub> le nombre d'erreurs majeures et ICrit. le nombre d'erreurs critiques.

Notre métrique MQM prend en compte la gravité-ité des erreurs identifiées dans l'hypothèse MT, conduisant à une mesure plus fine que HTER ou DA. Lorsqu'ils sont utilisés dans nos expériences, ces val-les valeurs ont été divisées par 100 et tronquées à 0. Comme dans la section 3.1, nous avons construit un ensemble de données d'entraînement  $D = \{s_i, h_i, r_i, y_i\}_{i=1}^N$  où  $s_i$  désigne le texte source,  $h_i$  l'hypothèse MT,  $r_i$  la traduction de référence, et  $y_i$  le MQM score for l'hypothèse  $h_i$ .

## 4 Expériences

Nous formons deux versions du modèle Estimator de-décrit dans la section 2.3: un qui régresse sur HTER (COMET-R) et un qui régresse sur HTER trained with the QT21 corpus, and un autre qui régresse sur notre mise en œuvre propriétaire-tation de MQM (COMET-MQM) formé avec le corpus interne MQM. Pour le Classement de Traduction modèle, décrit dans la section 2.4, nous nous entraînons avec le Corpus WMT de 2017 et 2018 (COMET-OMET-RANG). Dans cette section, nous introduisons la formation

configuration pour ces modèles et évaluation correspondante-configuration de l'installation.

### 4.1 Configuration de la Formation

Les deux versions des Estimateurs (COMET-HTER/MQM) partagent la même configuration d'entraînement et hyper-paramètres (les détails sont inclus dans l'Ap-pendices). Pour la formation, nous chargeons le pré-entraîné encodeur et initialiser à la fois la couche de pooling et le régresseur feed-forward. Tandis que le layer-wise les scalaires de la couche de pooling sont initialement définis à zéro, les poids provenant de la propagation avant sont ini-tialisés aléatoirement. Pendant la formation, nous divisons le paramètres du modèle en deux groupes : les paramètres de l'encodeur pa-ramètres, qui incluent le modèle de l'encodeur et le scalaires à partir de ces paramètres du régresseur, qui inclure les paramètres du flux d'information en avant supérieur réseau. Nous appliquons un dégel progressif et une discrimination-taux d'apprentissage innés (Howard et Ruder, 2018), ce qui signifie que le modèle d'encodeur est gelé pour un époque tandis que le feed-forward est optimisé avec un taux d'apprentissage de 3e-5. À la première époque, le tout le modèle est affiné mais le taux d'apprentissage pour les paramètres de l'encodeur sont réglés à 1e-5 afin de éviter l'oubli catastrophique.

En contraste avec les deux Estimateurs, pour le COMET-R nous affignons dès le départ. De plus, puisque ce modèle n'ajoute rien nouveaux paramètres au-dessus de XLM-RoBERTa (base) autre que les scalaires de couche de pooling, nous utilisons un seul taux d'apprentissage de 1e-5 pour l'ensemble du modèle.

### 4.2 Configuration de l'évaluation

Nous utilisons les données de test et la configuration du WMT 2019 Tâche Partagée de Métriques (Ma et al., 2019) afin in order to comparez les modèles COMET avec les performances supérieures-soumissions de la tâche partagée et autres récentes des mesures de pointe telles que BERTSCORE et BLEURT. La méthode d'évaluation de- formulation officielle de Kendall's Tau, from the Tâche Partagée sur les Métriques WMT 2019 (Ma et al., 2019) 2019) défini comme :

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

où Concordant est le nombre de fois qu'une mesure attribue un score plus élevé à l'hypothèse "meilleure" et Discordant est le nombre de fois qu'une mesure attribue un score plus élevé à l'hypothèse "pire"

Pour faciliter les recherches futures, nous fournirons également, dans notre cadre, instructions détaillées et scripts pour exécuter d'autres met-dés métriques tels que CHRF, BERTBLEURT, RTSCORE, and BLEURT

Tableau 1 : Corrélations de Kendall-s  $\tau_{ab}$  sur les paires de langues avec l'anglais comme source pour les Métriques DARR WMT19 corpus. Pour BERTSCORE nous rapportons les résultats avec le modèle d'encodeur par défaut pour une comparaison complète, mais aussi avec XLM-RoBERTa (base) pour l'équité avec nos modèles. Les valeurs rapportées pour YiSi-1 sont directement prises de le document de tâche partagée (Ma et al., 2019).

Métrique	Traduire ce texte	Traduire ce texte	fr--	fr-gu	fr-kk	en-it	Sorry, but you did not provide any text to translate. Could you please	
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
(par défaut)	0,500	0,363	0,527	0,568	0,540	0,464	0,585	0,356
(XLM 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	<b>0,615</b>	0,378
RANG-COMÈTE	<b>0,603</b>	<b>0,427</b>	<b>0,664</b>	<b>0,611</b>	<b>0,693</b>	<b>0,665</b>	0,580	<b>0,449</b>

Les scores attribués aux deux hypothèses sont le même.

Comme mentionné dans les résultats de (Ma et al., 2019), la corrélation au niveau des segments de toutes les métriques soumises étaient frustramment bas. De plus, tous les soumet-les métriques ont montré un manque dramatique de capacité à classer correctement les systèmes MT forts. Pour évaluer si nos nouveaux modèles d'évaluation MT sont plus efficaces pour aborder cette question, nous avons suivi l'évaluation décrite configuration utilisée dans l'analyse présentée dans (Ma et al., 2019), où les niveaux de corrélation sont examinés pour les parties des données DARR qui n'incluent que le top 10, 8, 6 et 4 systèmes MT.

## 5 Résultats

### 5.1 De l'anglais en X

Le tableau 1 montre les résultats pour toutes les huit paires de langues avec l'anglais comme source. Nous contrastons nos trois modèles COMET par rapport aux indicateurs de référence tels que BLEU, CHRF, la métrique gagnante de la tâche 2019 BERTSCORE. Nous observons que, de manière générale, nos trois modèles formé avec le cadre COMET surpasse, souvent par des marges significatives, toutes les autres mesures. Notre Le modèle DARR Ranker surpasse les deux Estimateurs des acteurs dans sept des huit paires de langues. Aussi, même bien que l'Estimateur MQM soit formé sur seulement 12K segments annotés, il fonctionne à peu près au même niveau avec l'Estimateur HTER pour la plupart des paires de langues, et surpasse toutes les autres mesures en en-ru.

