

Evaluating Large Language Models in Arabic Tasks: A Survey of Benchmarks, Methods, and Gaps

Anonymous ACL submission

Abstract

This survey provides the first systematic review of Arabic benchmarks that target the evaluation of LLMs with Arabic support. It showcases an analysis of 40+ Arabic evaluation benchmarks across NLP tasks, knowledge domains, cultural understanding, and specialized capabilities. A taxonomy organizing benchmarks into four categories is proposed: Knowledge, NLP Tasks, Culture and Dialects, and Target-Specific evaluations. The analysis reveals significant progress in Arabic benchmark diversity while identifying critical gaps: limited temporal evaluation, insufficient multi-turn dialogue assessment, and cultural misalignment in translated datasets. Three primary approaches are examined: native collections, translated collections, and synthetically generated collections. Their trade-offs regarding authenticity, scale, and cost are discussed. This work serves as a comprehensive reference for Arabic NLP researchers, providing insights into benchmark methodologies, reproducibility standards, and evaluation metrics while offering recommendations for future development.

1 Introduction

In recent years, Large Language Models (LLMs) have achieved major advances in natural language understanding and reasoning, moving closer to the vision of AGI (Naveed et al., 2025). Since the transformer era, LLMs have demonstrated strong multilingual capabilities beyond English (Huang et al., 2025). Arabic, spoken by nearly 500 million people worldwide¹, underscores the importance of multilingual evaluation. Consequently, both Arabic-specific and multilingual LLMs with Arabic support have been released in open- and closed-source settings (Al-Khalifa et al., 2025). Robust benchmarks are essential for systematic

¹World Bank (2024). Population, Total – Arab World. <https://data.worldbank.org/indicator/SP.POP.TOTL?locations=1A>

evaluation, and in parallel with progress in Arabic-capable LLMs, substantial effort has been devoted to Arabic benchmarking.

Arabic benchmark development faces distinct challenges. Data scarcity and limited diversity in Arabic web content (Al-Khalifa et al., 2025) increase the cost of dataset creation. To address this, researchers rely on translation from English, synthetic data generation using LLMs, and native Arabic data collection. Each approach involves trade-offs: translation often leads to cultural misalignment, while synthetic data risks bias and circular evaluation. Both require extensive human validation and remain vulnerable to cultural misalignment (Nacar et al., 2025).

These challenges are amplified by Arabic’s linguistic diversity across Modern Standard Arabic and numerous regional dialects (Keleg et al., 2025). Hence, Arabic demands dedicated benchmarking analysis due to its rich morphology, more than 20 dialects that function almost as separate languages, strong cultural sensitivity that limits translation-based evaluation, persistent data scarcity, and a fragmented evaluation landscape absent in English-centric benchmarks.

While recent surveys have analyzed Arabic-capable LLMs from various perspectives (Mashaabi et al., 2024; Al-Khalifa et al., 2025; Rhei and Roussinov, 2025), none provides a focused and comprehensive treatment of Arabic benchmarks. This survey fills that gap by systematically reviewing Arabic evaluation datasets and methodologies. Our contributions are threefold: (1) we introduce a taxonomy organizing benchmarks into four categories; (2) we analyze over 40 Arabic benchmarks; and (3) we examine evaluation practices, tools, trends, and critical gaps. We release a companion repository² consolidating benchmark datasets, code, and frameworks.

²Anonymous-github-Link

079	The paper is organized as follows. Section 2 re-	3 Benchmarks Taxonomy	128
080	views Arabic-supporting LLMs and related surveys.	We reviewed 40+ existing Arabic benchmarks	129
081	Section 3 introduces the taxonomy, with Sections	and constructed a taxonomy that captures various	130
082	4–7 detailing each category. Section 8 discusses	themes and categories, as depicted in Figure 1. Be-	131
083	open challenges, and Section 9 concludes.	fore introducing the taxonomy, we present the re-	132
084	2 Background	search methodology and the criteria used for bench-	133
085	This section provides the necessary context for un-	mark inclusion.	134
086	derstanding the Arabic benchmarking landscape for	3.1 Inclusion Criteria	135
087	LLMs. We first categorize existing LLMs that sup-	While surveying the current state of benchmarks	136
088	port Arabic or are Arabic-specific, by their training	used to evaluate LLMs in Arabic tasks, we devised	137
089	approaches, then review related work to position	inclusion criteria that determined which bench-	138
090	our contribution within the broader literature.	marks would be discussed in this survey paper. All	139
091	2.1 LLMs with Arabic Support	existing works mentioned fall under one of the fol-	140
092	State-of-the-art LLMs with Arabic support can be	lowing categories, which we discuss below:	141
093	categorized into three types. Native models are	<ul style="list-style-type: none"> • Existing Arabic benchmarks introduced in 	142
094	trained from scratch exclusively on Arabic data, en-	academic papers that evaluate specific Arabic	143
095	abling Arabic-only interaction. Examples include	capabilities in LLMs. Public availability of	144
096	Jais (Sengupta et al., 2023), ArabianGPT (Koubaa	the benchmark dataset or evaluation pipeline	145
097	et al., 2024), and AraGPT (Antoun et al., 2021).	is not a requirement for inclusion.	146
098	Multilingual models support multiple languages	<ul style="list-style-type: none"> • Arabic benchmarks that lack a detailed techni- 	147
099	including Arabic, such as Qwen3 (Yang et al.,	cal report but are well established in the field,	148
100	2025), Gemma3 (Team et al., 2025b), and Llama	as evidenced by their use in evaluating re-	149
101	(Grattafiori et al., 2024), as well as closed-source	leased LLMs or their presence in well-known	150
102	models like ChatGPT (OpenAI, 2023) and Claude	leaderboards.	151
103	(Anthropic, 2024). Adapted Arabic models apply	<ul style="list-style-type: none"> • Arabic subsets of established multilingual 	152
104	continued pretraining or supervised fine-tuning to	benchmarks, such as MMLU and EXAMS.	153
105	existing multilingual models to enhance Arabic per-	3.2 Taxonomy	154
106	formance, exemplified by AceGPT (Huang et al.,	Our proposed taxonomy organizes benchmarks into	155
107	2024), SILMA (Team, 2024), Fanar (Team et al.,	four categories:	156
108	2025a), and Falcon-Arabic (TII, 2025).	Knowledge includes benchmarks evaluating gen-	157
109	2.2 Related Work	eral knowledge and STEM capabilities, along with	158
110	Several surveys have examined LLMs that have	domain-specific benchmarks in fields such as law	159
111	Arabic capabilities. Mashaabi et al. (2024) pro-	and medicine.	160
112	vided an in-depth discussion of pretraining and fine-	Natural Language Processing (NLP) encom-	161
113	tuning data for LLMs with Arabic support, includ-	passes early task-specific benchmarks and com-	162
114	ing dialectal coverage, and listed available models	prehensive multi-task benchmarks, reflecting the	163
115	with details on accessibility and reproducibility.	evolution from narrow task evaluation to unified	164
116	Rhel and Roussinov (2025) surveyed pretrained	assessment across diverse dialects and domains.	165
117	Arabic LLMs with focus on classical NLP applica-	Culture and Dialects groups benchmarks assess-	166
118	tions and benchmarks. Most recently, Al-Khalifa	ing cultural knowledge and dialect understanding,	167
119	et al. (2025) presented the historical evolution of	addressing the essential property of cultural aware-	168
120	Arabic NLP, common pretraining and fine-tuning	ness in LLMs with Arabic capabilities.	169
121	strategies, and current research trends and chal-	Target-Specific covers benchmarks designed to	170
122	lenges. While Al-Khalifa et al. (2025) briefly dis-	assess particular LLM properties such as safety,	171
123	cussed a subset of existing benchmarks, no compre-	hallucination detection, instruction-following, and	172
124	hensive survey of Arabic benchmarks exists. This	vision capabilities.	173
125	work fills that gap by systematically reviewing eval-	This taxonomy emerged from analyzing com-	174
126	uation techniques and benchmarking datasets for	mon patterns across benchmarks and reflects the	175
127	LLMs with Arabic support.		

evolution from task-specific evaluation to comprehensive assessment. In the following sections, we describe each category and we provide a thorough table (Tables 6) listing all discussed benchmarks with their characteristics.

4 Knowledge

Paper	Primary Topic (Type)	Total
MMLU_ar	General Knowledge (MCQ)	15k
EXAMS_ar	General Knowledge (MCQ)	0.5k
ArabicMMLU	General Knowledge (MCQ)	14.57k
AraSTEM	STEM (MCQ)	11.63k
GAT	Linguistic Ability (MCQ)	0.56k
Qiyas	Linguistics & Math (MCQ)	2.4k
ArabLegalEval	Law (MCQ, GEN)	26.4k
AraMed	Medical (GEN)	270k
MizanQA	Law (MCQ)	1.7k
Fann or Flop	Poetry (GEN)	6.9k
GATmath, GATLc	Linguistics & Math (MCQ)	9k
3LM	STEM & Coding (MCQ)	3.1k
Arabic-GSM8K	Mathematics (GEN)	8.9k
MedArabiQ	Medical (MCQ)	0.7k
Hajji-FAQ	Religious QA (QA)	1.2k

Table 1: Arabic knowledge benchmarks summarized at a high level.

This section examines benchmarks evaluating LLMs’ acquired knowledge and reasoning capabilities, covering both general and STEM topics as well as specialized domains such as law, medicine, and poetry, as summarized in Table 1.

4.1 General and STEM

Multilingual efforts produced early general knowledge benchmarks with Arabic components. MMLU_ar (Hendrycks et al., 2020) comprises 14,079 human-translated MCQs spanning 57 subjects across difficulty levels, while EXAMS_ar (Hardalov et al., 2020) contains 562 high-school exam questions covering physics, chemistry, and biology. However, these suffer from translation concerns or limited scale.

