

مقدمة لنماذج المحوَّلات

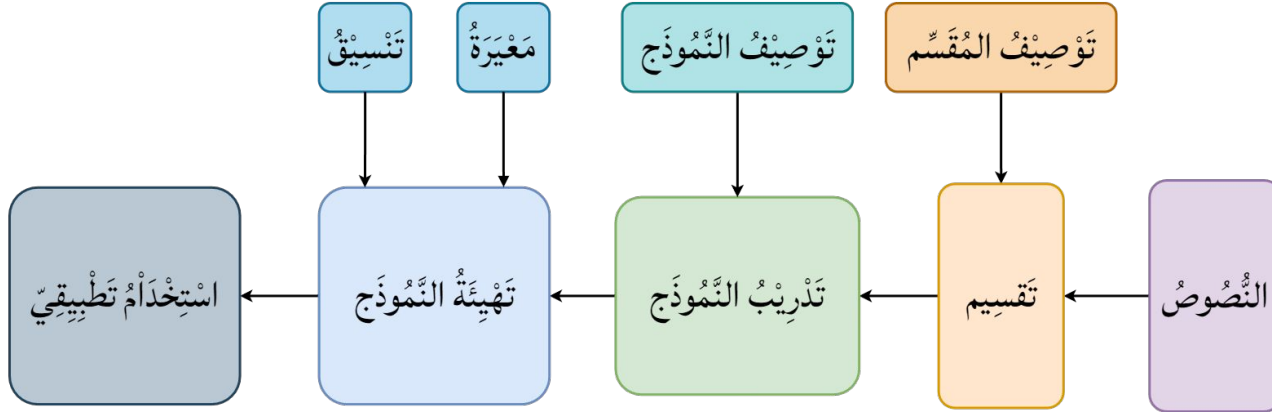
سري السباعي

فهرس العرض

1. مقدمة أوليّة عن العرض والمُصطلحات
2. تصوّر عام عن النماذج اللغوية
3. مقدمة عن أشكال نماذج المحوّلّات الثلاثة
4. شرح المحول الفاكّ المؤكّد
 - 4.1. البيانات وتمثيلها
 - 4.2. قُبيلَ الطبقة الجامعة
 - 4.3. الطبقة الجامعة
 - 4.4. بُعيد الطبقة الجامعة
5. معيرة النماذج

مقدمة عن العرض

المقصد من هذا العرض تصوير أُسُس النماذج اللغوية من أصلها إلى النماذج الحديثة مع ذكر أشهر النماذج والمدونات العربية. وتبيين أهم المصطلحات في المجال وترجمتها. وكذا الإحالة على مصادر مفيدة لمن أراد التوسع فما مجلسنا هذا إلا قطرة من يَمٍ وقليلًا من جم.



أهم المصطلحات المُستعملة

Transformers	المحولات	Language model	نموذج لغوي
Decoder	مفكك	FineTune	معيرة (ضبط دقيق)
Encoder	مشفر	Understand	فهم
Corpus\Dataset	مدونة (المُدخلات)	Generative	توليد
Tokenizer	مُقَسِّم (مُجزّء)	Token	كُليمة (جزء لغوي)

الحوسبيّات



اللسانيات

أصل العرض

المحاولات التوليدية العربية بنظرة تقنية رياضية بحثية مختصرة

سري بن تيسير السباعي

مهندس باحث في معمل الروبوتات وإنترنت الأشياء في جامعة الأمير سلطان في السعودية

مقدمة

الحمد لله وصلاة وسلاما على رسول الله أما بعد: فهذه ورقة معاصرة ألخص فيها أسس نماذج المحاولات المولدة فأكبر التشفير¹ مخصصا منها ما كان متعلقا بالعربية. قسمتها على ثلاث مقالات كل مقالة فصول وكل فصل أبواب وقبلهن مقدمة في ترجمة المصطلحات المستعملة في الورقة. أما تفصيل المقالات الثلاث:

فأولهن: في الكلام عن المدونات النصية العربية أنواعها وطرق استنصاحها² وتجميعها ثم نتطرق إلى طرق تمثيلها الرقمية للنموذج، والثانية: نصفت فيها التركيب الرياضي وهيكل النموذج المحول فالك التشفير وأنواع بناء مكوناته واختلافها وألمح إلى مبحث سبب تعلم هذه النماذج، مختتما بذكر أسماء أشهر النماذج اللغوية العربية.

والثالثة في مميزات النماذج اللغوية الكبيرة واستعمالها على بيانات مخصوصة والتطبيقات المختلفة لذلك.

ويتخلل الفصول تنبيهات وفوائد حسب المناسب ان شاء الله.

واعلم أن هذه الورقة لم تُبَيَّن فيها نية الاستقصاء ولا همة الجمع (فهذا يحتاج كتابا أو أكثر) وإنما غايتها التنبيه والإلماح لعدد كبير من المصادر والمجالات والأفكار البحثية وربطها بالعربية نسأل الله العون وحسن الاختيار وصلى الله وسلم على نبيِّنا محمد وعلى آله وصحبه أجمعين.