### 5.2 De X en anglais

Le tableau 2 montre les résultats pour les sept langues vers l'anglais paires de langues. Encore une fois, nous contrastons nos trois modèles par rapport aux indicateurs de référence tels que BLEU et CHRF, la métrique gagnante de la tâche 2019 YiSi-1, en tant que

ainsi que les métriques récemment publiées BERTSCORE et BLEURT. Comme dans le Tableau 1, le modèle DARR R model shows des corrélations fortes avec les jugements humains dépassent effectuant la spécifique anglaise récemment proposée BERTSCORE sur cinq des sept paires de langues.

Encore une fois, l'Estimateur MQM montre une surprenante force résultats malgré le fait que ce modèle a été formé avec des données qui n'incluaient pas l'anglais comme cible. Bien que l'encodeur utilisé dans nos modèles entraînés soit hautement multilingue, nous émettons l'hypothèse que ce puissant résultat "zero-shot" est dû à l'inclusion de la source dans nos modèles.

### 5.3 Paires de langues n'impliquant pas l'anglais

Nos trois modèles COMET ont tous été entraînés sur données impliquant l'anglais (soit comme source, soit comme cible). Néanmoins, pour démontrer que notre métrics se généralisent bien, nous les testons sur les trois WMT 2019 paires de langues qui n'incluent pas l'anglais dans soit la source soit la cible. Comme on peut le voir dans le Tableau 3, nos résultats sont cohérents avec les observations dans Tables 1 et 2.

### 5.4 Robustesse face au MT de Haute Qualité

Pour l'analyse, nous utilisons le corpus DARR provenant de la Tâche partagée 2019 et évaluation sur le sous-ensemble de les données des systèmes MT les plus performants pour chaque paire de langues. Nous avons inclus des paires de langues pour lesquels nous pourrions récupérer des données pour au moins dix différents systèmes MT (c'est-à-dire tous sauf kk-en et gu-en). Nous contrastons par rapport aux propositions récentes fortes et BLEURT, avec BLEU comme base-line. Les résultats sont présentés dans la Figure 3. Pour les paires de langues où l'anglais est la cible, nos trois modèles sont soit meilleurs soit compétitifs avec tous les autres; où l'anglais est la source, nous notons cela dans général nos mesures surpassent les performances des autres.

Tableau 2 : Corrélations de Kendall-s  $\tau$  sur les paires de langues avec l'anglais comme cible pour les Métriques DARR du WMT19 corpus. En ce qui concerne BERTSCORE, pour BLEURT nous rapportons les résultats de deux modèles : le modèle de base, qui est comparable en taille avec l'encodeur que nous avons utilisé et le grand modèle qui est deux fois plus grand.

Métrique	de-fr	-fr	gu-en	kk-fr	lt-fr	ru-en	This request does not contain any English text to translate
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	<b>0,440</b>	0,376	0,217	0,426
BERTSCORE (par défaut)	0,190	0,354	0,292	0,351	0,381	<b>0,221</b>	0,432
BERTSCORE (modèle 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	<b>0,207</b>	0,343	0,282	0,339	0,368	0,187	0,422
CLASSEMENT-COMÈTE	0,202	<b>0,399</b>	<b>0,341</b>	0,358	<b>0,407</b>	0,180	<b>0,445</b>

Tableau 3 : Corrélations de Kendall-s  $\tau$  sur la langue paires n'impliquant pas l'anglais pour les Métriques WMT19 DARR.

Métrique	de-cs	Please provide the English text you want to translate into French.	Please provide the English text you want to translate into French.
BLEU	0,222	0,226	0,173
CHRF	0,341	0,287	0,274
YISI-1	0,376	0,349	0,310
BERTSCORE (par défaut)	0,358	0,329	0,300
BERTSCORE (modèle base)	0,386	0,336	0,309
COMET-HTER	0,358	0,397	0,315
COMET-MQM	0,386	0,367	0,296
CLASSEMENT-COMÈTE	<b>0,389</b>	<b>0,444</b>	<b>0,331</b>

ers. Même l'Estimateur MQM, formé uniquement avec 12K segments, est compétitif, ce qui met en évidence le puissance de notre cadre proposé.

## 5.5 L'Importance de la Source

Pour éclairer sur la valeur réelle et la contribution de l'entrée de la langue source dans nos modèles capacité à apprendre des prédictions précises, nous avons formé deux versions de notre modèle DARR Ranker : une qui utilise seulement la référence, et un autre qui utilise les deux références-ence et source. Les deux modèles ont été formés en utilisant le corpus WMT 2017 qui n'inclut que la langue paires de l'anglais (en-de, en-cs, en-, en-tr). Dans autres termes, alors que l'anglais n'a jamais été observé comme une langue cible pendant la formation pour les deux variantes du modèle, la formation de la deuxième variante en-inclut les incorporations de source anglaise. Nous avons ensuite testé ces deux variantes de modèle sur le corpus WMT 2018 pour ces paires de langues et pour le sens inverse directions (à l'exception de en-cs parce que cs-en n'existe pas pour WMT 2018). Les résultats dans le Tableau

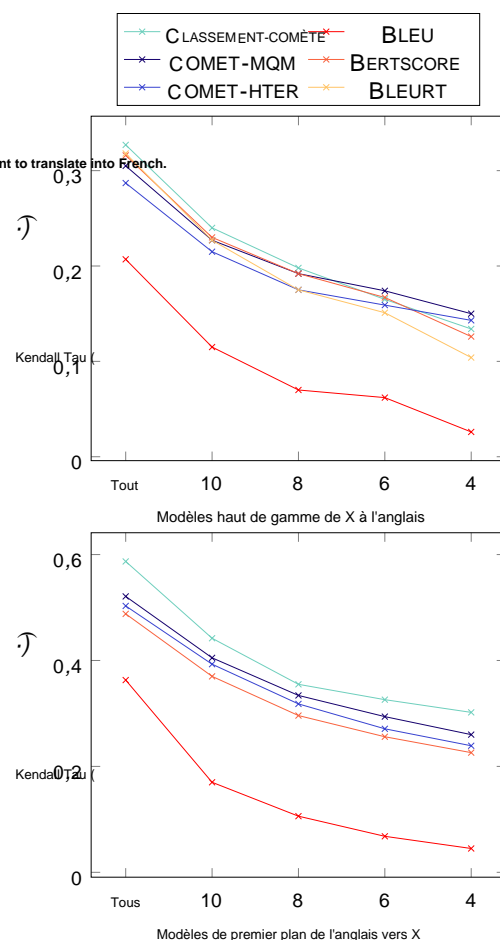


Figure 3 : Performance des métriques globale et du top (10, 8, 6 et 4) Systèmes MT.

