

LARA BENCH: Benchmarking Arabic AI with Large Language Models

Anonymous ACL submission

Abstract

Recent advancements in Large Language Models (LLMs) have significantly influenced the landscape of language and speech research. Despite this progress, these models lack specific benchmarking against state-of-the-art (SOTA) models tailored to particular languages and tasks. *LARA BENCH* addresses this gap for Arabic Natural Language Processing (NLP) and Speech Processing tasks, including sequence tagging and content classification across different domains. We utilized models such as GPT-3.5-turbo, GPT4, BLOOMZ, Whisper, and USM, employing zero and few-shot learning techniques to tackle 33 distinct tasks across 61 publicly available datasets. This involved 98 experimental setups, encompassing ~296K data points, ~46 hours of speech, and 30 sentences for Text-to-Speech (TTS). This effort resulted in 330+ sets of experiments. Our analysis focused on measuring the performance gap between SOTA models and LLMs. The overarching trend observed was that SOTA models generally outperformed LLMs in zero-shot learning, with a few exceptions. Notably, larger computational models with few-shot learning techniques managed to reduce these performance gaps. Our findings provide valuable insights into the applicability of LLMs for Arabic NLP and speech processing tasks.

1 Introduction

Generative Pre-trained Transformer (GPT) models are examples of large language models (LLMs) trained on massive datasets and using hundreds of millions of parameters. Several LLMs have been recently released for use through APIs or pre-trained models and have demonstrated a high level of coherence in generating content in response to specific user tasks. However, quality assessments of released LLMs generally lack references to previous research and comparison with state-of-the-art (SOTA) methods that the research community has

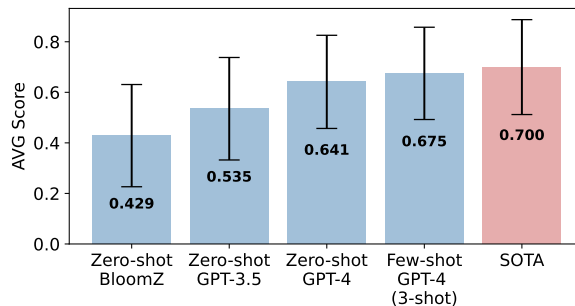


Figure 1: Average performance of the models as compared to SOTA across 21 unique NLP tasks and 31 testing setups.

used for systematic evaluation and monitoring of scientific progress for various languages and tasks.

Several research initiatives have evaluated these large models' performance on standard NLP and speech processing tasks. The HELM project (Liang et al., 2022) assessed English LLMs across various metrics and scenarios. BIG-Bench (Srivastava et al., 2022) introduced a large-scale evaluation with 214 tasks, including low-resource languages. GPT2.5 (Radford et al., 2019), ChatGPT (OpenAI, 2023), and BLOOM (Scao et al., 2022), were recently evaluated by Bang et al. (2023); Ahuja et al. (2023); Hendy et al. (2023); Khondaker et al. (2023). Large speech models such as Whisper (Radford et al., 2022) and Universal Speech Model (USM) (Zhang et al., 2023) were also explored for speech recognition and translation tasks. Initiatives such as SUPERB (Yang et al., 2021) were introduced to support benchmarking tools and leaderboards for several speech-related tasks. Bubeck et al. (2023) explored GPT-4's capabilities to determine if it surpasses mere memorization, possessing a profound and adaptable comprehension of concepts, skills, and domains. Their results indicate that GPT-4 demonstrates a higher level of general intelligence compared to its predecessors.

LARA BENCH study fulfills an important objective of assessing the LLMs capabilities for supporting

070 Arabic language processing tasks, for Modern Stan- 121
071 dard Arabic (MSA) and dialectal Arabic (DA), at 122
072 the same level of depth and breadth as for English 123
073 tasks. Our evaluation involves 61 publicly avail- 124
074 able datasets and 98 test setups used to perform 125
075 and evaluate language processing tasks in both 126
076 MSA and dialectal Arabic across various genres 127
077 (e.g., news articles, tweets, meetings, telephony, 128
078 and broadcast content). Our evaluation focuses on 129
079 assessing the capabilities of GPT-3.5-turbo, GPT- 130
080 4, and BLOOMZ (176B) for NLP tasks, and of 131
081 Whisper (Large, 1.55B) and USM (2B) for Speech 132
082 processing, in both zero and few-shot settings. We 133
083 investigate: (i) *can LLMs effectively perform Ara- 134
084 bic NLP and Speech processing tasks without prior 135
085 task-specific knowledge (zero-shot)?* (ii) *how does 136
086 performance vary across tasks with different com- 137
087 plexities in zero- and few-shot settings?* (iii) *how 138
088 do LLMs compare to current SOTA models, and 139
089 are open LLMs as effective as the commercially 140
090 available (closed) models?* Our investigation re- 141
091 veals unique insights about LLMs’ performance on 142
092 Arabic NLP and Speech tasks: 143

093 **A. Zero-shot Multi-task Performer.** GPT-4 out- 144
094 performs GPT-3.5 and BLOOMZ in majority of 145
095 the NLP tasks (see Figure 1). However, a large per- 146
096 formance gap between GPT-4 and SOTA models 147
097 remains due to the higher quality SOTA models. 148
098 For speech tasks, USM outperforms all the Whisper 149
099 variants and performs comparably with SOTA. 150

100 **B. Few-shot and SOTA.** GPT-4 reduces the per- 151
101 formance gap with SOTA in the few-shot (only 152
102 3-shots) setting (see Figure 1). This significant per- 153
103 formance gain is noticed for almost all tasks, par- 154
104 ticularly for more complex semantic and question- 155
105 answering tasks compared to syntactic and seg- 156
106 mentation tasks. Similarly, whisper models exhibit 157
107 promising results in speech recognition with just 158
108 2 hours of speech examples in few-shot finetuning. 159
109 Open models (BLOOMZ and Whisper) performed 160
110 poorly w.r.t. to their commercially available coun- 161
111 terparts. However, fine-tuning with more instruc- 162
112 tions may help these open models to achieve closer 163
113 performance to SOTA and other closed LLMs. 164

114 **C. MSA vs Dialect.** The gaps in LLMs’ perfor- 165
115 mance between MSA and dialectal datasets (e.g., 166
116 for machine translation (MT) and speech recogni- 167
117 tion task) are more pronounced, indicating ineffec- 168
118 tiveness of LLMs for under-represented dialects. 169

119 **D. Hallucination and Data Contamination.**

120 GPT models, specially GPT-3.5, suffer from the

hallucination problem. We noticed the model 121
insert extra information (e.g., for MT task with 122
Bible dataset) from its parametric memory. 123

Benchmarking LLMs raises concerns about their 124
exposure to existing datasets. In our study, we uti- 125
lized datasets that were released after the cut-off 126
date of GPT’s training (September 2021). More- 127
over, we applied a prompt-based method with tai- 128
lored instructions (Golchin and Surdeanu, 2023) 129
on nine datasets using GPT-4, to determine if these 130
datasets are present in the GPT-4 model. Our ex- 131
periments yielded no examples from these datasets. 132

To the best of our knowledge, our *LAraBench* 133
study is the *first* comprehensive effort that includes 134
commercial (close) and open source LLMs and 135
evaluates zero- and few-shot settings for a wide 136
range of Arabic NLP and Speech tasks. It is the 137
first to include the evaluation of Whisper and USM 138
models for Arabic ASR and the *first* to report bench- 139
marks for a standard Arabic TTS generative model. 140
All resources and findings of the *LAraBench* study 141
will be made publicly available to the community 142
to scale up the effort.¹ 143

144 **2 Tasks and Datasets**

145 The LAraBench study was designed with an ambi- 146
tious goal of empowering the research community 147
and practitioners with the most comprehensive eval- 148
uation of LLMs for Arabic NLP and Speech tasks 149
to date. It includes 61 publicly available datasets to 150
support 9 task groups² discussed below. We briefly 151
describe each task and refer to Appendix A for a 152
comprehensive description of tasks and datasets.

Word Segmentation, Syntax and Information 153
Extraction. We explore six sequence tagging 154
tasks: i) word segmentation, ii) POS-tagging, iii) 155
lemmatization, iv) diacritization, v) parsing, and 156
vi) named-entity recognition (NER), using pub- 157
licly available datasets. We also include a dialect 158
identification task (e.g., Egyptian dialect) since vo- 159
cabulary, pronunciation, grammar, and idiomatic 160
expressions vary across dialects. For our bench- 161
marking we used QADI (Abdelali et al., 2021) and 162
ADI (in-house) datasets. 163

Machine Translation (MT). Machine transla- 164
tion of Arabic is challenging due to morphological 165
complexity and dialectal variations. We experiment 166

¹<http://anonymous.com/>

²Our task categorization is derived from the taxonomy outlined in the list of tracks established by ACL 2022.

with AraBench (Sajjad et al., 2020), an extensive suite of data from diverse genres and dialects.

Sentiment, Stylistic and Emotion Analysis.

These tasks involve understanding and analyzing aspects of human expression and communication. We benchmark Sentiment Analysis, Emotion Recognition, Stance Detection, and Sarcasm Detection with datasets from Elmadany et al. (2018), Mohammad et al. (2018), Chouigui et al. (2017), and Abu Farha et al. (2021), respectively.

News Categorization. This task involves classification of news articles into pre-defined categories or topics (Sebastiani, 2002). We support benchmark evaluations using SANAD news article corpus (Einea et al., 2019) and ASND social media dataset (Chowdhury et al., 2020b).

Demographic Attributes. Demographic information, including gender, age, and country of origin, hold significant value across various applications such as population analysis. We include datasets that enable experimentation with tasks of identifying country, gender (Mubarak et al., 2022) and location (Mubarak and Hassan, 2021).

Ethics and NLP: Factuality, Disinformation and Harmful Content Detection.

These tasks have emerged as critical areas within the field of NLP. We support benchmarking of several *detection* tasks, such as: i) Propaganda (Alam et al., 2022b), ii) Factuality using the datasets Baly et al. (2018a); Alam et al. (2021b); Khouja (2020), iii) Harmful content (Nakov et al., 2022b), iv) Offensive language (Zampieri et al., 2020), and v) Hate speech (Mubarak et al., 2021a).

Semantics. This task group includes Semantic Textual Similarity (STS) and Natural Language Inference (NLI). We benchmark STS using two datasets: SemEval-2017 STS task (Cer et al., 2017a) and similarity in Arabic question pairs, as explored by Seelawi et al. (2019). For the XNLI task, we used the translated version of Arabic from XNLI corpus (Conneau et al., 2018).

Question Answering (QA). For the QA task, we employed ARCD (Mozannar et al., 2019), MLQA (Lewis et al., 2019), TyDiQA (Clark et al., 2020), and XQuAD (Artetxe et al., 2020) datasets.

Speech Processing. We evaluate the large speech models on two tasks: speech recognition (ASR) and text-to-speech (TTS) synthesis. For ASR, we

include datasets varying domain and dialects, e.g. MGB2 (Ali et al., 2016), QASR.CS (Mubarak et al., 2021b) and ESCWA.CS (Ali et al., 2021a). For TTS, we evaluated with in-house 30 test sentences, covering diverse topics (e.g., education, health).

3 Methodology

For benchmarking of Arabic NLP and Speech processing tasks, we use zero- and few-shot learning involving GPT-3.5-Turbo, GPT-4 and BLOOMZ for NLP, and Whisper (small, medium, and large), USM and Amazon Polly for Speech. We also compared LLM’s performance with the respective SOTA models.

The use and evaluation of LLMs involve prompting and post-processing of output to extract the expected content. Therefore, for each task, we explored a number of prompts, guided by the same instruction and format as recommended in the Azure OpenAI Studio Chat playground, and PromptSource (Bach et al., 2022). After obtaining a reasonable prompt, we used it to complete the evaluation using the task and modality-specific API services, e.g., OpenAI API from Azure for NLP tasks and Google’s USM API for Speech tasks. For BLOOMZ, we set up on-premises hosting and use.

We based our model selection on factors like performance, language support, and accessibility. For NLP tasks, we chose OpenAI models because they consistently outperformed others for English tasks. Initially, we used GPT-3.5 and later transitioned to GPT-4 when it became available. Limited budget and lack of Arabic language support led us to avoid other closed models. Among open models, we selected BLOOMZ because it’s a large multilingual model, including 4% Arabic content.³ For ASR, we chose Whisper and USM due to their excellent performance in recent studies.

3.1 Models and Prompts for NLP Tasks

Zero-shot Setup. For tasks with GPT-3.5-Turbo, GPT-4 and BLOOMZ, we use zero-shot prompting giving natural language instructions describing the task and specify the expected output. We allow the LLM itself to build context that narrows the inference space and produces more accurate output.

Few-shot Setup. In order to explore the maximum potential of specific LLMs, e.g., GPT-4

³New models such as JAIS (Sengupta et al., 2023) and AceGPT (Huang et al., 2023) have been released as we speak and we leave their benchmarking for the future.

model, we used available training data to select few-shot examples and provide context for the task. For a few tasks and datasets (e.g., location, name to country), training sets are either private or not available and therefore they could not be included in our few-shot experiments. We used maximal marginal relevance-based (MMR) selection to construct example sets that are deemed relevant and diverse (Carbonell and Goldstein, 1998), following the proven method by Ye et al. (2022). The MMR method computes the similarity between a test example and the example pool (e.g., training dataset) and selects m examples (shots). We apply MMR on top of embeddings of multilingual sentence-transformers (Reimers and Gurevych, 2019). In our few-shot investigation, we performed experiments on all tasks and datasets using only 3-shots to primarily reduce computational and API expenses. Additionally, we expanded our analysis to include 3, 5, and 10 instances across seven distinct datasets drawn from various task categories. More details are provided in Section F.2 of the Appendix.

Prompts Design. Prompt design is a complex and iterative process that present challenges due to the unknown representation of information within LLMs and a need for different types of outputs across tasks, e.g., token classification vs. sentence classification. The instructions expressed in our prompts were in English, including the content examples in Arabic. In Appendix C, we provide examples of prompts for different tasks. We also examined Arabic instructions in our study, to understand the effect of native language prompts. For this set of experiments we selected seven datasets from seven different task groups. More details can be found in Section F.3 (Appendix).

Post Processing. Outputs of LLMs are post-processed to enable automatic comparison with gold standard labels. Depending on the task, this may include mapping prefixes, or filtering tokens. For example, for POS tagging, the tags ‘*preposition*’, ‘*P*’, ‘*PRP*’, ‘*حرف جر*’, are mapped onto *PREP*. For NER, the model switches the tag of the prediction i.e., B-PER predicts as PER-B, and therefore requires remapping of the NER tags.

3.2 Models and Prompts for Speech Tasks

We use zero- and few-shot settings to benchmark large speech models. For ASR, we use three Whisper models (OpenAI) – small, medium, and large,

and the USM model (Google). For the details of the models, see Table B.2 in Appendix. We compare these large models to SOTA: supervised KANARI⁴ conformer-based (Chowdhury et al., 2021) offline and RNN-T based streaming ASR.⁵ For the TTS task, we compare two public systems: Amazon Polly TTS engine⁶ and KANARI TTS system.⁷

Zero-shot Setup. For zero-shot setup, we use the initial (or pre-trained) weights of Whisper and API of USM models with a goal to benchmark the performances of these LLMs in different domains, for different Arabic dialects, and for code-switching with no domain knowledge. As a prompt to the model, we passed only a language flag.

Few-shot Setup. Under this setup, we fine-tune Whisper (small and large) with 2 hours of domain-specific speech data and compare it with the SOTA models trained from scratch with 3K hours of speech.

ASR Post Processing. ASR is evaluated based on word error rate (WER) that aligns the model’s output with reference transcription and penalizes the output based on insertion, deletion, and substitution errors. The measure is unable to disambiguate code-switching and minor formatting differences introduced by multilingual scripts or non-standardized orthography. Hence, post-processing is a crucial component. We normalized ‘alif’, ‘ya’ and ta-marbuta’, and adapted a minimalist Global Mapping File (GLM) (Chowdhury et al., 2021) to transliterate common words and handle rendering mismatch. Thus keeping room for further improvement with more enhanced post-processing.

3.3 Random Baseline

We also calculated a random baseline for the NLP tasks (further details can be found in Appendix, Section F.1). The aim is to determine if the LLMs predictions are not merely the result of chance. It also serves as a lower limit to be expected for each task.

3.4 Evaluation Metrics

To measure the performance of each task, we followed current state-of-art references and used the metric reported in the respective work. This includes: Accuracy (Acc), F1 (macro, micro, and

⁴<https://fenek.ai/>

⁵<https://arabicasr.kanari.ai/>

⁶<https://aws.amazon.com/polly/>

⁷<https://arabic tts.kanari.ai/>

weighted), word error rate (WER), Jaccard Similarity (JS), Pearson Correlation (PC), and mean opinion score (MOS) for naturalness, intelligibility and diacritization. We report average MOS (10-point Likert scale) from 3 native-annotators.

4 Results and Discussion

In Tables 1, 2, 3 and 4, we report the results of different NLP and Speech related tasks. In the below sections, we summarize the results and challenges specific to the task groups.

4.1 NLP Tasks

In Table 1, we report the random baseline, GPT-3.5, GPT-4 (zero-shot and few-shot), and BLOOMZ and compare them to SOTA.⁸ In almost all tasks, models outperform random baseline, indicating that the predictions of the models are not by chance.

Word Segmentation, Syntax and Information Extraction. As Table 1 shows, for almost all tasks in this group, the performance is significantly below SOTA performance. For example, the difference between SOTA and GPT-4 (zero-shot) ranges from 19.4% (NER) to 73.8% (segmentation).

Machine Translation. Table 2 reports MT results by averaging them dialect-wise for different datasets. Appendix F.6 reports detailed results. The results indicate the short-coming of LLMs when explored with standard and dialectal Arabic.

Sentiment, Stylistic and Emotion Analysis. In the second group of Table 1, we report results for sentiment, emotion, stance and sarcasm detection mainly over tweets. We observe that performance gap significantly reduced between GPT-4 (best of zero- and few-shot) vs. SOTA compared to GPT-3.5 vs. SOTA, 4.77% vs 15.14%, respectively. For sarcasm detection task with ArSarcasm dataset, GPT-4 even outperformed SOTA by 4.41%.

News Categorization. Table 1 shows that performance gap reduced significantly ranging from 15.63% to 6.44% for GPT-3.5 to GPT-4, respectively. Low performance on tweet dataset (ASND) might be due to the higher number of class labels.

Demographic/Protected Attributes. Among the three tasks in this group, two of them (“name info”

⁸Note that some results are missing either due to the unavailability of training data (marked with NA) or the incapability of the Bloomz model (marked with ‡).

and “location” identification) demonstrate a significant performance improvement (by 3.62%) over the SOTA results, using GPT-4 model.

Ethics and NLP: Factuality, Disinformation and Harmful Content Detection. Across eleven tasks, the performance gap significantly reduced with GPT-4 model, however in some tasks, model’s performance is significantly lower than the SOTA. For example, for factuality with COVID-19 disinfo. dataset, GPT-4 model’s performance is 33% lower than the SOTA, even though performances of GPT-4 significantly improved compared to GPT-3.5. This task is generally challenging requiring deep contextual analysis and reasoning abilities, and domain knowledge in many of the cases. With a few demonstrations (only 3-shots) may not be enough to determine the factuality of the content.

