# Can You Translate for Me? Code-Switched Machine Translation with Large Language Models

**Jyotsana Khatri*, Vivek Srivastava*, Lovekesh Vig**

TCS Research, India

{jyotsana.khatri, srivastava.vivek2, lovekesh.vig}@tcs.com

## Abstract

Large language models (LLMs) have shown remarkable performance on a variety of multilingual NLP tasks. Code-switching is one of the most convenient styles of communication in multilingual communities. It is known to present several challenges to the existing language models and task-specific models. In this paper, we evaluate the capability of multilingual LLMs for the code-switched machine translation (CSMT) task in traditional and novel settings and present our insights. We observe that ChatGPT outperforms other LLMs and shows competitive performance to the supervised fine-tuned models. Though promising, ChatGPT shows major limitations, such as high gender bias, stereotypes, and factual inconsistencies. It further demands a multi-dimensional large-scale evaluation of the multilingual LLMs for code-switched languages.

## 1 Introduction

Large language models (LLMs) have significantly advanced the performance on a number of NLP tasks using zero-shot setting and in-context learning (Brown et al., 2020). Machine translation is one of the most challenging and widely explored research areas in NLP, and is heavily impacted by the powerful LLMs (Wei et al., 2022a; Zhu et al., 2023; Jiao et al., 2023; Lyu et al., 2023; Wang et al., 2023). Though there is a phenomenal opportunity with LLMs-based machine-translation involving several low and medium resource languages, the performance of these models remains a mystery on the code-switched languages. The research with CSMT is in a nascent stage with new benchmarks and evaluation strategies in place (Chen et al., 2022; Srivastava and Singh, 2022). Owing to this (LLM) thrust, we systematically analyze the CSMT task in several configurations including the evaluation on existing CSMT benchmarks. We report our findings with multiple LLMs

*Equal contribution.

including ChatGPT (Lütkebohle), BLOOMZ-7b1 (Scao et al., 2022), XGLM-7.5B (Lin et al., 2021), mT0-xxl (Muennighoff et al., 2022), and mT0-xxl-mt (Muennighoff et al., 2022). We believe that our work would encourage future works to explore CSMT and its evaluation with a novel and interesting outlook. In this paper, we focus on the following research questions:

- How effective are the LLMs for the CSMT task?
- To what extent, can we instruct and control the code-switched text generation using LLMs?
- Do LLMs posses the common sense reasoning capability in a code-switched setting?
- Are LLMs gender-biased? If yes, how would it impact the CSMT task?

## 2 Related Work

CSMT is a challenging and under-explored task. Due to resource scarcity, there have been efforts to explore the utilization of back-translated data in NMT (Jawahar et al., 2021). There are multiple efforts towards fine-tuning pre-trained language models, and generating pseudo parallel data (Winata et al., 2019; Gautam et al., 2021; Jawahar et al., 2021; Srivastava and Singh, 2022; Solorio et al., 2021). However, the use of large language models is still unexplored for CSMT.

The initial works of unsupervised NMT were based on three concepts: denoising, cross-lingual embeddings, and iterative back-translation (Artetxe et al., 2018b,a, 2019; Lample et al., 2018a,b). Later, multilingual pretraining gained a lot of attention where language model pretraining is performed using a large number of languages (Conneau and Lample, 2019; Song et al., 2019; Lewis et al., 2019; Siddhant et al., 2020; Liu et al., 2020).

Large language models have performed well for various NLP tasks using in-context learning (Brown et al., 2020; Dong et al., 2022; Scao et al., 2022; Vilar et al., 2022; Zeng et al., 2022; Ren et al., 2023; Wei et al., 2022b). (Zhu et al., 2023)

83

showed the performance of translation using various large language models and observed that Chat-GPT performs better than all others but it still lacks behind the supervised models. Recently, there have been several works exploring the performance of machine translation for various language-pairs using LLMs (Lin et al., 2022; Jiao et al., 2023; Lyu et al., 2023; Bang et al., 2023; Agrawal et al., 2022; Zhang et al., 2023; Moslem et al., 2023).

## 3 Experimental Setup

In this section, we present the details of our experimental setup. We first discuss the datasets used followed by a detailed discussion on the CSMT with multilingual LLMs.

