



Fann or Flop: A Multigenre, Multiera Benchmark for Arabic Poetry Understanding in LLMs

Wafa Alghallabi^{1†} Ritesh Thawkar^{1†} Sara Ghaboura^{1†} Ketan More^{1†}
Omkar Thawakar^{1†} Hisham Cholakkal¹ Salman Khan^{1,2} Rao Muhammad Anwer^{1,3}

¹Mohamed bin Zayed University of AI, ²Australian National University, ³Aalto University

{wafa.alghallabi, sara.ghaboura, omkar.thawakar}@mbzuai.ac.ae

<https://mbzuai-oryx.github.io/FannOrFlop/>

Abstract

Arabic poetry is one of the richest and most culturally rooted forms of expression in the Arabic language, known for its layered meanings, stylistic diversity, and deep historical continuity. Although large language models (LLMs) have demonstrated strong performance across languages and tasks, their ability to understand Arabic poetry remains largely unexplored. In this work, we introduce *Fann or Flop*, the first benchmark designed to assess the comprehension of Arabic poetry by LLMs in 12 historical eras, covering 14 core poetic genres and a variety of metrical forms, from classical structures to contemporary free verse. The benchmark comprises a curated corpus of poems with explanations that assess semantic understanding, metaphor interpretation, prosodic awareness, and cultural context. We argue that poetic comprehension offers a strong indicator for testing how good the LLM understands classical Arabic through Arabic poetry. Unlike surface-level tasks, this domain demands deeper interpretive reasoning and cultural sensitivity. Our evaluation of state-of-the-art LLMs shows that most models struggle with poetic understanding despite strong results on standard Arabic benchmarks. We release *Fann or Flop*¹ along with the evaluation suite² as an open-source resource to enable rigorous evaluation and advancement for Arabic language models.

1 Introduction

Arabic is among the world’s most lexically rich languages, with a vocabulary exceeding 12.3 million words—far surpassing that of most modern languages (AlSuyuti, 15th Century; Andrews, 2024). A single word can convey multiple meanings, varied pronunciations, and diverse interpretations, reflecting the language’s profound semantic com-

¹<https://huggingface.co/datasets/omkarthawakar/FannOrFlop>

²<https://github.com/mbzuai-oryx/FannOrFlop>

[†]Equal contribution.



Figure 1: **Chronological Wheel of Arabic Poetic Eras.** This circular taxonomy visualizes the evolution of Arabic poetry across 12 major historical eras, from the Pre-Islamic and Transitional periods through the Abbasid, Andalusian, and Mamluk dynasties, up to the Modern era. The layout reflects both temporal flow and the rich cultural shifts that shaped poetic expression. Detailed taxonomy by genre, meter, and notable poets presented in Table 2.

plexity. Despite its official status in 27 countries—ranking third in global geopolitical presence (wikipedia, 2025)—only a fraction of this lexicon remains in common use today.

To unify communication across its many dialects, Modern Standard Arabic (MSA) emerged in the late 19th and early 20th centuries as a formal register (oussama, 2024). Today, it is the primary language of education, media, and governance in the Arab world. Although linguists distinguish Classical Arabic (CA) from MSA, native speakers generally view them as a unified formal variety (wikipediaArabic, 2025). Nevertheless, even the most comprehensive Arabic dictionaries—such as Lisan al-Arab (Manzur, 14th Century), Taj al-Lugha (al Jawhari, 10th Century), and

Feature	AQMAR	Tafsir	Ashaar	Ara Bench	Arabic SQuAD	ARCD	AraBERT Collection	CAMeL Corpus	Tash keela	PADIC	MADAR	Fann or Flop
Dialectal Variety	✗	✗	✗	✓	✗	✗	✓	✗	✗	✓	✓	✓
Poetic Device Annotation	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✓
Verses/Sentence-Level Annotation	✗	✓	✓	✗	✗	✗	✗	✗	✗	✗	✗	✓
Temporal/Historical Context	✗	✓	✓	✗	✗	✗	✗	✗	✗	✗	✗	✓
QA-Style Task Format	✗	✗	✗	✗	✓	✓	✓	✗	✗	✗	✗	✓
Open-Source	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

Table 1: **Comparison of key Arabic NLP datasets.** Existing Arabic NLP resources typically address isolated features such as dialectal coverage, QA formats, or classical text processing. In contrast, **Fann or Flop** uniquely integrates multiple underrepresented dimensions (i.e. dialectal diversity, poetic device annotation, verse-level granularity, temporal grounding, and a QA-style evaluation format) positioning it as the first comprehensive benchmark for Arabic poetry understanding. AQMAR (Mohit et al., 2012), Tafsir (Ahmed et al., 2022), Ashaar (Alyafeai et al., 2023), AraBench (Sajjad et al., 2020), Arabic-SQuAD (Mozannar et al., 2019), ARCD (Mozannar et al., 2019), AraBERT Collection (Antoun et al., 2020), CAMeL Corpus (Abdul-Mageed et al., 2020; Khalifa et al., 2018), Tashkeela (Zerrouki and Balla, 2017), PADIC (Meftouh et al., 2015), MADAR (Bouamor et al., 2018).

al-Mu‘jam al-Mu‘asir (alsharekh, 2019; Ar-Riyadh, 2025)—cover only a small portion of the historical corpus, revealing the inherent challenges of Arabic lexicography and sociolinguistic narrowing of usage. Within this broader linguistic context, Arabic poetry has served as a repository of cultural and intellectual expression from the older era to the modern time. Poetic forms such as long odes (qasida), lyrical love poems (ghazal), elegies (ritha’), strophic songs (muwashsha), and vernacular verse (zajal) are marked by distinct metrical, rhetorical, and performative characteristics. While contemporary poets explore free verse and modernist motifs, classical forms continue to exert a strong literary and cultural influence.

Recent advances in LLM, such as GPT (Chen et al., 2025), LLaMA (Touvron et al., 2023), AceGPT (Huang et al., 2023), Jais (Sengupta et al., 2023), and Falcon (Malartic et al., 2024), have demonstrated impressive multilingual capabilities, including Arabic. However, most Arabic natural language processing (NLP) benchmarks focus on tasks such as sentiment analysis, question answering, or recognition of named entities (Antoun et al., 2020; Abdul-Mageed et al., 2021; Obeid et al., 2020), typically in MSA or dialectal prose. These benchmarks often miss the linguistic depth and cultural nuances that are inherent in Arabic poetry. As LLMs are increasingly evaluated for their ability to handle complex linguistic phenomena, such as metaphor, figurative language, and stylistic nuance, their limitations become evident (Liu et al., 2022; Bisk et al., 2020). The FLUTE benchmark (Chakrabarty et al., 2022) and the FigLang 2024 workshop (FIGLANG202, 2024) have reaffirmed that non-literal language understanding remains a

significant challenge. This challenge is particularly acute in Arabic, where poetry is densely layered with intertextuality and cultural symbolism. Arabic poetry thus provides a uniquely demanding testbed for assessing deep linguistic in language models.

To address this gap, we introduce *Fann or Flop*, the first benchmark dedicated to evaluating LLMs’ understanding of Arabic poetry. Our benchmark comprises 6,984 poem-explanation pairs curated from 12 distinct historical poetic eras (see Figure 1), which can be broadly seen as spanning three major historical periods: pre-Islamic, classical, and contemporary. It covers 14 poetic genres and includes a range of metrical forms, as detailed in Table 1. Each sample is manually verified by native Arabic speakers with domain knowledge to ensure linguistic authenticity and interpretive accuracy. This rich and diverse collection makes Fann or Flop a reliable benchmark for evaluating deep cultural and literary reasoning in Arabic NLP. Figure 2 represents the examples from our proposed Fann or Flop dataset, showcasing the diversity of eras, genres, and poetic styles covered.

Our goal is to provide a diagnostic on how well your language model understands and interprets genuine Classical Arabic. Unlike general text, poetry requires sensitivity to rhetorical devices, metrical patterns, and sociohistorical context, making it a rigorous and culturally grounded indicator of language proficiency. Our benchmark serves as a clear indicator of whether a model has truly been exposed to and internalized high-quality Arabic content, offering a focused lens into its cultural and linguistic depth. We evaluate a range of open-source and commercial LLMs using Fann or Flop and find that, despite strong performance on con-

Era	Approx. Years	Genres (Theme)	Meter	Notable Poets
Pre-Islamic (Jahiliyyah)	Until 610 CE	Satire, Separation, Wisdom	At-Tawil, Al-Kamel, Al-Basit	Imru al-Qays, Antarah ibn Shaddad, Zuhayr ibn Abi Sulma
Transitional Poets (Mukhadramun)	Late 6th – Early 7th c.	Praise, Apology, Religious	Ar-Rojz, Ar-Ramel	Hassan ibn Thabit, Labid ibn Rabi'a, Al-Khansa
Islamic	610–661 CE	Religious, Wisdom, Patience	Al-Madid, Al-Kamel	Abu Sallama Al-Arhabi, Onayf Ibn Kitra
Umayyad	661–750 CE	Love, Satire, Political	At-Tawil, Al-Wafer, As-Sari'	Jarir, al-Farazdaq, al-Akhtal
Abbasid	750–1258 CE	Praise, Elegy, Wisdom	Al-Basit, Kamel, Al-Monsareh, Al-Moktadab	Abu Nuwas, al-Mutanabbi, al-Buhturi, Abu Tammam
Fatimid	909–1171 CE	Religious, Praise, Sadness	Ar-Rojz, Al-Mutakareb	Ibn Hayus, Abu al-Ala al-Ma'arri
Andalusian	756–1492 CE	Love, Longing, Wisdom	Mowachah, Al-Mowaliya, Al-Mohtath	Ibn Sahl Al-Andalusi, Ibn Zaydun, Ibn Khafaja
Ayyubid	1171–1250 CE	Religious, Praise, Elegy	Al-Kamel, Al-Khafif	Ibn al-Farid, Mohyiddine Bin Arabi
Mamluk	1250–1517 CE	Wisdom, Praise, Religious	Al-Wafer, Ar-Rojz	Bahaa'eddine Zuhair, Safiyueddine Alhilli
Between the Two Dynasties	1258–1517 CE	Religious, Wisdom, Reproach	Al-Mutadarek, Ar-Ramel	Bashar bn Burd
Ottoman	1517–1800 CE	Religious, Love, General	Al-Kamel, Al-khafif	Bnt Al-Shahna, Ibn Razka
Modern	19th c. – Present	Nationalism, Love, Social Justice	Free Meter	Ahmad Shawqi, Hafeth Ibrahim

Table 2: **Taxonomy of Arabic Poetic Eras with Genre and Meter Coverage.** This table provides a structured overview of 12 major eras in Arabic poetic history, detailing their approximate chronological spans, the most prominent poetic themes (genres) representative of each era, the dominant metrical patterns (Arabic *buhūr*) used in poetic composition, and notable poets who exemplify the literary character of their time. The genre column highlights recurring thematic concerns such as satire, elegy, love, nationalism, and religious devotion, while the meter column showcases the classical metrical forms like *At-Tawil*, *Al-Kamel*, and *Ar-Rojz*, along with innovations such as free verse in the modern period. This taxonomy reflects the dynamic interplay between form, content, and historical context in shaping Arabic poetic expression.

ventional Arabic tasks, most models struggle with the interpretive depth required by poetry. These findings highlight the need for culturally informed benchmarks that better reflect the depth and diversity of Arabic. We release Fann or Flop as an open-source resource to support the development and evaluation of Arabic-capable language models.

2 The Fann or Flop Dataset

2.1 Dataset Taxonomy

To capture the linguistic, historical, and thematic richness of Arabic poetry, we construct an expert-verified taxonomy that organizes poems across both form and era. As illustrated in Figure 1 and detailed

in Table 2, the taxonomy traces 12 distinct poetic eras, from the pre-Islamic period to modern times, encompassing 14 genres that capture the dominant styles, concerns, and historical contexts of each era. It illustrates how poetic expression evolved over the centuries.

This structured framework was carefully reviewed and validated by scholars specializing in Arabic language and literature to ensure both linguistic accuracy and contextual relevance. Their expertise helped align the taxonomy with established literary traditions while accommodating the nuances of classical and modern poetic forms. Beyond its utility for literary and philological analysis, the taxonomy serves as a robust foundation for

Pre-Islamic (الجاهلي / قبل الإسلام)		Umayyad (الأموي)		Abbasid (العباسي)	
الشاعر: امرؤ القيس بن خُجَر		الشاعر: الفرزدق		الشاعر: أبو نواس	
النوع: غزل	بحر: الطويل	النوع: هجاء	بحر: الطويل	النوع: الحكمة	بحر: الرمل
فقا نك من ذكرى حبيب ومَنزل فَوضِعَ فالْمُفَرَّاءَ لَمْ يَغْفُرْ رَسْمُهَا تَرَى بَعْدَ الزَّامِ في عَرْضِهَا بِسَقَطِ اللّوى بَيْنَ الدُّخُولِ فَحُومِلَ	بِمَا نَسَجْتَهَا مِنْ جُلُوبٍ وَشَمَالٍ وَقَبِيعَاتِهَا كَأَنَّهُ حَبْ قُلْفَلٍ	فيسري فأني أرض قومك إني وَأَنِّي عَلَى سَعْدٍ بِمَا هِيَ أَهْلُهُ عِظَامُ الْمُقَارِي يَأْمُنُ الْجَارُ فَيُجْعِلُهَا أَرَى خُفْيَةَ خَوْفَاءَ جَمًّا فُلُوفُهَا	أرى خُفْيَةَ خَوْفَاءَ جَمًّا فُلُوفُهَا وَحِزْرَ أَحَادِيثِ الْغَرِيبِ صَدُوقُهَا إِذَا مَا الثُّرَيَّا أَخْلَفَتْهَا يُرِوْقُهَا	كُلُّ نَاعٍ فَسَيُنْعَى كُلُّ مَذْخُورٍ سَيُنْفَى لَيْسَ غَيْرُ اللَّهِ يَبْقَى إِنَّ شَيْئًا قَدْ كُفِينَا	كُلُّ بَالِكٍ فَسَيُنْكِي كُلُّ مَذْكَورٍ سَيُنْكِي مَنْ عَلَا قَالَتْهُ أَعْلَى هُ لَهْ نَبْسِي وَنَبْشِي
هذه القصيدة رثاء مؤثر، تُعبر عن شوق وحنين عميقين للماضي الجميل، وخاصة للذكرى حبيب ومكان عزيزين على قلب الشاعر. سَأشرح كل بيت على حدة:		تحدث هذه القصيدة عن قبيلة طواعة، وتصف حالها وخصائصها، وتُعبّر عن استنكار الشاعر لسلوك بعض أفرادها. سَأشرح كل بيت على حدة:		تُعدّ هذه القصيدة وصفًا بديعًا لعين جميلة، وتُبرز مهارة الشاعر في استخدام التشبيهات والاستعارات. سَأشرح كل بيت على حدة:	
البيت الأول: يبدأ الشاعر بدعوة للبقاء على ذكرى حبيب ومَنزل كان يقيم فيه. "سَقَطَ اللّوى" هو مكانٌ مُحدّد، يُرجّح أن يكون اسم مكان أو منطقة جبلية، "بين الدخول فحومل" توضيح لموقعه بين مرتفعين أو مكانين متقابلين.		البيت الأول: فيسري فأني أرض قومك إني: نداءٌ للشاعرة (أو المتحدثة) تدعو فيه إلى زيارة أرض قبيلة طواعة، وتصفها بأنها كالنمّ بالنسبة لها، مُعبرة عن علاقة قوية وارتباط عاطفي. أرى خُفْيَةَ خَوْفَاءَ جَمًّا فُلُوفُهَا: تصف الشاعرة منطقة من الأرض (خُفْيَةَ) بأنها واسعة وخصبة (خَوْفَاءَ جَمًّا)، وأن ثمارها وفيرة (فُلُوفُهَا).		البيت الأول: كل من ينعي ميتًا (ناعٍ) سوف يُنعى هو أيضًا يوفًا ما، وكل من ينكي ميتًا (بالك) سوف يُنكى هو أيضًا. هذا البيت يُشير إلى حتمية الموت الذي يطال الجميع بلا استثناء.	
البيت الثاني: يصف الشاعر المنزل بقوله "فَوضِعَ"، أي أنه واضح المعالم. "فالْمُفَرَّاءَ" تشير إلى الطريق المؤدية إلى المنزل، "لم يعف رسمها" أي أن آثار المنزل لا تزال باقية، "لما نسجتها من جنوب وشمال" يشير إلى اتساع المنزل وبنائه المتين الذي استمر رغم مرور الزمن.		البيت الثاني: وَأَنِّي عَلَى سَعْدٍ بِمَا هِيَ أَهْلُهُ: تثنى الشاعرة على طيبة أهل القبيلة وكرمهم، وتؤكد أنهم أهل سعادة (سَعْدٍ) بسبب أخلاقهم الحميدة. وَحِزْرَ أَحَادِيثِ الْغَرِيبِ صَدُوقُهَا: تُشير إلى أن أهل القبيلة يتميزون بالصدق والأمانة، وأن كلامهم صادق (صَدُوقُهَا)، وهذا ما يجعلهم محل ثقة لدى الغرابة.		البيت الثاني: كل ما يُخزن من مال أو متاع (مَذْخُور) سيؤول وينتهي، وكل من يُذكر ويشهر (مَذْكَور) سَيُنْسَى مع مرور الزمن. هذا البيت يُبرز زوال الممتلكات والشهرة.	
البيت الثالث: يصف الشاعر المنزل من جديد، "تَعَزَّ الزَّامُ" أي آثار البناء، "في عَرْضِهَا" أي في ساحاته الواسعة، "وقبعاها كأنه حب قفل" يشير إلى انتشار آثار البناء كحبوب الفلفل المنتشرة.		البيت الثالث: عِظَامُ الْمُقَارِي يَأْمُنُ الْجَارُ فَيُجْعِلُهَا: تشبيه يُبرز أمان القبيلة، حيث تُشبه بيوتها (المقاري)		البيت الثالث: لا يبقى إلا الله وحده، فمن اعتقد أنه علا أو تفوق، فإن الله تعالى هو الأعلى والأعظم منه. هذا البيت يؤكد على أزلية الله وبقائه، وعلى عظمته وتفوقه على كل شيء.	
البيت الرابع: "إن شئنا (أي: أمرًا ما من أمور الحياة) قد كلفنا به، ولبيدنا الأمر نسعى ونشقى في سبيل تحقيقه. هذا البيت يُشير إلى مشقة الحياة وسعي الإنسان لتحقيق أهدافه.		البيت الثالث: عِظَامُ الْمُقَارِي يَأْمُنُ الْجَارُ فَيُجْعِلُهَا: تشبيه يُبرز أمان القبيلة، حيث تُشبه بيوتها (المقاري)		البيت الرابع: إن شئنا (أي: أمرًا ما من أمور الحياة) قد كلفنا به، ولبيدنا الأمر نسعى ونشقى في سبيل تحقيقه. هذا البيت يُشير إلى مشقة الحياة وسعي الإنسان لتحقيق أهدافه.	