وما من كاتبٍ إلا سيفنى ... ويُبقِي الدهر ما كتبت يده

فلا تكتب بكفك غير شيء ... يسرُّك في القيامة أن تراه³

نظرة عامة عن النماذج اللغوية

طريقة ومنهج لتمثيل اللغة حوسبياً إما **لتوليد أو فهم** بعض **خصائصها اللغوية** حسبَ مرتبتها (بدايةً من التحليل النطقي وصولاً إلى التحليل المُعجمي الشامل الفهم الدلالي).

فمن الأمثلة على الفهم المَبني على قواعد ثابتة نظام دردشة "إِلزّا" (الخمسينات) وكذا أنظمة فهم أقسام الكلام العربي وتصريفها

وأما أمثلة التوليد فلم تظهر إلا بعد طفرة **التعلّم الإحصائي** بأن يتوقّع النموذج احتمال الكُليمة التالية بناءً على سياق محدد

أمثلة على الطُّرق القديمة

البحث بالكلمات

البرمجة المباشرة

باستعمال "الأنماط الطبيعية"
REGEX

المُتَجَهات العَدِّيَّة

التحليل الصرفي للجذور
والجذوع

معجم تركيبى للتوليد والتحليل
الصرفي والنحوي

الشبكات العصبيَّة الصِّرفَة

الشبكات التكراريَّة

المُتَجَهات الكلماتيَّة

أهم المحطات التاريخية باختصار

2017 المحولات

2022 متمم-3

2013 الكلمجات
word2vec

1997 الشبكات التكرارية
RNN



متمم

أنواع المحولات الثلاثة

متمم (مفكك)

مُترجم (مُشفر فاك)

متفهم (المُشفر)

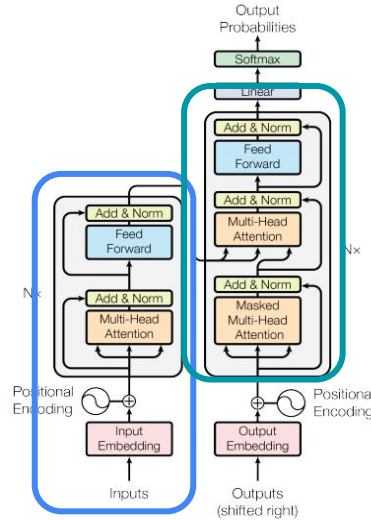
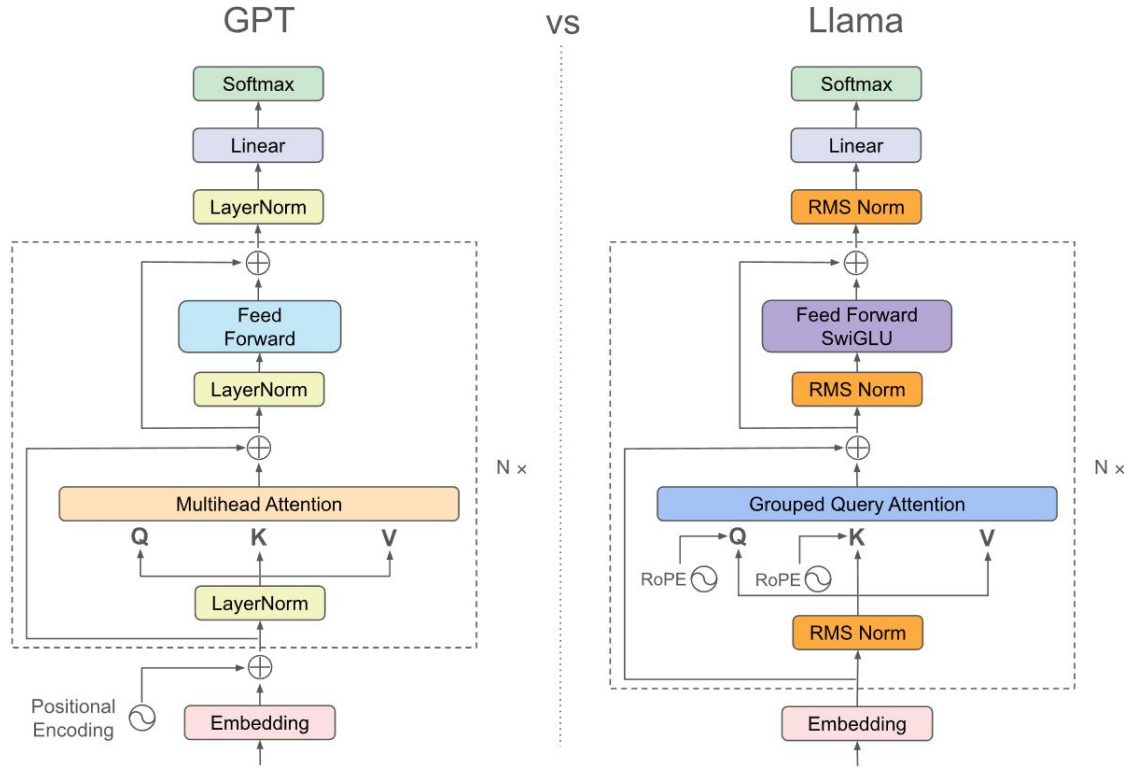


Figure 1: The Transformer - model architecture.