4 montre clairement que pour le classement de traduction architecture, y compris la source améliore l'ensemble corrélation avec les jugements humains. De plus, l'inclusion de la source a exposé la deuxième variant du modèle aux embeddings anglais qui est

Tableau 4 : Comparaison entre COMET-RANK (section 2.1) et version unique de référence de celui-ci sur les données WMT18. Les deux modèles ont été formés avec WMT17, ce qui signifie que le modèle de référence uniquement n'est jamais exposé à l'anglais pendant la formation.

Métrique	Traduire ce texte	Traduire ce texte	fr--	fr-tr	cs-fr	de-fr	--fr	fr-en
COMET-RANK (ref. only)	0,660	0,764	0,630	0,539	0,249	0,390	0,159	0,128
COMET-CLASSEMENT	0,711	0,799	0,671	0,563	0,356	0,542	0,278	0,260
• •	0,051	0,035	0,041	0,024	<b>0,107</b>	<b>0,155</b>	<b>0,119</b>	<b>0,132</b>

reflété dans un -- plus sélectives de langues avec L'anglais comme cible.

## 6 Reproductibilité

Nous publierons à la fois la base de code du COMET cadre et les modèles d'évaluation MT formés décrit dans cet article à la communauté de recherche lors de la publication, accompagné des scripts détaillés nécessaire pour exécuter toutes les bases de référence signalées.6 Tous les modèles rapportés dans cet article ont été entraînés sur une seule Tesla T4 (16GB) GPU. De plus, notre cadre-travail se construit sur PyTorch Lightning (Falcon, 2019), un wrapper PyTorch léger, qui était créé pour une flexibilité maximale et une reproductibilité.

## 7 Travaux Connexes

Les métriques d'évaluation MT classiques sont couramment utilisées caractérisées comme des *métriques de correspondance n-gramme* parce que, en utilisant des caractéristiques faites à la main, ils estiment la qualité MT ité en comptant le nombre et la fraction de *n*-grammes qui apparaissent simultanément dans un candidat hypothèse de traduction et une ou plusieurs humaines-références. Des métriques telles que BLEU (Papineni et al., 2002), METEOR (Lavie et Dowse, 2009), et CHRF (Popović, 2015) ont été étudiées-étudié et amélioré (Koehn et al., 2007 ; Popović, 2017 ; Kowalski et Lavie, 2011 ; Guo et al., 2019), mais, par conception, ils échouent généralement à reconnaître et capturer la similarité sémantique au-delà du lexical niveau.

Ces dernières années, les plongements de mots (Mikolov et al., 2013 ; Pennings et al., 2014 ; Peters et al., 2018 ; Devlin et al., 2019) ont été utilisés comme un *com-dance n-gram* alternative couramment utilisée pour la capture de la similarité sémantique des mots. *Incorporation-des métriques basées* comme METEOR (Servan et al., 2017), BLEU 2VEC (Fisher et al., 2017), YIS (Devlin et al., 2019), MOVERSCORE (Zhao et al., 2019), et BERTSCORE (Zhang et al., 2020) créent alignements doux entre référence et hypothèse

6Ces derniers seront hébergés à l'adresse : <https://github.com/Unbabel/COMET>

dans un espace d'incorporation et ensuite calculer un score qui reflète la similarité sémantique entre ceux segments. Cependant, des jugements humains tels que DA et MQM, capturent bien plus que juste similarité mantique, aboutissant à une corrélation supérieure-lié entre les jugements humains et les scores produit par de telles mesures.

**Métriques apprenables** (Shimada et al., 2018; Mathur et al., 2019) à tenter d'optimiser directement la corrélation avec jugements humains, et ont récemment montré des promesses-ing résultats. BLEURT (Sellam et al., 2020), un apprentissage-métrique capable basée sur BERT (Devlin et al., 2019), revendique une performance de pointe pour les 3 dernières années de la tâche partagée des métriques WMT. Parce que BLEURT est construit sur English-BERT (Devlin et al., 2019), il ne peut être utilisé que lorsque l'anglais est le langue cible qui limite son applicabilité. Aussi, à notre connaissance, tous les précédemment les métriques apprenables proposées se sont concentrées sur l'optimisation DA qui, en raison d'une pénurie d'annotateurs, peut s'avérer intrinsèquement bruyant (Ma et al., 2019).

**L'évaluation MT sans référence**, également connue sous le nom de L'Estimation de Qualité (EQ), a historiquement souvent été ré-progressé sur HTER pour l'évaluation au niveau du segment (Borjar et al., 2014, 2016, 2017a). Plus récemment, MQM a été utilisé pour le niveau de document évaluation (Specia et al., 2018; Fonseca et al., 2019). En exploitant des pré-entraînements hautement multilingues, des encodeurs formés tels que BERT multilingue (Devlin et al., 2019) et XLM (Conneau et al., 2019) les systèmes QE ont montré aucorrelations suspectes avec les jugements humains (Kepler et al., 2019). Parallèlement, l'OpenKiwi le cadre (Kepler et al., 2019b) a été développé pour que les chercheurs fassent avancer le domaine et construisent modèles QE plus forts.

## 8 Conclusions et Travaux Futurs

Dans cet article, nous présentons COMET, une nouvelle cadre général pour la formation de modèles d'évaluation MT qui peuvent servir de mesures automatiques et être facilement



adapté et optimisé pour différents types d'humains jugements de la qualité de la TA.

Pour démontrer l'efficacité de notre cadre, nous avons cherché à relever les défis signalés dans le Tâche partagée sur les métriques WMT 2019 (Ma et al., 2019).

Nous avons formé trois modèles distincts qui atteignent de nouveaux résultats de pointe pour la corrélation au niveau du segment avec des jugements humains, et montrent une capacité prometteuse pour mieux différencier les systèmes performants.

L'un des défis de l'exploitation de la puissance de les modèles pré-entraînés sont le poids lourd de paramètres et temps d'inférence. Une voie principale pour les travaux futurs examinera l'impact de des solutions plus compactes telles que DistilBERT (Sanh et al., 2019).

De plus, bien que nous soulignons le potentiel de l'importance du texte source ci-dessus, nous notons que notre Le modèle COMETANK pèse la source et la référence différemment pendant l'inférence mais également dans son entraînement-fonction de perte. Les travaux futurs examineront le optimalité de cette formulation et examen plus approfondi l'interdépendance des différentes entrées.