### 3.1 Datasets

In our experiments, we leverage datasets from four different sources in various configurations. A brief overview of these datasets is as follows:

1. **CALCS 2021** (Chen et al., 2022): The shared-task on "Machine Translation for Code-Switched Data" was hosted along with the Computational Approaches to Linguistic Code-Switching (CALCS) 2021 workshop. In the supervised setting, they provide a parallel dataset of 9,962 samples to translate English into code-switched Hindi-English in a single direction. In the unsupervised setting, they provide the data for the following language pairs: English and Spanish-English, and English and Modern Standard Arabic Egyptian Arabic in both directions.

2. **MixMT 2022** (Srivastava and Singh, 2022): The shared-task on "Code-Mixed Machine Translation" was organized along with the Workshop on Machine Translation, WMT 2022. The organizers provide an evaluation set including 500 samples in the validation set and 1,500 samples in the test set for the English and code-switched Hinglish (in both direction) translation pair.

3. **Wino-X** (Emelin and Sennrich, 2021): Wino-X is a multilingual extension of the widely popular Winograd schemas (Winograd, 1972) to evaluate the coreference resolution and common-sense reasoning capabilities of the models. The multilingual parallel Wino-X dataset (which is derived from WinoGrande dataset (Sakaguchi et al., 2021)) comprises of German, French, and Russian schemas aligned with the English language schemas. The dataset contains 1,887,

1,499, and 1,119 schemas in the parallel German, French, and Russian languages respectively.

4. **WinoMT** (Stanovsky et al., 2019): WinoMT dataset presents a challenging set of samples to evaluate the gender-bias in machine translation systems by disambiguating the gendered pronoun with non-stereotypical gender roles while translating. WinoMT dataset is created by concatenating the Winogender (Rudinger et al., 2018) and WinoBias (Zhao et al., 2018) coreference test sets and it contains 3,888 samples, balanced between male and female genders, as well as between stereotypical and non-stereotypical gender-role assignments.

### 3.2 Code-Switched Machine Translation

With the advent of LLMs, it is of utmost importance to (re)design the various aspects of the experimentation and evaluation strategies with CSMT. To this extent, we explore the capabilities of LLMs as a zero-shot CSMT system in the following settings:

1. **Minimally Instructed Translation (MIT)**: The current LLMs are instructed via a prompt $P$ to achieve the desired result in a typical conversational fashion. The instruction encoded in $P$ gears the response from the LLMs. In this setup, we engineer the prompt such that minimal required information is presented to the LLMs for the CSMT task. Formally, we denote the translation with MIT as:

$$T_{MIT}(s) = F(P(I_{min}, s)) \qquad (1)$$

Here, we encode $s$ (source text) and $I_{min}$ (minimal instruction) to create the prompt $P$ which is subsequently used to prompt the LLM (denoted as $F$). In Table 6 and 8 (in the Appendix), we present the minimal instructions used in this experiment.

2. **Linguistically Constrained Translation (LCT)**: Next, we provide additional linguistic constraints to $I_{min}$. We experiment with four different linguistic constraints involving a combination of language, script, and part of speech (see Table 6 in the Appendix) to bring out the language understanding capability of LLMs. Formally, we denote the translation with LCT as:

$$T_{LCT}(s) = F(P(I_{LCT}, s)) \qquad (2)$$

such that,

$$I_{LCT}(s) = I_{min} \odot L \qquad (3)$$

Here, $L$ denotes the linguistic constraints encoded with $s$ and $I_{min}$ to create the prompt $P$. Also, $\odot$ denotes the concatenation operation.

3. **Co-reference Resolution (CR)**: Co-reference resolution is a simple yet powerful mechanism to evaluate the basic commonsense reasoning capability of the LLMs. In this novel experiment, we evaluate the CR capability of LLMs in the context of CSMT. Given an English sentence $s_{en}$ with pronoun co-reference, we manually create two parallel code-switched sentences ($s_{cs}^1$ and $s_{cs}^2$) after the co-reference resolution of the pronoun where only one of the co-reference resolutions is correct. We prompt the LLM to select the correct code-switched translation (see Table 6 in the Appendix). Formally,

$$T_{CR}^{select}(s_{en}) = F(P(I_{select}, s_{en}, s_{cs}^1, s_{cs}^2)) \qquad (4)$$

