

# MuDRiC: Multi-Dialect Reasoning for Arabic Commonsense Validation

Anonymous ACL submission

## Abstract

Commonsense validation evaluates whether a sentence aligns with everyday human understanding, a critical capability for developing robust natural language understanding systems. While substantial progress has been made in English, the task remains underexplored in Arabic, particularly given its rich linguistic diversity. Existing Arabic resources have primarily focused on Modern Standard Arabic (MSA), leaving regional dialects underrepresented despite their prevalence in spoken contexts. To bridge this gap, we present two key contributions. We introduce MuDRiC, an extended Arabic commonsense dataset incorporating multiple dialects. To the best of our knowledge, this is the first Arabic multi-dialect commonsense reasoning dataset. We further propose a novel method adapting Graph Convolutional Networks (GCNs) to Arabic commonsense reasoning, which enhances semantic relationship modeling for improved commonsense validation. Our experimental results demonstrate that this approach consistently outperforms the baseline of direct language model fine-tuning. Overall, our work enhances Arabic natural language understanding by providing a foundational dataset and a new method for handling its complex variations. Data and code are available at this [link](#).

## 1 Introduction

Common sense reasoning is a fundamental task in natural language processing (NLP), enabling machines to interpret and generate text in ways that align with human intuition (Sap et al., 2020). It is critical for AI systems to make plausible inferences about the world and engage in human-like conversation (Davis and Marcus, 2015). However, conversational commonsense often involves implicit social norms, cultural references (Sadallah et al., 2025), and pragmatic reasoning that vary across dialects. Despite progress in English (Levesque et al., 2012; Sap et al., 2019a; Talmor et al., 2019)

and other high-resource languages, common sense reasoning remains challenging for languages with dialectal diversity, such as Arabic, primarily due to severe scarcity of annotated data for dialects.

Most existing Arabic common sense benchmarks focus exclusively on Modern Standard Arabic (MSA), neglecting the rich diversity of Arabic dialects (Lamsiyah et al., 2025; Sadallah et al., 2025; Khaled et al., 2023). Dialects such as Egyptian, Gulf, Levantine, and Moroccan dominate everyday communication across the Arab world. Beyond lexical or grammatical variation, these dialects encode fine-grained regional cultural knowledge, making dialectal commonsense reasoning a culturally grounded and challenging task. As a result, models trained solely on MSA often fail to generalize to dialectal content. To address this gap, we introduce the first multi-dialect Arabic commonsense dataset balanced across Egyptian, Gulf, Levantine, and Moroccan dialects.

In terms of approaches to address Arabic common sense tasks, prior work heavily relies on MSA-centered models, e.g., AraBERT (Antoun et al., 2020), which perform barely above chance on dialectal data (Lamsiyah et al., 2025; Khaled et al., 2023). Dialect-specific models such as MARBERT (Abdul-Mageed et al., 2021) also show weak performance due to nuanced differences between dialects. More related work in Appendix A.

We propose integrating base language models with graph-based augmentation to capture deeper semantic relationships. This integration of graph-based methods significantly enhances cross-dialect robustness. To summarize our main contributions:

- We introduce MuDRiC: the first multi-dialect Arabic common sense benchmark, enabling more inclusive and robust Arabic NLP systems.
- We propose graph-based augmentation training strategy to enhance performance on dialectal data.

## 2 Dataset

**Task Description and Formulation** Given a single sentence, the task aims to identify whether it is reasonable (labeled as 1) or non-reasonable (labeled as 0), based on its alignment with common sense. We cast commonsense validation as a binary classification task to provide a unified and simple formulation across datasets with different original structures. This setting allows us to evaluate the commonsense plausibility of a single sentence independently, rather than relying on relative comparisons between candidates and facilitates scaling the task to multiple Arabic dialects.

**MSA** We use two established Modern Standard Arabic (MSA) datasets for commonsense validation: the Arabic Dataset for Commonsense Validation (ADCV; Tawalbeh and Al-Smadi (2020)) and ArabicSense (Lamsiyah et al., 2025). ADCV contains 11,000 instances, each consisting of a pair of sentences, where the task is to select the more commonsensical option. We convert this setup into a binary classification task by separating each sentence pair into two individual sentences, assigning label 1 to the original correct (reasonable) sentence and label 0 to its incorrect counterpart. This results in 22,000 labeled samples. ArabicSense includes 5,650 multiple-choice instances with two candidate sentences per instance, one of which is commonsensical. We apply the same conversion strategy, assigning labels accordingly, yielding 11,288 MSA samples after removing duplication.

**Dialects Extension** Based on the MSA datasets above, we translate them into four Arabic dialects including Egyptian, Moroccan, Gulf and Levantine using GPT-4o (OpenAI, 2024). The statistical distribution of the extended datasets is summarized in Table 1, while Table 2 presents sample MSA sentences from ADCV alongside their corresponding dialectal translations. Prompting details can be found in Appendix B.

The final composite dataset ensures a parallel representation across four major Arabic dialect families. This addresses a critical gap in Arabic NLP, where previous benchmarks have been limited to either Modern Standard Arabic or isolated dialectal efforts without systematic comparison.

