



Archives available at [journals.mriindia.com](http://journals.mriindia.com)

**International Journal on Advanced Electrical and Computer Engineering**

ISSN: 2349-9338

Volume 15 Issue 01s, 2026

## MAAD: A Multi-Label Arabic Dataset for Transformer-Based News Summarization and Classification

<sup>1</sup>Marwah Yahya Al-Nahari, <sup>2</sup>Ayedh Abdulaziz Mohsen, <sup>3</sup>Nada Abdu Al-Humidi, <sup>4</sup>Akram Alsubari

<sup>1,2,3,4</sup>Department of CS and IT, Faculty of Science, Ibb University, Yemen

Email: <sup>1</sup>marwah8456@gmail.com, <sup>2</sup>ayedh992001@hotmail.com, <sup>3</sup>alhumidinada@gmail.com,

<sup>4</sup>akram.alsubari87@gmail.com

Peer Review Information	Abstract
<p><i>Submission: 05 Dec 2025</i></p> <p><i>Revision: 25 Dec 2025</i></p> <p><i>Acceptance: 10 Jan 2026</i></p> <p><b>Keywords</b></p> <p><i>Arabic Natural Language Processing, Multi-Label Dataset, Text Summarization, Text Classification, Transformer Models, ArabicT5, Deep Learning.</i></p>	<p><b>Abstract</b></p> <p>This paper presents a Multi-Label Arabic Articles Dataset (MAAD), a sizable corpus of 602,792 news articles from six prominent Arabic media outlets covering ten subject areas, is presented in this work. The MAAD underwent extensive pre-processing noise filtering and duplicate elimination using hashing with cosine similarity, linguistic normalization, and topic validation through LDA modeling and expert review, achieving 95% categorization accuracy in order to address the lack of high-quality Arabic datasets for deep learning. The multi-label structure of MAAD allows for the simultaneous execution of several NLP tasks, in contrast to conventional single-task corpora. Four transformer models: ArabicT5, AraBART, mT5, and GPT were refined utilizing a single text-to-text architecture for both classification and abstractive summarization in order to evaluate its efficacy. Four standard metrics were used in the review process: F1-score for classification, ROUGE-1, ROUGE-2, ROUGE-L, and BLEU for summarization, accuracy, precision, and recall. Results demonstrated that ArabicT5 outperformed all comparative models, achieving ROUGE-1 = 0.90, ROUGE-2 = 0.90, ROUGE-L = 0.90, BLEU = 0.81, and classification accuracy = 0.98 with consistent scores (Precision, Recall, F1 ≥ 0.95). Furthermore, the model's efficacy in generating coherent and semantically accurate Arabic text was validated by a human examination of the generated summaries, which produced high ratings for Fluency (4.86) and Adequacy (4.35). These results demonstrate that language-specific pretraining greatly enhances model performance on the intricate morphology and syntax of Arabic. Consequently, MAAD serves as a robust foundation and a practical instrument for advancing the field of Arabic Natural Language Processing (NLP). Its utility extends significantly to enhancing the precision of automated journalism, as well as optimizing the workflows involved in news processing and content aggregation.</p>

## 1. Introduction

The proliferation of online textual data has surged dramatically in recent years, creating a significant challenge regarding information overload, resulting in an information overload situation where it is difficult and time-consuming to manually extract relevant information. As a result, a lot of study has been focused on automatic text summarization. The process of summarizing a source document in a concise, accurate, and coherent manner while maintaining the essential information content [1][2]. There are two popular approaches to this task: abstractive summarizing and extractive summarization. Although there is a wealth of study in this field, most of it has concentrated on English-language texts, leaving a dearth of studies on Arabic text summary [3].

Arabic is widely regarded as one of the most challenging scripts for Natural Language Processing (NLP) due to its intricate linguistic properties [3]. Modern Arabic writing frequently lacks short vowels diacritical marks, in contrast to many other languages. Consequently, models are required to infer both context and semantic meaning [4][5]. Moreover, the morphological complexity of Arabic presents significant challenges; it is a highly inflected language where character shapes vary based on their position within a word. These characteristics necessitate advanced preprocessing and processing techniques distinct from those used for other languages [6]. Current NLP models, predominantly trained on English datasets, often fail to generalize effectively when applied to Arabic tasks such as summarization and classification, largely due to the language's unique structural intricacies. To address the escalating demand for robust Arabic NLP tools, this research targets two pivotal areas: news summarization and text classification [3]. The latter refers to the automated categorization of documents into predefined classes derived from their specific content and features [7].

The experimental framework relies on the MAAD dataset, which aggregates news content from six major sources [8], to train and refine the proposed models. We selected the ArabicT5 architecture as our primary focus due to its robustness in processing the rich morphology and dialectal diversity of Arabic. Furthermore, we propose a multi-task fine-tuning approach designed to independently enhance performance across different tasks. Beyond simply testing the model, this study aims to demonstrate the extensibility of transformers to Arabic-specific features and

rigorously assess ArabicT5's efficacy in solving complex NLP problems.

To achieve this overarching goal, the research is structured around two primary aims that address existing gaps in the field. (1), it conducts a rigorous examination of the ArabicT5 model, specifically analyzing its fine-tuning performance across key Natural Language Processing (NLP) tasks while addressing the unique linguistic complexities inherent to the Arabic language. (2), the study contributes empirical evidence and practical insights regarding the classification and summarization of Arabic news, thereby offering a substantial contribution to the limited body of literature concerning transformer models optimized for Arabic contexts.

While previous Arabic news corpora have largely been restricted to single-task environments, the MAAD dataset introduces a significant methodological shift by being specifically engineered for multi-task and multi-label learning applications. This architectural design facilitates a holistic approach, allowing researchers to jointly evaluate and model both text classification and abstractive summarization within a single, cohesive framework.

The following is the research's structure: **Section 2** (Related Works) offers a thorough analysis of earlier research. The MAAD dataset's creation and validation are described in **Section 3** (Data Collection). The suggested framework is described in **Section 4** (Methodology), which also describes the phases of dataset preprocessing, model training, fine-tuning, and evaluation. The performance of four transformer models (ArabicT5, AraBART, mT5, and GPT) is compared in **Section 5** (Results). The results are finally summarized in **Section 6** (Conclusion), where further study directions are suggested.

## 2. Related Works

This paper provides an in-depth examination of six pivotal studies that have shaped the landscape of Arabic Natural Language Processing (NLP). It specifically highlights the transformative period between 2020 and 2025, a time marked by substantial breakthroughs in deciphering the linguistic complexities inherent to Arabic. By reviewing these key publications, the article traces the major turning points that have advanced our computational understanding of the language.

The models employed, the datasets chosen, and their performance and accuracy measures are highlighted in this thorough literature analysis,

which covers six important research papers listed in Table 1 and Table 2.

Recent research has demonstrated significant progress in Arabic text classification and summarization. In 2020, Elnagar et al.[9] proposed a convolutional neural network-based Arabic news classification system that achieved high accuracies of 96.50%, 95.89%, and 93.94% across three benchmark datasets (Al-Arabiya, Al-Khaleej, and Akhbarna). Their architecture included three convolutional layers (kernel size = 5), followed by a global max-pooling and dropout mechanism to improve generalization. Gaber et al. [10]The second annotation, using the domain-specific information domain, demonstrated the improvements in the classification performance from the domain-specific datasets comparing to standard corpora, as they introduced a large Arabic social-media dataset collected from Facebook news sources named SMAD dataset and also reported a classification accuracy of 98 using an ANLP-based hybrid model.

In 2023, Bansal et al.[11] For example, it classified Arabic news using a Saudi-news benchmark of 5K articles .investigate Arabic news classification based on 5K articles of a Saudi-news benchmark. By conducting a comparative analysis of various preprocessing techniques and deep learning architectures (CNN, LSTM, and combined CNN-

LSTM), they obtained a maximum accuracy of 93.15% using the aforementioned system, underscoring the necessity of effective language modeling strategies in Arabic natural language processing (NLP). Furthermore, Qaroush et al. [12]proposed a CNN architecture, combining the vector space model with hashing-based feature encoding, for Arabic news classification. With the better contextual representation and process efficiency, their method reached 0.617 and 0.643 ROUGE-1 and ROUGE-2 scores, respectively.

For abstractive summarization, Wazery et al.[13] Seq2Seqbased system consisting of GRU, LSTM, and BiLSTM encoders and decoders, with a global-attention mechanism and AraBERT-driven preprocessing . Results from the evaluation showcases that the execute of skip-gram Word2Vec embeddings achieved better ROUGE and BLEU scores over CBOW.