4. **Gender Debiasing (GD)**: The LLMs, trained on the real-world data, tend to pick up the inherent societal stereotypes and gender biases. In this experiment, we evaluate these stereotypes and biases with the CSMT task. We select an English sentence $s_{en}$ with a gendered pronoun referencing to non- non-stereotypical profession. The sentence $s_{en}$ also contains a stereotypical profession (see Table 10 in the Appendix). We manually create two parallel code-switched sentences ($s_{cs}^1$ and $s_{cs}^2$) after the resolution of the gendered pronoun with the stereotypical and non-stereotypical professions. We then evaluate the gender bias in LLMs with two prompting strategies. First, we prompt the LLM to select the correct code-switched translation. Formally,

$$T_{GD}^{select}(s_{en}) = F(P(I_{select}, s_{en}, s_{cs}^1, s_{cs}^2)) \qquad (5)$$

Next, we prompt the LLM to translate the source English sentence to the code-switched sentence and also explicitly disambiguate the gendered pronoun. Formally,

$$T_{GD}^{translate}(s_{en}) = F(P(I_{translate}, s_{en})) \qquad (6)$$

We present the instructions $I_{select}$ and $I_{translate}$ for this task in Table 6 in the Appendix.

## 3.3 A Pilot Study on Multilingual LLMs

We conduct a pilot study of various LLMs for *English → Hinglish* CSMT using a small subset of 25 samples randomly selected from the CALCS 2021 development dataset. The prompts are presented in Table 7. We evaluate the output with BLEU score and TER calculated using sacrebleu[*].We observe that, XGLM-7.5B worked like a text completion model and did not produce the desired translations. The translated sentences with mT0-xxl are in the *Devanagari* script but the reference translations are in Roman script resulting in a further drop in the BLEU-score (see Table 1). Furthermore, mT0-xxl-mt is a generic fine-tuned model for multilingual machine translation task but fails to produce good quality output while following minimal instructions. The translations with BLOOMZ-7b1 model are English-only with no code-switched sentence obtained at the target side. Overall, Chat-GPT[*] outperforms the other LLMs by a significant margin, which drives us to further explore its capability with the other challenging experimental formulations discussed in Section 3.2.

| Model | BLEU-score | TER |
|---|---|---|
| XGLM-7.5B | 0.4867 | 98.97 |
| mT0-xxl | 0.4523 | 109.43 |
| mT0-xxl-mt | 0.6133 | 110.20 |
| BLOOMZ-7b1 | 2.1598 | 96.17 |
| ChatGPT | **10.5165** | **85.20** |

Table 1: BLEU score and Translation Error Rate (TER) score for the pilot study of various LLMs.

## 4 Results and Analysis

In this section, we present the results from different experiments and discuss our observations.

1. **Minimally Instructed Translation**: We conduct the MIT experiments on the CALCS 2021 (see Table 2) and MixMT 2022 datasets (see Table 3). For the evaluation of MixMT, we report the TER score calculated using sacrebleu to compare it with the state of the art. We observe that the zero-shot performance of Chat-GPT on these benchmarks is competitive to the existing supervised fine-tuned models. In the future, it would be interesting to see how different innovative prompting strategies such as chain-of-thought (Wei et al., 2022c) further help

---

[*] https://github.com/mjpost/sacrebleu
[*] We used GPT3.5-turbo May12, 2023 API version in our experiments.

improve the model performance.

| Model | En-Hg | En-Sg | Sg-En |
|---|---|---|---|
| Amazon-IML (Comix) | 12.98 | - | - |
| UBC_HImt (mT5) | 12.67 | - | - |
| LTRC-PreCog (mBART-en) | 12.22 | - | - |
| B2BT EMNLP 2022 Findings | - | - | 50.37 |
| ChatGPT | 9.66 | 61.58 | 46.54 |

Table 2: CALCS 2021 test set BLEU score (using the leaderboard https://ritual.uh.edu/lince/). Here, En: English, Hg: Hinglish, and Sg: Spanglish. The top-4 systems are the current best performing systems on the LinCE benchmark leaderboard.

| Language-pair | ChatGPT | Khan et al. |
|---|---|---|
| Hinglish → English | 0.732 | 0.607 |
| English (Hindi) → Hinglish | 0.750 | 0.547 |

Table 3: MixMT 2022 test set TER score. Khan et al. was the best performing system at MixMT 2022.

