Mitigating the Language Mismatch and Repetition Issues in LLM-based Machine Translation via Model Editing

Anonymous ACL submission

Abstract

Large Language Models have recently revolutionized the NLP field, while they still fall short in some specific down-stream tasks. In the work, we focus on utilizing LLMs to perform machine translation, where we observe that two patterns of errors frequently occur and drastically affect the translation quality: language mismatch and repetition. The work sets out to explore the potential for mitigating these two issues by leveraging model editing methods, e.g., by locating FFN neurons or something that are responsible for the errors and deactivating them in the inference time. We find that directly applying such methods either limited effect on the targeted errors or has significant negative side-effect on the general translation quality, indicating that the located components may also be crucial for ensuring machine translation with LLMs on the rails. To this end, we propose to refine the located components by fetching the intersection of the locating results under different language settings, filtering out the aforementioned information that is irrelevant to targeted errors. The experiment results empirically demonstrate that our methods can effectively reduce the language mismatch and repetition ratios and meanwhile enhance or keep the general translation quality in most cases.

1 Introduction

001

017

024

037

041

Pre-trained Large Language Models (LLMs) are natural machine translators with in-context learning(Brown et al., 2020; Touvron et al., 2023; Vilar et al., 2023; Bawden and Yvon, 2023; Zhang et al., 2023a), while they still fall behind specialized Machine Translation (MT) systems like NLLB(Koishekenov et al., 2023). Previous studies utilise In-Context Learning (Agrawal et al., 2023) (ICL), instruction tuning(Xu et al., 2023; Alves et al., 2023) and post-editing methods(Jiao et al., 2023; Ki and Carpuat, 2024; Raunak et al., 2023)



Figure 1: The illustration of the language mismatch error (a) and the repetition error (b).

042

044

045

047

048

050

054

059

060

061

062

063

064

065

067

to improve the translation quality. Are there any issues that were ignored in previous studies hindering the LLM-based machine translation from further development? In this work, we identify two issues in the LLM-based machine translation: Language Mismatch and Repetition (as shown in Figure 1). We check the occurrence of these errors and find that: (1) they are common errors in the whole translation set (e.g., in the en \rightarrow de setting, language mismatch occurs in over 40% cases with zero-shot prompting); (2) they are severe errors for machine translation systems (e.g., repetition errors usually lead to a over 50% decrease on the BLEU score).

Nonetheless, the inherent reason for these errors still remains unclear, let alone patching them. In recent research works on model editing (Dai et al., 2022; Meng et al., 2022; Todd et al., 2023), they typically leverage analyzing tools like causal mediation analysis (Pearl, 2014; Vig et al., 2020), integrated gradient attribution (Sundararajan et al., 2017) to locate important component units (e.g., FFN neurons, attention heads and stuff) that are highly responsible for specific behavior patterns of LLMs, and then precisely control these behaviors by manipulating the located components (e.g.,

amplifying and suppressing the activation values 068 of neurons). Inspired by these works, we ask a 069 research question: Can we leverage model editing 070 methods to mitigate aforementioned language mismatch and repetition issues? We set out to adapt two widely-used model editing technique, Function Vectors (Todd et al., 2023) (FV) and Knowledge Neurons (Dai et al., 2022) (KN), to MT scenarios in an aim to locate error-relevant component units 076 inside LLMs. However, our empirical results show 077 that directly adapting FV and KN either has limited effect on the targeted errors or has significant sideeffect on the general translation quality, which indicates that the located component units may be not only responsible for targeted error patterns but also crucial for ensuring machine translation with LLMs on the rails and hence manipulate them could result in affecting the general translation behavior.

> We then aim to filter out the error-irrelevant components from the located results. A natural hypothesis is that the location for the important errorrelevant modules is supposed to be independent to translation language settings. Comparing the locating results under the different translation language settings (de \rightarrow en, en \rightarrow de, zh \rightarrow en and en \rightarrow zh), we do observe that a proportion of located component units are shared across different language settings. Grounded on this observation, we propose to refine the located components by fetching the intersection of the locating results under different language settings. The empirical results demonstrate that the modified methods can effectively reduce the language mismatch and repetition ratios and meanwhile keep or enhance the general translation quality in most cases.

094

100

101

102

103

104

106

108

109

110

111

112

113

114 115

116

117

118

119

Our main contributions are three-fold: (1) We identify two patterns of errors in LLM-based MT that frequently occur and badly affect the translation quality: language mismatch and repetition. (2) We investigate the potential for leveraging model editing methods (FV and KN) to reduce these errors. We find that directly adapting the editing methods either has limited effect on the targeted errors or has significant side-effect on the general translation quality. (3) We propose to refine the located modules by fetching the intersection of the locating results under different language settings. We show that with the refined locating results we could arouse the potential for editing methods to handle the language mismatch and repetition errors and meanwhile enhance or keep the general translation quality in most cases. The performance of our

methods could sometimes be comparable with traditional methods that adapt LLMs to MT tasks (e.g., 5-Shot ICL, LoRA and Full-FineTuning) without additional requirements like long-context prompting and fine-tuning. 120

121

122

123

124

125

126

127

128

129

130

131

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

169

2 Related Work

Large Language Models for Machine Translation One surprising ability of LLMs is that they are natural machine translators with Zero-Shot or One-Shot prompt(Brown et al., 2020; Touvron et al., 2023; Vilar et al., 2023; Bawden and Yvon, 2023; Robinson et al., 2023; Zhang et al., 2023a). However, there is still a gap(Xu et al., 2023) between pre-trained LLM and large-scale NMT systems like NLLB(Koishekenov et al., 2023) on the machine translation task. To bridge this gap, previous studies utilise in-context learning(Moslem et al., 2023; Agrawal et al., 2023; Bawden and Yvon, 2023; Vilar et al., 2023), model tuning(Xu et al., 2023; Alves et al., 2023; Zhang et al., 2023b), and interaction with annotation methods(Jiao et al., 2023; Ki and Carpuat, 2024) to improve the translation quality. Even though LLM has achieved massive success in machine translation(Kocmi et al., 2023a), some of the issues from LLM itself may challenge machine translation, such as Hallucination(Bang et al., 2023). Meanwhile, these problems from LLM are challenging to detect only with MT metrics. Alves et al. (2023) find few-shot tuning can improve the translation quality based on MT metrics(Papineni et al., 2002; Rei et al., 2022) but detect the machine translation hallucination with a case-based hallucination design.

Locating Based Model Editing Precisely locating a small set of important modules (e.g., neurons (Dai et al., 2022), hidden states (Todd et al., 2023), MHSA (Li et al., 2024) and MLP (Meng et al., 2022) outputs) and editing their values to steer large-scale models toward assumed behaviours (e.g., updating factual associations (Meng et al., 2022; Hase et al., 2023), detoxifying (Wang et al., 2024), decreasing hallucination (Li et al., 2023), switching languages (Tang et al., 2024) and patching reasoning errors (Li et al., 2024)) is a recently emerging paradigm. Nonetheless, such techniques are still largely under-explored in context of MT. In this work, we investigate the potential for adapting two representative locating based editing approaches (specifically, function vectors (Todd et al., 2023) and knowledge neurons (Dai et al., 2022)) 170 171

172

173

199

to the MT scenario to mitigate its two fundamental but crucial issues: language mismatch and repetition (Zhang et al., 2021).

3 Preliminary

In this section, we detail the data preparation process, including the data source, prompt template, and dataset construction. Additionally, we provide information about the model and the evaluation metrics used to support the ensuing experiments.