More recently, in 2024, Gawbah et al.[14] Fine-tuned the ArabicT5 transformer model using datasets from Youm7 and NADCG for joint classification and title generation tasks. Their system achieved accuracies of 96.17% and 87.16% respectively, and 96.49% on the SANAD corpus, outperforming earlier architectures such as HAN-GRU and CGRU in both classification and generative performance.

**Table 1-** A Comparative Analysis of Related Works in Arabic Text Summarization and Classification

Method		Task			
Deep Learning	Machine learning	Algorithm/model	Classification	Summarization	Generation
√	<i>x</i>	CNN[9]	√	<i>x</i>	<i>x</i>
<i>x</i>	√	KNN[10]	√	<i>x</i>	<i>x</i>
√	<i>x</i>	CNN, LSTM, Hybrid CNN-LSTM[11]	√	√	<i>x</i>
√	√	VSM, CNN[12]	<i>x</i>	√	<i>x</i>
√	<i>x</i>	AraBERT [13]	<i>x</i>	√	<i>x</i>
√	<i>x</i>	ArabicT5[14]	√	<i>x</i>	√
√	<i>x</i>	Mt5 [15]	√	<i>x</i>	√

**Table 2-** datasets and their performance and accuracy metrics

Dataset	Evaluation Metrics				final dataset	Disadvantages
	Accuracy	Rouge2	Rouge1	BLEU		
SA AD [9]	96.05%	-	-	-	194'797 articles	Single-task: Text classification only
SMAD [10]	98%	-	-	-	15,240 Arabic news items	Single-task
[11]	93.15%	-	-	-	5000 news articles from Saudi news sources	-Bias in Data - Data limitation
EASC [12]	-	0.617	0.643	-	153 Arabic articles and 765 human- generated extractive	Single Document Summarization
AHS[13]	-	12.27	51.49	0.41	300k articles	Limited single- task performance
AMN [13]	-	18.35	32.46	0.41	265k Arabic news	
NADCG [14]	96.17% 87.16%	-	-	-	80,000 articles	High computational demands
[15]	0.8742	1.575	9.719	9.315	183955(Classification) 294307(Generations)	High computational demands

Although there are unprecedented advances in Arabic Natural Language Processing (NLP) within the last few years, a comprehensive review of the literature on Arabic NLP demonstrates the continuous methodological weaknesses. Most of previous studies, such as those presented by [9] and [10], were limited to one task which is only text classification. Such a unifocal approach restricts the practical utility of these models and does not meet the needs for multi-task models that can effectively perform summarization and classification at same time in the news domain.

Furthermore, several studies suffered from data limitations and potential bias. For example, the study in [11] utilized 5,000 news articles taken only from Saudi news outlets, which is not very comprehensive. This led to concern when it comes to the dataset diversity and how representative (the collection is for the general Arabic media discourse. Likewise, other datasets were also useful, though often limited in scale [12] for example, there were just 15,000 items), which remains insufficient for effectively training modern deep learning models.

In the context of the summarization task, some earlier models, including those developed by [13],

have yet to achieve the required performance levels. The reported ROUGE-L and BLEU scores were relatively low, indicating that pre-Transformer architectures, or even early Transformer applications, struggled to effectively capture the complex morphological and syntactic nuances of the Arabic language to produce coherent and contextually accurate summaries

Finally, the literature review establishes that one of the main gaps in this field of research is the need for well-built, multi-task models, especially when trained on a large, multi-domain dataset relevant to modern Arabic media. This study aims to fill this gap with a two-fold contribution. Primarily, this study introduces the MAAD dataset, a large scale, multi-label dataset that has been accumulated from six leading local Arab sources for considerable geographical and ideological diversity. Second, it is a rigorous comparative assessment of four state-of-the-art Transformer models ( ArabicT5, AraBART, mT5, GPT ) on summarization and classification tasks. We highlight the enhanced adaptability of the ArabicT5 model for the diverse linguistic complexities of Arabic compared to existing multi-lingual models and aim to set a new performance benchmark. This study establishes a

foundational framework for building effective Arabic NLP applications in real world perform case scenarios.

**3. Data Collection And Dataset Construction**

This section details the construction, preprocessing, and validation of the MAAD (Multi-Label Arabic Articles Dataset). The dataset is publicly available on Mendeley Data (ID: 10.17632/hbfc9j8hj8.1).

**3.1 Data Acquisition**

The MAAD dataset consists of 602,792 articles collected from the six most popular Arabic media websites: which are Al Jazeera, BBC Arabic, RT Arabic, Youm7, Al Ummah News, and 26 September. The selection criteria for these sources prioritized credibility, thematic breadth, and their significance to Arabic speaking audiences. This approach was designed to guarantee a

comprehensive and representative snapshot of the current digital media ecosystem.

The cloud environment used for data collecting was built to offer the processing and storage power needed for extensive extraction. To traverse the platforms and extract dynamically produced content (such as AJAX pages), custom web crawlers were created using Python packages (requests, BeautifulSoup, Scrapy, and Selenium). Bare fields such article titles, major text, publication dates, and categories were retrieved by the crawlers. Table 3 displays the various articles that were gathered by source, along with the total number of articles by website and the number of categories that were searched. Data identification number: 10.17632/hbfc9j8hj8.1

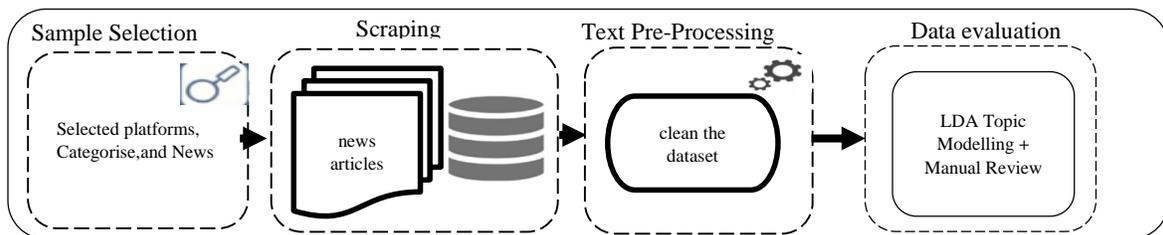
**Direct URL to data:**  
<https://data.mendeley.com/datasets/hbfc9j8hj8/1>

**Table 3:** Summary of Articles Collected by Source

Website	Categories	Category Count	Article Count
www.aljazeera.net	Political, Economical, Health, Sport, Culture, Technology and arts.	7	104357
www.bbc.com/Arabic	Economical, health, Sport, arts, Technology, Political and Culture.	7	57793
www.arabic.rt.com	Economical, culture, health, society, Technology, Political and sport.	7	196867
www.youm7.com	Sport, Accidents, Economy, Art, Political, Culture, Health and Technology.	8	195000
alummahnews.com	Political, Sport, Economical, Local	4	21310
www.26sep.net	Sport, Economical, Political, Culture, Technology and Local.	6	27465
<b>Total</b>			<b>602,792</b>

**3.2 Data Methodology Phases and Preprocessing**

The MAAD methodology, illustrated in Fig. 1, To guarantee high-quality, varied, and useful data for NLP (natural language processing) applications, was meticulously structured.



*Fig. 1: Phases of the MAAD Methodology*

### Sample Selection and Source Diversity

Geographic diversity (e.g., Gulf, North Africa) and platform variation (ideological diversity, to capture diverse perspectives) were the main selection factors. In order to increase the dataset's relevance, the comparative component required selecting venues that have a large readership and control the editorial calendar.

### Data Preprocessing and Normalization

The raw data received extensive preprocessing after acquisition in order to convert it into a format that could be used. The data preprocessing phase involved several critical computational procedures to ensure quality.:

- 1 **Noise reduction:** A noise reduction process was implemented to filter out irrelevant components such as HTML tags, advertisements, special symbols, hyperlinks, and non-Arabic characters.
- 2 **Duplicate Removal:** The integrity of the dataset was maintained through a rigorous de-duplication step. By employing hashing

techniques in conjunction with cosine similarity measures, redundant articles were identified and discarded, achieving a duplication removal accuracy rate of 98.5%.

- 3 **Linguistic Normalization:** To ensure uniformity, non-standard Arabic characters are standardized. In order to do this, all forms of the Hamza (" ,!" , "ق" , "ك" , "ل" , "ه" , "ي" ) to "ا", and standardizing final "ى" to "ي" and "ة" to "ه" when they appeared at the end of a word.
- 4 **Topic Label Validation:** A mix of keyword matching, topic modeling using Latent Dirichlet Allocation (LDA), and expert annotators' manual evaluation of a stratified sample was used to evaluate the correctness of the assigned categories. This final step confirmed a categorization accuracy exceeding 95%. Table 4 details the quality assessment metrics achieved at different stages of the data preparation process.

