# Quality Estimation based feedback training for improving pronoun translation

**Anonymous ACL submission** 

#### Abstract

Pronoun translation is a longstanding challenge in neural machine translation (NMT), often requiring inter-sentential context to ensure linguistic accuracy. To address this, we introduce ProNMT, a novel framework designed to enhance pronoun and overall translation quality in context-aware machine translation systems. ProNMT leverages Quality Estimation (QE) models and a unique Pronoun Generation Likelihood-Based Feedback mechanism to iteratively fine-tune pre-trained NMT models without relying on extensive human annotations. The framework combines OE scores with pronoun-specific rewards to guide training, ensuring improved handling of linguistic nuances. Extensive experiments demonstrate significant gains in pronoun translation accuracy and general translation quality across multiple metrics. ProNMT offers an efficient, scalable, and context-aware approach to improving NMT systems, particularly in translating context-dependent elements like pronouns.

## 1 Introduction

011

013

017

019

021

037

041

Document translation is a critical application of machine translation (MT), facilitating cross-lingual transfer of knowledge. In an increasingly interconnected world, multilingual communication is essential to ensure equitable access to information and services. Despite advances in neural machine translation (NMT) models and the rise of large language models (LLMs), document translation remains a challenging and relevant area of research (Sun et al., 2022; Wang et al., 2023). Traditional MT systems often process sentences independently, which can lead to inconsistencies in terminology and style across a document.

Document translation addresses these limitations by leveraging contextual information across sentences, paragraphs, or entire documents to produce more coherent and accurate translations. For instance, incorporating document-level context improves the handling of anaphora (Voita et al., 2018), lexical disambiguation, and stylistic consistency. 042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

079

081

Various techniques have been employed in literature to fine-tune MT models, with the most prominent ones being: a) Supervised fine-tuning (SFT), which uses labeled data to fine-tune the model in a supervised fashion, and b) Reinforcement Learning from human feedback (RLHF) which is based on optimizing a reward function based on expert judge rankings (human preferences) and using it with Proximal Policy Optimization (PPO) to finetune the model. One of the major drawbacks of both methods is the explicit dependence on human experts to either label huge datasets or rank candidate translations. Furthermore, training PPO involves tuning a large set of hyperparameters and loading multiple models (reference, critic, and reward), which comes at the expense of computational power and expansive memory resources. Although less stable and faster than SFT, RLHF using PPO has shown superior performance in aligning models (Ramamurthy et al., 2023).

Various attempts have been made to integrate feedback to improve the quality of translations in the field of neural machine translation (NMT) as well. Although a few works employ real but limited human feedback (Kreutzer et al., 2018a, Kreutzer et al., 2018b), others focus on using similarity scores between candidates and reference translation as a simulated human feedback. Quality Estimation (QE) models have recently been proven to be an adept proxy for real human feedback-based reward models (He et al., 2024). These QE models, facilitated by the advent of more human evaluation data and better language models (Rei et al., 2020), provide a numerical score to indicate the quality of candidate translation. Our proposed framework is based on exploiting these QE model evaluations to assist the feedback training process iteratively, bypassing the requirement to perform human evaluations since it is very costly in most cases.

158

159

160

161

162

163

164

168

132

133

Despite substantial progress in various areas related to neural machine translation (NMT), the task of translating pronouns has always been inherently difficult for MT models due to dependence on intersentential context for their translation. For example, see the case presented in Fig 1. In both languages, a pronoun in the second sentence refers to advertising. Hence, when the second sentence is translated from English to German, the translation of the pronoun *it* is ambiguous without the previous sentence. This calls for the need for a framework that can help fine-tune MT models to give better overall and pronoun translation scores.

084

096

100

101

102

103

104

105

106

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

Existing methods for document translation often incorporate techniques such as architectural-level modifications to include document-level context, concatenation of context sentences, and multitasking approaches. These methods aim to enhance the translation quality by leveraging the additional contextual information available in documents. However, most models are traditionally trained using reference translations, without explicitly focusing on key areas where document-level machine translation (MT) systems excel—such as improving the translation of pronouns, which are often challenging in standard sentence-level MT systems.