Semantics: The results for various semantic tasks reported in Table 1 indicate that the performance on three out of the four tasks surpasses the SOTA, with an overall improvement of 7%.

Question answering (QA): Results on four QA datasets (Table 1) show that for three of them, GPT-4 achieved higher performance than SOTA with an overall improvement of 4.42%.

4.2 Speech Recognition and Synthesis

In Table 3, we reported the performance of ASR using different datasets and models. We observed that USM outperforms Whisper in all datasets in both zero and few-shot setting. The USM model performs comparably to standard task- and domain-specific ASR systems and is better equipped to handle cross-language and dialectal code-switching data from unseen domains compared to the SOTAs and Whispers few-shot finetuned model.

Both the subjective and objective evaluations for the TTS are reported Table 4. The results show that KANARI models outperformed Amazon Polly significantly in objective evaluation (WER). Subjective scores show KANARI is better in naturalness and diacritization. With almost similar performance in intelligibility.

5 Findings

NLP Model Performances. Our comprehensive study highlights the disparities in performance of LLMs – GPT-3.5 and GPT-4, as compared to SOTA models, in zero and few-shot settings. GPT-3.5

Task Name	Dataset	Metric	Random Baseline	BLOOMZ	Zero-shot GPT-3.5	Zero-shot GPT-4	Few-Shot GPT-4 (3-shot)	SOTA
Word Segmentation, Syntax and Information Extraction								
Segmentation	WikiNews	Acc	0.272	‡	0.195	0.252	0.927	0.989 (Darwish and Mubarak, 2016)
Segmentation	Samih et al. (2017)	Acc _{AVG}	0.309	‡	0.283	0.372	0.850	0.931 Samih et al. (2017)
Lemmatization	WikiNews	Acc	0.348	‡	0.471	0.397	NA	0.973 (Mubarak, 2018)
Diacritization	WikiNews	WER	0.963	‡	0.308	0.420	0.237	0.045 (Mubarak et al., 2019)
Diacritization	Darwish et al. (2018)	WER	0.999	‡	0.928	0.899	0.994	0.031 (Darwish et al., 2018)
POS	WikiNews	Acc	0.030	‡	0.231	0.479	0.367	0.953 (Darwish et al., 2017c)
POS	Samih et al. (2017)	Acc	0.036	‡	0.073	0.511	0.323	0.892 Samih et al. (2017)
POS	‡GLUE (Arabic)	Acc	0.032	‡	0.159	0.402	0.524	0.686 (Liang et al., 2020a)
Parsing	Conll2006	UAS	0.001	‡	0.239	0.504	0.551	0.796 (Lei et al., 2014)
NER	ANERcorp	F1 _{Macro}	0.008	‡	0.210	0.355	0.420	0.886 (Gridach, 2018)
NER	Aqmar	F1 _{Macro}	0.007	‡	0.230	0.365	0.390	0.690 (Schneider et al., 2012)
NER	QASR	F1 _{Macro}	0.009	‡	0.208	0.504	NA	0.698 (Mubarak et al., 2021b)
Dialect	QADI	F1 _{Macro}	0.052	0.067	0.149	0.243	NA	0.600 (Abdelali et al., 2021)
Dialect	ADI	F1 _{Macro}	0.092	0.098	0.169	0.229	0.260	0.26/0.57 (lexical/acoustic) (In-house)
Sentiment, Stylistic and Emotion Analysis								
Sentiment	ArSAS	F1 _{Macro}	0.222	0.251	0.550	0.569	0.598	0.758 (Hassan et al., 2021)
Emotion	SemEval18-Task1	JS	0.167	0.142	0.395	0.373	0.489	0.541 (Hassan et al., 2022)
Stance	Unified-FC	F1 _{Macro}	0.193	0.235	0.232	0.495	0.358	0.558 (Baly et al., 2018b)
Stance	ANS	F1 _{Macro}	0.281	0.223	0.620	0.762	0.721	0.767 (Khouja, 2020)
Sarcasm	ArSarcasm	F1 _(POS)	0.240	0.286	0.465	0.400	0.504	0.460 (Farha and Magdy, 2020)
Sarcasm	ArSarcasm-2	F1 _(POS)	0.333	0.436	0.537	0.573	0.537	0.623 (Alharbi and Lee, 2021)
News Categorization								
News Cat.	ASND	F1 _{Macro}	0.048	0.371	0.512	0.667	0.594	0.770 (Chowdhury et al., 2020b)
News Cat.	SANAD/Akhbarona	Acc	0.142	0.582	0.730	0.877	0.892	0.940 (Elnagar et al., 2020)
News Cat.	SANAD/AlArabiya	Acc	0.144	0.716	0.922	0.921	0.925	0.974 (Elnagar et al., 2020)
News Cat.	SANAD/AIKhaleej	Acc	0.142	0.738	0.864	0.911	0.899	0.969 (Elnagar et al., 2020)
Demographic Attributes								
Name Info	ASAD	F1 _{Weighted}	0.014	‡	0.570	0.629	NA	0.530 (Under review)
Location	UL2C	F1 _{Macro}	0.027	0.118	0.339	0.735	NA	0.881 (Mubarak and Hassan, 2021)
Gender	Arap-Tweet	F1 _{Macro}	0.521	0.532	0.883	0.868	0.980	0.821 (Mubarak et al., 2022)
Ethics and NLP: Factuality, Disinformation and Harmful Content Detection								
Offensive lang.	OffensEval2020	F1 _{Macro}	0.454	0.533	0.460	0.623	0.874	0.905 (Mubarak et al., 2020b)
Hate Speech	OSACT2020	F1 _{Macro}	0.376	0.503	0.430	0.669	0.644	0.823 (Mubarak et al., 2020c)
Adult Content	ASAD	F1 _{Macro}	0.421	0.513	0.460	0.727	0.832	0.889 (Mubarak et al., 2021a)
Spam	ASAD	F1 _{Macro}	0.405	0.152	0.440	0.745	NA	0.989 (Hassan et al., 2021)
Subjectivity	In-house	F1 _{Macro}	0.496	0.428	0.670	0.677	0.745	0.730 (In-house)
Propaganda	WANLP22	F1 _{Micro}	0.139	0.108	0.353	0.472	0.537	0.649 (Samir et al., 2022)
Check-worthy	CT-CWT-22	F1 _(POS)	0.398	0.431	0.526	0.560	0.554	0.628 (Du et al., 2022)
Factuality	COVID-19 Disinfo.	F1 _{Weighted}	0.582	0.749	0.393	0.485	0.491	0.831 (Alam et al., 2021b)
Factuality	Unified-FC	F1 _{Macro}	0.464	0.460	0.306	0.581	0.621	In-house
Factuality	ANS	F1 _{Macro}	0.505	0.550	0.252	0.539	0.704	0.643 (Khouja, 2020)
Claim	CT-CWT-22	Acc	0.498	0.532	0.703	0.587	0.686	0.570 (Eyuboglu et al., 2022)
Harmful content	CT-CWT-22	F1 _(POS)	0.269	0.144	0.471	0.533	0.494	0.557 (Bilel et al., 2022)
Attention-worthy	CT-CWT-22	F1 _{Weighted}	0.125	0.148	0.258	0.257	0.412	0.206 (Nakov et al., 2022a)
Semantics								
STS	STS2017-Track 1	PC	0.005	0.537	0.799	0.813	0.809	0.754 (Cer et al., 2017b)
STS	STS2017-Track 2	PC	-0.136	0.512	0.828	0.848	0.857	0.749 (Cer et al., 2017b)
STS QS (Q2Q)	Mawdoo3 Q2Q	F1 _{Micro}	0.491	0.910	0.816	0.895	0.935	0.959 (Seelawi et al., 2019)
XNLI (Arabic)	XNLI	Acc	0.332	0.500	0.489	0.753	0.774	0.713 (Artetxe et al., 2020)
Question answering (QA)								
QA	ARCD	F1 _(EM)	0.085	0.368	0.502	0.705	0.704	0.613 (Mozannar et al., 2019)
QA	MLQA	F1 _(EM)	0.066	0.377	0.376	0.620	0.653	0.548 (Lewis et al., 2019)
QA	TyDi QA	F1 _(EM)	0.111	0.456	0.480	0.744	0.739	0.820 (Clark et al., 2020)
QA	XQuAD	F1 _(EM)	0.047	0.367	0.442	0.729	0.722	0.665 (Artetxe et al., 2020)

Table 1: Results on NLP tasks. QS: Question similarity, PC: Pearson Correlation, JS: Jaccard Similarity, EM: Exact match, POS: positive class. Best result per row is **boldfaced**. NA: experiments could not be performed due to a lack of training data. BLOOMZ does not understand some tasks at all as marked with ‡ symbol.

445 exhibits a significant performance gap when compared to SOTA. However, GPT-4 manages to narrow this gap to some extent and even outperforms the SOTA models in high-level abstract tasks such as STS, QA, claim detection, news categorization, 448 demographic attributes, and XNLI. Moreover, GPT-4 449 outperforms GPT-3.5 across all tasks. However, 450 451

452 it remains a challenge for GPT-4 to surpass SOTA 453 performance consistently in sequence tagging (es- 454 pecially syntactic and segmentation) tasks. The 455 performance of BLOOMZ is significantly lower 456 than SOTA and GPT models, and in some cases 457 lower than random baseline. The performances of 458 both open and close models are heavily dependent 459

Dataset	Dialect	#Sent.	BloomZ	Zero-shot GPT-3.5	Zero-shot GPT-4	SOTA
APT	LEV	1000	11.38	18.55	17.77	21.90
APT	Nile	1000	12.95	21.58	18.99	22.60
MADAR	Gulf	16000	32.34	34.60	36.18	32.46
MADAR	LEV	12000	31.36	33.42	35.24	32.45
MADAR	MGR	14000	23.59	23.91	27.83	23.14
MADAR	MSA	2000	42.33	37.55	37.67	43.40
MADAR	Nile	8000	34.87	36.97	37.93	35.15
MDC	LEV	3000	10.00	17.38	16.05	17.63
MDC	MGR	1000	8.28	14.46	14.20	13.90
MDC	MSA	1000	15.75	21.05	19.34	20.40
Media	Gulf	467	14.22	22.68	22.76	19.60
Media	LEV	250	7.54	17.65	16.65	16.80
Media	MGR	526	4.87	11.58	10.20	9.60
Media	MSA	1258	20.66	35.34	33.57	32.65
Bible	MGR	1200	17.09	16.72	15.29	29.00
Bible	MSA	1200	22.91	22.08	17.53	31.20

Table 2: BLEU score on MT using zero-shot prompts. #Sent: number of test set sentences. SOTA results are reported in (Sajjad et al., 2020).

on the *effective prompt* and implementing appropriate *post-processing techniques*. Overall, these findings indicate the potential of GPT-4 as a *multi-task model* without heavily relying on task-specific resources, particularly in zero/few-shot settings.

The *few-shot results* across seven different datasets show an average improvement of 0.656 (0-shot) to 0.721 (10-shot) indicating the promise of few-shot learning, as depicted in Figure 2 (in Appendix), with individual results are reported in Table 9 (in Appendix).

The use of *native language prompts* with GPT-4 in a zero-shot context highlighted the role played by the prompt language, as we observed increased performance in three out of seven datasets compared to their counterparts with English prompts while two underperformed, and one showed equivalent performance (see Table 10 in Appendix).

When evaluating these LLMs in *multi-dialectal* settings, the performance gap between MSA and dialectal test sets becomes more evident. For example, in both the GPT-models, we noticed a large discrepancy in the POS accuracy of 0.810 versus 0.379 on MSA and dialects respectively. Similarly, for the dialect identification we notice a significant difference between the SOTA acoustic and lexical model with respect to LLMs results.

From the average *performance gap between semantic and syntactic tasks*, as reported in Table 11 (in Appendix), we noticed the discrepancy in semantic tasks is much lower than in syntactic tasks, across the three LLMs. This suggests that these models might be better equipped at encoding and expressing semantic information than in pinpointing specific syntactic phenomena in their inputs.

Dataset dom./dial.	Models	Zero-Shot	N-Shot (2hrs)	SOTA
MGB2 Broadcast/MSA	W.S	46.70	36.8	
	W.M	33.00	-	O: 11.4
	W.Lv2	26.20	18.8	S:11.9
	USM	15.70	N/A	
MGB3 Broadcast/EGY	W.S	83.20	77.5	
	W.M	65.90	-	O: 21.4
	W.Lv2	55.60	44.6	S: 26.70
	USM	22.10	N/A	
MGB5 Broadcast/MOR	W.S	135.20	114.6	
	W.M	116.90	-	O: 44.1
	W.Lv2	89.40	85.5	S:49.20
	USM	51.20	N/A	
QASR.CS Broadcast/Mixed	W.S	63.60	-	
	W.M	48.90	-	O: 23.4
	W.Lv2	37.90	31.2 ⁺	S: 24.90
	USM	27.80	N/A	
DACS Broadcast /MSA-EGY	W.S	61.90	-	
	W.M	48.70	-	O: 15.9
	W.Lv2	34.20	30.4 ⁺	S: 21.3
	USM	14.30	N/A	
ESCWA.CS Meeting/Mixed	W.S	101.50	-	
	W.M	69.30	-	O: 49.8
	W.Lv2	60.00	53.6 ⁺	S:48.00
	USM	45.70	N/A	
CallHome Telephony/EGY	W.S	155.90	152.9	
	W.M	113.70	-	O: 45.8*
	W.Lv2	78.70	64.6	S: 50.90
	USM	54.20	N/A	

Table 3: Reported WER (\downarrow) on ASR in zero and few-shot setup and domain-specific ASR setup. W.S,M,Lv2 stands for OpenAI Whisper small, medium and Largev2 model. O: represent offline; S: streaming ASR; * represent the model’s input is 8kHz sampling rate and Offline model was re-trained to accommodate telephony data. ⁺ represent model fine-tuned with 2hrs of MGB2-data.

Moreover, these performance gaps can also be linked to *undesirable hallucination*. In particular, during the MT for the Bible, results reveal an interesting phenomenon. It appears that the GPT models, particularly GPT-3.5-turbo, tend to hallucinate and insert additional content in their responses.

Is the data contaminated? We have used some datasets for evaluation that are released after the cut-off date of ChatGPT training, which include subjectivity, propaganda, check worthiness, factuality (CT-CWT-22), harmful content, and attention worthiness. Moreover, we experiment with nine datasets using the tailored instructions approach proposed by Golchin and Surdeanu (2023) revealing that GPT-4 could not produce any example from these datasets. Thus, we can confirm that the models have not been contaminated with such datasets. More details in Appendix F.5.

Speech Model Performances: We observed the performance of these models is heavily depen-

Model	Subjective (MOS) ↑			Objective ↓	
	Diac.	Natur.	Intel.	WER	CER
Amazon	8.2	8.3	9.8	5.2	1.0
KANARI	9.5	8.6	9.8	3.7	1.2

Table 4: Evaluation for Arabic TTS. Diac.: Diacritization, Natur.: Naturalness, Intel.: Intelligibility.

514 dependent on the architecture parameters. USM model
515 performs comparably with SOTA for MSA. Both
516 Whisper (and its variants) and USM show a per-
517 formance gap when dealing with dialects specially
518 Moroccan dialect. Fine-tuning the open model
519 (Whisper Largev2) with only 2 hours of speech
520 data bridges the performance gap significantly, in-
521 dicating the potential to be a robust and strong
522 foundation model. Our observation also suggests
523 that USM model is better equipped to handle code-
524 switching phenomena in spoken utterance than the
525 supervised large transformer models.

526 6 Related Work

527 **Models for NLP:** Since the inception of the trans-
528 former architecture (Vaswani et al., 2017), there
529 have been efforts to develop larger models with
530 its variants such as BERT (Devlin et al., 2019),
531 RoBERTa (Liu et al., 2019), XLM-RoBERTa (Con-
532 neau et al., 2020), GPT models (Radford et al.,
533 2018, 2019; Ouyang et al., 2022) among others.

534 Such advancements have led to the development
535 LLMs with parameter sizes exceeding 100 bil-
536 lion, which are pre-trained on massive datasets.
537 Examples of LLMs include Megatron (Shoeybi
538 et al., 2019), GPT-3 (Brown et al., 2020), GPT-
539 Jurassic (Lieber et al., 2021), OPT-175B (Zhang
540 et al., 2022), and Bloom (Scao et al., 2022). This
541 unprecedented scale enabled new capabilities that
542 address the zero-shot and multilingual tasks learn-
543 ing. ChatGPT (GPT-3.5) and its subsequent model
544 GPT-4 is the latest development in NLP that have
545 addressed many limitations of prior LLMs and en-
546 abled us to perform diverse tasks (OpenAI, 2023).
547 The ability of LLMs to solve various tasks can be at-
548 tributed to the meticulous design of prompts, which
549 enable the generation of desired responses (Wei
550 et al., 2022; Shin et al., 2020).

551 **Models for Speech Processing:** When handling
552 complex audio/speech data, LLMs face signifi-
553 cant challenges. However, with the advent of
554 self-supervised learning, models like Wav2vec,
555 WavLM, and Whisper have been leading in address-
556 ing these challenges (Baevski et al., 2019, 2020;
557 Chen et al., 2022; Radford et al., 2022). More re-

cent developments like the Universal Speech Model
(USM) and VALL-E have demonstrated superior
capabilities in ASR and zero-shot TTS tasks, re-
spectively (Zhang et al., 2023; Wang et al., 2023).

558 **LLMs Benchmarking:** Since the release of Chat-
559 GPT, there have been efforts to evaluate the perfor-
560 mance of LLMs on standard NLP tasks (Bubeck
561 et al., 2023; Bang et al., 2023; Ahuja et al., 2023;
562 Hendy et al., 2023). Liang et al. (2022) conducted
563 a comprehensive assessment of LLMs for English.
564 It encompassed various metrics such as accuracy,
565 calibration, toxicity, and efficiency, along with 42
566 scenarios involving 30 prominent language models.
567 **Benchmarks on Arabic:** The complexity and lin-
568 guistic diversity of Arabic have led to a limited
569 number of benchmarks for language tasks, such as
570 ORCA (Elmadany et al., 2022), ALUE (Seelawi
571 et al., 2021), ArBERT (Abdul-Mageed et al., 2021),
572 and AraBench (Sajjad et al., 2020).

573 **LAraBench:** To the best of our knowledge, our
574 study represents the first comprehensive Arabic
575 language benchmarking effort exploring GPT-3.5
576 (zero-shot), GPT-4 (zero- and few-shot), BLOOMZ
577 (zero-shot), and Speech models like Whisper and
578 USM. Our evaluation spans a broad array of LLMs,
579 tasks, and datasets, distinguishing it from prior
580 benchmarks in terms of task and dataset diversity,
581 test setup, modalities (text, speech), and state-of-
582 the-art comparisons. Table 8 (Appendix E), pro-
583 vides a detailed comparison.