Data Source We choose three high-resource lan-179 guages: English, Chinese, German which show 180 good performance on MT tasks(Robinson et al., 181 2023). For the detailed language setting, we in-182 clude two language pairs: English-Chinese and English-German, and four translation directions: en \rightarrow de, de \rightarrow en, en \rightarrow zh and zh \rightarrow en (where en, de, zh represent English, German and Chinese, re-186 spectively). In the data choice, we use the human-187 made dataset from general MT tasks of WMT21, 188 WMT22 and WMT23¹ to ensure both high data quality and flexible data domain. These data make the machine translation approach a real-life usage 191 to help us understand the current state of machine 192 translation tasks using LLMs. 193

194**Prompt Template**For machine translation tasks,195a widely-adopted (Zhang et al., 2023a; Bawden and196Yvon, 2023; Vilar et al., 2023) K-Shot In-Context197Learning (ICL) prompt template (taking the lan-198guage setting of en \rightarrow zh for an example) is:

English : $src_1 \setminus n$ Chinese : $tgt_1 \setminus n$

English : $src_K \setminus n$ Chinese : $tgt_K \setminus n$ English : $src_q \setminus n$ Chinese :

200 Where (src_i, tgt_i) refers to the *i*-th in-context trans-201 lation exemplar $(src_i \text{ refers to a sentence of source}$ 202 language and tgt_i refers to the corresponding sen-203 tence of target language.). src_q refers to the real 204 sentecne of source language that needs to be trans-205 lated. We call this prompt template *Lang Prompt* 206 and regard it as the default prompt template for the 207 follow-up experiments in this paper.

208Dataset ConstructionIn the data construction209part, we construct the \mathcal{D}_{exps} (data from WMT21)210to provide the ICL exemplars used in the K-Shot211prompt for machine translation tasks. We use the

WMT22 data as the \mathcal{D}_{train} to fine-tune a model or locate the critical parts in an LLM for model editing methods. For the testing and validation, we construct the \mathcal{D}_{test} (data from WMT23) for various modifications (e.g. fine-tuning(Devlin et al., 2019) or model editing methods(Todd et al., 2023; Dai et al., 2022)). (Please refer to Appendix A for detailed dataset information)

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

Model To support the in-depth exploration and analysis of how the two kinds of errors happen. We use LLaMA2-7B as our backbone language model to implement the machine translation task and further adaptation(Touvron et al., 2023).

Evaluation Metrics For the machine translation metrics, we consider the overlapping-based metrics BLEU(Papineni et al., 2002) and neural-based metrics COMET(Rei et al., 2022) to evaluate the translation quality (For a detailed toolkit and detection process, please refer to Appendix B).

4 Language Mismatch and Repetition Error in LLM-MT

In our initial experiments, we observe that LLMbased machine translation struggles with the following two types of common errors. One is **Language Mismatch**, referring to the language of the translation result is not the target language. For example, In the en \rightarrow zh machine translation, the target language is Chinese while the language of generated sentence is still English. Another is **Repetition**, referring to a substring is generated repeatedly until the end of the generation. To evaluate these errors, we additionally introduce two metrics: Language **Mismatch Ratio** (LMR) (the percentage of cases occurring the language mismatch error) and **R**epetition **Ratio** (RR) (the percentage of cases occurring the repetition error).

Language mismatch and repetition error are common and crucial In this paragraph, we aim to find the ratio of language mismatch and repetition error in Zero-Shot and One-Shot. For detailed language settings, we consider $en \rightarrow de$, $de \rightarrow en$, $en \rightarrow zh$, and $zh \rightarrow en$. We use the \mathcal{D}_{test} and \mathcal{D}_{exps} as the test set and prompt exam source, respectively. We choose LER and RR to represent the ratio of language mismatch and repetition in a setting (e.g. $en \rightarrow de$ (Zero-Shot)). For translation quality evaluation, we choose the BLEU (Papineni et al., 2002) as the metrics. We observe that language mismatch is frequent in Zero-Shot and seldom in One-Shot.

¹https://github.com/wmt-conference/wmtX-newssystems. $X \in \{21, 22, 23\}$

Repetition error cases in One-Shot are without language mismatch but combined with language mismatch in Zero-Shot. Based on our observation, we do experiments and analysis in Zero-Shot for the language mismatch and in One-Shot for the repetition error.

261

262

269

270

271

275

276

278

279

282

283

286

293

To explore the relation between the above errors and translation quality, we split the translation results into four sets to evaluate the BLEU performance after error detection. The four sets include two error sets: language mismatch set and repetition error set, one regular set (where instance without both errors), and one Origin set that includes all cases. The results of Table 1 illustrate: (1) the gap between the regular set and the original set shows both language mismatch and repetition error hurt the translation quality; (2) Language mismatch is the main reason for the low performance in Zero-Shot; (3) Even though we observe a low repetition ratio in One-Shot, the gap between repetition set and regular set shows that repetition is a severe error in the original set; (4) The performance gap between regular and error cases indicates a direct way to improve the translation quality by eliminating these errors.

Setting	$\mathbf{L}(\downarrow)$	OB (↑)	$LB(\uparrow)$	$\mathbf{RB}(\uparrow)$
$zh \rightarrow en(\mathbf{Z})$	0.0486	17.13	8.77	17.60
en \rightarrow zh (Z)	0.3269	16.34	3.13	25.29
en \rightarrow de (Z)	0.4524	12.61	1.65	21.86
$de \rightarrow en(\mathbf{Z})$	0.0219	35.34	23.23	35.66
Setting	R (↓)	OB (↑)	RRB (↑)	RB (↑)
zh→en (O)	0.0035	18.87	2.13	19.06
en \rightarrow zh (O)	0.0146	27.78	2.08	29.47
$en \rightarrow de$ (O)	0.0141	24.97	12.64	25.86
$de \rightarrow en$ (O)	0.0018	36.54	6.10	36.71

Table 1: The correlation between error ratio and BLEU. (Z) represents the Zero-Shot prompting, and (O) represents the One-Shot prompting. L: language mismatch ratio; **R**: repetition ratio; **OB**: The BLEU on the original set; **LB**: The BLEU on the language mismatch set; **RRB**: The BLEU on the repetition error set; **RB**: The BLEU on the regular set.

5 Can we mitigate language mismatch and repetition via model editing?

In this section, We aim to investigate the potential for leveraging model editing methods (Dai et al., 2022; Meng et al., 2022; Todd et al., 2023) to precisely mitigate the aforementioned two severe issues in MT: language mismatch and repetition. We mainly focus on two widely-used model editing methods: Function Vectors (FV) (Todd et al., 2023) and Knowledge Neurons (KN) (Dai et al., 2022), for both of them are representative (i.e., Causal Mediation Analysis (Meng et al., 2022; Pearl, 2014) for FV and Integrated Gradient Attribution (Qi et al., 2019; Lundstrom et al., 2022) for KN) and influential (Bai et al., 2024; Hojel et al., 2024; Niu et al., 2024; Chen et al., 2024). In the following paragraphs, we adapt the idea of FV (corresponding to Machine translation vectors) and KN (corresponding to Machine translation neurons and Repetition neurons) to MT scenarios, in an aim to both enhance the LLMs' understanding to MT (for both language mismatch and repetition errors) and their specific ability to handle repetition errors.