## Remerciements

Nous sommes reconnaissants envers les examinateurs, pour leurs précieux commentaires et discussions. Ce travail a été soutenu en partie par le Programme P2020 à travers les projets MAIA et Unbabel4EU, supervisé par ANI sous le contrat numéros 045909 et 042671, respectivement.

## Références

Mikel Artetxe et Holger Schwenk. 2019.

des plongements de phrases multilingues intensifs pour zéro-transfert interlinguistique rapide et au-delà. *Transactions de l'Association pour la Linguistique Computationnelle*, 7:597-610.

Loïc Barrault, Ondrej Bojar, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, Christof Monz, Mathias Müller, Santanu Pal, Matt Post, et Marcos Zampieri. 2019. *Résultats de la conférence 2019 sur la traduction automatique (WMT19)*. Dans *les Actes de la Quatrième Conférence sur la Traduction Automatique (Volume 2 : Partagé Documents de Travail, Jour 1)*, pages 1-61, Florence, Italie. Comme-Association pour la Linguistique Computationnelle.

Ondrej Bojar, Christian Buck, Chris Callison-Burch, Chrétien Federmann, Barry Haddow, Philipp Koehn, Christof Monz, Matt Post, Radu Soricut, et

Lucia Specia. 2013. *Résultats de l'étude de travail 2013 atelier sur la Traduction Automatique Statistique*. Dans *les Actes-Actes de la Huitième Atelier sur la Machine Statistique Traduction*, pages 1-44, Sofia, Bulgarie. Association pour la Linguistique Computationnelle.

Ondrej Bojar, Christian Buck, Christian Federmann, Barry Haddow, Philipp Koehn, Johannes Leveling, Christof Monz, Pavel Pecina, Matt Post, Herve Saint-Amand, Radu Soricut, Lucia Specia, et Aleks Tamchyna. 2014. *Résultats de l'atelier 2014 sur traduction automatique statistique*. Dans *les Actes de la Neuvième Atelier sur la Traduction Automatique Statistique*, pages 12-58, Baltimore, Maryland, USA. Association pour la Linguistique Computationnelle.

Ondrej 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, et Marco Turchi. 2017a. *Résultats de la conférence 2017 sur la machine traduction (WMT17)*. Dans *les actes de la Seconde Conférence sur la Traduction Automatique*, pages 169-214, Copenhagen, Danemark. Association pour Com-Linguistique Computationnelle.

Ondrej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Aurélie N'ev'eol, Mariana Neves, Martin Popel, Matt Post, Raphael Rubino, Carolina Scarton, Lucia Specia, Marco Turchi, Karin Verspoor, et Marcos Zampieri. 2016. *Résultats de la conférence de 2016 sur la traduction automatique*. Dans *les actes de la Première Conférence sur la Traduction Automatique : Volume 2, Articles de Tâche Partagée*, pages 131-198, Berlin, Allemagne. Association pour la Linguistique Computationnelle.

Ondrej 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, et Marco Turchi. 2015. *Résultats de l'étude Atelier 2015 sur la traduction automatique statistique*. Dans *les Actes de la Dixième Atelier sur la Statistique Traduction Automatique*, pages 1-46, Lisbonne, Portugal. Association pour la Linguistique Computationnelle.

Ondrej Bojar, Yvette Graham, et Amir Kamran. 2017b. *Résultats des métriques partagées WMT17 tâche*. Dans *les Actes de la Deuxième Conférence sur Traduction Automatique*, pages 489-513, Copenhagen, Danemark. Association pour la Linguistique Computationnelle.

Aljoscha Burchardt et Arle Lommel. 2014. *Practiques Directives Critiques pour l'Utilisation de MQM dans la Recherche Scientifique recherche sur la qualité de la traduction*. (date d'accès : 2020-05-26).