**Quality Control** We apply multiple quality control steps to ensure the reliability of the translated dataset. First, each dialectal translation is automat-

Source Dataset	MSA Samples	Dialectal Samples
ADCV	22,000	88,000
ArabicSense	11,288	45,152
<b>Total</b>	<b>33,288</b>	<b>133,152</b>

Table 1: Statistical distribution of datasets, with each MSA extending to four dialects.

ically verified using Gemini 2.5 Flash by jointly evaluating the translated sentence and its original MSA counterpart (Prompt in Appendix B). Samples flagged as incorrect by Gemini are subsequently reviewed by native-speaker annotators. In total, approximately 8,580 samples (5.2% of the dataset) were flagged; among these, 27% were confirmed to be incorrect (corresponding to roughly 1.4% of the full dataset) and were corrected by the annotators.

As an additional validation step to the original source datasets, we randomly sample 500 instances and ask two independent annotators to verify the correctness of their commonsense labels (reasonable vs. non-reasonable). All sampled instances were found to be correctly labeled.

## 3 Methodology

We explore (i) graph-based augmentation to inject relational structure and (ii) domain-adversarial training (Appendix D) to encourage dialect-invariant representations.

### 3.1 Graph-based Language Model Reps

Inspired by prior work that integrates graph encoders with Transformer models (Jiawei et al., 2020; Zhibin et al., 2020), we augment pretrained Masked Language Models (MLMs) with a Graph Convolutional Network (GCN) that encodes local word-level relations and surface morphological cues. This is motivated by the continuum of Arabic dialects and their non-standard orthography, which introduce substantial spelling and morphological variation. As a result, sequence-based fine-tuning becomes brittle, and dialect-invariant objectives are less reliable (Sha’ban and Habash, 2025; Bhatia et al., 2025). The graph encoder connects related variants via message passing, complementing contextual semantics and improving cross-dialect commonsense validation.

**Pipeline Overview** We first construct a semantic graph for each input instance to capture relational

MSA Text	Egyptian	Gulf	Moroccan	Levantine	Label
لا أحد يريد العيش مع الفئران (No one wants to live with rats)	محدث عايز يعيش مع الفئران	ما في أحد بيبي يعيش مع الفئران	حتى واحد ما بغا يعيش مع الفئران	ما حدا بده يعيش مع الفئران	1
تقوم جورجيا تك بتدريب التنين (Georgia Tech trains dragons)	جورجيا تك بتدرب التنين	جورجيا تك تقوم بتدريب التنين	جورجيا تك كيدزبو التنين	جورجيا تك عم تدرب التنين	0

Table 2: MSA and dialectal examples from ADCV. Label 1 = reasonable and label 0 = non-reasonable sentence.

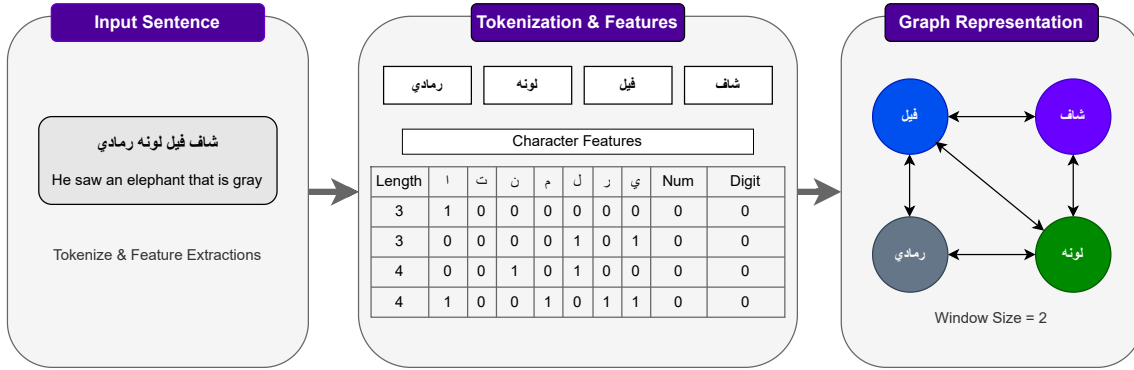


Figure 1: The creation process of the graph representation for a sentence.

dependencies between entities. This graph is processed using GCNs to produce a fixed-length graph representation. In parallel, a BERT-based encoder generates a contextualized text embedding, which is projected to the same dimensionality as the graph representation. The two embeddings are then fused via concatenation. This fusion enriches the model by combining token-level semantic representations with higher-order structural information, enabling more effective reasoning over the input.

**Graph Representation** To build the graph representation (Figure 1), each input text is first tokenized into words. A co-occurrence graph is then constructed, where nodes correspond to unique words and undirected edges connect words appearing within a fixed-size sliding window. Each node is initialized with a handcrafted feature vector derived from word-level statistics, including word length and Arabic-specific morphological indicators (e.g., character counts and digit-related features). This results in lightweight and informative node representations that reflect the surface morphology and character patterns in the Arabic text.

The resulting graph is then processed by a multi-layer GCN with one hidden layer. The GCN layers propagate and aggregate features across the graph, enabling the model to learn contextual structural patterns. A global mean-pooling layer is then applied to extract a single fixed-length vector that summarizes the entire graph.