ArabicMMLU (Koto et al., 2024) addressed these limitations with 14,575 native Arabic MCQs curated from educational exams across Arab countries, covering all school levels plus university, spanning STEM, social sciences, humanities, and Arabic language understanding. AraSTEM (Mustapha et al., 2024) further specialized in STEM with 11,637 native MCQs which is 7,000 more STEM samples than ArabicMMLU though it remains unpublished despite evaluation results on open-weight models.

The 3LM suite (Boussaha et al., 2025) combines three benchmarks totaling 3,151 questions: 3LM_nat from Arabic STEM exams, 3LM_syn

synthetically generated via Yourbench’s pipeline (Shashidhar et al., 2025) using Qwen3-235B-A22B, and 3LM_code comprising Arabic translations of HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021) instructions and comments.

Several benchmarks leverage Saudi Arabia’s General Aptitude Test (GAT), which assesses verbal abilities (reading comprehension, contextual errors, sentence completion) and quantitative skills (arithmetic, algebra, geometry, data analysis). Early efforts (Alkaoud, 2024; Al-Khalifa and Al-Khalifa, 2024) used 456 and 2,407 GAT samples respectively but lacked reproducibility and scale, evaluating only GPT-3.5 and GPT-4 with limited shots. GATmath and GATLc (AlBallaa et al., 2025) addressed these issues with 7k and 9k samples respectively (16k total), publicly released with 5-shot evaluation on diverse Arabic and bilingual LLMs. Arabic-GSM8K (Omartificial-Intelligence-Space, 2025) provides human-validated translations of the established GSM8K benchmark (Cobbe et al., 2021), assessing middle-school mathematical reasoning through 5-shot settings.

4.2 Domain Knowledge

Legal Domain. ArabLegalEval (Hijazi et al., 2024) pioneered Arabic legal LLM evaluation using documents scraped from Saudi Arabia’s Ministry of Justice and Board of Experts websites. The benchmark employs three approaches: synthetic MCQ generation using GPT-4 and Claude-3-opus with in-context examples from ArabicMMLU’s law section, QA pairs from governmental FAQ sections, and machine-translated datasets from Legalbench (Guha et al., 2023) verified by legal experts. MIZANQA (Bahaj and Ghogho, 2025) extends legal evaluation to Moroccan law using MCQs from various law exams.

Poetry and Linguistics. Fann or Flop (Alghalabi et al., 2025) assesses poetry understanding through 6,984 poem-explanation pairs, evaluating metaphorical and figurative comprehension. Evaluation uses BLEU, chrF(++), BERTScore, and mDeBERTaV3 for character-level overlap and semantic alignment, with GPT-4o as judge for faithfulness, grammatical correctness, and interpretive depth.

Medical Domain. Al-Majmar et al. (2024) introduced 808k medical QA samples from the Altibbi patient-doctor forum³. AraMed (Alasmari et al., 2024) refined this to 270k high-quality samples

³<https://altibbi.com/>

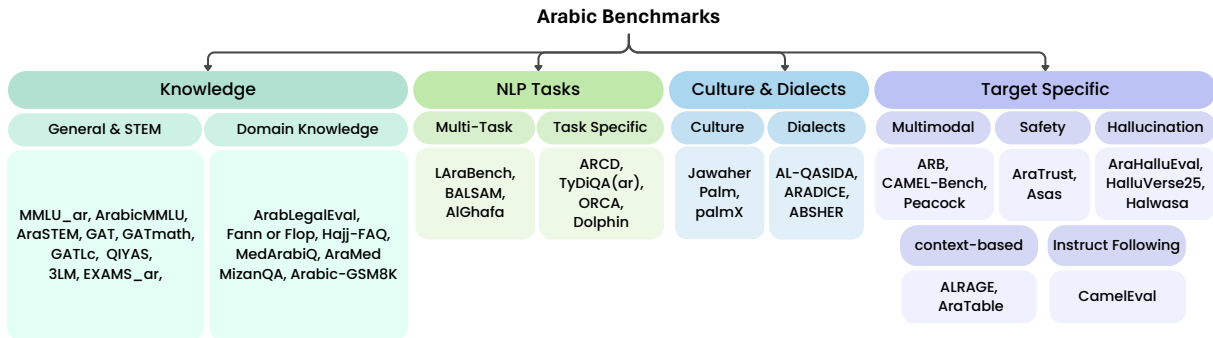


Figure 1: Taxonomy of Arabic Benchmarks.

based on vote counts. MedArabiQ (Daoud et al., 2025) further curated 100 AraMed samples, enhanced them through grammatical correction and LLM modification, and added medical exam questions (MCQs and fill-in-the-blank), yielding 700 samples evaluating Arabic medical knowledge.

Religious Domain. Hajj-FAQ (Aleid and Azmi, 2025) addresses religious knowledge through question-answering on Hajj fatwas (Islamic legal rulings), providing a benchmark for LLMs’ understanding of Islamic jurisprudence and pilgrimage-related guidance. This specialized dataset evaluates models’ ability to handle religiously and culturally sensitive content requiring both linguistic competence and domain-specific religious knowledge.

5 NLP Tasks

Paper	Primary Task (Type)	Total
ARCD	Reading Comprehension (MCQ)	1.4k
TyDiQA (Ar)	Reading Comprehension (GEN)	14k
ORCA	Multi-task NLP (BENCH)	588k
Dolphin	Multi-task NLP (BENCH)	2M
Alghafa	Reading Comprehension (MCQ)	33.2k
LArABench	Multi-task NLP (BENCH)	296k
BALSAM	Multi-task NLP (BENCH)	52k

Table 2: Arabic NLP benchmarks covering core language understanding and generation tasks.

LLM evaluation on fundamental NLP tasks has evolved from narrow, task-specific datasets to comprehensive, multi-dimensional benchmarks. Early efforts such as ARCD (Mozannar et al., 2019) for reading comprehension and TyDiQA (Clark et al., 2020) for question answering established foundational paradigms for the pre-LLM era. The landscape has since shifted toward unified benchmarks that assess generalist capabilities across understanding, generation, and reasoning tasks, addressing morphological complexity, dialectal variation, and domain diversity Table 2.

5.1 Comprehensive Multi-Task Benchmarks

LArABench (Abdelali et al., 2024) introduced one of the first systematic benchmarking suites for Arabic NLP and speech, evaluating general-purpose LLMs against task-specific models across 33 tasks and 61 datasets. By integrating earlier resources such as ARCD and dialect identification corpora, it bridges pre-LLM and LLM evaluation. Results from zero- and few-shot evaluation of GPT-3.5-turbo, GPT-4 (OpenAI et al., 2024), BLOOMZ (Muennighoff et al., 2023), Jais-13b-chat (Sengupta et al., 2023), and Whisper (Radford et al., 2023) show that specialized models outperform LLMs in zero-shot settings, while larger LLMs substantially narrow the gap under few-shot prompting.

BALSAM (Al-Matham et al., 2025) provides a community-driven benchmark emphasizing instruction-following across 78 tasks and 14 categories (52K examples), spanning core NLP tasks such as summarization, QA, translation, and reasoning. It mitigates contamination via blind test sets, offers an integrated leaderboard, and demonstrates that LLM-as-judge aligns more closely with human judgments than traditional metrics. Drawing on mixed natural, translated, and synthetic data from sources including xP3 (Muennighoff et al., 2023), PromptSource (Bach et al., 2022), SuperNaturalInstructions (Wang et al., 2022), TruthfulQA (Lin et al., 2022), and newly curated datasets, BALSAM shows that large closed-source models outperform smaller Arabic-centric models, with performance influenced by tokenization, data scale, and Arabic-specific tuning. Identified limitations include residual cultural misalignment and limited evaluation of multi-turn dialogue and hallucination.

5.2 Task-Specific Benchmarks

Understanding and generation capabilities require evaluation beyond broad benchmark suites. ORCA (Elmadany et al., 2023) targets Arabic natural language understanding by consolidating 60 datasets into seven task categories, covering classification, sequence labeling, semantic similarity, inference, disambiguation, and question answering across MSA and dialects. It evaluates 18 multilingual and Arabic-specific models using task-appropriate metrics and provides a public leaderboard with detailed metadata, supporting reproducible comparison and tracking progress toward modern LLM evaluation.

Dolphin (Nagoudi et al., 2023) focuses on Arabic natural language generation across nine tasks, including generation, paraphrasing, summarization, transliteration, and grammatical error correction, spanning MSA and dialects. Aggregating 15 datasets, it evaluates encoder–decoder and decoder-only models using standard generation metrics and task-specific measures. Results highlight persistent gaps between general-purpose multilingual models and Arabic-finetuned baselines, though Dolphin predates instruction-tuned LLMs and relies on zero-shot evaluation, likely underestimating current performance.

Together with LAraBench and BALSAM, ORCA and Dolphin form complementary benchmarks covering understanding, generation, reasoning, and instruction-following across dialects and domains, providing a foundation for future Arabic LLM evaluation.

6 Culture and Dialects

Paper	Primary Focus (Type)	Total
Jawaher	Cultural Proverbs (GEN)	10k
PALM	Cultural Language Use (GEN)	17.4k
PalmX	Cultural Knowledge (MCQ)	6.4k
Commonsense	Cultural Commonsense (MCQ)	3.5k
AraDiCE	Dialect & Culture (BENCH)	81.8k
AL-QASIDA	Dialectal Analysis (BENCH)	–
Absher	Saudi Dialect & Culture (MCQ)	18k

Table 3: Benchmarks targeting Arabic culture, dialects, and region-specific knowledge.

Evaluation of Arabic cultural and dialectal understanding has evolved incrementally, with successive benchmarks addressing limitations of earlier efforts, as summarized in Table 3. Jawaher (Magdy et al., 2025) introduced one of the earliest culturally grounded resources by compiling 10,037 Arabic proverbs annotated with dialectal origin and

idiomatic meaning, enabling the evaluation of figurative reasoning in dialectal Arabic. However, its scope was limited to proverbial knowledge.