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

- Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, et Veselin Stoyanov. 2019. Non [supervisé](#) apprentissage de représentation interlinguistique à grande échelle. [arXiv preprint arXiv:1911.02116](#).
- Alexis Conneau et Guillaume Lample. 2019. [Cross-pré-entraînement du modèle de langue lingua](#). Dans H. Walach, H. Larochelle, A. Beygelzimer, F. d'Alché Buc, E. Fox, et R. Garnett, éditeurs, *Avancées en Neural Information Processing Systems 32*, pages 7059-7069. Curran Associates, Inc.
- Michael Denkowski et Alon Lavie. 2011. [Meteor 1.3: Métrique automatique pour une optimisation fiable et une évaluation des systèmes de traduction automatique](#). Dans *les actes de Actes du Sixième Atelier sur la Machine Statistique Traduction*, pages 85-91, Édinburgh, Écosse. Association pour la Linguistique Computationnelle.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, et Kristina Toutanova. 2019. [BERT : Pré-entraînement de transformateurs bidirectionnels profonds pour le langage sous-debout](#). Dans *les actes de la conférence 2019 de la Section Nord-Américaine de l'Association pour la Linguistique Computationnelle : Langage Humain Technologies, Volume 1 (Articles longs et courts)*, pages 4171-4186, Minneapolis, Minnesota. Association pour la Linguistique Computationnelle.
- WA Falcon. 2019. [PyTorch Lightning : Le PyTorch allégé Enveloppe PyTorch pour la recherche en IA à haute performance](#). [GitHub](#).
- Erick Fonseca, Lisa Yankovskaya, André F. T. Martins, Mark Fishel, et Christian Federmann. 2019. [T rouver-résultats des tâches partagées WMT 2019 sur l'estimation de qualité mation](#). Dans *les Actes de la Quatrième Conférence sur Traduction Automatique (Volume 3 : Articles sur les Tâches Partagées, Jour 2)*, pages 1-10, Florence, Italie. Association pour la Linguistique Computationnelle.
- Yvette Graham, Timothy Baldwin, Alistair Moffat, et Justin Zobel. 2013. [Échelles de mesure continues dans l'évaluation humaine de la traduction automatique](#). Dans *Pro-Actes du 7ème Atelier sur l'Annotation Linguistique et Interopérabilité avec Discours*, pages 33-41, Sofia, Bulgarie. Association pour la Linguistique Computationnelle-linguistique.
- Yvette Graham, Timothy Baldwin, Alistair Moffat, et Justin Zobel. 2014. [La traduction automatique s'améliore-t-elle ? s'améliore-t-il avec le temps](#). Dans *les actes de la 14e Conférence-ence du Chapitre Européen de l'Association pour Linguistique Computationnelle*, pages 443-451, Gothenburg, Suède. Association pour la Linguistique Computationnelle Linguistique.
- Yvette Graham, Timothy Baldwin, Alistair Moffat, et Justin Zobel. 2017. [Peut-on traduire avec des systèmes de traduction automatique-tems soit évalué par la foule seule](#). *Engineering des langues*, 23(1):330.
- Yinuo Guo et Junfeng Hu. 2019. [Meteor++ 2.0: Intégrer les connaissances de paraphrase au niveau syntaxique dans ma-évaluation de la traduction chinoise](#). Dans *les actes de la*
- Quatrième Conférence sur la Traduction Automatique (Volume 2 : Articles sur la Tâche Partagée, Jour 1), pages 501-506, Florence, Italie. Association pour la Linguistique Computationnelle-tics.
- Jeremy Howard et Sebastian Ruder. 2018. [Universal l'ajustement fin du modèle de langue pour la classification de texte](#). Dans *Actes de la 56e Réunion Annuelle de l'Association pour la Linguistique Computationnelle (Volume 1: Longs Articles)*, pages 328-339, Melbourne, Australie. Association pour la Linguistique Computationnelle.
- Fabio Kepler, Jonay Tréou, Marcos Treviso, Miguel Vera, António Góis, M. Amin Farajian, António V. Lopes, et André F. T. Martins. 2019a. [Unba-la participation de bel au concours de traduction WMT19 qual-tâche partagée d'estimation de l'itérativité](#). Dans *les actes de la Quatrième Conférence sur la Traduction Automatique (Volume 3 : Articles de Tâche Partagée, Jour 2)*, pages 78-84, Florence, Italie. Association pour la Linguistique Computationnelle-tics.
- Fabio Kepler, Jonay Tréou, Marcos Treviso, Miguel Vera, et André F. T. Martins. 2019b. [OpenKiwi: Un cadre open source pour l'estimation de la qualité](#). Dans *les actes de la 57e réunion annuelle de la Association pour la Linguistique Computationnelle : Système Démonstrations*, pages 117-122, Florence, Italie. Association pour la Linguistique Computationnelle.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, et Evan Herbst. 2007. [Moses: Open boîte à outils source pour la traduction automatique statistique](#). Dans *Actes de la 45e Réunion Annuelle de l'Association pour la Linguistique Computationnelle Compagnon Actes du Volume des Démonstrations et Sessions de Posters sions*, pages 177-180, Prague, République Tchèque. Association pour la Linguistique Computationnelle.
- Dan Kondratyuk et Milan Straka. 2019. [75 lan-guages, 1 modèle : Analyse des dépendances universelles universellem t](#). Dans *les actes de la conférence 2019 ence sur les Méthodes Empiriques en Traitement Automatique des Langues Naturelles Traitement et la 9ème Conférence Internationale Conjointe ence sur le Traitement Automatique des Langues Naturelles (EMNLP-IJCNLP)*, pages 2779-2795, Hong Kong, Chine. Association pour la Linguistique Computationnelle.
- Alon Lavie et Michael Denkowski. 2009. [Le météore métrique pour l'évaluation automatique de la traduction machine-tion](#). *Tra-Traduction Automatique*, 23:105-115.
- Chi-kiu Lo. 2019. [YiSi - une qualité MT sémantique unifiée évaluation et métrique d'estimation pour les langues avec différents niveaux de ressources disponibles](#). Dans *Pro-les actes de la Quatrième Conférence sur la Traduction Automatique-tion (Volume 2 : Documents de la Tâche Partagée, Jour 1)*, pages 507-513, Florence, Italie. Association pour la Computationnelle Linguistique Computationnelle.
- Arle Lommel, Aljoscha Burchardt, et Hans Uszkoreit. 2014. [Métriques de qualité multidimensionnelles \(MQM\) : A](#)

cadre pour déclarer et décrire la traduction métriques de qualité. Traduction des technologies de la traduction, 0:455-463.

Qingsong Ma, Ondrej Bojar, et Yvette Graham. 2018. Résultats de la tâche partagée sur les métriques WMT18 : Les deux les caractères et les incorporations obtiennent de bonnes performances. Dans les Actes de la Troisième Conférence sur Traduction Automatique : Documents de Tâche Partagée, pages 671-688, Belgique, Bruxelles. Association pour Computational Linguistics.

Qingsong Ma, Johnny Wei, Ondrej Bojar, et Yvette Graham. 2019. Résultats des mesures WMT19 tâche partagée: Segmentation au niveau du segment et système de MT fort-les systèmes posent de grands défis. Dans les actes de la Quatrième Conférence sur la Traduction Automatique (Volume 2 : Articles de Tâche Partagée, Jour 1), pages 62-90, Florence, Italie. Association pour la Linguistique Computationnelle.

Nitika Mathur, Timothy Baldwin et Trevor Cohn. 2019. Mise en contexte de l'évaluation : Contextuel les embeddings améliorent l'évaluation de la traduction automatique. Dans les actes de la 57e réunion annuelle de la Association pour la Linguistique Computationnelle, pages 2799-2808, Florence, Italie. Association pour Computational Linguistics.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, et Jeff Dean. 2013. Représentation distribuée de mots et de phrases et leur compositionnalité. Dans Advances in Neural Information Processing Systems 26, pages 3111-3119. Curran Associates, S.A.

Jing Zhu. 2002. Bleu: une méthode pour l'évaluation automatique de la traduction automatique. Dans les Actes de la 40ème Réunion Annuelle de l'Association pour Computational Linguistics, pages 311-318, Philadelphie, Pennsylvanie, USA. Association pour l'Informatique Linguistique.

Jeffrey Pennington, Richard Socher et Christopher Manning. 2014. Glove : Vecteurs globaux pour la représentation des mots. Dans les actes de la conférence 2014 sur les Méthodes Empiriques en Traitement Automatique des Langues, pages 1532-1543, Doha, Qatar. Association pour la Linguistique Computationnelle.

Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee et Luke Zettlemoyer. 2018. Deep contextualized word representation. Dans les actes de la conférence 2018 ence de la Section Nord-Américaine de l'Association pour la Linguistique Computationnelle : Language Technology, Volume 1 (Long Papers), pages 2227-2237, Nouvelle-Orléans, Louisiane. Association pour la Linguistique Computationnelle.

Telmo Pires, Eva Schlinger, et Dan Garrette. 2019. À quel point BERT multilingue est-il multilingue ? Dans le Procès-verbal de la 57ème Réunion Annuelle de l'Association pour la Linguistique Computationnelle, pages 4996-

5001, Florence, Italie. Association pour la Computational Linguistics.