**Semantic Representation by MLMs** In parallel, the input text is encoded using a BERT-based model. We extract the contextual embedding of the “[CLS]” token from the final hidden state layer, which serves as a summary representation of the input sequence. Both the graph and the BERT embeddings are projected into a shared fusion space using learnable linear projections.

**Fusion** To combine the two representation inputs, we employ a multi-head self-attention mechanism over the concatenated graph and MLM embeddings. This allows the model to dynamically weigh the contribution of each modality and to learn complex interactions between them. The output of the attention layer is flattened and passed through a feedforward classification head. Figure 2 and Algorithm 1 summarize the pipeline in Appendix C.

## 4 Experiments

### 4.1 Experimental Setup

**Baselines** We fine-tune three Arabic-centric pre-trained masked language models (MLMs), i.e., CAMELBERT-mix (Inoue et al., 2021), AraBERT and MARBERT as baselines, and evaluated on all dialects datasets. The three models were pretrained on large Arabic corpus, making them suitable for the task at hand. AraBERT focuses on Modern Standard Arabic. MARBERT emphasizes dialectal Arabic, incorporating substantial representation from various regional dialects, which can better

Methods	MSA	Egyptian	Gulf	Levantine	Moroccan	Avg.
AraBERTv2	65.88	68.33	65.78	65.36	62.09	65.34
AraBERTv2 + GCN	<b>68.15</b>	<b>70.00</b>	<b>66.63</b>	<b>67.03</b>	<b>63.24</b>	<b>67.01</b>
CAMeLBERT-mix	73.30	73.52	74.63	73.82	66.45	72.34
CAMeLBERT-mix + GCN	<b>74.21</b>	<b>75.61</b>	<b>76.24</b>	<b>76.70</b>	<b>68.87</b>	<b>74.28</b>
MARBERTv2	80.12	80.73	78.81	80.03	71.09	78.16
MARBERTv2 + GCN	<b>81.64</b>	<b>81.26</b>	<b>80.87</b>	<b>80.64</b>	<b>73.06</b>	<b>79.53</b>

Table 3: Accuracy of baselines and our method on MSA and dialects datasets in % (higher is better).

capture the dialectal characteristics of the dataset. Similar to MARBERT, CAMeLBERT-mix was pre-trained to dialectal Arabic in addition to modern standard Arabic and classical Arabic.

**Our Methods** We evaluated the effectiveness of graph-based representations of MLMs which fused GCN-based embeddings with masked language models’ contextual embeddings. This experimental setup allowed us to examine the hypothesis that graph-enhanced representations can improve downstream task performance.

**Data Split and Training Setups** In all experiments, we trained models using cross-entropy loss and optimized with AdamW (Loshchilov and Hutter, 2017), using a learning rate of  $2e-5$ , weight decay of 0.01, and a batch size of 128 for 3 epochs. We used 110K MuDRiC samples with ADCV as source dataset, balanced across both common-sense and dialect labels. The data were split into 70%/15%/15% train, development, and test sets. This split was fixed and reused across all experiments to ensure fair and consistent comparisons.

## 4.2 Results

Table 3 reports accuracy on MSA and four dialectal subsets. The *Avg. Dialects* column corresponds to the overall average across all five subsets (MSA + the four dialects).

**Which is the Best Base Model?** Across the table, performance follows a stable ranking: *MARBERTv2* > *CAMeLBERT-mix* > *AraBERTv2*. This ordering is consistent with the expected degree of *dialect exposure* during pretraining: MARBERTv2 is explicitly dialect-heavy, CAMeLBERT-mix is more balanced, and AraBERTv2 is comparatively more MSA-oriented. The implication is that dialectal generalization is primarily constrained by representation quality learned during pretraining; downstream methods help, but they do not compensate fully for a strong pretraining mismatch.

## GCN Consistently Improves Performance.

The consistent improvements across all models suggest that the GCN provides information that the transformer alone does not reliably capture, such as relational structure, lexical/semantic neighborhood effects, or instance-to-instance dependencies. Importantly, the gains appear larger (in relative terms) for the weaker backbones (AraBERTv2 and CAMeLBERT-mix), which supports the interpretation that GCN fusion helps *mitigate* dialect/domain mismatch by encouraging useful sharing across related samples or features. For MARBERTv2, the improvement is still present but more incremental, consistent with the idea that a dialect-rich backbone already captures much of the needed variation and the GCN mainly refines decision boundaries.

Overall, the table supports a clear conclusion: dialect-aware pretraining is the strongest driver of performance, and GCN fusion is a reliable enhancement.

## 5 Conclusion

This work presented two major contributions to Arabic NLP: (1) the creation of the first large-scale, multi-dialect common sense reasoning dataset, and (2) Enhanced Arabic Commonsense Reasoning methodology combining graph-based embeddings with pre-trained BERT-based models to enhance performance. By systematically expanding MSA commonsense reasoning benchmarks into four major dialects we established a crucial resource for evaluating dialect robustness. Our experiments demonstrated that neither MSA-focused models (e.g., AraBERT) nor dialect-pretrained models (e.g., MARBERT) alone suffice for reliable common sense classification across dialects. Instead, our hybrid graph-based approach to structured commonsense representation outperformed prior methods, setting a new benchmark for dialect-aware Arabic NLP.