PALM (Alwajih et al., 2025a) expanded cultural coverage through a large-scale, community-driven effort spanning all 22 Arab countries, covering diverse genres and topics and evaluating cultural inclusivity via perplexity and generation quality. PalmX (Alwajih et al., 2025b) further targeted deep understanding of Arabic and Islamic culture through MCQ-based evaluation in MSA, covering traditions, history, and religious practices. Arab-specific commonsense reasoning, absent from prior benchmarks, was later addressed by Commonsense Reasoning in Arab Culture (Sadallah et al., 2025), which introduced culturally grounded inference tasks.

Dialectal evaluation has followed a parallel trajectory. AraDiCE (Mousi et al., 2025) assessed dialect identification, generation, and cognitive reasoning in Egyptian, Levantine, and Gulf Arabic, largely via dialectal adaptations of existing benchmarks. Alqasida et al. (Robinson et al., 2025) proposed a broader evaluation framework spanning identification, comprehension, generation quality, and dialect–MSA translation. NADI 2024 (Abdul-Mageed et al., 2024) was employed for dialect identification and translation-based chat evaluation, while Absher (Al-Monef et al., 2025) focused specifically on Saudi dialect vocabulary, phrases, and proverbs.

Overall, despite meaningful progress, benchmarks targeting Arabic cultural alignment and dialectal diversity remain limited in scope and coverage. Future efforts must broaden dialectal representation and deepen cultural reasoning evaluation to better reflect real-world Arabic language use.

7 Target-Specific Tasks

Paper	Primary Focus (Type)	Total
CamelEval	Instruction Following (GEN)	1.6k
Halwasa	Hallucination Detection (GEN)	10k
Henna	Multimodal Understanding (GEN)	1.1k
CAMEL-Bench	Multimodal (BENCH)	29k
AraTrust	Safety & Trust (MCQ)	0.52k
Arabic Safeguard	Safety Evaluation (GEN)	5.8k
ALRAGE	Retrieval-Augmented QA (GEN)	21.2k
AraTable	Tabular Reasoning (GEN)	0.6k
AraHalluEval	Hallucination Detection (GEN)	1.5k
HalluVerse25	Hallucination Classification (CL)	0.8k
ARB	Multimodal Reasoning (GEN)	1.36k
ASAS	Safety Benchmark (GEN)	0.8k

Table 4: Target-specific Arabic benchmarks for safety, multimodality, reasoning, and robustness.

Beyond general NLP capabilities, LLMs with Arabic capabilities require evaluation on specialized tasks reflecting real-world deployment scenarios and LLM-specific challenges, as summarized in Table 4.

Instruction-Following. CamelEval (Qian et al., 2024) evaluates conversational abilities and instruction-following through 1,610 generative questions: 805 human-validated translations from AlpacaEval (Dubois et al., 2023) and 805 synthetically generated using GPT-4 from culturally grounded textbooks. Evaluation uses LLM-as-judge computing win rates for open-ended generation.

Context-Based Reasoning. ALRAGE (El Filali et al., 2025) targets RAG evaluation through 2.12K question-answer-context trios from 40 Arabic books, synthetically generated via Meta-Llama-3.1-70B and validated by native speakers. AraTable (Alshaikh et al., 2025) addresses structured tabular data through 41 tables (15 QA pairs each) evaluating direct answering, fact verification, and complex reasoning. Tables sourced from Wikipedia, government portals, and GPT-4o were verified by human experts, with evaluation combining accuracy and the Assisted Self-Deliberation (ASD) framework employing judge LLMs and human evaluation.

Hallucination Detection. HalluVerse25 (Abdjalil et al., 2025) provides fine-grained multilingual evaluation across entity, relation, and sentence hallucination types, offering 828 Arabic samples from Wikidata autobiographies where GPT-4 injected false hallucinations validated by human judges. AraHalluEval (Alansari and Luqman, 2025) focuses on Arabic through QA (300 samples from TyDiQA-GoldP-AR and translated TruthfulQA) and summarization (100 XLSum instances), distinguishing factuality and faithfulness hallucinations through manual annotation. Halwasa (Mubarak et al., 2024) provides larger-scale evaluation with 10K samples where ChatGPT and GPT-4 generated factual sentences from 1,000 SAMER Arabic Lexicon words, with hallucination indicators annotated by 200 annotators enabling sentence-level analysis.

Safety and Trustworthiness. Ashraf et al. (2025) introduced 5,799 questions spanning direct attacks, indirect attacks, and harmless requests with sensitive words, employing dual-perspective evaluation from governmental and oppositional viewpoints. ASAS (aiastrolabe, 2025) provides the first human-rated Arabic safety index for red-teaming frontier models, exposing persistent weaknesses

in top-performing systems. AraTrust (Alghamdi et al., 2024) assesses trustworthiness across nine dimensions through 522 human-written MCQs.

Multimodal Capabilities. Peacock (Alwajih et al., 2024) pioneered Arabic multimodal evaluation with Henna benchmark combining standard VQA/OCR prompts with focus on culture understanding. CAMEL-Bench (Ghaboura et al., 2024) provides large-scale evaluation across eight domains and 38 subdomains, revealing promising performance but substantial weaknesses in specialized domains requiring precise vision-text alignment. ARB (Ghaboura et al., 2025) advances step-by-step multimodal reasoning, emphasizing logical integration of visual and textual inputs over simple captioning.

8 Discussion

Category	Coverage
Q&A / Reading Comprehension	Moderate (5 datasets)
Translation & Multitask Generation	Moderate (4 datasets)
Reasoning & Multi-step Thinking	Limited (3 datasets)
STEM / Academic Evaluation	Strong (8 datasets)
Law / Legal Reasoning	Limited (2 datasets)
Poetry / Literature / Arts	Limited (3 datasets)
Cultural Alignment & Dialect Evaluation	Strong (7 datasets)
Commonsense & Cultural Reasoning	Limited (3 datasets)
Hallucination / Truthfulness	Limited (3 datasets)
Retrieval-Augmented / Contextual Tasks	Limited (3 datasets)

Table 5: Summary of Arabic NLP dataset coverage and key gaps. Coverage indicates the number of datasets per task category, while Key Gaps outline remaining needs such as standardization, scale, and cultural balance.

This section discusses key observations regarding reproducibility practices, community initiatives, methodological inconsistencies, and critical gaps requiring future attention.

8.1 Reproducibility and Evaluation Benchmarks

Reproducibility is fundamental to scientific progress, yet our survey reveals significant heterogeneity in benchmark accessibility and transparency. We classify benchmarks into three categories based on their openness:

Private Benchmarks release neither datasets nor evaluation pipelines, severely limiting community impact and making verification impossible. Examples include AraSTEM (11,637 samples) and Halwasa (10K samples), which remain unavailable despite published results.

Partially-Public Benchmarks release datasets but withhold evaluation code, hindering exact reproduc-

490	tion due to ambiguities in preprocessing, prompt-	ROUGE, and BLEU.	541
491	ing, and metric computation.		
492	Public Benchmarks release both datasets and com-	8.3 Methodological Inconsistencies and	542
493	plete pipelines, either through dedicated reposi-	Quality Issues	543
494	tries (LAraBench, 3LM) or integration with main-	Through detailed examination of released bench-	544
495	stream frameworks like lighteval (Habib et al.,	marks, we identified several concerning inconsis-	545
496	2023) and lm-eval-harness (Gao et al., 2024).	tencies and quality issues that affect evaluation va-	546
497	Our analysis reveals that approximately 25% of	lidity and cross-benchmark comparability.	547
498	surveyed benchmarks remain private or partially	Multiple-Choice Formatting. Benchmarks em-	548
499	public, limiting their impact towards Arabic NLP	ploy inconsistent option labeling: some use Latin	549
500	community. We strongly advocate for complete	letters while others use Arabic letters. This incon-	550
501	transparency as the community standard, recogniz-	sistency affects model performance as LLMs may	551
502	ing that data contamination concerns can be ad-	be more familiar with Latin alphabetic indices from	552
503	ressed through alternative mechanisms such as	English pretraining. Standardization is needed for	553
504	blind test sets (as employed by BALSAM), peri-	fair cross-benchmark comparison.	554
505	odic dataset refreshment, or held-out evaluation	Prompting Variations. Few-shot evaluation	555
506	techniques rather than complete privatization.	ranges from 0-shot to 5-shot with limited system-	556
507		atic investigation of optimal settings per task type.	557
508	8.2 Common Metrics	Prompt phrasing varies from formal MSA to con-	558
509	Over the course of this survey, we identified the	versational styles, potentially affecting responses.	559
510	most commonly used metrics, which we list and	Quality Control Issues. Manual inspection re-	560
511	discuss below:	vealed typos, grammatical errors, and formatting	561
512	Accuracy This was a common approach in MCQ	inconsistencies in several released datasets. Some	562
513	benchmarks, with most inferring accuracy based	benchmarks contain culturally inappropriate con-	563
514	on the log-probabilities of choice indices. Other	tent despite claims of cultural alignment. These	564
515	benchmarks used joint log-probabilities over the	quality issues undermine benchmark validity and	565
516	actual choice text, which typically leads to lower	highlight the need for rigorous review processes.	566
517	scores. Accuracy based on log-probabilities has	LLM-as-Judge Validation. While BALSAM re-	567
518	been widely used since the rise of LLMs, but it	ports 0.824-0.977 correlation with human judg-	568
519	is less applicable to instruction-tuned models and	ments, most benchmarks adopting LLM-as-judge	569
520	cannot be used to evaluate closed-source models.	lack validation specifically in Arabic contexts.	570
521	LLM-as-Judge For generative tasks, several ex-	Judge model selection, prompt sensitivity, and po-	571
522	isting Arabic benchmarks rely on a closed-source	tential circular evaluation when judges resemble	572
523	model acting as an expert to evaluate LLM out-	evaluated models require systematic investigation.	573
524	puts, typically by prompting the model to score		
525	responses using a predefined prompt. Despite be-	8.4 Arabic LLM Leaderboards	574
526	ing adopted by nine benchmarks, this approach	Centralized leaderboards serve critical functions:	575
527	suffers from high API costs, expert bias, incons-	establishing performance baselines, enabling fair	576
528	istent standards, and limited reproducibility. To	model comparisons, tracking field progress over	577
529	mitigate some of these issues, a small number of	time, and guiding model selection for practition-	578
530	benchmarks incorporated human validation on a	ers. Several initiatives have emerged to fulfill these	579
531	limited subset to increase trust in the used expert.	roles for LLMs with Arabic support, each with	580
532	Human-as-Judge In at least two of the reported	distinct philosophies and design choices.	581
533	benchmarks, a Human-as-Judge approach was	The Open Arabic LLM Leaderboard (OALL)	582
534	adopted, in which human evaluators manually as-	(El Filali et al., 2024a) pioneered open-source	583
535	sess model responses. This approach represents	rankings, initially using Alghafa, EXAMS_ar, and	584
536	the most reliable method for obtaining qualitative	MMLU_ar. OALL v2 (El Filali et al., 2025) transi-	585
537	signals about LLM performance. However, it is	tioned to native benchmarks (ArabicMMLU, AL-	586
538	both intensive in terms of time and labor, making	RAGE, AraTrust, MadinahQA), reflecting commu-	587
539	it difficult to apply for large scale datasets.	nity consensus against translated content. BAL-	588
540	Apart from the aforementioned metrics, com-	SAM (Al-Matham et al., 2025) offers compre-	589
	monly known NLP metrics were used such F1, EM,	hensive evaluation across 78 tasks (52K samples)	590

