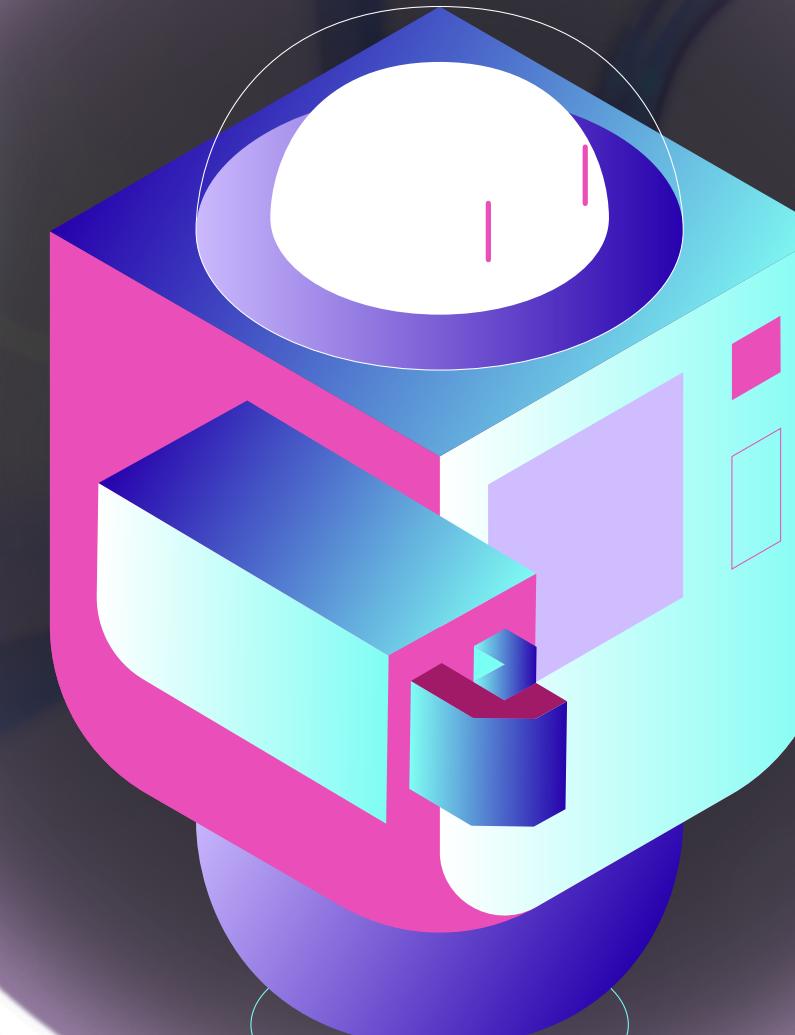
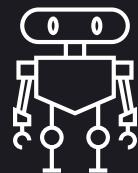


# ARABIC QUESTION ANSWERING BASED ON LLM

Arabic language, Question Answering,  
Natural Language Processing, AraElectra,  
Fine-Tuning, DuckDuckGo, NLP