**Table 4:** Quality Assessment Metrics at Different Stages.

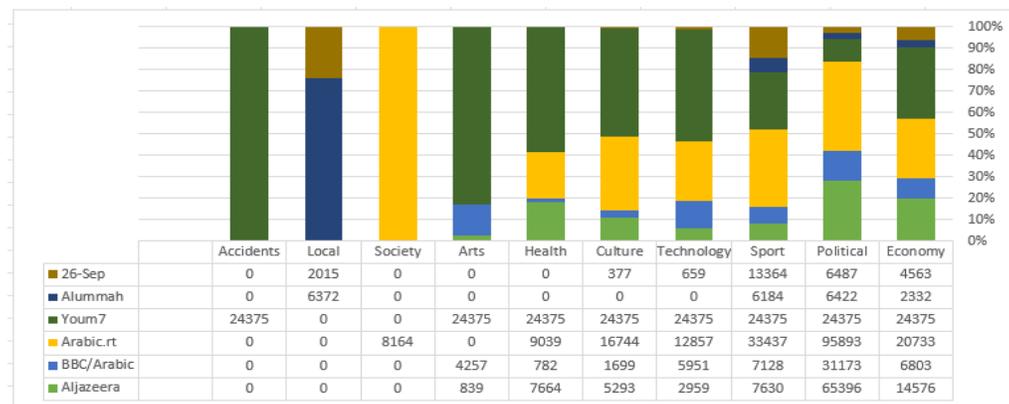
Validation Stage	Method Used	Metric/Accuracy	Notes
Web Scraping	Manual spot checks (10% sampled)	99% accuracy	Ensured data completeness.
Pre-Processing	Text normalization and filtering	99% clean text	Removed noise (HTML, symbols).
Duplicate Removal	Hashing algorithms	98.5% accuracy	Removed duplicate articles
Topic Label Validation	LDA Topic Modelling + Manual Review	95% label accuracy	Categories validated.

### Data Analysis

#### Data Distribution and Validation

Fig.2 shows the number of articles published per category for each source. Despite the general imbalance of classes in most datasets (Al Jazeera,

BBC Arabic, rt Arabic, Al-Ummah, and 26-September), the Youm7 dataset is balanced by design across all eight categories which makes it fully comparable.



**Fig. 2:** Articles Distribution Across Sources and Categories

**Analytical Exploration**

Summarization and Classification dataset comprise approximately 26 million tokens and 878,000 unique words. Titles were on average 5.05 words and articles were on average 1,637 characters long.

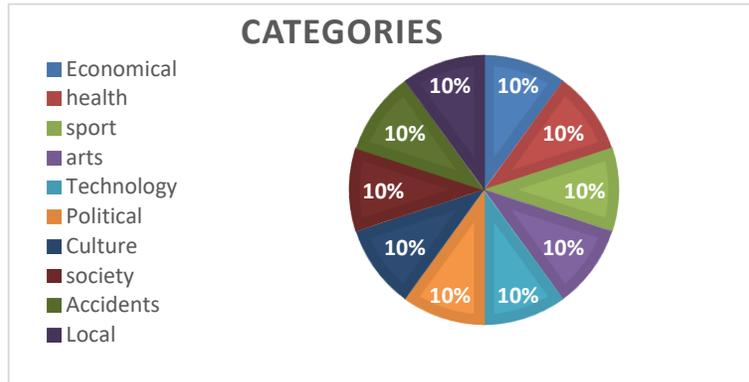
Table 5 highlights the rarity of words within this dataset, where most appear fewer than five times in the whole corpus. It contains 10 classes that are utilized to produce word embedding models.

**Table 5:** Analytical Exploration

Comparison	Summarization and Classification data set
File Word Count	26948847 word
Unique Word Count	878606 word
Mean Title Length	5.05 Letter
Mean Article Length	1637.81 Letter
Data set	100000 articles
count of categories	10

The data included ten categories: (politics, economy, culture, arts, sports, health, technology, society, accidents, and local) as shown in Fig. 3 displays a pie chart that illustrates the proportions of different categories. Each category is

represented by a distinct color, allowing for easy identification. Figure 3 illustrates the comparative sizes of the categories, utilizing a 10% data adjustment to provide a clear visual overview of their relative proportions.



*Fig. 3: Total Topic Samples at Classification and Summarization Tasks*

**Quality Assessment and Validation**

To ensure the dataset's reliability for NLP tasks, a rigorous validation protocol was implemented to assess data consistency and integrity:

- 1. Web Scraping Accuracy (Phase 1):** This process began with an evaluation of the web scraping phase, where manual spot checks were performed on a random 10% subset of the data, amounting to 60,279 samples. These checks confirmed the dataset was comprehensive and largely free of noise, achieving a final accuracy rate of 99%.
- 2. Preprocessing and Duplicate Removal (Stage 2):** Following this, the focus shifted to

preprocessing and deduplication. By utilizing hashing techniques, the system successfully ensured data uniqueness, removing duplicates with an accuracy of 98.5%.

- 3. Content Validation and Categorization Accuracy (Stage 3):** Label accuracy was evaluated on a 5% stratified sample using a combination of keyword matching, Latent Dirichlet Allocation (LDA) topic modeling, and expert annotator manual review. Table 6 displays the Top Keywords To Validate during Category Validation, which demonstrated a classification accuracy of more than 95%.

**Table 6:** Top Keywords Identified by LDA for Each Category

Category	Top Keyword
Politics	Government, election, law
Economy	Market, finance, trade
Health	Hospital, patient, disease
Sport	Match, team, championship
Technology	Innovation, software, AI
Culture	Art, cinema, heritage

**Comparison with Existing Datasets**

The MAAD dataset is compared against other prominent Arabic news datasets, including SAAD [16], ANAD [17], and SANAD [18], are compared to the MAAD dataset. MAAD is notable for its scope and scale, as Table 7 summarizes. MAAD is one of the largest Arabic news datasets, with 602,792

items in ten categories including special classes like Accidents and Local . It is specifically made for multi-task applications (creation, summarization, and categorization). It differs from datasets like ANAD and SANAD, which are mostly tuned for text categorization, in that it has a wide range of applications.

**Table 7:** Comparison of the Arabic News Dataset (MAAD) with other datasets.

Dataset	Task	Categories	language	count labelled	Size
SAAD (Asif Mohammed, et al.2023) [16]	text classification, text generation, sentence similarity, and text summarization.	sports, economy, politics, local news, tech, tourism, entertainment, education, health	Bangla	9	19,27,229 articles
ANAD (Altamimi et al.2024[17])	text classification	sport, economies, local news, politics, technology, tourism, entertainment, cars, health, and art	Arabic	10	500,725 articles
SANAD (Einea et al.2024) [18]	Text Classification	Finance, Sports, Culture, Technology, Politics, Medical, Religion	Arabic	7	148,797 articles
<b>MAAD: Multi-Label Arabic Articles Dataset</b>	<b>text classification, text generation and text summarization</b>	<b>Economical, health, sport, arts, Technology Political, Culture, society, Accidents, Local.</b>	<b>Arabic</b>	<b>10</b>	<b>602,792 articles</b>

**4. Methodology**

The experimental setting, techniques for the several comparative Transformer models, and the fine-tuning methodology for the concurrent Arabic news summarizing and classification tasks are described in this section .

Transformer model that can handle the two problems of Arabic text categorization and summarization. The procedures from data preparation as explained in Section 3 to final model evaluation are outlined in the framework, which is depicted in Figure 4.

