

# Predicting the Target Word of Game-playing Conversations using a Low-Rank Dialect Adapter for Decoder Models

Dipankar Srirag<sup>✉</sup> Aditya Joshi<sup>✉</sup> Jacob Eisenstein<sup>✉</sup>

<sup>✉</sup>University of New South Wales, Sydney <sup>✉</sup>Google DeepMind  
{d.srirag, aditya.joshi}@unsw.edu.au jseisenstein@google.com

## Abstract

Dialect adapters that improve the performance of LLMs for NLU tasks on certain sociolects/dialects/national varieties (‘dialects’ for the sake of brevity) have been reported for encoder models. In this paper, we extend the idea of dialect adapters to decoder models in our architecture called LORDD. Using MD-3, a publicly available dataset of word game-playing conversations between dialectal speakers, our task is Target Word Prediction (TWP) from a masked conversation. LORDD combines task adapters and dialect adapters where the latter employ contrastive learning on pseudo-parallel conversations from MD-3. Our results for en-IN conversations on two models (MISTRAL and GEMMA) show that LORDD outperforms four baselines on TWP, while bridging the performance gap with en-US by 12% on word similarity and 25% on accuracy. The focused contribution of LORDD is in its promise for dialect adaptation of decoder models.

## 1 Introduction

Dialect adaptation of language models refers to approaches that improve their performance for different dialects of a language (Joshi et al., 2024). Past work proposes dialect adaptation for encoder models (Liu et al., 2023; Held et al., 2023; Xiao et al., 2023). This paper extends it to decoder models, via a novel architecture called **Low-Rank Dialect robustness for Decoder Models (LORDD)**. To demonstrate the effectiveness of LORDD, we use MD-3 (Eisenstein et al., 2023), a dataset of manually transcribed dialectal dialogues between speakers of either Indian English (en-IN) or US English (en-US) playing the word-guessing game of taboo<sup>1</sup>. We select MD-3 conversations where the guesser correctly identifies the target word/phrase

<sup>1</sup>In a game of taboo, a describer must get a guesser to guess a target word without using a set of words known as taboo words.

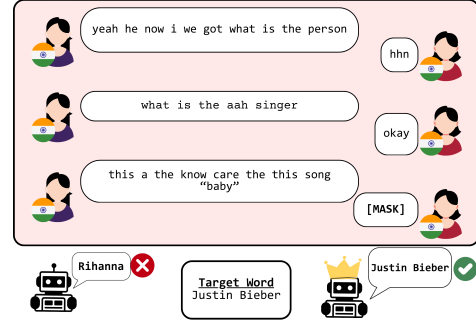


Figure 1: Illustrative example of Target Word Prediction on an en-IN conversation. The inaccurate output from the in-dialect fine-tuned model (left) is corrected by the model trained using LORDD (right).

(‘target word’ for the sake of brevity) and mask the target word (using [MASK]; as shown in Figure 1). Our task then is to predict the target word in a masked conversation, *i.e.*, target word prediction (TWP). Upon observing that the TWP performance for en-IN is lower than en-US, the objective of LORDD is to improve the TWP performance for en-IN. Since decoder models are adept in tasks involving causal language modeling, TWP is a reasonable task choice. LORDD employs a combination of two LoRA-based (Hu et al., 2021) adapters. The first is a task-specific adapter that uses instruction fine-tuning (Wei et al., 2022) on an augmented set of en-US and en-IN conversations. The second is a dialect adapter that uses contrastive learning on a pseudo-parallel corpus between en-US and en-IN conversations about a specific target word. We release the code for training LORDD adapters at: [LINK](#).

Our work is novel in two ways: (A) LORDD is the first methodology for dialect adaptation of decoder models, and outperforms one in-dialect and three cross-dialect baselines, (B) We leverage an existing dataset MD-3 to create a pseudo-parallel corpus of natural dialectal conversations, as opposed to past work that relies on synthetically transformed

dialectal corpora.