In this work, we introduce **ProNMT**, a novel framework designed to exploit the capabilities of context-aware MT models to improve translation quality, with a particular emphasis on pronoun translation. ProNMT is built upon two core ideas:

- 1. **Quality Estimation (QE)**: This component evaluates the translation quality of pronouns in the output by estimating how well the generated pronouns align with contextual and linguistic expectations.
- 2. **Pronoun Generation Likelihood-Based Feedback**: This is a unique training mechanism where feedback is provided based on the likelihood of generating correct pronouns during translation.

We define "pronoun generation likelihood" as 124 the probability assigned by the model to a pronoun 125 token, given the source sentence and the previously 126 127 generated tokens. This metric serves as a proxy for assessing the quality of pronoun translation, as 128 detailed in Section 3.2. By incorporating this like-129 lihood into our feedback loop, we aim to guide the 130 model toward improved pronoun handling while 131

simultaneously enhancing the overall translation quality.

This paper makes the following key contributions:

- 1. **Pronoun-Focused Feedback Mechanism**: We propose a feedback mechanism specifically tailored to enhance pronoun translation in context-aware MT systems. This mechanism integrates pronoun generation likelihood as a measurable and actionable metric during training.
- 2. **Context-Aware Quality Enhancement**: We exploit document-level context to improve both pronoun translation and the general quality of translations, demonstrating how a targeted approach to pronoun handling can benefit overall translation performance.
- 3. Novel Framework for Fine-Tuning Pre-Trained MT Models: ProNMT provides a practical framework for fine-tuning pretrained MT models, leveraging quality estimation and feedback-driven training processes to address long-standing challenges in document-level MT.
- 4. Evaluation and Results: We empirically validate our approach through extensive experiments, highlighting significant improvements in pronoun translation and overall translation quality compared to existing methods.

By addressing the specific challenges of pronoun translation and leveraging document-level context, ProNMT sets a new benchmark for improving the quality of translations in pre-trained machine translation models.

- $\begin{array}{c} {\bf EN} \quad \mbox{This advertising cheapens all women.} \\ It cheapens every one of us and our daughters. \end{array}$
- **DE** Eine solche Art von <u>Werbung</u> erniedrigt alle Frauen. *Sie* erniedrigt uns alle und unsere Töchter.

Figure 1: Example illustrating the inter-sentential dependence for pronoun translation. Pronouns of interest are *in italics*, and the antecedents they refer to are <u>underlined</u>. Data taken from Europarl EN <-> DE dataset.

## 2 Related Works

Incorporating context is generally better than context-agnostic models (Sim Smith, 2017). The

258

259

261

262

263

264

265

267

268

221

222

223

primary methods to incorporate contexts often 169 use either concatenation (Tiedemann and Scherrer, 170 2017; Junczys-Dowmunt, 2019) or multi-encoder-171 based approaches. Multi-encoder architectures, 172 while helps to achieve better results, similar results could be obtained by passing random con-174 text instead of actual context in the additional en-175 coder (Li et al., 2020). Appicharla et al. (2024) 176 explored multi-task learning (MTL) in contextaware NMT by explicitly modeling context en-178 coding to enhance sensitivity to context choice. 179 Experiments on German-English language pairs 180 showed that the MTL approach outperformed 181 concatenation-based and multi-encoder DocNMT 182 models in low-resource settings. However, they ob-183 served that MTL models struggled to generate the source from the context, suggesting that available document-level parallel corpora may not be sufficiently context-aware. (Wang et al., 2020) used pre-187 vious three sentences during pre-training of Crosslingual Language model Pre-training. Translating pronouns accurately in Neural Machine Translation (NMT) systems remains a significant challenge, primarily due to the necessity of utilizing inter-192 sentential context. Similarly, (Appicharla et al., 193 2023) investigated the impact of different context settings on pronoun translation accuracy. They 195 trained multi-encoder models using previous sentences, random sentences, and a mix of both as con-197 text, evaluating their performance on the ContraPro 198 test set. Their models performed well even with 199 random context, indicating that the models were somewhat agnostic to the specific context provided. Voita et al. (2018) observed that using documentlevel context helps in better pronoun translation. While the previous works observed the effective-204 ness of context in translating pronouns, the effect of pronoun translation as one of the objective is yet to be explored. 207 Our work differs from these approaches by introducing **ProNMT**, a framework that leverages Qual-