584 7 Conclusion and Future Studies

585 This study is the *first* large-scale benchmark that
586 brings together both Arabic Speech and NLP tasks
587 under the same study. We report the performance
588 of LLMs for a variety of tasks covering different
589 domains and dialects. Our study also considers
590 tasks with a wide range of complexity ranging from
591 token to text classification, different application
592 settings, NER to sentiment, factuality and disinforma-
593 tion, ASR, and TTS among others. We evaluate
594 33 tasks and 61 datasets with 98 test setups, which
595 are very prominent for Arabic AI. We compare and
596 report the performance of each task and dataset
597 with SOTA, which will enable the community and
598 practitioners of large language models to decide on
599 their uses of these models. Future work aims to in-
600 vestigate open models and explore ways to reduce
601 the performance gap with SOTA; enhance prompts
602 for better performance; and expand datasets and
603 tasks studied.

608 Limitations

609 The main focus of this study was to benchmark
610 large language models for Arabic NLP and Speech
611 tasks. Given that this is a work in progress, there
612 are currently some limitations. In this edition, we
613 evaluated several large models: ChatGPT, USM,
614 and Whisper models and compared them to SOTA.
615 We plan to extend our study by adding other models
616 such Bard, Claude, MMS, and other open multi-
617 lingual models that have Arabic. In this work, we
618 benchmarked 61 datasets with 98 test setups for 33
619 tasks. However, we did not benchmark all avail-
620 able data sets. For example, the study reported in
621 (Elmadany et al., 2022) benchmarked 19 sentiment
622 datasets, whereas we only covered one. It is also
623 possible that we missed many other Arabic NLP
624 and Speech tasks, which we will attempt to cover in
625 the future. Our current results are highly dependent
626 on prompt design. Additional efforts on prompt
627 engineering could potentially improve the results.

628 In addition, performance may vary depending
629 on the version of the models we used.⁹ For GPTs,
630 we utilized gpt-3.5-turbo-0301 and gpt-4-0314 ver-
631 sions for our NLP tasks. To ensure transparency
632 and reproducibility, we are committed to sharing all
633 our experimental resources, including prompts and
634 parameter details. This will facilitate the easy repli-
635 cation of our results using the provided pipeline
636 and the fixed model versions. The same princi-
637 ple extends to our speech models. We have taken
638 steps to maintain versioning not only for the models
639 themselves but also for the prompts used. This en-
640 sures that our work remains reproducible for future
641 researchers in the field.

642 **Potential Risk** We do not oversee any potential
643 risk that can result from our study.

644 Ethics Statement

645 We used publicly available and in-house developed
646 datasets in our study. Any biases are unintended.

647 References

648 Ahmed Abdelali, Mohammed Attia, Younes Samih, Ka-
649 reem Darwish, and Hamdy Mubarak. 2019. *Dia-*
650 *critization of maghrebi Arabic sub-dialects*. *arXiv*
651 *preprint arXiv:1810.06619*.

652 Ahmed Abdelali, Hamdy Mubarak, Younes Samih,
653 Sabit Hassan, and Kareem Darwish. 2020. Ara-

654 bic dialect identification in the wild. *arXiv preprint*
655 *arXiv:2005.06557*.

Ahmed Abdelali, Hamdy Mubarak, Younes Samih,
656 Sabit Hassan, and Kareem Darwish. 2021. *QADI:*
657 *Arabic dialect identification in the wild*. In *Proceed-*
658 *ings of the Sixth Arabic Natural Language Process-*
659 *ing Workshop*, pages 1–10, Kyiv, Ukraine (Virtual).
660 Association for Computational Linguistics. 661

Muhammad Abdul-Mageed, AbdelRahim Elmadany,
662 et al. 2021. Arbert & marbert: Deep bidirectional
663 transformers for Arabic. In *Proceedings of the 59th*
664 *Annual Meeting of the Association for Computational*
665 *Linguistics and the 11th International Joint Confer-*
666 *ence on Natural Language Processing (Volume 1:*
667 *Long Papers)*, pages 7088–7105. 668

Ibrahim Abu Farha and Walid Magdy. 2020. *From*
669 *Arabic sentiment analysis to sarcasm detection: The*
670 *ArSarcasm dataset*. In *Proceedings of the 4th Work-*
671 *shop on Open-Source Arabic Corpora and Process-*
672 *ing Tools, with a Shared Task on Offensive Language*
673 *Detection*, pages 32–39, Marseille, France. European
674 Language Resource Association. 675

Ibrahim Abu Farha, Wajdi Zaghouni, and Walid Magdy.
676 2021. Overview of the wanlp 2021 shared task on
677 sarcasm and sentiment detection in arabic. In *Pro-*
678 *ceedings of the Sixth Arabic Natural Language Pro-*
679 *cessing Workshop*. 680

Kabir Ahuja, Rishav Hada, Millicent Ochieng, Prachi
681 Jain, Harshita Diddee, Samuel Maina, Tanuja Ganu,
682 Sameer Segal, Maxamed Axmed, Kalika Bali, et al.
683 2023. MEGA: Multilingual evaluation of generative
684 ai. *arXiv preprint arXiv:2303.12528*. 685

Firoj Alam, Stefano Cresci, Tanmoy Chakraborty, Fab-
686 rizio Silvestri, Dimiter Dimitrov, Giovanni Da San
687 Martino, Shaden Shaar, Hamed Firooz, and Preslav
688 Nakov. 2022a. *A survey on multimodal disinfor-*
689 *mation detection*. In *Proceedings of the 29th Inter-*
690 *national Conference on Computational Linguistics,*
691 *COLING '22*, pages 6625–6643, Gyeongju, Republic
692 of Korea. 693

Firoj Alam, Fahim Dalvi, Shaden Shaar, Nadir Durrani,
694 Hamdy Mubarak, Alex Nikolov, Giovanni Da San
695 Martino, Ahmed Abdelali, Hassan Sajjad, Kareem
696 Darwish, and Preslav Nakov. 2021a. Fighting the
697 COVID-19 infodemic in social media: A holistic
698 perspective and a call to arms. In *Proceedings of the*
699 *International AAAI Conference on Web and Social*
700 *Media, ICWSM '21*, pages 913–922. 701

Firoj Alam, Hamdy Mubarak, Wajdi Zaghouni, Gio-
702 vanni Da San Martino, and Preslav Nakov. 2022b.
703 *Overview of the WANLP 2022 shared task on propa-*
704 *ganda detection in Arabic*. pages 108–118. 705

Firoj Alam, Shaden Shaar, Fahim Dalvi, Hassan Saj-
706 jad, Alex Nikolov, Hamdy Mubarak, Giovanni
707 Da San Martino, Ahmed Abdelali, Nadir Durrani,
708 Kareem Darwish, Abdulaziz Al-Homaid, Wajdi Za-
709 ghouni, Tommaso Caselli, Gijs Danoe, Friso Stolk,
710 Britt Bruntink, and Preslav Nakov. 2021b. *Fighting*
711 *the COVID-19 infodemic: Modeling the perspective*
712 *of journalists, fact-checkers, social media platforms,*
713 *policy makers, and the society*. In *Findings of the*
714

⁹<https://platform.openai.com/docs/models/overview>

715	<i>Association for Computational Linguistics: EMNLP 2021</i> , pages 611–649, Punta Cana, Dominican Republic. Association for Computational Linguistics.	774
716		775
717		776
718	Abeer ALDayel and Walid Magdy. 2021. Stance detection on social media: State of the art and trends . In <i>Information Processing & Management</i> , 58(4):102597.	777
719		778
720		779
721	Abdullah I Alharbi and Mark Lee. 2021. Multi-task learning using a combination of contextualised and static word embeddings for arabic sarcasm detection and sentiment analysis. In <i>Proceedings of the Sixth Arabic Natural Language Processing Workshop</i> , pages 318–322.	780
722		781
723		782
724		783
725		784
726		785
727	Ahmed Ali, Peter Bell, James Glass, Yacine Messaoui, Hamdy Mubarak, Steve Renals, and Yifan Zhang. 2016. The MGB-2 challenge: Arabic multi-dialect broadcast media recognition. In <i>2016 IEEE Spoken Language Technology Workshop (SLT)</i> , pages 279–284. IEEE.	786
728		787
729		788
730		789
731		790
732		791
733	Ahmed Ali, Shammur Chowdhury, Amir Hussein, and Yasser Hifny. 2021a. Arabic code-switching speech recognition using monolingual data. <i>arXiv preprint arXiv:2107.01573</i> .	792
734		793
735		794
736		795
737	Ahmed Ali, Suwon Shon, Younes Samih, Hamdy Mubarak, Ahmed Abdelali, James Glass, Steve Renals, and Khalid Choukri. 2019. The MGB-5 challenge: Recognition and dialect identification of dialectal Arabic speech. In <i>2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)</i> , pages 1026–1033. IEEE.	796
738		797
739		798
740		799
741		800
742		801
743		802
744	Ahmed Ali, Stephan Vogel, and Steve Renals. 2017. Speech recognition challenge in the wild: Arabic MGB-3. In <i>2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)</i> , pages 316–322. IEEE.	803
745		804
746		805
747		806
748		807
749	Zien Sheikh Ali, Watheq Mansour, Tamer Elsayed, and Abdulaziz Al-Ali. 2021b. Arafacts: the first large Arabic dataset of naturally occurring claims. In <i>Proceedings of the Sixth Arabic Natural Language Processing Workshop</i> , pages 231–236.	808
750		809
751		810
752		811
753		812
754	Francesco Antici, Luca Bolognini, Matteo Antonio Inajetovic, Bogdan Ivasiuk, Andrea Galassi, and Federico Ruggeri. 2021. Subjectivita: An italian corpus for subjectivity detection in newspapers. In <i>Experimental IR Meets Multilinguality, Multimodality, and Interaction</i> , pages 40–52, Cham. Springer International Publishing.	813
755		814
756		815
757		816
758		817
759		818
760		819
761	Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020. On the cross-lingual transferability of monolingual representations. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4623–4637.	820
762		821
763		822
764		823
765		824
766	Stephen Bach, Victor Sanh, Zheng Xin Yong, Albert Webson, Colin Raffel, Nihal V Nayak, Abheesht Sharma, Taewoon Kim, M Saiful Bari, Thibault Févry, et al. 2022. Promptsources: An integrated development environment and repository for natural language prompts. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: System Demonstrations</i> , pages 93–104.	825
767		826
768		827
769		828
770		829
771		830
772		831
773		832
	Alexei Baeovski, Steffen Schneider, and Michael Auli. 2019. vq-wav2vec: Self-supervised learning of discrete speech representations. <i>arXiv preprint arXiv:1910.05453</i> .	833
		834
	Alexei Baeovski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. <i>Advances in neural information processing systems</i> , 33:12449–12460.	
	Ramy Baly, Georgi Karadzhov, Dimitar Alexandrov, James Glass, and Preslav Nakov. 2018a. Predicting factuality of reporting and bias of news media sources . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3528–3539, Brussels, Belgium. Association for Computational Linguistics.	
	Ramy Baly, Mitra Mohtarami, James Glass, Lluís Màrquez, Alessandro Moschitti, and Preslav Nakov. 2018b. Integrating stance detection and fact checking in a unified corpus. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> .	
	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. <i>arXiv preprint arXiv:2302.04023</i> .	
	Alberto Barrón-Cedeño, Firoj Alam, Tommaso Caselli, Giovanni Da San Martino, Tamer Elsayed, Andrea Galassi, Fatima Haouari, Federico Ruggeri, Julia Maria Struß, Rabindra Nath Nandi, et al. 2023. The clef-2023 checkthat! lab: Checkworthiness, subjectivity, political bias, factuality, and authority. In <i>Advances in Information Retrieval: 45th European Conference on Information Retrieval, ECIR 2023, Dublin, Ireland, April 2–6, 2023, Proceedings, Part III</i> , pages 506–517. Springer.	
	Yassine Benajiba and Paolo Rosso. 2007. ANERsys 2.0: Conquering the NER task for the Arabic language by combining the maximum entropy with pos-tag information. In <i>IJCAI</i> , pages 1814–1823.	
	Yassine Benajiba, Paolo Rosso, and José Miguel BenedíRuiz. 2007. ANERsys: An arabic named entity recognition system based on maximum entropy. In <i>Computational Linguistics and Intelligent Text Processing</i> , pages 143–153, Berlin, Heidelberg. Springer Berlin Heidelberg.	
	Douglas Biber and Edward Finegan. 1988. Adverbial stance types in english. <i>Discourse processes</i> , 11(1):1–34.	
	Taboubi Bilel, Ben Nessir Mohamed Aziz, and Hatem Haddad. 2022. iCompass at CheckThat! 2022: ARBERT and AraBERT for Arabic checkworthy tweet identification. In <i>Working Notes of CLEF 2022 - Conference and Labs of the Evaluation Forum, CLEF '2022</i> , Bologna, Italy.	
	Houda Bouamor, Nizar Habash, and Kemal Oflazer. 2014. A multidialectal parallel corpus of Arabic. In <i>LREC</i> , pages 1240–1245.	

835	Houda Bouamor, Nizar Habash, Mohammad Salameh, Wajdi Zaghouani, Owen Rambow, Dana Abdulrahim, Ossama Obeid, Salam Khalifa, Fadhil Eryani, Alexander Erdmann, et al. 2018. The MADAR Arabic Dialect Corpus and Lexicon. In <i>LREC</i> .	895
836		896
837		897
838		898
839		899
840	Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference . In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.	900
841		901
842		902
843		903
844		904
845		905
846		906
847	David A Broniatowski, Amelia M Jamison, SiHua Qi, Lulwah AlKulaib, Tao Chen, Adrian Benton, Sandra C Quinn, and Mark Dredze. 2018. Weaponized health communication: Twitter bots and Russian trolls amplify the vaccine debate . <i>American journal of public health</i> , 108(10):1378–1384.	907
848		908
849		909
850		910
851		911
852		912
853	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners . <i>Advances in Neural Information Processing Systems</i> .	913
854		914
855		915
856		916
857		917
858		918
859		919
860		920
861		921
862		922
863		923
864		924
865	Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. 2023. Sparks of artificial general intelligence: Early experiments with GPT-4 . Technical report, Microsoft Research.	925
866		926
867		927
868		928
869		929
870		930
871		931
872	Sabine Buchholz and Erwin Marsi. 2006. CoNLL-X shared task on multilingual dependency parsing. In <i>Proceedings of the tenth conference on computational natural language learning (CoNLL-X)</i> , pages 149–164.	932
873		933
874		934
875		935
876		936
877	Jaime Carbonell and Jade Goldstein. 1998. The use of mmr, diversity-based reranking for reordering documents and producing summaries. In <i>Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval</i> , pages 335–336.	937
878		938
879		939
880		940
881		941
882		942
883	Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017a. SemEval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation . In <i>Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)</i> , pages 1–14, Vancouver, Canada. Association for Computational Linguistics.	943
884		944
885		945
886		946
887		947
888		948
889		949
890	Daniel Cer, Mona Diab, Eneko E. Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017b. SemEval-2017 Task 1: Semantic Textual Similarity Multilingual and Cross-lingual Focused Evaluation . In <i>Proceedings of the 11th International Workshop on Semantic Evalu-</i>	950
891		951
892		952
893		953
894		954
	<i>ation (SemEval-2017)</i> , SemEval-2017, pages 1–14, Vancouver, Canada.	955
	Gullal Singh Cheema, Sherzod Hakimov, Abdul Sittar, Eric Müller-Budack, Christian Otto, and Ralph Ewerth. 2022. MM-claims: A dataset for multimodal claim detection in social media . In <i>Findings of the Association for Computational Linguistics: NAACL 2022</i> , pages 962–979, Seattle, United States. Association for Computational Linguistics.	956
	Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, et al. 2022. Wavlm: Large-scale self-supervised pre-training for full stack speech processing. <i>IEEE Journal of Selected Topics in Signal Processing</i> , 16(6):1505–1518.	957
	Amina Chouigui, Oussama Ben Khiroun, and Bilel Elayeb. 2017. ANT corpus: An Arabic news text collection for textual classification. In <i>2017 IEEE/ACS 14th International Conference on Computer Systems and Applications (AICCSA)</i> , pages 135–142. IEEE.	958
	Shammur A Chowdhury, Younes Samih, Mohamed El-desouki, and Ahmed Ali. 2020a. Effects of dialectal code-switching on speech modules: A study using egyptian Arabic broadcast speech. <i>Proc. Interspeech</i> .	959
	Shammur Absar Chowdhury, Ahmed Abdelali, Kareem Darwish, Jung Soon-Gyo, Joni Salminen, and Bernard J Jansen. 2020b. Improving Arabic text categorization using transformer training diversification. In <i>Proceedings of the fifth arabic natural language processing workshop</i> , pages 226–236.	960
	Shammur Absar Chowdhury, Amir Hussein, Ahmed Abdelali, and Ahmed Ali. 2021. Towards one model to rule all: Multilingual strategy for dialectal code-switching Arabic asr. <i>Proc. Interspeech</i> .	961
	Shammur Absar Chowdhury, Hamdy Mubarak, Ahmed Abdelali, Soon-gyo Jung, Bernard J Jansen, and Joni Salminen. 2020c. A multi-platform arabic news comment dataset for offensive language detection. In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> , pages 6203–6212.	962
	Jonathan H. Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. 2020. Tydi qa: A benchmark for information-seeking question answering in typologically diverse languages. <i>Transactions of the Association for Computational Linguistics</i> .	963
	Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL ’20</i> , pages 8440–8451, Online. Association for Computational Linguistics.	964
	Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. XNLI: Evaluating cross-lingual sentence representations. In <i>Proceedings of</i>	965