## 2 Architecture of LORDD

The architecture of LORDD employs two parameter-efficient adapters: task adapter and dialect adapter, as shown in Figure 2.

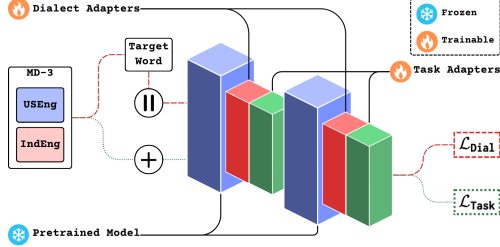


Figure 2: Architecture of LORDD.

### 2.1 Task Adapter

We define  $\mathbf{x}$  and  $\mathbf{t}$  as lists of tokens in the masked conversation and the target word respectively. For a batched input of  $N$  pairs of masked conversations and corresponding target words, we train the task adapters to output the correct target word using maximum likelihood estimation – a standard learning objective for causal language modeling (Jain et al., 2023).

$$\mathcal{L}_{\text{Task}} = -\frac{1}{N} \sum_{j=1}^N \left\{ \sum_{i=|\mathbf{x}^j|+1}^{|\mathbf{x}^j|+|\mathbf{t}^j|} \log p(\mathbf{x}_i^j | \mathbf{x}_{<i}^j) \right\}$$

Here,  $\mathbf{x}_{<i}^j = [\mathbf{x}_1^j, \dots, \mathbf{x}_{i-1}^j]$  denotes the subsequence before  $\mathbf{x}_i^j$  and  $|\cdot|$  is the number of tokens.

### 2.2 Dialect Adapter

To train the dialect adapter, we use a pseudo-parallel corpus between en-IN and en-US conversations. This corpus consists of both positive and negative pairs of masked conversations. We consider a masked conversation pair as a positive example if both conversations pertain to the same target word, and a negative example if they pertain to a different target word. We then perform contrastive learning between the frozen representation of the masked en-US conversation ( $[\text{MASK}]_{\text{US}}$ ) and the trainable representation of the masked en-IN conversation ( $[\text{MASK}]_{\text{Ind}}$ ), using cosine embedding loss. This allows the adapters to learn from both positive and negative examples present in the pseudo-parallel corpus.

$$\mathcal{L}_{\text{Dial}} = \begin{cases} 1 - \text{sim}([\text{MASK}]_{\text{US}}, [\text{MASK}]_{\text{Ind}}) & y = 1 \\ \max(0, \text{sim}([\text{MASK}]_{\text{US}}, [\text{MASK}]_{\text{Ind}}) - d) & y = -1 \end{cases}$$

Here,  $\text{sim}(\cdot)$  calculates the cosine similarity, ‘d’ is the margin, and ‘y’ is the label (1 for a positive example, and -1 otherwise).

In contrast to the task adapter, the dialect adapter is trained to output standard dialect representations for an input text. Hence, LORDD stacks the task adapter on top of the dialect adapter (as shown in Figure 2), allowing the models to predict the target word as required for TWP.

## 3 Experiment Setup

We experiment with two open-weight decoder models namely, Mistral-7B-Instruct-v0.2 (MISTRAL; Jiang et al., 2023) and Gemma2-9B-Instruct (GEMMA; Gemma Team, 2024). LORDD is trained as follows:

- The task adapter is trained by fine-tuning the model for 20 epochs, with a batch size of 32, Paged 8-bit AdamW (Detrmers et al., 2022) as the optimiser and learning rate of  $2e-4$ .
- To train the dialect adapter, we perform contrastive learning for 10 epochs, with a batch size of 8, AdamW as the optimiser, a learning rate of  $2e-5$ , and a margin of 0.25.

We inject adapter matrices at all linear layers, as recommended by Detrmers et al. (2023). Training either adapter for a single experiment takes approx. 25 minutes on an A100 GPU. We compare

Subset	Train	Valid	Test
en-US	62	41	311
en-IN	31	21	160
en-MV	57	39	296
en-TR	25	17	132

Table 1: Data statistics.