ity Estimation (QE) models and a Pronoun Genera-210 tion Likelihood-Based Feedback mechanism to iter-211 atively fine-tune pre-trained NMT models. Unlike 212 previous methods that rely on modeling context 213 through auxiliary tasks or synthetic data, ProNMT 214 bypasses the need for extensive human annotations 215 216 and explicitly focuses on improving pronoun translation within context-aware systems. By integrat-217 ing QE scores with pronoun-specific rewards, our 218 method effectively guides the training process to 219 enhance pronoun handling and overall translation 220

quality, offering a scalable solution to longstanding challenges in document-level translation.

## 3 Methodology

#### 3.1 Framework

Given a pre-trained MT model  $M_0(x, \theta_0)$  with initial parameters  $\theta_0$ , which generates an output ybased on multinomial sampling with underlying distribution  $p_{M_0}(y|x, \theta_0)$ , our aim is to guide the model to generate better translations with a focus on pronouns using a QE-based reward function r(x, y). Note that the QE-based reward function does a reference-free estimation of the translation quality. We define the optimization objective as:

$$\max_{\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim p_M(y|x;\theta)} r(x, y).$$
(1)

In each iteration *i*, we choose a batch of sentence pairs  $\{X_i, Y_i\}$  of size *B*. For each sample input  $x_p \in X_i$ , we then generate *k* candidate translations  $y_1^p, y_2^p, \dots, y_k^p$ . Then, we extract the pronoun of interest from each candidate  $y_j^p \forall j$  and calculate the reward as defined in (2). This helps our framework to choose the best proxy for human feedback (best candidate translation, say  $y_j^p$ ) and update the model parameters using iterative supervised fine-tuning (SFT). This iterative training process is presented in Algorithm 1.

## 3.2 Reward function

The reward function used for training is a linear combination of pronoun-based and translationbased reward metrics. For a given pair of sentences  $(x_p, y_p^t)$  containing the target pronoun token  $y_{pi}^t$ and a generated candidate translation  $y_k^p$ , we assess overall translation quality using the COMET model "wmt21-comet-ge-da" that employs a referencefree evaluation approach and is built on the XLM-R architecture (Rei et al., 2020). This gives us  $R_{translation}$ , a normalized score between -1 and 1, where 1 means perfect translation. Next, to assess the pronoun translation reward  $R_{PGL}$ , we identify the pronoun token  $y_{kj}^p$  in candidate translation  $y_k^p$ and its "Pronoun Generation Likelihood (PGL)" defined as:  $P(y_{kj}^p | \mathbf{x}_{\mathbf{p}}, \mathbf{y}_{k_{1:j-1}}^p; \theta)$ . If the pronoun token matches with that in reference translation, then we set the pronoun reward to PGL itself. If it does not, then it is set to -PGL. Lastly, if no pronoun token in present in the candidate in the first place,  $R_{PGL}$  is set to 0. We are now in a position to define the overall reward  $r(x, y_i)$  in Equation 2.

$$r(x, y_i) = \beta \cdot R_{PGL} + \alpha R_{translation} \quad (2)$$

where

$$R_{PGL} = \begin{cases} PGL & \text{if } y_{kj}^p = y_{pi}^t, \\ -PGL & \text{if } y_{kj}^p ! = y_{pi}^t, \\ 0 \text{ otherwise} \end{cases}$$

## Algorithm 1 ProNMT

**Require:** Training set  $\mathcal{X}$ , reward function r(x, y), initial model  $M_0 = P(y|x; \theta_0)$ , batch size B, temperature T, the number of candidates k1: for iteration i in  $0, 1, \ldots, N-1$  do  $D_i \leftarrow \text{SampleBatch}(\mathcal{X}, b)$ 2:  $\mathcal{B} \leftarrow \emptyset$ 3: 4: for each  $x \in D_i$  do  $y_1,\ldots,y_k \sim P_T(y|x;\theta_i)$ 5:  $y^* \leftarrow \arg \max_{y_j \in \{y_1, \dots, y_k\}} r(x, y_j)$  $\mathcal{B} \leftarrow \mathcal{B} \cup \{(x, y^*)\}$ 6: 7: end for 8: Fine-tune  $\theta_i$  on  $\mathcal{B}$  to obtain  $M_{i+1} =$ 9:  $P(y|x;\theta_{i+1})$ 10: end for