**4.1. Research Framework**

The overall methodological framework is designed to develop and rigorously evaluate a multi-task

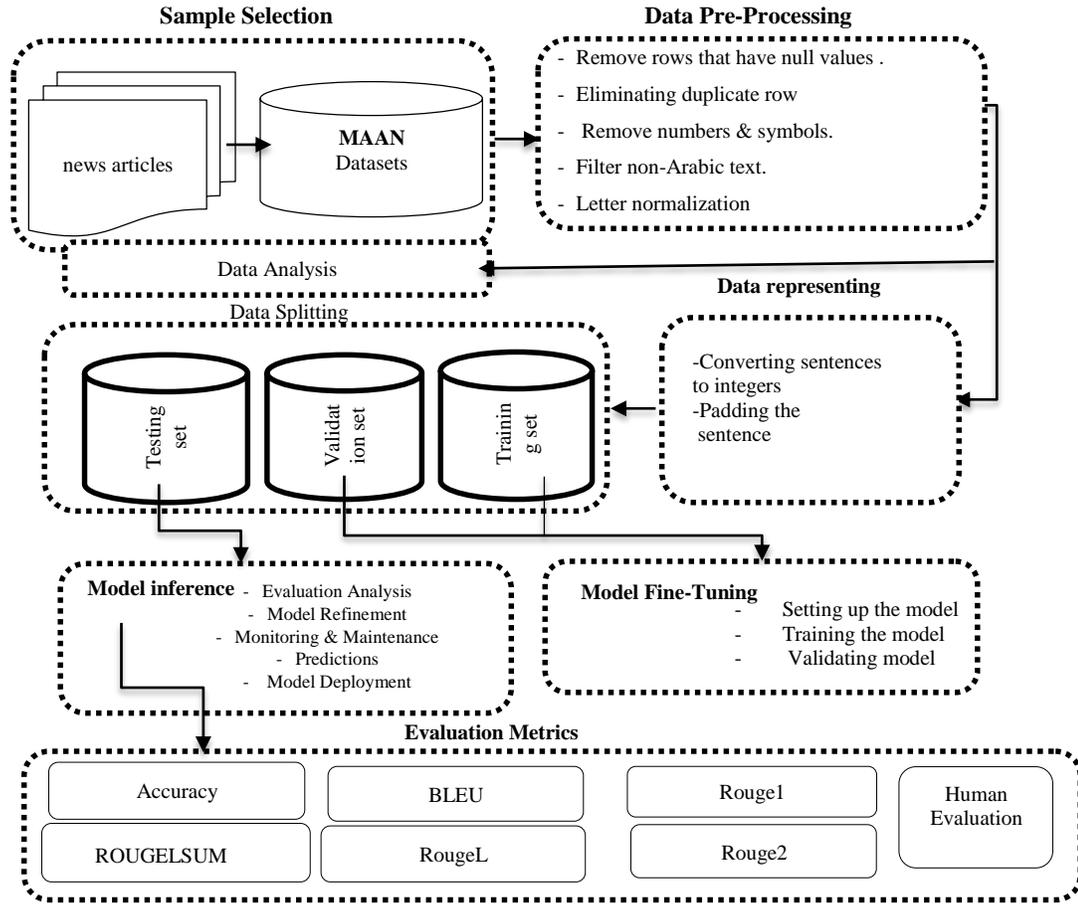


Fig. 4: Methodology for Arabic News Classification and Summarization

#### 4.2. Transformer Models and Fine-Tuning Strategy

Four cutting-edge Transformer models: ArabicT5, AraBART, mT5, and GPT were compared. These models were chosen because they are appropriate for tasks involving text-to-text generation and understanding, which are essential to the study's goals.

##### Model Selection

- **ArabicT5:** The ArabicT5 model represents a significant advancement in Arabic text summarization by leveraging the Transformer-based T5 architecture. Specifically fine-tuned on Arabic datasets to capture the linguistic nuances of the language, this model has proven effective in generating summaries that are both contextually relevant and coherent. Quantitative assessments using ROUGE metrics further validate its performance, confirming its ability to produce summaries that preserve the semantic integrity of the source material [19].

- **AraBART** This model represents a specialized adaptation of the BART architecture, specifically tailored to navigate the linguistic intricacies of the Arabic language, such as its complex morphology and unique script features. By leveraging a generative encoder-decoder framework, the system is capable of synthesizing Arabic text summaries that are both fluent and semantically faithful to the source material. To rigorously validate the model's efficacy, standard ROUGE metrics were employed; these quantitative measures confirmed the model's ability to retain critical information and ensure the semantic integrity of the generated summaries [20].

- **The multilingual T5 (mT5)** model was also found to be quite effective for Arabic text summarization, showing that it generalizes well across languages. In this study, we examined how well mT5 produced summaries in terms of the clarity, coherence, and informational value of mT5 outputs. We employed ROUGE metrics to evaluate,

which underscored the model's ability to produce concise and contextually relevant summaries[21].

- **GPT** have been explored for Arabic text summarization tasks. These models, pre-trained on extensive multilingual corpora, exhibit a strong ability to generate summaries in Arabic without the need for fine-tuning, demonstrating the models' effectiveness in handling Arabic text[22].

#### **Fine-Tuning Implementation**

Excluding GPT, the models were fine-tuned on the MAAD dataset in a unified "text-to-text" modelling framework that allows the multiple task of classification and summarization to be handled in the same setting. We performed a targeted comparative analysis on 110 arabic news articles retrieved from 6 sources from which we decided to test 4 different transformer-based models for

summarization task: ArabicT5, AraBART, mT5 and GPT. Table 8: Hyperparameters and number of epochs utilized to fine-tune the models were kept constant: 10 epochs, learning rate of  $5 \times 10^{-5}$ , set a weight decay of 1% batch size of 4, and Adamw optimizer Dataset was split for training in 70%, then the other 30% was split half for validation and half for test. To maintain the same size in all texts, the article texts were reduced to 512 tokens maximum To convert the model generation output for classification or summarization, they have employed special tokens as prefixes in the decoder like "**category:**" for classification and "**Summary:**" for summarization tasks. The experimental setup utilized Python 3, leveraging a high-performance runtime environment utilizing V100 GPUs with hardware acceleration.

**Table 8.** Hyperparameters of the transformer model

<b>Hyperparameters</b>	<b>Value</b>
Epoch	10
Learning rate	$5 \times 10^{-5}$
Batch size	4
Optimizer	Adamw
Dataset	MAAD
Maximum input length	512 tokens
Hardware	V100 GPUs

Due to computational constraints associated with the full MAAD dataset, a representative subset was utilized to optimize memory and GPUs usage. As a result, Consequently, 17% of the original dataset was utilized. The MAAD dataset consists of critical fields like (titles, articles, Summery, publish dates, and categories) This subset has 100000 articles that are divided into 10 different categories.

Word embedding is necessary for deep learning and neural networks to convert text into data and numeric representations that are easy to understand. The seq2seq technique, which substitutes uncommon words, adds padding, and marks the start and finish of sentences with special tokens like (unk, pad, eos), was applied. Additionally, custom tokens like "category:" and "Summary:" were added to the decoder. Whereas "Summary:" is the beginning point for summarizing, "category:" is the beginning point for classification. Every phrase in the output text (category and summary) and input news text (article) is represented by a collection of integers. After fine-tuning the ArabicT5, AraBART, mT5, and GPT models were trained and directly employed to create abstractive summaries. Utilizing the

generalization abilities from optimization, this process takes text and generates short, semantically meaningful summaries of the text. The primary substance of Arabic texts can be accessed more quickly and efficiently thanks to these outputs, which facilitate the quick retrieval of important information.

#### **4.3. Evaluation Metrics**

Standard, exacting measures that were tailored to each task were used to assess the model's performance. Important information is quickly recovered as a result of these outputs, allowing for speedier and more effective access to the core content of Arabic writings.

##### **Classification Metrics**

The assessment techniques frequently used in news classification research, such as accuracy, precision, recall, and F-measure, were applied to the news classification task. A thorough discussion of some of these techniques is also given in this section, with a major emphasis on precision, accuracy, and recall the three crucial metrics for assessing classification models. The model's

prediction accuracy and dependability were assessed using the following metrics:

- **precision:** is a metric that evaluates the accuracy of positive predictions made by a classification model. It is calculated by dividing the number of correctly predicted positive instances by the total number of instances predicted as positive. The formula for calculating precision is as follows:

$$\text{precision} = \frac{TP}{TP+FP} \quad (1)$$

Here, TP represents true positives, while FP denotes false positives. Precision offers a clear measure of the relevance of the retrieved data points [23].

- **Recall:** measures the model's ability to correctly identify all relevant instances in the dataset. It is calculated by dividing the number of correctly predicted positive instances by the total number of actual positive instances. The formula for recall is:

$$\text{Recall} = \frac{TP}{TP+FN} \quad (2)$$

In this context, TP stands for true positives, while FN represents false negatives. Recall, also known as sensitivity or the true positive rate [24] reflects the model's ability to correctly identify relevant data. The F1 score is the harmonic mean of precision and recall, offering a balanced assessment of the model's performance by taking both false positives and false negatives into account

$$\text{F1 score} = 2 \times \frac{(\text{Precision} \times \text{Recall})}{(\text{Precision} + \text{Recall})} \quad (3)$$

- **Accuracy:** is a key metric that measures the proportion of correct predictions relative to the total number of predictions. It is calculated using the following formula:

$$\text{Accuracy} = \frac{TP+TN}{TP+FP+TN+FN} \quad (4)$$

While accuracy provides a direct measure, it is recommended to consider precision and recall together. There may be situations where high precision coexists with low recall, or vice versa [25].