954			
955		<i>the 2018 Conference on Empirical Methods in Natural Language Processing, EMNLP '18</i> , pages 2475–2485.	
956			
957	Kareem Darwish, Ahmed Abdelali, Hamdy Mubarak, Younes Samih, and Mohammed Attia. 2018. Diacritization of moroccan and tunisian Arabic dialects: A crf approach. <i>OSACT</i> , 3:62.		
958			
959			
960			
961	Kareem Darwish, Dimitar Alexandrov, Preslav Nakov, and Yelena Mejova. 2017a. Seminar users in the Arabic twitter sphere. In <i>Social Informatics: 9th International Conference, SocInfo 2017, Oxford, UK, September 13-15, 2017, Proceedings, Part I 9</i> , pages 91–108. Springer.		
962			
963			
964			
965			
966			
967	Kareem Darwish and Hamdy Mubarak. 2016. Farasa: A new fast and accurate Arabic word segmenter. In <i>Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)</i> , pages 1070–1074.		
968			
969			
970			
971			
972	Kareem Darwish, Hamdy Mubarak, and Ahmed Abdelali. 2017b. Arabic diacritization: Stats, rules, and hacks . In <i>Proceedings of the Third Arabic Natural Language Processing Workshop</i> , pages 9–17, Valencia, Spain. Association for Computational Linguistics.		
973			
974			
975			
976			
977			
978	Kareem Darwish, Hamdy Mubarak, Ahmed Abdelali, and Mohamed Eldesouki. 2017c. Arabic POS tagging: Don't abandon feature engineering just yet. In <i>Proceedings of the third arabic natural language processing workshop</i> , pages 130–137.		
979			
980			
981			
982			
983	Thomas Davidson, Dana Warmusley, Michael Macy, and Ingmar Weber. 2017. Automated hate speech detection and the problem of offensive language. <i>Proceedings of the International AAAI Conference on Web and Social Media</i> , 11(1):512–515.		
984			
985			
986			
987			
988	Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo Ko, Alan Cowen, Gaurav Nemade, and Sujith Ravi. 2020. GoEmotions: A dataset of fine-grained emotions . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4040–4054, Online. Association for Computational Linguistics.		
989			
990			
991			
992			
993			
994			
995	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4171–4186.		
996			
997			
998			
999			
1000			
1001			
1002			
1003	Dimitar Dimitrov, Bishr Bin Ali, Shaden Shaar, Firoj Alam, Fabrizio Silvestri, Hamed Firooz, Preslav Nakov, and Giovanni Da San Martino. 2021. Detecting propaganda techniques in memes . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 6603–6617, Online. Association for Computational Linguistics.		
1004			
1005			
1006			
1007			
1008			
1009			
1010			
1011			
1012	Mingzhe Du, Sujatha Das Gollapalli, and See-Kiong Ng. 2022. Nus-ids at checkthat! 2022: identifying	check-worthiness of tweets using checkthat5. <i>Working Notes of CLEF</i> .	1014 1015
1013			
		Omar Einea, Ashraf Elnagar, and Ridhwan Al Debsi. 2019. SANAD: Single-label Arabic news articles dataset for automatic text categorization. <i>Data in brief</i> , 25:104076.	1016 1017 1018 1019
		Paul Ekman. 1971. Universals and cultural differences in facial expressions of emotion. In <i>Nebraska symposium on motivation</i> . University of Nebraska Press.	1020 1021 1022
		AbdelRahim Elmadany, El Moatez Billah Nagoudi, and Muhammad Abdul-Mageed. 2022. Orca: A challenging benchmark for Arabic language understanding. <i>arXiv preprint arXiv:2212.10758</i> .	1023 1024 1025 1026
		AbdelRahim A. Elmadany, Hamdy Mubarak, and Walid Magdy. 2018. ArSAS: An Arabic speech-act and sentiment corpus of tweets. <i>OSACT</i> , 3:20.	1027 1028 1029
		Ashraf Elnagar, Ridhwan Al-Debsi, and Omar Einea. 2020. Arabic text classification using deep learning models. <i>Information Processing & Management</i> , 57(1):102121.	1030 1031 1032 1033
		A. Etman and A. A. Louis Beex. 2015. Language and dialect identification: A survey . In <i>2015 SAI Intelligent Systems Conference (IntelliSys)</i> , pages 220–231.	1034 1035 1036
		Ahmet Bahadir Eyuboglu, Mustafa Bora Arslan, Ekrem Sonmezer, and Mucahid Kutlu. 2022. TOBB ETU at CheckThat! 2022: detecting attention-worthy and harmful tweets and check-worthy claims. In <i>Working Notes of CLEF 2022 - Conference and Labs of the Evaluation Forum, CLEF '2022, Bologna, Italy</i> .	1037 1038 1039 1040 1041 1042
		Ibrahim Abu Farha and Walid Magdy. 2020. From arabic sentiment analysis to sarcasm detection: The arsarcasm dataset. In <i>Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools, with a Shared Task on Offensive Language Detection</i> , pages 32–39.	1043 1044 1045 1046 1047 1048
		Hongyu Gao, Yan Chen, Kathy Lee, Diana Palsetia, and Alok N. Choudhary. 2012. Towards online spam filtering in social networks. In <i>Network and Distributed System Security Symposium, NDSS '12</i> , pages 1–16.	1049 1050 1051 1052
		Razan Ghanem, Hasan Erbay, and Khaled Bakour. 2023. Contents-based spam detection on social networks using roberta embedding and stacked blstm. <i>SN Computer Science</i> , 4(4):380.	1053 1054 1055 1056
		Shahriar Golchin and Mihai Surdeanu. 2023. Time travel in llms: Tracing data contamination in large language models. <i>arXiv preprint arXiv:2308.08493</i> .	1057 1058 1059
		Raul Gomez, Jaume Gibert, Lluís Gomez, and Dimosthenis Karatzas. 2020. Exploring hate speech detection in multimodal publications. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer vision</i> , pages 1470–1478.	1060 1061 1062 1063 1064
		Mourad Gridach. 2018. Deep learning approach for arabic named entity recognition. In <i>Computational Linguistics and Intelligent Text Processing: 17th International Conference, CICLing 2016, Konya, Turkey, April 3–9, 2016, Revised Selected Papers, Part I 17</i> , pages 439–451. Springer.	1065 1066 1067 1068 1069 1070

1071	Jan Hajic, Otakar Smrz, Petr Zemánek, Jan Šnidauf, and Emanuel Beška. 2004. Prague Arabic dependency treebank: Development in data and tools. In <i>Proc. of the NEMLAR Intern. Conf. on Arabic Language Resources and Tools</i> , volume 1.	1130
1072		1131
1073		1132
1074		1133
1075		1134
1076	Sabit Hassan, Hamdy Mubarak, Ahmed Abdelali, and Kareem Darwish. 2021. Asad: Arabic social media analytics and understanding. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations</i> , pages 113–118.	1135
1077		1136
1078		1137
1079		1138
1080		1139
1081		1140
1082	Sabit Hassan, Shaden Shaar, and Kareem Darwish. 2022. Cross-lingual emotion detection. In <i>Proceedings of the Thirteenth Language Resources and Evaluation Conference</i> , pages 6948–6958.	1141
1083		1142
1084		1143
1085		1144
1086	Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are GPT models at machine translation? a comprehensive evaluation. <i>arXiv preprint arXiv:2302.09210</i> .	1145
1087		1146
1088		1147
1089		1148
1090		1149
1091		1150
1092	Huang Huang, Fei Yu, Jianqing Zhu, Xuening Sun, Hao Cheng, Dingjie Song, Zhihong Chen, Abdulmohsen Alharthi, Bang An, Ziche Liu, et al. 2023. Acegpt, localizing large language models in arabic. <i>arXiv preprint arXiv:2309.12053</i> .	1151
1093		1152
1094		1153
1095		1154
1096		1155
1097	Fatemah Husain and Ozlem Uzuner. 2021. A survey of offensive language detection for the Arabic language. <i>ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)</i> , 20(1):1–44.	1156
1098		1157
1099		1158
1100		1159
1101	Md Tawkat Islam Khondaker, Abdul Waheed, El Moatez Billah Nagoudi, and Muhammad Abdul-Mageed. 2023. GPTAraEval: A comprehensive evaluation of chatgpt on arabic nlp. <i>arXiv preprint arXiv:2305.14976</i> .	1160
1102		1161
1103		1162
1104		1163
1105		1164
1106	Jude Khouja. 2020. Stance prediction and claim verification: An Arabic perspective . In <i>Proceedings of the Third Workshop on Fact Extraction and VERification (FEVER)</i> , pages 8–17, Online. Association for Computational Linguistics.	1165
1107		1166
1108		1167
1109		1168
1110		1169
1111	Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Pratik Ringshia, and Davide Testuggine. 2020. The hateful memes challenge: Detecting hate speech in multimodal memes . In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 2611–2624.	1170
1112		1171
1113		1172
1114		1173
1115		1174
1116		1175
1117	Lev Konstantinovskiy, Oliver Price, Mevan Babakar, and Arkaitz Zubiaga. 2021. Toward automated factchecking: Developing an annotation schema and benchmark for consistent automated claim detection . <i>Digital Threats: Research and Practice</i> , 2(2).	1176
1118		1177
1119		1178
1120		1179
1121		1180
1122	Dilek Küçük and Fazli Can. 2020. Stance detection: A survey. <i>ACM Computing Surveys (CSUR)</i> , 53(1):1–37.	1181
1123		1182
1124		1183
1125	Gaurav Kumar, Yuan Cao, Ryan Cotterell, Chris Callison-Burch, Daniel Povey, and Sanjeev Khudanpur. 2014. Translations of the callhome Egyptian Arabic corpus for conversational speech translation . In <i>Proceedings of the 11th International Workshop on Spoken Language Translation: Papers</i> , pages 244–248, Lake Tahoe, California.	1184
1126		1185
1127		1186
1128		1187
1129		1188
		1189
	Tao Lei, Yu Xin, Yuan Zhang, Regina Barzilay, and Tommi Jaakkola. 2014. Low-rank tensors for scoring dependency structures. In <i>Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1381–1391.	1190
		1191
	Patrick Lewis, Barlas Oğuz, Ruty Rinott, Sebastian Riedel, and Holger Schwenk. 2019. Mlqa: Evaluating cross-lingual extractive question answering. <i>arXiv preprint arXiv:1910.07475</i> .	1192
		1193
	Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yan Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. 2022. Holistic evaluation of language models. <i>arXiv preprint arXiv:2211.09110</i> .	1194
		1195
	Yaobo Liang, Nan Duan, Yeyun Gong, Ning Wu, Fenfei Guo, Weizhen Qi, Ming Gong, Linjun Shou, Daxin Jiang, Guihong Cao, Xiaodong Fan, Ruofei Zhang, Rahul Agrawal, Edward Cui, Sining Wei, Taroan Bharti, Ying Qiao, Jiun-Hung Chen, Winnie Wu, Shuguang Liu, Fan Yang, Daniel Campos, Rangan Majumder, and Ming Zhou. 2020a. XGLUE: A new benchmark dataset for cross-lingual pre-training, understanding and generation. <i>arXiv</i> , abs/2004.01401.	1196
		1197
	Yaobo Liang, Nan Duan, Yeyun Gong, Ning Wu, Fenfei Guo, Weizhen Qi, Ming Gong, Linjun Shou, Daxin Jiang, Guihong Cao, et al. 2020b. Xglue: A new benchmark datasetfor cross-lingual pre-training, understanding and generation. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 6008–6018.	1198
		1199
	Opher Lieber, Or Sharir, Barak Lenz, and Yoav Shoham. 2021. Jurassic-1: Technical details and evaluation. <i>White Paper. AI21 Labs</i> , 1.	1200
		1201
	Bing Liu and Lei Zhang. 2012. A survey of opinion mining and sentiment analysis. In <i>Mining text data</i> , pages 415–463. Springer.	1202
		1203
	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	1204
		1205
	Clyde R. Miller. 1939. The Techniques of Propaganda. From “How to Detect and Analyze Propaganda,” an address given at Town Hall. The Center for learning.	1206
		1207
	Saif Mohammad, Felipe Bravo-Marquez, Mohammad Salameh, and Svetlana Kiritchenko. 2018. Semeval-2018 task 1: Affect in tweets. In <i>Proceedings of the 12th international workshop on semantic evaluation</i> , pages 1–17.	1208
		1209
	Hussein Mozannar, Karl El Hajal, Elie Maamary, and Hazem Hajj. 2019. Neural Arabic question answering. <i>arXiv preprint arXiv:1906.05394</i> .	1210
		1211
	Hamdy Mubarak. 2018. Build fast and accurate lemmatization for Arabic . In <i>Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)</i> , Miyazaki, Japan. European Language Resources Association (ELRA).	1212
		1213
		1214
		1215
		1216
		1217
		1218
		1219

1190	Hamdy Mubarak, Ahmed Abdelali, Sabit Hassan, and Kareem Darwish. 2020a. Spam detection on Arabic twitter. In <i>Social Informatics: 12th International Conference, SocInfo 2020, Pisa, Italy, October 6–9, 2020, Proceedings 12</i> , pages 237–251. Springer.	1251
1191		1252
1192		1253
1193		1254
1194		1255
1195	Hamdy Mubarak, Ahmed Abdelali, Hassan Sajjad, Younes Samih, and Kareem Darwish. 2019. Highly effective Arabic diacritization using sequence to sequence modeling. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 2390–2395.	1256
1196		1257
1197		1258
1198		1259
1199		1260
1200		1261
1201		1262
1202		1263
1203	Hamdy Mubarak, Shammur Absar Chowdhury, and Firoj Alam. 2022. ArabGend: Gender analysis and inference on Arabic Twitter . In <i>Proceedings of the Eighth Workshop on Noisy User-generated Text (W-NUT 2022)</i> , pages 124–135, Gyeongju, Republic of Korea. Association for Computational Linguistics.	1264
1204		1265
1205		1266
1206		1267
1207		1268
1208		1269
1209	Hamdy Mubarak, Kareem Darwish, Walid Magdy, Tamer Elsayed, and Hend Al-Khalifa. 2020b. Overview of OSACT4 Arabic offensive language detection shared task . In <i>Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools, with a Shared Task on Offensive Language Detection</i> , pages 48–52, Marseille, France. European Language Resource Association.	1270
1210		1271
1211		1272
1212		1273
1213		1274
1214		1275
1215		1276
1216		1277
1217	Hamdy Mubarak, Kareem Darwish, Walid Magdy, Tamer Elsayed, and Hend Al-Khalifa. 2020c. Overview of osact4 Arabic offensive language detection shared task. In <i>Proceedings of the 4th Workshop on open-source arabic corpora and processing tools, with a shared task on offensive language detection</i> , pages 48–52.	1278
1218		1279
1219		1280
1220		1281
1221		1282
1222		1283
1223		1284
1224	Hamdy Mubarak and Sabit Hassan. 2021. UI2c: Mapping user locations to countries on Arabic twitter. In <i>Proceedings of the Sixth Arabic Natural Language Processing Workshop</i> , pages 145–153.	1285
1225		1286
1226		1287
1227		1288
1228	Hamdy Mubarak, Sabit Hassan, and Ahmed Abdelali. 2021a. Adult content detection on Arabic twitter: Analysis and experiments. In <i>Proceedings of the Sixth Arabic Natural Language Processing Workshop</i> , pages 136–144.	1289
1229		1290
1230		1291
1231		1292
1232		1293
1233	Hamdy Mubarak, Amir Hussein, Shammur Absar Chowdhury, and Ahmed Ali. 2021b. QASR: QCRI Aljazeera speech resource—a large scale annotated Arabic speech corpus. <i>arXiv preprint arXiv:2106.13000</i> .	1294
1234		1295
1235		1296
1236		1297
1237		1298
1238		1299
1239	Mahmoud Nabil, Mohamed Aly, and Amir Atiya. 2015. Astd: Arabic sentiment tweets dataset. In <i>Proceedings of the 2015 conference on empirical methods in natural language processing</i> , pages 2515–2519.	1300
1240		1301
1241		1302
1242	Preslav Nakov, Alberto Barrón-Cedeño, Giovanni Da San Martino, Firoj Alam, Rubén Míguez, Tommaso Caselli, Mucahid Kutlu, Wajdi Zaghouni, Chengkai Li, Shaden Shaar, Hamdy Mubarak, Alex Nikolov, Yavuz Selim Kartal, and Javier Beltrán. 2022a. Overview of the CLEF-2022 CheckThat! lab task 1 on identifying relevant claims in tweets. In <i>Working Notes of CLEF 2022—Conference and Labs of the Evaluation Forum, CLEF ’2022</i> .	1303
1243		1304
1244		1305
1245		1306
1246		1307
1247		1308
1248		1309
1249		1310
1250		1310
	Preslav Nakov, Alberto Barrón-Cedeño, Giovanni da San Martino, Firoj Alam, Julia Maria Struß, Thomas Mandl, Rubén Míguez, Tommaso Caselli, Mucahid Kutlu, Wajdi Zaghouni, et al. 2022b. Overview of the clef-2022 checkthat! lab on fighting the covid-19 infodemic and fake news detection. In <i>Experimental IR Meets Multilinguality, Multimodality, and Interaction: 13th International Conference of the CLEF Association, CLEF 2022, Bologna, Italy, September 5–8, 2022, Proceedings</i> , pages 495–520. Springer.	1251
		1252
		1253
		1254
		1255
		1256
		1257
		1258
		1259
		1260
		1261
	Preslav Nakov, David Corney, Maram Hasanain, Firoj Alam, Tamer Elsayed, Alberto Barrón-Cedeño, Paolo Papotti, Shaden Shaar, and Giovanni Da San Martino. 2021. Automated fact-checking for assisting human fact-checkers. In <i>Proceedings of the 30th International Joint Conference on Artificial Intelligence, IJ-CAI ’21</i> , pages 4551–4558.	1262
		1263
		1264
		1265
		1266
		1267
		1268
	OpenAI. 2023. GPT-4 technical report . Technical report, OpenAI.	1269
		1270
	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in Neural Information Processing Systems</i> , 35:27730–27744.	1271
		1272
		1273
		1274
		1275
		1276
	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. Robust speech recognition via large-scale weak supervision. <i>arXiv preprint arXiv:2212.04356</i> .	1277
		1278
		1279
		1280
	Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving Language Understanding by Generative Pre-Training". Technical report, Open AI.	1281
		1282
		1283
		1284
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	1285
		1286
		1287
		1288
	Revanth Gangi Reddy, Sai Chetan Chinthakindi, Zhenhailong Wang, Yi Fung, Kathryn Conger, Ahmed Elsayed, Martha Palmer, Preslav Nakov, Eduard Hovy, Kevin Small, et al. 2022. Newsclaims: A new benchmark for claim detection from news with attribute knowledge. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 6002–6018.	1289
		1290
		1291
		1292
		1293
		1294
		1295
		1296
	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. 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 3982–3992.	1297
		1298
		1299
		1300
		1301
		1302
	Sara Rosenthal, Noura Farra, and Preslav Nakov. 2017. Semeval-2017 task 4: Sentiment analysis in twitter. In <i>Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)</i> , pages 502–518.	1303
		1304
		1305
		1306
		1307
	Hassan Sajjad, Ahmed Abdelali, Nadir Durrani, and Fahim Dalvi. 2020. AraBench: Benchmarking dialectal Arabic-English machine translation. In <i>Pro-</i>	1308
		1309
		1310