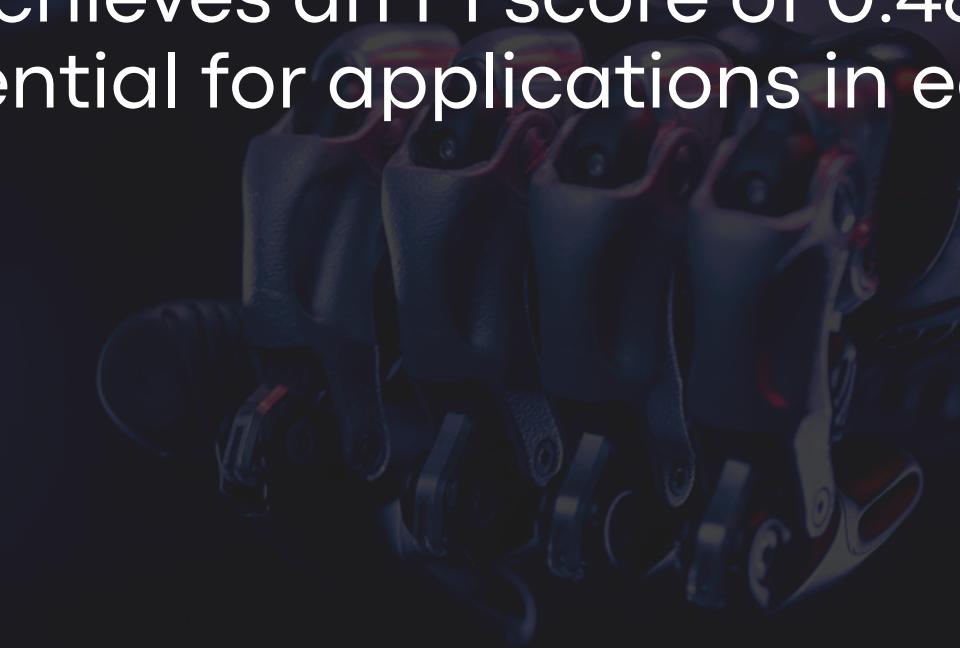
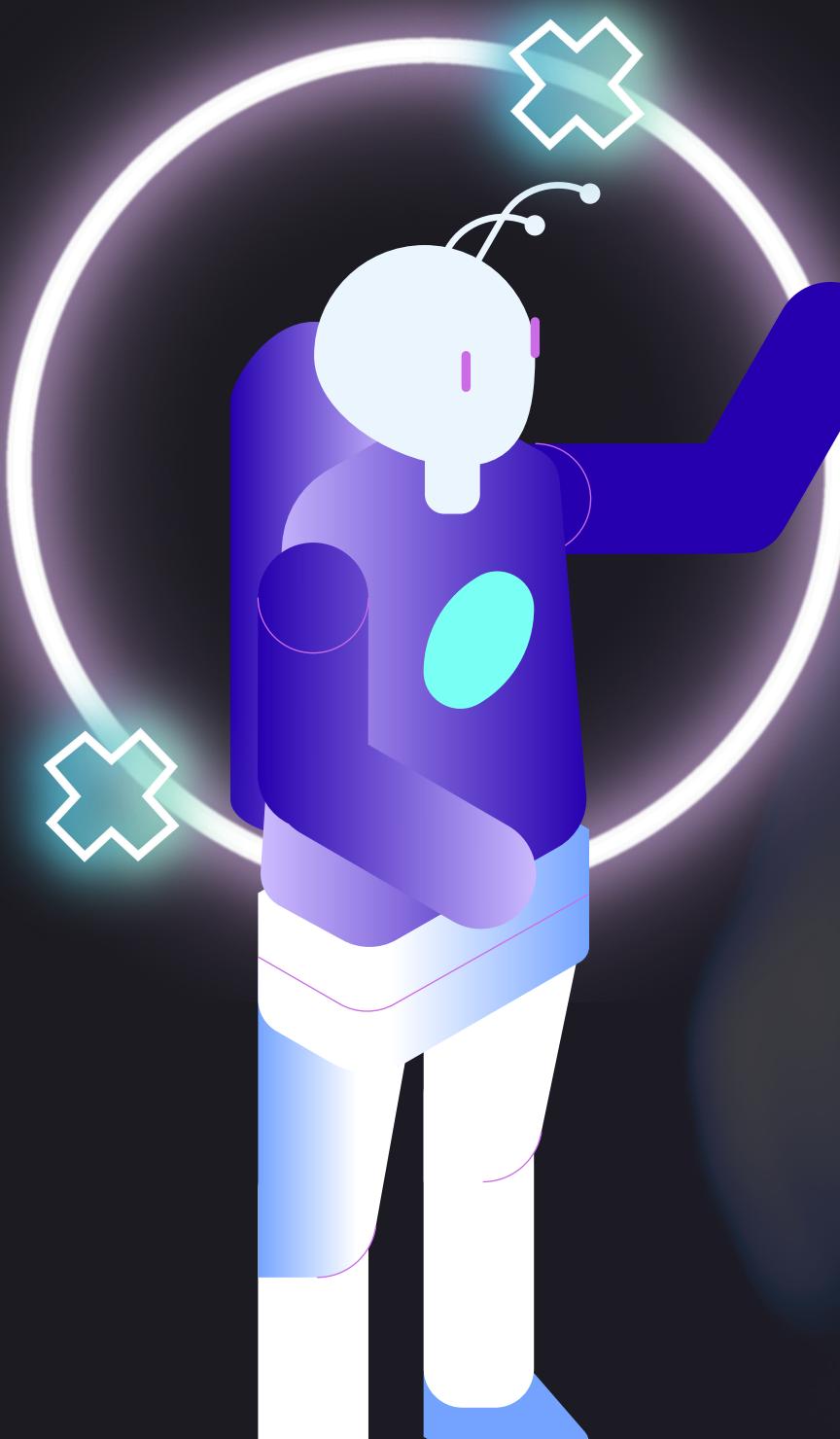


# DEVELOPMENT OF THE SYSTEM

This presentation outlines the development of an advanced Arabic Question Answering (QA) system designed to address the scarcity of robust QA tools for Arabic.

By fine-tuning the AraElectra model and integrating NLP techniques like Named Entity Recognition (NER), query reformulation, and web based retrieval via the DuckDuckGo API.

the system delivers contextually accurate answers. Evaluated on the AAQAD dataset, it achieves an F1 score of 0.4814 and accuracy of 0.6372, showing strong potential for applications in education, research, and customer service



# WHY ARABIC QA MATTERS

## CONTEXT

Arabic QA systems lag behind English due to linguistic complexities (e.g., morphology, diacritics) and limited annotated data

## GOAL

Develop a smart QA system to extract precise answers from web documents or user inputs, enhancing information access for Arabic speakers.

## APPLICATIONS

Education, research, customer service, leveraging growing Arabic online content.

## KEY FEATURES

Supports bilingual input (Arabic/English), multiple context modes (predefined, manual, web-based), and a user-friendly Gradio interface.