Andalusian (الأندلسي)		Islamic (الإسلامي)		Transitional (المخضرمون)	
الشاعر: ابن خفاجة		الشاعر: عيسى بن فاتك		الشاعر: لبيد بن ربيعة	
النوع: صبر	بحر: الطويل	النوع: سياسية	بحر: الوافر	النوع: فراق	بحر: الرجز
كفاني شكوى أن أرى المجد شاكياً وَحَسْبُ الزَّيْزَانِ أَنْ تَرَاني بَاكِتاً أَدَارِي فَوَادٍ يَصُدُّعُ الصُّنْدُوقَ زَفَرَةً وَزَجَّعَ أَنْيَبَ يَحْلُبُ السَّمْعَ سَاجِياً وَكَيْفَ أَوَارِي مِنْ أَوَارٍ وَجُدْتِي لَهُ صَادِرًا عَنْ مَثَلِ الْمَاءِ صَادِتاً		فَلَمَّا أَصْبَحُوا صَبَّأُوا وَقَامُوا إِلَى الْجُرَدِ الْعَتَاقِ شَسُومِينَا فَلَمَّا اسْتَجْمَعُوا حَمَلُوا عَلَيْنِهِمْ فَطَلَّ ذُو الْجَعَالِ يَفْتَلُونَا بَقِيَّةَ يَوْمِهِمْ حَتَّى آتَاهُمْ سَوَادُ اللَّيْلِ فِيهِ يَرَاوَعُونَا		يَا هَرِمًا وَأَنْتَ أَهْلُ عَدَلٍ إِنْ وَرَدَ الْأَخْوَصُ مَاءً قَبْلِي لَيْذَهَبَ أَهْلُهُ بِأَهْلِي لَا تَجْعَلَنَّ شَكْلِيْمٌ وَشَكْلِي وَنَسَلُ آبَائِهِمْ وَنَسَلِي لَقَدْ تَهَيَّئْتُ عَنْ سَفَاهِ الْجَهْلِ	
تحدث القصيدة عن شاعر يعاني من ضائقة نفسية عميقة، يتجلى ذلك من خلال صور شعرية مؤثرة تعبر عن حزنه والألم. سنشرح كل بيت على حدة:		تصوّر القصيدة معركةً بين فئة قليلة مؤمنة وجيش كبير، حيث تنتصر الفئة بثباتها وإيمانها، وتدين طاعة الظالمين وتعلي من شأن العدالة الإلهية. سَأشرح كل بيت على حدة:		تُعبّر هذه القصيدة عن رفض قاطع للظلم والاحتجاج، وتُظهر شجاعة الشاعر في مواجهة الجور، وتُفصّل رحيله عن بيئة فاسدة. سَأفصّل شرح كل بيت على حدة:	
البيت الأول: يشعر الشاعر بأن رؤية المجد نفسه حزناً كافية لتكون شكواه، ولا يحتاج لأكثر من ذلك. حزنه الشخصي يضافي حزن المجد ذاته، بل يكفيه أن يراه الناس يبكي. "المجد" هنا رمز لما هو عظيم أو مرغوب، وشكواه تدل على ضياع معالمة أو انحصار.		البيت الأول: يصف هذا البيت بداية يوم المعركة، حيث قام القوم بأداء الصلاة مع شروق الشمس، ثم توجهوا إلى ساحة القتال وهم مستعدون ومنظمون. وتشير عبارة "الجرّد العتاق" إلى الخيول النشيطة الأصبغة، مما يرمز إلى الجهوية والعزيمة.		البيت الأول: يا هَرِمًا: نداءٌ موجّه إلى القاضي أو الحاكم (هرم: قاضي أو حاكم يُقصّد به هنا رمزاً للعدل). وَأَنْتَ أَهْلُ عَدَلٍ: تأكيدٌ على أنّ الحاكم يُفترض به أن يكون عادلاً. إِنْ وَرَدَ الْأَخْوَصُ مَاءً قَبْلِي: تشبيه يُعبّر عن الظلم. "الأخوص" هو الشخص الذي يسبق غيره في الحصول على الماء أو الحقوق، بمعنى أن هناك من سبق الشاعر في الحصول على حقه. البيت يُشير إلى وجود ظلم وقع على الشاعر، حيث تمّ إعطاء حقوقه لمن هو أقل استحقاقاً.	
البيت الثاني: يُخفي الشاعر قلبه المتألم الذي يتسبب في آتني حزين يُخرج الدموع بغزارة. "يَصُدُّعُ الصُّنْدُوقَ زَفَرَةً" صورة بليغة تُظهر شدة الألم الذي يُمرّق صدره، و"يَحْلُبُ الدَّمْعَ سَاجِياً" تشبيهٌ دقيق لسيلان الدموع الغزير.		البيت الثاني: يوضح هذا البيت هجوم الجيش الكبير على المجموعة القليلة، حيث كان الهجوم مفاجئاً وعنيفاً. و"ذُو الْجَعَالِ" هم أصحاب السيف، مما يدل على شدة القتال وكثرة القتلى في صفوف المؤمنين.		البيت الثاني: لَيْذَهَبَ أَهْلُهُ بِأَهْلِي: دُعاءٌ على الظالم بأن يذهب هو ومن معه مع أهله الشاعر، بمعنى أن الشاعر لا يرغب في الارتباط أو المشاركة معهم في شيء. لَا تَجْعَلَنَّ شَكْلِيْمٌ وَشَكْلِي: طلّت واضحٌ من الحاكم بعدم الجمع بينه وبين الظالمين، وعدم ربط مصيره بمصيرهم، فهو يرفض تماماً التشابه أو الارتباط بهم....	
البيت الثالث: يتساءل الشاعر كيف يُخفي حزنه وهو غارق فيه، كمن يخر من ينبوع الماء، لا يستطيع إخفاء البلل. الصورة تُبرز استحالة إخفاء الألم المتغلغل في كيانه.		البيت الثالث: استمر القتال طوال النهار حتى دخل الليل، حيث استغل الطرف المهاجم ظلمة الليل للفرار والتمويه والهرب....			

Figure 2: **Representative Poetic Samples Across Arabic Literary Eras.** This figure presents curated excerpts from Arabic poems spanning key historical eras, illustrating the evolution of language, themes, and stylistic expression. The Pre-Islamic sample reflects tribal valor and rhetorical precision; the Umayyad excerpt captures satire and social commentary; the Abbasid example highlights philosophical reflection and refined metaphorical use; the Transitional era showcases a poetic voice confronting injustice and advocating moral clarity; and the Andalusian selection reveals emotional openness and psychological depth through lyrical expression. Together, these samples provide insight into how Arabic poetry has adapted to diverse historical, cultural, and ideological contexts. Refer to Appendix F.1, Figure 9 for the GPT-4o-generated English translations of the Arabic poetic samples.

computational modeling. It enables more precise automatic genre classification and facilitates temporal contextualization across different eras of Arabic poetry, thereby supporting culturally informed and interpretable Arabic NLP research.

2.2 Data Collection

We curated Arabic poems from a well-established digital archive³, which hosts a broad spectrum of poets, genres, and historical periods. A custom web scraper was developed to extract the poem texts along with associated metadata, including

³<https://arabic-poetry.net>

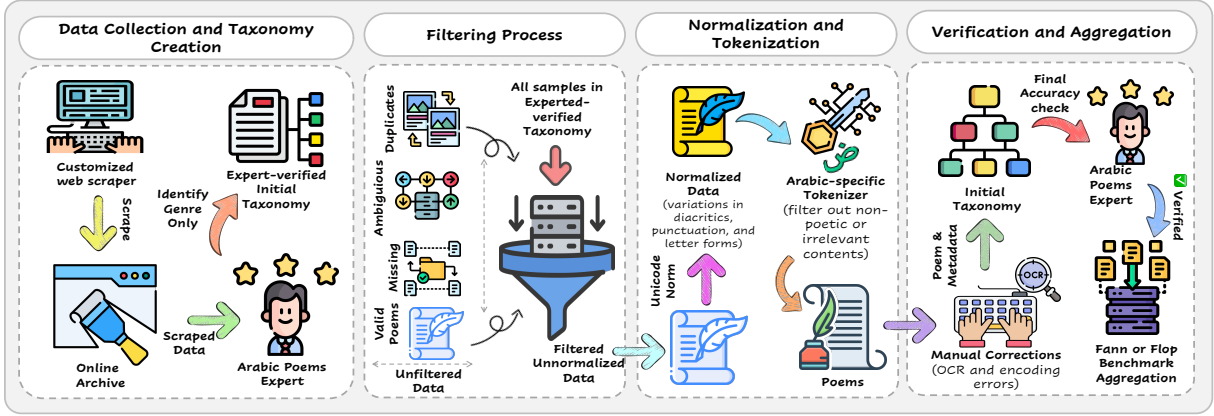


Figure 3: **Fann or Flop Pipeline.** Fann or Flop is built out of the multi-stage pipeline. It begins with scraping Arabic poems from a trusted online archive using a custom web scraper. Extracted poems are matched to an initial expert-verified taxonomy and filtered to remove duplicates, ambiguous metadata, and invalid entries. The filtered texts then undergo normalization (e.g., unifying diacritics, punctuation, and letter forms) and Arabic-specific tokenization, with non-poetic or irrelevant content excluded. Manual corrections are applied to fix OCR and encoding errors. In the final stage, linguistic experts verify each sample to ensure proper alignment with genre and era labels.

the poet’s name, historical era, genre, meter, and poem’s title. The resulting dataset extends across 12 distinct eras, from the pre-Islamic period to the modern era, and reflects a diverse range of poetic forms and styles. All entries were curated following our expert-verified taxonomy (see Table 2), ensuring consistency across genre and era classifications. This structured approach preserves both the linguistic richness and historical specificity of Arabic poetry, providing a valuable resource for research in both Arabic NLP and digital literary studies. By aligning each sample with a well-defined literary context, the dataset enables more accurate model evaluation and supports culturally grounded language understanding.

2.3 Data Filtering and Verification

To ensure data quality, consistency, and alignment with the expert-defined taxonomy, we applied a multi-step filtering and verification pipeline, illustrated in Figure 3. The process consisted of the following stages:

- **Duplicate and Metadata Filtering:** Starting with a collection of over 10,000 Arabic poems, we removed duplicate entries and discarded those with missing or ambiguous metadata, such as unknown poets or unspecified historical eras, resulting in a curated dataset of 6,984 high-quality poems.
- **Unicode Normalization:** All poems were standardized using Unicode normalization to

address orthographic inconsistencies common in Arabic, including variations in diacritics, punctuation, and letter forms (e.g., alternate representations of *alif* and *ta marbuta*).

- **Text Tokenization and Content Filtering:** We applied an Arabic-specific tokenizer to segment the text accurately. Non-poetic or irrelevant content, such as editorial comments, footnotes, and prose fragments, was automatically excluded.
- **Manual Correction of Encoding Errors:** A sample subset of poems was manually reviewed to correct common OCR and encoding issues that were not resolved through automated preprocessing.
- **Expert Validation of Labels:** All genre and era annotations were reviewed by Arabic language and literature experts. This validation step ensured that each poem was accurately categorized in accordance with the taxonomy introduced in Section 2.1.

