Efficient k-Nearest-Neighbor Machine Translation with Dynamic Retrieval

Anonymous ACL submission

Abstract

To achieve non-parametric NMT domain adaptation, k-Nearest-Neighbor Machine Translation (kNN-MT) constructs an external datastore to store domain-specific translation knowledge, which derives a kNN distribution to interpolate the prediction distribution of the NMT model via a linear interpolation coefficient λ . Despite its success, kNN retrieval at each timestep leads to substantial time overhead. To address this issue, dominant studies resort to kNN-MT with adaptive retrieval (kNN-MT-AR), which dy-011 namically estimates λ and skips kNN retrieval if λ is less than a fixed threshold. Unfortunately, 014 kNN-MT-AR does not yield satisfactory results. In this paper, we first conduct a preliminary study to reveal two key limitations of kNN-MT-AR: 1) the optimization gap leads to inaccurate estimation of λ for determining kNN retrieval 019 skipping, and 2) using a fixed threshold fails to accommodate the dynamic demands for kNNretrieval at different timesteps. To mitigate these limitations, we then propose kNN-MT with dynamic retrieval (kNN-MT-DR) that significantly extends vanilla kNN-MT in two aspects. Firstly, we equip kNN-MT with a MLPbased classifier for determining whether to skip kNN retrieval at each timestep. Particularly, we explore several carefully-designed scalar features to fully exert the potential of the classifier. Secondly, we propose a timestep-aware threshold adjustment method to dynamically generate the threshold, which further improves the efficiency of our model. Experimental results on the widely-used datasets demonstrate the effectiveness and generality of our model.¹

1 Introduction

041

As an effective paradigm for non-parametric domain adaptation, *k*-Nearest-Neighbor Machine Translation (*k*NN-MT) (Khandelwal et al., 2020) derives from *k*-Nearest-Neighbor Language Model (*k*NN-LM) (Khandelwal et al., 2019) and has garnered much attention recently. Typically, kNN-MT introduces translation knowledge stored in an external datastore to enhance the NMT model, which can conveniently achieve non-parametric domain adaptation by changing external datastores.

In *k*NN-MT, a datastore containing key-value pairs is first constructed with an off-the-shelf NMT model, where the key is the decoder representation and the value corresponds to its target token. During translation, the current decoder representation is used as a query to retrieve *k* nearest pairs from the datastore, where retrieved values are converted into a probability distribution. Finally, via a linear interpolation coefficient λ , this distribution is used to adjust the prediction distribution of the NMT model. In spite of success, retrieving at each timestep incurs substantial time overhead, which becomes considerable as the datastore expands.

То address this drawback, researchers have proposed two categories of approaches: 1) datastore compression that improves retrieval efficiency by reducing the size of datastores (Martins et al., 2022a; Meng et al., 2022; Wang et al., 2022; Dai et al., 2023; Zhu et al., 2022; Deguchi et al., 2023); 2) retrieval reduction that skips some kNN retrieval to speed up decoding. In this regard, the most representative work is kNN-MT with adaptive retrieval (kNN-MT-AR) (Martins et al., 2022a) that skips kNN retrieval when the coefficient λ is less than a fixed threshold α . However, kNN-MT-AR does not achieve desired results as reported in (Martins et al., 2022a).

In this work, we mainly focus on the studies of *retrieval reduction*, which is compatible with the other type of studies. To this end, we first reimplement *k*NN-MT-AR (Martins et al., 2022a) and conduct a preliminary study to analyze its limitations. Through in-depth analyses, we show that 1) the optimization gap leads to inaccurate estimation of λ for determining *k*NN retrieval skipping; 2) with the increase in timesteps, the demand for *k*NN 042

043

044

047

048

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

078

079

081

082

¹We will release our code upon the acceptance of our paper.

> 104 105

103

106 107

108 109

110

112

113 114

115 116

117

118 119

1

121

123 124

122

125

126

127

128

129

130

131

retrieval diminishes, which proves challenging for the fixed threshold α to handle effectively.

To overcome the above defects, we then significantly extend the vanilla kNN-MT into kNN-MT with dynamic retrieval (kNN-MT-DR), which accelerates the model decoding in two aspects. Concretely, instead of relying on the interpolation coefficient λ , we introduce a MLP-based classifier to explicitly determine whether to skip kNN retrieval as a binary classification task. Particularly, instead of using the decoder representation as the input of the classifier, we explore several carefullydesigned scalar features to fully exert the potential of the classifier. Besides, we propose a timestepaware threshold adjustment method to dynamically generate the threshold, so as to further improve the efficiency of our model.

To summarize, main contributions of our work include the following four aspects:

 Through in-depth analyses, we conclude two defects of kNN-MT-AR: the optimization gap leads to inaccurate estimation of λ for kNN retrieval skipping, and a fixed threshold is unable to effectively handle the varying demands of kNN retrieval at different timesteps.

• We propose to equip kNN-MT with an explicit classifier to determine whether to skip kNN retrieval, where carefully-designed features enable our model to achieve a better balance between model acceleration and performance.

• We propose a timestep-aware threshold adjustment method to further improve the efficiency of our model.

• Empirical evaluations on the multi-domain datasets validate the effectiveness of our model, as well as its compatibility with datastore compression methods.