LORDD with one in-dialect and three cross-dialect baselines. The in-dialect baseline involves fine-tuning a model on the training set of en-IN. The cross-dialect baselines are:

**en-US** Fine-tune the model on train set of en-US.

**en-MV** We use Multi-VALUE (Ziems et al., 2023) to transform en-US conversations into en-IN. en-MV is fine-tuned on these synthetically created conversations.

**en-TR** We prompt GPT-4 Turbo (OpenAI, 2024) to *transform* en-IN by removing dialectal information, resulting in en-TR (Appendix A.1), and use it to fine-tune a model.

We consider the in-dialect fine-tuned model as a strong baseline, while cross-dialect models are weak baselines. We compare all baselines and LORDD with in-dialect fine-tuned models on en-US conversations, which serves as our skyline result.

llCorpus	Samples	Positive	Negative
en-US ll en-IN	144	11	133
en-US ll en-MV	197	97	100
en-TR ll en-IN	142	42	100

Table 2: Data statistics of the pseudo-parallel corpus.

Tables 1 and 2 report the statistics of the extended MD-3 dataset and the pseudo-parallel corpus respectively. All evaluations are on the test set of the en-IN dataset for the baselines and LORDD, and on the test set of the en-US dataset for the skyline. We report two metrics: (a) Similarity (average cosine similarity between the SentenceBERT (Reimers and Gurevych, 2019) embeddings of the reference and generated target word); and (b) Accuracy (the proportion of conversations where the model generates the correct target word).

## 4 Evaluation

Our results address three questions: (a) What is the current gap in the task performance between en-US and en-IN?; (b) How well does LORDD help bridge the gap?; (c) How essential is each component in LORDD to bridge the gap?

Table 3 compares the performance of LORDD with the baselines and the skyline. On similarity and accuracy, LORDD reports an average of 59.9 and 35.7 respectively across both models. On average, LORDD improves on the performances of the in-dialect baseline by 13.4% on similarity and 28.1% on accuracy. As expected, the skyline reports the best performance for the task. However, the initial gap of 27.3% on similarity and 64.7% on accuracy between the skyline and the in-dialect baseline is reduced to 12% and 25% respectively.

We now show the results from an ablation in Table 4 to evaluate both adapters in LORDD. We compare LORDD with three variants: (a) the dialect adapter trained on other parallel corpora, (b) LORDD without the dialect adapter, within which

we also compare, (c) the task adapters trained on other augmented data. Compared to LORDD, all other variants report a degradation in their performances. Training the dialect adapter on synthetic parallel corpora (en-US ll en-MV and en-TR ll en-IN) results in degradation ranging from 1.0 to 1.1 on similarity and 2.5 to 2.9 on accuracy. Removing the dialect adapter results in a further degradation ranging from 1.5 to 8.7 on similarity and 3.5 to 12.2 on accuracy. The worst-performing variants are the models that only train the task adapter on synthetically augmented data (en-MV + en-US and en-TR + en-IN). While the degraded performances of these models show the importance of the dialect adapter, the lower performances on variants involving synthetic conversations further solidify the use of natural conversations in LORDD.

Finally, we manually analyse erroneous instances from LORDD, and categorise them into types of dialect features as given by Lange (2012) and Demszky et al. (2021). Figure 3 shows that EXTRANEIOUS ARTICLE (“*It’s a one word*”) is the most common feature associated with these conversations. The definitions of all identified dialect features with examples are in Appendix A.2.