## 310 Limitations & Future Work

311 While our work advances dialect-aware common  
312 sense reasoning, several limitations warrant discus-  
313 sion: the dialectal data generation process relied on  
314 GPT-4o for translation, which may introduce sub-  
315 tle semantic shifts or stylistic inconsistencies com-  
316 pared to naturally occurring dialectal speech, and  
317 while we implemented quality checks, the absence  
318 of large-scale human validation leaves room for po-  
319 tential noise, particularly in idiomatic expressions  
320 requiring deep cultural familiarity; the framework  
321 treats all dialects as equally distinct from MSA,  
322 overlooking gradient dialectal relationships. For in-  
323 stance, Levantine Arabic shares more lexical over-  
324 lap with MSA than Moroccan Arabic, potentially  
325 leading to uneven generalization where linguisti-  
326 cally closer dialects benefit implicitly; the binary  
327 labeling scheme (reasonable vs. non-reasonable)  
328 oversimplifies the continuum of common sense  
329 plausibility, failing to capture partially valid or  
330 context-dependent interpretations; moreover, the  
331 focus on four major dialects excludes dozens of  
332 other Arabic varieties, risking the marginalization  
333 of less common dialects like Sudanese or Yemeni  
334 Arabic, an area future work should address.

335 Future work will prioritize reducing the per-  
336 sistent Moroccan gap via dialect-specific adapta-  
337 tion (e.g., continued pretraining, normalization,  
338 or dialect-aware modules). We also plan to ex-  
339 pand coverage to additional underrepresented var-  
340 ieties (e.g., Sudanese, Yemeni, Algerian, Iraqi),  
341 test broader diglossic/code-switched settings (e.g.,  
342 Moroccan Arabic–French), and extend evaluation  
343 beyond sentence-level to contextual or multi-turn  
344 commonsense reasoning.

## 345 Ethical Statement

346 **Data License** A primary ethical consideration in  
347 our work is the licensing and provenance of the data  
348 used. Our dataset builds upon two publicly avail-  
349 able resources: the Arabic Dataset for Common-  
350 sense Validation and ArabicSense, both of which  
351 have been released for research purposes with ap-  
352 propriate usage permissions. To ensure compliance  
353 with licensing constraints, we generated novel di-  
354 alectal variants derived from the Modern Standard  
355 Arabic (MSA) instances provided in the original  
356 datasets. This approach ensures that all newly cre-  
357 ated content remains consistent with the intended  
358 research scope of the original licenses and miti-  
359 gates potential concerns related to data reuse and

redistribution.

**Biased Language** As the dialectal variants were  
generated using GPT-4o, we rely on the model’s  
built-in safety mechanisms to not generate outputs  
that may contain biased, offensive, or contextually  
inappropriate language.

## Positive Impact of Commonsense Validation

Our work advances existing methods and datasets  
for Arabic commonsense validation by introduc-  
ing dialectal variants and exploring novel model-  
ing approaches within this domain. We believe  
that enhancing commonsense understanding across  
Arabic dialects can contribute meaningfully to real-  
world applications such as fake news detection,  
fact-checking, and mitigating the spread of mis-  
leading or harmful content.

## References

- Muhammad Abdul-Mageed, AbdelRahim Elmadany,  
and El Moatez Billah Nagoudi. 2021. **ARBERT &  
MARBERT: Deep bidirectional transformers for Ara-  
bic**. In *Proceedings of the 59th Annual Meeting of the  
Association for Computational Linguistics and the  
11th International Joint Conference on Natural Lan-  
guage Processing (Volume 1: Long Papers)*, pages  
7088–7105, Online. Association for Computational  
Linguistics.
- Norah Alshahrani, Saied Alshahrani, Esmā Wali, and  
Jeanna Matthews. 2024. **Arabic synonym BERT-  
based adversarial examples for text classification**. In  
*Proceedings of the 18th Conference of the European  
Chapter of the Association for Computational Lin-  
guistics: Student Research Workshop*, pages 137–  
147, St. Julian’s, Malta. Association for Computa-  
tional Linguistics.
- Wissam Antoun, Fady Baly, and Hazem Hajj. 2020.  
**AraBERT: Transformer-based model for Arabic lan-  
guage understanding**. In *Proceedings of the 4th Work-  
shop on Open-Source Arabic Corpora and Process-  
ing Tools, with a Shared Task on Offensive Language  
Detection*, pages 9–15, Marseille, France. European  
Language Resource Association.
- Gagan Bhatia, El Moatez Billah Nagoudi, Abdellah  
El Mekki, Fakhreddin Alwajih, and Muhammad  
Abdul-Mageed. 2025. **Swan and ArabicMTEB:  
Dialect-aware, Arabic-centric, cross-lingual, and  
cross-cultural embedding models and benchmarks**.  
In *Findings of the Association for Computational  
Linguistics: NAACL 2025*, pages 4654–4670, Al-  
buquerque, New Mexico. Association for Computa-  
tional Linguistics.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng  
Gao, and Yejin Choi. 2020. **PIQA: reasoning about  
physical commonsense in natural language**. In *The*