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

337

338

339

341

342

5.1 Machine Translation Vectors

In the original FV paper, the authors argue that the key information of a task (e.g., task \mathcal{T} , demonstrated with few-shot In-Context Learning (ICL) exemplars in the prompt) is compactly represented and transported in a small set of attention heads in LLMs. The original definition for FV (of task \mathcal{T}) is the summation of these located head vectors. FV can be directly added into the "residual stream" (Elhage et al., 2021) of forwarding computation of Transformer-based (Vaswani et al., 2017) LLMs to help them perform ideal behaviour of task \mathcal{T} . Therefore, a natural question is: *Can* we use FV to enhance LLMs' understanding to MT and mitigate aforementioned language mismatch and repetition issues? We use Ten-Shot ICL prompts \mathcal{P} (the template of machine translation prompts is the Lang Prompt (Zhang et al., 2023a) in Section 3.) to locate important attention heads, where the data are sampled from \mathcal{D}_{train} . For brevity, we denote the normal Ten-Shot ICL input (omitting language signs, i.e., "English", "Chinese" and "German") as: inp = $[(src_1, tgt_1), (src_2, tgt_2), ..., (src_{10}, tgt_{10}), src_q]$ $\in \mathcal{P}$, where *src* and *tgt* refer to sentences of source and target languages respectively; index $1 \sim 10$ refers to ten ICL exemplars and q refers to "query" (the real source sentence that requires to be translated.). On its basis, we construct the shuffled version of the original ICL input: inp = $[(src_1, tgt_1), (src_2, tgt_2), ..., (src_{10}, tgt_{10}), src_a],$ where for each ICL exemplar (src_k, tgt_k) , $k \in [1..10]$, the target sentence $tgt_k \neq tgt_k$.

Extracting machine translation vectors First, we locate attention heads that are important to the MT with a Causal Mediation procedure: (1) Send both *inp* and *inp* to the same LLMs (denoting the model as θ), (2) Fetch both probabilities of predicting the ground-truth target sentence tgt_q from models with the above two inputs: $p_{\theta}(tgt_q|inp)$ and $p_{\theta}(tgt_q|inp)$, (3) Adopt intervention: replacing a single attention head output (e.g., for the *i*-th head in the *j*-th layer, the denotation is h_{i}^{i} .) in the shuffled run with inp with the attention head output at the same place (h_i^i) in the clear run with inp, (4) Calculate the Causal Indirect Effect (CIE $(h_j^i \rightarrow h_j^i | inp)$) of the intervention: $p_{\theta}(tgt_q|\widetilde{inp},\widetilde{h_j^i}\to h_j^i)-p_{\theta}(tgt_q|\widetilde{inp})$ and (5) Calculate the Average Indirect Effect for head h_i^i :
$$\begin{split} \text{AIE}(h^i_j) &= \mathop{\mathbb{E}}_{inp\in\mathcal{P}}[\text{CIE}(\widetilde{h^i_j} \to h^i_j | inp)]. \end{split}$$
 The AIE values for all heads in LLaMA2-7B

The AIE values for all heads in LLaMA2-7B under the language settings² of "de \rightarrow en" and "en \rightarrow zh" are depicted in Figure2. We observe that



Figure 2: Heatmaps of AIE values for attention heads in LLaMA2-7B for de \rightarrow en setting (a) and en \rightarrow zh setting (b). x-axis and y-axis refer to the layer and head. Brighter color refers to the head with larger AIE value.

for machine translation there are sparsely a few heads of which the corresponding AIE values strikingly stand out among 1024 heads. We select top-32 heads (according to their AIE values, denoted as \mathcal{H}) to extract FV in the follow-up experiments.

Let $h_j^i(inp)$ denote the output of attention head h_j^i given the input prompt *inp*. Following Todd et al. (2023), we extract the machine translation vector with a specific language setting (e.g., zh \rightarrow en) $\mathcal{V}_{zh\rightarrow en}$ with the following formula:

$$\mathcal{V}_{\text{zh}\to\text{en}} = \mathop{\mathbb{E}}_{inp\in\mathcal{P}_{\text{zh}\to\text{en}}} \left[\sum_{h_{i}^{i}\in\mathcal{H}} h_{j}^{i}(inp)\right]$$
(1)

Editing LLMs via machine translation vectors We directly add the extracted the machine transla-

Zero-Shot	$L(\downarrow)$	B(†)	C (†)
LLaMA2-7B	0.0486	17.1288	0.722
+MT vectors	-72.84%	-37.35%	-1.84%
+MT neurons	-18.72%	-4.28%	-0.15%
One-Shot	$\mathbf{R}(\downarrow)$	B(†)	C (†)
LLaMA2-7B	0.0035	18.8714	0.7376
+MT vectors	482.86%	-23.07%	-1.68%
+MT neurons	0.0%	-0.35%	-0.03%
+RP neurons	-8.57%	0.07%	0.0%

Table 2: Performance of LLaMA2-7B (and after applying model editing methods) on \mathcal{D}_{test} (under the language setting of $zh \rightarrow en$). *Zero-Shot* and *One-Shot* refer to that using zero-shot prompt (for language mismatch errors) and one-shot prompt (for repetition errors) for MT. For evaluation metrics, L: Language mismatch ratio; **R**: Repetition ratio; **B**: BLEU and **C**: COMET22DA, where **B** and **C** mainly evaluate the general translation quality. For plain LLaMA2-7B, the results are absolute values; for LLaMA2-7B with editing methods, the results are relative **improvement percentages**.

tion vector to the "residual stream" (being aligned with the original paper, at 11-th layer for LLaMA2-7B) in the forwarding process. The performance of LLaMA2-7B (e.g., under the language setting of $zh\rightarrow en$.) after adopting machine translation vectors are posted in Table 2.

376

377

378

379

380

381

382

383

384

385

388

389

390

391

392

393

394

395

396

398

399

400

401

402

403

404

405

406

We observe that leveraging machine translation vectors (+*MT vectors*) can (1) reduce the language mismatch errors to a large extent (-72.84%) while simultaneously (2) introduce more repetition errors (+482.86%) and (3) do harm to the general translation quality: -37.35% (Zero-Shot) and -23.07% (One-Shot) for BLEU.

5.2 Machine Translation Neurons and Repetition Neurons

The original KN technique is constructed on the basis of Geva et al. (2021)'s observation that the Feed-Forward Networks (FFNs) in the Transformer (Vaswani et al., 2017) can be viewed as keyvalue memories, which can store the encodings of a set of patterns (e.g., sentence end with specific words, specific topics and factual knowledge.). Given a specific pattern, KN can be used to locate a small set of neurons in FFNs that are the attribution of this pattern and manipulate the expression of this pattern by amplifying, suppressing or erasing the activation values of the located neurons. We ask a natural question: Can we use KN to locate and manipulate skilled neurons that are responsible for MT or the repetition error pattern? In the MT scenarios, We denote the input prompt inp (also

374

375

343

344

357

²Due to the page limit, We post experiment results only under part of the language settings results in the main text. For the rest language settings, we post them in Appendix C, Similarly hereinafter.

407omitting language sign) as $[src_q]$ (Zero-Shot) or408 $[(src_0, tgt_0), src_q]$ (One-Shot) and the correspond-409ing output as tgt_q , where the (src_0, tgt_0) is the ICL410exemplar (sampled from \mathcal{D}_{exps}) and (src_q, tgt_q) is411the "query", the real case used for locating neurons412(sampled from \mathcal{D}_{train}) or testing edited models413(sampled from \mathcal{D}_{test}).

Locating Important Neurons for MT We ran-414 domly sample a token t in each tgt_q and use t 415 to split tgt_q into two parts: $tgt_q = (tgt_q, tgt_q)$ 416 $(t \in t\vec{gt_q})$. To fully model the MT and mean-417 while restrict the computation, we focus on the 418 probability of $p(t|inp^+)$, where t refers to the 419 first token of $t\vec{gt_q}$ and inp^+ refers to the con-420 catenation of inp and $t \overset{\leftarrow}{gt_q}$. Focusing on a single 421 neuron $w_i^{(l)}$ (*i*-th intermediate neuron in the *j*-th 422 FFN), we denote its activation value as $\overline{w_i}^{(l)}$. Then 423 we can introduce this variable into $p(t|inp^+)$ as 424 $p(t|inp^+, w_i^{(l)} = \overline{w_i}^{(l)}) \triangleq f(\overline{w_i}^{(l)})$ (fixing t and 425 inp^+ , the probability can be viewed as an objective 426 function whose only variable is the value of neuron 427 $w_{i}^{(l)}$). We calculate the attribution score of neuron 428 $w_i^{(l)}$ by Integrated Gradient (Sundararajan et al., 429 2017): 430

431

432

433

434

435

$$\operatorname{Attr}(w_i^{(l)}|f) = \overline{w_i}^{(l)} \int_{\alpha=0}^1 \frac{\partial f(\overline{w_i}^{(l)})}{\partial w_i^{(l)}} d\alpha.$$
(2)

We calculate the mean value of the attribution scores for each neuron with 2,000 examples in \mathcal{D}_{train} and select top-5 neurons as Machine Translation neurons (*MT neurons*).