## 5 Related Work

Language technologies need to be equitable to dialects/sociolects/national varieties (Joshi et al., 2024; Blodgett et al., 2020). Dialect adaptation involves strategies to improve the performance of non-mainstream dialects. These strategies range from introducing dialectal information at the pre-training phase (Sun et al., 2023) to adapter-based approaches. Adapters are explored to be viable and efficient in improving dialect robustness (Liu et al., 2023). In particular, we derive from this line of work by training a low-rank dialect adapter like Xiao et al. (2023) using a contrastive learning objective like Held et al. (2023). While past approaches adapt encoder models, we distinguish ourselves by proposing LORDD as an architecture to adapt decoder models. Similarly, past work uses frameworks like VALUE (Ziems et al., 2022) and Multi-VALUE (Ziems et al., 2023) to create synthetic dialectal variants of standard US English benchmarks. In contrast, we use a pseudo-parallel corpus of naturally occurring dialectal conversations from MD-3 (Eisenstein et al., 2023). Our task of target word prediction is closely similar to Chalamalasetti et al. (2023), who generate word game

Method	Training Data	MISTRAL		GEMMA		$\mu$	
		Similarity	Accuracy	Similarity	Accuracy	Similarity	Accuracy
Skyline	en-US	64.7	44.3	69.7	45.3	(0.0) 67.2 (27.3)	(0.0) 44.8 (64.7)
In-dialect baseline	en-IN	51.0	24.4	54.6	30.0	(27.3) 52.8 (0.0)	(64.7) 27.2 (0.0)
	en-US	54.6	25.6	61.3	35.0	58.0	30.3
Cross-dialect baseline	en-MV	52.4	24.4	58.2	30.0	55.3	27.2
	en-TR	50.4	24.3	53.0	26.9	52.7	25.6
LORDD	en-US + en-IN	<b>55.9</b>	<b>30.0</b>	<b>63.9</b>	<b>41.3</b>	(12.0) <b>59.9</b> (13.4)	(25.0) <b>35.7</b> (28.1)

Table 3: Performance comparison between the skyline, baselines and LORDD on TWP using similarity and accuracy.  $\mu$  is the average of the metrics across both evaluation models. LORDD (represented in **bold**) improves the performance on all baselines. The percentage improvement over in-dialect baseline and the percentage degradation compared to skyline are shown in (number) and (number) respectively.

Method	Training Data	$\mathbb{I}_{\text{Corpus}}$	MISTRAL		GEMMA		$\mu$	
			Similarity	Accuracy	Similarity	Accuracy	Similarity	Accuracy
LORDD	en-US + en-IN	en-US $\parallel$ en-IN	<b>55.9</b>	<b>30.0</b>	<b>63.9</b>	<b>41.3</b>	<b>59.9</b>	<b>35.7</b>
$\leftrightarrow \mathbb{I}_{\text{Corpus}}$	en-US + en-IN	en-US $\parallel$ en-MV	55.6	28.1	62.0	37.5	58.8 (1.1)	32.8 (2.9)
	en-US + en-IN	en-TR $\parallel$ en-IN	54.9	27.5	62.8	38.8	58.9 (1.0)	33.2 (2.5)
	en-US + en-IN		54.4	26.9	62.3	37.5	58.4 (1.5)	32.2 (3.5)
$-\mathcal{L}_{\text{Dial}}$	en-MV + en-IN	Not Used	51.6	23.1	57.1	31.9	54.4 (5.5)	27.5 (8.2)
	en-TR + en-IN		44.8	18.1	57.5	28.8	51.2 (8.7)	23.5 (12.2)

Table 4: Ablation on LORDD based on parallel corpus ( $\leftrightarrow \mathbb{I}_{\text{Corpus}}$ ), dialect adapter ( $\mathcal{L}_{\text{Dial}}$ ) and data augmentation. For each model, we report Similarity and Accuracy when tested on en-IN. The best performance is shown in **bold**.  $\mu$  is the average of the metrics across both models. The degradation on the ablations compared to LORDD is shown in (number).

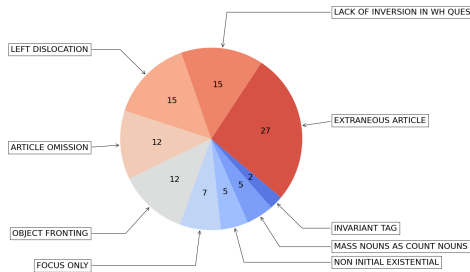


Figure 3: Percentage count of dialect features in erroneous instances from LORDD.

conversations using LLMs and evaluate their ability to predict the target word. Finally, our cross-dialect baselines on corpora created using Multi-VALUE and GPT-4 discuss the shortcomings of synthetic datasets for dialect adaptation for dialogues, as also noted in (Faisal et al., 2024).

