

PaINLP at AR-MS Shared Task: Guidelines Paper for Arabic Manuscript Transcription

Mutaz Ayesh

Cardiff University
AyeshMA@cardiff.ac.uk

Abstract

This paper describes the guidelines that the PaINLP team recursively developed and followed during the transcription of the assigned batch, as part of the AR-MS Shared Task. The team, which consists of a single experienced transcriber, has manually transcribed 500 images of lines from the Omar Al-Saleh Memoir Collection.

Keywords: transcription guidelines, OCR, Arabic transcription

1. Introduction

This document presents the transcription guidelines that were synthesized and followed as part of the Arabic Manuscript Understanding Shared Task organized by Nakba-NLP 2026 (Hamoud et al., 2026). The team, which consists of a single experienced annotator, received a subset of the Omar Al-Saleh Memoir Collection to transcribe. The aim of the task was to support Arabic optical character recognition (OCR), and specifically handwritten text recognition (HTR).

The subset contained 500 images of lines, spanning 13 years of mostly handwritten material (1953-1965). Few sentences from the subset were typed. Images were line-level crops, not full pages. The team, however, also received the full context of these lines, which proved useful.

The methodology can be described as one of strict alignment with the handwritten lines, especially on the character level. Orthographic details were preserved and normalization or editorial or linguistic corrections were avoided whenever possible. This strict alignment also encompassed preserving incomplete or truncated words at the beginning or end of the line-based images.

Metric	Value
Total Lines / Segments	500
Total Characters	38,252
Total Words	6,757
Vocabulary Size (Unique Words)	4,118
Average Words per Line	13.51
Edge Cases	55

Table 1: Summary statistics of the batch.

2. Guidelines and Examples

There were four main guidelines that guided the transcription process, which were provided by the organizers: (1) transcribing the text exactly as it

appears; (2) preserving the original spelling and character forms of the author; (3) preserving punctuation marks, numbers, and diacritics in their original forms; and (4) providing the best interpretation for unclear, damaged, or illegible text. Multiple instructions were extracted accordingly.

Footnotes indicate the identifiers corresponding to each line crop, usually formatted as year_page_line.

2.1. Arabic-Indic Digits

Arabic-Indic digits were preserved as in the source images, with no normalization to Western numerals (0–9), as shown in Table 2.

Transcribed line	Notes
١٤ ¹ صفر	"14 Safar".
١٣٨٢ ² هجرية	"1382 Hijri".

Table 2: Examples illustrating the preservation of Arabic-Indic digits in the transcription.

2.2. Typos

Table 3 provides examples of writing errors that were preserved in the transcription.

Transcribed line	Notes
(٣) دراسة النواحي العسكرية في الليجان المختصة ونضع مخططاً يستند على هذه الدعائم ³	داسة instead of دراسة due to the <i>alef</i> clearly preceding the <i>ra'</i> .
الحراسة الليلية والنهارية مدر العام ولا يسمح لهؤلاء الافراد المدربين ان يغادروا البلاد الا ⁴	The unnecessary <i>madda</i> diacritic in الحراسة; مدر instead of مدار; and هؤلاء instead of هؤلاء.

Table 3: Instances of uncorrected typos.

¹ID: 732294a1c8ef4d55ad023c877c098e26-0004-37

²ID: 90aa2437fec34a18bf1fad0830294b33-0017-08

³ID: 1962_p086_l0022

⁴ID: 1959_p134_l0017

2.3. Incomplete Words

Incomplete words, whether in the beginning or end of the line were transcribed. When necessary, the *tatwīl* character was used to ensure that the truncated word matches the handwritten text, as shown in Table 4.

Transcribed line	Notes
قرأت الجرائد المحلية وتبعت الاحداث والحوادث تستغرب ولا تتكاد تصدق فر ⁵	The incomplete word here is فر at the end of the line.
اني اكزه التدخين وارى فيه الضرر وانما اخفف واقلل واتيته ع ⁶	The <i>tatwīl</i> character was used due to the word-initial <i>ayn</i> .