#### **4** Experiments

#### 4.1 Data for training

While a contrastive test-suite to assess pronoun translation quality called Contrapro (Müller et al., 2018) is available, we found it to contain relatively shorter and easy-to-translate source sentences which makes it trivial for fine-tuning a pretrained MT model (see appendix section B). This motivated us to design our own Europarl-based filtered dataset. We start preparing our training data by adopting Europarl our base EN <-> DE sentence corpus. It contains 1920209 sentence pairs of English and German languages. To make the dataset suited for pronoun translation we adopt the following filtering process: for each pair of sentences (s, t) in English and German, extract iff

- s contains the English pronoun it, and t contains a German pronoun that is third person singular (er, sie or es), as indicated by their part-of-speech tags.
- Those pronouns are aligned to each other.

Note that we only consider the pronoun "it" and its German translations "er" (Masculine), "sie" (Feminine) or "es" (Neutral) due to the crucial dependence on source or target side context for its translation quality. This corpus filtering process is also crucial to reduce noise in QE-based feedback training. If all pronouns are considered in the filtered one-pronoun sentences, we observed the training to be noisy due to the possible ambiguity in pronoun translation, i.e. more than one candidate pronoun translations may be valid (see appendix section A.1 for details). Further, checking alignment of pronouns is also important as "it" may correspond to different german pronouns like "ist","das" or "dies" as well (Müller et al., 2018). This filtering process reduces the dataset to 117834 sentences containing "es", 17447 sentences containing "er" and 39439 sentences containing "sie". To tackle class imbalance in the filtered dataset, we sample 15000 sentences from each class to create our final dataset used for training. The final train set contains 42750 sentences, test set contains 1500 sentences and validation set contains 750 instances. For our document-level experiment, we preprocess the source sentence  $x_i$  in the following format: "<context>  $x_{i-1}$  <\context>  $x_i$ ". Note that the same seed was set for shuffling data instances while creating train, test and val sets, for both with and without context experiments. This was done to maintain uniformity in performing assessments.

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

329

330

331

332

333

334

335

336

337

338

339

341

342

343

344

346

347

#### 4.2 Training details

We begin the training process by considering a base pre-trained model. For our experiments, we chose the distilled 600M parameter variant of NLLB-200. In this model, we deploy an iterative Supervised Fine Tuning (SFT) trainer using the trl package from transformers library. This helps us to fine tune the model parameters on the reward-chosen candidates iteratively for each mini-batch. We choose the batch size B = 4 for our experiments. For each input sentence  $x_p$  in a mini-batch, we generate k = 10 candidate translation  $(y_1^p, y_2^p \dots y_{10}^p)$ through multinomial sampling. The "most-suited" candidate  $y_i^c$  is then greedily chosen for each sample in the mini-batch and fed into the Iterative-SFT Trainer to update the model parameters. We run the trainer for a maximum of 700 iterations with a learning rate of 6e-5 and gradient accumulation steps = 16. For assessing the model in a documentlevel setting, we first fine-tuned the model on a subset of our filtered dataset using Huggingface's trainer class. For fine-tuning, we train the model for 3 epochs with the hyperparameters mentioned in Table 1. This is done so that model learns to use the context but not output it in translation during generation. We perform SFT training for without-

270

272

273

274

275

276

279

280

281

284

285

286

291

296

348 context model on same hyperparameters as well,
349 although not required, it helps to normalise test350 ing scores across with and without-context trained
351 models.

| Hyperparameter              | Value              |  |
|-----------------------------|--------------------|--|
| Learning Rate               | $8 \times 10^{-3}$ |  |
| Per Device Train Batch Size | 8                  |  |
| Gradient Accumulation Steps | 16                 |  |
| Eval Accumulation Steps     | 16                 |  |
| Per Device Eval Batch Size  | 8                  |  |
| Weight Decay                | 0.01               |  |
| Number of Train Epochs      | 10                 |  |

Table 1: Hyperparameters for fine-tuning process

#### 4.3 Evaluation

354

361

367

371

372

374

379

383

We chose to assess the model trained through ProNMT during training and testing. For the testing evaluation, we assess the model checkpoint with best validation rewards.