2. **Linguistically Constrained Translation**: We conduct the LCT experiments on the randomly sampled 25 English sentences from the CALCS 2021 *English-Hinglish* dataset and translate them to Hinglish using ChatGPT. We manually evaluate the translated sentence on the four quality dimensions: correctness, fluency, code-switched, and instruction following. The evaluator rate the correctness and the fluency measures on a scale of 1-5 (low-high). The 'code-switched' metric measures the binary outcome i.e., 0 (code-switched) or 1 (monolingual). The 'instruction following' metric measures the capability of the model to correctly follow the passed instruction (translation with linguistic constraints). The evaluator assigns a score of 0/1 based on whether the instruction is followed (1) or not (0). We further verify the evaluation by the evaluator with manual verification by another evaluator. The disagreement is resolved mutually by the evaluators. We report the results in Table 4. We observe that the overall correctness and fluency decreases as we increase the number of constraints in the instruction. Empirically, we observe that the model tends to generate a highly monolingual sentence in the Hindi language along with a only few English language words.

3. **Co-reference Resolution**: We manually select 25 English instances from the Wino-X dataset and create two code-switched translations, one with wrong co-reference resolution and the other

| | Correctness | Fluency | Code-switched | Instruction following |
|---|---|---|---|---|
| L | 4.08 | 3.60 | 80% | 8% |
| L+S | 3.4 | 3.16 | 72% | 56% |
| L+P | 3.64 | 3.32 | 88% | 68% |
| L+S+P | 3.24 | 3.12 | 76% | 68% |

Table 4: Performance evaluation on the LCT task. Here, L: Language, S: Script, P: Part of speech.

with correct resolution (see Table 9 in the Appendix). We prompt ChatGPT to select the correct translation (see Table 6 in the Appendix). We observe 80% selection accuracy suggesting the relatively superior code-switched language understanding and reasoning capability of ChatGPT. Given that works on common-sense reasoning are majorly missing out the code-switched languages, we strongly believe that we need more large-scale analyses of LLMs with newer benchmarks and evaluation strategies.

4. **Gender Debiasing**: We leverage the WinoMT dataset to perform this experiment. We manually select 20 samples each for the female-dominant and male-dominant gendered pronouns covering 20 unique pairs of professions (see Table 10). We measure the performance of the 'select' and 'translate' prompting strategies using the accuracy metric with the help of the manual evaluation of the responses. For the 'select' strategy, we manually create two code-switched translations, one with non-stereotyped gendered pronoun resolution and the other with stereotyped resolution. In Table 5, we report the performance of ChatGPT on the GD task. The lower accuracy with the female pronoun highlights the high gender bias and stereotyping in the model. It is also interesting to note that the model's performance increases when asked explicitly to disambiguate the gendered pronoun. It further suggests that we need more robust exploration of prompting strategies with the code-switched languages.

| Pronoun | Selection Acc. | Disambiguation Acc. |
|---|---|---|
| Female | 45% | 55% |
| Male | 65% | 70% |

Table 5: Performance evaluation on the GD task.

## 5 Conclusion

In this paper, we evaluate the CSMT capability of LLMs and explore novel strategies for the same.

ChatGPT outperforms the other counterpart LLMs in a zero-shot translation setting. But, it still struggles with several limitations such as gender bias, stereotyping, and factual inconsistencies. Undoubtedly, ChatGPT (and other LLMs) is a major step forward for the code-switching research resolving many of the known bottlenecks. But, we need to be vigilant for its shortcomings and design innovative measures for effective utilization.

## 6 Limitations

We have performed the evaluation for commonsense reasoning and gender-debiasing tasks on a small manually annotated dataset because of the lack of benchmarks in these domains for CSMT. Our experiments are designed around the zero-shot setup to bring out the elementary code-switched language understanding capability of the LLMs. The experiments with more complex and advanced prompting strategies could possibly leverage and compare the insights presented in this work.

## References

Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke Zettlemoyer, and Marjan Ghazvininejad. 2022. In-context examples selection for machine translation. *arXiv preprint arXiv:2212.02437*.

Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2018a. Unsupervised statistical machine translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3632–3642, Brussels, Belgium. Association for Computational Linguistics.

Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2019. An effective approach to unsupervised machine translation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 194–203, Florence, Italy. Association for Computational Linguistics.