2 Related Work

Datastore Compression. In this aspect, the size of the datastore for kNN retrieval is decreased to make retrieval efficient. For example, Martins et al. (2022a) compress the datastore by greedily merging neighboring pairs that share the same values, and applying PCA algorithm (Wold et al., 1987) to reduce the dimension of stored keys. Meanwhile, Zhu et al. (2022) prune the datastore based on the concept of local correctness, while Wang et al. (2022) presents a cluster-based compact network to condense the dimension of stored keys, coupled with a cluster-based pruning strategy to discard redundant pairs. Additionally, some studies opt for dynamically adopting more compact datastores. For instance, for each token in the input sentence, Meng et al. (2022) identify the relevant parallel sentences that contain this token and then collect corresponding word-aligned target tokens to construct a smaller datastore. Subsequently, Dai et al. (2023) conduct sentence-level retrieval and dynamically construct a compact datastore for each input sentence. With the same motivation, Deguchi et al. (2023) suggest retrieving target tokens from a subset of neighbor sentences related to the input sentence, where a look-up table based distance computation method is used to expedite retrieval. 132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

Retrieval Reduction. In this regard, some kNN retrieval is reduced to decrease time overhead for retrieval. For instance, Martins et al. (2022b) adopt chunk-wise kNN retrieval rather than timestepwise one, and Martins et al. (2022a) explore two approaches to reduce the frequency of kNN retrieval operations: 1) one introduces a caching mechanism to speed up decoding, where the cache mainly contains retrieved pairs from previous timesteps, and skip kNN retrieval if the distance between the query and any cached key is less than a predefined threshold; 2) the other proposes to conduct kNN retrieval when the interpolation coefficient λ is less than a predefined threshold α , which, however, does not achieve satisfactory results.

Our work mainly focuses on the second type of studies mentioned above. We first conduct a preliminary study to in-depth analyze two limitations of the λ -based kNN retrieval skipping. To address these limitations, we introduce a classifier to explicitly determine whether to skip kNN retrieval as a classification task. Notably, almost concurrently with our work, Shi et al. (2023) also use a classifier to speed up model decoding, sharing a similar motivation with ours. However, our work not only achieves better results, but also significantly differs from theirs in the following three aspects:

First, we explore several carefully-designed scalar features as the input for the classifier, which are crucial for achieving better performance. Second, when training the classifier, we adopt more reasonable criteria to construct training samples. To be specific, in addition to skipping retrieval when the target token ranks the 1st position in the NMT prediction distribution, we believe that the model should also skip when the target token can not be obtained through kNN retrieval. Finally,
based on the observation that the demand for kNN
retrieval diminishes as timesteps increase, we propose a timestep-aware threshold method to further
improve the efficiency of our model.

3 Preliminary Study

3.1 Background

188

190

192

194

195

196

198

204

205

206

207

210

211

212

213

214

215

216

Typically, given an off-the-shelf NMT model f_{θ} , a vanilla *k*NN-MT model is constructed through the following two stages:

Datastore Construction. At this stage, all parallel sentence pairs in the training corpus $C = \{(x, y)\}$ are first fed into the NMT model f_{θ} in a teacherforcing manner (Williams and Zipser, 1989). At each timestep *t*, the decoder representation h_t and its corresponding target token y_t are collected to form a key-value pair, which is then added to the key-value datastore $\mathcal{D} = \{(h_t, y_t) \mid \forall y_t \in \mathbf{y}, (\mathbf{x}, \mathbf{y})\},$ where $h_t = f_{\theta}(\mathbf{x}, \mathbf{y}_{< t})$.

Translating with Retrieved Pairs. During inference, the datastore is used to assist the NMT model. Specifically, the decoder representation \hat{h}_t is used as a query to retrieve k pairs $\mathcal{N}_t = \{(h_i, y_i)\}_{i=1}^k$ from \mathcal{D} , which are then converted into a probability distribution over the vocabulary, abbreviated as kNN distribution:

$$p_{k\mathrm{NN}}(\hat{y}_t | \boldsymbol{x}, \boldsymbol{y}_{< t}) \propto \sum_{(h_i, y_i) \in \mathcal{N}_t} \mathbb{1}_{(\hat{y}_t = y_i)} \exp(\frac{-d(h_i, \hat{h}_t)}{\tau}), \qquad (1)$$

where $\mathbb{1}_{(*)}$ is an indicator function, $d(h_i, \hat{h}_t)$ measures the Euclidean distance between the query \hat{h}_t and the key h_i , and τ is a predefined temperature. Finally, *k*NN-MT interpolates p_{kNN} with the prediction distribution p_{NMT} of the NMT model as a final translation distribution:

$$p(\hat{y}_t | \boldsymbol{x}, \boldsymbol{y}_{< t}) = \lambda p_{kNN} + (1 - \lambda) p_{NMT}, \quad (2)$$

217 where λ denotes a predefined interpolation coeffi-218 cient tuned on the validation set.

219kNN-MT with Adaptive RetrievalObviously,220the retrieval of kNN-MT at each timestep incurs221significant time overhead. To address this limita-222tion, Martins et al. (2022a) follow (He et al., 2021)223to explore kNN-MT with adaptive retrieval (kNN-224MT-AR). Unlike the vanilla kNN-MT, they dynam-225ically estimate the interpolation coefficient λ using

α	IT	Koran	Law	Medical	Subtitles
0.25	0.27	0.14	0.04	0.12	0.50
0.50	0.50	0.54	0.26	0.43	0.60
0.75	0.51	0.59	0.40	0.42	0.60

Table 1: F1 scores of the λ -based kNN retrieval skipping of kNN-MT-AR (Martins et al., 2022a) on the test sets.

a light MLP network, which takes several neural and count-based features as the input. Then, they not only interpolate the kNN and NMT prediction distributions with λ , but also skip kNN retrieval when λ is less than a fixed threshold α . During training, they minimize the cross-entropy (CE) loss over the interpolated translation distribution. Unfortunately, extensive results on several commonlyused datasets indicate that kNN-MT-AR does not achieve satisfactory results.