1311			1370
1312		<i>ceedings of the 28th International Conference on Computational Linguistics</i> , pages 5094–5107.	1371
1313	Younes Samih, Mohamed Eldesouki, Mohammed Attia,		1372
1314	Kareem Darwish, Ahmed Abdelali, Hamdy Mubarak,		1373
1315	and Laura Kallmeyer. 2017. Learning from relatives:		1374
1316	Unified dialectal Arabic segmentation. In <i>Proceed-</i>		1375
1317	<i>ings of the 21st Conference on Computational Nat-</i>		1376
1318	<i>ural Language Learning (CoNLL 2017)</i> , pages 432–		1377
1319	441.		1378
1320	Ahmed Samir, Abu Bakr Soliman, Mohamed Ibrahim,		1379
1321	Laila Hesham, and Samhaa R El-Beltagy. 2022.		1380
1322	Ngu_cnlp at wanlp 2022 shared task: Propaganda		1381
1323	detection in arabic. <i>WANLP 2022</i> , page 545.		1382
1324	Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi,		1383
1325	and Noah A Smith. 2019. The risk of racial bias in		1384
1326	hate speech detection. In <i>Proceedings of the 57th</i>		1385
1327	<i>annual meeting of the association for computational</i>		1386
1328	<i>linguistics</i> , pages 1668–1678.		1387
1329	Teven Le Scao, Angela Fan, Christopher Akiki, El-		1388
1330	lie Pavlick, Suzana Ilić, Daniel Hesslow, Roman		1389
1331	Castagné, Alexandra Sasha Luccioni, François Yvon,		1390
1332	Matthias Gallé, et al. 2022. BLOOM: A 176b-		1391
1333	parameter open-access multilingual language model.		1392
1334	<i>arXiv preprint arXiv:2211.05100</i> .		1393
1335	Anna Schmidt and Michael Wiegand. 2017. A survey		1394
1336	on hate speech detection using natural language pro-		1395
1337	cessing. In <i>Proceedings of the fifth international</i>		1396
1338	<i>workshop on natural language processing for social</i>		1397
1339	<i>media</i> , pages 1–10.		1398
1340	Nathan Schneider, Behrang Mohit, Kemal Oflazer, and		1399
1341	Noah A Smith. 2012. Coarse lexical semantic an-		1400
1342	notation with supersenses: an arabic case study. In		1401
1343	<i>Proceedings of the 50th Annual Meeting of the As-</i>		1402
1344	<i>sociation for Computational Linguistics (Volume 2:</i>		1403
1345	<i>Short Papers)</i> , pages 253–258.		1404
1346	Fabrizio Sebastiani. 2002. Machine learning in auto-		1405
1347	mated text categorization. <i>ACM computing surveys</i>		1406
1348	(<i>CSUR</i>), 34(1):1–47.		1407
1349	Haitham Seelawi, Ahmad Mustafa, Hesham Al-		1408
1350	Bataineh, Wael Farhan, and Hussein T Al-Natsheh.		1409
1351	2019. Nsurl-2019 task 8: Semantic question simi-		1410
1352	larity in Arabic. In <i>Proceedings of the First Inter-</i>		1411
1353	<i>national Workshop on NLP Solutions for Under Re-</i>		1412
1354	<i>sourced Languages (NSURL 2019) co-located with</i>		1413
1355	<i>ICNLSP 2019-Short Papers</i> , pages 1–8.		1414
1356	Haitham Seelawi, Ibraheem Tuffaha, Mahmoud Gzawi,		1415
1357	Wael Farhan, Bashar Talafha, Riham Badawi, Ziad		1416
1358	Sober, Oday Al-Dweik, Abed Alhakim Freihat, and		1417
1359	Hussein Al-Natsheh. 2021. Alue: Arabic language		1418
1360	understanding evaluation. In <i>Proceedings of the</i>		1419
1361	<i>Sixth Arabic Natural Language Processing Workshop</i> ,		1420
1362	pages 173–184.		1421
1363	Neha Sengupta, Sunil Kumar Sahu, Bokang Jia,		1422
1364	Satheesh Katipomu, Haonan Li, Fajri Koto,		1423
1365	Osama Mohammed Afzal, Samta Kamboj, Onkar		1424
1366	Pandit, Rahul Pal, et al. 2023. Jais and jais-chat:		1425
1367	Arabic-centric foundation and instruction-tuned open		1426
1368	generative large language models. <i>arXiv preprint</i>		1427
1369	<i>arXiv:2308.16149</i> .		1428
	Shaden Shaar, Maram Hasanain, Bayan Hamdan,		1429
	Zien Sheikh Ali, Fatima Haouari, Alex Nikolov,		1430
	Mucahid Kutlu, Yavuz Selim Kartal, Firoj Alam,		1431
	Giovanni Da San Martino, Alberto Barrón-Cedeño,		1432
	Rubén Míguez, Javier Beltrán, Tamer Elsayed, and		1433
	Preslav Nakov. 2021. Overview of the CLEF-2021		1434
	CheckThat! lab task 1 on check-worthiness estima-		1435
	tion in tweets and political debates. In <i>2021 Working</i>		1436
	<i>Notes of CLEF - Conference and Labs of the Evalua-</i>		1437
	<i>tion Forum</i> .		1438
	Shivam Sharma, Firoj Alam, Md. Shad Akhtar, Dimitar		1439
	Dimitrov, Giovanni Da San Martino, Hamed Firooz,		1440
	Alon Halevy, Fabrizio Silvestri, Preslav Nakov, and		1441
	Tanmoy Chakraborty. 2022. Detecting and under-		1442
	standing harmful memes: A survey . In <i>Proceedings</i>		1443
	<i>of the Thirty-First International Joint Conference on</i>		1444
	<i>Artificial Intelligence, IJCAI '22</i> , pages 5597–5606,		1445
	Vienna, Austria. International Joint Conferences on		1446
	Artificial Intelligence Organization. Survey Track.		1447
	Taylor Shin, Yasaman Razeghi, Robert L Logan IV,		1448
	Eric Wallace, and Sameer Singh. 2020. Autoprompt:		1449
	Eliciting knowledge from language models with		1450
	automatically generated prompts. <i>arXiv preprint</i>		1451
	<i>arXiv:2010.15980</i> .		1452
	Mohammad Shoeybi, Mostofa Patwary, Raul Puri,		1453
	Patrick LeGresley, Jared Casper, and Bryan Catan-		1454
	zaro. 2019. Megatron-lm: Training multi-billion		1455
	parameter language models using model parallelism.		1456
	<i>arXiv preprint arXiv:1909.08053</i> .		1457
	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao,		1458
	Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch,		1459
	Adam R Brown, Adam Santoro, Aditya Gupta,		1460
	Adrià Garriga-Alonso, et al. 2022. Beyond the		1461
	imitation game: Quantifying and extrapolating the		1462
	capabilities of language models. <i>arXiv preprint</i>		1463
	<i>arXiv:2206.04615</i> .		1464
	Md Tawkat, Islam Khondaker, Abdul Waheed,		1465
	El Moatez Billah Nagoudi, and Muhammad Abdul-		1466
	Mageed. 2023. GPTAraEval: A comprehensive eval-		1467
	uation of chatgpt on arabic nlp. <i>arXiv e-prints</i> , pages		1468
	arXiv–2305.		1469
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob		1470
	Uszkoreit, Llion Jones, Aidan N Gomez, ukasz		1471
	Kaiser, and Illia Polosukhin. 2017. Attention is all		1472
	you need. <i>Advances in neural information processing</i>		1473
	<i>systems</i> , 30.		1474
	Chengyi Wang, Sanyuan Chen, Yu Wu, Ziqiang Zhang,		1475
	Long Zhou, Shujie Liu, Zhuo Chen, Yanqing Liu,		1476
	Huaming Wang, Jinyu Li, et al. 2023. Neural codec		1477
	language models are zero-shot text to speech synthe-		1478
	sizers. <i>arXiv preprint arXiv:2301.02111</i> .		1479
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten		1480
	Bosma, Fei Xia, Ed H Chi, Quoc V Le, Denny Zhou,		1481
	et al. 2022. Chain-of-thought prompting elicits rea-		1482
	soning in large language models. In <i>Advances in</i>		1483
	<i>Neural Information Processing Systems</i> .		1484
	Shu-wen Yang, Po-Han Chi, Yung-Sung Chuang,		1485
	Cheng-I Jeff Lai, Kushal Lakhotia, Yist Y Lin,		1486
	Andy T Liu, Jiatong Shi, Xuankai Chang, Guan-		1487
	Ting Lin, et al. 2021. Superb: Speech processing		1488

1430 universal performance benchmark. *arXiv preprint*
1431 *arXiv:2105.01051*.

1432 Xi Ye, Srinivasan Iyer, Asli Celikyilmaz, Ves Stoyanov, Greg Durrett, and Ramakanth Pasunuru. 2022.
1433 Complementary explanations for effective in-context
1434 learning. *arXiv preprint arXiv:2211.13892*.
1435

1436 Marcos Zampieri, Preslav Nakov, Sara Rosenthal, Pepa
1437 Atanasova, Georgi Karadzhov, Hamdy Mubarak,
1438 Leon Derczynski, Zeses Pitenis, and Çağrı Çöltekin.
1439 2020. Semeval-2020 task 12: Multilingual offensive
1440 language identification in social media (offenseval
1441 2020). *arXiv preprint arXiv:2006.07235*.

1442 Rabih Zbib, Erika Malchiodi, Jacob Devlin, David Stal-
1443 lard, Spyros Matsoukas, Richard Schwartz, John
1444 Makhoul, Omar Zaidan, and Chris Callison-Burch.
1445 2012. Machine translation of Arabic dialects. In *Pro-*
1446 *ceedings of the 2012 conference of the north amer-*
1447 *ican chapter of the association for computational*
1448 *linguistics: Human language technologies*, pages 49–
1449 59.

1450 Daniel Zeman, Joakim Nivre, Mitchell Abrams, Elia
1451 Ackermann, Noëmi Aepli, Hamid Aghaei, and
1452 R Ziane. 2020. Universal dependencies 2.5. *LIN-*
1453 *DAT/CLARIAHCZ digital library at the Institute of*
1454 *Formal and Applied Linguistics (UFAL), Faculty of*
1455 *Mathematics and Physics, Charles University*. url:
1456 <http://hdl.handle.net/11234/1-3226>.

1457 Lei Zhang, Shuai Wang, and Bing Liu. 2018. Deep
1458 learning for sentiment analysis: A survey. *Wiley*
1459 *Interdisciplinary Reviews: Data Mining and Knowl-*
1460 *edge Discovery*, 8(4):e1253.

1461 Susan Zhang, Stephen Roller, Naman Goyal, Mikel
1462 Artetxe, Moya Chen, Shuohui Chen, Christopher De-
1463 wan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022.
1464 Opt: Open pre-trained transformer language models.
1465 *arXiv preprint arXiv:2205.01068*.

1466 Yu Zhang, Wei Han, James Qin, Yongqiang Wang,
1467 Ankur Bapna, Zhehuai Chen, Nanxin Chen, Bo Li,
1468 Vera Axelrod, Gary Wang, Zhong Meng, Ke Hu, An-
1469 drew Rosenberg, Rohit Prabhavalkar, Daniel S. Park,
1470 Parisa Haghani, Jason Riesa, Ginger Perng, Hagen
1471 Soltau, Trevor Strohman, Bhuvana Ramabhadran,
1472 Tara Sainath, Pedro Moreno, Chung-Cheng Chiu, Jo-
1473 han Schalkwyk, Françoise Beaufays, and Yonghui
1474 Wu. 2023. *Google usm: Scaling automatic speech*
1475 *recognition beyond 100 languages*. *arXiv preprint*
1476 *arXiv:2303.01037*.

1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521

Appendix

A Tasks and Datasets

In this section, we discuss the tasks and the associated datasets by grouping them based on ACL-2022 track.¹⁰ In Tables 5 and 6, we provide a summarized description of the test sets used for evaluating textual and speech processing tasks, respectively.

A.1 Word Segmentation, Syntax and Information Extraction

A.1.1 Segmentation

Segmentation is an important problem for language like Arabic, which is rich with bound morphemes that change the tense of verbs, or represent pronouns and prepositions in nouns. It is a building block for NLP tasks such as search, part-of-speech tagging, parsing, and machine translation. The idea is segmenting Arabic words into prefixes, stems, and suffixes, which can facilitate many other tasks.

Datasets

WikiNews For modern standard Arabic (MSA), we used the WikiNews dataset of (Darwish and Mubarak, 2016) which comprises 70 news articles in politics, economics, health, science and technology, sports, arts, and culture. The dataset has 400 sentences (18,271 words) in total.

Tweets For the dialectal Arabic, we used the dataset in (Samih et al., 2017), which provides 1400 tweets in Egyptian, Gulf, Levantine, and Maghrebi dialects for a total of 25,708 annotated words .

A.1.2 Part-Of-Speech (POS) Tagging

Part-of-speech (POS) is one of the fundamental components in the NLP pipeline. It helps in extracting higher-level information such as named entities, discourse, and syntactic parsing.

Datasets

WikiNews We used for this task the WikiNews dataset tagged for POS (Darwish et al., 2017c) for modern standard Arabic.

Tweets For POS tagging with noisy texts and different dialects we used the same dataset reported in (Samih et al., 2017) (see §A.1.1).

XGLUE We also used the Arabic part of XGLUE benchmark (Liang et al., 2020b) for POS tagging, which uses a subset of Universal Dependencies Treebanks (v2.5) (Zeman et al., 2020).

¹⁰<https://www.2022.aclweb.org/callpapers>

A.1.3 Lemmatization

Lemmatization is another component in the NLP pipeline, which reduces words to their base or root form, known as a lemma. It takes into consideration the morphological analysis of the words, which uses the context and POS to convert a word to its simplest form. This task differs from segmentation which only separates a word stem from prefixes and suffixes. In contrast, lemmatization requires returning the lexicon entry for a certain word, which may depend on POS tagging.

Dataset We used WikiNews dataset tagged for lemmas (Mubarak, 2018) (see §A.1.1 for the details of the dataset).

A.1.4 Diacritization

Diacritization involves assigning the diacritics to each letter in an Arabic word within a sentence. Diacritical marks indicate the correct pronunciation and meaning of the written Arabic words. For example, different word diacritizations could transform a noun into a verb or vice versa.

Datasets

WikiNews We use a dataset of modern standard Arabic from (Mubarak et al., 2019) that comprises fully diacritized WikiNews corpus (Darwish et al., 2017b).

Bibles This dataset includes translations of the New Testament into two Maghrebi sub-dialects: Moroccan and Tunisian (Darwish et al., 2018; Abdelali et al., 2019).

A.1.5 Parsing

Dependency parsing is the task of identifying syntactical and grammatical relations among the words in a sentence. These dependencies result in a hierarchical tree representation that captures the structure of the sentence at different levels.

Dataset For this task we used the Arabic part of CoNLL-X 2006 shared tasks on dependency parsing (Buchholz and Marsi, 2006), which has 4,990 scoring tokens and uses the Prague Arabic Dependency Treebank (Hajic et al., 2004).

A.1.6 Named-Entity Recognition (NER)

This task involves identifying and classifying the words in a sentence that are proper names, names of places, entities like organizations or products, amongst other things. This depends on understanding the context and the relations of a word or a

1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568

Dataset	Task	Domain	Test Set Size
Word Segmentation, Syntax and Information Extraction			
WikiNews	Segmentation	News articles (MSA)	400 sentences
Samih et al. (2017)	Segmentation	Tweets (Dialects: EGY, LEV, GLF, MGR)	70 X 4 dialects
WikiNews	Lemmatization	News articles (MSA)	400 sentences
WikiNews	Diacritization	News articles (MSA)	400 sentences
Darwish et al. (2018)	Diacritization	Sentences (Dialects: Moroccan, Tunisian)	1,640 X 2 dialects
WikiNews	POS	News articles (MSA)	400 sentences
Samih et al. (2017)	POS	Tweets (Dialects: EGY, LEV, GLF, MGR)	70 X 4 dialects
XGLUE (Arabic)	POS	Web, Wikipedia	680 sentences
Conll2006	Parsing	MSA	146 sentences
ANERcorp	NER	News articles	924 sentences
AQMAR	NER	Wikipedia	1,976 sentences
QASR	NER	Transcripts	7,906 segments
QADI	Dialect	Tweets	3,797
ADI	Dialect	Transcripts	751
Sentiment, Stylistic and Emotion Analysis			
ArSAS	Sentiment	Tweets	4,213
SemEval2018-Task1	Emotion	Tweets (Dialectal)	1,518
Unified-FC	Stance	News articles	3,042 claim-article pairs
ANS	Stance	News articles	379 headline pairs
ArSarcasm	Sarcasm	Tweets	2,110
ArSarcasm-2	Sarcasm	Tweets	3,000
News Categorization			
ASND	News Cat.	Posts*	1,103
SANAD/Akhbarona	News Cat.	News articles	7,843
SANAD/AlArabiya	News Cat.	News articles	7,125
SANAD/AlKhaleej	News Cat.	News articles	4,550
Demographic Attributes			
ASAD	Name Info	Wikidata	80,130
UL2C	Location	User loc. (Twitter)	28,317
Arap-Tweet	Gender	Usernames (Twitter)	640
Ethics in NLP: Factuality, Disinformation and Harmful Content Detection			
OffensEval2020	Offensive lang.	Tweets (Dialectal)	2,000
OSACT2020	Hate Speech	Tweets (Dialectal)	2,000
ASAD	Adult Content	Tweets (Dialectal)	10,000
ASAD	Spam	Tweets (Dialectal)	28,383
In-house	Subjectivity	News articles	297 sentences
WANLP23	Propaganda	Tweets	323
CT-CWT-22	Checkworthiness	Tweets (COVID19)	680
COVID19 Disinfo.	Factuality	Tweets	996
Unified-FC	Factuality	News articles	422 claims
ANS	Factuality	News articles	456 headlines
CT-CWT-22	Claim	Tweets (COVID19)	1,248
CT-CWT-22	Harmful content	Tweets (COVID19)	1,201
CT-CWT-22	Attention-worthy	Tweets (COVID19)	1,186
Semantic Textual Similarity (STS)			
STS2017-Track 1	STS	Image captions	250 sentence pairs
STS2017-Track 2	STS	Image captions	250 sentence pairs
Mawdo3 Q2Q	STS QS (Q2Q)	Questions	3,715 question pairs
XNLI	XNLI	ANC	5,010 sentence pairs
Question Answering (QA)			
ARCD	QA	Wikipedia	702 questions
MLQA	QA	Wikipedia	5,335 questions
TyDi QA	QA	Wikipedia	921 questions
XQuAD	QA	Wikipedia	1,190 questions

Table 5: Summary on test sets and their sizes used in evaluation for the different textual tasks. **ANC**: American National Corpus. **Posts***: posts from Twitter, Youtube and Facebook. **News Cat.**: News Categorization