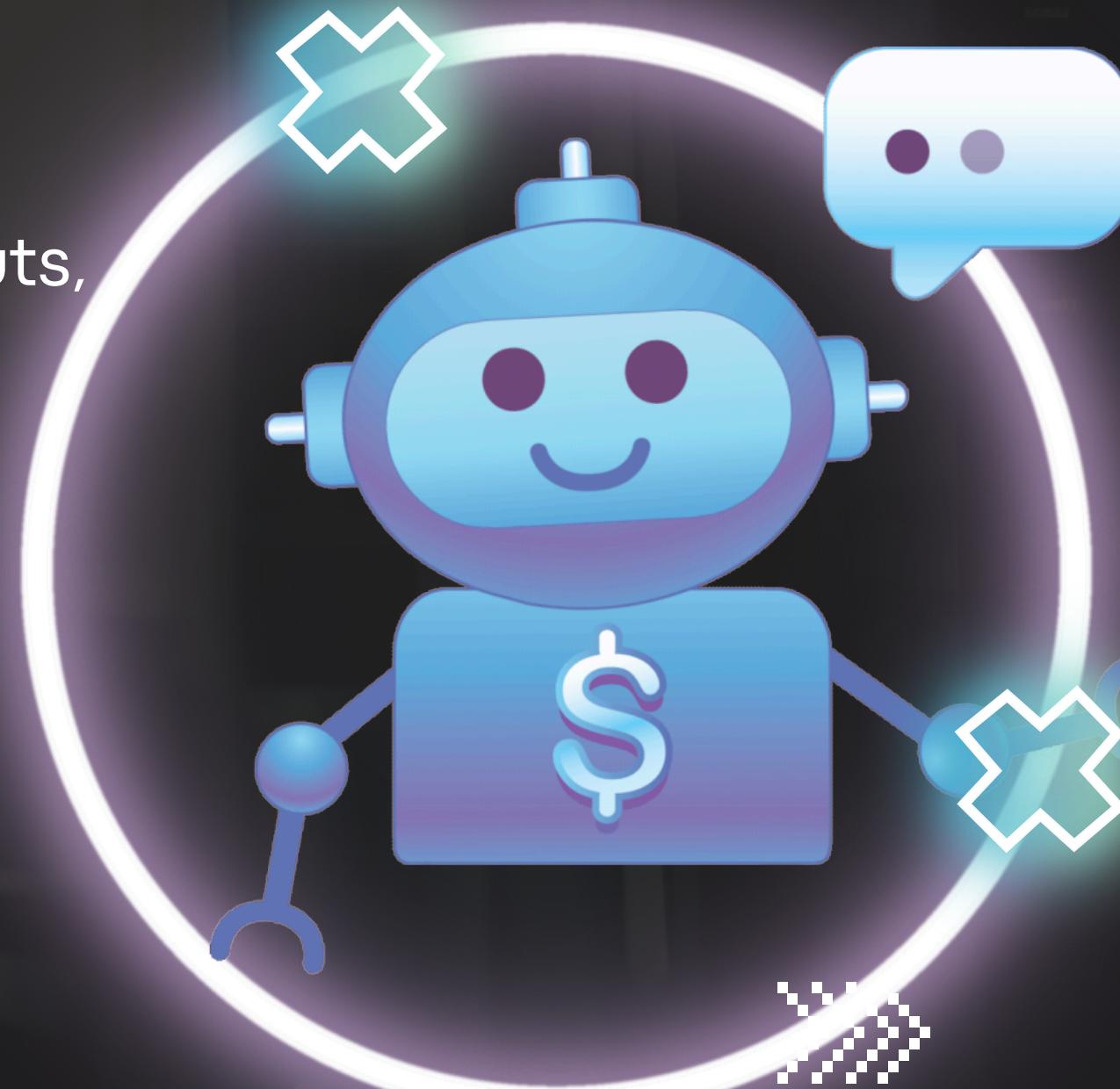


# HOW THE SYSTEM WORKS

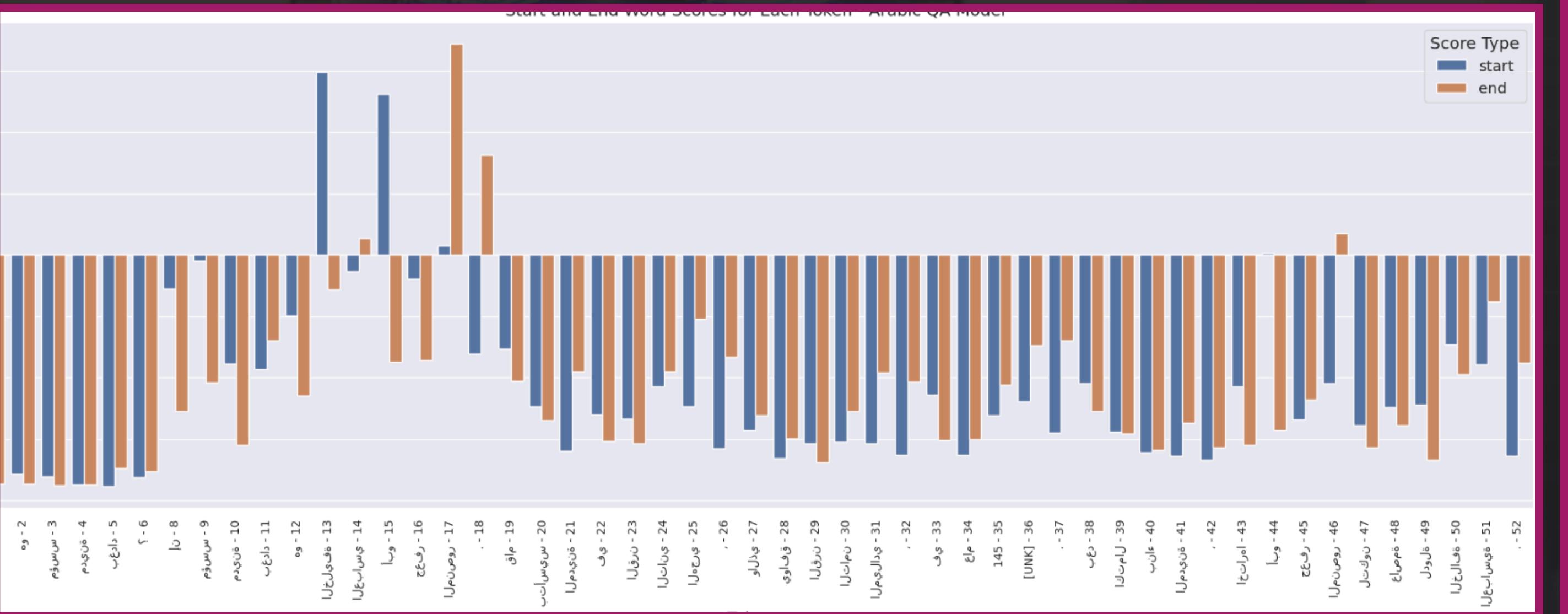
Query Processing: Preprocesses questions using tokenization, NER, and query reformulation to confirm user intent.