3 Fann or Flop Benchmark Evaluation

Evaluation Framework. To assess the quality of LLM-generated explanations for Arabic poetry, we adopt a multi-layered evaluation pipeline. This framework integrates (i) automatic lexical metrics, (ii) semantic and entailment modeling, (iii) human expert annotation (interpretive depth), and (iv)

	Model	BLEU	chrF(++)	BERTScore	Textual Entailment	Faithfulness/ Consistency	Fluency/ Grammaticality	Interpretive- Depth
Closed	GPT-4o-2024-08-06 (OpenAI, 2024)	0.0395	0.2882	0.6410	0.6775	3.92 (± 0.99)	4.96 (± 0.20)	7.52
	GPT-4o-mini-2024-07-18 (OpenAI, 2024)	0.0395	0.2542	0.6124	0.4383	2.91 (± 0.75)	4.28 (± 0.57)	7.50
	Gemini-2.5-Flash (AI, 2025b)	0.0153	0.2618	0.6319	0.7475	4.25 (± 1.00)	4.98 (± 0.16)	7.22
	Gemini-2.0-Flash (AI, 2025a)	0.0395	0.2618	0.6393	0.7154	3.99 (± 1.04)	4.95 (± 0.22)	6.50
	Gemini-1.5-Pro (Reid et al., 2024)	0.0395	0.2618	0.6333	0.6180	3.59 (± 1.00)	4.80 (± 0.41)	5.38
	Fanar-Star (Team et al., 2025)	0.0138	0.1538	0.5677	0.6468	2.16 (± 0.92)	3.40 (± 0.76)	2.88
Open	Deepseek-V3 (Liu et al., 2024)	0.0395	0.2771	0.6335	0.5117	3.36 (± 0.91)	4.98 (± 0.16)	4.75
	Deepseek-R1 (Guo et al., 2025)	0.0395	0.2771	0.6335	0.5117	3.38 (± 0.92)	4.98 (± 0.16)	4.25
	Llama-3.3-70B (Meta AI, 2024)	0.0153	0.2618	0.6393	0.5364	2.51 (± 0.90)	3.37 (± 0.73)	7.20
	Qwen-3 (Team, 2025)	0.0296	0.2837	0.6158	0.6468	3.98 (± 0.90)	4.73 (± 0.45)	6.50
	Aya-Expanse (Dang et al., 2024)	0.0329	0.2771	0.6328	0.6468	3.76 (± 0.90)	4.68 (± 0.47)	5.88
	Jais (Sengupta et al., 2023)	0.0312	0.2698	0.6245	0.6023	3.21 (± 0.88)	4.35 (± 0.52)	5.35
	ALLaM-7B (Bari et al., 2024)	0.0119	0.0463	0.5375	0.5997	1.32 (± 0.62)	2.11 (± 0.89)	3.12
	AceGPT-v2-70B-Chat (Huang et al., 2023)	0.0402	0.0412	0.5759	0.6061	2.52 (± 0.91)	3.46 (± 0.95)	4.12

Table 3: **Comparison of closed and open-source models on the Arabic poem understanding task using both automatic and human evaluations.** BLEU, chrF(++), and BERTScore capture lexical and semantic similarity with reference explanations, while textual entailment assesses factual alignment. Human evaluation includes interpretive depth, while faithfulness and fluency are automatically judged using GPT-4o as a reference grader. Closed models like GPT-4o and Gemini-2.5-Flash achieve strong overall performance, while open models such as Deepseek-V3 and Aya-Expanse show promising consistency and interpretability. This benchmark highlights the potential of open models and the need for deeper cultural reasoning in Arabic poetic understanding.

	Model	Pre-Islamic	Transitional	Early Islamic	Umayyad	Abbasid	Fatimid
Closed	GPT-4o-2024-08-06 (OpenAI, 2024)	0.6285	0.6304	0.6341	0.6285	0.6421	0.6398
	GPT-4o-mini-2024-07-18 (OpenAI, 2024)	0.5980	0.6060	0.6134	0.5998	0.6125	0.6127
	Gemini-2.5-Flash (AI, 2025b)	0.6245	0.6264	0.6286	0.6253	0.6326	0.6282
	Gemini-2.0-Flash (AI, 2025a)	0.6290	0.6303	0.6326	0.6312	0.6404	0.6373
	Gemini-1.5-Pro (Reid et al., 2024)	0.6255	0.6293	0.6223	0.6278	0.6338	0.6307
	Fanar-Star (Team et al., 2025)	0.5694	0.5749	0.5695	0.5696	0.5720	0.5666
Open	Deepseek-V3 (Liu et al., 2024)	0.6225	0.6303	0.6311	0.6263	0.6313	0.6330
	Deepseek-R1 (Guo et al., 2025)	0.6271	0.6296	0.6321	0.6247	0.6324	0.6359
	Llama-3.3-70B (Meta AI, 2024)	0.5705	0.5703	0.5701	0.5668	0.5831	0.5719
	Qwen-3 (Team, 2025)	0.6111	0.6152	0.6129	0.6136	0.6164	0.6145
	Aya-Expanse (Dang et al., 2024)	0.6214	0.6232	0.6220	0.6232	0.6343	0.6294
	Jais (Sengupta et al., 2023)	0.6172	0.6218	0.6241	0.6183	0.6285	0.6239
	ALLaM-7B (Bari et al., 2024)	0.5786	0.5826	0.5917	0.5790	0.5862	0.5799
	AceGPT-v2-70B-Chat (Huang et al., 2023)	0.6194	0.6246	0.6329	0.6213	0.6261	0.6225

	Model	Andalusian	Ayyubid	Mamluk	Between Dynasties	Ottoman	Modern
Closed	GPT-4o-2024-08-06 (OpenAI, 2024)	0.6386	0.6440	0.6563	0.6440	0.6510	0.6487
	GPT-4o-mini-2024-07-18 (OpenAI, 2024)	0.6151	0.6167	0.6273	0.6176	0.6202	0.6140
	Gemini-2.5-Flash (AI, 2025b)	0.6297	0.6340	0.6421	0.6336	0.6415	0.6341
	Gemini-2.0-Flash (AI, 2025a)	0.6346	0.6409	0.6533	0.6414	0.6504	0.6441
	Gemini-1.5-Pro (Reid et al., 2024)	0.6313	0.6349	0.6409	0.6355	0.6443	0.6387
	Fanar-Star (Team et al., 2025)	0.5746	0.5684	0.5569	0.5831	0.5586	0.5392
Open	Deepseek-V3 (Liu et al., 2024)	0.6337	0.6404	0.6482	0.6393	0.6404	0.6368
	Deepseek-R1 (Guo et al., 2025)	0.6353	0.6404	0.6509	0.6408	0.6423	0.6373
	Llama-3.3-70B (Meta AI, 2024)	0.5791	0.5755	0.5935	0.5854	0.5797	0.5794
	Qwen-3 (Team, 2025)	0.6153	0.6163	0.6189	0.6160	0.6242	0.6149
	Aya-Expanse (Dang et al., 2024)	0.6289	0.6366	0.6475	0.6367	0.6393	0.6398
	Jais-30B-v3 (Sengupta et al., 2023)	0.6279	0.6321	0.6413	0.6307	0.6348	0.6316
	ALLaM-7B (Bari et al., 2024)	0.5876	0.5925	0.6004	0.5884	0.5933	0.5864
	AceGPT-v2-70B-Chat (Huang et al., 2023)	0.6168	0.6280	0.6466	0.6212	0.6205	0.6265

Table 4: **Era-wise Evaluation using BERTScore.** Model-wise performance breakdown using BERTScore evaluation across different Arabic poetic eras, evaluating understanding and generation quality within historical and stylistic contexts. The eras span from Pre-Islamic to Modern periods, offering a fine-grained analysis of model capabilities across evolving linguistic and cultural expressions. This table highlights gaps in temporal generalization and cultural grounding, motivating the need for era-aware training and evaluation in Arabic literary modeling.

LLM-as-Judge scoring. Together, these complementary layers capture both surface-level fidelity and the deeper interpretive demands of poetic understanding.

3.1 Automatic and Semantic Metrics

For automatic evaluation, we compute BLEU (Papineni et al., 2002) and chrF(++) (Popović, 2017) scores to quantify semantic and character-level overlap between model outputs and actual poem explanation references. While useful for consistency checks, these metrics are limited in capturing the nuanced variation allowed in literary interpretation.

To assess semantic alignment, we employ BERTScore (Zhang et al., 2019), leveraging Arabic-pretrained transformers such as AraBERT (Antoun et al., 2020) to quantify the semantic similarity between model-generated explanations and human-authored references. In addition, we incorporate Textual Entailment (TE) analysis using mDeBERTaV3 (He et al., 2021), a multilingual model fine-tuned for Natural Language Inference (NLI). This enables us to evaluate whether the reference explanation logically entails the generated output. Unlike surface-level similarity metrics, this approach provides a deeper measure of semantic consistency, capturing whether the generated interpretation remains faithful to the intended meaning of the expert-authored reference, even when expressed using different lexical or syntactic forms.

3.2 Human Evaluation and Inter-Annotator Agreement

To capture interpretive and literary nuance beyond automated metrics, we performed human evaluation on model-generated explanations. Annotators used a rubric-based scale (0–10) grounded in poetic analysis, comprising the following criteria:

- **Literal Comprehension (0–1):** Does the explanation correctly reflect the surface meaning of the poem?
- **Thematic and Emotional Depth (0–2):** Does it convey underlying themes, sentiment, or tone (e.g., longing, satire, mysticism)?
- **Cultural and Historical Appropriateness (0–2):** Does it demonstrate awareness of cultural, religious, or historical context?
- **Stylistic Sensitivity (0–3):** Does it acknowledge rhetorical and literary features such as

metaphor, figurative language, rhythm, or imagery?

- **Expressiveness and Coherence (0–2):** Is the explanation clear, well-articulated, and stylistically appropriate in Arabic?

Annotators. Three native Arabic speakers with expertise in Arabic literature and linguistics voluntarily applied the evaluation rubric, reflecting a shared commitment to advancing Arabic linguistic research and strengthening LLM evaluation in culturally specific domains.

Rubric Guidelines. The evaluation rubric was collaboratively refined by the annotators to capture the cultural, literary, and linguistic dimensions of Arabic poetry understanding, ensuring a balanced and consistently applied framework (see Table 17 for a detailed breakdown).

Inter-Annotator Agreement (IAA). To evaluate the reliability of the interpretive depth scoring, 10% of the data was randomly sampled and Krippendorff’s α was computed across the three annotators, resulting in an IAA of 87.7%. This level of agreement underscores the consistency of the scoring process and validates the robustness of the interpretive depth metric.

3.3 LLM-as-Judge Evaluation

Based on human-annotated interpretive depth scores (Section 3.2), GPT-4o (OpenAI, 2024) demonstrated the strongest performance among all open- and closed-source models (Table 3), thereby motivating its adoption as the LLM-as-judge. In this role, it evaluates explanation quality across two key dimensions: *Faithfulness/Consistency*—whether the explanation accurately reflects the poem without introducing hallucinated or misleading content—and *Fluency/Grammaticality*—whether it is grammatically correct, clear, and natural in Arabic.

4 Results and Analysis:

Table 3 presents the performance of both closed and open-source models on Arabic poem understanding using a combination of automatic metrics (BLEU, chrF(++), BERTScore, Textual Entailment, faithfulness and fluency) and human evaluation such as interpretive depth analysis. These metrics collectively assess the quality, relevance, and clarity of model-generated explanations when interpreting Arabic poetry.