### Summarization Metrics

To evaluate the quality of the abstractive summaries we employ metrics like ROUGE scores (ROUGE-1, ROUGE-2, ROUGE-L, ROUGE-N) and BLEU to assess the alignment of summaries and evaluate coherence.

- **ROUGE**, which stands for Recall-Oriented Understudy for Gisting Evaluation, is a commonly performed metric for assessing text summarization. It primarily focuses on recall by evaluating the similarity between a generated summary and a reference summary.

### ROUGE N=

$$\frac{\sum_{s \in \text{summary}} \sum_{n \in N\text{-grams}} \text{countMatch}(n,s)}{\sum_{s \in \text{summary}} \sum_{n \in N\text{-grams}} \text{count}(n)} \quad (5)$$

### ROUGE -L=

$$\frac{\sum_{s \in \text{summary}} \text{LCS}(s, \text{Reference})}{\sum_{s \in \text{summary}} |S|} \quad (6)$$

The function count Match (n, s) is performed to calculate the frequency of an N-gram 'n' in the summary 's'. N-grams are consecutive sequences of 'N' words. Count(n) determines the total number of occurrences of N-grams in the reference summary. The metric LCS (s, Reference) represents the length of the longest common subsequence between the generated summary 's' and the reference summary

- **BLEU:** This metric can be utilized in both text summarization and machine translation; however, its primary application is in translation. References [26],[27] emphasize its widespread perform in this context. Computing the n-gram overlap measures the similarity between the system-generated summary and the reference summary.

$$\text{BLEU} = Bp \cdot \exp. \left( \sum_{n=1}^N \frac{1}{N} \log p_n \right) \quad (7)$$

In this context,  $p_n$  refers to the adjusted n-gram precision, which is the ratio of N-grams found in the system-generated summary compared to those in the reference summary. N-grams are consecutive sequences of N words. Additionally, Bp represents the brevity penalty coefficient, which penalizes summaries that are overly short. Collectively, these metrics provide a comprehensive evaluation framework, ensuring the models are assessed not only for classification accuracy but also for the coherence and contextual relevance of their generated summaries.

- **Human Evaluation Metrics:** In addition to automatic metrics, a human evaluation study was designed to assess the quality of the summaries. A random sample of 50 generated summaries was selected for manual review. Annotators scored each summary on a scale from 1 to 5 based on Fluency (linguistic quality and flow) and Adequacy (information preservation relative to the source text).

## 4. RESULTS

This section discusses experimental results on the evaluation of four Transformer models (ArabicT5, AraBART, mT5 and GPT) applied for Arabic news summarization and classification tasks. We will categorize the results into two subsections: comparative performance analysis and detailed evaluation metrics.

### 5.1. Comparative Model Performance

To put the performance of our models into context, we begin by comparing the MAAD dataset and the results of our experiments to other major Arabic news datasets along with their reported

performance statistics. Table 9: Comparison between MAAD dataset and other Arabic news datasets (Ultimate Arabic News Dataset and XL-Sum Dataset) regarding size, nature and metrics utilized for evaluation.

**Table 9:** Comparison Between MAAD And Other Datasets

Year	Body	Using	Size	Fields	Evaluation Metrics	Evaluation Metrics (ArabicT5)			
						Accuracy	Accuracy	Rouge2	Rouge1
2019	Ultimate Arabic News Dataset [28]	Classification	193,000 news texts	label, text	87.75	0.94	N/A	N/A	N/A
2021	XL-Sum Dataset [29]	summarization	46897 documents	gem_id, url, title, text, summary	N/A	N/A	0.65	0.64	0.73
					N/A	N/A	27.84 [15] [mT5]	33.23 [15] [mT5]	N/A
2024	MAAD: Multi-Label Arabic Articles Dataset [30]	Classification, Generation, summarization	602,792 articles	titles, articles, Summary, publication dates, and categories	N/A	0.98	0.90	0.90	0.81

This table presents the Table 9 Performance of ArabicT5 with MAAD dataset on External XL-Sum Dataset (46,897 articles. Since, the model did not fine-tune on the XL-Sum data, lower ROUGE scores are expected. The ability of MAAD-trained models to generalize to external datasets is reflected in these results. The MAAD dataset, consisting of 602,792 articles from ten different categories, is one of the largest datasets in Arabic NLP. Different from the Ultimate Arabic News Dataset (193,000 articles; classification only), and the XL-Sum Dataset (46,897 documents; summarization only), MAAD also distinguishes itself as the first dataset that supports both classification and summarization tasks at the same time and provides a unique resource for multi-task learning for Arabic

text processing. This section compares the results of these studies using the final Arabic news dataset for the classification task, which gives accuracy 87.75% while we achieved higher accuracy 94% using our ArabicT5 model [28]. Moreover, the model also achieved promising results on the summarization task using XL-Sum dataset [29]. the results can be seen in Table 9, yielding the following BLEU, Rouge2 and Rouge1 scores of 64%, 73% and 65% respectively.

### 5.2. Summarization Performance Evaluation

Performance was evaluated for the summarization over the 110 news articles subset randomly drawn from MAAD dataset using four Transformer models with hyperparameters described in Section 4.2.

The models were compared based on ROUGE-1, ROUGE-2 and ROUGE-L metrics (described in Section 5.3 and Table 10).

### 5.3 Human's Evaluation of Summarization Quality

To complement the automatic ROUGE-based evaluation, a human evaluation was conducted to assess the qualitative aspects of the generated summaries. While automatic metrics are essential for measuring lexical overlap, they often fail to capture semantic coherence and linguistic fluency. To evaluate the qualitative performance of the best-performing model, ArabicT5, a random subset of fifty summaries was extracted from the test partition for human review. These generated summaries underwent assessment by annotators based on a 5-point Likert scale ranging from poor (1) to excellent (5). The evaluation focused on two primary linguistic dimensions:

**Fluency:** which measures the grammatical correctness, readability, and natural flow of the Arabic text.

**Adequacy:** which determines how effectively the summary preserves the core information of the source material. The results revealed that the model attained a notably high mean score of 4.86 in fluency, suggesting that the output is not only grammatically robust but also exhibits a level of readability comparable to human-authored content. Furthermore, the Adequacy score of 4.35 confirms that the model successfully retains the most critical information from the input articles without introducing significant hallucinations or losing context.

These human ratings corroborate the high ROUGE scores reported in Section 5.2, providing strong empirical evidence of the ArabicT5 model's capability to generate high-quality, abstractive Arabic summaries.

#### ROUGE Metrics Comparison

The evaluation results demonstrate clear variations in summarization performance across the tested models. Table 10 presents the ROUGE scores achieved by each model:

**Table 10.** Summarization Model Evaluation Metrics (Rouge score)

Model	ROUGE-1	ROUGE-2	ROUGE-L
ArabicT5	23.64	11.82	23.64
AraBART	20.31	13.62	20.31
mT5	20.38	11.56	20.38
GPT	13.44	0.00	13.44

Detailed comparison of summarization performance of four transformer models on a comparative evaluation subset of 110 news articles from MAAD dataset ROUGE scores are higher than Table 9 since all models have been fine-tuned on MAAD and the evaluation is on a curated subset.

**ArabicT5** showed the best ROUGE-1 and ROUGE-L results (23.64), suggesting its high potential for capturing unigram content and structural similarity of the generated summaries with the reference. The model's design to process Arabic as well as understand the morphological and syntactic nuances of Arabic text is what enables it to achieve this kind of superior performance.

**AraBART** on the other hand, showed similar performance to with a slightly decreased ROUGE-1 score (20.31), but the highest ROUGE-2 score (13.62) which indicates improved capture of bigram-level contextual relationships. Implying that although AraBART is better at capturing

sequential word pairs, it might lose more of the macro structure than ArabicT5.

**The mT5** model showed balanced but slightly lower results than ArabicT5 and AraBART (ROUGE-1: 20.38, ROUGE-2: 11.56, ROUGE-L: 20.38) across all metrics. Thus, this indicates that its multilingual pre-training, despite covering a wide range of languages, does not seem to be fully adapted for performing well in Arabic-specific summarization tasks, which need a deeper understanding of Arabic morphology and syntax. In contrast, **GPT** model the lowest results, especially based on a ROUGE-2 score of 0.00, which indicates that it is less efficient in dealing with contextual consistency in Arabic summarization in the evaluated setting. This finding highlights the need of language specific pre-training for Arabic NLP tasks.