with private test sets preventing contamination, including both closed and open-source models. IL-MAAM (Nacar et al., 2025) specializes in culturally aligned evaluation using refined ArabicMMLU by ensuring religious sensitivity and social norms. The AraGen benchmark (El Filali et al., 2024b) and its associated leaderboard adopt a 3C3H evaluation metric (Correctness, Completeness, Conciseness, Helpfulness, Honesty, Harmlessness), using LLM-as-judge, and combine it with a dynamic, blind-testing approaches to push for robust and fair benchmarking of Arabic capabilities in LLMs.

8.5 Critical Gaps

Despite substantial progress, several critical areas remain underexplored or entirely absent from current Arabic benchmarking efforts. Table 5 summarizes the coverage of datasets per category as well as key gaps.

Underexplored Areas. Current benchmarks lack temporal evaluation, multi-turn dialogue assessment, code-switching evaluation (Arabic-English/French mixing), low-resource dialect coverage (Sudanese, Mauritanian), pragmatic understanding (sarcasm, indirect speech), and specialized domains beyond law/medicine (education, journalism, technical documentation).

Methodological Challenges. Data contamination threatens validity as training corpora expand. Static benchmarks fail to capture temporal degradation. Cultural alignment lacks standardized metrics. Heavy reliance on LLM-as-judge requires more rigorous Arabic-specific validation.

Dataset Size Matters. As Figure 2 shows, benchmark sizes range from under 500 samples (Qiyas, EXAMS_ar) to over 2 million (Dolphin). This dramatic variability is consequential: small benchmarks provide weaker statistical signals and are more susceptible to overfitting, while aggregate scores averaging across diverse benchmark sizes can be misleading if not properly weighted. Evaluation interpretation must account for dataset scale.

8.6 Recommendations

To address the identified gaps, future Arabic benchmarks should:

Prioritize native Arabic data with mandatory cultural and dialectal review to reduce translation bias.

Ensure full reproducibility by releasing evaluation code, prompts, and scoring procedures within standardized frameworks.

Control contamination using blind test sets and automated overlap detection.

Standardize evaluation protocols by fixing prompts, answer formats, and decoding settings to enable fair comparison.

Validate automated judges through calibration against diverse Arabic human annotations.

Expand dialectal coverage by defining target dialect sets and minimum sample sizes.

Strengthen dataset quality control via multi-pass annotation and transparent documentation.

Incorporate temporal evaluation through versioned or rolling benchmarks to measure performance drift.

These steps are essential for building Arabic benchmarks that are culturally grounded, reproducible, and robust over time.

9 Conclusion

This survey comprehensively reviews 40+ Arabic benchmarks, providing the first systematic taxonomy across NLP tasks, knowledge domains, cultural understanding, and specialized capabilities. The field has evolved from task-specific datasets (ARCD, TyDiQA) to comprehensive benchmarks (BALSAM, LAraBench, ArabicMMLU) addressing understanding, generation, and reasoning.

Progress includes increased dialectal coverage, cultural grounding, and domain diversity. However, persistent challenges remain: data contamination, cultural misalignment in translations, insufficient coverage of temporal reasoning and multi-turn dialogue, and inconsistent reproducibility. Our recommendations emphasize reproducibility, cultural authenticity, dialectal inclusivity, and methodological rigor. As Arabic capabilities in LLMs advance, evaluation must authentically capture Arabic’s linguistic complexity and cultural richness. This survey serves as a comprehensive reference, guiding development of more robust, culturally aligned LLM evaluation on Arabic capabilities.

Limitations

Despite our comprehensive review, this survey has limitations. Rapidly evolving Arabic LLM benchmarking means recent or proprietary benchmarks may be excluded, particularly those without public documentation. Our review covers benchmarks through early 2025, so some information may quickly become outdated. Several benchmarks

(e.g., AraSTEM, Halwasa, Qiyas) remain unavailable, limiting analysis to published descriptions.

We focus on evaluation benchmarks rather than training datasets, emphasizing text-based and multimodal tasks while excluding speech-only benchmarks. Assessments of benchmark quality and cultural alignment reflect our perspectives; native speakers from other regions may judge differently. Finally, our recommendations are informed opinions, not community consensus, and classifications of benchmarks as public or private reflect the information available at the time of writing.

References

Samir Abdaljalil, Hasan Kurban, and Erchin Serpedin. 2025. [Halluverse25: Fine-grained multilingual benchmark dataset for llm hallucinations](#). *Preprint*, arXiv:2503.07833.

Ahmed Abdelali, Hamdy Mubarak, Shammur Chowdhury, Maram Hasanain, Basel Mousi, Sabri Boughorbel, Samir Abdaljalil, Yassine El Kheir, Daniel Izham, Fahim Dalvi, Majd Hawasly, Nizi Nazar, Youssef Elshahawy, Ahmed Ali, Nadir Durrani, Natasa Milic-Frayling, and Firoj Alam. 2024. [LaraBench: Benchmarking Arabic AI with large language models](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 487–520, St. Julian’s, Malta. Association for Computational Linguistics.

Muhammad Abdul-Mageed, Amr Keleg, AbdelRahim Elmadany, Chiyu Zhang, Injy Hamed, Walid Magdy, Houda Bouamor, and Nizar Habash. 2024. Nadi 2024: The fifth nuanced arabic dialect identification shared task. *arXiv preprint arXiv:2407.04910*.

aiastrolabe. 2025. Redteaming frontier llms with aiastrolabe arabic safety index (asas). Accessed: Oct. 3, 2025.

Shahad Al-Khalifa and Hend Al-Khalifa. 2024. The qiyas benchmark: Measuring chatgpt mathematical and language understanding in arabic. *arXiv preprint arXiv:2407.00146*.

Shahad Al-Khalifa, Nadir Durrani, Hend Al-Khalifa, and Firoj Alam. 2025. The landscape of arabic large language models. *Communications of the ACM*.

Nashwan Ahmed Al-Majmar, Hezam Gawbah, and Akram Alsubari. 2024. Ahd: Arabic healthcare dataset. *Data in Brief*, 56:110855.

Rawan Al-Matham, Kareem Darwish, Raghad Al-Rasheed, Waad Alshammari, Muneera Alhoshan, Amal Almazrua, Asma Al Wazrah, Mais Alheraki, Firoj Alam, Preslav Nakov, Norah Alzahrani, Eman alBilali, Nizar Habash, Abdelrahman El-Sheikh, Muhammad Elmallah, Haonan Li, Hamdy Mubarak,

Mohamed Anwar, Zaid Alyafeai, and 24 others. 2025. [Balsam: A platform for benchmarking arabic large language models](#). *Preprint*, arXiv:2507.22603.

Renad Al-Monef, Hassan Alhuzali, Nora Alturayef, and Ashwag Alasmari. 2025. [Absher: A benchmark for evaluating large language models’ understanding of saudi dialects](#). *Preprint*, arXiv:2507.10216.

Aisha Alansari and Hamzah Luqman. 2025. [Arahallueval: A fine-grained hallucination evaluation framework for arabic llms](#). *Preprint*, arXiv:2509.04656.

Ashwag Alasmari, Sarah Alhumoud, and Waad Alshammari. 2024. Aramed: Arabic medical question answering using pretrained transformer language models. In *Proceedings of the 6th Workshop on Open-Source Arabic Corpora and Processing Tools (OS-ACT) with Shared Tasks on Arabic LLMs Hallucination and Dialect to MSA Machine Translation@ LREC-COLING 2024*, pages 50–56.

Safa AlBallaa, Nora AlTwaresh, Abdulmalik AlSalman, and Sultan Alfarhood. 2025. Gatmath and gatlc: Comprehensive benchmarks for evaluating arabic large language models. *PLoS One*, 20(9):e0329129.

Hayfa A Aleid and Aqil M Azmi. 2025. Hajj-fqa: A benchmark arabic dataset for developing question-answering systems on hajj fatwas: H. aleid and a. azmi. *Journal of King Saud University Computer and Information Sciences*, 37(6):135.

Wafa Alghallabi, Ritesh Thawkar, Sara Ghaboura, Ketan More, Omkar Thawakar, Hisham Cholakkal, Salman Khan, and Rao Muhammad Anwer. 2025. Fann or flop: A multigenre, multiera benchmark for arabic poetry understanding in llms. *arXiv preprint arXiv:2505.18152*.

Emad A. Alghamdi, Reem I. Masoud, Deema Alnuhait, Afnan Y. Alomairi, Ahmed Ashraf, and Mohamed Zaytoon. 2024. [Aratrust: An evaluation of trustworthiness for llms in arabic](#). *Preprint*, arXiv:2403.09017.

Mohamed Alkaoud. 2024. A bilingual benchmark for evaluating large language models. *PeerJ Computer Science*, 10:e1893.