Poem	GPT-4o	Gemini 2.5 Flash
<p>أَبْعَدَ الْحَارِثَ الْمَلِكَ ابْنَ عَمْرٍو لَهُ مَلِكُ الْعِرَاقِ إِلَى عُثْمَانَ</p> <p>مُجَاوَرَةً بَنِي شَمْجَى بْنِ جَرْمٍ هُوَأَنَا مَا أُتِيحُ مِنَ الْهُوَانِ</p> <p>وَيَمْنَعُنْهَا بَنُو شَمْجَى بْنِ جَرْمٍ مَعِيزُهُمْ خَنَّاكَ ذَا الْخَنَانِ</p>	<p>لَهُ مَلِكُ الْعِرَاقِ إِلَى عُثْمَانَ أَبْعَدَ الْحَارِثَ الْمَلِكَ ابْنَ عَمْرٍو يبدأ الشاعر بذكر الحارث بن عمرو، الملك الذي يمتد ملكه من العراق إلى عمان، مما يعكس عظمة سلطته واتساع نفوذه. يستخدم الشاعر هنا أسلوب الفخر والتعظيم من خلال ذكر المسافات الجغرافية الشاسعة التي يسيطر عليها الملك، مما يضفي على الحارث هالة من القوة والهيبة. كما أن استخدام كلمة "بعد" يوحي باستمرارية تأثير الحارث حتى بعد رحيله، مما يعزز من مكانته التاريخية.</p> <p>هُوَأَنَا مَا أُتِيحُ مِنَ الْهُوَانِ مُجَاوَرَةً بَنِي شَمْجَى بْنِ جَرْمٍ يشير الشاعر إلى مجاورة بني شمجى بن جرم، ويصف هذه المجاورة بأنها هوان، مما يعكس التوتر والصراع بين القبائل. يستخدم الشاعر هنا الكناية عن الضعف والذل الذي قد يصيب القبيلة نتيجة هذه المجاورة، ويبرز من خلال ذلك التحديات التي تواجهها القبائل في الحفاظ على مكانتها وقوتها. كما أن استخدام كلمة "ما أتيح" يوحي بأن هذا الهوان ليس اختيارياً بل مفروضاً عليهم.</p> <p>مَعِيزُهُمْ خَنَّاكَ ذَا الْخَنَانِ وَيَمْنَعُنْهَا بَنُو شَمْجَى بْنِ جَرْمٍ يصف الشاعر كيف أن بني شمجى بن جرم يمنعون القبيلة من حقوقها، مستخدماً كلمة "معيزهم" كاستعارة للثروة أو الموارد التي يحتفظون بها لأنفسهم. كما أن استخدام "خَنَّانك ذا الخنان" يعكس التوسل والرجاء، مما يضفي على البيت طابعاً من الاستعطاف والطلب. هذا الأسلوب البلاغي يعبر عن التوتر بين القبائل والحاجة إلى تدخل أو مساعدة خارجية لحل النزاعات.</p>	<p>لَهُ مَلِكُ الْعِرَاقِ إِلَى عُثْمَانَ أَبْعَدَ الْحَارِثَ الْمَلِكَ ابْنَ عَمْرٍو يستهل الشاعر بأسلوب الاستفهام الإنكاري "أبعد"، مستكبراً ومستعظماً الحال الذي آلت إليه الأمور بعد وفاة الملك الحارث بن عمرو، ويذكر سعة ملكه وامتداد سلطانه من العراق إلى عمان، وهي كناية عن عظمته وقوته التي يضعها في مقابلة مع الوضع الراهن الذي سيفصح عنه لاحقاً، مؤكداً على التباين الشديد بين الماضي والحاضر.</p> <p>هُوَأَنَا مَا أُتِيحُ مِنَ الْهُوَانِ مُجَاوَرَةً بَنِي شَمْجَى بْنِ جَرْمٍ يصف الشاعر الحال بعد الملك العظيم، وهو مجاورة قبيلة بني شمجى بن جرم، ويصف هذه المجاورة بأنها "هواناً"، مؤكداً على شدة الذل والمهانة التي لحقت بهم، ويستخدم تعبير "ما أتيح من الهوان" للمبالغة في وصف حجم الهوان الذي وقع عليهم، وكأنه أقصى ما يمكن أن يقع من ذل، وهو بذلك يقرن بين العز السابق والذل اللاحق.</p> <p>مَعِيزُهُمْ خَنَّاكَ ذَا الْخَنَانِ وَيَمْنَعُنْهَا بَنُو شَمْجَى بْنِ جَرْمٍ يستمر في وصف الهوان بذكر مثال عملي، وهو أن بني شمجى بن جرم يمنعون شيئاً ما (ربما الإبل أو المرعى أو غير ذلك مما يعود على القوم بالخير) أو يمنعونهم من معيزهم هم، ويختتم البيت بعبارة "خَنَّانك ذا الخنان"، وهي صيغة استرحام أو تحسر وتجعج على الحال، وتدل على بلوغ الذل منتهاه حتى أصبحوا في موضع من يحتاج إلى الشفقة والرحمة من قبيلة كانت أقل شأنًا منهم في الماضي.</p>
Original Explanation	Qwen 3	Aya-expanse
<p>لَهُ مَلِكُ الْعِرَاقِ إِلَى عُثْمَانَ أَبْعَدَ الْحَارِثَ الْمَلِكَ ابْنَ عَمْرٍو أبعد: حرف استفهام يدل على التعجب والاستنكار. الحارث الملك ابن عمرو غفرو: يشير إلى حارث بن عمرو، وهو ملك ذو سلطان عظيم. لَهُ مَلِكُ الْعِرَاقِ إِلَى عُثْمَانَ: يحدد نطاق ملك حارث، ممتداً من العراق في الشمال إلى عمان في الجنوب، مما يدل على قوة ملكه وشموه. يُعِزُّ هذا البيت عن المفاجأة والدهشة من هزيمة ملك بهذا القدر من النفوذ.</p> <p>هُوَأَنَا مَا أُتِيحُ مِنَ الْهُوَانِ مُجَاوَرَةً بَنِي شَمْجَى بْنِ جَرْمٍ مُجَاوَرَةً: ظرف زمان، أي: بعد أن كان مجاوراً. بَنِي شَمْجَى بْنِ جَرْمٍ: قبيلة بني شمجى بن جرم، وهي القبيلة التي هزمت حارث بن عمرو. هُوَأَنَا مَا أُتِيحُ مِنَ الْهُوَانِ: هوأنا شديد، لم يُز مثله من قبل، أي هزيمة ساحقة ومذلّة. يبرز هذا البيت شدة الهزيمة التي لحقت بحارث.</p> <p>مَعِيزُهُمْ خَنَّاكَ ذَا الْخَنَانِ وَيَمْنَعُنْهَا بَنُو شَمْجَى بْنِ جَرْمٍ وَيَمْنَعُنْهَا: أي يمنعون حارث بن عمرو من استعادة ملكه. بَنُو شَمْجَى بْنِ جَرْمٍ: تكرر ذكر القبيلة الفائزة لتأكيد سيطرتها. مَعِيزُهُمْ خَنَّاكَ ذَا الْخَنَانِ: (معيزة: هنا بمعنى نصرهم أو انتصارهم). "خَنَّانك ذا الخنان" يستخدم هنا مجازاً، ربما يقصد به القوة والسلطان اللذان أظهرهما بني شمجى، أو ربما يتحدث الشاعر عن قوة حارث نفسه (التي أصبحت بيد خصومه الآن). البيت يبرز استمرار سيطرة بني شمجى على ما كان ملكاً لحارث. باختصار، القصيدة تصف هزيمة حارث بن عمرو، ملك عظيم، على يد بني شمجى بن جرم، وتبرز شدة الهزيمة واستمرار سيطرة القبيلة المنتصرة. وتظهر القصيدة استعراب الشاعر من هزيمة ملك بهذا الحجم.</p>	<p>لَهُ مَلِكُ الْعِرَاقِ إِلَى عُثْمَانَ أَبْعَدَ الْحَارِثَ الْمَلِكَ ابْنَ عَمْرٍو يتساءل الشاعر في هذا البيت عن الحارث الملك بن عمرو، الذي كان يملك العراق حتى عمان، مستخدماً أسلوب الاستفهام الإنكاري لتأكيد بعده عن مكانته السابقة، ويعكس ذلك تغيراً في الوضع السياسي أو الاجتماعي. كلمة "ملك" تدل على عظمته، بينما ذكر "العراق إلى عمان" يوسع من نطاق سلطته، مما يبرز عظمته السابقة بالمقارنة مع وضعه الحالي.</p> <p>هُوَأَنَا مَا أُتِيحُ مِنَ الْهُوَانِ مُجَاوَرَةً بَنِي شَمْجَى بْنِ جَرْمٍ يصف الشاعر في هذا البيت المجاورة بين الحارث وبني شمجى بن جرم، مستخدماً مجاز الاستعارة في "مجاورة" لتصوير القرب الجغرافي والسياسي، بينما "هواناً" تظهر الذل الذي أصاب الحارث، و"ما أتيح من الهوان" يبرز شدة هذا الذل باستخدام أسلوب المبالغة. الجالس بين "مجاورة" و"هواناً" يضفي تناغماً لغوياً يعزز المعنى.</p> <p>مَعِيزُهُمْ خَنَّاكَ ذَا الْخَنَانِ وَيَمْنَعُنْهَا بَنُو شَمْجَى بْنِ جَرْمٍ يُظْهِر هذا البيت منع بني شمجى بن جرم للحارث من معيزتهم، مستخدمين "خَنَّانك ذا الخنان" كاستعارة مكينة للتعبير عن رحمتهم التي تتحول إلى سلاح ضده. التشبيه الضمني بين الخنان والقوة يبرز تناقضاً درامياً، بينما الطباق بين "منع" و"خنان" يبرز تناقض القوة والضعف، مما يُظْهِر تفوق القبيلة على الملك السابق.</p>	<p>لَهُ مَلِكُ الْعِرَاقِ إِلَى عُثْمَانَ أَبْعَدَ الْحَارِثَ الْمَلِكَ ابْنَ عَمْرٍو يبدأ الشاعر بذكر الحارث الملك بن عمرو، الذي كان يملك مملكة واسعة تمتد من العراق إلى عمان. المعنى الحرفي هو تحديد هوية الحارث وتوضيح نطاق حكمه. لكن المعنى المجازي يتجاوز ذلك، حيث يضع الأساس لفكرة التناقض بين المجد والهوان، حيث كان الحارث ملكاً عظيماً ولكنه أصبح لاحقاً في موقف ضعيف.</p> <p>هُوَأَنَا مَا أُتِيحُ مِنَ الْهُوَانِ مُجَاوَرَةً بَنِي شَمْجَى بْنِ جَرْمٍ هنا، يصف الشاعر حالة الحارث بعد أن أصبح جارا لبني شمجى بن جرم. الكلمة "هواناً" تشير إلى حالة الضعف والمهانة التي أتاحت له، مما يدل على التناقض الصارخ مع ماضيه المجيد. يستخدم الشاعر الطباق بين "مجاورة" و"هواناً" لتعزيز هذا التناقض.</p> <p>مَعِيزُهُمْ خَنَّاكَ ذَا الْخَنَانِ وَيَمْنَعُنْهَا بَنُو شَمْجَى بْنِ جَرْمٍ في هذه الأبيات، يُعْرِب الشاعر عن منع بني شمجى بن جرم للحارث من دخول أرضهم. عبارة "معيزهم خَنَّانك ذا الخنان" هي كناية عن كرمهم وخنائهم، حيث يمنعون أرضهم من أن تصبح ملاذاً للحارث، مما يدل على تعاطفهم العميق. الاستعارة الضمنية هنا هي تشبيه الأرض بالملاد الأمن.</p>

Figure 4: **Qualitative Comparison of Model-Generated Explanations for a Single Arabic Poem.** This figure presents a representative Arabic poem alongside its original human-written explanation and corresponding verse-by-verse explanations generated by four different language models. The comparison highlights how each model interprets the poem’s rhetorical devices, imagery, and thematic depth relative to the gold explanation. This qualitative analysis illustrates variations in faithfulness, fluency, and literary sensitivity, offering insight into each model’s ability to handle nuanced Arabic poetic language and convey its intended meaning.

Overall, closed models such as GPT-4o and Gemini-2.5-Flash achieve consistently strong scores across both automatic and human evaluations. Notably, Gemini-2.5-Flash attains the highest textual entailment score (0.7475), along with high fluency and faithfulness scores, indicating strong alignment with poetic content and natural language clarity. GPT-4o also performs well across all dimensions, with the highest BERTScore and a strong balance of semantic coherence and linguistic quality. Among open models, Deepseek-V3, Aya-Expanse, and Qwen-3 show competitive performance, especially in fluency and textual entailment. However, models like ALLaM-7B and

AceGPT-v2 lag significantly in both lexical and semantic overlap, as well as in human-judged fluency and consistency.

A key insight from this evaluation is that most state-of-the-art models perform well on content expressed in Modern Standard Arabic (MSA) but struggle with the classical forms and linguistic intricacies present in historical and poetic Arabic. Despite high scores in generic semantic metrics, many models fail to capture deeper cultural and metaphorical meanings embedded in traditional Arabic poetry. Our analysis highlights the importance of domain-specific evaluation for literary and cultural tasks. It also underscores the need for

building or fine-tuning models that are more sensitive to classical Arabic forms. The gap between fluency and interpretive depth in some models suggests that future research should focus not just on surface-level correctness but also on deeper reasoning and cultural grounding. Such efforts are essential for advancing Arabic NLP in creative and heritage-preserving applications.

Table 4 shows era-wise performance of closed and open-source models on Arabic poem understanding using BERTScore, which captures semantic similarity with human explanations. Closed models like GPT-4o and Gemini variants perform consistently well, especially on modern and recent historical eras. In contrast, open models such as Deepseek-V3 and Aya-Expanse perform reasonably on some eras but struggle with older poetic forms like Pre-Islamic and Umayyad due to their complex language and cultural depth. This highlights that while current models are effective on MSA, they face challenges with classical Arabic. A complementary analysis using Textual Entailment is included in the Appendix (refer Table 18), further supporting these findings.

Additionally, Figure 4 shows a qualitative comparison of model-generated explanations for a classical Arabic poem. It compares outputs from GPT-4o, Gemini 2.5 Flash, Qwen 3, and Aya-Expanse against a human-written explanation. The figure highlights differences in faithfulness, fluency, and interpretive depth, showing how well each model captures the poem’s meaning, style, and literary richness. This example clearly illustrates the strengths of advanced models like GPT-4o in understanding nuanced poetic language.

5 Conclusion

Arabic poetry represents one of the richest and most culturally nuanced forms of expression within the Arabic language, characterized by layered meanings, stylistic diversity, and deep historical roots. In this paper, we introduced Fann or Flop, the first benchmark specifically developed to evaluate the capabilities of LLMs in understanding Arabic poetry across 12 historical eras, spanning from pre-Islamic to contemporary periods, and encompassing a broad spectrum of poetic genres and metrical forms. Our benchmark includes carefully curated diagnostic questions aimed at assessing semantic comprehension, metaphorical interpretation, prosodic awareness, and sensitivity to cultural

contexts. Through extensive evaluation, we demonstrated that despite strong performances on standard Arabic language tasks, state-of-the-art LLMs consistently struggle with the interpretative and culturally embedded dimensions of Arabic poetic texts. By releasing Fann or Flop as an open-source resource, we aim to encourage further research, promote rigorous assessment methodologies, and support advancements in linguistically and culturally rich Arabic language modeling.

6 Limitations and Societal Impact

While Fann or Flop provides a rigorous framework for evaluating LLMs’ understanding of Arabic poetry, it has several limitations. The benchmark covers only a portion of the broader Arabic poetic tradition, as some poems could not be included due to missing metadata, unclear authorship, or lack of reliable era or genre annotations. Additionally, poetry often invites multiple valid interpretations, which current evaluation metrics may not fully capture, even with expert-curated references. Expanding the dataset to include more diverse annotations, as well as dialectal and regional poetic forms, remains a key area for future work.

On the societal front, this benchmark contributes to the preservation and computational accessibility of Arabic literary heritage by positioning poetry as a meaningful testbed for language understanding. By promoting the development of culturally informed and linguistically grounded models, Fann or Flop encourages more inclusive and context-sensitive NLP. Nonetheless, as with any system trained on culturally rich and potentially sensitive material, there is a risk of misinterpretation or misuse. Ensuring transparency, human oversight, and responsible deployment is essential to safeguard the ethical impact of this work, especially in educational, literary, and public-facing applications.

References

- Muhammad Abdul-Mageed, Shady Elbassuoni, Jad Doughman, AbdelRahim Elmadany, El Moatez Billah Nagoudi, Yorgo Zoughby, Ahmad Shaher, Iskander Gaba, Ahmed Helal, and Mohammed El-Razzaz. 2021. *DiaLex: A benchmark for evaluating multi-dialectal Arabic word embeddings*. In *Proceedings of the Sixth Arabic Natural Language Processing Workshop*, pages 11–20, Kyiv, Ukraine (Virtual). Association for Computational Linguistics.
- Muhammad Abdul-Mageed, AbdelRahim Elmadany, and El Moatez Billah Nagoudi. 2020. Arbert &