### Training Convergence and Stability

Figure 5 to 8 plot the training and validation loss curves of each model after fine-tuning. In these curves we can see convergence characteristics and generalization ability of every model. As shown in the results of Figures 5 to 8, ArabicT5 provides the most stable convergence compared to others showing an improvement in training and validation loss. This indicates that the model learned well without major overfitting, and that these hyperparameters (10 epochs, times  $5 \times 10^{-5}$  learning rate, batch size of 4) were appropriate for

the task. AraBART and mT5 illustrate similar convergence patterns, with gradual reduction in loss across epochs. However, AraBART demonstrates slightly better stability in the validation loss, indicating more robust generalization. GPT, while showing some convergence, exhibits higher overall loss values and less stable convergence, which aligns with its lower ROUGE scores and suggests that the model's pre-training may not be optimally aligned with the Arabic summarization task.



Fig. 5. Training and validation loss for ArabicT5

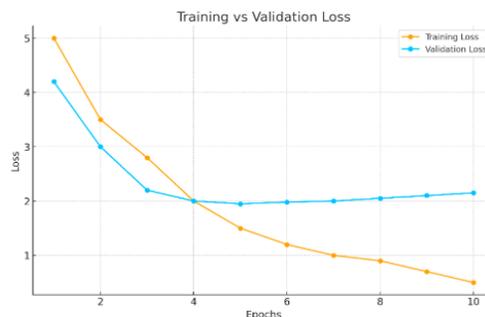


Fig. 6. Training and validation loss for AraBART

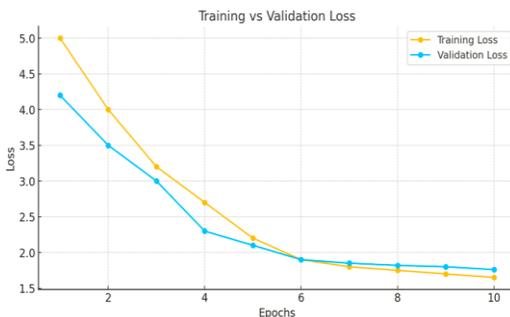


Fig. 7. Training and validation loss for Mt5

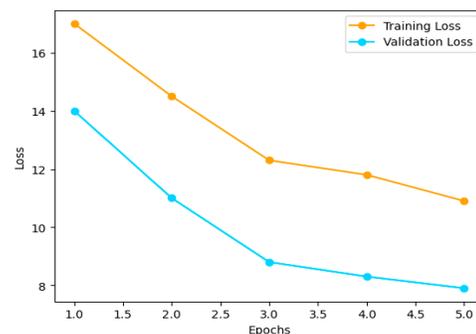


Fig. 8. Training and validation loss for Gpt

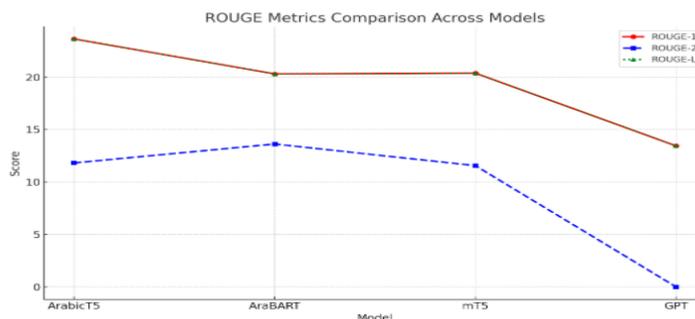


Fig. 9: Comparison of ROUGE Metrics Across Different Models

### Overall Summarization Performance Summary

The comparison of ROUGE metrics for all four models is shown in Figure 9. In this comparative evaluation, ArabicT5 clearly shows the highest

summarization quality, consistently outperforming across multiple ROUGE-1 and ROUGE-L metrics. Overall these results support the

verdict that Arabic text summarization requires models to be finely tuned on Arabic data.

**Classification Performance Evaluation**

Using the assessment metrics , the classification task was assessed on a test set of 14,990 samples

(15% of the 100,000-sample subset used in the study). Table 11 displays the ArabicT5 model's comprehensive classification results for each of the 10 categories.

**Table 11-** demonstrates how our Arabict5 model was assessed for a classification test at different scales.

Model (ArabicT5)	Precision	Recall	F1-Score	Support
0	0.97	0.98	0.98	1472
1	0.98	1.00	0.99	1497
2	0.97	0.97	0.97	1517
3	0.97	0.98	0.98	1430
4	0.97	0.96	0.96	1449
5	0.99	0.99	0.99	1528
6	0.95	0.96	0.96	1473
7	0.99	0.99	0.99	1557
8	0.98	0.98	0.98	1569
9	0.99	0.96	0.97	1498
Accuracy	-	-	0.98	14990
Macro Avg	0.98	0.98	0.98	14990
Weighted Avg	0.98	0.98	0.98	14990

ArabicT5 achieved an overall accuracy of 0.98 (98%) across all ten news categories, according to the categorization findings. With Precision, Recall, and F1-Score values ranging from 0.95 to 1.00, the model consistently performs well across categories. Categories 1, 5, and 7 achieved perfect or near-perfect Recall scores (1.00 and 0.99), indicating exceptional capability in identifying all relevant instances of these categories.

The macro-averaged and weighted averaged metrics (both 0.98) confirming that the model's performance across categories are balanced and no significant performance degradation occurs on minority categories. The uniformity of

performance through out the categories, indicates the non-homogeneity of the multi-task learning method and the homogeneous category distribution of the MAAD dataset.

**Detailed Summarization Evaluation with Extended Training**

In order to additionally verify the stability of the ArabicT5 model, we performed long training with varying models sizes and epochs settings. The results of the summarization performance metrics (ROUGE-1, ROUGE-2, ROUGE-L and BLEU) for both ArabicT5-small and ArabicT5-base are shown in Tab. 12.

**Table 12:** Similar experimental results of a summarization approach on data

task	Model	Vocab	Epochs	Rouge1	Rouge2	ROUGESUM	BLEU	Max length input	Max length target
summarization of Article	ArabicT5-small	32K	1	0.88	0.88	0.88	0.80	512	64
	ArabicT5-base	32K	3	0.90	0.90	0.90	0.81	512	64

This table presents the detailed results of extended training experiments on the full MAAD training set (100,000 articles). The highest ROUGE scores represent the optimal model performance achieved through comprehensive hyperparameter tuning.

### Training Dynamics and Generalization

Figure 10 illustrates the training and validation loss curves for the ArabicT5-base model across

three epochs. The results illustrate that the training loss gradually increases, while the validation loss decreases slightly and then stabilizes. This pattern shows that important patterns are still being learned by the model without substantial overfitting. Improved generalization ability and steady convergence during the training phase are demonstrated by the declining validation loss.



Fig. 10: Training Loss and Validation Loss

The stability of the validation loss after epoch 2 suggests that the model has reached an optimal point where further training provides diminishing returns. These results confirm that training over three epochs is the optimal approach for this specific task, effectively balancing computational resource usage with high model performance. As illustrated in Table 13, a qualitative analysis of the summaries generated by each model from the

same Arabic news source reveals distinct variations in fluency, information retention, and summarization depth. When compared to the human-written baseline, ArabicT5 stands out for producing the most coherent and balanced summary. Meanwhile, AraBART excels in narrative fluidity, mT5 offers a smooth interpretation with moderate factual accuracy, and GPT produces concise, yet comprehensible, outputs.