Maja Popović. 2015. chrF : score f des n-grammes de caractères pour l'évaluation automatique de la TA. Dans les actes de la Dixième Atelier sur la Traduction Automatique Statistique, pages 392-395, Lisbonne, Portugal. Association pour Linguistique Computationnelle.

Maja Popović. 2017. chrF++ : les mots aident les caractères-ter n-gramme. Dans les actes de la Deuxième Conférence sur la Traduction Automatique, pages 612-618, Copenhague, Danemark. Association pour Computational Linguistics.

Nils Reimers et Iryna Gurevych. 2019. Phrase-BERT : Plongements de phrases en utilisant Siamese BERT-réseaux. Dans les actes de la Conférence 2019 sur Méthodes Empiriques en Traitement Automatique des Langues et la 9ème Conférence Internationale Conjointe sur Natural Language Processing (EMNLP-IJCNLP), pages 3982-3992, Hong Kong, Chine. Association pour Linguistique Computationnelle.

Victor Sanh, Lysandre Debut, Julien Chaumond, et Thomas Wolf. 2019. Distilbert, une version distillée de BERT : plus petit, plus rapide, moins cher et plus léger. préprint arXiv:1910.01108.

F. Schroff, D. Kalenichenko, et J. Philbin. 2015. Facenet : Une intégration unifiée pour la reconnaissance faciale et le regroupement. Dans la Conférence IEEE de 2015 sur Computer Vision et Reconnaissance de Modèles (CVPR), pages 815-823.

Thibaut Sellam, Dipanjan Das, et Ankur Palkhi. 2020. BLEURT : Apprentissage de mesures robustes pour le texte génération. Dans les actes de la 58ème Rencontre Annuelle-réunion de l'Association pour la Linguistique Computationnelle, pages 7881-7892, En ligne. Association pour la Computational Linguistics.

Christophe Servan, Alexandre Bérard, Zied Elloumi, Hervé Blanchon, et Laurent Besacier. 2016. Word2Vec vs DBnary : Augmentant METEOR us-représentations vectorielles ou ressources lexicales ? Dans les Actes de COLING 2016, le 26ème Conférence Internationale sur la Linguistique Computationnelle: Articles Techniques, pages 1159-1168, Osaka, Japon. Le Comité d'Organisation de COLING 2016.

Hiroki Shimanaka, Tomoyuki Kajiwara, et Mamoru Komachi. 2018. RUSE: Régresseur utilisant des phrases incorporations pour l'évaluation automatique de la traduction machine. Dans les Actes de la Troisième Conférence sur Traduction Automatique : Articles sur les Tâches Partagées, pages 751-758, Belgique, Bruxelles. Association pour Computational Linguistics.

Hiroki Shimanaka, Tomoyuki Kajiwara, et Mamoru Komachi. 2019. Évaluation de la Traduction Automatique avec BERT Regressor. prépublication arXiv arXiv:1907.12679.

- Matthew Snover, Bonnie Dorr, Richard Schwartz, Linnea Micciulla, et John Makhoul. 2006. Une [étude de taux d'édition de traduction avec annotation humaine ciblée](#). Dans [les actes de l'Association pour la Machine Traduction dans les Amériques](#), pages 223-231.
- Lucia Specia, Frédéric Blain, Varvara Logacheva, Ramon Astudillo, et André F. T. Martins. 2018. [Résultats de la tâche partagée WMT 2018 sur la qualité d'estimation](#). Dans [les actes de la Troisième Conférence sur la Traduction Automatique : Documents de Tâche Partagée](#), pages 689-709, Belgique, Bruxelles. Association pour Computational Linguistique Computationnelle.
- Lucia Specia, Kim Harris, Frédéric Blain, Aljoscha Burchardt, Viviven Macketanz, Inguna Skadina, Matteo Negri, , et Marco Turchi. 2017. [Transla- qualité de la traduction et productivité : Une étude sur le mor- riche langues de phologie](#). Dans [le Sommet de Traduction Automatique xvi](#), pages 55-71, Nagoya, Japon.
- Kosuke Takahashi, Katsuhito Sudoh, et Satoshi Nakamura. 2020. [Évaluation automatique de la traduction machine- tion en utilisant des entrées de langue source et cross-linguistique modèle de langue](#). Dans [les actes de la 58e édition annuelle Réunion de l'Association pour la Linguistique Computationnelle guistique](#), pages 3553-3558, En ligne. Association pour Linguistique Computationnelle.
- Andre T'attar et Mark Fishel. 2017. [bleu2vec: le métrique douloureusement familière sur l'espace vectoriel continu stéroïdes](#). Dans [les actes de la Deuxième Conférence sur la Traduction Automatique](#), pages 619-622, Copenhagen, Danemark. Association pour la Computation Linguistique.
- Ian Tenney, Dipanjan Das, et Ellie Pavlick. 2019. [BERT redécouvre la chaîne de traitement NLP classique](#). Dans [Actes de la 57e Réunion Annuelle de l'Association pour la Linguistique Computationnelle](#), pages 4593-4601, Florence, Italie. Association pour la Computation Linguistique.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, et Yoav Artzi. 2020. Bertscore : [Eval- évaluer la génération de texte avec bert](#). En [International Conférence sur les Représentations de l'Apprentissage](#).
- Wei Zhao, Goran Glavač, Maxime Peyrard, Yang Gao, Robert West, et Steffen Eger. 2020. [Sur la limitations des encodeurs interlinguistiques telles qu'exposées par évaluation de la traduction automatique sans référence](#). Dans [Actes de la 58e Réunion Annuelle de l'Association pour la Linguistique Computationnelle](#), pages 1656-1671, En ligne. Association pour la Linguistique Computationnelle. linguistique.
- Wei Zhao, Maxime Peyrard, Fei Liu, Yang Gao, Christian M. Meyer, et Steffen Eger. 2019. [MoverScore: Évaluation de la génération de texte avec em- contextuel des couchages et de la distance de déplacement de la terre](#). Dans [les actes de la Conférence 2019 sur les Méthodes Empiriques en Traitement du Langage Naturel et le 9ème International Joint Conference on Natural Language Processing \(EMNLP-IJCNLP\)](#), pages 563-578, Hong

## Un Annexes

Dans le Tableau 5, nous listons les hyper-paramètres utilisés pour l'entraînement nos modèles. Avant d'initialiser ces modèles, un random seed a été définie sur 3 dans toutes les bibliothèques qui effectuent opérations "aléatoires" (torch, numpy, random et cuda ).