Context Retrieval: Supports predefined texts, user-provided inputs, or real-time web searches via DuckDuckGo API, ranked by relevance.

**AraElectra** : A transformer-based model optimized for Arabic, using discriminative pre-training for efficiency.



# TOKENIZER AND VISUALIZING SCORES



[CLS]	306 (Question)
583 (Question)	هو (Question)
10611 (Question)	مؤسس (Question)
1171 (Question)	مدينة (Question)
3112 (Question)	بغداد (Question)
105 (Question)	? (Question)
[SEP]	3 (Question)
476 (Context)	إن (Context)
10611 (Context)	مؤسس (Context)
1171 (Context)	مدينة (Context)
3112 (Context)	بغداد (Context)
583 (Context)	هو (Context)
12280 (Context)	الخليفة (Context)
29743 (Context)	العباسي (Context)
1195 (Context)	أبو (Context)
9501 (Context)	جعفر (Context)
17647 (Context)	المنصور (Context)
20 (Context)	.
1178 (Context)	قام (Context)
22077 (Context)	باتيس (Context)
1665 (Context)	المدينة (Context)
305 (Context)	في (Context)
2890 (Context)	القرن (Context)
1161 (Context)	الثاني (Context)
27231 (Context)	الهجري (Context)
' (Context)	103 (Context)
1619 (Context)	والذى (Context)
12423 (Context)	يواقب (Context)
2890 (Context)	القرن (Context)
6687 (Context)	الثامن (Context)
22923 (Context)	الميلادي (Context)
' (Context)	103 (Context)
305 (Context)	في (Context)
515 (Context)	علم (Context)
145 (Context)	23045 (Context)
[UNK] (Context)	1 (Context)
' (Context)	20 (Context)
446 (Context)	بعد (Context)
14166 (Context)	اكمال (Context)
1745 (Context)	بناء (Context)

# HOW THE SYSTEM WORKS CONT...

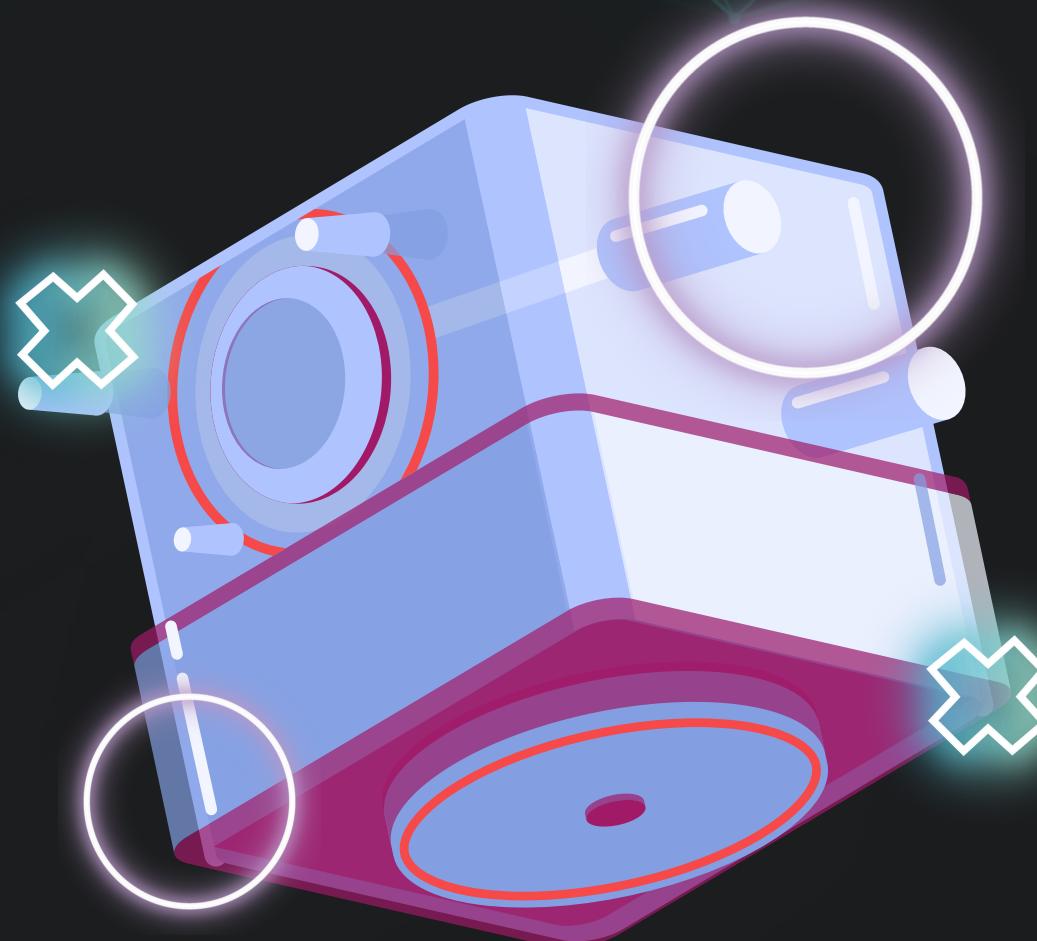
Answer Generation: Uses fine-tuned AraElectra model with prompt engineering (four strategies: direct, formal specificity, detailed precision, analytical completeness) for accurate answer spans.