Table 13. Comparison models summary

Model	Summary
Original text	مع دخول اليوم الـ734 من حرب الإبادة على الشعب الفلسطيني، أعلن الرئيس الأميركي والوسطاء التوصل إلى اتفاق يمهد لإنهاء الحرب في غزة. وأكدت حركة حماس أن الاتفاق يقضي بإنهاء الحرب على غزة، وانسحاب الاحتلال منها، ودخول المساعدات، وتبادل الأسرى، مؤكدة أنها سلمت قوائم الأسرى الفلسطينيين الذين سيفرج عنهم بموجب الاتفاق. من جهته، قال المتحدث باسم الخارجية القطرية إنه تم الاتفاق في مفاوضات شرم الشيخ على بنود تنفيذ المرحلة الأولى لاتفاق وقف إطلاق النار بغزة، مؤكداً أن الاتفاق سيؤدي إلى وقف الحرب وإطلاق المحتجزين الإسرائيليين والأسرى الفلسطينيين. وقالت القناة 12 الإسرائيلية إن وقف إطلاق النار يدخل حيز التنفيذ الساعة 12 ظهر اليوم بالتوقيت المحلي (9.00 صباحاً بتوقيت غرينيتش) بعد توقيع الاتفاق. وأعلن الجيش الإسرائيلي أن قواته بدأت استعداداتها لتنفيذ الاتفاق بناء على توجيهات المستوى السياسي وتقييم الموقف، مضيفة أنها تواصل الانتشار في الميدان والاستعداد لأي تطورات عملياتية محتملة. ميدانياً، نفذت طائرات الاحتلال فجر اليوم قصفاً جويًا ومدفعيًا على مدينتي غزة وخان يونس، رغم الإعلان عن الاتفاق، كما قصف الاحتلال بالمدفعية محيط شارع الصناعة، جنوب غرب حي الصبرة بمدينة غزة، فيما أفادت وسائل إعلام إسرائيلية بإصابة جنود في عملية للمقاومة. وفي الضفة الغربية المحتلة، تجددت الاعتقالات والاعتقالات فجر اليوم، بعد ساعات من استشهاد فلسطيني وإصابة آخرين برصاص المستوطنين وقوات الاحتلال.
ArabicT5	أعلنت الولايات المتحدة والوسطاء عن اتفاق ينهي الحرب على غزة، يتضمن وقف إطلاق النار، وانسحاب قوات الاحتلال، وتبادل الأسرى، والسماح بدخول المساعدات. حركة حماس أكدت التزامها بالاتفاق وتسليم قوائم الأسرى، بينما أشار المتحدث باسم الخارجية القطرية إلى أن المرحلة الأولى ستبدأ

	اليوم الساعة 12 ظهرًا. ورغم الاتفاق، واصل الاحتلال قصفه لمدينتي غزة وخان يونس، فيما شهدت الضفة الغربية اقتحامات واعتقالات جديدة.
AraBART	تم التوصل في مفاوضات شرم الشيخ إلى اتفاق شامل لإنهاء الحرب على غزة وبدء تنفيذ وقف إطلاق النار ظهر اليوم. الاتفاق يشمل انسحاب الاحتلال وتبادل الأسرى وفتح المعابر لدخول المساعدات. ورغم الإعلان، نفذت إسرائيل قصفًا جديدًا على غزة وخان يونس، بينما استمرت الاعتقالات في الضفة الغربية، ما يثير تساؤلات حول التزام الاحتلال بالاتفاق.
mT5	أعلنت الأطراف الوسيطة التوصل لاتفاق لوقف الحرب في غزة يشمل تبادل الأسرى ودخول المساعدات. من المقرر أن يبدأ وقف إطلاق النار عند الساعة 12 ظهرًا. ومع ذلك، استمرت الغارات الإسرائيلية على غزة وخان يونس، ووقعت مواجهات واعتقالات في الضفة الغربية.
GPT	توصلت الأطراف بوساطة دولية إلى اتفاق لوقف الحرب في غزة يبدأ ظهر اليوم، يشمل وقف النار، تبادل الأسرى، وفتح المساعدات. رغم ذلك، واصلت إسرائيل قصفها على غزة وخان يونس وعملياتها في الضفة الغربية.

### Synthesis of Results

The ArabicT5 model's efficacy for Arabic news processing is demonstrated by the thorough examination across both summarization and classification tasks. The key findings are as follows:

**1 Summarization Performance:** ArabicT5 achieved superior ROUGE scores (ROUGE-1 = 0.90, ROUGE-2 = 0.90, ROUGE-L = 0.90) and BLEU score of 0.81, indicating strong capability in generating coherent, semantically accurate summaries that preserve both content and structure.

**2 Classification Performance:** ArabicT5 achieved 98% accuracy across ten news categories, with consistent performance across all categories (Precision, Recall, and F1-Score all  $\geq 0.95$ ), demonstrating robust multi-class classification capability.

**3 Model Comparison:** Among the four models evaluated, ArabicT5 consistently outperformed AraBART, mT5, and GPT across all metrics, confirming the advantage of Arabic-specific pre-training for Arabic NLP tasks.

**4 Training Efficiency:** The model exhibited stable convergence with appropriate hyperparameters (10 epochs for comparative evaluation, 3 epochs for extended training), achieving high performance without significant overfitting.

The results underscore the significant potential of transfer learning and text-to-text architectures in advancing Arabic natural language processing, particularly within multi-task environments that require simultaneous categorization and summarization. The superior performance exhibited by ArabicT5 not only validates the scientific utility of the MAAD dataset but also establishes it as a vital, comprehensive resource for future research in Arabic news processing, effectively addressing a critical gap in the field.

### 5. Conclusion

Through the development of the Multi-Label Arabic Articles Dataset (MAAD) and a comparative analysis of advanced Transformer models for simultaneous news summarization and classification, this research establishes a robust framework for advancing Arabic Natural Language Processing. By contributing a substantial, high-quality dataset for general domains and demonstrating the empirical advantages of language-specific models in handling the linguistic complexities of Arabic, this study effectively addresses a significant deficiency in existing Arabic NLP literature.

#### 5.1. Key Contributions

This study makes three main contributions: This study presents three primary contributions to the field. First, we introduced the MAAD Dataset, a rigorously verified collection comprising 602,792 news articles sourced from six prominent Arabic media outlets across ten distinct categories. To ensure high quality, the data underwent a comprehensive preprocessing phase that combined automated techniques such as keyword matching and LDA topic modeling with manual expert verification, resulting in a classification accuracy exceeding 95%. A key distinction of this dataset is its design for multi-task learning applications, including classification, summarization, and generation, unlike many existing Arabic datasets restricted to single-task optimization. Second, we established a robust comparative evaluation framework to assess four advanced Transformer models: ArabicT5, AraBART, mT5, and GPT. By utilizing a unified experimental environment with standardized hyperparameters and metrics, this framework highlights the performance differences between multilingual architectures and those pre-trained specifically on Arabic data. Finally, empirical evidence from our experiments demonstrates that

the ArabicT5 model consistently outperforms its counterparts across the evaluated tasks.

Summarization: ROUGE-1 = 0.90, ROUGE-2 = 0.90, ROUGE-L = 0.90, and BLEU = 0.81. These results were further validated by human evaluation, where the model achieved 4.86 in Fluency and 4.35 in Adequacy.

- Classification: Overall accuracy of 0.98 (98%) with consistent performance across all ten categories (Precision, Recall, and F1-Score all  $\geq$  0.95)

### 5.2. Implications for Arabic NLP

The findings of this study have significant implications for the field of Arabic Natural Language Processing:

1. The consistent effectiveness of ArabicT5 across all evaluation metrics underscores the critical value of language-specific pre-training within Arabic NLP. By specifically addressing the unique morphological intricacies, syntactic frameworks, and semantic subtleties of the Arabic language, this model demonstrates superior performance compared to both multilingual alternatives like mT5 and general-purpose large language models such as GPT within this context.
2. This finding reinforces the promise that domain and language-specific pre-training is key to achieving the best performance on specialized terminology/ domain-specific language.
3. The unique application of an individual Transformer model for both summarization and classification exemplify the feasibility and utility of multi-task learning methods for Arabic NLP. While there is a direct computation efficiency by this approach, indirect efficiency is also available here which allows transferring knowledge from simple tasks to some extent, and may help in enhancing generalization and robustness.
4. The excellent results obtained on MAAD together with its large size and more diverse category representation further confirms the quality and applicability of the dataset emerging new higher-end DL models capable of both detection and accuracy. Due to the detailed and well-encompassed nature of the dataset, we believe it will serve as a benchmark for future research in the field of Arabic NLP.

### Comparison with Existing Work

This study's findings demonstrate a substantial advancement relative to prior work in Arabic Natural Language Processing. Specifically, when evaluated on the MAAD dataset, the ArabicT5 model achieved a classification accuracy of 98%, surpassing the 87.75% benchmark previously

established by the Ultimate Arabic News Dataset [15]. Furthermore, the efficacy of the proposed methodology is reinforced by the high ROUGE scores achieved 0.90 for both ROUGE-1 and ROUGE-2 which indicate a significant improvement over earlier summarization techniques.

### Limitations and Future Directions

Although this study significantly advances Arabic NLP, several though this study significantly advances Arabic natural language processing, there are a number of limitations and future research opportunities that should be discussed:

#### Limitations:

1. Although representative, the sample of 110 articles for summarization and 14,990 articles for categorization used in this comparative analysis cannot fully capture the range and diversity of Arabic news across all genres and styles.
2. The work only addresses Modern Standard Arabic (MSA) and ignores dialectal issues, which continue to be a significant obstacle to Arabic NLP.
3. Since the evaluation metrics (ROUGE and BLEU) are mostly focused on recall, they might not adequately account for all facets of summarization quality, including readability and factual consistency.
4. Future study intends to extend this to a larger-scale assessment using inter-annotator agreement analysis to better examine the model's generative nuances, even though a human evaluation was carried out on a sample of 50 articles with encouraging results.