- marbert: Deep bidirectional transformers for arabic. *arXiv preprint arXiv:2101.01785*.
- Sajawel Ahmed, Rob van der Goot, Misbahur Rehman, Carl Kruse, Ömer Özsoy, Alexander Mehler, and Gemma Roig. 2022. [Tafsir dataset: A novel multi-task benchmark for named entity recognition and topic modeling in classical Arabic literature](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3753–3768, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Google AI. 2025a. [Gemini 2.0 flash](#). Large language model, accessed May 20, 2025.
- Google AI. 2025b. [Gemini 2.5 flash](#). Large language model (Preview), accessed May 20, 2025.
- Abu Nasr al Jawhari. 10th Century. Taj al-lugha wa sihah al-arabiya - Wikipedia — [en.wikipedia.org. https://en.wikipedia.org/wiki/Abu_Nasr_al-Jawhari](https://en.wikipedia.org/wiki/Abu_Nasr_al-Jawhari). [Accessed 06-05-2025].
- alsharekh. 2019. Al-mujam al muaser— lexicon.alsharekh.org. <https://lexicon.alsharekh.org/>. [Accessed 06-05-2025].
- AlSuyuti. 15th Century. Al-mizhar fi eulum allughat wa'anwaeiha - Wikipedia — [en.wikipedia.org. https://en.wikipedia.org/wiki/Al-Suyuti](https://en.wikipedia.org/wiki/Al-Suyuti). [Accessed 06-05-2025].
- Zaid Alyafeai, Maged S Al-Shaibani, and Moataz Ahmed. 2023. Ashaar: Automatic analysis and generation of arabic poetry using deep learning approaches. *arXiv preprint arXiv:2307.06218*.
- Toni Andrews. 2024. Is Arabic The Richest Language In Words? - Interpreters & Translators, Inc. — [ititranslates.com. https://ititranslates.com/is-arabic-the-richest-language-in-words/](https://ititranslates.com/is-arabic-the-richest-language-in-words/). [Accessed 06-05-2025].
- Wissam Antoun, Fady Baly, and Hazem Hajj. 2020. Arabert: Transformer-based model for arabic language understanding. *arXiv preprint arXiv:2003.00104*.
- Mujam Ar-Riyadh. 2025. Mujam ar-riyadh— dictionary.ksaa.gov.sa. <https://dictionary.ksaa.gov.sa/>. [Accessed 06-05-2025].
- M Saiful Bari, Yazeed Alnumay, Norah A. Alzahrani, Nouf M. Alotaibi, Hisham A. Alyahya, Sultan Al-Rashed, Faisal A. Mirza, Shaykhah Z. Alsubaie, Hassan A. Alahmed, Ghadah Alabduljabbar, Raghad Alkhathran, Yousef Almushayqih, Raneem Alnajim, Salman Alsubaihi, Maryam Al Mansour, Majed Al-rubaian, Ali Alammari, Zaki Alawami, Abdulmohsen Al-Thubaity, Ahmed Abdelali, Jeril Kuriakose, Abdalghani Abujabal, Nora Al-Twairesh, Areeb Alowisheq, and Haidar Khan. 2024. [Allam: Large language models for arabic and english](#). *Preprint, arXiv:2407.15390*.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439.
- Houda Bouamor, Nizar Habash, Mohammad Salameh, Wajdi Zaghoulani, Owen Rambow, Dana Abdulrahim, Ossama Obeid, Salam Khalifa, Fadhil Eryani, Alexander Erdmann, et al. 2018. The madar arabic dialect corpus and lexicon. In *Proceedings of the eleventh international conference on language resources and evaluation (LREC 2018)*.
- Tuhin Chakrabarty, Arkadiy Saakyan, Debanjan Ghosh, and Smaranda Muresan. 2022. Flute: Figurative language understanding through textual explanations. *arXiv preprint arXiv:2205.12404*.
- Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2025. Sharegpt4v: Improving large multi-modal models with better captions. In *European Conference on Computer Vision*, pages 370–387. Springer.
- John Dang, Shivalika Singh, Daniel D'souza, Arash Ahmadian, Alejandro Salamanca, Madeline Smith, Aidan Peppin, Sungjin Hong, Manoj Govindassamy, Terrence Zhao, et al. 2024. Aya expanse: Combining research breakthroughs for a new multilingual frontier. *arXiv preprint arXiv:2412.04261*.
- FIGLANG2024. 2024. FigLang2024 — [sites.google.com. https://sites.google.com/view/figlang2024](https://sites.google.com/view/figlang2024). [Accessed 07-05-2025].
- Giuseppe Gallipoli and Luca Cagliero. 2025. It is not a piece of cake for gpt: Explaining textual entailment recognition in the presence of figurative language. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 9656–9674.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2021. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing. *arXiv preprint arXiv:2111.09543*.
- Huang Huang, Fei Yu, Jianqing Zhu, Xuening Sun, Hao Cheng, Dingjie Song, Zhihong Chen, Abdulmohsen Alharthi, Bang An, Juncai He, et al. 2023. Acegpt, localizing large language models in arabic. *arXiv preprint arXiv:2309.12053*.
- Salam Khalifa, Nizar Habash, Fadhil Eryani, Ossama Obeid, Dana Abdulrahim, and Meera Al Kaabi. 2018. [A morphologically annotated corpus of emirati Arabic](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*.
- Emmy Liu, Chen Cui, Kenneth Zheng, and Graham Neubig. 2022. Testing the ability of language models to interpret figurative language. *arXiv preprint arXiv:2204.12632*.
- Quentin Malartic, Nilabhra Roy Chowdhury, Ruxandra Cojocaru, Mugariya Farooq, Giulia Campesan, Yasser Abdelaziz Dahou Djilali, Sanath Narayan, Ankit Singh, Maksim Velikanov, Basma El Amel Boussaha, et al. 2024. Falcon2-11b technical report. *arXiv preprint arXiv:2407.14885*.
- Ibn Manzur. 14th Century. Lisan al-arab — en.wikipedia.org. https://en.wikipedia.org/wiki/Ibn_Manzur. [Accessed 06-05-2025].
- Karima Meftouh, Salima Harrat, Salma Jamoussi, Mourad Abbas, and Kamel Smaili. 2015. Machine translation experiments on padic: A parallel arabic dialect corpus. In *Proceedings of the 29th Pacific Asia conference on language, information and computation*, pages 26–34.
- Meta AI. 2024. [Llama 3.2: Revolutionizing edge ai and vision with open, customizable models](#).
- Behrang Mohit, Nathan Schneider, Rishav Bhowmick, Kemal Oflazer, and Noah A Smith. 2012. Recall-oriented learning of named entities in arabic wikipedia. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, pages 162–173.
- Hussein Mozannar, Karl El Hajal, Elie Maamary, and Hazem Hajj. 2019. Neural arabic question answering. *arXiv preprint arXiv:1906.05394*.
- Ossama Obeid, Nasser Zalmout, Salam Khalifa, Dima Taji, Mai Oudah, Bashar Alhafni, Go Inoue, Fadhl Eryani, Alexander Erdmann, and Nizar Habash. 2020. Camel tools: An open source python toolkit for arabic natural language processing. In *Proceedings of the twelfth language resources and evaluation conference*, pages 7022–7032.
- Susanna Olivero. 2024. *Figurative Language Understanding based on Large Language Models*. Ph.D. thesis, Politecnico di Torino.
- OpenAI. 2024. [Gpt-4o mini: advancing cost-efficient intelligence](#).
- OpenAI. 2024. [Gpt-4o system card](#). *Preprint*, arXiv:2410.21276.
- oussama. 2024. Modern Standard Arabic – The Missing Glossary - — blog.jarrousse.org. <https://blog.jarrousse.org/2024/03/27/modern-standard-arabic-the-missing-glossary/>. [Accessed 06-05-2025].
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Maja Popović. 2017. [chrF++: words helping character n-grams](#). In *Proceedings of the Second Conference on Machine Translation*, pages 612–618, Copenhagen, Denmark. Association for Computational Linguistics.
- Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, and et al. 2024. [Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context](#). *Preprint*, arXiv:2403.05530.
- Hassan Sajjad, Ahmed Abdelali, Nadir Durrani, and Fahim Dalvi. 2020. Arabench: Benchmarking dialectal arabic-english machine translation. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5094–5107.
- Neha Sengupta, Sunil Kumar Sahu, Bokang Jia, Satheesh Katipomu, Haonan Li, Fajri Koto, William Marshall, Gurpreet Gosal, Cynthia Liu, Zhiming Chen, et al. 2023. Jais and jais-chat: Arabic-centric foundation and instruction-tuned open generative large language models. *arXiv preprint arXiv:2308.16149*.
- Fanar Team, Ummar Abbas, Mohammad Shahmeer Ahmad, Firoj Alam, Enes Altinisik, Ehsannedin Asgari, Yazan Boshmaf, Sabri Boughorbel, Sanjay Chawla, Shammur Chowdhury, Fahim Dalvi, Kareem Darwish, Nadir Durrani, Mohamed Elfeky, Ahmed Elmagarmid, Mohamed Eltabakh, Masoomali Fatehkia, Anastasios Fragkopoulous, Maram Hasanain, Majd Hawasly, Mus’ab Husaini, Soon-Gyo Jung, Ji Kim Lucas, Walid Magdy, Safa Messaoud, Abubakr Mohamed, Tasnim Mohiuddin, Basel Mousi, Hamdy Mubarak, Ahmad Musleh, Zan Naeem, Mourad Ouzani, Dorde Popovic, Amin Sadeghi, Husrev Taha Sencar, Mohammed Shinoy, Omar Sinan, Yifan Zhang, Ahmed Ali, Yassine El Kheir, Xiaosong Ma, and Chaoyi Ruan. 2025. [Fanar: An arabic-centric multimodal generative ai platform](#). *Preprint*, arXiv:2501.13944.
- Qwen Team. 2025. [Qwen3](#).
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- wikipedia. 2025. ar.wikipedia.org. https://ar.wikipedia.org/wiki/%D8%A7%D9%84%D9%84%D8%BA%D8%A9_%D8%A7%D9%84%D8%B9%D8%B1%D8%A8%D9%8A%D8%A9. [Accessed 06-05-2025].

- wikipediaArabic. 2025. Varieties of Arabic - Wikipedia — en.wikipedia.org. https://en.wikipedia.org/wiki/Varieties_of_Arabic. [Accessed 06-05-2025].
- Taha Zerrouki and Amar Balla. 2017. Tashkeela: Novel corpus of arabic vocalized texts, data for auto-diacritization systems. *Data in brief*, 11:147.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.
- Cheng Zhao, Bin Wang, and Zhen Wang. 2024. Understanding literary texts by llms: A case study of ancient chinese poetry. *arXiv preprint arXiv:2409.00060*.

A Appendix

This appendix provides supplementary material to support our study of Arabic poetry understanding in language models. It includes four key sections: (1) a brief overview of related work in Arabic NLP, highlighting recent progress in benchmark development and the specific gaps our work addresses; (2) detailed dataset statistics, including token distribution, genre coverage, and temporal representation across poetic eras; (3) additional details on the prompts used for model generation and evaluation; and (4) a selection of qualitative examples from the *Fann or Flop* benchmark that illustrate its richness and the interpretive challenges it presents. Together, these components underscore the linguistic, historical, and cultural depth of our dataset and evaluation framework.

B Related Work

Understanding Arabic poetry computationally intersects with multiple subfields of NLP, including language modeling, data set construction, figurative language interpretation, and the evaluation of cultural knowledge. To contextualize our contribution, we review prior work across two key domains: Arabic NLP benchmarks and poetry understanding in LLMs.

B.1 Arabic NLP Benchmarks

Over the past decade, Arabic NLP has advanced considerably with the introduction of large-scale benchmarks such as SOQAL (Arabic-SQuAD and ARCD) (Mozannar et al., 2019), AraBench (Sajjad et al., 2020), and the AraBERT Collection (Antoun et al., 2020). These benchmarks cover essential tasks such as sentiment analysis, named entity recognition (NER), and question answering, and typically support both MSA and dialectal varieties. However, they largely overlook CA, which remains underrepresented in the main resources. Consequently, while models trained on these datasets perform well on surface-level tasks, they lack the depth to assess cultural, rhetorical, and literary understanding, especially in classical poetic contexts.

Additional resources such as the CAMEL corpus (Abdul-Mageed et al., 2020; Khalifa et al., 2018), Tashkeela (Zerrouki and Balla, 2017), PADIC (Meftouh et al., 2015), and MADAR (Bouamor et al., 2018) have enriched the field through morphologically annotated corpora, diacritized texts, and dialectal content. However,

these datasets are primarily designed for structural tasks such as morphological disambiguation or dialect identification, without engaging the semantic or figurative dimensions of the poetic language.

More recently, efforts have extended Arabic NLP to the literary and religious domains. The Tafsir dataset (Ahmed et al., 2022) introduces a benchmark derived from *Tafsir al-Tabari*, including NER and topic modeling in CA. AQMAR (Mohit et al., 2012) targets recall-oriented NER in Arabic Wikipedia, offering annotations across standard and domain-specific entity types. Although both datasets engage with classical Arabic and semantic granularity, they do not address poetry or the interpretive challenges it poses.

Among the most directly relevant efforts is Ashaar (Alyafeai et al., 2023), the first large-scale Arabic poetry dataset. It includes tasks such as meter classification, era identification, and poet recognition, along with descriptive metadata. Despite its contributions to computational poetics, Ashaar lacks verse-level annotation, rhetorical device modeling, question-answer style interpretation, and historical contextualization, limiting its ability to evaluate deeper poetic reasoning in language models.

B.2 Poetry Understanding in NLP

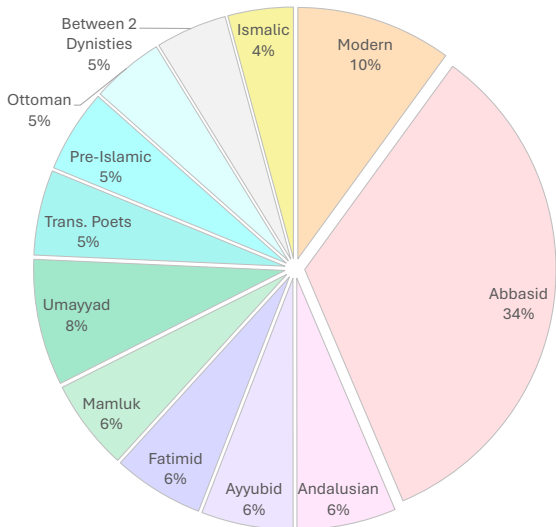
Outside Arabic, poetry and figurative language have emerged as valuable testbeds for assessing the reasoning of LLM (Liu et al., 2022; Bisk et al., 2020; Olivero, 2024). Benchmarks like FLUTE (Chakrabarty et al., 2022) and the FigLang shared tasks (FIGLANG202, 2024) reveal persistent challenges in handling metaphor, simile, and symbolic expression. Recent works (Gallipoli and Cagliero, 2025; Zhao et al., 2024) further expose the limitations of LLMs in interpreting literary texts, including complex poetic structures and non-literal meaning. Despite Arabic’s longstanding poetic legacy, this evaluation line remains largely unexplored for Arabic, leaving a notable gap in culturally grounded reasoning tasks.

Fann or Flop addresses this gap by combining a chronological taxonomy of Arabic poetry with interpretive question-answering. It spans 12 eras and integrates dialectal variation, rhetorical analysis, historical context, and verse-level annotation. As summarized in Table 1, no existing benchmark offers this breadth of poetic features, positioning Fann or Flop as the first comprehensive diagnostic tool for evaluating Arabic poetic understanding in LLMs.

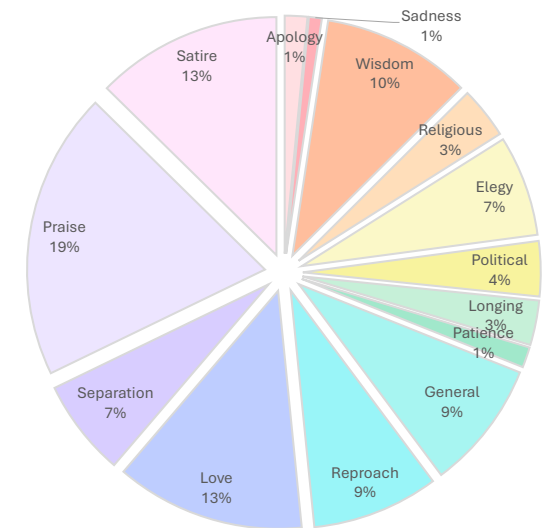
C Fann or Flop Data Statistics

To better characterize the distributional properties of our curated Arabic poetry dataset, we present a series of descriptive statistics that cover both historical and thematic dimensions.

These include the distribution of poems across major eras (Figure 5a), the overall distribution of poetic genres (Figure 5b), and a genre-by-era breakdown (Figure 6).



(a) **Distribution of poems by historical era.** The chart shows the proportion of poems collected from each era. Abbasid, Modern, and Andalusian periods are the most represented, reflecting their central role in Arabic literary production.



(b) **Distribution of poems by genre.** This chart shows the proportion of poetic genres across the dataset. Praise, Satire, and Love dominate the distribution, while genres such as Apology and Sadness appear less frequently.

Figure 5: **Era and Genre Statistics.** Subfigure (a) displays the distribution of poems across historical eras, while subfigure (b) shows the overall genre distribution across the dataset.

Complementing these visualizations, we also include detailed per-era tables listing the most represented poets and the number of poems attributed to each. Together, these statistics contextualize the coverage of the dataset and support downstream applications such as genre classification, diachronic literary analysis, and poet-specific modeling.

Tables 5 to 16 provide a breakdown of the number of poems attributed to prominent Arabic poets across different historical eras. Each table is dedicated to one era:

Table 5: Pre-Islamic era; Table 6 Transitional (Early Islamic) period; Table 7: Modern era; Table 8: Islamic era; Table 9: Umayyad era; Table 10: Abbasid era; Table 11: Between 2 Dynasties; Table 12: Fatimid Dynasty; Table 13: Andalusian era; Table 14: Ayybid era. Table 15: Mamluk Dynasty; Table 16: Ottoman era.

Era		Pre-Islamic
Poet	Poems	
Aws ibn Hajar	35	
al-Samaw'al	12	
al-Sulayk ibn al-Sulaka	7	
Imru' al-Qais ibn Hujr	34	
Zuhayr ibn Abi Sulma	48	
Salama ibn Jandal	14	
Tarfah ibn al-Abd	26	
Urwah ibn al-Ward al-Absi	31	
Ubayd ibn al-Abras	40	
Amr ibn Qami'a	21	
Amr ibn Kulthum	24	
Antarah ibn Shaddad	82	
Total	374	

Table 5: Poem counts for major poets from the Pre-Islamic era.

Era		Transitional Poet
Poet	Poems	
Al-Hadira	7	
Al-Hutay'a	95	
Al-Khansa	92	
Hassan ibn Thabit	74	
Amir ibn al-Tufayl	41	
Amr ibn Barraqa	5	
Labid ibn Rabi'a	69	
Total	383	

Table 6: Poem counts for major poets from the Early-Islamic Transitional period.

Era		Modern
Poet	Poems	
Ahmed Shawqi	460	
Hafiz Ibrahim	240	
Total	700	

Table 7: Poem counts for major poets from the Modern era.