Ebtesam Almazrouei, Ruxandra Cojocaru, Michele Baldo, Quentin Malartic, Hamza Alobeidli, Daniele Mazzotta, Guilherme Penedo, Giulia Campesan, Murgariya Farooq, Maittha Alhammedi, Julien Launay, and Badreddine Noune. 2023. [AlGhafa evaluation benchmark for Arabic language models](#). In *Proceedings of ArabicNLP 2023*, pages 244–275, Singapore (Hybrid). Association for Computational Linguistics.

Rana Alshaikh, Israa Alghanmi, and Shelan Jeawak. 2025. [Aratable: Benchmarking llms’ reasoning and understanding of arabic tabular data](#). *Preprint*, arXiv:2507.18442.

Fakhraddin Alwajih, Abdellah El Mekki, Samar Mohamed Magdy, AbdelRahim A. Elmadany, Omer

795	Nacar, El Moatez Billah Nagoudi, Reem Abdel-Salam, Hanin Atwany, Youssef Nafea, Abdulfattah Mohammed Yahya, Rahaf Alhamouri, Hamzah A. Alsayadi, Hiba Zayed, Sara Shatnawi, Serry Sibae, Yasir Ech-chammakhy, Walid Al-Dhabyani, Marwa Mohamed Ali, Imen Jarraya, and 25 others. 2025a. Palm: A culturally inclusive and linguistically diverse dataset for Arabic LLMs . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 32871–32894, Vienna, Austria. Association for Computational Linguistics.	851
796		852
797		853
798		854
799		855
800		856
801		857
802		858
803		
804		
805		
806		
807	Fakhraddin Alwajih, Abdellah El Mekki, Hamdy Mubarak, Majd Hawasly, Abubakr Mohamed, and Muhammad Abdul-Mageed. 2025b. Palmx 2025: The first shared task on benchmarking llms on arabic and islamic culture . <i>Preprint</i> , arXiv:2509.02550.	
808		
809		
810		
811		
812	Fakhraddin Alwajih, El Moatez Billah Nagoudi, Gagan Bhatia, Abdelrahman Mohamed, and Muhammad Abdul-Mageed. 2024. Peacock: A family of arabic multimodal large language models and benchmarks . <i>Preprint</i> , arXiv:2403.01031.	
813		
814		
815		
816		
817	Anthropic. 2024. Introducing the next generation of claude .	
818		
819	Wissam Antoun, Fady Baly, and Hazem Hajj. 2021. AraGPT2: Pre-trained transformer for Arabic language generation . In <i>Proceedings of the Sixth Arabic Natural Language Processing Workshop</i> , pages 196–207, Kyiv, Ukraine (Virtual). Association for Computational Linguistics.	
820		
821		
822		
823		
824		
825	Yasser Ashraf, Yuxia Wang, Bin Gu, Preslav Nakov, and Timothy Baldwin. 2025. Arabic dataset for llm safeguard evaluation . <i>Preprint</i> , arXiv:2410.17040.	
826		
827		
828	Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. 2021. Program synthesis with large language models . <i>Preprint</i> , arXiv:2108.07732.	
829		
830		
831		
832		
833	Stephen H. Bach, Victor Sanh, Zheng-Xin Yong, Albert Webson, Colin Raffel, Nihal V. Nayak, Abheesht Sharma, Taewoon Kim, M Saiful Bari, Thibault Fevry, Zaid Alyafeai, Manan Dey, Andrea Santilli, Zhiqing Sun, Srulik Ben-David, Canwen Xu, Gunjan Chhablani, Han Wang, Jason Alan Fries, and 8 others. 2022. Promptsources: An integrated development environment and repository for natural language prompts . <i>Preprint</i> , arXiv:2202.01279.	
834		
835		
836		
837		
838		
839		
840		
841		
842	Adil Bahaj and Mounir Ghogho. 2025. Mizanqa: Benchmarking large language models on moroccan legal question answering . <i>arXiv preprint arXiv:2508.16357</i> .	
843		
844		
845		
846	Basma El Amel Boussaha, Leen AlQadi, Mugariya Farooq, Shaikha Alsuwaidi, Giulia Campesan, Ahmed Alzubaidi, Mohammed Alyafeai, and Hakim Hacid. 2025. 3lm: Bridging arabic, stem, and code through benchmarking . <i>arXiv preprint arXiv:2507.15850</i> .	
847		
848		
849		
850		
	Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, and 39 others. 2021. Evaluating large language models trained on code . <i>Preprint</i> , arXiv:2107.03374.	851
		852
		853
		854
		855
		856
		857
		858
	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>Preprint</i> , arXiv:2003.05002.	859
		860
		861
		862
		863
		864
	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems . <i>arXiv preprint arXiv:2110.14168</i> .	865
		866
		867
		868
		869
		870
	Mouath Abu Daoud, Chaimae Abouzahir, Leen Kharouf, Walid Al-Eisawi, Nizar Habash, and Farah E Shamout. 2025. Medarabiq: Benchmarking large language models on arabic medical tasks . <i>arXiv preprint arXiv:2505.03427</i> .	871
		872
		873
		874
		875
	Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Alpaca-farm: A simulation framework for methods that learn from human feedback . <i>Preprint</i> , arXiv:2305.14387.	876
		877
		878
		879
		880
	Ali El Filali, Hamza Alobeidli, Clémentine Fourier, Basma El Amel Boussaha, Ruxandra Cojocaru, Nathan Habib, and Hakim Hacid. 2024a. Open arabic llm leaderboard . https://huggingface.co/spaces/OALL/Open-Arabic-LLM-Leaderboard .	881
		882
		883
		884
		885
	Ali El Filali, Manel ALOUI, Tarique Husaain, Ahmed Alzubaidi, Basma El Amel Boussaha, Ruxandra Cojocaru, Clémentine Fourier, Nathan Habib, and Hakim Hacid. 2025. The open arabic llm leaderboard 2 . https://huggingface.co/spaces/OALL/Open-Arabic-LLM-Leaderboard .	886
		887
		888
		889
		890
		891
	Ali El Filali, Neha Sengupta, Abouelseoud, Preslav Nakov, and Clémentine Fourier. 2024b. Rethinking llm evaluation with 3c3h: Aragen benchmark and leaderboard . https://huggingface.co/blog/leaderboard-3c3h-aragen .	892
		893
		894
		895
		896
	AbdelRahim Elmadany, El Moatez Billah Nagoudi, and Muhammad Abdul-Mageed. 2023. Orca: A challenging benchmark for arabic language understanding . <i>Preprint</i> , arXiv:2212.10758.	897
		898
		899
		900
	Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, and 5 others. 2024. The language model evaluation harness .	901
		902
		903
		904
		905
		906
		907
		908

909	Sara Ghaboura, Ahmed Heakl, Omkar Thawakar, Ali Alharthi, Ines Riahi, Abduljalil Saif, Jorma Laaksonen, Fahad S. Khan, Salman Khan, and Rao M. Anwer. 2024. Camel-bench: A comprehensive arabic lmm benchmark . <i>Preprint</i> , arXiv:2410.18976.	966
910		967
911		968
912		
913		
914	Sara Ghaboura, Ketan More, Wafa Alghallabi, Omkar Thawakar, Jorma Laaksonen, Hisham Cholakkal, Salman Khan, and Rao Muhammad Anwer. 2025. Arb: A comprehensive arabic multimodal reasoning benchmark . <i>Preprint</i> , arXiv:2505.17021.	
915		
916		
917		
918		
919	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. The llama 3 herd of models . <i>Preprint</i> , arXiv:2407.21783.	
920		
921		
922		
923		
924		
925		
926		
927	Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher Ré, Adam Chilton, Aditya Narayana, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel N. Rockmore, Diego Zambrano, Dmitry Talisman, Enam Hoque, Faiz Surani, Frank Fagan, Galit Sarfaty, Gregory M. Dickinson, Haggai Porat, Jason Hegland, and 21 others. 2023. Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language models . <i>Preprint</i> , arXiv:2308.11462.	
928		
929		
930		
931		
932		
933		
934		
935		
936		
937	Nathan Habib, Clémentine Fourrier, Hynek Kydlíček, Thomas Wolf, and Lewis Tunstall. 2023. Lighteval: A lightweight framework for llm evaluation .	
938		
939		
940	Momchil Hardalov, Todor Mihaylov, Dimitrina Zlatkova, Yoan Dinkov, Ivan Koychev, and Preslav Nakov. 2020. Exams: A multi-subject high school examinations dataset for cross-lingual and multilingual question answering . <i>arXiv preprint arXiv:2011.03080</i> .	
941		
942		
943		
944		
945		
946	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding . <i>arXiv preprint arXiv:2009.03300</i> .	
947		
948		
949		
950	Faris Hijazi, Somayah AlHarbi, Abdulaziz AlHussein, Harethah Abu Shairah, Reem AlZahrani, Hebah Al-Shamlan, Omar Knio, and George Turkiyyah. 2024. Arablegaleval: A multitask benchmark for assessing arabic legal knowledge in large language models . <i>arXiv preprint arXiv:2408.07983</i> .	
951		
952		
953		
954		
955		
956	Huang Huang, Fei Yu, Jianqing Zhu, Xuening Sun, Hao Cheng, Dingjie Song, Zhihong Chen, Abdulmohsen Alharthi, Bang An, Juncai He, Ziche Liu, Zhiyi Zhang, Junying Chen, Jianquan Li, Benyou Wang, Lian Zhang, Ruoyu Sun, Xiang Wan, Haizhou Li, and Jinchao Xu. 2024. Acegpt, localizing large language models in arabic . <i>Preprint</i> , arXiv:2309.12053.	
957		
958		
959		
960		
961		
962		
963	Kaiyu Huang, Fengran Mo, Xinyu Zhang, Hongliang Li, You Li, Yuanchi Zhang, Weijian Yi, Yulong Mao, Jincheng Liu, Yuzhuang Xu, Jinan Xu, Jian-Yun Nie, and Yang Liu. 2025. A survey on large language models with multilingualism: Recent advances and new frontiers . <i>Preprint</i> , arXiv:2405.10936.	966
964		967
965		968
	Amr Keleg, Sharon Goldwater, and Walid Magdy. 2025. Revisiting common assumptions about arabic dialects in nlp . <i>arXiv preprint arXiv:2505.21816</i> .	969
		970
		971
	Fajri Koto, Haonan Li, Sara Shatnawi, Jad Doughman, Abdelrahman Sadallah, Aisha Alraeesi, Khalid Al-mubarak, Zaid Alyafeai, Neha Sengupta, Shady Shehata, Nizar Habash, Preslav Nakov, and Timothy Baldwin. 2024. ArabicMMLU: Assessing massive multitask language understanding in Arabic . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 5622–5640, Bangkok, Thailand. Association for Computational Linguistics.	972
		973
		974
		975
		976
		977
		978
		979
		980
	Anis Koubaa, Adel Ammar, Lahouari Ghouti, Omar Najar, and Serry Sibae. 2024. Arabiangpt: Native arabic gpt-based large language model . <i>Preprint</i> , arXiv:2402.15313.	981
		982
		983
		984
	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods . <i>Preprint</i> , arXiv:2109.07958.	985
		986
		987
	Samar Mohamed Magdy, Sang Yun Kwon, Fakhraddin Alwajih, Safaa Taher Abdelfadil, Shady Shehata, and Muhammad Abdul-Mageed. 2025. JAWAHER: A multidialectal dataset of Arabic proverbs for LLM benchmarking . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 12320–12341, Albuquerque, New Mexico. Association for Computational Linguistics.	988
		989
		990
		991
		992
		993
		994
		995
		996
		997
	Malak Mashaabi, Shahad Al-Khalifa, and Hend Al-Khalifa. 2024. A survey of large language models for arabic language and its dialects . <i>arXiv preprint arXiv:2410.20238</i> .	998
		999
		1000
		1001
	Basel Mousi, Nadir Durrani, Fatema Ahmad, Md. Arid Hasan, Maram Hasanain, Tameem Kabbani, Fahim Dalvi, Shammur Absar Chowdhury, and Firoj Alam. 2025. AraDiCE: Benchmarks for dialectal and cultural capabilities in LLMs . In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 4186–4218, Abu Dhabi, UAE. Association for Computational Linguistics.	1002
		1003
		1004
		1005
		1006
		1007
		1008
		1009
	Hussein Mozannar, Karl El Hajal, Elie Maamary, and Hazem M. Hajj. 2019. Neural arabic question answering . <i>CoRR</i> , abs/1906.05394.	1010
		1011
		1012
	Hamdy Mubarak, Hend Al-Khalifa, and Khaloud Suliman Alkhalefah. 2024. Halwasa: Quantify and analyze hallucinations in large language models: Arabic as a case study . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 8008–8015, Torino, Italia. ELRA and ICCL.	1013
		1014
		1015
		1016
		1017
		1018
		1019
		1020