User Interface: Gradio-based, bilingual, interactive platform for input and answer display.



# BUILDING THE ARABIC QA SYSTEM

Query Processing: Handles Arabic linguistic challenges (diacritics, morphology) via AraElectra's tokenizer and NER for entity-focused answers.



Context Retrieval: Uses DuckDuckGo API for web searches, ranking documents by keyword overlap and semantic similarity.

Answer Generation: Fine-tuned AraElectra model predicts answer spans, enhanced by prompt engineering and rule-based post-processing for grammar and cultural fit.

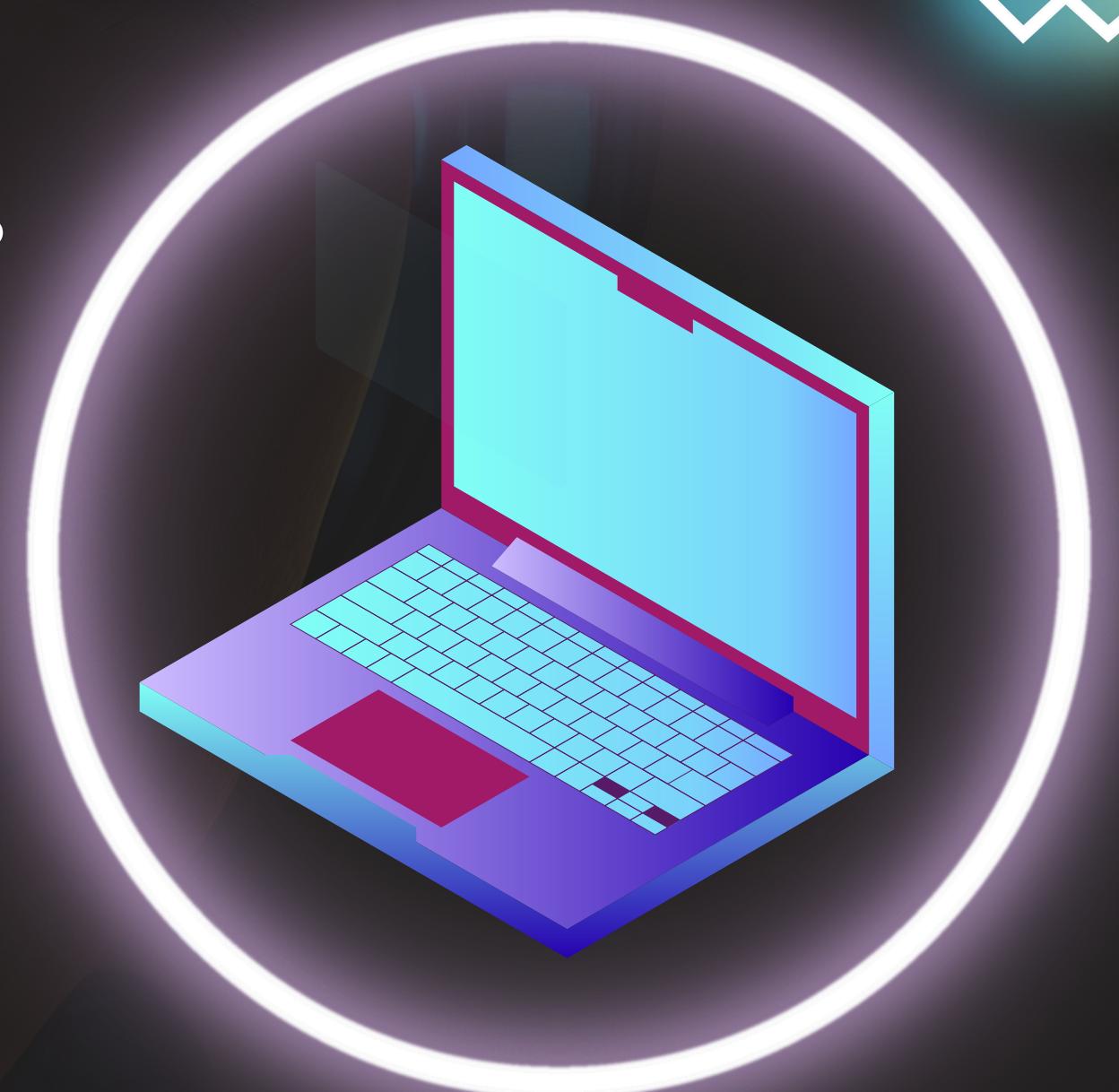
# IMPLEMENTATION

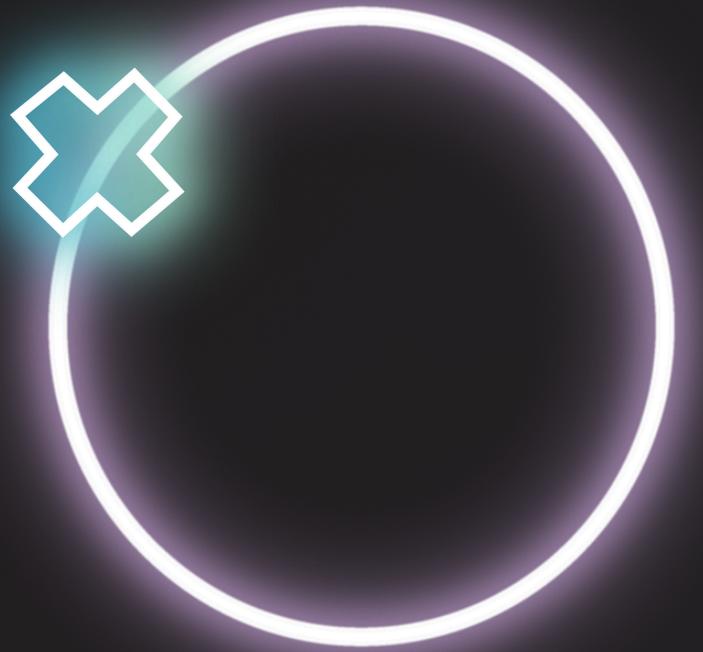
Tools: Python, Hugging Face Transformers, Gradio, DuckDuckGo API.