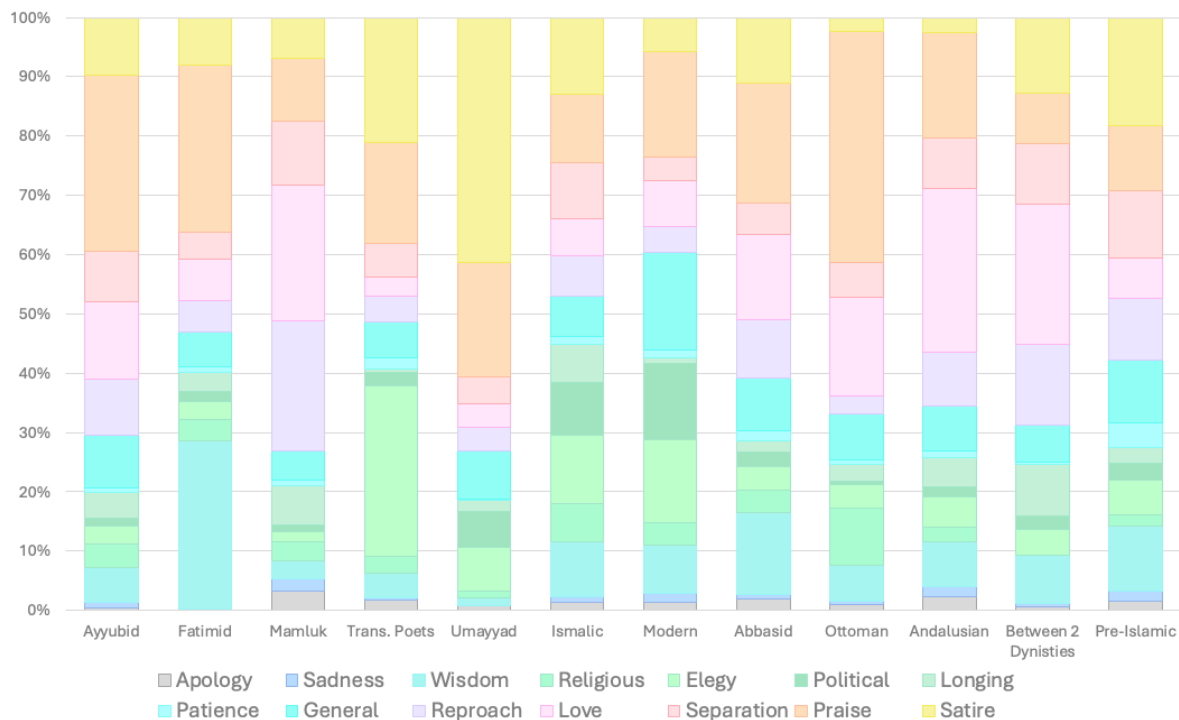


Figure 6: **Genre distribution across historical eras.** This stacked bar chart illustrates how poetic themes evolved across different dynasties. It highlights patterns such as the prominence of Praise and Satire during the Abbasid and Umayyad eras, and the diverse thematic expression in Modern poetry.

Era	Islamic
Poet	Poems
Abu Muhammad al-Faq'asi	28
Al-Akraa bin Muath Al-Kushairi	16
Asmaa Bin Kharja El-Fazari	4
Aasha Taroud	4
Khuzaima ibn Thabit al-Ansari	14
Khalid ibn al-Walid	7
Az-Zzubayr bin Al-Awam	4
As-Samhari Al-Okliy	11
Al-Ghitamish Al-Dabbi	4
Abd al-Rahman ibn Abi Bakr al-Siddiq	3
Jubaiha' al-Ashja'i	7
Habib ibn Khidrah al-Hilali	6
Ka'b ibn Mashhur al-Makhbali	13
Mas'ud al-Mazini	3
Satira al-USaybiyya	4
Ziyad ibn Abihi	3
Ziyad ibn Hanzala al-Tamimi	4
Murrah ibn Junada	3
Atika bint Zayd	6
Abd al-Aziz ibn Zararah al-Kalabi	6
Urwa ibn Hizam	6
Ali ibn al-Husayn	6
Amr ibn al-'As	26
Amra bint Mirdas	4
Other Poet with 1 or 2 poems	102
Total	294

Table 8: Poem counts for major poets from the Islamic era.

Era	Umayyad Era
Poet	Poems
Al-Akhtal	136
Jarir	239
Al-Farazdak	178
Ubaydallah ibn al-Ruqayyat	15
Total	568

Table 9: Poem counts for major poets from the Umayyad Era.

Era	Abbasid
Poet	Poems
Abu al-Atahiya	362
Abu Firas al-Hamdani	129
Abu Nuwas	701
Abu Tammam	29
Ibn al-Rumi	227
Imam al-Shafi'i	20
al-Buhturi	601
al-Mutanabbi	273
Total	2342

Table 10: Poem counts for major poets from the Abbasid era.

Era	Between Dynasties
Poet	Poems
Bashar bin Bord	321
Total	321

Table 11: Poem counts for major poets from the Between 2 Dynasties.

Era	Fatimid Dynasty
Poet	Poems
Abu al-Ala al-Ma'arri	183
Ibn Hayyus	120
Arqala al-Kalbi	106
Total	409

Table 12: Poem counts for major poets from the Fatimid Dynasty.

Era	Andalusian
Poet	Poems
Abu Ishaq al-Albiri	38
Ibn Khafaja	225
Ibn Zaydun	146
Ibn Sahl al-Andalusi	37
Total	446

Table 13: Poem counts for major poets from the Andalusian era.

Era	Ayyubid Dynasty
Poet	Poems
Ibn al-Farid	35
Sibt Ibn al-Tawawidhi	291
Muhyiddin Ibn Arabi	87
Total	413

Table 14: Poem counts for major poets from the Ayyubid Dynasty.

Era	Mamluk Dynasty
Poet	Poems
Baha al-Din Zuhayr	368
Safiyy al-Din al-Hilli	40
Total	408

Table 15: Poem counts for major poets from the Mamluk Dynasty.

Era	Ottoman
Poet	Poems
Abu al-Ma'ali al-Talawi	75
Ibn Razka	19
Ibn Matuq al-Musawi	74
al-Kawkabani	2
Bint al-Shuhna	2
Abd al-Rahman al-Musili	58
Muhammad al-Isba'i	31
Muhammad al-Sharafi al-Safaqsi	65
Total	326

Table 16: Poem counts for major poets from the Ottoman era.

Prompt

You are an academic expert in Arabic literature and poetry analysis. Your task is to provide a deep linguistic, rhetorical, and literary explanation for every verse of the Arabic poem provided below, regardless of its style or period.

The poem to be analyzed is: **<POEM_CONTENT>**

Please adhere strictly to the following guidelines:

1. **Carefully read the entire poem to understand its deep meaning and global message**, general theme, purpose, emotional tone, and cultural or historical context, which should inform your verse-by-verse analysis.

2. **For each verse** (or each paired line, if the poem follows a two-hemistich structure), write one cohesive paragraph in formal Arabic that integrates:

- **The literal meaning** (المعنى الحرفي): What the poet is directly saying.
- **The figurative, symbolic, or rhetorical meaning** (المعنى المجازي): Emotional or intellectual connotations, artistic imagery, and rhetorical devices, such as التثنية, الاستعارة, الكناية, etc. You must explicitly name these forms and explain their function within the verse's meaning and effect.

3. **The explanation must be:**

- **Unified** : Do not separate the literal and figurative meanings. Present them in a single integrated paragraph.
- **Focused on the verse alone** : Do not reference other verses explicitly within your explanation. However, you may consider the overall poem global message and meaning implicitly for interpretive accuracy.

4. **Avoid generic commentary** (e.g., "the poet expresses love or sadness"). Instead, your analysis must refer directly to the language and structure of the verse and explain how the meaning is conveyed.

5. **If there are multiple possible interpretations, choose the clearest and most textually supported one**, without mentioning alternative readings.

6. Each paragraph should be between **2 to 4 well-formed sentences**, reflecting literary insight and academic precision.

7. **Ensure** that every verse of the poem is explained. **Do not skip or summarize verses.**

8. **Generate a global explanation** of the whole poem piece to support your verse explanation and refer to it to ensure consistency and connectivity and coherence from start to end. The explanation should not exceed a few lines.

9. **Return your output** in the following strict and valid JSON format only:

```

{
  "No extra text, no explanations, no markdown:":
  "json":
  {
    "explanation": "<full poem summary>",

    [
      {
        "verse": "<full text of verse 1>",
        "explanation": "<full Arabic explanation for verse 1>"
      },

      {
        "verse": "<full text of verse 2>",
        "explanation": "<full Arabic explanation for verse 2>"
      },

      //... continue for all remaining verses ...
    ]
  }
}

```

Figure 7: **The verse-level explanation prompt used for evaluation.** This prompt instructs the model to produce detailed verse-by-verse explanations in Arabic. It guides the model to integrate both literal and figurative interpretations, explicitly name rhetorical devices (e.g., metaphor, personification, paronomasia). The prompt enforces coherence, academic rigor, and structural consistency by requiring output in a strict JSON format.

D Annotator Rubric Guidelines

Table 17 presents the scoring rubric used by annotators. Each annotator assigns a single value for every criterion, and the criterion scores are summed to produce an individual interpretive depth score. The final reported score is the average of these values across annotators.

Criterion	Score	Guideline
Literal Comprehens.	0	Misinterprets or ignores the poem’s literal meaning; adds unrelated content.
	1	Correctly reflects the surface meaning of the poem.
Thematic & Emotional Depth	0	Fails to capture themes, sentiment, or tone.
	1	Identifies a theme or emotion only superficially.
	2	Clearly conveys central themes and emotions with nuance.
Cultural & Historical Appropriateness	0	Culturally/historically inaccurate; ignores context.
	1	Shows partial or limited awareness of context.
	2	Demonstrates accurate and insightful cultural/historical awareness.
Stylistic Sensitivity	0	No recognition of rhetorical/literary features.
	1	Mentions some features, but inaccurately or superficially.
	2	Correctly identifies key features, but lacks depth.
	3	Nuanced discussion of literary devices (e.g., metaphor, rhythm, imagery).
Expressiveness & Coherence	0	Incoherent, disorganized, or ungrammatical.
	1	Understandable but lacks fluency or cohesion.
	2	Clear, well-structured, and stylistically appropriate in Arabic.

Table 17: Annotator evaluation rubric of model-generated explanations. Each criterion is scored independently, with higher values indicating better performance.

E Prompts Used

E.1 Model Generation Prompt

To generate verse-level explanations suitable for evaluating both open- and closed-source models, we developed a carefully optimized generation prompt. The prompt design followed an iterative

and augmented process. Initially, we used a simple bilingual (in Arabic and English) prompt asking for explanations. Based on early outputs, which tended to capture local semantic meaning but lacked coherence and global context, we progressively refined the prompt to elicit more structured and connected responses.

Through multiple rounds of testing, expert evaluation, and prompt engineering, we incorporated explicit instructions to address both local (verse-specific) and global (poem-wide) interpretive elements as support. This enhancement significantly improved the quality of the generated explanations, resulting in outputs that were more coherent, context-aware, and semantically aligned with the original verses.

After extensive comparison, expert reviewers favored the English version of the prompt over its Arabic counterpart, as it more consistently achieved local-global alignment and produced well-connected, high-quality explanations. This final version of the English prompt (Figure 7) was adopted for all subsequent evaluations.

E.2 Model Evaluation Prompt

To ensure consistent and reliable automatic LLM-Judge evaluation of model-generated poem explanations, we designed a clear and structured system prompt (see Figure 8). The prompt positions the evaluator as an expert Arabic linguist and literary critic, responsible for assessing AI-generated verse-by-verse explanations against ground-truth references.

Each poem is evaluated on two key dimensions: Faithfulness/Consistency, which measures how accurately the explanation reflects the verse’s intended meaning, and Fluency/Grammaticality, which assesses the quality of the generated text in Modern Standard Arabic. Annotators assign a score from 1 to 5 for each criterion based on the overall performance across all verses, without providing per-verse feedback or open-ended commentary.

The prompt ensures simplicity, objectivity, and high inter-annotator agreement, making it well-suited for evaluating poetic reasoning in culturally rich and linguistically nuanced contexts like Arabic poetry.

F Additional Examples: Qualitative, Translated, and Quantitative Insights

In the following section, we present a more detailed evaluation of the Textual-Entailment (Refer

System Prompt used to evaluate the poem explanation for Faithfulness and Fluency metric

```
You are an expert Arabic linguist and literary evaluator.

Your task is to evaluate a full Arabic poem's verse-by-verse explanations.
You will compare ground-truth (human-written) explanations with
generated explanations from an AI model.

You will judge each verse explanation based on the following two criteria:

---

### Evaluation Criteria (per verse)

1. Faithfulness / Consistency:
  Is the generated explanation consistent with the meaning of the verse?
  - Score 5: Deeply faithful to the verse's content
  - Score 3: General alignment but loses poetic imagery
  - Score 1: Misinterprets or invents meaning

2. Fluency / Grammaticality:
  Is the generated explanation well-formed Modern Standard Arabic?
  - Score 5: Fluent, grammatically correct
  - Score 3: Understandable with minor issues
  - Score 1: Awkward, incomplete, or ungrammatical

---

### What You Will Receive

You will receive for each poem:
- 'poem_title'
- "ground_truth": a list of objects { "v": <int>, "text": <string> }
- "generated": a list of objects with the same v indices

---

### What You Must Do

- Compare all verses together and assign a single score of 1-5 for each
  criterion.
- Do not provide per-verse scores or any comments.

Then:
- Calculate average scores for the whole poem
- Provide an 'overall_score' (1-5) that reflects your judgment across all
  verses

---

Do NOT provide any comments or rationale.
Respond with valid JSON only in this format:

### Output Format (in JSON)

{
  "faithfulness_score": <1-5>,
  "fluency_score": <1-5>,
  "overall_score": <1-5>
}
```

Figure 8: System prompt used for LLM-Judge evaluation of verse-by-verse poem explanations. LLM (OpenAI, 2024) compare AI-generated outputs with original explanations and assign overall scores for faithfulness to meaning and fluency in MSA, following clearly defined criteria. The structured format ensures consistency and reliability across evaluations.

	Model	Pre-Islamic	Transitional	Early Islamic	Umayyad	Abbasid	Fatimid
Closed	GPT-4o-2024-08-06 (OpenAI, 2024)	0.6425	0.6502	0.7116	0.6166	0.6699	0.7050
	GPT-4o-mini-2024-07-18 (OpenAI, 2024)	0.4355	0.4789	0.5436	0.4200	0.4266	0.4532
	Gemini-2.5-Flash (AI, 2025b)	0.7275	0.7308	0.7527	0.7112	0.7417	0.7542
	Gemini-2.0-Flash (AI, 2025a)	0.6908	0.7156	0.7458	0.6798	0.7033	0.7462
	Gemini-1.5-Pro (Reid et al., 2024)	0.6004	0.6372	0.6497	0.6312	0.6035	0.6502
	Fanar-Star (Team et al., 2025)	0.6142	0.6354	0.6621	0.5900	0.6413	0.6717
Open	Deepseek-V3 (Liu et al., 2024)	0.5066	0.5875	0.6174	0.5482	0.4736	0.5581
	Deepseek-R1 (Guo et al., 2025)	0.5066	0.5875	0.6174	0.5482	0.4736	0.5581
	Llama-3.3-70B (Meta AI, 2024)	0.5456	0.5469	0.5747	0.5211	0.5341	0.5387
	Qwen-3 (Team, 2025)	0.6142	0.6354	0.6621	0.5900	0.6413	0.6717
	Aya-Expanse (Dang et al., 2024)	0.6142	0.6354	0.6621	0.5900	0.6413	0.6717
	ALLaM-7B (Bari et al., 2024)	0.5619	0.5630	0.6037	0.5844	0.5848	0.6158
	Jais (Sengupta et al., 2023)	0.6124	0.6289	0.6482	0.6047	0.6295	0.6421
	AceGPT-v2-70B-Chat (Huang et al., 2023)	0.5851	0.5656	0.6104	0.5770	0.6119	0.6095

	Model	Andalusian	Ayyubid	Mamluk	Between Dynasties	Ottoman	Modern
Closed	GPT-4o-2024-08-06 (OpenAI, 2024)	0.7128	0.6774	0.7393	0.6656	0.7379	0.6843
	GPT-4o-mini-2024-07-18 (OpenAI, 2024)	0.4869	0.4303	0.4507	0.4240	0.4836	0.3988
	Gemini-2.5-Flash (AI, 2025b)	0.7778	0.7416	0.7866	0.7398	0.7994	0.7544
	Gemini-2.0-Flash (AI, 2025a)	0.7527	0.7320	0.7698	0.7164	0.7585	0.6951
	Gemini-1.5-Pro (Reid et al., 2024)	0.6710	0.6074	0.6377	0.5971	0.6441	0.5965
	Fanar-Star (Team et al., 2025)	0.6749	0.6454	0.7105	0.6342	0.7151	0.6429
Open	Deepseek-V3 (Liu et al., 2024)	0.5927	0.5065	0.5448	0.4929	0.5226	0.4705
	Deepseek-R1 (Guo et al., 2025)	0.5927	0.5065	0.5448	0.4929	0.5226	0.4705
	Llama-3.3-70B (Meta AI, 2024)	0.5873	0.5221	0.5849	0.5129	0.5712	0.4897
	Qwen-3 (Team, 2025)	0.6749	0.6454	0.7105	0.6342	0.7151	0.6429
	Aya-Expanse (Dang et al., 2024)	0.6749	0.6454	0.7105	0.6342	0.7151	0.6429
	ALLaM-7B (Bari et al., 2024)	0.5892	0.6044	0.6736	0.5905	0.6556	0.6302
	Jais (Sengupta et al., 2023)	0.6540	0.6399	0.6812	0.6183	0.6625	0.6348
	AceGPT-v2-70B-Chat (Huang et al., 2023)	0.6215	0.6131	0.6683	0.5681	0.6273	0.6044

Table 18: **Era-wise Evaluation using Textual Entailment (TE)**. Era-wise performance of closed and open-source models on the Arabic poem understanding task, measured using the Textual Entailment metric. This metric evaluates how well the model-generated explanation logically aligns with the original poem content. The results are grouped across key historical eras, from Pre-Islamic to Modern, allowing a fine-grained view of model strengths and limitations across time periods. Closed models such as GPT-4o and Gemini variants demonstrate consistently high entailment across most eras, while select open models like Deepseek-V3 and Aya-Expanse show promising results in specific historical contexts. This analysis highlights the importance of temporal generalization and cultural grounding in building robust Arabic literary reasoning models.