1021	Niklas Muennighoff, Thomas Wang, Lintang Sutawika,	Zhaozhi Qian, Faroq Altam, Muhammad Alqurishi, and	1077
1022	Adam Roberts, Stella Biderman, Teven Le Scao,	Riad Souissi. 2024. Camelevel: Advancing cultur-	1078
1023	M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hai-	ally aligned arabic language models and benchmarks.	1079
1024	ley Schoelkopf, Xiangru Tang, Dragomir Radev,	<i>Preprint</i> , arXiv:2409.12623.	1080
1025	Alham Fikri Aji, Khalid Almubarak, Samuel Al-		
1026	banie, Zaid Alyafeai, Albert Webson, Edward Raff,	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brock-	1081
1027	and Colin Raffel. 2023. Crosslingual generaliza-	man, Christine McLeavey, and Ilya Sutskever. 2023.	1082
1028	tion through multitask finetuning . In <i>Proceedings</i>	Robust speech recognition via large-scale weak super-	1083
1029	<i>of the 61st Annual Meeting of the Association for</i>	vision. In <i>Proceedings of the 40th International Con-</i>	1084
1030	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	<i>ference on Machine Learning</i> , ICML'23. JMLR.org.	1085
1031	pages 15991–16111, Toronto, Canada. Association		
1032	for Computational Linguistics.	Haneh Rhel and Dmitri Roussinov. 2025. Large lan-	1086
		guage models and arabic content: a review. <i>arXiv</i>	1087
1033	Ahmad Mustapha, Hadi Al-Khansa, Hadi Al-Mubasher,	<i>preprint arXiv:2505.08004</i> .	1088
1034	Aya Mourad, Ranam Hamoud, Hasan El-Husseini,		
1035	Marwah Al-Sakkaf, and Mariette Awad. 2024.	Nathaniel R. Robinson, Shahd Abdelmoneim, Kelly	1089
1036	Arastem: A native arabic multiple choice question	Marchisio, and Sebastian Ruder. 2025. Al-qasida:	1090
1037	benchmark for evaluating llms knowledge in stem	Analyzing llm quality and accuracy systematically in	1091
1038	subjects. <i>arXiv preprint arXiv:2501.00559</i> .	dialectal arabic . In <i>Findings of the Association for</i>	1092
		<i>Computational Linguistics: ACL 2025</i> .	1093
1039	Omer Nacar, Serry Taiseer Sibae, Samar Ahmed, Safa	Abdelrahman Sadallah, Junior Cedric Tonga, Khalid	1094
1040	Ben Atitallah, Adel Ammar, Yasser Alhabashi, Ab-	Almubarak, Saeed Almheiri, Farah Atif, Chatrine	1095
1041	dulrahman S. Al-Batati, Arwa Alsehibani, Nour Qan-	Qwaider, Karima Kadaoui, Sara Shatnawi, Yaser	1096
1042	dos, Omar Elshehy, Mohamed Abdelkader, and Anis	Alesh, and Fajri Koto. 2025. Commonsense reason-	1097
1043	Koubaa. 2025. Towards inclusive Arabic LLMs: A	ing in Arab culture . In <i>Proceedings of the 63rd</i>	1098
1044	culturally aligned benchmark in Arabic large lan-	<i>Annual Meeting of the Association for Computational</i>	1099
1045	guage model evaluation . In <i>Proceedings of the First</i>	<i>Linguistics (Volume 1: Long Papers)</i> , pages 7695–	1100
1046	<i>Workshop on Language Models for Low-Resource</i>	7710, Vienna, Austria. Association for Computa-	1101
1047	<i>Languages</i> , pages 387–401, Abu Dhabi, United Arab	tional Linguistics.	1102
1048	Emirates. Association for Computational Linguistics.		
1049	El Moatez Billah Nagoudi, AbdelRahim Elmadany,	Neha Sengupta, Sunil Kumar Sahu, Bokang Jia,	1103
1050	Ahmed Oumar El-Shangiti, and Muhammad Abdul-	Satheesh Katipomu, Haonan Li, Fajri Koto, William	1104
1051	Mageed. 2023. Dolphin: A challenging and diverse	Marshall, Gurpreet Gosal, Cynthia Liu, Zhim-	1105
1052	benchmark for Arabic NLG . In <i>Findings of the As-</i>	ing Chen, Osama Mohammed Afzal, Samta Kam-	1106
1053	<i>sociation for Computational Linguistics: EMNLP</i>	boj, Onkar Pandit, Rahul Pal, Lalit Pradhan,	1107
1054	<i>2023</i> , pages 1404–1422, Singapore. Association for	Zain Muhammad Mujahid, Massa Baali, Xudong	1108
1055	Computational Linguistics.	Han, Sondos Mahmoud Bsharat, and 13 others. 2023.	1109
		Jais and jais-chat: Arabic-centric foundation and	1110
1056	Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad	instruction-tuned open generative large language	1111
1057	Saqib, Saeed Anwar, Muhammad Usman, Naveed	models . <i>Preprint</i> , arXiv:2308.16149.	1112
1058	Akhtar, Nick Barnes, and Ajmal Mian. 2025. A com-	Sumuk Shashidhar, Clémentine Fourier, Alina Lo-	1113
1059	prehensive overview of large language models. <i>ACM</i>	zovskia, Thomas Wolf, Gokhan Tur, and Dilek	1114
1060	<i>Transactions on Intelligent Systems and Technology</i> ,	Hakkani-Tür. 2025. Yourbench: Easy custom evalua-	1115
1061	16(5):1–72.	tion sets for everyone . <i>Preprint</i> , arXiv:2504.01833.	1116
1062	Omartificial-Intelligence-Space. 2025. Ara-	Fanar Team, Ummar Abbas, Mohammad Shahmeer Ah-	1117
1063	abic gsm8k: Arabic grade school math	mad, Firoj Alam, Enes Altinisik, Ehsannedin Asgari,	1118
1064	dataset. https://huggingface.co/datasets/	Yazan Boshmaf, Sabri Boughorbel, Sanjay Chawla,	1119
1065	Omartificial-Intelligence-Space/	Shammur Chowdhury, Fahim Dalvi, Kareem Dar-	1120
1066	Arabic-gsm8k-v2 .	wish, Nadir Durrani, Mohamed Elfeky, Ahmed El-	1121
1067	OpenAI. 2023. Chatgpt. https://chat.openai.com/	magarmid, Mohamed Eltabakh, Masoomali Fatehkia,	1122
1068	chat . Accessed: 2 October 2025.	Anastasios Fragkopoulos, Maram Hasanain, and 23	1123
		others. 2025a. Fanar: An arabic-centric multimodal	1124
1069	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,	generative ai platform . <i>Preprint</i> , arXiv:2501.13944.	1125
1070	Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-	Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya	1126
1071	man, Diogo Almeida, Janko Altmenschmidt, Sam Alt-	Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin,	1127
1072	man, Shyamal Anadkat, Red Avila, Igor Babuschkin,	Tatiana Matejovicova, Alexandre Ramé, Morgane	1128
1073	Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim-	Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey	1129
1074	ing Bao, Mohammad Bavarian, Jeff Belgum, and	Cideron, Jean bastien Grill, Sabela Ramos, Edouard	1130
1075	262 others. 2024. Gpt-4 technical report . <i>Preprint</i> ,	Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev,	1131
1076	arXiv:2303.08774.	and 197 others. 2025b. Gemma 3 technical report .	1132
		<i>Preprint</i> , arXiv:2503.19786.	1133