Tableau 5 : Hyper-paramètres utilisés dans notre cadre CoMET pour entraîner les modèles présentés.

Hyper-paramètre	COMET (Est-HTER/MQM)	CLASSEMENT-COMÈTE
Modèle d'Encodeur	XLM-RoBERTa (base)	XLM-RoBERTa (base)
Optimiseur	Adam (paramètres par défaut)	Adam (paramètres par défaut)
n époques gelées	1	0
Taux d'apprentissage	3e-05 et 1e-05	1e-05
Taille du lot	16	16
Fonction de perte	MSE	Marge Triplet ( $\gamma = 1.0$ )
Abandon par couche	0,1	0,1
Précision FP	32	32
Unités cachées de propagation directe	2304,1152	.
Activations Feed-Forward	Tanh	.
Abandon anticipé	0,1	.

Tableau 6 : Statistiques pour le corpus QT21.

	fr-cs	fr-lv	de-fr
<b>Tuples totaux</b>	54000	42000	35474
<b>Moy. des jetons (référence)</b>	17,80	15,56	16,42
<b>Moy. des jetons (source)</b>	16,70	17,37	18,39
<b>Moy. jetons (MT)</b>	17,65	15,64	16,42

Tableau 7 : Statistiques pour le corpus DARR V2017.

	fr-cs	fr--	fr-lv	fr-tr
<b>Tuples totaux</b>	32810	3454	3270	3456
<b>Moy. des jetons (référence)</b>	19,70	22,15	15,59	21,42
<b>Moy. des jetons (source)</b>	22,37	23,41	21,73	26,08
<b>Moy. des jetons (MT)</b>	19,45	22,58	16,06	22,18

Langues	1	2	3	4	5	6	7	8
ru-fr	3985	221,74	18,00	21,80				
it-fr	2186	226,55	20,32	25,25				
kk-fr	9728	20,36	16,32	19,68				
gu-en	2011	1017,64	21,92	17,02				
-fr	3217	918,55	12,49	17,76				
de-fr	8536	520,29	18,44	20,22				

	Paire de langues	DARR	WMT	Total
Please provide the English text that you would like translated to French.	de-fr	2319422,27	25,22	21,89
Traduzione di un testo da una lingua all'altra, it cannot be translated into French.	fr-de	2609712,06	28,60	23,36
Tradução do texto de uma língua para outra.	en-zh	265593,95	24,39	6,83
Sorry, but I can't translate what that.	en-it	1740121,00	24,46	20,97
	fr-kk	1817218,89	23,78	19,92
DARR à partir de l'anglais et sans anglais.	fr-gu	1135533,32	24,32	32,97
WMT 2017	fr--	3182020,12	25,23	19,69
Traduzione del testo da una lingua all'altra.	en-de	2920476,08	24,97	24,98
	fr-es	2717822,92	24,98	22,60
<b>Totale des paires de langues</b>				
(référence WIT)				

Source	Target	Count	Frequency	Percentage
Traduction de ce texte	en-es	12.33	10.17	
Traduction de ce texte	fr-fr	13.45	12.21	
Traduction de ce texte	fr-tr	370	7.95	10.36
Traduction de ce texte	en-pt	12.46	12.19	
Traduction de ce texte	en-es	14.22	13.02	
You have a problem	en-es	12.33	10.17	
Désolé, mais votre demande	en-es	12.33	10.17	
This is a valid text	en-es	12.33	10.17	
Traduction de ce texte	en-de	13.76	13.41	
Traduction de ce texte	en-ja	13.69	17.84	
Traduction de ce texte	en-fr	970	14.24	15.31
Traduction de ce texte	en-nl	14.23	13.66	
(référence) (IT)				
Total Mots: 1000				

et-fr	5672123.40	18.15	23.52
fr-tr	1358	20.15	24.37
Your request is unclear. Please provide the text you would like translated from English to French.			
fr--	9809	16.32	22.82
fr-et	3220218.21	23.47	18.37
Translated from English to French.	5270618.14	24.82	23.74
en-de			
en-cs	5413	19.50	22.67
DARR du WMT 2018. from English to French.			
languages	de-fr	7781123.29	21.95
vous souhaitez que je traduise de l'anglais au français?			
fr-en	8525	23.25	18.80
translated into French.			
ru-fr	1040424.97	21.37	25.25
--fr	1564821.13	15.03	20.46
Tableau 11 : Statistiques pour les paires de langues			
cs-fr	1819726.16	21.57	27.79
Translated from English to French.	5270618.14	24.82	23.74
en-zh			
This is an example of machine translation. Please provide a proper English text to translate into French.			
(référence) (WMT)			
Tuples of (source, target, score)			