3.2 Limitations of *k*NN-MT-AR.

In this subsection, we conduct a preliminary study to explore the limitations of kNN-MT-AR. We strictly follow the settings of (Martins et al., 2022a) to re-implement their kNN-MT-AR, and then conduct two groups of experiments on the commonlyused multi-domain datasets released by Aharoni and Goldberg (2020).

As reported by Martins et al. (2022a), dynamically determining whether to skip kNN retrieval based on λ leads to significant performance degradation. In the first group of experiments, to further provide evidence of this conclusion, we perform decoding on the test sets in a teacher-forcing manner and analyze the F1 scores of λ -based kNN retrieval skipping. As shown in Table 1, F1 scores remain relatively low no matter which thresholds and datasets are used.

For the above results, we believe that there are two reasons leading to the inaccurate estimation of λ , which in turn makes λ unsuitable for deciding whether to skip kNN retrieval. In addition to lacking the information of kNN distribution for λ estimation², we believe that **the optimization objective of minimizing the CE loss over the translation distribution may be unsuitable to train an accurate** λ **estimator for determining** kNN **retrieval skipping**. To verify this claim, we consider whether to skip kNN retrieval as a standard binary classification task and use a binary CE loss to train

²Due to the consideration of model efficiency, kNN-MT-AR do not exploit the kNN retrieval information to estimate λ , which has been shown to be effective in previous studies (Zheng et al., 2021; Jiang et al., 2022).



Figure 1: The changes of BLEU improvements between adjacent intervals. [0,5] means that kNN-MT only conducts retrieval when timestep ranges from 0 to 5.

a classifier for λ estimation. Note that this classifier is also based on MLP and contains the same input as kNN-MT-AR. To avoid description confusion, we denote the λ trained by kNN-MT-AR and the above binary CE loss as Tran- λ and Bina- λ , respectively. Then, we calculate the average absolute value of the difference between Bina- λ and Tran- λ at all timesteps. The statistical results show that the average difference is 0.1495, and 29.12% of timesteps exhibit a difference exceeding 0.2. These findings indicate significant differences between Bina- λ and Tran- λ .

266

270

271

272

275

277

278

279

290

291

299

301

In the second group of experiments, we conduct experiments with vanilla kNN-MT to explore the impact of kNN retrieval during different timestep intervals. Specifically, we limit the model to only perform kNN retrieval in specific timestep intervals, where each interval starts from 0 and increases by 5 timesteps in length. From Figure 1, we observe that with the increase in timesteps, the performance gain caused by kNN retrieval gradually decreases across all datasets. This observation reveals that **the demand for** kNN retrieval varies **at different timesteps, which can not be handled** well by the fixed threshold α in kNN-MT-AR.

In summary, the above two defects seriously limit the practicality of kNN-MT-AR. Therefore, it is of great significance to explore more effective skipping kNN retrieval methods for kNN-MT.

4 Our Model

In this section, we significantly extend kNN-MT into kNN-MT-DR in the following two aspects.

4.1 Classifier for Determining *k*NN Retrieval Skipping

Unlike kNN-MT-AR leveraging the interpolation coefficient λ for determining whether to skip kNN

retrieval, we directly equip kNN-MT with a binary classifier to determine whether to skip at each timestep. This classifier comprises a two-layer MLP network with ReLU activation. At timestep t, we conduct kNN retrieval only if the prediction probability of the classifier on conducting kNN retrieval exceeds a timestep-aware threshold α_t , otherwise we will directly skip kNN retrieval. In the following, we will discuss the classifier, which involves the construction of training samples, input features, and the training objective. 302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

321

322

323

324

325

327

328

329

331

333

334

335

336

338

339

340

341

342

343

346

347

348

351

Construction of Training Samples To train the classifier, one crucial step is to construct training samples. In this regard, within the exploration of kNN-LM, He et al. (2021) propose to construct training examples with two distinct labels, namely, "conducting retrieval" and "skipping retrieval", by comparing the prediction probabilities of kNN and NMT distributions on the target token y_t : when $p_{kNN}(y_t)$ is greater than $p_{NMT}(y_t)$, then kNN retrieval should be conducted, otherwise it can be skipped. However, such a criterion still leads to a lot of redundant kNN retrieval. For example, when the target token y_t has the highest probability in the NMT prediction distribution, there is no need to perform kNN retrieval, even if $p_{kNN}(y_t) \ge p_{NMT}(y_t)$. Taking the IT validation set as an example, 69.8% of timesteps satisfy $p_{\text{kNN}}(y_t) \ge p_{\text{NMT}}(y_t)$, among which 77.9% of the timesteps have y_t ranking the 1st position in the NMT prediction distribution. Based on the above analysis, we traverse the parallel sentence pairs in the validation set, and collect various information at each timestep to construct training samples according to the following criteria:

- kNN retrieval should be skipped if one of the two conditions is satisfied: 1) yt ranks the 1st position in the NMT prediction distribution, and 2) yt does not appear in the pairs obtained via kNN retrieval. Obviously, kNN retrieval yields no benefit in both conditions.
- kNN retrieval should be conducted if yt is not the top-1 token in the NMT prediction distribution and it occurs in the kNN retrieval pairs. In this situation, conducting kNN distribution has the potential to improve translation.