413	<i>Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020</i> , pages 7432–7439. AAAI Press.		
414			
415			
416			
417			
418			
419			
420	Ernest Davis and Gary Marcus. 2015. Commonsense reasoning and commonsense knowledge in artificial intelligence. <i>Commun. ACM</i> , 58(9):92–103.		
421			
422			
423	Li Du, Xiao Ding, Kai Xiong, Ting Liu, and Bing Qin. 2022. e-CARE: a new dataset for exploring explainable causal reasoning. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 432–446, Dublin, Ireland. Association for Computational Linguistics.		
424			
425			
426			
427			
428			
429			
430	Javid Ebrahimi, Hao Yang, and Wei Zhang. 2021. How does adversarial fine-tuning benefit bert? <i>CoRR</i> , abs/2108.13602.		
431			
432			
433	Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. 2016. Domain-adversarial training of neural networks.		
434			
435			
436			
437	Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. <i>Neural Comput.</i> , 9(8):1735–1780.		
438			
439			
440	Jena D. Hwang, Chandra Bhagavatula, Ronan Le Bras, Jeff Da, Keisuke Sakaguchi, Antoine Bosselut, and Yejin Choi. 2021. (comet-) atomic 2020: On symbolic and neural commonsense knowledge graphs. In <i>Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021</i> , pages 6384–6392. AAAI Press.		
441			
442			
443			
444			
445			
446			
447			
448			
449			
450			
451	Go Inoue, Bashar Alhafni, Nurpeis Baimukan, Houda Bouamor, and Nizar Habash. 2021. The interplay of variant, size, and task type in Arabic pre-trained language models. In <i>Proceedings of the Sixth Arabic Natural Language Processing Workshop</i> , pages 92–104, Kyiv, Ukraine (Virtual). Association for Computational Linguistics.		
452			
453			
454			
455			
456			
457			
458	Mete Ismayilzada, Debjit Paul, Syrielle Montariol, Mor Geva, and Antoine Bosselut. 2023. CRoW: Benchmarking commonsense reasoning in real-world tasks. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 9785–9821, Singapore. Association for Computational Linguistics.		
459			
460			
461			
462			
463			
464			
465	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L��lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth��e Lacroix, and William El Sayed. 2023. Mistral 7b. <i>CoRR</i> , abs/2310.06825.		
466			
467			
468			
469			
470			
471			
472			
	Zhang Jiawei, Zhang Haopeng, Xia Congying, and Sun Li. 2020. GRAPH-BERT: Only attention is needed for learning graph representations.		473 474 475
	Akbar Karimi, Leonardo Rossi, and Andrea Prati. 2021. Adversarial training for aspect-based sentiment analysis with bert. In <i>2020 25th International Conference on Pattern Recognition (ICPR)</i> , pages 8797–8803.		476 477 478 479
	M Moneb Khaled, Aghyad Al Sayadi, and Ashraf Elngar. 2023. Commonsense validation and explanation in arabic text: A comparative study using arabic bert models. In <i>2023 24th International Arab Conference on Information Technology (ACIT)</i> , pages 1–6.		480 481 482 483 484
	Thomas N. Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. In <i>5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings</i> . OpenReview.net.		485 486 487 488 489 490
	Salima Lamsiyah, Kamyar Zeinalipour, Samir El amary, Matthias Brust, Marco Maggini, Pascal Bouvry, and Christoph Schommer. 2025. ArabicSense: A benchmark for evaluating commonsense reasoning in Arabic with large language models. In <i>Proceedings of the 4th Workshop on Arabic Corpus Linguistics (WACL-4)</i> , pages 1–11, Abu Dhabi, UAE. Association for Computational Linguistics.		491 492 493 494 495 496 497 498
	Hector J. Levesque, Ernest Davis, and Leora Morgenstern. 2012. The winograd schema challenge. In <i>Principles of Knowledge Representation and Reasoning: Proceedings of the Thirteenth International Conference, KR 2012, Rome, Italy, June 10-14, 2012</i> . AAAI Press.		499 500 501 502 503 504
	Bill Yuchen Lin, Xinyue Chen, Jamin Chen, and Xiang Ren. 2019. KagNet: Knowledge-aware graph networks for commonsense reasoning. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2829–2839, Hong Kong, China. Association for Computational Linguistics.		505 506 507 508 509 510 511 512
	Bill Yuchen Lin, Wangchunshu Zhou, Ming Shen, Pei Zhou, Chandra Bhagavatula, Yejin Choi, and Xiang Ren. 2020. CommonGen: A constrained text generation challenge for generative commonsense reasoning. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 1823–1840, Online. Association for Computational Linguistics.		513 514 515 516 517 518 519
	Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. In <i>International Conference on Learning Representations</i> .		520 521 522
	OpenAI. 2024. Gpt-4o system card. <i>CoRR</i> , abs/2410.21276.		523 524
	Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. How multilingual is multilingual BERT? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4996–5001, Florence, Italy. Association for Computational Linguistics.		525 526 527 528 529 530
	Abdelrahman Sadallah, Junior Cedric Tonga, Khalid Almubarak, Saeed Almheiri, Farah Atif, Chatrine		531 532