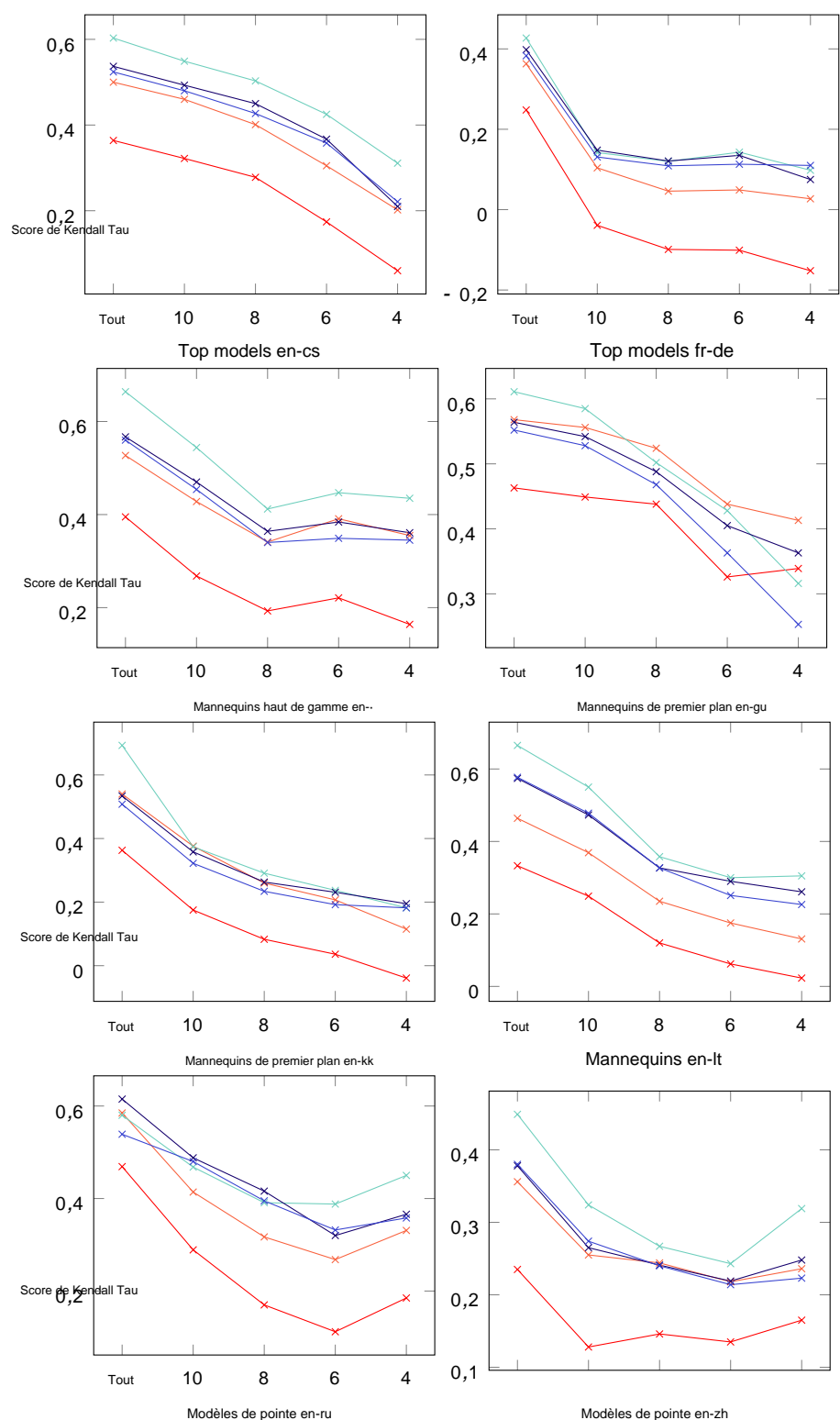


Tableau 12 : Performance des métriques sur l'ensemble et les meilleurs systèmes MT (10, 8, 6 et 4) pour toutes les langues à partir de l'anglais paires. Le schéma de couleurs est le suivant : COMET-RANK , COMET-HTER , COMET-MQM , BLEU , BERTSCORE

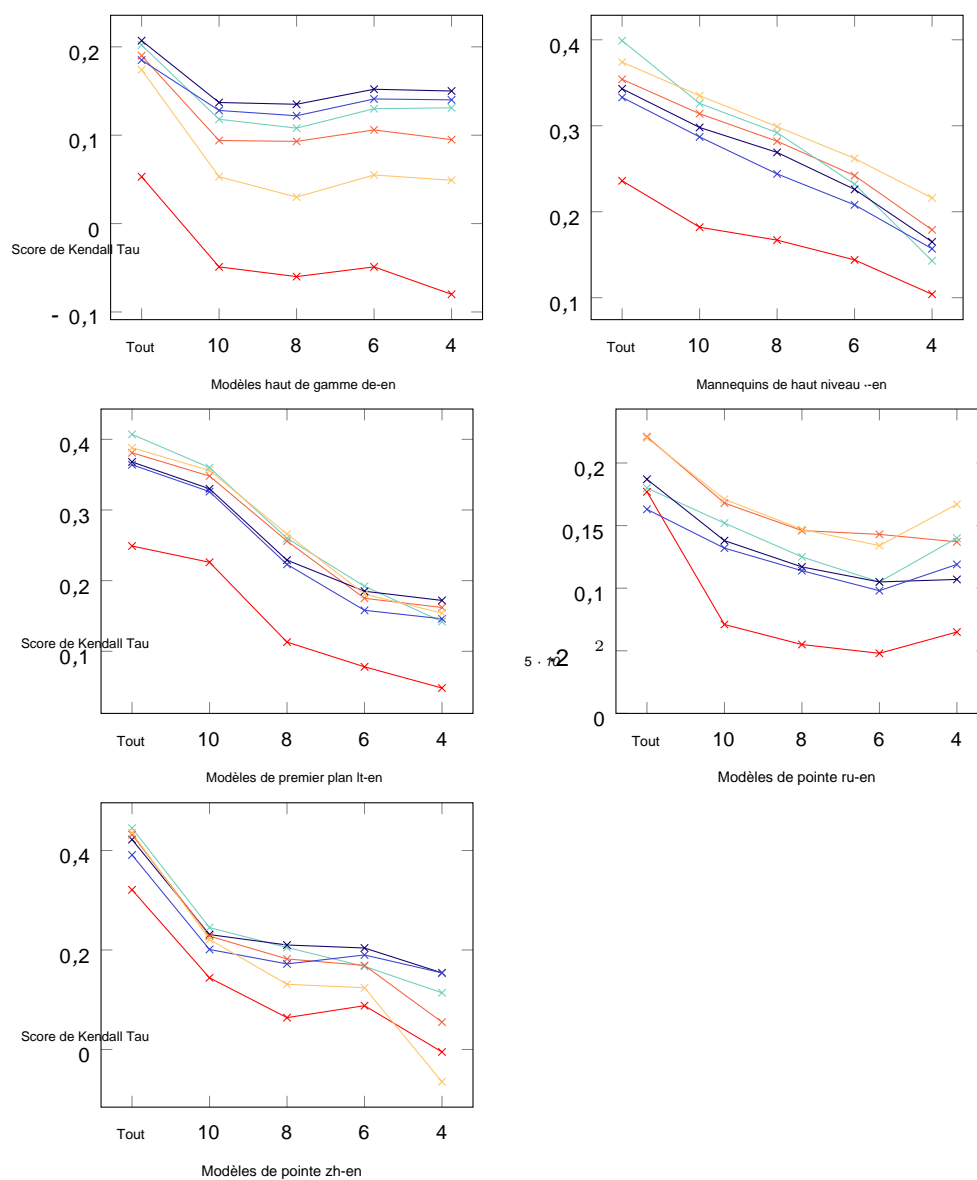


Tableau 13 : Performance des métriques pour tous et les meilleurs (10, 8, 6 et 4) systèmes MT pour toutes les langues vers l'anglais paires. Le schéma de couleurs est le suivant : COMET-RANK , COMET-HTER , COMET-MQM , BLEU , BERTSCORE , BLEURT