1569	collection of words in a sentence, and is key to	roccan ¹¹ and Tunisian ¹² dialects (Abdelali et al.,	1614
1570	tasks such as question answering.	2019).	1615
1571	Datasets	Media Dataset The dataset consists of 7.5 hours	1616
1572	ANERCorp We used the test corpus of the AN-	of recordings collected from five public broadcast-	1617
1573	ERCorp dataset (Benajiba et al., 2007; Benajiba	ing channels that cover programs with Maghrebi,	1618
1574	and Rosso, 2007), which contains 316 articles,	Lebanese, Omani dialects, and MSA with genres	1619
1575	150,286 tokens and 32,114 types, and classifies	involving movies, news reports, and cultural pro-	1620
1576	words into one of four classes (organization, loca-	grams. The recordings were transcribed and trans-	1621
1577	tion, person and miscellaneous), we used the test	lated by a professional translation house (Sajjad	1622
1578	split of the dataset for our evaluation.	et al., 2020).	1623
1579	AQMAR The dataset is developed as an evalua-	A.3 Dialect Identification	1624
1580	tion suite for the named entity recognition task in	Dialect is defined as the speaker’s grammatical, lex-	1625
1581	Arabic. It consists of a collection of 28 Wikipedia	ical, and phonological variation in pronunciation	1626
1582	articles with 74,000 tokens. We consider the arti-	(Etman and Beex, 2015). Automatic Dialect Identifi-	1627
1583	cles corresponding to the test split for our evalua-	cation (ADI) has become an important research	1628
1584	tion. (Schneider et al., 2012).	area in order to improve certain applications and	1629
1585	QASR The QASR dataset consists of 70k words	services, such as ASR and many downstream NLP	1630
1586	extracted from 2,000 hours of transcribed Arabic	tasks.	1631
1587	speech (Mubarak et al., 2021b).	Dataset For this task, we used the QADI dataset	1632
1588	A.2 Machine Translation (MT)	containing a wide range of country-level Arabic	1633
1589	The machine translation evaluation set is a rich	dialects covering 18 different countries in the Mid-	1634
1590	set that covers a variety of Arabic in addition to	dle East and North Africa region (Abdelali et al.,	1635
1591	the Modern Standard Arabic (MSA). The genera-	2020). It consists of 540,590 tweets from 2,525	1636
1592	of the evaluation set also cover formal, informal,	users.	1637
1593	speech, and other modalities. These types and va-	A.4 Sentiment, Stylistic and Emotion Analysis	1638
1594	rieties allowed us to assess the system and reveal	A.4.1 Sentiment Analysis	1639
1595	its potential and limitations. For this study, we fo-	Sentiment analysis has been an active research area	1640
1596	ocused on translating Arabic to English and used the	and aims to analyze people’s sentiment or opin-	1641
1597	datasets discussed below.	ion toward entities such as topics, events, individ-	1642
1598	Datasets	uals, issues, services, products, organizations, and	1643
1599	MADAR Corpus This dataset consists of 2,000	their attributes (Liu and Zhang, 2012; Zhang et al.,	1644
1600	sentences from the BTEC corpus translated to mod-	2018). This task involves classifying the content	1645
1601	ern standard Arabic and four major dialects from	into sentiment labels such as positive, neutral, and	1646
1602	15 countries (Bouamor et al., 2018).	negative.	1647
1603	(Zbib et al., 2012) It is collected from the Arabic-	Dataset ArSAS dataset consists of 21k Arabic	1648
1604	Dialect/English Parallel Text (APT), which consists	tweets covering multiple topics that were collected,	1649
1605	of 2,000 sentences with 3.5 million tokens of trans-	prepared, and annotated for six different classes of	1650
1606	lated dialectal Arabic (Zbib et al., 2012).	speech-act labels and four sentiment classes (El-	1651
1607	Multi-dialectal Parallel Corpus of Arabic	madany et al., 2018). For the experiments, we used	1652
1608	(MDC) This dataset also consists of 2,000 sen-	only sentiment labels from this dataset.	1653
1609	tences in Egyptian, Palestinian, Syrian, Jordanian,	A.4.2 Emotion Recognition	1654
1610	and Tunisian dialects and their English counter-	Emotion recognition is the task of categorizing dif-	1655
1611	parts (Bouamor et al., 2014).	ferent types of content (e.g., text, speech, and vi-	1656
1612	The Bible It consists of 8.2k parallel sentences	sual) in different emotion labels (six basic emo-	1657
1613	translated into modern standard Arabic, and to Mo-		

¹¹The Morocco Bible Society <https://www.biblesociety.ma>

¹²The United Bible Societies <https://www.bible.com>

1658	tions (Ekman, 1971) or more fine-grained categories (Demszky et al., 2020)).		1705
1659			1706
1660	Dataset For the emotion recognition tasks we		1707
1661	used SemEval-2018 Task 1: Affect in Tweets (Mo-		1708
1662	hammad et al., 2018). The task is defined as classi-		1709
1663	fying a tweet as one or more of the eleven emotion		1710
1664	labels, which is annotated as a multilabel (pres-		1711
1665	ence/absence of 11 emotions) annotation setting.		1712
1666			1713
1667	A.4.3 Stance Detection		1714
1668	Stance is defined as the expression of the speaker’s		1715
1669	view and judgment toward a given argument or		1716
1670	statement (Biber and Finegan, 1988). Given that		1717
1671	the social media platforms allow users to consume		1718
1672	and disseminate information by expressing their		1719
1673	views, enabling them to obtain instant feedback		1720
1674	and explore others’ views, it is important to char-		1721
1675	acterize a stance expressed in a given content. Au-		1722
1676	tomatic stance detection also allows for assessing		1723
1677	public opinion on social media, particularly on dif-		1724
1678	ferent social and political issues such as abortion,		1725
1679	climate change, and feminism, on which people ex-		1726
1680	press supportive or opposing opinions (ALDayel		1727
1681	and Magdy, 2021; Küçük and Can, 2020). The task		1728
1682	involves “classification as the stance of the pro-		1729
1683	ducer of a piece of text, towards a target as either		1730
1684	one of the three classes: {support, against, neither}		1731
1685	or {agree, disagree, discuss, or unrelated}” (Küçük		1732
1686	and Can, 2020).		1733
1687			1734
1688	Datasets		1735
1689			1736
1690	Unified-FC dataset consists of claims collected		1737
1691	from Verify.sy (false claims) and Reuters (true		1738
1692	claims), which resulted in 422 claims. Based		1739
1693	on these claims documents are collected using		1740
1694	Google custom search API and filtered by com-		1741
1695	puting claim-documents similarity (Baly et al.,		1742
1696	2018b). This approach resulted in 3,042 claim-		1743
1697	documents pairs, which are then annotated for		1744
1698	stance (agree, disagree, discuss, unrelated) by Ap-		1745
1699	pen crowd-sourcing platform.		1746
1700			1747
1701	ANS Khouja (2020) developed a dataset by first		1748
1702	sampling news titles from Arabic News Texts		1749
1703	(ANT) corpus (Chouigui et al., 2017) and then gen-		1750
1704	erating true and false claims. From these claims		1751
	stance (three classes – agree, disagree, other) is		1752
	annotated from a pair of sentences using Amazon		1753
	Mechanical Turk and Upwork. The dataset consists		
	of 3,786 claim-reference pairs.		
	ArSarcasm Abu Farha and Magdy (2020) de-		
	veloped an Arabic sarcasm detection dataset. The		
	dataset was created using previously available Ara-		
	bic sentiment analysis datasets (Rosenthal et al.,		
	2017; Nabil et al., 2015) and adds sarcasm and di-		
	allect labels to them. The dataset contains 10,547		
	tweets, 1,682 of which are sarcastic. The training		
	set contains 8,437 tweets, while the test set contains		
	2,110 tweets.		
	ArSarcasm-v2 This dataset is an extension of the		
	original ArSarcasm dataset published along with		
	the paper (Farha and Magdy, 2020). ArSarcasm-		
	v2 consists of ArSarcasm along with portions of		
	DAICT corpus and some new tweets. Each tweet		
	was annotated for sarcasm, sentiment and dialect.		
	The final dataset consists of 15,548 tweets divided		
	into 12,548 training tweets and 3,000 testing tweets.		
	ArSarcasm-v2 was used and released as a part of		
	the shared task on sarcasm detection and sentiment		
	analysis in Arabic.		
	A.5 News Categorization		
	News text categorization was a popular task in the		
	earlier days of NLP research (Sebastiani, 2002).		
	The idea of to assign a category $C = \{c_1, \dots, c_n\}$		
	to a document $D = \{d_1, \dots, d_n\}$. For the news		
	categorization the D is a set of news articles and		
	C is a set of predefined categories. Most often a		
	news article can be categorized into more than one		
	category and the models are trained in a multilabel		
	setting. While earlier work mostly focused on news		
	article, however, lately it has been used for the		
	categorization of tweets in which news articles are		
	shared as a part of a tweet.		
	Datasets		
	Social Media Posts ASND is a News Tweets		
	dataset (Chowdhury et al., 2020b), collected from		
	Aljazeera news channel accounts on Twitter, Face-		
	book, and YouTube. The dataset consists of twelve		
	categories such as art-and-entertainment, business-		
	and-economy, crime-war-conflict, education, envi-		
	ronment, health, human-rights-press-freedom, poli-		
	tics, science-and-technology, spiritual, sports, and		
	(xii) others. We used the test split from each dataset		
	for the evaluation.		
	Arabic News SANAD corpus is a large col-		
	lection of Arabic news articles collected from		
	Akhbarona, AlKhaleej, and AlArabiya (Einea et al.,		
	2019). The dataset has separate collections gather-		
	ed from different news media, each of which has		

1754	six news categories; namely culture, finance, medi-	or opinions. Otherwise, the sentence is considered	1798
1755	cal, politics, sports and technology.	objective (Antici et al., 2021). Given that the identi-	1799
1756	A.6 Demographic/Protected Attributes	fication of subjectivity is subjective itself, therefore,	1800
1757	Demographic information (e.g., gender, age, coun-	it poses challenges in the annotation process by the	1801
1758	try of origin) are useful in many different appli-	annotator. The complexity lies due to the different	1802
1759	cations such as understanding population charac-	levels of expertise by the annotators, different in-	1803
1760	teristics, personalized advertising, socio-cultural	terpretations and their conscious and unconscious	1804
1761	studies, etc. Demographic information helps gov-	bias towards the content they annotate. The content	1805
1762	ernments, businesses, and organizations understand	can be text (e.g., sentence, article), image or multi-	1806
1763	their target audiences, and plan accordingly.	modal content, consisting of opinionated, factual	1807
1764	A.6.1 Gender	or non-factual content. The annotation typically	1808
1765	Gender analysis can reveal important differences	has been done using two labels, objective (OBJ)	1809
1766	between male and female users such as topics of	and subjective (SUBJ).	1810
1767	interest, gender gap, preferences, etc.	Dataset The dataset consists of sentences curated	1811
1768	Dataset We used the ArabGend test set, which	from news articles. The dataset has been developed	1812
1769	contains 1,000 names collected from Twit-	based on the existing AraFacts dataset (Ali et al.,	1813
1770	ter (divided equally between males and fe-	2021b) that contains claims verified by Arabic fac-	1814
1771	males) (Mubarak et al., 2022).	checking websites, and each claim is associated	1815
1772	A.6.2 Location	with web pages propagating or negating the claim.	1816
1773	Identifying user locations is useful for many appli-	The news articles are collected from different news	1817
1774	cations such as author profiling, dialect identifica-	media. News articles were automatically parsed,	1818
1775	tion, recommendation systems, etc. Often, users	split into sentences and filtered poorly-formatted	1819
1776	on social media platforms, such as Twitter, declare	sentences using a rule-based approach. The dataset	1820
1777	their locations in noisy ways, and mapping these	has been released as a part of Task 2 of CLEF2023	1821
1778	locations to countries is a challenging task.	CheckThat Lab (Barrón-Cedeño et al., 2023).	1822
1779	Dataset We used the UL2C dataset, which con-	A.7.2 Propaganda Detection	1823
1780	tains 28K unique locations, as written by Arabic	Propaganda can be defined as a form of commu-	1824
1781	Twitter users, and their mappings to Arab coun-	nication that aims to influence the opinions or the	1825
1782	tries (Mubarak and Hassan, 2021).	actions of people towards a specific goal; this is	1826
1783	A.6.3 Name Info	achieved utilizing well-defined rhetorical and psy-	1827
1784	Names contain important information about our	chological devices (Dimitrov et al., 2021). In differ-	1828
1785	identities and demographic characteristics, includ-	ent communication channels, propaganda (persua-	1829
1786	ing factors like gender, nationality, and ethnicity.	sion techniques) is conveyed through the use of di-	1830
1787	The purpose of this task is to predict the country of	verse techniques (Miller, 1939), which range from	1831
1788	origin of a person name giving only their names.	leveraging the emotions of the audience, such as	1832
1789	Dataset We used an in-house dataset for mapping	using <i>emotional technique</i> or logical fallacies such	1833
1790	person names to World countries extracted from	as <i>straw man</i> (misrepresenting someone’s opinion),	1834
1791	Wikipedia. ¹³	hidden <i>ad-hominem fallacies</i> , and <i>red herring</i> (pre-	1835
1792	A.7 Ethics and NLP: Factuality,	sending irrelevant data).	1836
1793	Disinformation and Harmful content	Dataset The dataset used for this study consists	1837
1794	detection	of Arabic tweets (Alam et al., 2022b) posted by	1838
1795	A.7.1 Subjectivity Identification	different news media from Arab countries such as	1839
1796	A sentence is considered subjective when it is based	Al Arabiya and Sky News Arabia from UAE, Al	1840
1797	on – or influenced by – personal feelings, tastes,	Jazeera, and Al Sharq from Qatar, and from five	1841
		international Arabic news sources Al-Hurra News,	1842
		BBC Arabic, CNN Arabic, France 24, and Russia	1843
		Today. The final annotated dataset consists of 930	1844
		tweets. Alam et al. (2022b) formulated the task as	1845
		a multilabel and multiclass span level classification	1846
		task. For this study, we used the multilabel setup.	1847

¹³Paper is under revision.

1848 A.7.3 Check-worthiness Detection

1849 Fact-checking is a time-consuming and complex
1850 process, and it often takes effort to determine
1851 whether a claim is important to check, irrespective
1852 of its potential to be misleading or not. Check-
1853 worthiness detection is the first step and a criti-
1854 cal component of fact-checking systems (Nakov
1855 et al., 2021) and the aim is to facilitate manual
1856 fact-checking efforts by prioritizing the claims for
1857 the fact-checkers. Research on check-worthiness
1858 includes check-worthiness detection/ranking from
1859 political speeches, debates, and social media posts
1860 (Nakov et al., 2022a; Shaar et al., 2021). A check-
1861 worthy claim is usually defined by its importance
1862 to the public and journalists, and whether it can
1863 cause harm to an individual, organization, and/or
1864 society.

1865 **Dataset** For this study, we used the Arabic subset
1866 of the dataset released with Task 1A (Arabic) of the
1867 CLEF2022 CheckThat Lab (Nakov et al., 2022b)
1868 The dataset consists of 4,121 annotated tweets. The
1869 Arabic tweets were collected using keywords re-
1870 lated to COVID-19, vaccines, and politics.

1871 A.7.4 Claim Detection

1872 Information shared in the mainstream and social
1873 media often contains misleading content. Claim de-
1874 tection has become an important problem in order
1875 to mitigate misinformation and disinformation in
1876 those media channels. A factual (verifiable) claim
1877 is a sentence claiming that something is true, and
1878 this can be verified using factually verifiable in-
1879 formation such as statistics, specific examples, or
1880 personal testimony (Konstantinovskiy et al., 2021).
1881 Research on claim detection includes social media
1882 posts – text modality (Alam et al., 2021b), multi-
1883 modality (Cheema et al., 2022) and news (Reddy
1884 et al., 2022).

1885 Datasets

1886 **CT-CWT-22-Claim** We used the Arabic sub-
1887 set of the dataset released with Task 1B of the
1888 CLEF2022 CheckThat Lab (Nakov et al., 2022a).
1889 The dataset has been annotated using a multi-
1890 question annotation schema (Alam et al., 2021a),
1891 which consists of tweets collected using COVID-
1892 19 related keywords. The dataset contains 6,214
1893 tweets (Nakov et al., 2022b).

1894 **ANS** (Khouja, 2020) This dataset consists of
1895 4,547 true and false claims, which was developed
1896 based on Arabic News Texts (ANT) corpus. A

sample of articles was modified to generate true
and false claims using crowdsourcing.

1897 A.7.5 Attention-worthiness Detection

1898 In social media most often people tweet by blaming
1899 authorities, providing advice, and/or call for
1900 action. It might be important for the policy mak-
1901 ers to respond to those posts. The purpose of this
1902 task is to categorize such information into one of
1903 the following categories: *not interesting, not sure,*
1904 *harmfulness, other, blames authorities, contains*
1905 *advice, calls for action, discusses action taken, dis-*
1906 *cusses cure, asks a question.*

1907 **Dataset** For this task, we used a subset of the
1908 dataset Task 1D of the CLEF2022 CheckThat
1909 Lab (Nakov et al., 2022a), which contains 6,140
1910 annotated tweets.

1911 A.7.6 Factuality Detection

1912 Fact-checking has emerged as an important re-
1913 search topic due to a large amount of fake news, ru-
1914 mors, and conspiracy theories that are spreading in
1915 different social media channels to manipulate peo-
1916 ple’s opinions or to influence the outcome of major
1917 events such as political elections (Darwish et al.,
1918 2017a; Baly et al., 2018b). While fact-checking has
1919 largely been done by manual fact-checker due to
1920 the reliability, however, that does not scale well as
1921 the enormous amount of information shared online
1922 every day. Therefore, an automatic fact-checking
1923 system is important and it has been used for fa-
1924 cilitating human fact-checker (Nakov et al., 2021).
1925 The task typically involves assessing the level of
1926 factual correctness in a news article, media outlets,
1927 or social media posts. The content is generally
1928 judged to be of high, low, or mixed factual correct-
1929 ness, seven-point Likert scale^{14,15} or just binary
1930 labels {yes, no} (Baly et al., 2018a; Alam et al.,
1931 2021b).

1932 Datasets

1933 **News Articles** We used the dataset developed
1934 by Baly et al. (2018a) in which false claims are
1935 extracted from verify-sy¹⁶ and true claims are
1936 extracted from <http://ara.reuters.com>. The
1937 dataset consists of 3,042 documents.

1938 **Tweets** For the claim detection from tweets, we
1939 used the same dataset (Alam et al., 2021b) dis-
1940 cussed in A.7.4. As mentioned earlier, this dataset

¹⁴<https://mediabiasfactcheck.com>

¹⁵<https://allsides.com>

¹⁶<http://www.verify-sy.com>

1943
1944
1945
1946
1947
1948

1949

1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966

1967
1968
1969

1970
1971
1972
1973
1974
1975
1976
1977
1978
1979

1980
1981
1982
1983
1984
1985

1986
1987
1988

was annotated using a multi-questions annotation schema in which one of the questions was “does the tweet appear to contain false information?”. Based on the answer to this question factuality label of the tweet has been defined. The Arabic dataset contains a total of 4,966 tweets.

A.7.7 Harmful Content Detection

For the harmful content detection we adopted the task proposed in (Alam et al., 2021b; Nakov et al., 2022b) though the research on harmful content detection also include identifying or detecting offensive, hate-speech, cyberbullying, violence, racist, misogynistic and sexist content (Sharma et al., 2022; Alam et al., 2022a). For some of the those harmful content detection tasks we addressed them separately and discussed in the below sections. Alam et al. (2021b); Nakov et al. (2022b) proposed the as in the context of tweets and idea was to detect whether the content of the tweet aims to and can negatively affect society as a whole, specific person(s), company(s), product(s), or spread rumors about them. The content intends to harm or *weaponize the information*¹⁷ (Broniatowski et al., 2018).

Dataset We used the Arabic dataset proposed in (Nakov et al., 2022b), which consists of a total of 6,155 annotated tweets.

A.7.8 Offensive Language Detection

The use of offensive language in social media has become a major problem, which can lead to real-world violence (Husain and Uzuner, 2021; Sap et al., 2019). This literature for offensive language detection mainly focused on social media content and addressing for variety of languages. The task is mainly defined as whether the content (e.g., text, image, or multimodal) is offensive or not (Chowdhury et al., 2020c).