Dataset: AAQAD, split 80% training, 10% validation, 10% test.

Fine-Tuning: 5 epochs, AdamW optimizer, 2e-5 learning rate, cross-entropy loss.

Data Augmentation: Synonym replacement, back-translation for robustness.





# PERFORMANCE AND INSIGHTS

Dataset: AAQAD-test, Arabic question-answer pairs.

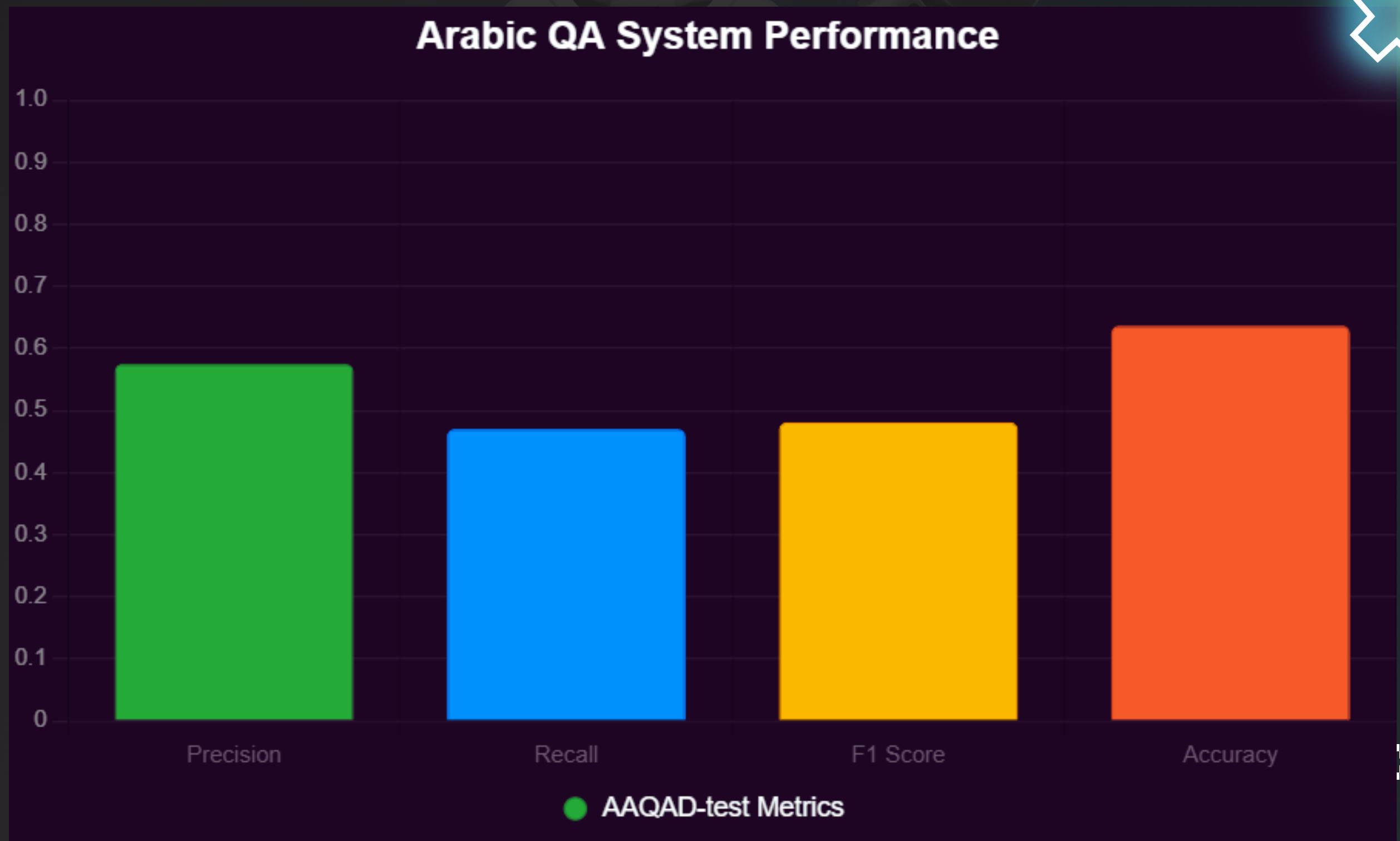
## METRICS

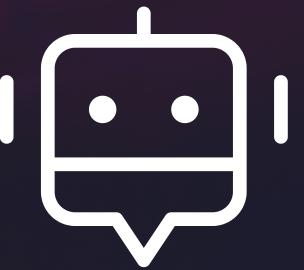
PRECISION: 0.5754

RECALL: 0.4708

F1 SCORE: 0.4814

ACCURACY: 0.6372





# SAMPLE PERFORMANCE

Factoid question (e.g., “What is the capital of Iraq?”):  
Perfect scores ( $F1=1.0$ )

Complex questions (e.g., space exploration): Mixed results,  
e.g., missed “Blue Origin” for reusable rockets, vague  
Artemis program answers.