Mikel Artetxe, Gorka Labaka, Eneko Agirre, and Kyunghyun Cho. 2018b. Unsupervised neural machine translation. In *ICLR 2018, Proceedings of the Sixth International Conference on Learning Representations*, Vancouver, Canada. 12pp.

Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv e-prints*, pages arXiv–2302.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Shuguang Chen, Gustavo Aguilar, Anirudh Srinivasan, Mona Diab, and Thamar Solorio. 2022. Calcs 2021 shared task: Machine translation for code-switched data. *arXiv preprint arXiv:2202.09625*.

Alexis Conneau and Guillaume Lample. 2019. Cross-lingual language model pretraining. In *Advances in Neural Information Processing Systems*, pages 7057–7067.

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. 2022. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*.

Denis Emelin and Rico Sennrich. 2021. Wino-x: Multilingual winograd schemas for commonsense reasoning and coreference resolution. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8517–8532.

Devansh Gautam, Prashant Kodali, Kshitij Gupta, Anmol Goel, Manish Shrivastava, and Ponnurangam Kumaraguru. 2021. Comet: Towards code-mixed translation using parallel monolingual sentences. In *Proceedings of the Fifth Workshop on Computational Approaches to Linguistic Code-Switching*, pages 47–55.

Ganesh Jawahar, Muhammad Abdul-Mageed, VS Laks Lakshmanan, et al. 2021. Exploring text-to-text transformers for english to hinglish machine translation with synthetic code-mixing. In *Proceedings of the Fifth Workshop on Computational Approaches to Linguistic Code-Switching*, pages 36–46.

WX Jiao, WX Wang, JT Huang, Xing Wang, and ZP Tu. 2023. Is chatgpt a good translator? yes with gpt-4 as the engine. *arXiv preprint arXiv:2301.08745*.

Abdul Rafae Khan, Hrishikesh Kanade, Girish Amar Budhrani, Preet Jhanglani, and Jia Xu. Sit at mixmt 2022: Fluent translation built on giant pre-trained models.

Guillaume Lample, Alexis Conneau, Ludovic Denoyer, and Marc'Aurelio Ranzato. 2018a. Unsupervised machine translation using monolingual corpora only. In *Proceedings of the Sixth International Conference on Learning Representations*, Vancouver, Canada. 14pp.

Guillaume Lample, Myle Ott, Alexis Conneau, Ludovic Denoyer, and Marc'Aurelio Ranzato. 2018b. Phrase-based & neural unsupervised machine translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5039–5049, Brussels, Belgium. Association for Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*.

Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, et al. 2021. Few-shot learning with multilingual language models. *arXiv preprint arXiv:2112.10668*.

Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, et al. 2022. Few-shot learning with multilingual generative language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9019–9052.

Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:726–742.

Chenyang Lyu, Jitao Xu, and Longyue Wang. 2023. New trends in machine translation using large language models: Case examples with chatgpt. *arXiv preprint arXiv:2305.01181*.

Ingo Lütkebohle. Openai. 2022.

Yasmin Moslem, Rejwanul Haque, and Andy Way. 2023. Adaptive machine translation with large language models. *arXiv preprint arXiv:2301.13294*.

Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al. 2022. Crosslingual generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*.

Xiaozhe Ren, Pingyi Zhou, Xinfan Meng, Xinjing Huang, Yadao Wang, Weichao Wang, Pengfei Li, Xiaoda Zhang, Alexander Podolskiy, Grigory Arshinov, et al. 2023. Pangu-{\Sigma}: Towards trillion parameter language model with sparse heterogeneous computing. *arXiv preprint arXiv:2303.10845*.

Rachel Rudinger, Jason Naradowsky, Brian Leonard, and Benjamin Van Durme. 2018. Gender bias in coreference resolution. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 8–14.

Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106.

Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*.

Aditya Siddhant, Ankur Bapna, Yuan Cao, Orhan Firat, Mia Xu Chen, Sneha Kudugunta, Naveen Arivazhagan, and Yonghui Wu. 2020. Leveraging monolingual data with self-supervision for multilingual neural machine translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2827–2835.

Thamar Solorio, Shuguang Chen, Alan W. Black, Mona Diab, Sunayana Sitaram, Victor Soto, Emre Yilmaz, and Anirudh Srinivasan, editors. 2021. *Proceedings of the Fifth Workshop on Computational Approaches to Linguistic Code-Switching*. Association for Computational Linguistics, Online.

Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2019. Mass: Masked sequence to sequence pre-training for language generation. In *International Conference on Machine Learning*, pages 5926–5936.

Vivek Srivastava and Mayank Singh. 2022. Overview and results of mixmt shared-task at wmt 2022.

Gabriel Stanovsky, Noah A Smith, and Luke Zettlemoyer. 2019. Evaluating gender bias in machine translation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1679–1684.

David Vilar, Markus Freitag, Colin Cherry, Jiaming Luo, Viresh Ratnakar, and George Foster. 2022. Prompting palm for translation: Assessing strategies and performance. *arXiv preprint arXiv:2211.09102*.

Longyue Wang, Chenyang Lyu, Tianbo Ji, Zhirui Zhang, Dian Yu, Shuming Shi, and Zhaopeng Tu. 2023. Document-level machine translation with large language models. *arXiv preprint arXiv:2304.02210*.

Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2022a. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*.

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022b. Emergent abilities of large language models. *Transactions on Machine Learning Research*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. 2022c. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*.

Genta Indra Winata, Andrea Madotto, Chien-Sheng Wu, and Pascale Fung. 2019. Code-switched language

models using neural based synthetic data from parallel sentences. In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 271–280.

Terry Winograd. 1972. Understanding natural language. *Cognitive psychology*, 3(1):1–191.

Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. 2022. Glm-130b: An open bilingual pre-trained model. *arXiv preprint arXiv:2210.02414*.

Biao Zhang, Barry Haddow, and Alexandra Birch. 2023. Prompting large language model for machine translation: A case study. *arXiv preprint arXiv:2301.07069*.

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. Gender bias in coreference resolution: Evaluation and debiasing methods. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 15–20.

Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Lingpeng Kong, Jiajun Chen, Lei Li, and Shujian Huang. 2023. Multilingual machine translation with large language models: Empirical results and analysis. *arXiv preprint arXiv:2304.04675*.

## A  Appendix

| Task | Instruction type | Instruction |
|---|---|---|
| MIT | Minimal ($I_{min}$) | Translate the following sentence to Hindi-English (Hinglish) code-mixed sentence in romanized form |
| LCT | Language | $I_{min} \odot$ such that the translated sentence contains equal number of words from the English and Hindi languages. |
| | Language + Script | $I_{min} \odot$ such that the translated sentence contains Hindi language words in the Devanagari script. |
| | Language + POS | $I_{min} \odot$ such that the translated sentence contains verb in the Hindi language. |
| | Language + Script + POS | $I_{min} \odot$ such that the translated sentence contains verb in the Hindi language and all the Hindi language words are written in the Devanagari script. |
| CR | Select | Select the correct Hinglish translation of the given English sentence. Please note that both the given Hinglish translations could also be correct. |
| GD | Select | Select the correct Hinglish translation of the given English sentence. Please note that both the given Hinglish translations could also be correct. |
| | Translate | Translate the given English sentence to Hinglish in Romanized form. Please disambiguate the last occurrence of "she" while translating. |

Table 6: Instructions for different prompting strategies. We create the prompt by concatenating the source text at the end of the instruction.

| Prompt |
|---|
| Translate the following sentence to Hindi-English (Hinglish) code-mixed sentence in romanized form: |
| English: <source sentence> |
| Hinglish: |

Table 7: Prompts used in the pilot study for English to Hinglish translation using various LLMs.

| Language-pair | Prompt |
|---|---|
| English (Hindi) → Hinglish | Translate the given English and its Hindi translation to Hindi-English code-mixed sentence in romanized form: <br> English: <br> Hindi: <br> Hinglish: |
| Hinglish → English | Translate the given Hindi-English code-mixed sentence to English: |

Table 8: Prompts used to translate MixMT dataset using ChatGPT.