Table 18) metric across the 12 historical eras, comparing both open-source and closed-source models on this dimension. To support comprehensive engagement and a clearer understanding of the data evaluated, we also include selected English translations of Arabic poetic samples, as well as additional qualitative Arabic examples. These additions offer deeper insight into the linguistic diversity, thematic range, and overall quality of the dataset used in our analysis.

F.1 English Translated Qualitative Samples

To support accessibility and improve cross-linguistic understanding, we translated selected Arabic poetry samples shown in Figure 2 into English. For this task, we used GPT-4o, which provided deeper context-aware translations that more

accurately capture the figurative and rhetorical nuances of the original verse, outperforming basic literal tools such as Google Translate. These translations allow non-Arabic speakers to more easily engage with the literary richness, emotional depth, and stylistic variety explored in our research.

F.2 Additional Qualitative Samples

To further showcase the dataset’s richness, we present additional qualitative samples spanning diverse historical periods, poetic genres (e.g., satire, elegy, political verse), and metrical patterns (See Figure 10). These examples were selected to demonstrate the stylistic, thematic, and rhetorical variety encountered in our evaluation.

Abbasid (العباسي)		Omayyad (الأموي)		Pre-Islamic (الجاهلي / قبل الإسلام)	
Poet: Abu Nuwas		Poet: Al-Farazdak		Poet: Imro' Al Quais	
Meter: Ar-Ramel	Genre: Wisdom	Meter: At-Tawil	Genre: Satire	Meter: At-Tawil	Genre: Love (Ghazal)
<p>Everyone who weeps—will one day be wept over</p> <p>Everyone remembered—will one day be forgotten</p> <p>Whoever rises—God is higher still</p> <p>It is for Him we strive and suffer</p>	<p>Everyone who mourns—will one day be mourned</p> <p>Everything stored—will one day perish</p> <p>None remains but God alone</p> <p>Indeed, some matters have already been taken care of for us</p>	<p>So go forth—my motherland is your people's land</p> <p>And praise Sa'd for what it deserves</p> <p>When the Pleiades fail to bring rain, their lightning still flashes in vain.</p>	<p>I see a vast, hollow land rich in fruit</p> <p>The best words from strangers are honest ones</p> <p>The solid dwellings neighbors feel no fear nearby</p>	<p>Stop, let us weep over the memory of a beloved and a home</p> <p>it became clear and the reading did not fade</p> <p>You see the droppings of wild deer scattered across its plains</p>	<p>In the valley of Al-Luwa, between Ad-Dakhul and Hawmal</p> <p>Because it was woven from the south and the north</p> <p>And its lowlands are like black peppercorns</p>
<p>This poem offers a profound depiction of existential truths, using vivid metaphors and rhetorical elegance. Here's a verse-by-verse explanation:</p> <p>Verse One: Everyone who mourns the dead will one day be mourned, and every weeper will himself be wept over. The verse reflects the inevitability of death, which spares no one.</p> <p>Verse Two: All that is hoarded—wealth or possessions—will eventually perish, and all who are remembered will one day be forgotten. A meditation on the impermanence of material wealth and fame.</p> <p>Verse Three: Only God remains. Whoever believes they have risen or triumphed, God is far above them. A declaration of God's eternal supremacy and the fleeting nature of human status.</p> <p>Verse Four: There are matters we have been entrusted with; for them, we strive and suffer. A reflection on the hardship of life and the human pursuit of purpose.</p>		<p>This poem speaks of the Tu'ah tribe, describing their character, values, and the poet's disapproval of certain behaviors among them. Here's a brief breakdown of each verse:</p> <p>Verse One: <i>So go forth—my motherland is your people's land.</i> A call from the poet to a woman (likely a poetess), inviting her to visit the land of Tu'ah, likening it to a mother—symbolizing deep emotional connection. <i>see a vast, hollow land rich in fruit.</i> Describes a fertile, expansive area with abundant yield.</p> <p>Verse Two: <i>And praise Sa'd for what it deserves.</i> The poet praises the tribe's goodness and noble character. <i>The best words from strangers are honest ones.</i> Emphasizes the tribe's reputation for truthfulness, which earns them trust even from outsiders.</p> <p>Verse Three: <i>The solid dwellings—neighbors feel no fear nearby.</i> A metaphor highlighting the tribe's hospitality and the safety they offer; their homes are places where others feel secure.</p>		<p>This poem is a heartfelt elegy reflecting the poet's deep yearning for a lost beloved and a cherished place. I will explain each verse individually:</p> <p>Verse One: The poet opens with a call to weep over the memory of a beloved and a once-inhabited home. <i>Sigt al-Luwa</i> is likely a mountainous site, located between <i>al-Dukhul</i> and <i>Hawmal</i>—two opposing landmarks.</p> <p>Verse Two: He describes the home's enduring traces—<i>Tawdih</i> and <i>al-Miqra</i>—still visible despite time, shaped by winds from the south and north, signifying the home's resilience and vastness.</p> <p>Verse Three: The scene is completed with scattered deer droppings in its courtyards and lowlands, likened to peppercorns—evoking a vivid image of abandonment and lingering traces of past life.</p>	
Transitional (المختلطة)		Islamic (الإسلامي)		Andalusian (الأندلسي)	
Poet: Labid Bin Rabi'a		Poet: Issa Bin Fatek		Poet: Ibn Khafaja	
Meter: Ar-Rojz	Genre: Separation	Meter: Al-Wafer	Genre: Politics	Meter: At-Tawil	Genre: Patience
<p>O Harem, you who are known for justice</p> <p>For his people would destroy my people</p> <p>(Between) their lineage and mine—there is no comparison.</p>	<p>If al-Ahwas reaches the water before me</p> <p>Do not equate their nature with mine</p> <p>I have warned against the foolishness of ignorance.</p>	<p>At daybreak, they prayed and rose to stand firm</p> <p>When they gathered in force, they charged at their foes</p> <p>For the rest of the day—until it reached them</p>	<p>They mounted their noble, swift steeds, well-trained and ready</p> <p>The sword-bearers kept on slaughtering us relentlessly.</p> <p>The darkness of night, in which they slipped away and evaded us</p>	<p>My complaint is complete when I see glory itself in sorrow</p> <p>I conceal a heart that cracks the chest with every sigh</p> <p>And how could I hide this burning fire I've found within me</p>	<p>It is enough of a tragedy that you see me in tears</p> <p>The echo of groans draws out calm, flowing tears</p> <p>Thirsting even as it emerges from the very source of water</p>
<p>This poem expresses a firm rejection of injustice and bias, showcasing the poet's courage in confronting corruption and distancing himself from a morally degraded environment. I will explain each verse individually:</p> <p>Verse One: The poet addresses <i>Haram</i>, symbolizing the judge or ruler, affirming his role as a man of justice. The line <i>"If al-Ahwas reaches the water before me"</i> metaphorically points to injustice, where someone less deserving claims rights ahead of the poet. It criticizes the misallocation of entitlements and the poet's exclusion.</p> <p>Verse Two: The poet prays that the oppressor and his kin be separated from his own people, reflecting a complete disassociation. In saying <i>"Do not equate their nature with mine,"</i> he urges the ruler not to compare or unite him with the wrongdoers, firmly rejecting any resemblance, affiliation, or shared fate.</p>		<p>The poem narrates a battle between a small group of faithful warriors and a much larger opposing force. Through faith, discipline, and resilience, the minority emerges victorious, while the poem denounces tyranny and elevates divine justice. I will explain each verse individually:</p> <p>Verse 1: At dawn, the believers begin their day with prayer before mounting their noble steeds (<i>jurd al-itaq</i>), symbolizing their preparedness and noble resolve.</p> <p>Verse 2: The larger army launches a fierce and sudden assault. The mention of <i>"sword-bearers"</i> (<i>dhaw al-ja'ith</i>) underscores the violence and high casualties inflicted upon the faithful.</p> <p>Verse 3: The battle rages until nightfall, when the attackers use the cover of darkness to escape—highlighting the spiritual and moral triumph of the steadfast believers.</p>		<p>This poem expresses the poet's profound emotional turmoil, using powerful imagery to convey his grief. Below is a verse-by-verse analysis:</p> <p>Verse 1: The poet finds it enough to see glory itself lamenting—his personal sorrow is a reflection of the decline of noble values. "Glory" serves as a symbol of lost greatness, making his grief all the more universal.</p> <p>Verse 2: He hides a heart overwhelmed by pain, whose sighs are so intense they "split the chest" and bring tears flowing like milk—a vivid metaphor for suppressed anguish.</p> <p>Verse 3: Immersed in sorrow, he questions how he could possibly conceal it—like one soaked by a spring who cannot pretend to be dry. The image stresses the impossibility of hiding deep emotional wounds.</p>	

Figure 9: Translated Samples. This figure presents English translations of the Arabic samples shown in Figure 2. The translations are included to facilitate understanding and accessibility for non-Arabic speakers, allowing broader engagement with the poetic content without requiring prior knowledge of Arabic.