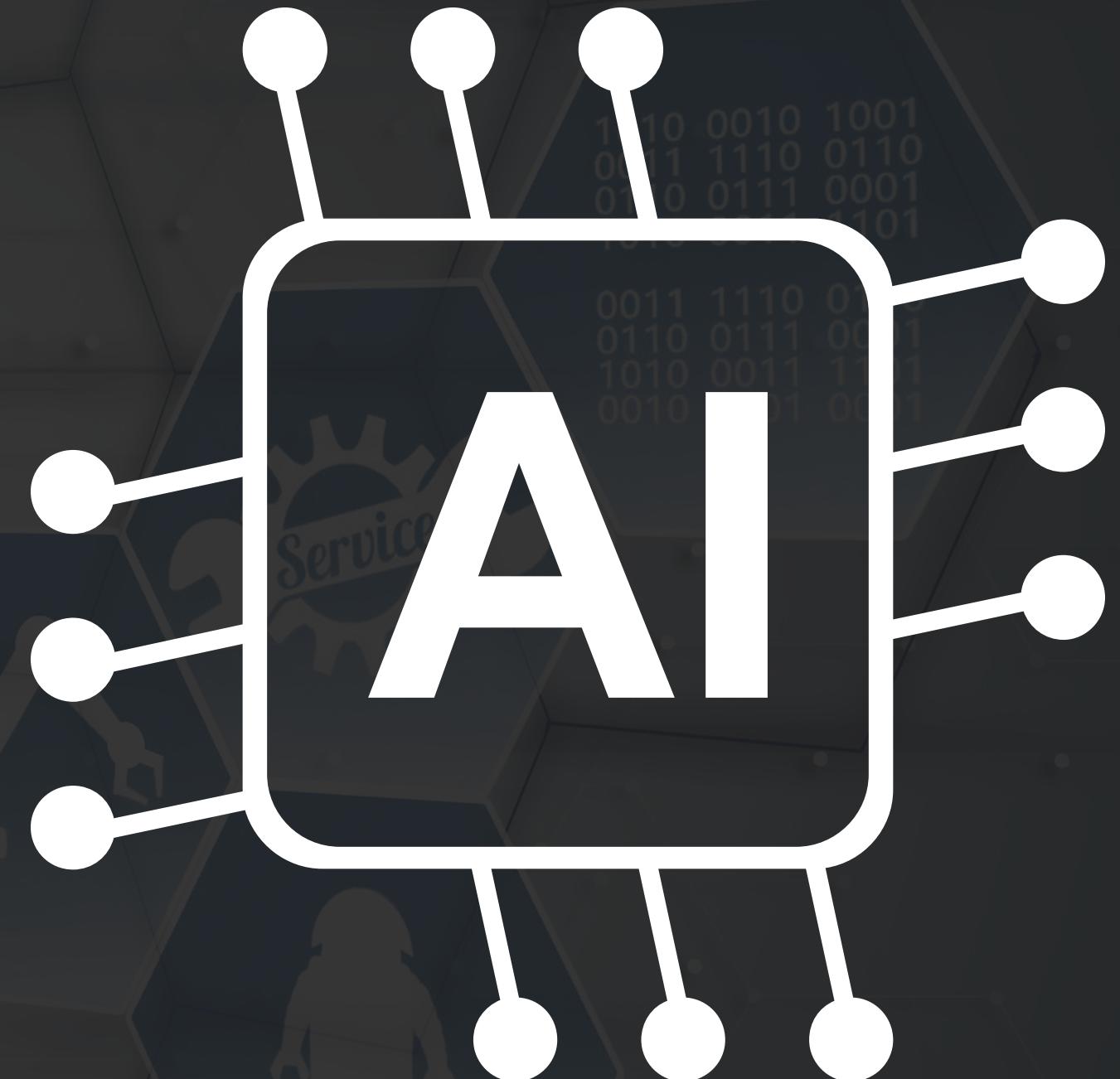


# TESTING

Unit Testing: Validated tokenizer, answer generation

Integration Testing: Ensured seamless query-to-answer pipeline.

User Acceptance: Confirmed usability via Gradio interface, strong for factoid queries, weaker for inferential ones.



# EXAMPLES

## Technology Context:

في القرن العشرين، تطورت تكنولوجيا المعلومات بشكل مذهل. تم تطوير الكمبيوترات والإنترنت وأجهزة الهواتف المحمولة وغيرها من الأجهزة الإلكترونية التي أصبحت جزءاً أساسياً من حياتنا اليومية. بفضل هذه التكنولوجيا، أصبحت الاتصالات أسرع وأسهل، وزادت سرعة نقل المعلومات والبيانات. كما تطورت البرامج والتطبيقات التي توفر العديد من الخدمات والوظائف المقيدة للأفراد والشركات.

ما هي التكنولوجيا التي ظهرت في القرن العشرين؟

Context length: 361 characters

Query has 72 tokens.

Segment A (Question) tokens: 11

Segment B (Context) tokens: 61

Answer start position: 16

Answer end position: 17

Answer tokens: ['تكنولوجيا', '،', 'المعلومات']

Answer: "تكنولوجيا المعلومات"

Confidence: 0.8793

ما الأجهزة الإلكترونية التي أصبحت جزءاً من حياتنا اليومية؟

Context length: 361 characters

Query has 73 tokens.

Segment A (Question) tokens: 12

Segment B (Context) tokens: 61

Answer start position: 24

Answer end position: 29

Answer tokens: ['الكمبيوتر', '،', '#ات', '،', 'والإنترنت', '،', 'أجهزة', '،', 'الهواتف', '،', 'المحمولة']

Answer: "الكمبيوترات والإنترنت وأجهزة الهواتف المحمولة"

Confidence: 0.8472

ما الفوائد التي يمكن أن تتحققها هذه البرامج والتطبيقات للأفراد والشركات؟

Context length: 361 characters

Query has 75 tokens.