Locating Important Neurons for Repetition 436 We first collect all of examples that occur the repeti-437 tion error. For a specific input prompt *inp*, the com-438 pletion of LLMs y can be divided into the follow-439 ing several parts: $y = [y_{norm}, y_{repe}, y_{repe}, y_{rest}]$, 440 where y_{norm} refers to the normal generation part 441 (except for the first-time generation of y_{repe}), y_{repe} 442 refers to the minimal repetition unit (the first y_{repe} 443 here is supposed to be treated as normal gener-444 445 ation) and y_{rest} (the follow-up generation after the second-time generation of y_{repe}). To con-446 centrate on the repetition error, we construct a 447 new input prompt $inp_{repe} = [inp, y_{norm}, y_{repe}]$ 448 and focus on the probability of $p(y_{repe}|inp_{repe})$. 449 Similar to the MT neurons part, we define neu-450 ron $w_i^{(l)}$, its value $\overline{w_i}^{(l)}$, its objective function 451 $p(y_{repe}|inp_{repe}, w_i^{(l)} = \overline{w_i}^{(l)}) \triangleq f_{repe}(\overline{w_i}^{(l)})$ and 452

its attribution score $\operatorname{Attr}(w_i^{(l)}|f_{repe})$ (repetition attribution score). A natural concern here is that the objective function $f_{repe}(\overline{w_i}^{(l)})$ might model the pattern of generating y_{repe} rather than the repetition error pattern. To exclude this concern, we additionally set a comparison objective function $f_{compare} = p(y_{repe}|[inp, y_{norm}], w_i^{(l)} = \overline{w_i}^{(l)})$ to model the first-time generation (normal generation) of y_{repe} . With $f_{compare}$, we can also get the attribution score $Attr(w_i^{(l)}|f_{compare})$ (comparison attribution score) of neuron $w_i^{(l)}$. We calculate the mean values of repetition and comparison attribution scores separately for each neuron $w_i^{(l)}$ with all of the cases in \mathcal{D}_{train} that occur the repetition error. We separately select top-300 neurons according to mean repetition and comparison attribution score, denoting the fetched sets as \mathcal{N}_{repe} and $\mathcal{N}_{compare}$. We select 5 neurons with the largest repetition attribution scores from $\mathcal{N}_{repe} \setminus \mathcal{N}_{compare}$ as the Repetition Neurons (RP neurons).

453

454

455

456

457

458

459 460

461

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

Editing LLMs via MT neurons and RP neurons For MT neurons, we edit LLMs by amplifying the activation values of these neurons (set the new values to be twice the original ones). For *RP* neurons, we edit LLMs by *erasing* the activation values of these neurons (set the new values to be zero). The performance of LLaMA2-7B (e.g., under the language setting of zh->en.) after adopting MT neurons and RP neurons are posted in Table 2. We observe that (1) adopting MT neurons can indeed help reduce language mismatch ratio to some extent(-18.72%) while also bring small negative side-effect to the translation quality (-4.28%for the BLEU score), (2) adopting MT neurons nearly have no effect on the repetition ratio and (3) adopting RP neurons can reduce the repetition ratio slightly (-8.57%) without affecting the metrics (BLEU and COMET22DA) of evaluating general translation quality.

Hence a short response to the question of this section is that *Directly leveraging model editing methods either has limited effect on errors (MT neurons and RP neurons) or significant negative side-effect on general translation quality (MT vectors)*. Nonetheless, we do observe the potential for mitigating the aforementioned errors with editing methods.



(a) Language Mismatch Ratio (b) COMET22DA

Figure 3: Performance ((a) for the decrease percentage of LMR; (b) for the improvement percentage of COMET22DA) of intervention (blue bars) with language settings of $zh \rightarrow en$, $en \rightarrow zh$ and $de \rightarrow en$ on the heads located with the language setting of $en \rightarrow de$. The red bars (comparison group) refer to the results for intervention on random heads of the same number.

6 Modifications to FV and KN in MT scenarios

In section, we mainly discuss our modifications (Section 6.1) to FV and KN methods (Section 5) to release their potential for better mitigating the language mismatch errors, repetition errors and even improving the general translation quality. Besides, we present systematical evaluation results for the modified editing methods and baselines in Section 6.2.

6.1 Modifications

Previous empirical results (Section 5) show that MT vectors are more effective to reduce language mismatch errors in comparison with MT neurons 513 while the RP neurons are more promising for handling repetition errors, suggesting that the inher-515 ent mechanisms for the recognition of target lan-516 guage and generating strings repeatedly locate in 517 MHSA heads and FFN intermediate neurons, re-518 spectively. To this end, in the follow-up experi-519 ments, we concentrate on modifying MT vectors to handle language mismatch errors and RP neu-521 rons to handle repetition errors. Our first modifica-522 tion is based on a natural hypothesis: The location 523 for the important modules inside LLMs that are 524 responsible for target language recognition and repetition errors is supposed to be independent to language settings. The hypothesis can also be verified to some extent by the important head locat-528 ing experiments depicted in Figure 4, where results 530 for different language settings ($zh \rightarrow en$, $en \rightarrow zh$, $en \rightarrow de$ and $de \rightarrow en$) share a large proportion of 531 top heads. Moreover, we locate top-12 important attention heads in LLaMA2-7B under the language setting of en \rightarrow de and apply *MT vectors* to 534

Zero-Shot	$L(\downarrow)$	B (↑)	C (↑)
LLaMA2-7B	0.0486	17.1288	0.722
+MTV	-92.46%	-0.81%	2.65%
+MTV-I	-80.15%	53.5%	15.51%
+MTV-I-D	-86.12%	$\mathbf{76.82\%}$	16.02 %
One-Shot	$\mathbf{R}(\downarrow)$	B (↑)	C(↑)
LLaMA2-7B	0.0035	18.8714	0.7376
+RPN	-8.57%	0.07%	0.0 %
+RPN-I	$-\mathbf{25.71\%}$	0.51 %	-0.04%

Table 3: Performance of LLaMA2-7B (and after applying model editing methods) on \mathcal{D}_{test} (under the language settings of $zh \rightarrow en$ for *Zero-Shot* and $zh \rightarrow en$ for **One-Shot**). Other notations and abbreviations are following Table 2.