#### 4.3.1 Training evaluation:

During training, we test model's progress in each iteration by calculating  $R_{translation}$  and  $R_{PGL}$  and averaging it over the mini-batch. These plots are presented in Fig. 3 for with-context model training and Fig. 2 for the without-context case. Moreover, in every 100 iterations, we perform a validation analysis during training. We calculate the average training reward and the cross-entropy loss in the validation set. We use these scores to keep track of model's training and to choose the model checkpoint with the highest training reward on validation set.

#### 4.3.2 Testing evaluation:

We evaluate the best model checkpoint found during training based on several open-source benchmarks as listed below:

• **BLEU** (Bilingual Evaluation Understudy)

• **COMET** (Crosslingual Optimized Metric for Evaluation of Translation): We employ two versions of COMET. For direct assessment, we use the *Unbabel/wmt22-comet-da* model, which is fine-tuned on human evaluation data from the WMT22 Metrics Shared Task and name this *COMET* in the results. For quality estimation without reference translations, we use the 'wmt21-comet-qe-da' model and



Figure 2: Batchwise training translation reward (a) and pronoun reward (b) for without-context model.



Figure 3: Batchwise training translation reward (a) and pronoun reward (b) for with-context model.

| Method                | En⇒De |         |        |        |  |  |  |
|-----------------------|-------|---------|--------|--------|--|--|--|
|                       | COMET | BLEU    | QE     | PGL    |  |  |  |
| NLLB 600M - DISTILLED |       |         |        |        |  |  |  |
| NO CONTEXT            |       |         |        |        |  |  |  |
| $\alpha \beta$        |       |         |        |        |  |  |  |
| 1.1/#tokens           | 72.88 | 15.2200 | 0.1189 | 0.3769 |  |  |  |
| $1 \ 1/\#avg - len$   | 73.95 | 15.2630 | 0.1186 | 0.3672 |  |  |  |
| $1.2\;1/\#tokens$     | 63.30 | 10.8723 | 0.0069 | 0.3425 |  |  |  |
| $1.2 \ 1/\#avg-len$   | 66.03 | 11.5783 | 0.0598 | 0.3370 |  |  |  |
| WITH CONTEXT          |       |         |        |        |  |  |  |
| $\alpha \beta$        |       |         |        |        |  |  |  |
| 1.1/#tokens           | 78.66 | 21.2250 | 0.0127 | 0.8270 |  |  |  |
| 11/#avg - len         | 78.98 | 21.9542 | 0.0121 | 0.8543 |  |  |  |
| $1.2\;1/\#tokens$     | 81.92 | 26.9574 | 0.0362 | 0.4183 |  |  |  |
| $1.2\;1/\#avg-len$    | 66.03 | 11.5783 | 0.0598 | 0.337  |  |  |  |

Table 2: Translation evaluation on test set for En $\Rightarrow$ De direction under various combinations of  $\alpha$  and  $\beta$ , using WMT21-COMET-QE-DA as reward model. #tokens refers to number of tokens in the respective sentence, #avg-len refers to the calculated average token length across the source side dataset, calculated to be approximately 30. QE ( $R_{translation}$ ) and PGL ( $R_{PGL}$ ) refer to the respective average reward calculated on the test set.