Table 4: Examples showing the transcription of incomplete words.

2.4. Preservation of Original Orthography

The letter *ع* was transcribed in its original form even in cases where the standard orthography would require the dotted counterpart *ي*. This includes word-final instances typically written with *ي*, such as *لأنني*, *شؤني*, and *في*. Two representative examples are shown in Table 5.

Transcribed line	Notes
احضر الى ابن سعود العشوة صندوق سيجار وذهبت ليلا الى الكعبيني لقضاء مصلحة ابو كحيل 7 واحضر لي مشهور	إلى is usually spelled as يلى.
وزارة لبنان: اشيع ان رشيد كرامي استقال بعد مساع كثيرة ولا يهمني هذا النبأ ⁸	The standard spelling of يهمني is إس يهمني.

Table 5: Orthographic adherence to the handwritten text, as exemplified by the *ع*.

2.5. Whitespaces

Whitespace was preserved exactly as it appeared in the source images. No normalization or correction of spacing was applied, including instances where spaces occurred between words and punctuation marks. For example, the author occasionally inserted spaces before punctuation, which were retained in the transcription, as shown in Table 6. The second example contains extended spacing between two verses of a poem, which is a common formatting convention in Arabic poetry.

⁵ID: 1959_p151_l0015

⁶ID: 1962_p045_l0026

⁷ID: 1959_p158_l0079

⁸ID: 1958_p123_l0021

Transcribed line	Notes
اتنى هدية ازهار ملأت اماكنها واعدقت على النفس بهجة وراحة ⁹	A whitespace was inserted between the word and the full-stop to reflect the written text.
كما قالت تلك الشاعرة: عجبا عجبا له ابكيه ملاً مدامعي وأقول 10 لا شلت يمين القاتل.	Multiple white spaces were inserted in our transcription to reflect those found between the verses.

Table 6: Examples illustrating the preservation of whitespace in the transcription.

2.6. Diacritics

Diacritics were preserved exactly as they appeared in the source images, including marks such as *madda* and *shadda*. No normalization was applied to *hamzat* forms (e.g., *أ، إ، ؤ، ئ*) even in cases where their insertion or modification might have been grammatically expected, ensuring that the transcription reflects the text as written rather than a corrected version.

Transcribed line	Notes
ينظمون صفوفهم ويهيئون رجالهم لحرب صاحقة اخشى ان تؤدي 11 العقبي لضياح البقية	The <i>hamza</i> was preserved due to appearing in the original image.
ضباط الجيش سمو انفسهم "الضباط الاحرار" واصدروا نشرات وزعوها في اجواء 12 الجيش	The <i>hamza</i> was not introduced in any of the underlined words where it would be expected (انفسهم، اجواء، اصدروا، الاحرار).
عدوا منافساً الى السعودية في الخليج؟ انى ارى هذا من السياسة 13 الاستعمارية القديمة	This example shows a <i>shadda</i> and a <i>tanwīn</i> .

Table 7: Instances that required a Google search. The side search helped disambiguate the spelling.

3. Challenging Cases

3.1. Instances that Required a Side Search

Google Search was utilized exclusively for ambiguous handwritten text involving foreign words or named entities (see Table 8). To ensure the transcription remained faithful and free from external bias, search results were used solely to generate spelling "hypotheses" for unclear strokes, which

⁹ID: 1958_p087_l0005

¹⁰ID: 1962_p177_l0049

¹¹ID: 1955_p184_l0035

¹²ID: 1955_p176_l0036

¹³ID: 1956_p047_l0027

were then cross-referenced against the manuscript. A term was only finalized in the transcription if the visual evidence completely supported the external Google results. This ensured that the document itself dictated the final text rather than outside knowledge.

Transcribed line	Notes
اتى هزاع وزير الداخلية القائم الى رام الله فالتفت حوله ارباب المصالح ¹⁴ ودعا حنا العسوس.	حنا العسوس and هزاع are Palestinian figures.
حزب آخر الى مساعدته وقد سقط اساطين اعداء ديغول امثال رينو ومنديس فرانس ونغي موليه ¹⁵ وغيرهم	منديس فرانس, رينو, ديغول, and نغي موليه are French politicians.