LLaMA2-7B with these located heads under the 535 language settings of $zh \rightarrow en$, $en \rightarrow zh$ and $de \rightarrow en$. 536 The results of Zero-Shot translation are depicted in 537 Figure 3 (experimental group, blue bars). We addi-538 tionally randomly select 12 heads to apply MT vec-539 tors and the results (comparison group) are shown 540 with red bars. We observe that for both the lan-541 guage mismatch ratio and COMET22DA, the per-542 formance of experimental group largely exceeds 543 the performance of comparison group under all 544 three other language settings, indicating that the 545 attention heads located under a single language set-546 ting can transfer to other language settings. Given 547 these evidences, we propose our first modification 548 to both MT vectors and RP neurons: We firstly lo-549 cate attention heads or FFN neurons separately 550 for each language setting and then get the final 551 located results by intersecting the located results 552 for all of language settings. We denote the MT 553 vectors fetched by intersected attention heads as 554 MT Vectors-Intersection (MTV-I) and intersected 555 RP neurons as **RePetition Neurons-Intersection** 556 (RPN-I). We post the results for leveraging MTV-I and *RPN-I* under the language settings of $en \rightarrow de$ 558 and $zh \rightarrow en$ in Table 3. We observe that: (1) for 559 MTV-I, the decrease percentage of language mis-560 match error ratio (-80.15%) is slightly lower than 561 MTV (-92.46%) while improvement percentage 562 of the BLEU score (53.5%) and COMET22DA 563 score (15.51%) exceed *MTV* (-0.81% and 2.65%) 564 by a large margin and (2) for RPN-I, the decrease 565 percentage of repetition error ratio (-25.71%) is 566 much higher than RPN (-8.57%), suggesting that 567 intersection of different language settings can filter 568 attention heads and FFN neurons that are irrelevant 569 to language mismatch errors and repetition errors 570 out. On the basis of MTV-I, we propose another 571 slight modification: Firstly calculate the MTV-I, 572

- 510
- 511 512
- 514

	de-	→en	en–	→de	zh-	→en	en—	≻zh
Zero-Shot	$L(\downarrow)$	B (†)	$L(\downarrow)$	B (↑)	$L(\downarrow)$	B (†)	$L(\downarrow)$	B (†)
LLaMA2-7B	0.0219	35.3448	0.4524	12.6084	0.0486	17.1288	0.3269	16.3441
+5-Shot ICL	-74.89%	4.93%	-92.06%	101.27%	-50.0%	12.46%	-82.59%	76.9%
+LoRA	-83.56%	0.68%	-95.25%	115.24%	-79.22%	6.62%	-77.58%	82.62%
+Full-FT	-8.68%	2.25%	-62.69%	55.41%	-33.33%	3.15%	-66.23%	62.64%
+MTV-I-D	-33.33%	-0.53%	-86.12%	76.82%	-54.12%	-14.08%	-69.9%	24.64%
One-Shot	$\mathbf{R}(\downarrow)$	B (†)	$\mathbf{R}(\downarrow)$	B (†)	R (↓)	B (†)	$\mathbf{R}(\downarrow)$	B (†)
LLaMA2-7B	0.0018	36.5445	0.0141	24.9685	0.0035	18.8714	0.0146	27.7798
+5-Shot ICL	0.0%	1.49%	14.89%	1.63%	-14.29%	2.07%	-17.12%	4.08%
+LoRA	-77.78%	-9.47%	-74.47%	-2.39%	5.71%	0.07%	-10.27%	0.37%
+Full-FT	22.22%	1.26%	-25.53%	4.9%	-22.86%	2.5%	22.6%	4.47%
+RPN-I	-38.89%	0.74%	-27.66%	0.35%	-25.71%	0.51%	-19.18%	-0.23%

Table 4: Overall Performance of LLaMA2-7B (and after applying model editing methods) on \mathcal{D}_{test} under all language settings. Other notations and abbreviations are following Table 2.

then divide it evenly according to the number of the intersected attention heads and add them to those heads. We denote this manner of leveraging *MTV-I* as *MTV-I-Distributional* (*MTV-I-D*). We also post the results of leveraging *MTV-I-D* in Table 3, where the results demonstrate that *MTV-I-D* can further achiever better performance than *MTV-I* in terms of language mismatch ratio, BLEU and COMET22DA.

6.2 Overall Results

573

574

575

576

577

579

580

581

583 584

585

589

590

591

592

593

594

595

596

598

603

To make readers get a better sense of the LLMs edited with our methods (MTV-I-D and RPN-I), we show the overall evaluation results for both our methods and traditional adaptation methods, including 5-Shot In-Context Learning (Brown et al., 2020) (5-Shot ICL), Low Rank Adaptation Tuning (Hu et al., 2022) (LoRA) and Full parameter Supervised Fine-Tuning (Alves et al., 2023) (Full-FT) for LLM-based MT in Table 4. Due to the page limit, we only post the performance on the metrics of language mismatch error ratio, repetition error ratio and BLEU score (We find that performance on COMET22 score is highly aligned with BLEU score). We observe that: (1) Applying the modified editing methods, MTV-I-D and RPN-I can generally reduce the error ratios for both language mismatch (\mathbf{L}) and repetition (\mathbf{R}) to a large degree, (2) The negative side-effect on the general translation quality (BLEU score, B) is minor (except when applying MTV-I-D under the setting of $\mathbf{zh} \rightarrow \mathbf{en}$, with a -14.08% decrease percentage on BLEU score). It is noteworthy that applying *MTV-I-D* can even improve the general translation quality to a large extent on the settings of $en \rightarrow de$

(76.82%) and **en** \rightarrow **zh** (24.64%) and (3) The performance of *MTV-I-D* and *RPN-I* can sometimes be comparable with (and even surpass) the traditional methods that adapt LLMs to the MT tasks, without additional requirements like long-context prompting and fine-tuning. 607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

7 Conclusion

In the work we find that two types of errors, language mismatch and repetition, occur frequently when performing the machine translation tasks with LLMs, bringing severe negative effect on the translation quality. We investigate the potentials of leveraging model editing methods to mitigate these issues and find that directly adopting function vectors and knowledge neurons may either have limited improvement on the error ratio metrics or bring noteworthy negative effect on the metrics that evaluate general machine translation quality (e.g., BLEU score), which indicates that the located attention heads and FFN neurons might be too coarse to only affect the error ratios without hurting the translation quality. To this end, we propose to refine the located attention heads and neurons by fetching the intersection of the locating results under different language settings. Our empirical results suggest that the modified function vectors and knowledge neurons methods (MTV-I-D and RPN-I) can effectively reduce the aforementioned two types of errors and even bring a positive influence on the translation quality metrics in most settings, indicating that there indeed exist a small set of modules that are highly responsible for the language mismatch and the repetition errors meanwhile.

741

742

743

744

745

689

690

Limitations

640

642

647

651

655

665

672

673

674

676

677

678

679

683

684

Our work is based on open-source LLaMA series models (Touvron et al., 2023)³. However, the effectiveness of these findings on other models, such as the open-sourced Baichuan 2 (Yang et al., 2023) or the close-sourced GPT-4 (OpenAI, 2023), remains unknown.

The model editing methods used in this paper require computational resources proportional to the size of the large language model (LLM). When applying our methods to a larger model, more computational resources will be necessary to achieve improved results. Our focus is on high-resource language settings for machine translation (MT). However, the observations and conclusions may differ when applied to low-resource or non-English language pair settings (e.g., zh \rightarrow de machine translation tasks)

We utilise automatic metrics for error and machine translation (MT) evaluation in our measurements. However, employing human-involved evaluations (Kocmi et al., 2023b) can offer a more profound understanding of the machine translation task with large language models (LLMs).

4 Ethics Statement

This paper utilizes a pre-trained large language model, with its training data sourced from web corpora that have not undergone ethical filtering. Consequently, it is capable of generating toxic content in the machine translation task (Wen et al., 2023). Moreover, we do not filter the source data or translation output in our work. Future research may build on our results to enhance the model, and we advocate for incorporating content supervision to prevent the dissemination of toxic content.

References

- Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke Zettlemoyer, and Marjan Ghazvininejad. 2023. Incontext examples selection for machine translation. In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 8857–8873. Association for Computational Linguistics.
 - Duarte M. Alves, Nuno Miguel Guerreiro, João Alves, José Pombal, Ricardo Rei, José Guilherme Camargo de Souza, Pierre Colombo, and André F. T. Martins. 2023. Steering large language models for machine translation with finetuning and in-context learning.

In Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023, pages 11127–11148. Association for Computational Linguistics.

- Yu Bai, Heyan Huang, Cesare Spinoso-Di Piano, Marc-Antoine Rondeau, Sanxing Chen, Yang Gao, and Jackie Chi Kit Cheung. 2024. Identifying and analyzing task-encoding tokens in large language models. *Preprint*, arXiv:2401.11323.
- Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. 2023. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. In Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics, IJCNLP 2023 -Volume 1: Long Papers, Nusa Dua, Bali, November 1 - 4, 2023, pages 675–718. Association for Computational Linguistics.
- Rachel Bawden and François Yvon. 2023. Investigating the translation performance of a large multilingual language model: the case of BLOOM. In Proceedings of the 24th Annual Conference of the European Association for Machine Translation, EAMT 2023, Tampere, Finland, 12-15 June 2023, pages 157–170. European Association for Machine Translation.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.
- Yuheng Chen, Pengfei Cao, Yubo Chen, Kang Liu, and Jun Zhao. 2024. Journey to the center of the knowledge neurons: Discoveries of language-independent knowledge neurons and degenerate knowledge neurons. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):17817–17825.
- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. Knowledge neurons in pretrained transformers. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8493– 8502, Dublin, Ireland. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of

³https://llama.meta.com/llama3/

- 746 747 748 749 750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771 772 773 774 775 775
- 772 773 774 775 776 777 778 779 780 781 782 783 784

791 792 793

801 802 deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers), pages 4171–4186. Association for Computational Linguistics.

- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, et al. 2021. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 1:1.
- Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. Transformer feed-forward layers are keyvalue memories. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5484–5495, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Peter Hase, Mohit Bansal, Been Kim, and Asma Ghandeharioun. 2023. Does localization inform editing? surprising differences in causality-based localization vs. knowledge editing in language models. In *Thirtyseventh Conference on Neural Information Processing Systems.*
- Alberto Hojel, Yutong Bai, Trevor Darrell, Amir Globerson, and Amir Bar. 2024. Finding visual task vectors. *Preprint*, arXiv:2404.05729.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. Lora: Low-rank adaptation of large language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net.
- Wenxiang Jiao, Jen-tse Huang, Wenxuan Wang, Zhiwei He, Tian Liang, Xing Wang, Shuming Shi, and Zhaopeng Tu. 2023. Parrot: Translating during chat using large language models tuned with human translation and feedback. In *Findings of the Association* for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023, pages 15009–15020. Association for Computational Linguistics.
- Dayeon Ki and Marine Carpuat. 2024. Guiding large language models to post-edit machine translation with error annotations. *CoRR*, abs/2404.07851.
- Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondrej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popovic, and Mariya Shmatova. 2023a. Findings of the 2023 conference on machine translation (WMT23): Ilms are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation, WMT*

2023, Singapore, December 6-7, 2023, pages 1–42. Association for Computational Linguistics.

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

- Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023b. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.
- Yeskendir Koishekenov, Alexandre Berard, and Vassilina Nikoulina. 2023. Memory-efficient NLLB-200: language-specific expert pruning of a massively multilingual machine translation model. In *Proceedings* of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023, pages 3567–3585. Association for Computational Linguistics.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023. Inferencetime intervention: Eliciting truthful answers from a language model. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Zhaoyi Li, Gangwei Jiang, Hong Xie, Linqi Song, Defu Lian, and Ying Wei. 2024. Understanding and patching compositional reasoning in llms. *Preprint*, arXiv:2402.14328.
- Daniel D Lundstrom, Tianjian Huang, and Meisam Razaviyayn. 2022. A rigorous study of integrated gradients method and extensions to internal neuron attributions. In *International Conference on Machine Learning*, pages 14485–14508. PMLR.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in GPT. In Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022.
- Yasmin Moslem, Rejwanul Haque, John D. Kelleher, and Andy Way. 2023. Adaptive machine translation with large language models. In *Proceedings of the* 24th Annual Conference of the European Association for Machine Translation, EAMT 2023, Tampere, Finland, 12-15 June 2023, pages 227–237. European Association for Machine Translation.
- Jingcheng Niu, Andrew Liu, Zining Zhu, and Gerald Penn. 2024. What does the knowledge neuron thesis have to do with knowledge? In *The Twelfth International Conference on Learning Representations*.

OpenAI. 2023. GPT-4 technical report. CoRR, abs/2303.08774.

860

873

874

875

876

877

878

893

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th Annual Meeting of the Association for Computational Linguistics, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Judea Pearl. 2014. Interpretation and identification of causal mediation. *Psychological methods*, 19(4):459.
- Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186– 191, Brussels, Belgium. Association for Computational Linguistics.
 - Zhongang Qi, Saeed Khorram, and Fuxin Li. 2019. Visualizing deep networks by optimizing with integrated gradients. In *CVPR workshops*, volume 2, pages 1–4.
- Vikas Raunak, Amr Sharaf, Yiren Wang, Hany Hassan Awadalla, and Arul Menezes. 2023. Leveraging GPT-4 for automatic translation post-editing. In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 12009–12024. Association for Computational Linguistics.
- Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Nathaniel R. Robinson, Perez Ogayo, David R. Mortensen, and Graham Neubig. 2023. Chatgpt MT: competitive for high- (but not low-) resource languages. In Proceedings of the Eighth Conference on Machine Translation, WMT 2023, Singapore, December 6-7, 2023, pages 392–418. Association for Computational Linguistics.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML'17, page 3319–3328. JMLR.org.
- Tianyi Tang, Wenyang Luo, Haoyang Huang, Dongdong Zhang, Xiaolei Wang, Xin Zhao, Furu Wei, and Ji-Rong Wen. 2024. Language-specific neurons: The key to multilingual capabilities in large language models. *Preprint*, arXiv:2402.16438.
- Eric Todd, Millicent L. Li, Arnab Sen Sharma, Aaron Mueller, Byron C. Wallace, and David Bau. 2023.
 Function vectors in large language models. *CoRR*, abs/2310.15213.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and finetuned chat models. CoRR, abs/2307.09288.

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. 2020. Investigating gender bias in language models using causal mediation analysis. *Advances in neural information processing systems*, 33:12388– 12401.
- David Vilar, Markus Freitag, Colin Cherry, Jiaming Luo, Viresh Ratnakar, and George F. Foster. 2023. Prompting palm for translation: Assessing strategies and performance. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 15406–15427. Association for Computational Linguistics.
- Mengru Wang, Ningyu Zhang, Ziwen Xu, Zekun Xi, Shumin Deng, Yunzhi Yao, Qishen Zhang, Linyi Yang, Jindong Wang, and Huajun Chen. 2024. Detoxifying large language models via knowledge editing. *Preprint*, arXiv:2403.14472.
- Jiaxin Wen, Pei Ke, Hao Sun, Zhexin Zhang, Chengfei Li, Jinfeng Bai, and Minlie Huang. 2023. Unveiling the implicit toxicity in large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1322– 1338, Singapore. Association for Computational Linguistics.
- Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023. A paradigm shift in machine translation: Boosting translation performance of large language models. *CoRR*, abs/2309.11674.

Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, Fan Yang, Fei Deng, Feng Wang, Feng Liu, Guangwei Ai, Guosheng Dong, Haizhou Zhao, Hang Xu, Haoze Sun, Hongda Zhang, Hui Liu, Jiaming Ji, Jian Xie, Juntao Dai, Kun Fang, Lei Su, Liang Song, Lifeng Liu, Liyun Ru, Luyao Ma, Mang Wang, Mickel Liu, MingAn Lin, Nuolan Nie, Peidong Guo, Ruiyang Sun, Tao Zhang, Tianpeng Li, Tianyu Li, Wei Cheng, Weipeng Chen, Xiangrong Zeng, Xiaochuan Wang, Xiaoxi Chen, Xin Men, Xin Yu, Xuehai Pan, Yanjun Shen, Yiding Wang, Yiyu Li, Youxin Jiang, Yuchen Gao, Yupeng Zhang, Zenan Zhou, and Zhiying Wu. 2023. Baichuan 2: Open large-scale language models. CoRR, abs/2309.10305.