| Method                         | En⇒De  |       |         |         |        |  |  |
|--------------------------------|--------|-------|---------|---------|--------|--|--|
|                                | Loss   | COMET | BLEU    | QE      | PGL    |  |  |
| NLLB 600M - DISTILLED          |        |       |         |         |        |  |  |
| NO CONTEXT                     |        |       |         |         |        |  |  |
| SFT                            | 1.0507 | 56.68 | 7.8904  | -0.1154 | 0.0902 |  |  |
| ONLY PGL REWARD                | 2.2632 | 50.30 | 1.8290  | -0.1534 | 0.8998 |  |  |
| ONLY QE REWARD                 | 2.2521 | 58.30 | 2.6723  | 0.048   | 0.3798 |  |  |
| PRONMT ( $\alpha *, \beta *$ ) | 2.1593 | 73.95 | 15.2630 | 0.1186  | 0.3672 |  |  |
| WITH CONTEXT                   |        |       |         |         |        |  |  |
| BASELINE                       | 6.807  | 66.03 | 11.5783 | 0.0598  | 0.337  |  |  |
| SFT                            | 1.2814 | 81.19 | 25.568  | 0.0258  | 0.1406 |  |  |
| ONLY PGL REWARD                | 2.2340 | 16.79 | 0.0703  | -0.6024 | 0.9675 |  |  |
| ONLY QE REWARD                 | 3.606  | 80.06 | 21.5588 | 0.0303  | 0.2262 |  |  |
| PRONMT ( $\alpha *, \beta *$ ) | 1.1486 | 81.92 | 26.9574 | 0.0362  | 0.4183 |  |  |

Table 3: Translation performance comparison of best combination with reward baselines, using WMT21-COMET-QE-DA as reward model for EN $\Rightarrow$ DE direction. QE ( $R_{translation}$ ) and PGL ( $R_{PGL}$ ) refer to the respective average reward calculated on the test set.  $\alpha *$  and  $\beta *$  refer to the hyperparameter configurations with highest COMET and BLEU scores (highlighted in bold in Table 2).

present it under the *QE* column in the results. This QE model predicts the quality of a translation based solely on the source sentence and the translation hypothesis.

By combining these metrics, we obtain a comprehensive evaluation of the model's performance in terms of accuracy, fluency, and alignment with human judgments.

# 5 Results

384

386

392

397

We present the translation assessment of NLLB 600M - distilled when trained on different chosen configurations in Table 3. The evaluation results reveal that incorporating document-level context significantly enhances both pronoun-specific and general translation quality, as demonstrated by the superior performance of context-aware models across all metrics. For instance, in the EN $\rightarrow$ DE direction, the context-aware ProNMT achieved a COMET score of 81.92 and a BLEU score of 26.95, compared to 73.95 and 15.2630, respectively, for the context-agnostic model. Pronoun-specific rewards, particularly those leveraging Pronoun Generation Likelihood (PGL), led to notable improvements in pronoun handling, with the context-aware model achieving a PGL score of 0.4183 versus 0.3672 for the baseline. However, models trained solely with PGL rewards underperformed on overall translation metrics, highlighting the importance of balancing PGL with Quality Estimation (QE)

398

399

400

401

402

403

404

405

406

407

408

409

410

411

rewards. The combined use of QE and PGL re-413 wards, optimized with appropriate weight config-414 urations, yielded the best results, as evidenced by 415 consistent improvements in batch-wise rewards dur-416 ing training. Context-agnostic models struggled to 417 resolve inter-sentential dependencies, further un-418 derscoring the necessity of leveraging document-419 level context for coherent and accurate translations. 420 We observed two key trends in the evaluation re-421 sults: first, the inclusion of document-level con-499 text significantly enhances both pronoun-specific 423 and overall translation quality, as evidenced by the 494 context-aware ProNMT achieving higher scores 425 across metrics such as COMET (81.92 vs. 73.95) 426 and BLEU (26.95 vs. 15.2630) compared to its 427 context-agnostic counterpart. 428

> Second, ProNMT's ability to jointly optimize QE and pronoun-specific rewards led to consistent improvements in external metrics like COMET and BLEU, particularly when both reward components  $(\alpha \neq 0, \beta \neq 0)$  were employed. This balance of rewards proved crucial for achieving concurrent gains in general translation quality and accurate handling of linguistically challenging pronoun translations.

## 6 Conclusion

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

In this paper, we introduced ProNMT, a novel framework designed to address the longstanding challenges of pronoun translation in Neural Machine Translation (NMT) systems. By leveraging Quality Estimation (QE) models and the Pronoun Generation Likelihood-Based Feedback mechanism, ProNMT effectively improves both pronounspecific and overall translation quality without the need for extensive human annotations. Our method uniquely integrates QE-based evaluations with pronoun-specific rewards, guiding iterative fine-tuning processes that are scalable, efficient, and context-aware.