Table 8: Instances that required a Google search. The side search helped disambiguate the spelling.

3.2. Incomprehensible Text

When a segment of text was difficult to decipher, it was transcribed exactly as it appeared in the line image without attempting correction or normalization. In such cases, the surrounding page context (provided by the organizers) was used to aid interpretation. While contextual information helped resolve ambiguous words (Example 2 in Table 9), it did not clarify all instances (Example 1).

Transcribed line	Notes
اذا حز به الاغياء وهزه ما زعزع مجرى حياته العادية انقلب كل ما فيه من مرح الى ترح ¹⁶	The second and fourth words are incomprehensible.
حكما لا وزن له ولا اهمية. كم اصبر وكم اتفاضى ولكن الصبر انطوى ¹⁷ واليد قصيرة والناس هائمون	The words كم اصبر وكم اتفاضى looked ambiguous due to the lack of space between them in the source image. This required reading the lines before and after it to decipher the handwritten text.

Table 9: Examples of incomprehensible handwriting that was transcribed to the closest visually identifiable reading.

To maintain a transparent uncertainty policy, any text that remained incomprehensible was transcribed to its closest visually identifiable reading and explicitly underlined and annotated with an “incomprehensible” comment. This approach ensures that future readers can easily distinguish between

¹⁴ID: 1955_p110_l0009

¹⁵ID: 1962_p175_l0019

¹⁶ID: 1959_p107_l0061

¹⁷ID: 1956_p127_l0051

confident transcriptions and visual approximations. PaINLP would be delighted to share the annotated spreadsheet with the organizers.

3.3. Multiple Lines in one Image

In a small number of cases, the provided line images contained text written across more than one visual line. When this occurred, the transcription preserved the original structure by inserting a new-line to reflect the layout of the source image. In addition, the author occasionally positioned numbers (e.g., dates or times) beneath the noun they modify. In such cases, the numbers were transcribed inline immediately after the relevant word. Examples of these cases are shown in Table 10.

Transcribed line	Notes
جثمت في الفراش انتجاعا للراحة وحباً في الأبعاد عن الناس وهذه خير حكمة ادرعها.	Line break preserved.
دفعت حساب الكهرباء ثلاثة ¹⁸ دنانير و ٣٧٠ فلسا.	
كل مصيبة اليمة. ففي دقيقة واحدة اختفى عدة الوف من من اليرانيين. عند الساعة ٩,٢٠ من ليل السبت ¹⁹	The number ٩,٢٠ was transcribed right after the word “الساعة” despite appearing below it in the source image.

Table 10: Examples of multi-line text.

4. Conclusion

This brief guidelines document aimed to bring context and structure to the PaINLP team’s submission for Subtask 1 of the Arabic Manuscript Understanding Shared Task. The transcription process generally aimed to reflect the handwritten text as much as possible so that OCR and HTR systems trained or evaluated on this batch can be assessed under realistic conditions that preserve orthographic variation, diacritics, ligatures, and other manuscript-specific elements.

Although the transcription was produced by a single annotator, consistency was supported through iterative self-review, revisiting earlier decisions when recurring handwriting patterns became clear, marking uncertain segments for later verification, and repeatedly checking batch-level agreement on CodaBench against the expert-verified reference.

¹⁸ID: 1958_p139_l0028

¹⁹ID: 1962_p136_l0010

References

Hadi Hamoud, Ahmad Ali Chamseddine, Bilal Shalash, Firas Ben Abid, Mustafa Jarrar, Chadi Abou Chakra, Bernard Ghanem, and Fadi A. Zaraket. 2026. NAKBA NLP 2026: Shared Task on Arabic Handwritten Manuscript Understanding (Palestine Memory–Omar Al-Saleh Memoir). In *Proceedings of the 2nd International Workshop on Nakba Narratives as Language Resources (Nakba-NLP 2026), co-located with the Language Resources and Evaluation Conference (LREC 2026)*, Palma, Mallorca, Spain.