## Appendix

### A Related Work

#### A.1 Commonsense Reasoning Datasets

**Common Sense Reasoning in English** There have been many benchmarks for English commonsense reasoning, such as CommonSenseQA (Talmor et al., 2019), ComVe (Wang et al., 2020), ATOMIC (Sap et al., 2019a) and ATOMIC 2020 (Hwang et al., 2021). Within the broader scope of commonsense reasoning, several specialized subfields have emerged, each targeting distinct types of implicit human knowledge. Earlier work focused on pronoun coreference resolution in linguistic contexts (Levesque et al., 2012), physical commonsense reasoning (Bisk et al., 2020), social reasoning (Sap et al., 2019b), and causal reasoning (Du et al., 2022). Additional efforts have explored commonsense in natural language generation (Lin et al., 2020), as well as the integration of commonsense reasoning into real-world NLP tasks (Ismayilzada et al., 2023).

Despite these advancements, most research and benchmarks are centered around English, leaving many other languages, such as Arabic, under-resourced.

**CommonSense Reasoning in Arabic** Recent years have witnessed exploration of Arabic commonsense reasoning. Initial efforts focused on translating English commonsense benchmarks into Modern Standard Arabic (MSA) (Tawalbeh and Al-Smadi, 2020), or leveraging large language models (LLMs) to generate MSA data from seed data (Lamsiyah et al., 2025). However, these datasets lack cultural nuances of Arabic. Recent work by Sadallah et al. (2025) fills this gap by collecting a dataset covering cultures of 13 countries across the Gulf, Levant, North Africa, and the Nile Valley. Despite this advancement, their dataset remains restricted to MSA and does not encompass the rich linguistic and cultural diversity embedded in Arabic dialects. Therefore, we collect the first Arabic dialects commonsense reasoning benchmark, extending commonsense evaluation to Arabic dialects, aiming to capture more authentic and regionally grounded reasoning patterns.

#### A.2 Approaches for Commonsense Reasoning

Prior research has primarily focused on fine-tuning transformer-based models or employing LLMs for

commonsense validation and explanation generation, without introducing improved task-specific representations that could enhance performance. Tawalbeh and Al-Smadi (2020) fine-tuned BERT, USE, and ULMFit models for binary classification, selecting the more plausible sentence from a pair. More recently, Lamsiyah et al. (2025) evaluated a suite of BERT-based encoders, including AraBERTv2 (Antoun et al., 2020), ARBERT, MARBERTv2 (Abdul-Mageed et al., 2021), CaMeLBERT, and mBERT (Pires et al., 2019), on two classification tasks: (i) distinguishing commonsensical from nonsensical statements, and (ii) identifying the underlying reasoning behind nonsensicality. They also assessed causal LLMs including Mistral-7B (Jiang et al., 2023), LLaMA-3 (Touvron et al., 2023) and Gemma, on the two tasks above, along with the task (iii) generating natural language explanations for commonsense violations. These approaches lacked exploring better representation learning techniques to enhance the performance.

#### Integration of Adversarial Training with Encoder Transformer Models

Prior work has explored integrating adversarial training with transformer-based models. Karimi et al. (2021) introduced BERT Adversarial Training (BAT), which fine-tuned BERT and domain-specific BERT-PT using adversarial perturbations in the embedding space to improve robustness in Aspect-Based Sentiment Analysis (ABSA). Ebrahimi et al. (2021) showed that adversarial training can preserve BERT’s syntactic abilities, such as word order sensitivity and parsing, during fine-tuning, compared to standard fine-tuning. Additionally, it demonstrated how adversarial training prevented BERT from oversimplifying representations by reducing over-reliance on a few words, leading to better generalization.

In Arabic context, Alshahrani et al. (2024) conducted a synonym-based word-level adversarial attack on Arabic text classification models using a Masked Language Modeling (MLM) task with AraBERT. This attack replaces important words in the input text with semantically similar synonyms predicted by AraBERT to generate adversarial examples that can fool state-of-the-art classifiers. To ensure grammatical correctness, they utilize CaMeLBERT as a Part-of-Speech tagger to verify that the synonym replacements match the original word’s grammatical tags, maintaining sentence grammar.

We investigate the use of adversarial training across dialects as a means to learn more robust and generalized representations, thereby enhancing model performance and resilience across the diverse landscape of Arabic dialects.

**Integration of Graph-based Approaches with Encoder Transformer Models** Graph Neural Networks (GNNs) (Scarselli et al., 2009), and particularly Graph Convolutional Networks (GCNs) (Kipf and Welling, 2017) have gained significant attention for their ability to model relational and topological structures in data. Integrating graph-based structures with encoder-based Transformer models enables models to better grasp higher-level connections and contextual dependencies that are crucial for complex language understanding tasks like commonsense reasoning. For example, GraphBERT (Jiawei et al., 2020) introduced leveraging Transformer-style self-attention over linkless subgraphs, allowing it to learn graph representations without relying on explicit edge connections. This approach mitigates issues such as over-smoothing and enhances parallelizability. In contrast, VGCN-BERT (Zhibin et al., 2020) adopts a hybrid design, incorporating a vocabulary-level graph convolutional network (VGCN) into the BERT architecture. It constructs a global word co-occurrence graph and fuses the GCN-derived word representations with the BERT input embeddings, thereby enriching the model’s understanding of global corpus-level semantics. Both models demonstrate how graph-derived features, when fused effectively with Transformer encoders, can improve downstream tasks like text classification by fusing graph extracted morphological features with the token-level contextual embeddings. In the context of commonsense reasoning, (Lin et al., 2019) proposed KAGNet, a model that integrates GCNs with Long Short-Term Memory networks (LSTMs) (Hochreiter and Schmidhuber, 1997) to encode knowledge paths from external commonsense knowledge bases, thereby improving question answering performance through structured reasoning.