Extensive experimental evaluations demon-451 strated that ProNMT consistently outperforms base-452 line systems across multiple metrics, including 453 COMET, BLEU and QE models. Importantly, in-454 corporating document-level context significantly 455 enhanced the handling of linguistically complex el-456 ements, such as pronouns, while maintaining high 457 458 performance on general translation tasks. These results validate the framework's ability to address 459 both inter-sentential dependencies and broader doc-460 ument coherence in machine translation. 461

## References

Ramakrishna Appicharla, Baban Gain, Santanu Pal, and Asif Ekbal. 2023. A case study on context encoding in multi-encoder based document-level neural machine translation. In *Proceedings of Machine Translation Summit XIX, Vol. 1: Research Track*, pages 160–172, Macau SAR, China. Asia-Pacific Association for Machine Translation. 462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

- Ramakrishna Appicharla, Baban Gain, Santanu Pal, Asif Ekbal, and Pushpak Bhattacharyya. 2024. A case study on context-aware neural machine translation with multi-task learning. In *Proceedings of the* 25th Annual Conference of the European Association for Machine Translation (Volume 1), pages 246–257, Sheffield, UK. European Association for Machine Translation (EAMT).
- Zhiwei He, Xing Wang, Wenxiang Jiao, Zhuosheng Zhang, Rui Wang, Shuming Shi, and Zhaopeng Tu. 2024. Improving machine translation with human feedback: An exploration of quality estimation as a reward model. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 8164–8180, Mexico City, Mexico. Association for Computational Linguistics.
- Marcin Junczys-Dowmunt. 2019. Microsoft translator at WMT 2019: Towards large-scale document-level neural machine translation. In *Proceedings of the Fourth Conference on Machine Translation (Volume* 2: Shared Task Papers, Day 1), pages 225–233, Florence, Italy. Association for Computational Linguistics.
- Julia Kreutzer, Shahram Khadivi, Evgeny Matusov, and Stefan Riezler. 2018a. Can neural machine translation be improved with user feedback? In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 3 (Industry Papers), pages 92–105, New Orleans - Louisiana. Association for Computational Linguistics.
- Julia Kreutzer, Joshua Uyheng, and Stefan Riezler. 2018b. Reliability and learnability of human bandit feedback for sequence-to-sequence reinforcement learning. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1777–1788, Melbourne, Australia. Association for Computational Linguistics.
- Bei Li, Hui Liu, Ziyang Wang, Yufan Jiang, Tong Xiao, Jingbo Zhu, Tongran Liu, and Changliang Li. 2020. Does multi-encoder help? a case study on contextaware neural machine translation. In *Proceedings* of the 58th Annual Meeting of the Association for Computational Linguistics, pages 3512–3518, Online. Association for Computational Linguistics.
- Mathias Müller, Annette Rios, Elena Voita, and Rico Sennrich. 2018. A Large-Scale Test Set for the Evaluation of Context-Aware Pronoun Translation in Neu-

525

- 529 530 531 532 533
- 534 535

536 537 538

- 539 540
- 540 541 542
- 543 544
- 5 5 5

555 556

557 558

560 561

- 56
- 564
- 565
- 566 567

568

569 570

571

574

ral Machine Translation. In *WMT 2018*, Brussels, Belgium. Association for Computational Linguistics.