Dataset For this task, we used the dataset from the SemEval-2020 Task 12 (OffensEval 2020) (Zampieri et al., 2020), which consists of 10,000 tweets, collected from a set of 660k Arabic tweets containing the vocative particle (“yA” – O) from April 15 to May 6, 2019.

A.7.9 Hate Speech Detection

Davidson et al. (2017) defined hate speech as “as language that is used to expresses hatred towards a

¹⁷The use of information as a weapon to spread misinformation and mislead people.

targeted group or is intended to be derogatory, to humiliate, or to insult the members of the group”. The literature for hate speech detection defined the task as detecting hate vs. non-hate from different types of content such as text, image and multimodal (Schmidt and Wiegand, 2017; Kiela et al., 2020; Gomez et al., 2020).

Dataset For this task, we also used the OSACT 2020 dataset (Mubarak et al., 2020c), which consists of 10,000 tweets with annotated label hate-speech, not-hate-speech.

A.7.10 Adult Content Detection

Identifying this type of content is important for social media platforms to make a safe place for users. Especially this type of content poses a serious threat to other vulnerable groups (e.g., younger age groups). The task typically involves detecting and identifying whether the textual content contains sensitive/adult content or account that share such content.

Dataset We used the dataset discussed in (Mubarak et al., 2021a), which contains 10,000 tweets collected by first identifying Twitter accounts that post adult content. Tweets are manually annotated as adult and not-adult.

A.7.11 Spam Detection

Spam content in social media includes ads, malicious content, and any low-quality content (Ghanem et al., 2023). Spam detection is another important problem as such content may often annoy and mislead the users (Gao et al., 2012).

Dataset We used the dataset discussed in (Mubarak et al., 2020a) for Arabic spam detection which contains 28K tweets manually labeled as spam and not-spam.

A.8 Semantic textual similarity

A.8.1 Textual Similarity

Semantic textual similarity is a measure used to determine if two sentences are semantically equivalent. The task involves generating numerical similarity scores for pairs of sentences, with performance evaluated based on the Pearson correlation between machine-generated scores and human judgments (Cer et al., 2017a). Two tasks were conducted to gauge the similarity between 250 pairs of Arabic sentences, as well as Arabic-English sentence pairs.

1989
1990
1991
1992
1993
1994
1995

1996
1997
1998
1999

2000
2001
2002
2003
2004
2005
2006
2007
2008

2009
2010
2011
2012
2013

2014
2015
2016
2017
2018
2019

2020
2021
2022
2023

2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035

2036
2037
2038
2039
2040

2041
2042
2043

2044
2045
2046
2047
2048
2049

2050

2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064

2065
2066
2067
2068
2069
2070
2071
2072
2073

2074

2075
2076
2077
2078

Dataset We used SemEval-2017 Task 1 (Track 1: ar-ar and Track 2: ar-en) dataset (Cer et al., 2017a), which is a translated version (machine translation followed by post-editing by human) of SNLI dataset (Bowman et al., 2015).

A.8.2 Semantic Question Similarity

The idea of this task is to determine how similar two questions are in terms of their meaning.

Dataset We used Mawdoo3 Q2Q dataset (NSURL-2019 task 8: Semantic question similarity in Arabic), which consists of 15,712 annotated pairs of questions. Each pair is labeled as *no semantic similarity (0)* or *semantically similar(1)* (Seelawi et al., 2019).

A.8.3 Natural Language Inference (NLI)

The XNLI task, known as Cross-lingual Natural Language Inference (Conneau et al., 2018), is a widely used benchmark in the field of natural language processing (NLP). It involves determining the logical relationship between pairs of sentences written in different languages. Specifically, the task requires NLP models to determine whether a given hypothesis sentence is entailed, contradicted, or neutral in relation to a given premise sentence, across multiple languages. The XNLI task serves as a rigorous evaluation of the cross-lingual transfer capabilities of NLP models, assessing their ability to understand and reason in different languages within a multilingual context.

Dataset The dataset we used for this study is the translated version of Arabic from XNLI corpus (Conneau et al., 2018). For the annotation, 250 English sentences were selected from ten different sources and then asked the annotators to produce three hypotheses per sentence premise. The resulting premises and hypotheses are then translated into 15 languages and we used the Arabic version for this study.

A.9 Question Answering (QA)

This task involves answering questions in Arabic based on a given text¹⁸. For this task, we use four different datasets consisting of (passage, question, and answer) pairs.

¹⁸This task is also referred to as machine reading comprehension where the model is tested on its ability to extract answers from the given text

Datasets 2079

ARCD consists of 1,395 Arabic MSA questions posed by crowd-sourced workers along with the text segments from Arabic Wikipedia. We use the test set only for our evaluation. The test set consists of 78 articles, 234 paragraphs, and 702 questions (Mozannar et al., 2019). 2080
2081
2082
2083
2084
2085

MLQA comprises multilingual question-answer instances in 7 languages, *English, Arabic, Simplified Chinese, Hindi, German, Vietnamese* and *Spanish*. We used the Arabic QA pairs from this dataset, which consist of 2389 articles, 4646 paragraphs, and 5335 questions (Lewis et al., 2019). 2086
2087
2088
2089
2090
2091

TyDi QA comprises 11 languages with 204K question-answer pairs. We used the data provided for the *Gold Passage task* in which a passage that contains the answer is provided and the task is to predict the span that contains the answer. We used the Arabic split of the data which contains 921 articles, 921 paragraphs and 921 questions (Artetxe et al., 2020). 2092
2093
2094
2095
2096
2097
2098
2099

XQuAD comprises 240 paragraphs and 1190 question-answers pairs from the development set of SQuAD v1.1 with their professional translations into ten languages. *Hindi, Turkish, Arabic, Vietnamese, Thai, German, Greek, Russian, Spanish* and *Chinese*. We use the the Arabic split of the data which consists of 48 articles, 240 paragraphs, and 1190 questions (Artetxe et al., 2020). We used the sQuAD version of all datasets along with the official squad evaluation script. 2100
2101
2102
2103
2104
2105
2106
2107
2108
2109

A.10 Speech Processing 2110

For this study, we address the speech modalities in the context of large foundation models, and we evaluate the following two tasks in this edition: (i) automatic speech recognition (ASR); and (ii) text to speech (TTS) models. In future, we will scale the speech benchmark with speech translation (ST) and spoken Arabic dialect identification spoken (ADI). 2111
2112
2113
2114
2115
2116
2117
2118

A.10.1 Speech Recognition 2119

The primary objective of an ASR system is to transform spoken language into written text. The task itself is challenging due to the presence of variability in human speech, which can be affected by factors such as accent, speaking style, code-switching, environmental factors like channels, and background noise among others. Furthermore, the 2120
2121
2122
2123
2124
2125
2126

Dataset	Task	Domain	Size
MGB2	ASR	Broadcast (MSA)	9.57 hrs
MGB3	ASR	Broadcast (EGY)	5.78 hrs
MGB5	ASR	Broadcast (MOR)	1.40 hrs
QASR.CS	ASR	Broadcast (Mixed) → Code-switching	5.90 hrs
DACS	ASR	Broadcast (MSA-EGY) → Code-switching	1.50 hrs
ESCWA.CS	ASR	Meeting (Mixed DA - ENG) → Code-switching	2.80 hrs
CallHome	ASR	Telephony (EGY)	20 phone conversations
In-house	TTS	Mixed Topics (education, health, etc)	20 sentences

Table 6: Summary on test sets and their sizes used in evaluation for the speech processing tasks.

presence of language-related challenges, including complex morphology, unstandardized orthography, and a wide array of dialects as a primary mode of communication, adds a layer of complexity to the task. Therefore to properly benchmark Arabic ASR, we covered a wide range of domains encapsulating different speaking styles, dialects, and environments. For our study, we considered broadcast news, telephony, and meeting data for MSA, Egyptian, Moroccan Arabic, etc., in both monolingual and code-switching setups.

Datasets

MGB2 consists of 9.57 hours of multi-dialect speech data that was collected from Aljazeera TV programs and manually transcribed. The data consists of a mix of Modern Standard Arabic (MSA) and various dialects, including Egyptian, Levantine, Gulf, and North African (Ali et al., 2016).¹⁹

MGB3 is a collection of 5.78 hours of multi-genre speech data in Egyptian dialect. The data was collected from YouTube videos and manually transcribed (Ali et al., 2017).²⁰

MGB5 is a collection of 1.4 hours of speech data in Moroccan dialect. The data was collected from YouTube videos and manually transcribed (Ali et al., 2019).²¹

ESCWA.CS is a collection of 2.8 hours of speech code-switching corpus collected over two days of meetings of the United Nations Economic and Social Commission for West Asia (ESCWA) in 2019 (Chowdhury et al., 2021).²²

QASR.CS is a collection of 5.9 hours of code-switching extracted from the Arabic broadcast

news data (QASR) to test the system for code-switching. The dataset also includes some instances where the switch is between Arabic and French, however, this type of instance are very rare occurrence (Mubarak et al., 2021b).²³

DACS is a collection of ≈ 1.5 hours of broadcast speech designed to evaluate the performance of ASR for code-switching between MSA to Egyptian dialect and vice versa (Chowdhury et al., 2020a).²⁴

CallHome Egyptian is a speech corpus of telephone conversations between native speakers of Egyptian Arabic. It consists of 20 unscripted telephone conversations, each of which lasts between 5-30 minutes (Kumar et al., 2014).²⁵

A.10.2 Text to Speech

Speech Synthesis a.k.a text to speech (TTS) helps users to get the written output easier and in some cases faster. Most state-of-the-art end-to-end TTS systems comprise three modules: text front-end, acoustic model, and vocoder. However, there is ongoing research to combine acoustic models and vocoder in a single neural network. Text front-end module normalizes input text by converting digits, symbols, abbreviations, and acronyms into full words, processing words with special sounds, borrowed words, etc. This task is challenging in Arabic due to missing diacritics in modern texts as explained in A.1.4. Therefore, the Arabic front-end part of the TTS is responsible for restoring the missing diacritics and text normalization.

Dataset For MSA TTS, we create the first public test dataset, which comprises 30 sentences covering different topics such as psychology, education, health, etc. The average length for each sentence

¹⁹<https://arabicspeech.org/mgb2>

²⁰<https://arabicspeech.org/mgb3>

²¹<https://arabicspeech.org/mgb5>

²²<https://arabicspeech.org/escwa>

²³<https://arabicspeech.org/qasr>

²⁴https://github.com/qcri/Arabic_speech_code_switching

²⁵<https://catalog.ldc.upenn.edu/LDC97S45>

is 8 words. This data is used for objective and subjective evaluation for Arabic TTS.

B Model Parameters

B.1 NLP Models

We used gpt-3.5-turbo-0301 and gpt-4-0314 versions for our tasks. In addition we used Bloomz 176B 8-bit version.

B.2 Speech Models

In Table 7, we provide the details of the speech model parameters.

Model	Layers	Width	Heads	Parameters
W.Small	12	768	12	244M
W.Medium	24	1024	16	769M
W.Large-v2	32	1280	20	1550m
USM	32	1526	16	2B

Table 7: Model parameters and architecture for Large pretrained ASRs. W. stands for Open.AI’s Whisper (Radford et al., 2022) and USM is Universal Speech Model from Google (Zhang et al., 2023)

C Prompts

The performance of the model is highly dependent on the prompting strategy. Designing the best prompts for each task is challenging and required several iterations. In many tasks, the output was not consistent for all instances of the datasets. For example, in many cases the model provides the desired labels however, there are cases where the model output different kind of error messages: (i) it’s trained only on English and cannot handle Arabic texts, (ii) the response was filtered due to the prompt triggering Azure OpenAI’s content management policy, (iii) it often provided extra tokens or swapped the tag (B-PER to PER-B). These resulted in an extra layer of post-processing and filtering of the evaluation dataset. Moreover, from our initial exploration, we noticed that, compared to language-specific (Arabic) prompts, English prompts (task-description) provide superior performance. Our underlying hypothesis is that with English task-description the input representations shift toward the English space that allows the model to process and understand the input better, giving better performance.²⁶

²⁶Note this observation aligns with other multilingual low-resource language studies.

For the segmentation task, with our initial prompt, we realized that the output was not segmented based on linguistic information but rather more Byte-Pair Encoding (BPE) like encoding. Based on that prompt is further redesigned, which resulted in a better outcome.

For factuality, disinformation, and harmful content detection tasks, the challenges were different from other tasks. One notable example is the propaganda detection task. The task requires determining whether a text snippet contains propagandistic language, and if it does, the model should detect which propaganda technique is used from a pre-defined list of techniques. Even with our best efforts to design the prompt for this task, the model still produced very unexpected responses, sometimes incomplete names of propaganda techniques, or even techniques not among the provided list. Another challenge with designing prompts for these tasks, is the issue of a task’s subjectivity where providing a crisp-clear classification task definition to the model is not possible. As an example, one of our tasks is to evaluate whether a tweet is offensive towards a person or an entity. In many instances, the model predicted tweets to be offensive, while in reality they were descriptive of the tweet’s author mental or physical state, or they were just repeating common negative statements or Arabic proverbs not directed at anyone indicating the model’s understanding of offensiveness is not inline of our definition.

In the following sections, we report the prompts we used for different tasks.

C.1 Word Segmentation, Syntax and Information Extraction

Segmentation

A word can be composed of one root and one or multiple affixes. Segment the following sentence into its morphological constituents: {inputSentence}"+". The output format should be a list of tuples, where each tuple consists of a word from the input text and its segmented form joined by a + sign.

Named Entity Recognition

Task Description: You are working as a named entity recognition expert and your task is to label a given arabic text with named entity labels. Your task is to identify and label any named entities present in the text without any explanation. The named entity labels that you will be using are PER (person), LOC (location),

2282	ORG (organization), MISC (miscellaneous). You	put them back as they were.	2342
2283	may encounter multi-word entities, so make sure		
2284	to label each word of the entity with the		
2285	appropriate prefix ('B' for first word entity,		
2286	'I' for any non-initial word entity). For words		
2287	which are not part of any named entity, you		
2288	should return 'O'. Note: Your output format		
2289	should be a list of tuples, where each tuple		
2290	consists of a word from the input text and its		
2291	corresponding named entity label. Input:		
2293	{inputSentence}		
<hr/>			
2294	POS		
<hr/>			
2295	These are the segmentation and POS tags for a		
2296	sample sentence:		
2297			
2298	فيلم جاذبية يتصدر ترشيحات جوائز الأكاديمية البريطانية		
2299	لفنون الفيلم والتلفزيون		
2300	فيلم NOUN		
2301	جاذبية NOUN+NSUFF		
2302	يتصدر V		
2303	ترشيحات NOUN+NSUFF		
2304	جوائز NOUN		
2305	الأكاديمية ال + أكاديمي + ة DET+NOUN+NSUFF		
2306	البريطانية ال + بريطاني + ة DET+ADJ+NSUFF		
2307	ل فنون لفنون PREP+NOUN		
2308	الفيلم ال + فيلم DET+NOUN		
2309	و التلفزيون والتلفزيون CONJ+DET+NOUN		
2310			
2311	get the segmentation and POS tags for this		
2313	sentence: {inputSentence}		
<hr/>			
2314	Assign POS tag to each morphological segment		
2315	within each word. group the tags for each word		
2316	with +: {inputSentence}"+". The output should		
2317	be in the format: [{word: label}, {word: label}]		
2319			
<hr/>			
2320	Label the following sentence with its		
2321	corresponding PENN Treebank POS Labels.		
2322	sentence: {inputSentence}		
2323	labels:		
2324			
<hr/>			
2326	Lemmatization		
<hr/>			
2327	for every word in the following sentence, write		
2328	only the lemmas without diacritics in separate		
2329	lines without explanation:		
2330	{inputSentence}		
2331			
<hr/>			
2333	Diacritization		
<hr/>			
2334	Diacritize fully the following Arabic sentence:		
2335	{inputSentence}		
2339			
<hr/>			
2338	Vowelized the following sentence:		
2339	{inputSentence}. Words that can't be vowelized		
2340			
<hr/>			
	Parsing		2343
<hr/>			
	Given the following features (in order: ID,		2344
	Form, Lemma, CPostTag, POSTag, Features),		2345
	predict the Head of each token in the following		2346
	sentence, which is either a value of a related		2347
	ID or 0. A value of zero means the token		2348
	attaches to the virtual root node:		2349
	{inputSentence}		2350
			2352
<hr/>			
	Dialect Identification		2353
<hr/>			
	Write only the country code of the Arabic		2354
	country in which this sentence is written in		2355
	its dialect without any explanation? Write only		2356
	the country code in ISO 3166-1 alpha-2 format		2357
	without explanation. Write 'MSA' if the		2358
	sentence is written in Modern Standard Arabic.		2359
	sentence: {inputSentence}		2360
	code:		2361
			2363
<hr/>			
	C.2 Sentiment, Stylistic and Emotion Analysis		2364
<hr/>			
	Sentiment analysis		2365
<hr/>			
	Choose only one sentiment between: Positive,		2366
	Negative, Neutral, or Mixed for this sentence:		2367
	sentence: {inputSentence}		2368
	label:		2369
			2370
<hr/>			
	Emotion detection		2372
<hr/>			
	Predict all the possible emotions in the		2373
	following Arabic sentences without explanation		2374
	and put them in a Python list. List of emotions		2375
	are: anger, anticipation, disgust, fear, joy,		2376
	love, optimism, pessimism, sadness, surprise,		2377
	and trust		2378
	sentence: {inputSentence}		2379
	labels:		2380
			2382
<hr/>			
	C.3 Demographic/Protected Attributes		2383
<hr/>			
	Gender		2384
<hr/>			
	If the following person name can be considered		2385
	as male, write 'm' without explanation, and if		2386
	it can be considered as female, write 'f'		2387
	without explanation.		2388
	person name: {inputSentence}		2389
	label:		2390
			2392
<hr/>			
	Location		2393
<hr/>			
	Map the following locations to one of the Arab		2394
	countries. Write the country code in ISO 3166-1		2395
	alpha-2 format without explanation. If the		2396
	country is outside Arab countries, write		2397
	'OTHERS', and if the location cannot be mapped		2398
			2399