In this work, we integrate graph neural network-projected embeddings into transformer-based encoders, enriching contextual representations with global structural information which are critical to common sense validation.

## B Data Generation and Quality Control

**Dialect generation** We design a prompt tailored for accurate and meaning-preserved translation:

أنت خبير في اللهجات العربية. ترجم الجملة التالية  
إلى اللهجة {dialect} بدون تغيير المعنى:  
. {sentence}

The prompt translates to “You are an expert in Arabic dialects. Translate the following sentence to {dialect}: {sentence}”. This ensures that the intended meaning of each sentence remains intact while reflecting natural dialectal usage.

**Automatic Evaluation** We prompt Gemini 2.5 Flash to assess the faithfulness of dialectal translations as follows:

**System Prompt:** You are an expert Arabic linguist. Your task is to verify whether DIALECT\_TEXT is a correct translation of MSA\_TEXT from Modern Standard Arabic (MSA) into the specified Arabic dialect. A translation is correct if the meaning, events, entities, time references, and polarity in MSA\_TEXT are faithfully preserved in DIALECT\_TEXT, even if wording differs due to dialectal variation. Ignore minor spelling, punctuation, and orthographic differences. Do not allow additions, omissions, or factual changes. Output only one word: true or false. Do not explain your decision.

DIALECT: <dialect\_name>  
MSA\_TEXT: <msa\_text>  
DIALECT\_TEXT: <dialect\_text>  
Answer:

## C Graph Embeddings-based Encoder Transformer models

### C.1 Extended Explanation of methodology

Figure 2 and Algorithm 1 summarize the pipeline.

## D Domain Adversarial Training

### D.1 Method

We implement domain-adversarial training (DANN) (Ganin et al., 2016) to encourage dialect-invariant features during fine-tuning. A shared Transformer encoder produces a sentence representation  $\mathbf{h}$  from the final-layer [CLS] vector. We attach two MLP heads: (i) a task classifier for commonsense validation, and (ii) a dialect discriminator predicting the dialect label. The dialect discriminator receives  $\mathbf{h}$  through a Gradient Reversal Layer (GRL), which multiplies its backpropagated gradient by  $-\alpha$ , so the encoder is optimized to reduce task loss while making dialect prediction harder.

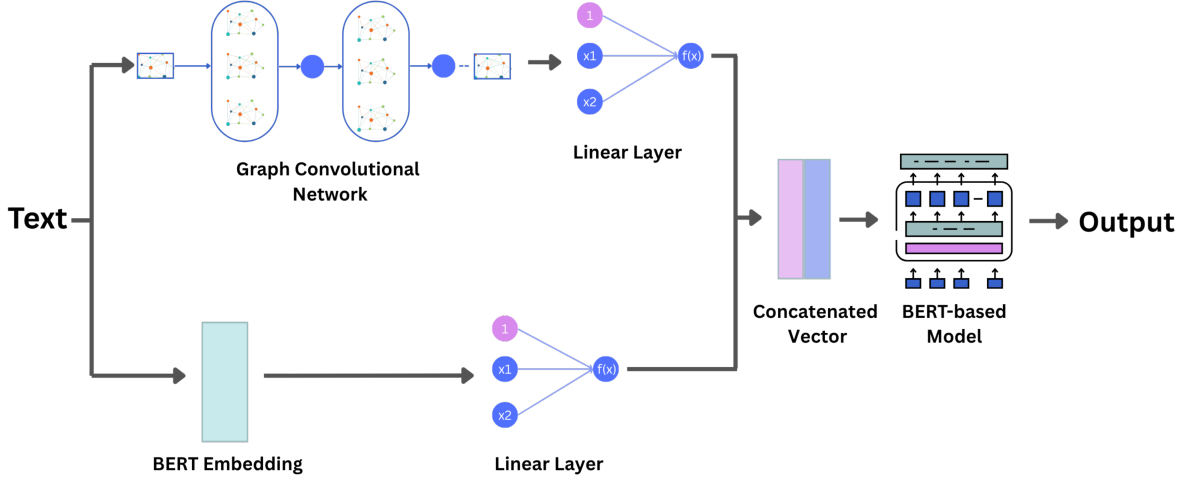


Figure 2: BERT Model with Graph Embeddings Fusion.