973

974

975

988

989

991

1002

1004

1006

1007

1008

1009

1010

1012

1013

1014

1015

1016

1017

1020

1021

1022

- Biao Zhang, Barry Haddow, and Alexandra Birch. 2023a. Prompting large language model for machine translation: A case study. In International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA, volume 202 of Proceedings of Machine Learning Research, pages 41092-41110. PMLR.
- Shaolei Zhang, Qingkai Fang, Zhuocheng Zhang, Zhengrui Ma, Yan Zhou, Langlin Huang, Mengyu Bu, Shangtong Gui, Yunji Chen, Xilin Chen, and Yang Feng. 2023b. Bayling: Bridging cross-lingual alignment and instruction following through interactive translation for large language models. CoRR, abs/2306.10968.
 - Ying Zhang, Hidetaka Kamigaito, Tatsuya Aoki, Hiroya Takamura, and Manabu Okumura. 2021. Generic mechanism for reducing repetitions in encoderdecoder models. In Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021), pages 1606–1615, Held Online. INCOMA Ltd.

Dataset Information Α

Table 5 shows the detailed data size for \mathcal{D}_{exps} , \mathcal{D}_{train} and \mathcal{D}_{test} . We use the WMT21 test set⁴ as the \mathcal{D}_{exps} , WMT22 test set⁵ as \mathcal{D}_{train} and WMT23 test set⁶ as \mathcal{D}_{test} .

The detailed data size for the K-shot (K =(0, 1, 5) setting is shown in Table 6. For all settings, we use the *lang prompt* as the prompt template (as shown in Section 3). For the Zero-Shot setting, we directly combine the source data with the *lang prompt.* For the One-Shot setting, we uniformly sample the data from \mathcal{D}_{exps} based on the length of the example source to alleviate the potential length

Setting	\mathcal{D}_{exps} Size	\mathcal{D}_{train} Size	\mathcal{D}_{test} size
en→de	1002	2037	557
$de{\rightarrow}en$	1000	1984	549
$en{\rightarrow}zh$	1002	2037	2074
$zh{\rightarrow}en$	1948	1875	1976

Table 5: Data size of \mathcal{D}_{exps} , \mathcal{D}_{train} , \mathcal{D}_{test} on four language settings.

bias from prompt example(Zhang et al., 2023a). We use the most natural selection method for the Five-Shot setting by randomly selecting five examples from \mathcal{D}_{exps} .

1024

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1041

1042

Setting	\mathcal{D}_0 Size	\mathcal{D}_1 Size	\mathcal{D}_5 size
en \rightarrow de (\mathcal{D}_{train})	2037	12222	2037
de \rightarrow en (\mathcal{D}_{train})	1984	9920	1984
en \rightarrow zh (\mathcal{D}_{train})	2037	12222	2037
$zh \rightarrow en (\mathcal{D}_{train})$	1875	11250	1875
en \rightarrow de (\mathcal{D}_{test})	557	3342	557
de \rightarrow en (\mathcal{D}_{test})	549	2745	549
en \rightarrow zh (\mathcal{D}_{test})	2074	12444	2074
$zh \rightarrow en (\mathcal{D}_{test})$	1976	11856	1976

Table 6: Data size of Zero-Shot (\mathcal{D}_0) , One-Shot (\mathcal{D}_1) and Five-Shot(\mathcal{D}_5) on four language settings. \mathcal{D}_{train} and \mathcal{D}_{test} represent the source data in the prompt.

Toolkits for evaluation B

For the language mismatch detection, we use the Polyglot toolkit⁷ to detect the language error. For repetition error, based on the definition of repetition error, we follow two rules to judge whether a translation result is repeated: (1) the generation length reaches the max_new_tokens setting⁸; (2) there exists a substring happening until the end of the generation. For the machine translation metrics, we use SacreBLEU(Post, 2018), Unbabel/wmt22comet-da⁹ and Unbabel/wmt22-cometkiwi-da¹⁰ to do evaluation.

The AIE values for all heads С

Figure 4 shows the AIE values of all heads of LLaMA2-7B on $en \rightarrow de$, $de \rightarrow en$, $en \rightarrow zh$ and $zh \rightarrow en$ settings.

⁴https://github.com/wmt-conference/ wmt21-news-systems

⁵https://github.com/wmt-conference/

wmt22-news-systems

⁶https://github.com/wmt-conference/ wmt23-news-systems

⁷https://github.com/aboSamoor/polyglot

⁸https://github.com/huggingface/tokenizers ⁹https://huggingface.co/Unbabel/

wmt22-comet-da

¹⁰https://huggingface.co/Unbabel/ wmt22-cometkiwi-da



Figure 4: Heatmaps of AIE values for attention heads in LLaMA2-7B for $en \rightarrow de$ setting (a), $de \rightarrow en$ setting (b), $en \rightarrow zh$ setting (c) and $zh \rightarrow en$ setting (d). The xaxis and y-axis refer to the layer and head, respectively. Brighter color refers to the head with larger AIE value.

D Results for direct adaptation

The complete results of direct adaptation on four language settings are shown in Table 7 (en \rightarrow de), 8 (de \rightarrow en), 9 (en \rightarrow zh) and 10 (zh \rightarrow de).

These tables show that the MT vectors can decrease the language mismatch ratio while the RP neurons help decrease repetition errors in all language settings.

Zero-Shot	$L(\downarrow)$	B(†)	C (†)
LLaMA2-7B	0.4524	12.6084	0.6113
+MT vectors	-92.46%	-0.81%	2.65%
+MT neurons	-11.1%	1.78%	0.15%
One-Shot	R (↓)	B (↑)	C (†)
LLaMA2-7B	0.0141	24.9685	0.7279
+MT vectors	487.94%	-39.11%	-10.87%
+MT neurons	4.26%	-1.05%	-1.06%
+RP neurons	-27.66%	0.77%	-0.3%

Table 7: Performance of LLaMA2-7B (and after applying model editing methods) on \mathcal{D}_{test} (under the language setting of $\mathbf{en} \rightarrow \mathbf{de}$). **Zero-Shot** and **One-Shot** refer to using a zero-shot prompt (for language mismatch errors) and one-shot prompt (for repetition errors) for MT tasks. For evaluation metrics, L: Language mismatch ratio; **R**: Repetition ratio; **B**: BLEU and **C**: COMET22DA, where **B** and **C** mainly evaluate the general translation quality. For plain LLaMA2-7B, the results are absolute values; for LLaMA2-7B with editing methods, the results are relative **improvement percentages**.

Zero-Shot	$L(\downarrow)$	B (↑)	C (†)
LLaMA2-7B	0.0219	35.3448	0.7836
+MT vectors	-74.89%	-33.85%	-5.53%
+MT neurons	8.22%	0.03%	0.23%
One-Shot	$\mathbf{R}(\downarrow)$	B(†)	C (↑)
LLaMA2-7B	0.0018	36.5445	0.7893
+MT vectors	727.78%	-33.62%	-4.38%
+MT neurons	22.22%	-0.35%	-0.11%
+RP neurons	%	%	%

Table 8: Performance of LLaMA2-7B (and after applying model editing methods) on \mathcal{D}_{test} (under the language setting of $\mathbf{de} \rightarrow \mathbf{en}$). The -- means the same result as the LLaMA2-7B since we do not detect any repetition on the training set under the same language setting. Notation and corresponding explanations can refer to Table 7.