1134	Silma Team. 2024. <i>Silma</i> .	A Summary Table	1154
1135	TII. 2025. <i>Falcon-arabic: A breakthrough in arabic language models</i> .	Table 6 provides a comprehensive overview of all the Arabic benchmarks reviewed in this survey. This appendix clarifies abbreviations and highlights key insights for interpreting the data.	1155
1136			1156
1137	Yizhong Wang, Swaroop Mishra, Pegah Alipoor-molabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Gary Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krma Doshi, Maitreya Patel, and 21 others. 2022. <i>Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks</i> . <i>Preprint</i> , arXiv:2204.07705.	A.1 Column Definitions	1157
1138			1158
1139		Sample Composition Columns:	1160
1140		• NAT: Natively authored Arabic samples (originally created in Arabic by native speakers)	1161
1141			1162
1142		• SYN: Synthetically generated samples (created using LLMs like GPT-4, Claude, or Llama)	1163
1143			1164
1144		• TRAN: Translated samples (translated from English or other languages to Arabic)	1165
1145			1166
1146		• Total: Total benchmark size	1167
1147	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. <i>Qwen3 technical report</i> . <i>Preprint</i> , arXiv:2505.09388.	Evaluation Details:	1168
1148			1169
1149		• Metrics: Evaluation metrics employed (e.g., Accuracy, F1, BLEU, ROUGE, LLM-as-Judge)	1170
1150			1171
1151		• Type: Question format - MCQ (Multiple-choice), GEN (Generative/open-ended), BENCH (Multi-task suite), CL (Classification), RS (Reasoning Steps)	1172
1152			1173
1153			1174
		Accessibility:	1175
		• PD (Public Dataset): Yes (publicly available), No (private/request-only), Partial (subset public)	1176
			1177
		• PR (Public Repository): Yes (complete evaluation code available), No (no public repository)	1178
			1179
			1180
		Notation: "-" (not available/applicable), "k" (thousands), "m" (millions), "DS" (datasets)	1181
			1182
		A.2 Key Insights from the Table	1183
		Dataset Composition Trends:	1184
		• Native vs. Translated: Early benchmarks (2020) relied heavily on translation (MMLU_ar: 100% translated), while recent benchmarks (2024-2025) prioritize native content (ArabicMMLU, GATmath: 100% native), reflecting community consensus on cultural authenticity.	1185
			1186
			1187
			1188
			1189
			1190
			1191
			1192
			1193
			1194

1195	• Synthetic Generation: Recent benchmarks increasingly employ synthetic generation (ALRAGE: 21.2k synthetic, ArabLegalEval: 10.58k synthetic), balancing scale with cost. However, purely native benchmarks like Jawaher (10k native) and Absher (18k native) demonstrate that large-scale native collection remains feasible.	1240
1196		1241
1197		1242
1198		1243
1199		
1200		
1201		
1202		
1203	• Hybrid Approaches: BALSAM exemplifies effective hybrid strategy (26k native, 1.7k synthetic, 24k translated), combining strengths of multiple approaches.	
1204		
1205		
1206		
1207	Scale Distribution:	
1208	• Large-Scale: Dolphin (2m samples) and LAraBench (296k) represent comprehensive multi-task benchmarks.	
1209		
1210		
1211	• Medium-Scale: Most benchmarks range from 1k-20k samples, balancing quality and coverage.	
1212		
1213		
1214	• Specialized: Domain-specific benchmarks (legal, medical) tend toward smaller, higher-quality datasets (1-7k samples).	
1215		
1216		
1217	Reproducibility Concerns:	
1218	• Fully Accessible: Only 21 of 41 benchmarks provide both public datasets and repositories.	
1219		
1220	• Dataset-Only: 8 benchmarks release datasets without evaluation code, hindering exact reproduction.	
1221		
1222		
1223	• Private: 9 benchmarks remain completely private such as AraST EM, Halwasa, CamelEval and Qiyas, including some with substantial contributions (AraST EM: 11.6k samples), severely limiting community impact.	
1224		
1225		
1226		
1227		
1228	• Critical Gap: The high proportion of inaccessible benchmarks impedes scientific progress and independent validation.	
1229		
1230		
1231	Evaluation Methodology Evolution:	
1232	• Traditional Metrics: Early benchmarks use standard metrics (Accuracy, F1, BLEU, ROUGE).	
1233		
1234		
1235	• LLM-as-Judge: Recent benchmarks increasingly adopt LLM-as-judge (BALSAM, ALRAGE, ArabLegalEval), particularly for generative tasks, though validation against human judgments remains limited.	
1236		
1237		
1238		
1239		
	• Human Evaluation: Hallucination benchmarks (AraHalluEval, Halwasa) maintain human-as-judge given the critical nature of the task, despite higher costs.	1244
	Task Coverage Patterns:	1244
	• NLP Tasks: Well-covered with multiple comprehensive benchmarks (ORCA, Dolphin, LAraBench, BALSAM).	1245
		1246
		1247
	• STEM: Substantial progress with native benchmarks (ArabicMMLU, AraST EM, GATmath) addressing earlier translation limitations.	1248
		1249
		1250
		1251
	• Culture & Dialects: Growing emphasis (7 benchmarks in 2025 alone), reflecting recognition of cultural alignment importance.	1252
		1253
		1254
	• Target-Specific: Emerging area with recent focus on hallucination detection and context-based reasoning.	1255
		1256
		1257
	Temporal Distribution: Most benchmarks were released in 2024-2025, indicating rapid field acceleration. However, this concentration also suggests potential redundancy in some areas (multiple GAT-based benchmarks) while other areas remain unexplored (temporal evaluation, code-switching, multi-turn dialogue).	1258
		1259
		1260
		1261
		1262
		1263
		1264
	Geographic and Cultural Bias: Several benchmarks focus on specific regional variants (Absher: Saudi dialect, MizanQA: Moroccan law, GATmath: Saudi standardized tests), highlighting both progress in regional representation and gaps in coverage of other Arabic-speaking regions (North Africa, Levant, Iraq).	1265
		1266
		1267
		1268
		1269
		1270
		1271
	B Arabic NLP Benchmarks Timeline	1272
	Figure 3 presents a timeline of Arabic benchmark releases from 2019 to 2025, categorized by our four-category taxonomy. The visualization reveals dramatic acceleration in benchmark development, with 82% of all benchmarks (34 of 41) released in 2024-2025 alone. This recent surge reflects growing recognition of Arabic LLM evaluation needs and demonstrates rapid field maturation across all categories, particularly in cultural/dialectal assessment where all 7 benchmarks emerged in 2025.	1273
		1274
		1275
		1276
		1277
		1278
		1279
		1280
		1281
		1282

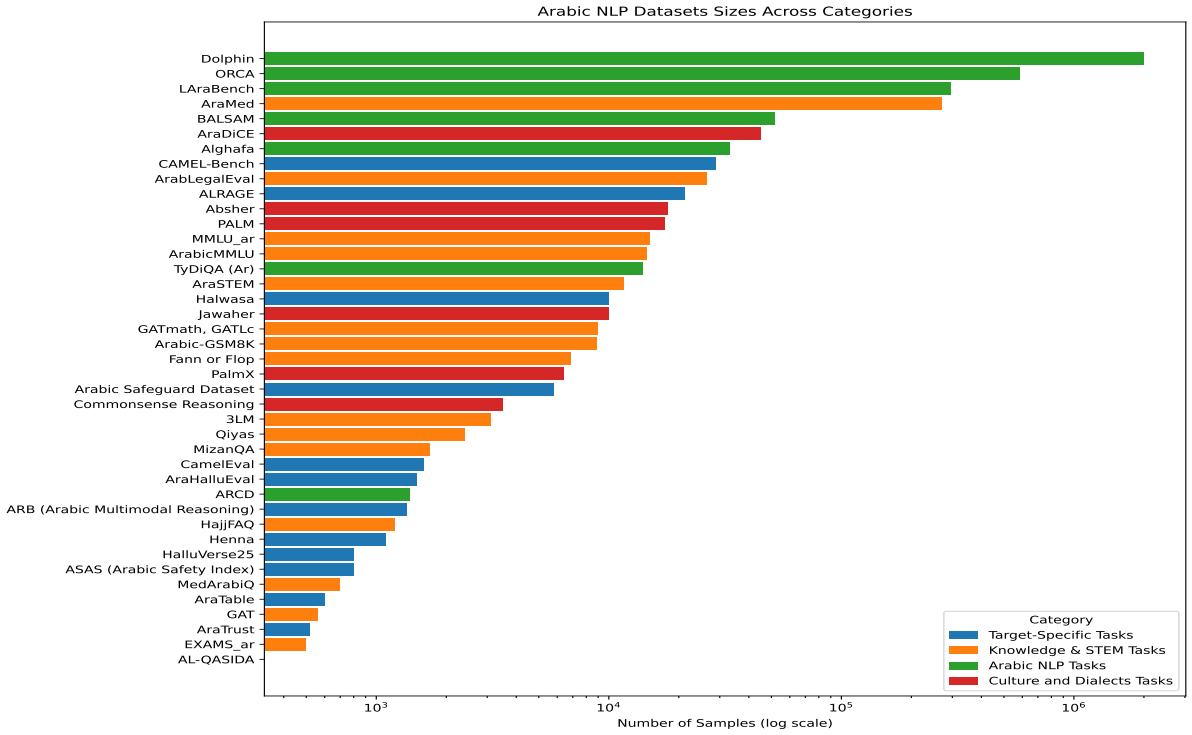


Figure 2: Arabic Datasets Sizes Across Categories. Size is log-scaled.