**Algorithm 1** The training algorithm for Graph Embeddings-based Encoder Transformer models.

```

1: Given:
2:  $\mathcal{D}_{\text{train}}$   $\triangleright$  Labeled corpora of text samples
3:  $\mathcal{T}$   $\triangleright$  Pretrained textual encoder (e.g., BERT)
4:  $\mathcal{G}$   $\triangleright$  Graph encoder (e.g., GCN)
5:  $\mathcal{F}$   $\triangleright$  Fusion mechanism (e.g., attention)
6:  $\mathcal{C}$   $\triangleright$  Classifier head
7:  $\theta$   $\triangleright$  Trainable parameters
8: Preparation:
9: for all  $(x, y) \in \mathcal{D}$  do
10:   Tokenize  $x \rightarrow \mathbf{t} \in \mathbb{R}^{L_h}$ 
11:   Convert  $x \rightarrow$  graph  $\mathcal{G}_x = (V, E, \mathbf{X})$ 
12: end for
13: Initialize:
14:  $\theta \leftarrow$  random or pretrained weights
15: for  $e = 1$  to  $E$  do
16:   Training Step:
17:   for all  $(x, y, \mathcal{G}_x) \in \mathcal{D}_{\text{train}}$  do
18:      $\mathbf{z}_t \leftarrow \mathcal{T}(x)$   $\triangleright$  Textual representation
19:      $\mathbf{z}_g \leftarrow \mathcal{G}(\mathcal{G}_x)$   $\triangleright$  Graph representation
20:      $\mathbf{z}_f \leftarrow \mathcal{F}(\mathbf{z}_t, \mathbf{z}_g)$   $\triangleright$  Fusion
21:      $\hat{y} \leftarrow \mathcal{C}(\mathbf{z}_f)$   $\triangleright$  Prediction
22:     Update  $\theta$  via  $\nabla_{\theta} \mathcal{L}(\hat{y}, y)$ 
23:   end for
24: end for

```

Both heads use the same architecture: Dropout  $\rightarrow$  Linear( $H \rightarrow 768$ )  $\rightarrow$  ReLU  $\rightarrow$  Dropout  $\rightarrow$  Linear( $768 \rightarrow C$ ), with dropout rate 0.1, where  $C=2$  for the main task and  $C=5$  for dialect prediction.

We optimize the combined objective:

$$\mathcal{L} = \mathcal{L}_{\text{main}} + \lambda \mathcal{L}_{\text{dial}}, \quad (1)$$

where both terms are cross-entropy losses and the GRL applies the adversarial signal to the shared encoder. We set  $\lambda=1.0$  and use a simple schedule for the GRL strength  $\alpha = \min(1, \frac{e+1}{5})$  as a function of epoch index  $e$ .

## D.2 Experimental Setup

We follow the same data split and base fine-tuning setup described in Section 4.1. Adversarial models are trained for 3 epochs using AdamW (learning rate  $2e-5$ , weight decay 0.01), with gradient clipping (max norm 1.0) and a linear warmup over the first 100 optimization steps (from 0.1lr to 1r, then constant). We select the best checkpoint using the dev weighted F1 of the main task and report results on the held-out test set.

## D.3 Results

Table 4 compares the baseline models against their domain-adversarial training variants. Overall, domain-adversarial training consistently leads to a substantial degradation in accuracy, indicating that enforcing domain invariance in this setting harms learning rather than improving cross-dialect generalization.

**Why Does Adversarial Training Fail?** All three backbones degrade under adversarial training, and the degradation is especially severe for AraBERTv2 (dropping to near chance). This pattern strongly suggests that dialect-specific cues are

Methods	MSA $\uparrow$	Egyptian $\uparrow$	Gulf $\uparrow$	Levantine $\uparrow$	Moroccan $\uparrow$	Avg. $\uparrow$
AraBERTv2	<b>65.88</b>	<b>68.33</b>	<b>65.78</b>	<b>65.36</b>	<b>62.09</b>	<b>65.34</b>
AraBERTv2 (Adv)	50.52	50.15	49.89	50.03	50.55	50.23
CAMeLBERT-mix	<b>73.30</b>	<b>73.52</b>	<b>74.63</b>	<b>73.82</b>	<b>66.45</b>	<b>72.34</b>
CAMeLBERT-mix (Adv)	66.52	67.12	67.26	68.91	64.18	66.80
MARBERTv2	<b>80.12</b>	<b>80.73</b>	<b>78.81</b>	<b>80.03</b>	<b>71.09</b>	<b>78.16</b>
MARBERTv2 (Adv)	79.58	79.00	77.57	78.52	70.24	76.98

Table 4: Accuracy (%) of baselines and adversarial training-based models on MSA and dialects datasets.

840 *not* purely nuisance variation for this task: enforcing  
841 dialect invariance likely removes information  
842 that is genuinely predictive (lexical, morphological,  
843 or orthographic markers correlated with the label).  
844 Another plausible contributor is optimization insta-  
845 bility: if the adversarial signal is too strong relative  
846 to the supervised signal, feature collapse can occur.  
847 An effect that would be amplified for a less dialect-  
848 robust backbone such as AraBERTv2. Practically,  
849 these results argue against treating dialect simply  
850 as a domain to be “erased”; a better direction may  
851 be *domain-aware* modeling (e.g., dialect embed-  
852 dings/adapters or mixture-of-experts) rather than  
853 domain-invariant representations. In contrast to ad-  
854 versarial training, adding GCN-based embeddings  
855 improves *every* backbone.