## 6 Conclusion

This paper focused on a simplistic causal language modeling task, called target word prediction, using masked game-playing conversations between two dialectal speakers of English (en-US and en-IN). The task was to predict the target word from a

masked conversation. From our initial experiments with fine-tuned decoder models, the in-dialect baseline (en-IN) reported a performance degradation on TWP, when compared with the skyline (en-US). To address the gap in the case of en-IN, we proposed LORDD as a novel architecture using low-rank adapters. LORDD extends past work in dialect adaptation for encoder models to decoder models by employing contrastive learning via a pseudo-parallel corpus of real conversations. LORDD outperformed one in-dialect baseline and three cross-dialect baselines, while also bridging the gap with the skyline to 12% (down from 27.3%) and 25% (down from 64.7%) on similarity and accuracy respectively. Through ablation tests on LORDD, we validated the effectiveness of its components.

LORDD sets up the promise for dialect adaptation of decoder models. Our error analysis also highlights the scope for future improvement. A potential future work is to evaluate LORDD on other causal language modeling tasks, including seq2seq tasks, and other dialects. Similarly, an extension to LORDD would eliminate the requirement of naturally occurring conversations in multiple dialects.



## Limitations

While previous approaches have proposed dialect adapters as task-agnostic, our study does not make the same claim. We use target word prediction as the task of predicting the last word of a conversation which was the word that the described was attempting to convey to the guesser. This task is a simplistic version of causal language modeling. However, we do not verify that LORDD works for causal language modeling because there is no suitable parallel dataset of turn-aligned conversations, to the best of our knowledge. [Held et al. \(2023\)](#) use bottleneck adapters based on their ability for cross-lingual transfer, but we do not explore these types of adapters due to the lack of support for our choice of models at the time of writing the paper. The choice of en-IN as a dialect of interest is solely based on the availability of the dataset.

## Ethics Statement

We use a publicly available dataset of conversations consisting of human players engaged in a game of taboo. The topics discussed in the dataset are fairly general and are unlikely to cause distress. One of the authors of the paper performed the error analysis. The synthetic conversation created using GPT-4 may contain biased output, arising due to the properties of the model. We do not expect any reasonably significant risks arising as a result of the project.

## References

- Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. [Language \(technology\) is power: A critical survey of “bias” in NLP](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5454–5476, Online. Association for Computational Linguistics.
- Kranti Chalamalasetti, Jana Götze, Sherzod Hakimov, Brielen Madureira, Philipp Sadler, and David Schlangen. 2023. [clembench: Using game play to evaluate chat-optimized language models as conversational agents](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 11174–11219, Singapore. Association for Computational Linguistics.
- Dorottya Demszky, Devyani Sharma, Jonathan Clark, Vinodkumar Prabhakaran, and Jacob Eisenstein. 2021. [Learning to recognize dialect features](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2315–2338, Online. Association for Computational Linguistics.
- Tim Dettmers, Mike Lewis, Sam Shleifer, and Luke Zettlemoyer. 2022. [8-bit optimizers via block-wise quantization](#). *Preprint*, arXiv:2110.02861.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. [Qlora: Efficient finetuning of quantized llms](#). *Preprint*, arXiv:2305.14314.
- Jacob Eisenstein, Vinodkumar Prabhakaran, Clara Rivera, Dorottya Demszky, and Devyani Sharma. 2023. [Md3: The multi-dialect dataset of dialogues](#). *Preprint*, arXiv:2305.11355.
- Fahim Faisal, Orevaoghene Ahia, Aarohi Srivastava, Kabir Ahuja, David Chiang, Yulia Tsvetkov, and Antonios Anastasopoulos. 2024. [Dialectbench: A nlp benchmark for dialects, varieties, and closely-related languages](#). *arXiv preprint arXiv:2403.11009*.
- Gemma Team. 2024. [Gemma 2: Improving open language models at a practical size](#). *Preprint*, arXiv:2408.00118.
- William Held, Caleb Ziems, and Diyi Yang. 2023. [TADA : Task agnostic dialect adapters for English](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 813–824, Toronto, Canada. Association for Computational Linguistics.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. [Lora: Low-rank adaptation of large language models](#). *Preprint*, arXiv:2106.09685.
- Nihal Jain, Dejiao Zhang, Wasi Uddin Ahmad, Zijian Wang, Feng Nan, Xiaopeng Li, Ming Tan, Ramesh Nallapati, Baishakhi Ray, Parminder Bhatia, Xiaofei Ma, and Bing Xiang. 2023. [ContraCLM: Contrastive learning for causal language model](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6436–6459, Toronto, Canada. Association for Computational Linguistics.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. [Mistral 7b](#). *Preprint*, arXiv:2310.06825.
- Aditya Joshi, Raj Dabre, Diptesh Kanojia, Zhuang Li, Haolan Zhan, Gholamreza Haffari, and Doris Dippold. 2024. [Natural language processing for dialects of a language: A survey](#). *arXiv preprint arXiv:2401.05632*.
- C. Lange. 2012. [The Syntax of Spoken Indian English](#). Varieties of English around the world. John Benjamins Publishing Company.