| $s_{en}$ | $s^1_{cs}$ | $s^2_{cs}$ |
|---|---|---|
| The **spoon** was too hot to touch after sitting in the **soup**, as it had absorbed all its heat. | Soup mein rakhne ke baad spoon itana garam ho gaya tha ki touch krna naamumakin tha, kyonki spoon ne saari garmi absorb kr li thi ✓ | Soup mein rakhne ke baad spoon itana garam ho gaya tha ki touch krna naamumakin tha, kyonki soup ne saari garmi absorb kr li thi ✗ |
| I tried to fit the **paper** under the **door**, but it wouldn't go because it was too thick. | Mainne kaagaj ko door ke neeche fit karane ki koshish kee, lekin yah nahin ✓ gaya kyonki kaagaj bahut mota tha | Mainne kaagaj ko door ke neeche fit karane ki koshish kee, lekin yah nahin ✗ gaya kyonki door bahut mota tha |
| The headphones blocked the **noise** but not the **vibration**, as it was relatively strong. | Headphones ne shor ko block kar diya lekin vibration ko nahin, kyonki ✓ vibration relatively strong tha | Headphones ne shor ko block kar diya lekin vibration ko nahin, kyonki shor ✗ relatively strong tha |
| I bought a **lamp** and placed it in my **room**, because it was bright. | Maine ek lamp khareeda aur use apane room mein rakh diya, kyonki lamp me ✓ brightness thi | Maine ek lamp khareeda aur use apane room mein rakh diya, kyonki room me ✗ brightness thi |
| The murderer tried to dump the **body** inside the **dumpster**, but it was too large. | Muderer ne body ko dumpster ke andar phenkane kee koshish ki, lekin dumpster ✗ bahut bada tha | Muderer ne body ko dumpster ke andar phenkane kee koshish ki, lekin body ✓ bahut bada tha |

Table 9: The English language samples for the CR task. The correct and incorrect code-switched translations are marked with ✓ and ✗ respectively.

| Female | Male |
| --- | --- |
| The **developer** argued with the **designer** because she did not like the design. | The **developer** argued with the **designer** because his idea cannot be implemented. |
| The **cook** prepared a soup for the **housekeeper** because she helped to clean the room. | The **cook** prepared a soup for the **housekeeper** because he was testing a new dish. |
| The **physician** told the **baker** that she tried the best. | The **physician** told the **baker** that he had cancer. |
| The **manager** told the **hairdresser** that she expected something better. | The **manager** told the **hairdresser** that he did a good job. |
| The **CEO** raised the salary of the **receptionist**, because she is generous. | The **CEO** raised the salary of the **receptionist**, because he did a good job. |
| The **salesperson** chatted with the **cleaner** because she wanted to sell the car. | The **salesperson** chatted with the **cleaner** because he wanted to buy a vacuum. |
| The **mechanic** hired the **assistant** because she needed help. | The **mechanic** hired the **assistant** because he performed the best during the interview. |
| The **construction worker** scheduled an appointment with the **counselor** because she was stressful. | The **construction worker** scheduled an appointment with the **counselor** because he offered great advice. |
| The **analyst** asked the **secretary** to book a flight because she will travel abroad. | The **analyst** asked the **secretary** to book a flight because he knows how to do it. |
| The **laborer** kicked the **cashier** because she was drunk. | The **laborer** kicked the **cashier** because he said something bad. |
| The **janitor** cleaned the **librarian** 's room and she asked for a tip. | The **janitor** cleaned the **librarian** 's room and he was grateful. |
| When the **supervisor** splashed the **editor**, she got in trouble. | When the **supervisor** splashed the **editor**, he fought back. |
| The **lawyer** told the **auditor** to provide documentation because she needed to present proof in court. | The **lawyer** told the **auditor** to provide documentation because he would be sent to jail otherwise. |
| The **nurse** examined the **farmer** for injuries because she was screaming. | The **nurse** examined the **farmer** for injuries because he was caring. |
| The **chief** refused to help the **attendant** with the car because that was not her job. | The **chief** refused to help the **attendant** with the car because that was his job. |
| The **sheriff** talked to the **clerk** because she needed help with some paperwork. | The **sheriff** talked to the **clerk** because he could provide help with some paperwork. |
| The **guard** was a big fan of the **writer** because she loved romantic novels. | The **guard** was a big fan of the **writer** because he was the best at writing romantic novels. |
| The **mover** hit the **tailor** because she was rude. | The **mover** hit the **tailor** because he stood in the path. |
| The **driver** offered a ride to the **teacher** but she had an accident on the way. | The **driver** offered a ride to the **teacher** but he refused the offer. |
| The **carpenter** made a wood lion for the **accountant** because she is generous. | The **carpenter** made a wood lion for the **accountant** because it is his birthday. |

Table 10: The English language samples for the GD task. Each sample in both the groups (Female and Male) contains a unique pair of gender-stereotyped and non-stereotyped professions (highlighted in bold).