#### Future Research Directions:

1. To improve the model's applicability to actual Arabic text processing contexts, future work should expand the MAAD dataset and model evaluation to include dialectal Arabic variations (Gulf, Levantine, and North African).
2. Investigating the potential for cross-lingual transfer learning between Arabic and other morphologically rich languages could yield insights into language-universal principles in NLP.
3. Beyond ROUGE and BLEU, a more thorough evaluation of summarization quality would be possible by combining human review with more complex artificial measures (such as BERTScore or METEOR).
4. Developing domain-specific variants of the ArabicT5 model for specialized domains (medical, legal, scientific) could further enhance performance in these critical application areas.
5. To demonstrate the practical viability of this research, future work should focus on embedding

the developed models within operational environments such as news aggregators, recommendation systems, and automated journalism platforms. This step is crucial not only for validating the utility of the models in real-world scenarios but also for uncovering specific challenges related to their deployment and integration.

### Final Remarks

This study shows that transfer learning and text-to-text Transformer models are very useful for improving Arabic NLP when they are appropriately modified and trained on high-quality Arabic data. Future Arabic text processing research will have a strong basis thanks to the MAAD dataset and the thorough evaluation mechanism provided in this work. Large-scale, high-quality corpora and the creation of language-specific models are still essential to guaranteeing accuracy, dependability, and scalability in both research and real-world applications as Arabic NLP develops.

Despite the complexity of the linguistic challenges in the language (such as rich morphology, diacritical variations, and diverse syntactic structures), ArabicT5 achieves impressive performance across multiple Arabic NLP tasks, confirming that language-specific approaches are still required for optimal performance in Arabic NLP tasks. In order to attain state-of-the-art results in Arabic language processing and create sophisticated Arabic NLP applications that better serve Arabic-speaking communities globally, future work can expand on current work and push its bounds.

### References

- [1] T. Uçkan and A. Karci, "Extractive multi-document text summarization based on graph independent sets," *Egyptian Informatics Journal* 21(3), 145–157 (2020).
- [2] M. Azam, et al., "Current trends and advances in extractive text summarization: A comprehensive review," *IEEE Access* (2025).
- [3] A. Farghaly and K. Shaalan, "Arabic natural language processing: Challenges and solutions," *ACM Transactions on Asian Language Information Processing (TALIP)* 8(4), 1–22 (2009).
- [4] S. Faizullah, et al., "A survey of OCR in Arabic language: applications, techniques, and challenges," *Applied Sciences* 13(7), 4584 (2023).
- [5] Gummadi, V. P. K. (2022). *MuleSoft API Manager: Comprehensive lifecycle management*. *Journal of Information Systems Engineering and Management*, 7(4), 1–9.
- [6] A. G. Al-Khulaidi and S. M. Yaseen, "Comparative analysis and evaluation of stemming and preprocessing techniques for Arabic text," *العلوم صندعاء جامعة مجلة*, 1(4) وال 1(4) (2023).
- [7] A. Elnagar, R. Al-Debsi, and O. Einea, "Arabic text classification using deep learning models," *Information Processing & Management* 57(1), 102121 (2020).
- [8] Aljazeera, BBC, Arabic.rt, Youm7, Alummahnews, 26sep: Available: <https://www.aljazeera.net/>, <https://www.bbc.com/arabic>, <https://arabic.rt.com/>, <https://www.youm7.com/>, <https://alummahnews.com>, <https://www.26sep.net/> (accessed Dec. 26, 2025).
- [9] A. M. Gaber, M. N. El-Din, and H. Moussa, "SMAD: Text classification of Arabic social media dataset for news sources," *SMAD* 12(10) (2021).
- [10] D. Bansal, N. Saini, and S. Saha, "DCBRTS: a classification-summarization approach for evolving tweet streams in multiobjective optimization framework," *IEEE Access* 9, 148325–148338 (2021).
- [11] A. Qaroush, et al., "An efficient single document Arabic text summarization using a combination of statistical and semantic features," *Journal of King Saud University-Computer and Information Sciences* 33(6), 677–692 (2021).
- [12] Y. M. Wazery, et al., "Abstractive Arabic text summarization based on deep learning," *Computational Intelligence and Neuroscience* (2022), Article ID 1566890.
- [13] H. Gawbah, et al., "Arabic News Classification and Generation Based on an Encoder-Decoder Transformer Model (ArabicT5)," in *Proceedings of the 1st International Conference on Emerging Technologies for Dependable Internet of Things (ICETI)*, IEEE, 2024.
- [14] A. A. Mohsen, M. Y. Al-Nahari, and A. Alsubari, "Classification and Generation of Arabic News Titles from Raw Text Based on an Encoder-Decoder Transformer Model (mT5)," 2024.
- [15] A. M. Saad, et al., "Bangla news article dataset," in J. Smith and K. Brown (eds.), *Proceedings of the International Conference on Data Science and Information Engineering (ICDSIE 2024)*,

- LNCS, vol. 9999, pp. 110874–110874, Springer, Heidelberg (2024).
- [16] M. Altamimi and A. M. Alayba, “ANAD: Arabic news article dataset,” in J. Smith and K. Brown (eds.), *Proceedings of the International Conference on Arabic Computational Linguistics (ICACL 2023)*, LNCS, vol. 9999, pp. 109460–109460, Springer, Heidelberg (2023).
- [17] O. Einea, A. Elnagar, and R. Al Debsi, “Sanad: Single-label Arabic news articles dataset for automatic text categorization,” in J. Smith and K. Brown (eds.), *Proceedings of the International Conference on Arabic Natural Language Processing (ICANLP 2019)*, LNCS, vol. 9999, pp. 104076–104076, Springer, Heidelberg (2019).
- [18] B. Alshemaimri, et al., “Summarizing Arabic Articles using Large Language Models,” *CS & IT Conference Proceedings* 14(10) (2024).
- [19] M. K. Eddine, et al., “Arabart: a pretrained Arabic sequence-to-sequence model for abstractive summarization,” arXiv preprint arXiv:2203.10945 (2022). Available: <https://huggingface.co/moussaKam/AraBAR T> (accessed Dec. 26, 2025).
- [20] L. Xue, et al., “mT5: A Massively Multilingual Pre-trained Text-to-Text Transformer,” Available: <https://pypi.org/project/langdetect/> (accessed May 05, 2024).
- [21] Z. Alyafeai, et al., “Taqyim: Evaluating Arabic NLP tasks using ChatGPT models,” arXiv preprint arXiv:2306.16322 (2023).
- [22] E. El Moatez Billah Nagoudi, A. Elmadany, and M. Abdul-Mageed, “AraT5: Text-to-text transformers for Arabic language generation,” arXiv preprint arXiv:2109.12068 (2021).
- [23] P. K. Mallick, S. Mishra, and G.-S. Chae, “Digital media news categorization using Bernoulli document model for web content convergence,” *Personal and Ubiquitous Computing* 27(3), 1087–1102 (2023).
- [24] D. J. Hemanth, D. Pelusi, and C. Vuppapapati (eds.), *Intelligent data communication technologies and Internet of things: Proceedings of ICICI 2021*, Springer Singapore (2022).
- [25] B. Jang, I. Kim, and J. W. Kim, “Word2vec convolutional neural networks for classification of news articles and tweets,” *PloS One* 14(8), e0220976 (2019).
- [26] A. Radford, et al., “Language models are unsupervised multitask learners,” *OpenAI blog* 1.8, 1–9 (2019).
- [27] M. Altamimi and A. M. Alayba, “ANAD: Arabic news article dataset,” *Data in Brief* 50, 109460 (2023).
- [28] M. A. R. Abdeen, et al., “A closer look at Arabic text classification,” *International Journal of Advanced Computer Science and Applications* 10(11) (2019).
- [29] T. Hasan, et al., “XL-sum: Large-scale multilingual abstractive summarization for 44 languages,” arXiv preprint arXiv:2106.13822 (2021).
- [30] Y. Al-Nahari, M. Marwah, A. Mohsen, N. Abdo Al-Humidi, and A. Alsubari, “MAAD: Multi-Label Arabic Articles Dataset,” *Mendeley Data*, V1, doi: 10.17632/hbfc9j8hj8.1 (2025).