<p>(غزلGhazal Love)</p> <p>أيا ما أحسبها مقلّة ولولا الملاحظة لم أعجب خلوقيّة في خلوقيها سؤداء من عتب الثعلب إذا نظر البازي عطفه كسّته شعاعاً على المنكب</p> <p>تُعَدّ هذه القصيدة وصفاً بديعاً لعين جميلة، وتُبرز مهارة الشاعر في استخدام التشبيهات والامتعارات. سألح كل بيت على حدة:</p> <p>البيت الأول: أيا ما أحسبها مقلّة ولولا الملاحظة لم أعجِبش: يبدأ الشاعر بمدح جمال العين قائلًا "أيا ما أحسبها مقلّة"، مُعرباً عن إعجابه الشديد بها. ثم يُضيف "ولو لا الملاحظة لم أعجب"، وهو يُشير إلى أن جمال العين ليس جمالاً عاديًا، بل هو جمال مُتقن ودقيق، يحتاج إلى ملاحظة مُدققة لتقدير جماله الكامل. فالجمال هنا ليس ظاهريًا فقط، بل فيه دقة وتفصيل تستحق التأمل.</p> <p>البيت الثاني: خلوقيّة في خلوقيها سؤداء من عتب الثعلِبش: يُشبه الشاعر بمدح سواد بؤبؤ العين بـ"سؤداء من عتب الثعلب"، وهو تشبيه بديع يُبرز سوادها الداكن وجمال لونه. كلمة "خلوقيّة" تعني حسنة الخلق، وتكريرها مع "في خلوقيها" يُبرز جمال خلق العين وجمال ما تحتويه من تفاصيل. يستخدم الشاعر عتب الثعلب -الذي يُعرف بلونه الأسود الداكن- لتوضيح عمق سواد بؤبؤ العين.</p> <p>البيت الثالث: إذا نظر البازي في عطفه كسّته شعاعاً على</p>	<p>(مدح Madh Praise)</p> <p>غَرَا نَمِرٌ وَقَادَ بَنِي تَمِيمٍ وَمَرَّ لَهُ الْيَأْمَانُ بِالسُّغُودِ فَقَفَّ الْغُلَّ عَنْ نَمِيمٍ بِنِ فُنُبٍ وَتَمِيمٌ فِي السَّلَاسِلِ وَالْفُيُودِ</p> <p>هذه القصيدة تتحدث عن غزو قام به نمر وقادته لبني تميم، وتصف تحريره لثَم من قنب من الأسر. سنشرح كل بيت على حدة:</p> <p>البيت الأول: غَرَا نَمِرٌ وَقَادَ بَنِي تَمِيمٍ وَمَرَّ لَهُ الْيَأْمَانُ بِالسُّغُودِ غَرَا: فعل "غزا" يدل على شن هجوم عسكري، ونمر هو اسم قائد الغزو. وقادَ بَنِي تَمِيمٍ: يبين هذا الشطر أن نمر قاد قبيلة بني تميم في غزوها. "قَادَ" تعني أدار وأمر بجمع القوات. وَمَرَّ لَهُ الْيَأْمَانُ بِالسُّغُودِ: هذا الشطر يصف مرور الأيام على نمر بالسُّغُود أي بالنجاح والتوفيق والانتصارات. "اليأمان" جمع "يوم" بمعنى الأيام المباركة أو الناجحة.</p> <p>البيت الثاني: فَقَفَّ الْغُلَّ عَنْ نَمِيمٍ بِنِ فُنُبٍ وَتَمِيمٌ فِي السَّلَاسِلِ وَالْفُيُودِ فَقَفَّ الْغُلَّ: الفعل "قَفَّ" يعني "توقف" أو "تراجع". "الغُلَّ" يُشير إلى القيد أو الأصفاد التي كانت تربطه. وَتَمِيمٌ فِي السَّلَاسِلِ وَالْفُيُودِ: هذا الشطر يوضح حالة تيم بن قنب قبل تحريره، حيث كان أسيرًا مقيدًا بالسلاسل والأغلال. "السلاسل" و"الفيود" كلمات مترادفة تُبرز شدة أسرِهِ واختصاره. القصيدة تصور بطولة نمر وقوة بني تميم في غزوه، وتبرز شجاعته بتحريرهم لأسير مهم كتيم بن قنب.</p>	<p>(حكمة Hikma Wisdom)</p> <p>أَعْنِ إِذْنِي تَهَبَّ الرِّيحَ رَهْوَ وَيَسِرِّي كَلِّمًا شَيْئًا الْغَمَامُ وَلَكِنَّ الْغَمَامَ لَهُ طِبَاعٌ تَجَسَّسُهُ بِهَا وَكَذَا الْكِرَامُ</p> <p>"القصيدة تتحدث عن قدرة الشاعر على التحكم في الطبيعة، أو على الأقل تصوره لذلك التحكم، وتُشرح حقيقة أن هذه القدرة محدودة بطبيعة الأشياء. سنشرح كل بيت على حدة:</p> <p>البيت الأول: أَعْنِ إِذْنِي تَهَبَّ الرِّيحَ رَهْوَ وَيَسِرِّي كَلِّمًا شَيْئًا الْغَمَامُ: يعني "إذني" أو "إبرادتي". يُعبر الشاعر عن قدرته على التحكم في الريح. تَهَبَّ الرِّيحَ رَهْوَ: الريح تهب بلطف وهدوء، مُشيرًا إلى درجة التحكم التي يمتلكها الشاعر. "رهوًا" تصف هدوء الريح وسلاسة حركتها. وَيَسِرِّي كَلِّمًا شَيْئًا الْغَمَامُ: وَيَسِرِّي الْغَمَامُ: السحاب) حيثما أضاء يُضفي الشاعر قدرة على التحكم في السحاب بالإضافة إلى الريح. "كَلِّمًا شَيْئًا" تعني "كلما أردت".</p> <p>البيت الثاني: وَلَكِنَّ الْغَمَامَ لَهُ طِبَاعٌ تَجَسَّسُهُ بِهَا وَكَذَا الْكِرَامُوَلَكِنَّ الْغَمَامَ لَهُ طِبَاعٌ: لكن الغمام له طبيعة الخاصة، أي أنه ليس خاضعًا تمامًا لإرادة الشاعر. تَجَسَّسُهُ بِهَا: تَجَسَّسَ: التفتاخر أو التباهي. هنا يُشير إلى أن الغمام يتفاخر بطبيعته، بمعنى أنه لا يخضع لإرادة الشاعر بشكل مطلق بل له قوانينه وخصائصه. وَكَذَا الْكِرَامُ: وكذا الكرام، أي الأشخاص الكرام ذوي المكانة العالية، لهم أيضًا طبيعتهم التي لا يُمكن إخضاعها بسهولة. يُشبه الشاعر</p>
<p>(رثاء Rithaa Elegy)</p> <p>جَزَعْتُ وَلَكِنْ مَا يَرُدُّ لِي الْجَزْءُ وَأَعُولْتُ لَوْ أَعْنِي الْعَوِيلُ وَلَوْ نَفَعَ بَا سَاكِي الْأَجْدَاتِ هَلْ لِي إِلَيْكُم عَلَى فَرِيكُم مَنَى الدَّهْرِ مُطْلَعُ فَوَائِهِ مَا أَبْقَى لِي الدَّهْرِ مَنَكُم حَبِيبًا وَلَا ذُخْرًا لَعَمْرِي وَلَا وَدَعَ</p> <p>تتحدث هذه القصيدة عن الشاعر وحزنه الشديد لفقد أحبابه، وتعبّر عن يأس عميق وفنوط من الزمن الذي سلب منه كل ما يملك من عزيز. سنشرح كل بيت على حدة:</p> <p>البيت الأول: جَزَعْتُ وَلَكِنْ مَا يَرُدُّ لِي الْجَزْءُ وَأَعُولْتُ لَوْ أَعْنِي الْعَوِيلُ وَلَوْ نَفَعَ: يبدأ الشاعر ببيان حالته النفسية: لقد جَزَعْتُ (أي حزنْتُ بشدة) ولكن الجزء لا يَرُدُّ لِي ما فقدته، وَأَعُولْتُ (أي بكيت بصوت عالٍ) ولكن العويل لا يُفيد ولا يُغيّر من الواقع المرّ. يعبر هذا البيت عن عجز الشاعر أمام مصيبتِهِ وقوة حزنه الذي يتجاوز أي علاج أو تسكين.</p> <p>البيت الثاني: بَا سَاكِي الْأَجْدَاتِ هَلْ لِي إِلَيْكُم عَلَى فَرِيكُم مَنَى الدَّهْرِ مُطْلَعُ: الدَّهْرِ مُطْلَعُ: يخاطب الشاعر هنا أحبابه الموتى (ساكِي الْأَجْدَاتِ)، ويسألهم عن إمكانية الوصول إليهم رغم بعدهم المادي. "على فريكم منى الدهر مُطلع" يعني منى الدهر مُطلع" يُشير إلى قُرْبهم الزماني منه في حياتهم، لكن الموت فصلٌ بينهم يفواصل لا تُفسي.</p> <p>البيت الثالث: فَوَائِهِ مَا أَبْقَى لِي الدَّهْرِ مَنَكُم حَبِيبًا وَلَا ذُخْرًا لَعَمْرِي وَلَا وَدَعَ: يقسم الشاعر بالله أنه لم يبقَ له من أحبابه شيء، لا حبيب يُواسيه، ولا ذُرَى (أي مال أو ثروة)، ولا ودّ (أي مودة وصداقة) يبقى من ذكرياتهم. يُبرز هذا البيت فقدان الشاعر لكل ما كان يُعزّزه ويُسعدُه في حياته.</p>	<p>(هجاء Hijaa Satire)</p> <p>إِرْفُضْ إِخْوَةَ مِنْ نَسَكٍ وَالزَّمْ مَوْدَةَ مِنْ فَتَكِ وَاهْضُ بِأَيْرِكَ مُنْعَطًا وَبِهِ فَطُوفٌ فِي السَّكَكِ فَإِذَا لَقِيتَ مُهْزِقًا أَحْوَى رَحِيمًا قَدْ تَسَكَّ</p> <p>هذه القصيدة تتم عن روح ثائرة رافضة للطلم والجيروت، وتدعو إلى التصدي له بقوة وحزم. سألح كل بيت على حدة:</p> <p>البيت الأول: إِرْفُضْ إِخْوَةَ مِنْ نَسَكٍ وَالزَّمْ مَوْدَةَ مِنْ فَتَكِ: الشاعر بتوجيه نصيحة قوية، فهو يرفض الأخوة مع من يُظهر النسك (أي التقوى الظاهرية والزهد المُصنّع) وذلك لأنّ هذا النسك قد يكون ستارًا يخفي الشرّ والبداع. ويُؤكد على أهمية الرمودة (المودة والصداقة) مع من هو قوَى وشجاع (الفتك). فهذا القوَى قد يكون حليفًا أفضل في مواجهة الطلم.</p> <p>البيت الثاني: وَاهْضُ بِأَيْرِكَ مُنْعَطًا وَبِهِ فَطُوفٌ فِي السَّكَكِ: يحثّ الشاعر على النهوض والانتلاق بقوة وشجاعة، باستخدام "أيرك" (أي سلاحك أو قوتك) بشكل واضح لا لبس فيه ("منعطًا" أي جليًا مظهرًا). ويُحثّ على التحول به إلى أعلى المراتب، أي على إظهار القوة وعدم الخوف من مواجهة الأعداء.</p> <p>البيت الثالث: فَإِذَا لَقِيتَ مُهْزِقًا أَحْوَى رَحِيمًا قَدْ تَسَكَّ: يصف الشاعر هنا نوعًا من الأعداء، فهم "مهزقون" (أي متظاهرون بالضعف والخضوع)، "أحواء" (أي ليون في ظاهريهم) و"رحيمون" (أي رقيقون في كالمهم)، لكنهم في الحقيقة "قد تَسَكَّوْا" (أي تظاهروا بالتقوى والزهد). وولاء هم أشدّ خطرًا من الأعداء الصريحين.</p>	<p>(مدح Madh Praise)</p> <p>على قدو أهل العزم تأتي الغزائم وتأتي على قدو الكرام المكارم وتُعظم في عين الصغبر صباها وتَصغُرُ في عين العظيم العظائم يُكَلِّفُ سَيْفُ الدُّوَلَةِ الْجَيْشَ هُمَةً وقد عَجَزَتْ عَنْهُ الْجِيُوشُ الْغُضَارُ</p> <p>"هذه تمجد القصيدة شجاعة سيف الدولة وقيادته الفذة، وتصور انتصاراته كملحمة بطولية تخدم المجد والدين. سنشرح كل بيت على حدة:</p> <p>البيت الأول: على قدر أهل العزم تأتي الغزائم وتأتي على قدر الكرام المكارم: يقرر الشاعر هنا علاقة طردية بين العزم والإرادة وبين تحقيق الأهداف. فالعزائم الكبيرة تأتي من أصحاب العزم الكبير. والمكارم العظيمة تصدر من الكرام ذوي النفوس الشريفة. العزم هنا هو السبب، والنجاح هو المسبب له.</p> <p>البيت الثاني: وتُعظم في عين الصغبر صباها وتَصغُرُ في عين العظيم العظائم: يوضح الشاعر هنا أنّ تقدير الأمور يختلف باختلاف الأشخاص. فالأشياء الصغيرة تبدو عظيمة في نظر الشخص الضليل، بينما تبدو الأشياء الكبيرة صغيرة في نظر الشخص العظيم ذي الطموحات الكبيرة والأهداف السامية. هذا يدل على ضيق أفق الصغير مقابل اتساع أفق العظيم.</p> <p>البيت الثالث: يُكَلِّفُ سَيْفُ الدُّوَلَةِ الْجَيْشَ هُمَةً وقد عَجَزَتْ عَنْهُ الْجِيُوشُ الْغُضَارُ: يشير الشاعر إلى قوة جيش "سيف الدولة" وقدرته على تحمل المهام الصعبة التي عجزت عنها جيوش أخرى ذات خبرة واسعة ("الغُضَارُ" تعني المحضرة في الحروب). يُبرز هذا البيت قدرة القيادة الفعالة على تحقيق ما عجز عنه غيرها.</p>
<p>(شوق Shawk Longing)</p> <p>تُفَرِّدُ دُمُوعِي بِشَوْقِي إِلَيْكَ وَيَشْدُدُّ قَلْبِي بِطَوْلِ الْكَزْبِ وَإِنِّي لَكَجُنْدٍ فِي الْجُبُودِ وَلَكِنَّ نَفْسِي تَأْتِي الْكَذِبِ وَإِنِّي عَلَيَّكَ لِجَارِي الدُّمُوعِ وَإِنِّي عَلَيَّكَ لَمَصَبٌ وَصِيبٌ</p> <p>في هذه القصيدة يُظهر الشاعر صراعًا بين محاولته إخفاء الحب وواقع مشاعره التي تفضيحه دموعه وقليه المثلث بالهم. فرغم أجهاده في الإنكار، إلا أن نفسه ترفض الزيف، وتفضحه شدة العشق وبكاؤه المستمر. سنشرح كل بيت على حدة:</p> <p>البيت الأول: تُفَرِّدُ دُمُوعِي بِشَوْقِي إِلَيْكَ وَيَشْدُدُّ قَلْبِي بِطَوْلِ الْكَزْبِ: يُعلن الشاعر عن شوقه من خلال دموعه التي لا يستطيع إخفاءها، وقليه يشهد على معاناته الطويلة. هذا البيت يُعبّر عن صدق المشاعر التي تنكشف رغم الصمت.</p> <p>البيت الثاني: وَإِنِّي لَكَجُنْدٍ فِي الْجُبُودِ وَلَكِنَّ نَفْسِي تَأْتِي الْكَذِبِ: يحاول الشاعر إنكار حبه وجهده في إخفاء مشاعره، لكن نفسه الصادقة ترفض الكذب، فتفضح حبه رغمًا عنه. البيت يُبرز صراعًا داخليًا بين العقل والوجدان.</p> <p>البيت الثالث: وَإِنِّي عَلَيَّكَ لِجَارِي الدُّمُوعِ وَإِنِّي عَلَيَّكَ لَمَصَبٌ وَصِيبٌ: يشير الشاعر إلى كثرة بكائه على الحبيب، فهو دائم البكاء (جاري الدموع)، و"مَصَبٌ" (يصب) و"صِيبٌ" (يصب) في العشق (صَبَّ وصِيبٌ)، مما يعكس شدة الحب والجزن.</p>	<p>(عتاب Itab Reproach)</p> <p>ألا ما لِسَيْفِ الدُّوَلَةِ التُّومُ عَاتِبَا فدأ الزور أمضى الشيوف فضايربا ومال إذا ما اشتقت أبصرت دونه تانيف لا أشتافها وسبايسبا وقد كان يدي تجلسي من سمانه أحداث فيها بذرما والكواكبا</p> <p>تُعبر الأبيات عن شوق الشاعر لمجالس العز والمكانة التي كان ينعم بها في ظل سيف الدولة، حيث يعصفه بأمر من السيوف ويعاتب غيابه. ويظهر الجنين إلى الأماكن التي كانت تجمعهم به، ويتجلى فيها الجمال والنجوم. سألح كل بيت على حدة:</p> <p>البيت الأول: بدأ أشجع الشاعر بعتاب موجه إلى سيف الدولة، مُستأنًا عن سبب غيابه أو تفكره. ثم يمتدحه بكونه لناس وأمضى السيوف في المعارك، مما يُظهر الاحترام رغم العتاب.</p> <p>البيت الثاني: يعبر الشاعر عن اشتياقه لسيف الدولة، لكنه لا يجد، بل يرى أماكن فاحلة لا تثير الشوق، مما يُبرز الفجوة والحرمان التي يشعر بها بسبب غياب القرب والوصال.</p> <p>البيت الثالث: يتذكر الشاعر ماضيه الجميل من كان سيف الدولة يُقرّنه منه، ويجعله جليشًا في مجلسه الراقي الذي يُشبه فيه الحضور بالندر والكواكب، في إشارة إلى عظمة المكان والناس فيه.</p>	<p>(سياسي Siyasi Political)</p> <p>أَبْعَدُ الْحَارِثُ الْمَلِكُ ابْنَ عَمْرٍو لَهُ مُلْكُ الْعِرَاقِ إِلَى عُمَانَ مُجَاوِزَةُ بَنِي شَمْسٍ بِنِ جَزْمٍ هُوَانَا مَا أَبْيَعُ مِنْ الْهَوَانِ وَيَنْعَمُهَا بَنُو شَمْسٍ بِنِ جَزْمٍ مَعْرِفُومُ حَنَانِكَ ذَا الْهَنَانِ</p> <p>يعبر الشاعر عن استغرابه من سقوط ملكٍ جليل كالحارث بن عمرو، الذي امتد نفوذه من العراق إلى عُمان، ثم انتهى به المطاف إلى الدل والهوان. فقد غدا جازًا لقبيلة صغيرة منته من استعادة ملكه، في تصوير لانهيار شديد بعد عز، وختام مجازي يبرز مرارة الموقف. سنشرح كل بيت على حدة:</p> <p>البيت الأول: يعبّر الشاعر عن دهشته من أن يُهزم الملك الحارث بن عمرو، الذي كان يملك بلادًا واسعة تمتد من العراق إلى عُمان. يستنكر كيف ليك هذه القوة والهيبة أن يُصاب بالهوان.</p> <p>البيت الثاني: يوضح أن البريمة جاءت على يد قبيلة بني شَمْسٍ بن جرم، ويصفها بأنها من أشد أشكال المذلة التي قد تُصيب ملكًا، في تصوير لانهيار شديد بعد العزة والسلطان.</p> <p>البيت الثالث: يشير إلى أن القبيلة المنتصرة تواصل منع الحارث من استعادة ملكه، ويفخر الشاعر بهذا النصر. يختم بنداء "حنانك ذا الجنان"، إيا كدعاء أو تعبير مجازي عن سطوة الخصم واستمرار هيمنتِهِ.</p>

Figure 10: Fann or Flop Samples by Genre. Additional representative examples from the Fann or Flop benchmark, illustrating the diversity of genres covered, including Love (Ghazal), Praise (Madh), Wisdom (Hikma), Satire (Hijaa'), Elegy (Rithaa'), Reproach ('Itab), Political Poetry, and Longing (Shawq). Each example showcases a poetic excerpt alongside an interpretive breakdown highlighting figurative language, rhetorical devices, and thematic nuances. These curated samples reflect the benchmark's aim to evaluate models' nuanced understanding of Arabic poetic tradition.