Segment A (Question) tokens: 14

Segment B (Context) tokens: 61

Answer start position: 68

Answer end position: 69

Answer tokens: ['الخدمات', '،', 'والوظائف']

## Context Text (wrapped):

إن عاصمة العراق هي مدينة بغداد ، كانت تسمى قديماً مدينة السلام. تقع مدينة بغداد في وسط العراق، وهي تابعة لمحافظة بغداد. كانت بغداد قديماً العاصمة المركزية للخلافة العباسية، فالذي بني مدينة بغداد هو الخليفة العباسي أبو جعفر المنصور في القرن الثاني الهجري، والذي يوافق القرن الثامن الميلادي، وبعد اكتمال بناء المدينة، اتخذها أبو جعفر المنصور عاصمة لدولة الخلافة العباسية، وبقيت بغداد عاصمة الخلافة العباسية حتى جاء الخليفة العباسي المعتصم بالله، وقام ببناء مدينة سامراء وجعلها عاصمة الخلافة، ومدينة سامراء تقع شمال مدينة بغداد.

## Testing answer\_question function:

### 1. First Question:

Question: ما هي عاصمة العراق؟

Context length: 526 characters

Query has 105 tokens.

Segment A (Question) tokens: 7

Segment B (Context) tokens: 98

Answer start position: 12

Answer end position: 12

Answer tokens: ['بغداد']

Answer: "بغداد"

Confidence: 0.8378

### 2. Second Question:

Question: ماذَا كانت تسمى قديماً؟

Context length: 526 characters

Query has 105 tokens.

Segment A (Question) tokens: 7

Segment B (Context) tokens: 98

Answer start position: 17

Answer end position: 18

Answer tokens: ['مدينة', '،', 'السلام']

Answer: "مدينة السلام"

Confidence: 0.9730

## Technology context

## Baghdad context

# TESTING

```
--- Question 1 ---  
Question: في أي عام تم تأسيس مدينة بغداد؟  
Context length: 237 characters  
-----  
✖ Using Advanced Prompt Engineering...  
  
💡 Generated 4 different prompt strategies:  
• direct_None_clean: 145 (confidence: 0.6536)  
• formal_specificity_clean: 145 (confidence: 0.8449)  
• detailed_precision_structure: 145 (confidence: 0.7886)  
• analytical_completeness_segment: 145 (الفن الذي يجري: 0.1542)  
  
⭕ Best Answer Selected:  
Answer: 145  
Confidence: 0.8449  
Strategy Used: formal_specificity_clean  
  
--- Question 2 ---  
Question: في أي قرن تم بناء مدينة بغداد؟  
Context length: 237 characters  
-----  
✖ Using Advanced Prompt Engineering...  
  
💡 Generated 4 different prompt strategies:  
• direct_None_clean: 145 (الفن الذي يجري: 0.2651)  
• formal_specificity_clean: 145 (الفن الذي يجري: 0.2573)  
• detailed_precision_structure: 145 (الفن الذي يجري: 0.3797)  
• analytical_completeness_segment: 145 (الفن الذي يجري: 0.2906)  
  
⭕ Best Answer Selected:  
Answer: 145  
Confidence: 0.3797  
Strategy Used: detailed_precision_structure
```

010 1001  
110 0110  
111 0001  
011 1101  
110 0110  
111 0001  
011 1101  
001 0001



Prompt engineering strategies for Arabic QA

# INTERFACE

## نظام الإجابة على الأسئلة العربية

السؤال

ما هو أقرب كوكب للشمس

طريقة الإجابة

النص المحدد مسبقاً  كتابة نص بنفسى  بحث على الإنترنت

أجب

حالة المعالجة

تمت المعالجة بنجاح

نتيجة النظام

Question: ما هو أقرب كوكب للشمس  
Context length: 587 characters

Query has 165 tokens.

Segment A (Question) tokens: 7  
Segment B (Context) tokens: 158  
Answer start position: 29  
Answer end position: 30  
Answer tokens: ['عطا', '#رد', 'رد']  
Answer: "عطارد"  
Confidence: 0.5975



Gradio-based user interface for the Arabic QA system

# CURRENT CHALLENGES

Inferential Questions:  
Lower F1 score (0.4398)  
due to limited  
semantic reasoning

Long Contexts:  
Performance degrades  
due to memory  
constraints.

Complex Queries:  
Prompt engineering  
less effective for non-  
factoid questions.

Example: Space  
exploration questions  
showed incomplete or  
vague answers (e.g.,  
missed “Perseverance”  
for Mars rovers).



# ENHANCING THE SYSTEM

Long-Context Optimization: Use context chunking or sparse attention (AraBERT) to handle longer texts.

Inferential Questions: Integrate knowledge graphs or external knowledge bases for better reasoning.

Dynamic Prompts: Apply reinforcement learning for adaptive prompt generation.

Multilingual Support: Enable cross-lingual QA (Arabic-English).

User Interface: Add real-time feedback and voice query support for broader accessibility.



# KEY TAKEAWAYS

<b>Achievements</b>	:Built a robust Arabic QA system with AraElectra, achieving 0.4814 F1 and 0.6372 accuracy on AAQAD.
<b>Innovations</b>	:Integrated prompt engineering, web retrieval, and Gradio interface for usability.
<b>Impact</b>	:Enhances Arabic NLP, supporting education, research, and customer service.
<b>Future Potential</b>	:Scalable for multilingual and complex query applications.



# THANK YOU!