Input Features. Unlike *k*NN-MT-AR, which uses the decoder representation and vectors mapped by other scalar features as the input, we consider several carefully-designed scalar features

437

438

439

440

441

399

400

401

402

403

404

405

406

407

408

409

410

411

412

as the input for the classifier directly. By doing so, we reduce the input dimension, achieving effective training and enabling efficient inference. Here, we give detailed descriptions to these features:

354

362

375

376

378

379

396

397

- $p_{\text{NMT}}(\hat{y}_t)$: the probability of the top-1 predicted token \hat{y}_t in the NMT prediction distribution. The higher the prediction confidence of the NMT model, the more likely the \hat{y}_t to be the correct one. In this situation, kNN retrieval is more likely to be skipped.
 - $\|\hat{h}_t\|_2$: the L_2 norm of current decoder representation. Inspired by (Liu et al., 2020), we use the vector norm of the decoder representation \hat{h}_t to measure the translation difficulty at current timestep: the larger $\|\hat{h}_t\|_2$, the more difficult the translation is.
 - max(*attn*): the maximal weight of the crossattention in the last layer of the decoder during current decoding timestep. A large weight means that the NMT model is relatively certain about which source token to be translated. In this case, the translation difficulty is often relatively low.

Finally, these features are concatenated and normalized with batch normalization (Ioffe and Szegedy, 2015) before being the input for the classifier.

Classifier Training. To achieve efficient domain adaptation for NMT, we fix the parameters of NMT model and only update those of classifier during training. Following He et al. (2021), we select 90% of the validation set to train the classifier, and use the remaining 10% for validation. Then, according to the above criterion, we construct training samples with different labels at each timestep to train our classifier. Considering the significant imbalance between two classes of training samples³, we adopt Focal Loss (Lin et al., 2017) to train our classifier as follows:

$$\mathcal{L}(p_c) = -\alpha_c (1 - p_c)^{\gamma} log(p_c), \qquad (3)$$

where c=0/1 denotes the label of skipping/conducting kNN retrieval, p_c is the prediction probability of the classifier on the label c, α_c is a weighting factor controlling the balance between different kinds of samples, and γ is a hyper-parameter adjusting the impacts of loss functions of easy and hard samples (Lin et al., 2017).

4.2 Timestep-aware Threshold Adjustment

As analyzed in Section 3.2, the benefit of kNN retrieval diminishes with the increase in timesteps, indicating that using the fixed threshold α is not the most reasonable choice. To deal with this issue, we propose a timestep-aware threshold adjustment method to accommodate the varied demands of kNN retrieval. Formally, we heuristically define a dynamic threshold function specific to the timestep:

$$\alpha_t = \alpha_{\min} + \operatorname{clip}(\frac{t}{T}; 0, 1)^2 \times (0.5 - \alpha_{\min}) \quad (4)$$

where $\operatorname{clip}(x; a, b)$ clamp x within the range of [a, b], t is the decoding timestep, α_{\min} is the lower limit of threshold, and T is the average length of sentences in the validation set. Apparently, with the increase of t, α_t will gradually increase until it reaches 0.5.

5 Experiments

5.1 Setup

Datasets. We conduct experiments on the multidomains dataset released by Aharoni and Goldberg (2020). The dataset comprises German-English parallel corpora across five domains: Koran, IT, Medical, Law, and Subtitles, the detailed statistics can be found in Appendix A. We employ Byte Pair Encoding (Sennrich et al., 2016) to split words into subwords. Finally, we use two metrics to evaluate the translation quality: SacreBLEU⁴ (Post, 2018) and COMET⁵ (Rei et al., 2020).

Model Configuration. We develop our model with kNN-BOX⁶ (Zhu et al., 2023) and use Faiss (Johnson et al., 2019) to build the datastore and search nearest neighbors. To ensure fair comparisons, we adopt the same settings as the previous study (Khandelwal et al., 2020). Concretely, we set the number of retrieved pairs to 8, the temperature τ to 100 for Koran and 10 for the other datasets, and λ to 0.7 for IT, Subtitles, 0.8 for the other datasets. We use a two-layer MLP network with ReLU activation (Agarap, 2018) to construct our classifier, of which hidden size is set to 32 because it is not sensitive in our model. Besides, we set the hyper-parameter α_{min} to 0.45 for Koran, Subtitles, 0.4 for the other datasets.⁷

³Through data analysis, we find that only 16.8% of training samples require kNN retrieval in the IT validation set.

⁴https://github.com/mjpost/sacrebleu

⁵https://github.com/unbabel/COMET

⁶https://github.com/NJUNLP/knn-box

⁷The details of tuning α_{\min} are reported in Appendix C.

Model	IT	Koran	Law	Medical	Subtitles	Average
Base NMT	38.35 / 82.74	16.26 / 72.04	45.48 / 85.66	39.99 / 83.13	29.27 / 79.76	33.87 / 80.67
Vanilla <i>k</i> NN-MT	45.83 / 85.19	20.37 / 72.30	61.16 / 87.46	54.22 / 84.73	31.28 / 80.13	42.57 / 81.96
k NN-MT-AR(α =0.25)	43.20 / 84.57	19.57 / 72.27	59.89 / 87.57	53.12 / 84.97	30.46 / 80.04	41.25 / 81.88
k NN-MT-AR(α =0.50)	41.19 / 84.05	17.23 / 72.25	58.83 / 87.50	51.22 / 84.69	29.45 / 79.87	39.58 / 81.67
k NN-MT-AR(α =0.75)	39.05 / 83.30	16.40 / 72.09	51.11 / 86.65	45.14 / 84.08	29.30 / 79.82	36.20 / 81.19
Faster kNN-MT	44.25 / 84.59	18.82 / 72.07	58.97 / 87.36	51.02 / 84.45	30.76 / 80.04	40.76 / 81.70
Ours	45.48 / 84.60	20.34 / 72.40	60.10 / 87.39	51.97 / 84.36	31.24 / 80.14	41.83 / 81.78

Table 2: BLEU / COMET scores of various models on the multi-domain test sets.