- Rajkumar Ramamurthy, Prithviraj Ammanabrolu, Kianté Brantley, Jack Hessel, Rafet Sifa, Christian Bauckhage, Hannaneh Hajishirzi, and Yejin Choi. 2023. Is reinforcement learning (not) for natural language processing: Benchmarks, baselines, and building blocks for natural language policy optimization. *Preprint*, arXiv:2210.01241.
- Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. Unbabel's participation in the WMT20 metrics shared task. In *Proceedings of the Fifth Conference on Machine Translation*, pages 911–920, Online. Association for Computational Linguistics.
- Karin Sim Smith. 2017. On integrating discourse in machine translation. In *Proceedings of the Third Workshop on Discourse in Machine Translation*, pages 110–121, Copenhagen, Denmark. Association for Computational Linguistics.
- Zewei Sun, Mingxuan Wang, Hao Zhou, Chengqi Zhao, Shujian Huang, Jiajun Chen, and Lei Li. 2022. Rethinking document-level neural machine translation. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 3537–3548, Dublin, Ireland. Association for Computational Linguistics.
- Jörg Tiedemann and Yves Scherrer. 2017. Neural machine translation with extended context. In *Proceedings of the Third Workshop on Discourse in Machine Translation*, pages 82–92, Copenhagen, Denmark. Association for Computational Linguistics.
- Elena Voita, Pavel Serdyukov, Rico Sennrich, and Ivan Titov. 2018. Context-aware neural machine translation learns anaphora resolution. In *Proceedings* of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1264–1274, Melbourne, Australia. Association for Computational Linguistics.
- Longyue Wang, Chenyang Lyu, Tianbo Ji, Zhirui Zhang, Dian Yu, Shuming Shi, and Zhaopeng Tu. 2023.
   Document-level machine translation with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16646–16661, Singapore. Association for Computational Linguistics.
- Longyue Wang, Zhaopeng Tu, Xing Wang, Li Ding, Liang Ding, and Shuming Shi. 2020. Tencent AI lab machine translation systems for WMT20 chat translation task. In *Proceedings of the Fifth Conference on Machine Translation*, pages 483–491, Online. Association for Computational Linguistics.

# A Appendix

# A.1 Training on all German pronouns

In our initial experiments, we tried to incorporate all German pronouns into our training framework. We considered the following consolidated list of pronouns: ['mein', 'uns', 'euer', 'ihnen', 'der', 'ihm', 'die', 'euch', 'diesen', 'unser', 'dem', 'denen', 'dieses', 'meinem', 'den', 'diese', 'du', 'seiner', 'meines', 'das', 'ich', 'deiner', 'dich', 'dir', 'meiner', 'meinen', 'es', 'meine', 'wir', 'sein', 'ihn', 'deren', 'diesem', 'sie', 'dessen', 'dieser', 'mich', 'ihr', 'mir', 'derer/deren', 'dein', 'ihrer', 'er']. Upon running ProNMT on this pronoun list, we get the results presented in Fig. 5. We attribute the noisy nature of the training curves to the noise introduced by the PGL reward and the possible ambiguity associated with translation in this scenario (Müller et al., 2018).



Figure 4: Batchwise training translation reward (a) and pronoun reward (b) for with-context model when no german pronouns are filtered.

# A.2 Source side token distribution in datasets

In this section, we will contrast the Europarl and Contrapro datasets with respect to the source side sentence token count distribution. Performing a basic statistical analysis, we get the following results.

# (i) Contrapro

| Mean: 12.533     | 59 |
|------------------|----|
| Median: 11.0     | 59 |
| Variance: 53.031 | 59 |

588

589

590

591

592

593

594

575

576

577

578

579

580

581

583

584

585

586

610

611

612

614

615

616

Standard Deviation: 7.282 Max length: 67 Min length: 2

# (ii) Europarl

Mean: 33.378 Median: 29.0 Variance: 404.219 Standard Deviation: 20.105 Max length: 212 Min length: 1



Figure 5: Token counts distribution for (a) Contrapro and (b) filtered Europarl dataset.

We hypothesize that the lower mean length and left-shifted distribution for the contrapro data set made the Quality estimation model scores saturated from the beginning of training, giving the model less room to learn. This was the reason we chose our own filtered Europarl dataset for training, testing, and validation.

## **B** Limitations

618 We identify the following limitations in our work:

The experiments were conducted only for
 NLLB 600M distilled variant. To assess the

robustness of our framework across MT models, we can expand the scope of our chosen models.

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

- The current framework only accommodates the translation of the pronoun "it". We can extend the framework to include nonoverlapping pronouns, i.e., pronouns whose translation is not ambiguous.
- We could not perform hyperparameter tuning for the SFT model training. The hyperparameters presented in Table 1 are default hyperparameters.
- We only consider the EN→DE direction in our experiments as it is considered to be a more difficult task than the opposite direction in MT.