2400	to any country in the world, write 'UNK'	Checkworthiness	2461
2401	without any explanation.		2462
2402	location: {inputSentence}	Classify the sentence as checkworthy or not	2463
2403	label:	checkworthy. Provide only the label.	2464
		sentence: {inputSentence}	2465
		label:	2466
2405	Name Info		
2406	Predict the country of citizenship of the	Claim detection	2468
2407	following person name. Write the country code		
2408	in ISO 3166-1 alpha-2 format without	Does this sentence contain a factual claim?	2469
2409	explanation.	Answer only by yes or no.	2470
2410	name: {inputSentence}	sentence: {inputSentence}	2471
2411	code:	label:	2472
2413			2473
2414	C.4 Ethics and NLP: Factuality,	Harmful content detection	2475
2415	Disinformation, Harmful content		
		Classify the following sentence as harmful or	2476
2416	Offensive Language	not harmful. Answer only by yes or no. Provide	2477
		only label.	2478
2417	If the following sentence is offensive, just	sentence: {inputSentence}	2479
2418	write "OFF", otherwise, just write "NOT_OFF"	label:	2480
2419	without explanation:		2481
2420	sentence: {inputSentence}		
2421	label:	Attention-worthy	2483
2423			
		Classify the sentence by whether it should get	2484
2424	Hate Speech	the attention of policymakers. Answer by yes or	2485
		no. If the predicted label is yes then classify	2486
2425	If the following sentence has hate speech, just	the sentence into one of the following	2487
2426	write "HS", otherwise, just write "NOT_HS"	categories: asks question, blame authorities,	2488
2427	without explanation:	calls for action, Harmful, contains advice,	2489
2428	sentence: {inputSentence}	discusses action taken, discusses cure, or	2490
2429	label:	other.	2491
2430		text: {input_sample}	2492
		label:	2493
			2494
2432	Adult Content		
		C.5 Semantics	2496
2433	Classify the following Arabic sentence as adult		
2434	language (the language used in adult	Semantic Textual Similarity	2497
2435	advertisement and porno advertisement) or not		
2436	adult language without illustration. In case	Given two sentences, produce a continuous	2498
2437	of adult language, just write "ADULT" without	valued similarity score on a scale from 0 to 5,	2499
2438	explanation, and in case of not adult	with 0 indicating that the semantics of the	2500
2439	language, just write "NOT_ADULT" without	sentences are completely independent and 5	2501
2440	explanation.	indicating semantic equivalence. The output	2502
2441	text: {inputSentence}	should be exactly in the form of a similarity	2503
2442	label:	score.	2504
2443		sentence 1: {inputSentence1}	2505
		sentence 2: {inputSentence2}	2506
		score:	2507
2445	Spam		2508
		Natural Language Inference	2510
2446	If the following sentence can be classified as		
2447	spam or contains an advertisemnt, write 'ADS'	You are provided with a premise and a	2511
2448	without explniation, otherwise write 'NOTADS'	hypothesis. Your task is to classify the	2512
2449	without explanantion.	hypothesis as true (entailment), false	2513
2450	sentence: {inputSentence}	(contradiction), or unknown (neutral) based on	2514
2451	label:	the given premise. The output should be true,	2515
2453		false or unknown.	2516
		premise: {inputSentence1}	2517
2454	Subjectivity	hypothesis: {inputSentence2}	2518
		output:	2519
2455	Classify the sentence as subjective or		2520
2456	objective. Provide only label.		
2457	sentence: {inputSentence}		
2458	label:		
2460			

2522	Classification (Question Similarity)	2575
2523	<hr/>	
2524	Are the following two questions semantically	2576
2525	similar? The output should be exactly either	2577
2526	yes or no.	
2527	question 1: {inputQuestion1}	2578
2528	question 2: {inputQuestion2}	2581
2530	label:	2582
	<hr/>	

2531	C.6 Question answering (QA)	2583
2532	<hr/>	
2533	Your task is to answer questions in Arabic	2584
2534	based on a given context.	
2535	Note: Your answers should be spans extracted	2585
2536	from the given context without any	2586
2537	illustrations.	2587
2538	You don't need to provide a complete answer.	2588
2539	context:{context}	2589
2540	question:{question}	2590
2541	answer:	2591
	<hr/>	

2543	D Post-processing	2592
2544	Post-processing was needed for almost all tasks	
2545	in order to match gold labels, which include refo-	2593
2546	rmating the output handling exceptions, missing	2594
2547	values, and unexpected values. Much like NLP	2595
2548	tasks, post-processing the transcription output from	2596
2549	the speech models is an important step. We noti-	2597
2550	cied that the performance of the Whisper models	2598
2551	is highly dependent on the post-processing. As the	2599
2552	models (Whisper family) are trained with massive	2600
2553	dataset created by weak supervision, the output is	2601
2554	quite noisy and needs extra care for post-processing.	2602
2555	In this study, we opt for a simple post-processing	2603
2556	pipeline so that the process is not overfitted to task-	2604
2557	based data styles.	2605

2558	E Benchmarks on Arabic: Details	2606
2559	Noteworthy among them is ORCA (Elmadany	
2560	et al., 2022), a large-scale benchmark that incor-	2607
2561	porates 60 diverse datasets organized into seven	2608
2562	comprehensive task clusters. This large-scale or-	2609
2563	ganization allows for a more in-depth and diverse	2610
2564	analysis of model performance across a multitude	2611
2565	of language tasks including but not limited to sen-	2612
2566	tence classification, text classification, structured	2613
2567	prediction, semantic similarity, natural language	2614
2568	inference, question-answering, and word sense dis-	2615
2569	ambiguation.	2616
2570	<i>AraBench</i> (Sajjad et al., 2020) is an evaluation	2617
2571	suite for dialectal Arabic-to-English machine trans-	2618
2572	lation. It offers a wide range of dialect categories	2619
2573	including 4 coarse, 15 fine-grained, and 25 city-	2620
2574	level dialects from various genres like media, chat,	2621

and travel. It also provides robust baselines that	2575
utilize different training methods like fine-tuning,	2576
back-translation, and data augmentation.	2577

The <i>ALUE</i> (Seelawi et al., 2021) benchmark of-	2578
fers 8 curated tasks and private evaluation datasets,	2579
covering areas like emotion classification, hate	2580
speech, and fine-grained dialect identification. Ara-	2581
bicBERT tops the performance in 7 of these 8 tasks,	2582
with evaluations also including BERT variants with	2583
AraVec and FastText models.	2584

<i>ARLUE</i> (Abdul-Mageed et al., 2021) bench-	2585
mark employs 42 datasets for six task clusters	2586
to evaluate multi-dialectal Arabic language under-	2587
standing, featuring BERT and XLM model variants.	2588
Fine-tuned models utilizing ARLUE lead the per-	2589
formance in all six clusters.	2590

As shown in Table 8, Our study provides a com-	2591
prehensive evaluation platform that advances the	2592
current benchmarks by presenting 33 distinct tasks	2593
over 61 datasets, which is the most extensive task	2594
coverage among current benchmarks. Unlike the	2595
AraBench, which focuses exclusively on Arabic-to-	2596
English translation tasks, and ALUE and ARLUE,	2597
which have a narrower task focus or a lesser num-	2598
ber of tasks, LAraBench provides a broader scope	2599
of evaluation tasks. This benchmark encompasses	2600
a multitude of language tasks that are paramount to	2601
understanding the robustness and generalizability	2602
of language models. Furthermore, LAraBench dis-	2603
tinguishes itself by not only including text modal-	2604
ity but also speech modality, thereby increasing	2605
the robustness and utility of our benchmark. Ad-	2606
ditionally, we successfully implemented GPT-3.5	2607
and GPT-4, demonstrating its compatibility with	2608
cutting-edge language models.	2609

Notably, the models employed in LAraBench	2610
have displayed comparable performance with the	2611
SOTA models, attesting to its robustness and high	2612
standard of evaluation. While SOTA models gen-	2613
erally outperform LLMs, our benchmark reveals	2614
that these LLMs can close the performance gap in	2615
certain tasks, particularly when increasing prompt	2616
complexity and transitioning from zero-shot to few-	2617
shot learning. This highlights LAraBench’s utility	2618
not only as a tool for model evaluation but also as	2619
an instrumental platform for identifying tasks under	2620
which LLMs might be able to match or even sur-	2621
pass SOTA performance. This benchmark serves	2622
as a challenging testbed for future language mod-	2623
els and contributes to the advancement of Arabic	2624
language understanding models.	2625

Reference	# tasks	# datasets	Fine-tuned Models	Zero-shot GPT-3.5	Few-shot GPT-3.5	Zero-shot GPT-4	Few-shot GPT-4	Zero-shot Bloomz	SOTA Comp.	Modality
AraBench (Sajjad et al., 2020)	1	6	Seq2Seq (transformer)	✗	✗	✗	✗	✗	✓	T, S
ARLUE (Abdul-Mageed et al., 2021)	13	42	ARBERT, MARBERT	✗	✗	✗	✗	✗	✓	T
ALUE (Seelawi et al., 2021)	8	8	AraBERT, mBERT, mBERT, ARBERT,	✗	✗	✗	✗	✗	✓	T
ORCA (Elmadany et al., 2022)	29	60	CamelBERT, MARBERT	✗	✗	✗	✗	✗	✓	T
GPTAraEval (Tawkat et al., 2023)	32	60	✗	✓	✓	✗	✓	✗	✗	T
LAraBench (Ours)	33	61	✗	✓	✗	✓	✓	✓	✓	T, S

Table 8: A comparison with prior studies. T: Text, S: Speech.

F Extended Experiments and Results

In this section, we provide extended versions of the results reported earlier in the paper.

F.1 Random Baseline

For different tasks, we used different approaches to compute random baseline, as discussed below.

- **Segmentation:** We first randomly decide how many segments a token should have (between 0, 1 and 2), and then randomly split the characters of that token into the chosen number of segments.
- **Lemmatization:** We first randomly decide the length of the lemma, and then randomly divide the remaining length between a prefix and suffix.
- **Diacritization:** we randomly choose between 9 choices for every character (8 diacritics and 1 choice for no diacritic).
- **QA:** Randomly select a span of tokens from the given context of each question.
- **Others (Multiclass and multilabel classification tasks):** For multiclass classification, we randomly assign a label to the test instance, with label selection based on the labels from the training set. For multilabel classification, which requires assigning multiple labels from a predefined set, both the number of labels and their selection were random, and these were assigned to the test instance.

F.2 Extended Few-shot Results

We conducted experiments using GPT-4 by incrementally increasing the number of shots. For this purpose, we chose one task from each of the seven groups listed in Table 1 in the paper. We tested the models using 3, 5, and 10 shots. For each task, we observed a general trend of increasing performance, with the exception of the gender task. On

average, performance improved from 0.656 in the 0-shot setting to 0.721 in the 10-shot setting. The results are presented in Table 9. To provide a clear overview of the comparison across different few-shot scenarios, we present the average performance in Figure 2.

Task Name	Metric	0-shot	3-shot	5-shot	10-shot
NER	Macro-F1	0.355	0.420	0.426	0.451
Sentiment	Macro-F1	0.569	0.598	0.619	0.639
News Cat.	Macro-F1	0.667	0.594	0.674	0.723
Gender	Macro-F1	0.868	0.980	0.931	0.937
Subjectivity	Macro-F1	0.677	0.745	0.740	0.771
XNLI (Arabic)	Acc	0.753	0.774	0.789	0.809
QA	F1 (exact match)	0.705	0.704	0.718	0.716
Average		0.656	0.688	0.700	0.721

Table 9: Results from few-shot experiments over seven tasks with GPT-4.

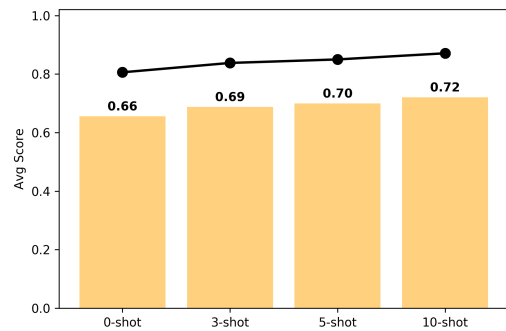


Figure 2: An average performance comparison (over seven tasks) of different few-shot experiments with GPT-4.

F.3 Native Language Prompts

We have conducted experiments using Arabic prompts for the *seven selected tasks*. The Arabic prompts were created by native Arabic speakers. The results are reported in Table 10. Using the Arabic prompts, three out of the seven tasks outperformed their counterparts that used English prompts, two underperformed, and one showed

equivalent performance. This finding partially supports the findings reported by Ahuja et al. (2023), which states that “the monolingual prompting setup outperforms the cross-lingual prompting strategy”. However, they also report that using Davinci-003, the English prompts yield better results than their translated version in the native language.

Task Name	Metric	English	Arabic
NER	Macro-F1	0.355	0.350
Sentiment	Macro-F1	0.569	0.547
News Cat.	Macro-F1	0.667	0.739
Gender	Macro-F1	0.868	0.892
Subjectivity	Macro-F1	0.677	0.725
XNLI (Arabic)	Acc	0.753	0.740
QA	F1 (exact match)	0.705	0.654
Average		0.656	0.664

Table 10: Results from GPT-4 using zero-shot prompts in both English and native languages.

F.4 Semantic vs. Syntactic Task Differences

We computed the performance difference between POS and MT, as shown in Table 11. The gap between SOTA and the three LLMs for POS (a syntactic task) is considerably larger than for MT (a semantic task). Moreover, the performance gap is much lower for semantic tasks compared to syntactic tasks, on average, across the three LLMs, as depicted in Table 11. This implies that these models might be better equipped to encode and express semantic information than to handle specific syntactic phenomena in their inputs.

	BLOOMZ	GPT-3.5	GPT-4	SOTA
Semantic				
MT	19.38	24.09	23.57	24.58
Semantics (STS, XNLI)	0.615	0.733	0.827	0.794
Syntactic				
POS	-	0.154	0.464	0.844
Parsing	-	0.239	0.504	0.796

Table 11: Average performance difference between semantic and syntactic tasks.

F.5 Data Contamination Assessment

The presence of test data from standard downstream NLP tasks in the training dataset of pre-trained LLMs’ may effect the evaluations. It is important to have blind test-sets to reliably assert that the models are not merely memorizing data patterns but have truly acquired the ability to generalize. Identifying whether the data has been contaminated or not is a challenging problem. In our

study, we have used the dataset that has been released after September 2021, which is a cut-off date for OpenAI’s GPT models.²⁷ The tasks include CT-CWT-22 tasks (Checkworthy, Claim, Harmful content, and Attention-worthy) introduced in 2022. Consequently, for these specific tasks, the potential for data contamination is none. Both GPT-3.5 and GPT-4 (in zero-shot and 3-shot scenarios) demonstrate results closely aligned with the state-of-the-art, mirroring trends seen in other 2021 test sets. In addition, the dataset for the subjectivity task is our in-house developed dataset, created at the end of 2022.

To further validate whether evaluation datasets have been exposed to the LLMs, we assessed various datasets using the methodology outlined in (Golchin and Surdeanu, 2023). This approach employs tailored instructions to ascertain if a model has encountered particular evaluation data. When applying this methodology to GPT-4, across a representative array of datasets, namely (1) Sentiment (ArSAS 2018), (2) Emotion (SemEval-2018 Task 1, Arabic), (3) Sarcasm (ArSarcasm-OSACT2020, ArSarcasm-v2-WANLP2021), (4) News Category (ASND 2020), (5) Gender (Arap-Tweet 2022), (6) Subjectivity (In-house 2022), (7) XNLI 2020 (Arabic), (8) Question Answering (XQuAD 2019), no instances were generated by GPT-4 from these datasets. For none of the 9 datasets and 8 tasks, GPT-4 could produce any example from it. Thus, based on these experiments, we can conclude that the Arabic datasets for different tasks are not included in the training data of GPT models.

F.6 Machine Translation (MT)

In Table 12, we report detailed results for MT, considering both dialect and city levels.

²⁷<https://platform.openai.com/docs/models/overview>

Dataset	Dialect	SC	City	#Sent	BloomZ	Zero-shot GPT-3.5	Zero-shot GPT-4	SOTA
APT	LEV	lv	-	1000	11.38	18.55	17.77	21.9
APT	Nile	eg	-	1000	12.95	21.58	18.99	22.6
MADAR	Gulf	iq	Baghdad	2000	30.99	32.47	34.83	29.1
MADAR	Gulf	iq	Basra	2000	29.63	32.92	34.72	29
MADAR	Gulf	iq	Mosul	2000	29.17	30.82	35.32	31.3
MADAR	Gulf	om	Muscat	2000	39.91	39.37	39.9	39.5
MADAR	Gulf	qa	Doha	2000	31.1	33.6	33.62	29.3
MADAR	Gulf	sa	Jeddah	2000	40.37	42.62	42.69	29.4
MADAR	Gulf	sa	Riyadh	2000	27.73	32.51	33.71	40.7
MADAR	Gulf	ye	Sana'a	2000	29.79	32.48	34.63	31.4
MADAR	LEV	jo	Amman	2000	35.56	35.09	36.24	35.1
MADAR	LEV	jo	Salt	2000	34.54	35.78	37.54	34.9
MADAR	LEV	lb	Beirut	2000	24.01	26.14	28.95	23.7
MADAR	LEV	ps	Jerusalem	2000	34.02	35.22	35.5	33.6
MADAR	LEV	sy	Aleppo	2000	30.92	34.09	35.47	34.3
MADAR	LEV	sy	Damascus	2000	29.1	34.19	37.74	33.1
MADAR	MGR	dz	Algiers	2000	23.13	22.43	25.95	21.3
MADAR	MGR	ly	Benghazi	2000	25.41	26.99	30.12	32
MADAR	MGR	ly	Tripoli	2000	30.05	32.82	38.63	25.9
MADAR	MGR	ma	Fes	2000	23.73	22.53	26.15	29.9
MADAR	MGR	ma	Rabat	2000	31.02	31.95	34.71	23.1
MADAR	MGR	tn	Sfax	2000	15	15.93	20.74	13.8
MADAR	MGR	tn	Tunis	2000	16.79	14.69	18.51	16
MADAR	MSA	ms	-	2000	42.33	37.55	37.67	43.4
MADAR	Nile	eg	Alexandria	2000	29.24	32.05	32.46	38.3
MADAR	Nile	eg	Aswan	2000	39.97	41.77	42.42	30.4
MADAR	Nile	eg	Cairo	2000	32.79	32.77	32.69	32.9
MADAR	Nile	sd	Khartoum	2000	37.48	41.27	44.13	39
MDC	LEV	jo	-	1000	10.43	17.75	16.96	17.7
MDC	LEV	ps	-	1000	9.32	15.72	14.22	15.3
MDC	LEV	sy	-	1000	10.24	18.66	16.96	19.9
MDC	MGR	tn	-	1000	8.28	14.46	14.2	13.9
MDC	MSA	ms	-	1000	15.75	21.05	19.34	20.4
Media	Gulf	om	-	467	14.22	22.68	22.76	19.6
Media	LEV	lb	-	250	7.54	17.65	16.65	16.8
Media	MGR	ma	-	526	4.87	11.58	10.2	9.6
Media	MSA	ms	-	637	22.14	37.87	34.41	29.7
Media	MSA	ms	-	621	19.17	32.8	32.73	35.6
QAraC	Gulf	qa	-	6713				16
Bible	MGR	ma	-	600	16.34	16.16	15.14	28.8
Bible	MGR	tn	-	600	17.83	17.27	15.43	29.2
Bible	MSA	ms	-	600	24.37	23.96	18.38	33.2
Bible	MSA	ms	-	600	21.44	20.2	16.68	29.2

Table 12: Results (BLEU score) on machine translation for different datasets using zero-shot prompts. #Sent. indicates number of sentences in test set. SOTA results are reported in (Sajjad et al., 2020).