442 **Baselines.** Our baselines include:

443

444

445 446

447

448

449

450

451

452

453

454

455

456 457

458

459

460

461

462

463

464

- **Base NMT** (Ng et al., 2019). Following Khandelwal et al. (2020), we use the WMT'19 German-English news translation task winner as the base NMT model.
- Vanilla *k*NN-MT (Khandelwal et al., 2020). It serves as a baseline, upon which we develop our model.
- kNN-MT-AR (Martins et al., 2022a). It performs retrieval only when the interpolation coefficient λ is less than a predefined threshold α. Note that it is our most important baseline. Particularly, we report the performance of kNN-MT-AR with α set to 0.25, 0.50, and 0.75, respectively.
- Faster *k*NN-MT (Shi et al., 2023). It is a concurrent work with ours, where a two-layer MLP network takes decoder representation as the input to determine whether to skip *k*NN retrieval at each timestep.

5.2 Main Results

To comprehensively evaluate various models, we report their translation quality and decoding speed.

Translation Quality. Table 2 presents BLEU and 465 COMET scores of various models on the multi-466 domain test sets. We observe that both kNN-MT-467 AR and Faster kNN-MT suffer from significant 468 performance declines compared to Vanilla kNN-469 MT, echoing with the results reported in previous 470 studies (Martins et al., 2022a; Shi et al., 2023). In 471 472 contrast, our model exhibits the least performance degradation. Specifically, our model achieves aver-473 age BLEU and COMET scores of 41.83 and 81.78 474 points, with only 0.74 and 0.18 points lower than 475 those of Vanilla kNN-MT, respectively. 476

477**Decoding Speed.** Model efficiency is a crucial478performance indicator for kNN-MT. As imple-

Model	IT	Koran	Law	Medical	Subtitles
]	Batch Size =	128		
Base NMT	3270.84	3912.95	3690.85	3152.59	4004.40
Vanilla kNN-MT	2584.31	3287.24	2300.23	2363.00	478.99
kNN-MT-AR	2724.76	3069.38	2241.93	2382.52	886.16
Faster kNN-MT	2912.67	3609.53	2923.79	2676.11	999.57
Ours	2944.38	3522.49	2933.76	2605.12	1002.13
		Batch Size	= 64		
Base NMT	3150.95	3730.90	3607.41	3111.54	3377.17
Vanilla kNN-MT	2506.85	2945.54	2252.18	2329.36	445.88
kNN-MT-AR	2789.59	2678.89	2125.88	2323.60	794.04
Faster kNN-MT	2783.62	3124.68	2726.92	2592.75	898.26
Ours	2798.39	3132.26	2755.02	2575.40	901.33
		Batch Size	= 32		
Base NMT	2559.84	2933.82	2995.43	2688.93	2635.05
Vanilla kNN-MT	2001.80	2360.50	1908.76	1955.65	408.54
kNN-MT-AR	2067.55	1792.74	1925.76	1694.34	676.28
Faster kNN-MT	2131.48	2432.76	2225.68	2047.19	735.26
Ours	2117.94	2392.60	2226.63	2031.51	737.85
		Batch Size	= 16		
Base NMT	1577.03	1878.36	1959.55	1737.23	1686.02
Vanilla kNN-MT	1378.65	1429.78	1318.55	1366.35	340.96
kNN-MT-AR	1369.49	1437.82	1244.21	1323.07	506.17
Faster kNN-MT	1396.32	1451.95	1455.46	1406.26	538.65
Ours	1441.66	1487.54	1472.04	1395.21	546.22
		Batch Size	= 1		
Base NMT	159.24	168.84	173.22	171.12	159.04
Vanilla kNN-MT	136.19	139.02	142.91	138.31	42.75
kNN-MT-AR	127.23	130.35	127.93	128.09	57.98
Faster k NN-MT	139.54	140.85	147.18	140.68	58.46
Ours	139.84	140.18	147.44	139.84	58.62

Table 3: Decoding speed (#Tok/Sec \uparrow) of various models using different batch sizes on the multi-domain test sets. Here, we only display the decoding speed of *k*NN-MT-AR(α =0.25), since *k*NN-MT-AR(α =0.5) and *k*NN-MT-AR(α =0.75) exhibit significant performance degradation, as reported in Table 2. All results are evaluated on an NVIDIA RTX A6000 GPU.

mented in previous studies (Zheng et al., 2021; Deguchi et al., 2023), we try different batch sizes: 1, 16, 32, 64 and 128, and then report the model efficiency using *"#Tok/Sec"*: the number of translation tokens generated by the model per second.

Experimental results are listed in Table 3. We have the following interesting findings: First, regardless of the batch size used, our model is more efficient than both Vanilla kNN-MT and kNN-MT-AR(α =0.25). Second, as the batch size increases,

487

488

479

569

570

571

572

573

574

575

576

577

528

529

Model	BLEU
Faster kNN-MT	44.25
Ours	45.48
Our Criteria⇒Conventional Criteria	43.90
Dynamic Threshold⇒Fixed Threshold	44.28
Focal Loss⇒Weighted CE Loss	44.79
$w/o p_{\rm NMT}(\hat{y}_t)$	$\bar{44.62}$
w/o $\ \hat{h}_t\ _2$	45.01
w/o $\max(Attn)$	45.12

Table 4: Ablation studies on the IT test set.