Zero-Shot	$L(\downarrow)$	B(†)	C(†)
LLaMA2-7B	0.3269	16.3441	0.6567
+MT vectors	-70.05%	18.2%	5.07%
+MT neurons	-5.32%	3.16%	0.35%
One-Shot	$\mathbf{R}(\downarrow)$	B (↑)	C (↑)
LLaMA2-7B	0.0146	27.7798	0.7444
+MT vectors	162.33%	-15.29%	-4.0%
+MT neurons	5.48%	-4.28%	-0.28%
+RP neurons	-4.11%	0.55%	0.05%

Table 9: Performance of LLaMA2-7B (and after applying model editing methods) on \mathcal{D}_{test} (under the language setting of **en** \rightarrow **zh**). *Zero-Shot* and *One-Shot* refer to using a zero-shot prompt (for language mismatch errors) and one-shot prompt (for repetition errors) for MT tasks. Notation and corresponding explanations can refer to Table 7.

Zero-Shot	$L(\downarrow)$	B(†)	C(†)
LLaMA2-7B	0.0486	17.1288	0.722
+MT vectors	-72.84%	-37.35%	-1.84%
+MT neurons	-18.72%	4.28%	-0.15%
One-Shot	$\mathbf{R}(\downarrow)$	B(†)	C (†)
LLaMA2-7B	0.0035	18.8714	0.7376
+MT vectors	482.86%	-23.07%	-1.68%
+MT neurons	0.0%	-0.35%	-0.03%
+RP neurons	-8.57%	0.07%	0.0%

Table 10: Performance of LLaMA2-7B (and after applying model editing methods) on \mathcal{D}_{test} (under the language setting of $\mathbf{zh} \rightarrow \mathbf{en}$). **Zero-Shot** and **One-Shot** refer to using a zero-shot prompt (for language mismatch errors) and one-shot prompt (for repetition errors) for MT tasks. Notation and corresponding explanations can refer to Table 7.

E Results for improved adaptation

Table 12, 11, 14 and 13 show the results for improved adaptation on $en \rightarrow de$, $de \rightarrow en$, $en \rightarrow zh$ and $zh \rightarrow en$ respectively.

Zero-Shot	$\mathbf{L}(\downarrow)$	$\mathbf{B}(\uparrow)$	$\mathbf{C}(\uparrow)$
LLaMA2-7B	0.0219	35.3448	0.7836
+MTV	-74.89%	-33.85%	0.0036 %
+MTV-I	-58.45%	-4.84%	-5.53%
+MTV-I-D	-33.33%	-0.53%	-0.22%
One-Shot	$\mathbf{R}(\downarrow)$	B (†)	C(†)
LLaMA2-7B	0.0018	36.5445	0.7893
+RPN	%	%	%
+RPN-I	%	%	%

Table 11: Performance of LLaMA2-7B (and after applying model editing methods) on \mathcal{D}_{test} (under the language settings of de—en for **Zero-Shot** and de—en for **One-Shot**). The – means the results is the same as the LLaMA2-7B since there is no repetition cases in the \mathcal{D}_{train} . Other notations and abbreviations following Table 7.

Zero-Shot	$L(\downarrow)$	B (↑)	C(†)
LLaMA2-7B	0.4524	12.6084	0.6113
+MTV	-92.46%	-0.81%	2.65%
+MTV-I	-80.15%	53.5%	15.51%
+MTV-I-D	-86.12%	$\mathbf{76.82\%}$	16.02 %
One-Shot	$\mathbf{R}(\downarrow)$	$\mathbf{B}(\uparrow)$	C(†)
LLaMA2-7B	0.0141	24.9685	0.7279
+RPN	-27.66%	0.77 %	-0.3%
+RPN-I	$-\mathbf{27.66\%}$	0.35%	-0.03%

Table 12: Performance of LLaMA2-7B (and after applying model editing methods) on \mathcal{D}_{test} (under the language settings of en \rightarrow de for **Zero-Shot** and en \rightarrow de for **One-Shot**). Other notations and abbreviations following Table 7.

Zero-Shot	$L(\downarrow)$	B (†)	C(†)
LLaMA2-7B	0.0486	17.1288	0.722
+MTV	-92.46%	-0.81%	2.65%
+MTV-I	-80.15%	53.5%	15.51%
+MTV-I-D	-86.12%	76.82 %	16.02 %
One-Shot	R (↓)	B(†)	C (†)
LLaMA2-7B	0.0035	18.8714	0.7376
+RPN	-8.57%	0.07%	0.0 %
+RPN-I	- 25.71 %	0.51 %	-0.04%

Table 13: Performance of LLaMA2-7B (and after applying model editing methods) on \mathcal{D}_{test} (under the language settings of zh \rightarrow en for **Zero-Shot** and zh \rightarrow en for **One-Shot**). Other notations and abbreviations are following Table 10.

Zero-Shot	$L(\downarrow)$	$\mathbf{B}(\uparrow)$	C(↑)
LLaMA2-7B	0.3269	16.3441	0.6567
+MTV	-70.05%	18.2%	5.07%
+MTV-I	-67.27%	19.08%	7.54%
+MTV-I-D	-69.9%	24.64 %	$\mathbf{8.82\%}$
One-Shot	$\mathbf{R}(\downarrow)$	$\mathbf{B}(\uparrow)$	C (†)
LLaMA2-7B	0.0146	27.7798	0.7444
+RPN	-4.11%	0.55%	0.05 %
+RPN-I	-19.18%	0.01%	-0.23%

Table 14: Performance of LLaMA2-7B (and after applying model editing methods) on \mathcal{D}_{test} (under the language settings of en \rightarrow zh for **Zero-Shot** and en \rightarrow zh for **One-Shot**). Other notations and abbreviations are following Table 7.

F Implementation Details

For all machine translation results on LLMs, we only maintain the first line of the generation as the translation result based on the format of *lang prompt*. In the real translation process, we use batch generation techniques (batch size = 4) and
set the maximum generation length of tokens to
400 with the Huggingface API¹¹ to do translations for any setting in this work.

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1071

1072

1073

1074

1075

1076

1077

1079

1080

1081

1082

Five-Shot For the Five-Shot setting, we directly use the D_5 on LLaMA2-7B to run machine translation task without intervention.

LoRA fine-tuning LoRA (Low-Rank Adaptation)(Hu et al., 2022) is a parameter-efficient tuning technique generally used in natural language processing. In our work, we use the LoRA(Hu et al., 2022) method to align the LLaMA2-7B model to the machine translation task. For the fine-tuning data, we combine the data of all language settings from \mathcal{D}_{train} into \mathcal{D}_0 and \mathcal{D}_1 for Zero-Shot setting and One-Shot setting respectively. Finally, we tune two LoRA models with the trl tool¹² with the selfsupervised tuning method. We train one epoch with a rank of 64 and a learning rate of $2e^{-4}$ for both Zero-Shot and One-Shot. We use one NVIDIA A100 80GB Tensor Core GPU card for the SFT training; either the Zero-Shot or One-Shot costs less than a half day.

Full fine-tuningWe use the same data and train-ing tool in the LoRA setting for full fine-tuning. In1083the training process, we use the bfloat16 precious1085to train the model on one NVIDIA A100 80GB1086

¹¹https://huggingface.co/

¹²https://github.com/huggingface/trl

Tensor Core GPU card for full fine-tuning with a 1087 lower learning rate $1e^{-6}$ compared to LoRA. 1088 We claim there is still room for improvement in 1089 the LoRA or Full fine-tuning methods. However, a 1090 complete understanding of the mismatch and repe-1091 1092 tition error should also be evaluated on large-scale data, which is one of the following steps for our 1093 research. 1094