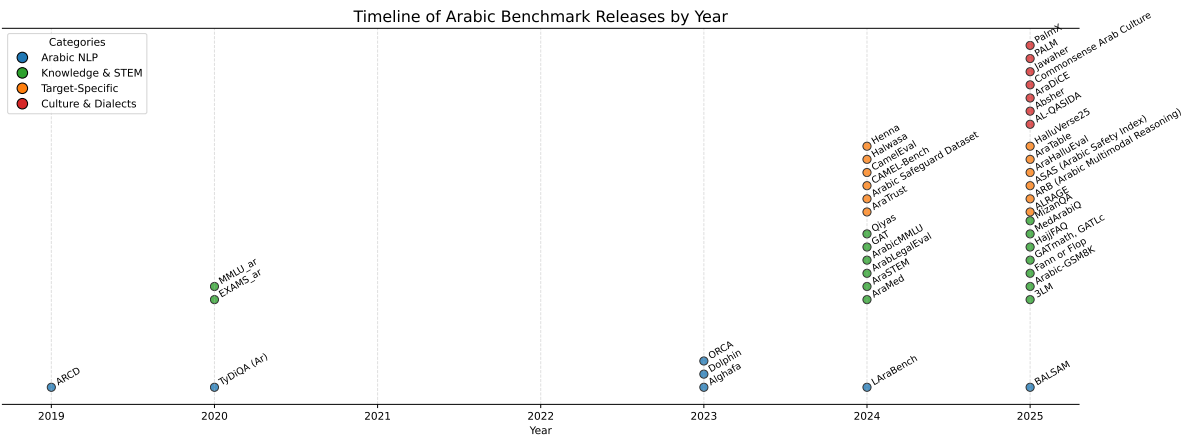


Figure 3: Timeline of Arabic benchmark releases (2019-2025) showing dramatic acceleration, with 82% of the benchmarks released in 2024-2025.

Year	Paper	Topic	NAT	SYN	TRAN	Total	Metrics	Type	PD	PR
Knowledge										
2020	MMLU_ar (Hendrycks et al., 2020)	Humanities • Social Science • STEM	-	-	15k	15k	Accuracy • Log-prob	MCQ	Yes	Yes
2020	EXAMS_ar (Hardalov et al., 2020)	Science • Social Science • Religion	0.5k	-	-	0.5k	Accuracy • Log-prob	MCQ	Yes	Yes
2024	ArabicMMLU (Koto et al., 2024)	Humanities • Social Science • STEM • Arabic Language	14.57k	-	-	14.57k	Accuracy • Log-prob	MCQ	Yes	Yes
2024	ArasTEM (Mustapha et al., 2024)	STEM	11.63k	-	-	11.63k	Accuracy • Log-prob	MCQ	No	No
2024	GAT (Alkaoud, 2024)	Linguistic Abilities	0.45k	-	-	0.56k	Accuracy	MCQ	No	No
2024	Oiyas (Al-Khalifa and Al-Khalifa, 2024)	Linguistic Abilities • Mathematics	2.4k	-	-	2.4k	Accuracy	MCQ	No	No
2024	ArabiLegalEval (Hijazi et al., 2024)	Law	79	10.58k	15.8	26.4k	Accuracy • LLM-as-Judge	MCQ, GEN	Yes	Yes
2024	ArabiMed (Alasmar et al., 2024)	Medical Domain • Healthcare	270k	-	-	270k	Accuracy • QA Metrics	GEN	Yes	Yes
2025	MizanQA (Bahaj and Ghogho, 2025)	Law	1.7k	-	-	1.7k	Accuracy • FI • ECE	MCQ	Yes	No
2025	Fann or Flop (Alghallabi et al., 2025)	Poetry	6.9k	-	-	6.9k	BLEU • chrF(++) • BERTScore	GEN	Yes	Yes
2025	GATmath, GATLc (AlBallaia et al., 2025)	Linguistic Abilities • Mathematics	9k	-	-	9k	Accuracy	MCQ	Yes	No
2025	3LM (Boussaha et al., 2025)	STEM • Coding	0.8k	1.74k	0.54k	3.1k	Accuracy • Log-prob	MCQ, GEN	Yes	Yes
2025	Arabic-GSM8K (Omarifical-Intelligence-Space, 2025)	Mathematics • Reasoning	-	-	8.9k	8.9k	Exact-Match	GEN	Yes	No
2025	MedArabiQ (Daoud et al., 2025)	Medical	-	0.7k	-	0.7k	Accuracy • BERTScore	MCQ	Yes	Yes
2025	Haji-FAQ (Alaid and Azmi, 2025)	Religious Domain • Islamic Jurisprudence	1.2k	-	-	1.2k	Accuracy • FI	QA	No	No
Arabic NLP Tasks										
2019	ARCD (Mozannar et al., 2019)	Reading Comprehension	1.4k	-	-	1.4k	EM • FI	MCQ	Yes	Yes
2020	TyDiQA (Ar) (Clark et al., 2020)	Reading Comprehension	14k	-	-	14k	EM • FI	GEN	Yes	Yes
2023	ORCA (Elmadany et al., 2023)	Sentiment • Text Classification • NER • QA • Paraphrase • NLI • Dialogue • Summarization • Paraphrasing • Writing	-	-	-	588k	Accuracy • FI • Pearson • EM	BENCH	Yes	No
2023	Dolphin (Nagoudi et al., 2023)	Dialect ID • MT	-	-	-	2m	BLEU • ROUGE • METEOR • BERTScore	BENCH	No	No
2023	Alghafa (Almazrouei et al., 2023)	QA • Reading Comprehension • Reasoning • Math • Commonsense	-	-	-	33.2k	Accuracy • Log-prob	MCQ	Yes	Yes
2024	LARA-Bench (Abdeljalil et al., 2024)	Sentiment • Topic Classification • NER • QA • Paraphrase • NLI • Dialogue ID • MT • Reasoning • Summarization	33.2k	-	-	296k	FI • BLEU • Task-specific	BENCH	Yes	Yes
2025	BALSAM (Al-Matham et al., 2025)	MT • Transliteration • DialectMT • Simplification • QREwrite • Paraphrase • Intent • GrammarCorr • GenderRewrite • TextClass • Sentiment • Sarcasm • DialectID • Command • Summarization • SubjectGen • AnswerExt • SeqTag	26k	1.7k	24k	52k	LLM-as-Judge • BLEU • ROUGE	BENCH	Partial	Yes
Culture and dialects' Tasks										
2025	Jawaher (Magdy et al., 2025)	Proverbs • Figurative QA	10k	-	-	10k	BLEURT • BERTScore • LLM-as-Judge • Human-as-Judge	GEN	Yes	No
2025	PALM (Alwajih et al., 2025a)	Prose • Dialogue • Cultural Expressions	17.4k	-	-	17.4k	LLM-as-Judge	GEN	Yes	Yes
2025	PalmX (Alwajih et al., 2025b)	Shared Task • Arabic Culture MCQs • Islamic Culture MCQs	-	6.4k	-	6.4k	Accuracy, Log-prob	MCQ	Yes	No
2025	Commonsense Reasoning in Arab Culture (Sadallah et al., 2025)	Cultural Commonsense QA	3.5k	-	-	3.5k	Accuracy • Log-prob	MCQ	Yes	Yes
2025	ArADICE (Mousi et al., 2025)	Dialect Identification/Generation/Translation • Cognitive Abilities on Dialect • Culture	41.8k	-	45k	81.8k	Accuracy	BENCH	Yes	Yes
2025	AL-QASIDA (Robinson et al., 2025)	Dialectal Analysis	-	-	-	-	Accuracy • Dialectal Error Rate	BENCH	No	No
2025	Absher (Al-Monef et al., 2025)	Saudi Dialect • Saudi Culture	18k	-	-	18k	Accuracy • FI	MCQ	No	No
Target-Specific Tasks										
2024	CamelEval (Qian et al., 2024)	Instruct-following	-	0.8k	0.8k	1.6k	Win-rate, LLM-as-Judge	GEN	No	No
2024	Hatwasa (Mubarak et al., 2024)	Hallucination	-	10k	-	10k	Human-as-Judge	GEN	No	No
2024	Henna (Alwajih et al., 2024)	Multimodal	-	1.1k	-	1.1k	LLM-as-Judge	VQA, OCR, GEN	Yes	Yes
2024	CAMEL-Bench (Ghaboura et al., 2024)	Multimodal	29.0k	-	-	29.0k	Accuracy, PMTS • LLM-as-Judge	VQA, OCR, RS, GEN	Yes	Yes
2024	Aratrust (Alghamdi et al., 2024)	Safety	0.52k	-	-	0.52k	Accuracy	MCQ	Yes	Yes
2024	Arabic Safeguard Dataset (Ashraf et al., 2025)	Safety	1.0k	-	4.8k	5.8k	LLM-as-Judge	GEN	Yes	Yes
2025	ALRAGE (El Filali et al., 2025)	Context-based (RAG)	-	21.2k	-	21.2k	LLM-as-Judge	GEN	Yes	Yes
2025	AratTable (Alshaikh et al., 2025)	Context-based (Tabular)	-	0.6k	-	0.6k	Accuracy • Human-as-Judge	GEN	Yes	Yes
2025	AratHalluEval (Alansari and Luqman, 2025)	Hallucination	0.4k	-	0.8k	1.5k	Human-as-Judge	GEN	Yes	Yes
2025	HalluVerse25 (Abdjalil et al., 2025)	Hallucination	-	0.8k	-	0.8k	Accuracy	CL	No	No
2025	ARB (Arabic Multimodal Reasoning) (Ghaboura et al., 2025)	Multimodal	1.36k	-	-	1.36k	BLEU • ROUGE • BERTScore • LabSE • LLM-as-Judge	VQA, RS	Yes	Yes
2025	ASAS (Arabic Safety Index) (Alastrolabe, 2025)	Safety	0.8k	-	-	0.8k	LLM-as-Judge	GEN	No	No

Table 6: Summary of Arabic Benchmarks. NAT: # of native samples, TRAN: # of translated samples, SYN: # of synthetic samples, PD: dataset is public, and PR: repo is public.