489 the efficiency advantage of our model becomes 490 more apparent. On most datasets, we find that the acceleration ratios of our model with large batch 491 sizes (64 or 128) are significantly higher than those 492 with small batch sizes (1 or 16). Finally, with the 493 increase of the datastore size, the efficiency ad-494 vantage of our model also becomes more signifi-495 496 cant. As analysed in Appendix A, the datastore in Subtitles contains the maximum number of pairs 497 while the datastore in Koran is the smallest. Cor-498 respondingly, our model has the most significant 499 acceleration effect on the Subtitles dataset, while 500 the acceleration effect on the Koran dataset is the 501 least significant. 502

Based on the above experimental results, we believe that compared with baselines, ours can achieve better balance between model performance degradation and acceleration.

5.3 Ablation Studies

503

504

507

508

510

511

512

513

514

515

516

517

519

521

523

524

525

527

Following previous studies (Zheng et al., 2021; Jiang et al., 2022), we compare our model with its variants on the IT test set. As shown in Table 4, we consider the following variants:

• Our Criteria
Conventional Criteria. As mentioned in Section 4.1, we adopt new criteria to determine whether kNN retrieval in training samples can be skipped. To verify the effectiveness of our criteria, we compare our criteria with the conventional criteria as mentioned in He et al. (2021): the kNN retrieval should be conducted if $p_{kNN}(y_t) \ge p_{NMT}(y_t)$, otherwise it can be skipped. We first report the proportion changes between two labels of training samples on the IT dataset. Using the conventional criteria, the proportion of training samples labeled as skipping retrieval is about 30.2%, which is significantly smaller than the proportion 83.2% in our criteria. Obviously, more kNN retrieval can be skipped

with our criteria. Second, we focus on the change of model performance. From Line 2, we observe that the conventional criteria leads to a significant performance degeneration, which strongly reveals the effectiveness of our critera.

- Dynamic Threshold⇒Fixed Threshold. We replace the proposed dynamic threshold α_t mentioned in Section 4.2 with the originally-used fixed threshold α=0.5 in this variant. As shown in Line 3, we observe that removing the dynamic threshold leads to a performance decline, demonstrating the effectiveness of our threshold adjustment method.
- Focal Loss⇒Weighted CE Loss. To make a fair comparison, we follow Shi et al. (2023) to adopt a weighted CE loss, which sets γ as 0 in Equation 3. Back to Table 4, we find that this variant is inferior to our model in terms of translation quality. However, it still surpasses Faster kNN-MT with a large margin, confirming the significant advantage of our model in translation quality.
- *w/o Input Features.* To verify the benefit of our carefully-designed features, we thoroughly construct several variants, each of which discards one kind of feature to train the classifier. As shown in Lines 6-8, all variants exhibit performance drops with varying degrees. Thus, we confirm all features are useful for our classifier.

5.4 Experiments on Adaptive *k*NN-MT

Adaptive kNN-MT (Zheng et al., 2021) is a widelyused variant of kNN-MT and significantly outperforms Vanilla kNN-MT in terms of performance. It introduces a meta-k network, a two-layer MLP incorporating distances and counts of all kNN retrieval pairs, to dynamically estimate λ . Our model can also utilize Adaptive kNN-MT as the base models. When using Adaptive kNN-MT as the base model, we dynamically estimate λ solely for timesteps considered to conduct kNN retrieval. Additionally, we explore the performance of Adaptive kNN-MT as the base model for kNN-MT-AR. To ensure fairness, we employ the λ of kNN-MT-AR to determine whether to skip kNN retrieval, and interpolate using the λ of Adaptive kNN-MT.

We also report the translation quality and decoding speed, as shown in the Table 5 and Table 6, respectively. Our model also demonstrate the least

Model	IT	Koran	Law	Medical	Subtitles	Average
Adaptive kNN-MT	47.26 / 85.99	20.15 / 73.22	62.68 / 88.07	56.49 / 85.25	31.49 / 80.25	43.61 / 82.56
+ k NN-MT-AR(α =0.25)	44.34 / 84.92	20.19 / 72.40	61.86 / 87.66	55.46 / 84.76	30.64 / 79.92	42.50 / 81.93
+ k NN-MT-AR(α =0.50)	41.34 / 84.51	17.04 / 72.05	59.71 / 87.37	52.33 / 84.59	29.37 / 79.83	39.96 / 81.67
+ k NN-MT-AR(α =0.75)	39.22 / 83.69	16.48 / 72.06	51.28 / 86.60	45.23 / 84.08	29.30 / 79.81	36.30 / 81.25
+ Faster kNN-MT	45.38 / 85.43	19.04 / 72.98	59.95 / 87.73	53.09 / 84.91	30.63 / 80.06	41.62 / 82.22
+ Ours	46.94 / 85.46	20.05 / 73.26	61.17 / 87.75	54.58 / 84.98	31.35 / 80.38	42.82 / 82.37

Table 5: BLEU / COMET scores of various models based on Adaptive kNN-MT.

Model	IT	Koran	Law	Medical	Subtitles
Adaptive kNN-MT	2583.92	3320.01	2292.75	2368.51	484.62
+ k NN-MT-AR(α =0.25)	2646.95	3098.34	2191.50	2235.01	873.98
+ Faster kNN-MT	2923.62	3665.24	2901.53	2733.55	952.27
+ Ours	2971.77	3569.44	2883.89	2712.45	1075.36

Table 6: Decoding speed (#Tok/Sec \uparrow) of various models based on Adaptive *k*NN-MT. Note that we also omit the results of *k*NN-MT-AR(α =0.5) and *k*NN-MT-AR(α =0.75). Here, we set the batch size as 128.