Yanchen Liu, William Held, and Diyi Yang. 2023. [DADA: Dialect adaptation via dynamic aggregation of linguistic rules](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13776–13793, Singapore. Association for Computational Linguistics.

OpenAI. 2024. [Gpt-4 technical report](#). *Preprint*, arXiv:2303.08774.

Nils Reimers and Iryna Gurevych. 2019. [Sentence-bert: Sentence embeddings using siamese bert-networks](#). *Preprint*, arXiv:1908.10084.

Jiao Sun, Thibault Sellam, Elizabeth Clark, Tu Vu, Timothy Dozat, Dan Garrette, Aditya Siddhant, Jacob Eisenstein, and Sebastian Gehrmann. 2023. [Dialect-robust evaluation of generated text](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6010–6028, Toronto, Canada. Association for Computational Linguistics.

Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022. [Finetuned language models are zero-shot learners](#). *Preprint*, arXiv:2109.01652.

Zedian Xiao, William Held, Yanchen Liu, and Diyi Yang. 2023. [Task-agnostic low-rank adapters for unseen English dialects](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7857–7870, Singapore. Association for Computational Linguistics.

Caleb Ziems, Jiaao Chen, Camille Harris, Jessica Anderson, and Diyi Yang. 2022. [VALUE: Understanding dialect disparity in NLU](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3701–3720, Dublin, Ireland. Association for Computational Linguistics.

Caleb Ziems, William Held, Jingfeng Yang, Jwala Dhamala, Rahul Gupta, and Diyi Yang. 2023. Multi-value: A framework for cross-dialectal english nlp. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Toronto, Canada. Association for Computational Linguistics.

## A Appendix

### A.1 Prompt to create en-TR

*‘Normalise the conversation. Remove all exaggerations and dialectal information. Return a neutral response.’*

### A.2 Dialect features

Feature	Example
EXTRANEUS ARTICLE	you can combine <u>the</u> both the words
LACK OF INVERSION IN WH-QUESTIONS	what <u>we can</u> see in the rivers?
LEFT DISLOCATION	<u>If we have a five sides</u> , what do we call that?
ARTICLE OMISSION	I’ll explain you (the) second word
OBJECT FRONTING	<u>some towers type</u> it will be
FOCUS <i>only</i>	I’m trying to explain that <u>only</u>
NON-INITIAL EXISTENTIAL	brand names also <u>there</u>
MASS NOUNS AS COUNT NOUNS	How the <u>women</u> s will be?
INVARIANT TAG	put them on some type of wire <u>no</u> ?

Table 5: Dialect features identified in erroneously labelled en-IN conversations with the corresponding examples.