Model	IT	Koran	Law	Medical	Subtitles	Average
PLAC	46.81 / 85.65	20.51 / 73.21	62.89 / 88.01	56.05 / 85.16	31.59 / 80.36	43.57 / 82.48
+ Ours	46.83 / 85.40	20.36 / 73.25	61.66 / 87.82	54.82 / 85.01	31.28 / 80.29	42.99 / 82.35
PCK	47.27 / 86.43	19.93 / 72.96	62.91 / 88.03	56.46 / 85.15	31.69 / 80.53	43.65 / 82.62
+ Ours	46.85 / 85.97	19.99 / 73.24	61.98 / 88.05	55.34 / 85.11	31.20 / 80.44	43.07 / 82.56

Table 7: BLEU / COMET scores of PLAC (Zhu et al., 2022) and PCK (Wang et al., 2022), alongside these integrated with ours.

Model	IT	Koran	Law	Medical	Subtitles
PLAC	2684.36	3398.53	2433.44	2383.00	749.49
+Ours	3027.95	3596.20	3025.14	2713.74	1461.30
РСК	2873.40	3535.19	2673.76	2617.73	979.52
+Ours	3072.21	3588.76	3009.64	2720.04	1801.97

Table 8: Decoding speed (#Tok/Sec \uparrow) of PLAC (Zhu et al., 2022) and PCK (Wang et al., 2022), alongside these integrated with ours. Here, we set the batch size as 128.

performance decline and achieve the most efficient decoding speed. Although Faster *k*NN-MT demonstrates comparable decoding speeds to ours, our model achieves superior performance.

5.5 Compatibility with Datastore Compression Methods

578

579

581

583

584

586

588

589

591

592

In this group of experiments, we choose PLAC (Zhu et al., 2022) and PCK (Wang et al., 2022) as the basic models for our compatibility experiment, both of which are derived from Adaptive *k*NN-MT. Typically, PLAC prunes the datastore by eliminating pairs with high knowledge margin values, while PCK introduces a cluster-based compact network to condense the dimension of stored keys and utilizes a cluster-based pruning strategy to discard redundant pairs.

Tables 7 and 8 report the translation quality and decoding speed, respectively. We can observe that our model can further improve the efficiency of these two models, with slight drops in translation quality. Thus, we confirm that ours is also compatible with both PLAC and PCK.

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

6 Conclusion and Future Work

In this work, we first in-depth analyze the limitations of kNN-MT-AR, and then significantly extend the vanilla kNN-MT to kNN-MT-DR in two aspects. First, we equip the model with a classifier to determine whether to skip kNN retrieval, where several carefully-designed scalar features are exploited to exert the potential of the classifier. Second, we propose a timestep-aware threshold adjustment method to further refine kNN retrieval skipping. Extensive experiments and analyses verify the effectiveness of our model.

Inspired by (Li et al., 2023), we will further improve our model by incorporating more source-side information into our classifier. Besides, we aim to generalize our model to kNN-LM (Khandelwal et al., 2019) and multilingual scenario (Stap and Monz, 2023), so as to validate its generalizability.

630

632

633

637

643

647

651

654

655

659

619 As our model integrates an additional classifier, there is an associated increase in time consumption. Notably, as the size of the datastore decreases, the 621 time overhead for kNN retrieval diminishes and classifier-related time cost becomes more apparent, 623 624 which results in a less pronounced acceleration in decoding. Besides, the experiments of decoding 625 speed are evaluated solely on a single computer, while the time overhead of kNN retrieval may differ across different hardware, yielding varied accel-629 eration results.

References

Limitations

- Abien Fred Agarap. 2018. Deep learning using rectified linear units (relu). *arXiv*.
- Roee Aharoni and Yoav Goldberg. 2020. Unsupervised domain clusters in pretrained language models. In *Proc. of ACL.*
- Yuhan Dai, Zhirui Zhang, Qiuzhi Liu, Qu Cui, Weihua Li, Yichao Du, and Tong Xu. 2023. Simple and scalable nearest neighbor machine translation. In *Proc. of ICLR*.
- Hiroyuki Deguchi, Taro Watanabe, Yusuke Matsui, Masao Utiyama, Hideki Tanaka, and Eiichiro Sumita.
 2023. Subset retrieval nearest neighbor machine translation. In *Proc. of ACL*.
- Junxian He, Graham Neubig, and Taylor Berg-Kirkpatrick. 2021. Efficient nearest neighbor language models. In *Proc. of EMNLP*.
- Sergey Ioffe and Christian Szegedy. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proc. of ICML*.
- Hui Jiang, Ziyao Lu, Fandong Meng, Chulun Zhou, Jie Zhou, Degen Huang, and Jinsong Su. 2022. Towards robust k-nearest-neighbor machine translation. In *Proc. of EMNLP*.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. *IEEE Trans. on Big Data*.
- Urvashi Khandelwal, Angela Fan, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. Nearest neighbor machine translation. In *Proc. of ICLR*.
- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2019. Generalization through memorization: Nearest neighbor language models. In *Proc. of ICLR*.
- Xuanhong Li, Peng Li, and Po Hu. 2023. Revisiting source context in nearest neighbor machine translation. In *Proc. of EMNLP*.

Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. 2017. Focal loss for dense object detection. In *Proc. of ICCV*. 667

668

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

- Xuebo Liu, Houtim Lai, Derek F. Wong, and Lidia S. Chao. 2020. Norm-based curriculum learning for neural machine translation. In *Proceedings of the* 58th Annual Meeting of the Association for Computational Linguistics, Online. Association for Computational Linguistics.
- Pedro Martins, Zita Marinho, and André FT Martins. 2022a. Efficient machine translation domain adaptation. In *Proc. of the 1st Workshop on Semiparametric Methods in NLP: Decoupling Logic from Knowledge.*
- Pedro Henrique Martins, Zita Marinho, and André FT Martins. 2022b. Chunk-based nearest neighbor machine translation. In *Proc. of EMNLP*.
- Yuxian Meng, Xiaoya Li, Xiayu Zheng, Fei Wu, Xiaofei Sun, Tianwei Zhang, and Jiwei Li. 2022. Fast nearest neighbor machine translation. In *Proc. of ACL Findings*.
- Nathan Ng, Kyra Yee, Alexei Baevski, Myle Ott, Michael Auli, and Sergey Edunov. 2019. Facebook FAIR's WMT19 news translation task submission. In *Proc. of MT*.
- Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proc. of MT*.
- Ricardo Rei, Craig Stewart, Ana C. Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proc. of EMNLP*.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proc. of ACL*.
- Xiangyu Shi, Yunlong Liang, Jinan Xu, and Yufeng Chen. 2023. Towards faster k-nearest-neighbor machine translation. *arXiv*.
- David Stap and Christof Monz. 2023. Multilingual *k*-nearest-neighbor machine translation. In *Proc. of EMNLP*.
- Dexin Wang, Kai Fan, Boxing Chen, and Deyi Xiong. 2022. Efficient cluster-based k-nearest-neighbor machine translation. *arXiv preprint arXiv:2204.06175*.
- Ronald J Williams and David Zipser. 1989. A learning algorithm for continually running fully recurrent neural networks. *Neural Computation*.
- Svante Wold, Kim Esbensen, and Paul Geladi. 1987. Principal component analysis. *Chemometrics and intelligent laboratory systems*.
- Xin Zheng, Zhirui Zhang, Junliang Guo, Shujian Huang, Boxing Chen, Weihua Luo, and Jiajun Chen. 2021. Adaptive nearest neighbor machine translation. In *Proc. of ACL*.

- 718
 719
 720
 721
 722
 723
 724
 725
- 72
- 728
- 729

731

733

736

738

740

741 742

743

747

748

749

751

755

- Wenhao Zhu, Shujian Huang, Yunzhe Lv, Xin Zheng, and Jiajun Chen. 2022. What knowledge is needed? towards explainable memory for knn-mt domain adaptation. *arXiv*.
- Wenhao Zhu, Qianfeng Zhao, Yunzhe Lv, Shujian Huang, Siheng Zhao, Sizhe Liu, and Jiajun Chen. 2023. knn-box: A unified framework for nearest neighbor generation. arXiv.

A Dataset Statistics

The number of parallel sentence pairs in different datasets and the sizes of the constructed datastores are shown in Table 9.

Dataset	IT	Koran	Law	Medical	Subtitles
Train	223K	18 K	467K	248K	14.46M
Valid	2K	2K	2K	2K	2K
Test	2K	2K	2K	2K	2K
Size	3.6M	0.5M	19.1M	6.9M	180.7M

Table 9: The statistics of datasets in different domains. We also list the size of the datastore, which is the number of stored pairs.

B Effect of Datastore Size

As analyzed in Section 5.2, our speed advantage becomes more significant with the increase of datastore size. To further verify this, we construct datastores of varying sizes by randomly deleting pairs from the original datastore, and employ the pruned datastores for kNN retrieval. The results of decoding speed on the Subtitles dataset are reported in Table 2. As expected, we observe that our model consistently surpasses kNN-MT, regardless of the datastore size. Furthermore, the efficiency advantage of our model over kNN-MT becomes more evident with the increase of datastore size. These results further confirm that the pronounced speed advantage of our model as the datastore expands.

C Hyper-Parameter Tuning

The performance and efficiency of our model is significantly impacted by the hyper-parameter α_{min} , and we tune α_{min} among the subset of $\{0.45, 0.40, 0.35\}$ on the validation set.

We report the BLEU score and #Tok/Sec, as shown in Table 10. As α_{min} decreases, the increment in BLEU score gradually diminishes, while the drop in decoding speed becomes more pronounced. So we set the hyper-parameter α_{min} to 0.45 for Koran, Subtitles, and 0.40 for other



Figure 2: Decoding speed(#Tok/Sec \uparrow) of Vanilla *k*NNMT and ours. Here, we set the batch size as 128.

datasets to achieve a balance between performance and efficiency.

757

758

759

760

761

Note that as the validation set is utilized in training the classifier network, there exists a potential risk of overfitting when tuning α_{min} , which may result in a suboptimal selection of α_{min} .

Datasets	0.45	0.40	0.35
IT	42.03 / 2978.73	42.30 / 2940.88	42.23 / 2878.88
Koran	19.53 / 3452.23	19.50 / 3415.35	19.58 / 3408.09
Law	58.66 / 3137.26	59.20 / 3097.68	59.31 / 3001.18
Medical	51.45 / 3155.22	51.75 / 3069.02	51.86 / 2989.31
Subtitles	32.05 / 1027.21	32.13 / 898.61	32.09 / 771.34

Table 10: BLEU \uparrow and #Tok/Sec \uparrow of our model on the multi-domain validation sets with different α_{min} . Here, we set the batch size as 128.