أشهر المفكّكات



Jais

AYA

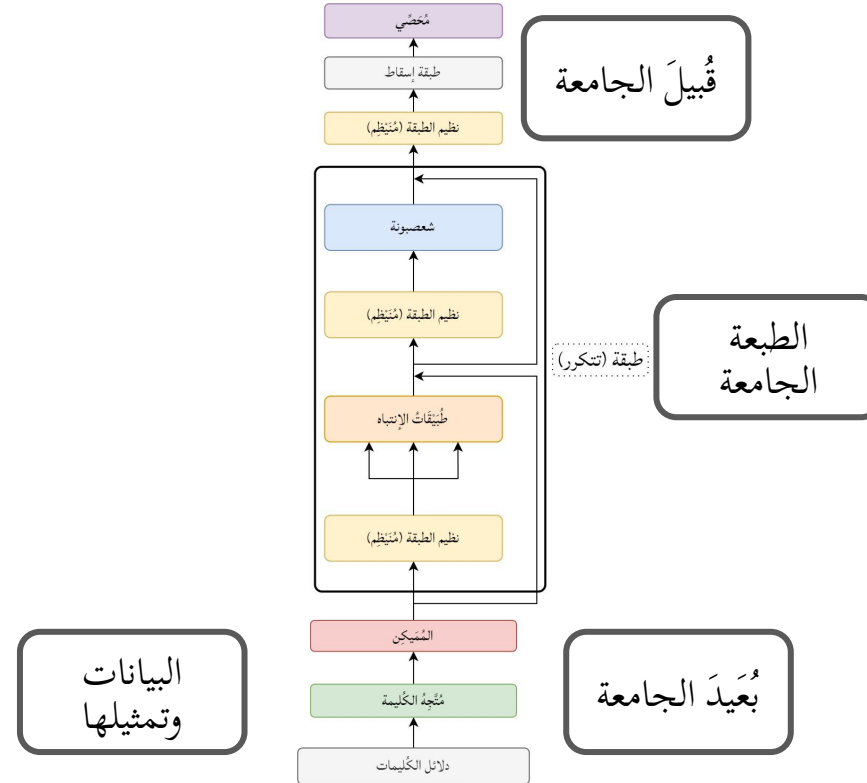
ALLAM

QWEN

ArabianGPT

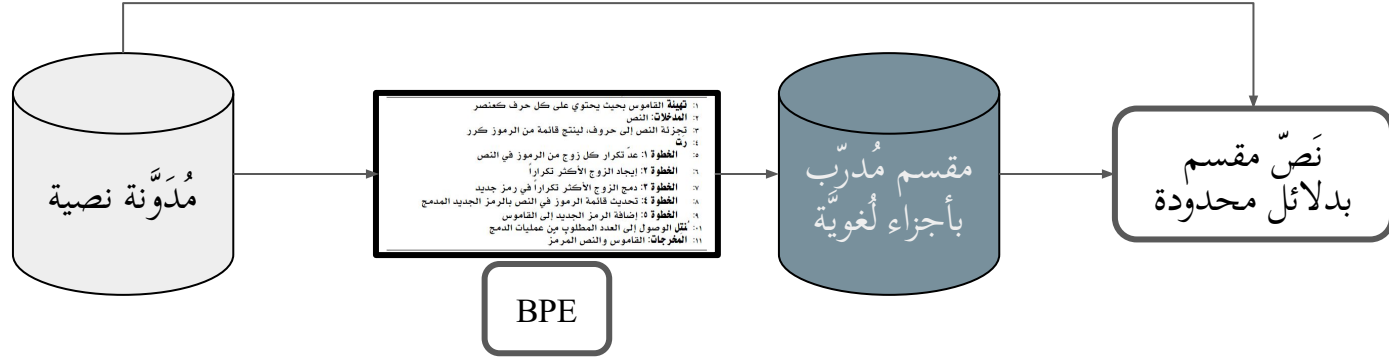
AraGPT

نموذج المحوّل الفاكّ



البيانات وتمثيلها

أشهر تمثيل للنصوص المكتوبة يكون باستعمال خوارزمية التزويج البايتي^١ وعليها بُنيت الباقيات وغالب النماذج العربية استعملتها



400 جيجا	مدونة CultureX
200 جيجا	مدونة ArabicText2022
-	مدونة 101 مليار

قُبيل الجامعة



$$E = \text{Embedding}(X) \quad E \in \mathbb{R}^{T \times d_{\text{model}}}$$

$$PE = \text{PositionalEncoding}(T, d_{\text{model}})$$

$$Z_0 = E + PE \quad Z_0 \in \mathbb{R}^{T \times d_{\text{model}}}$$

الطبقة الجامعة

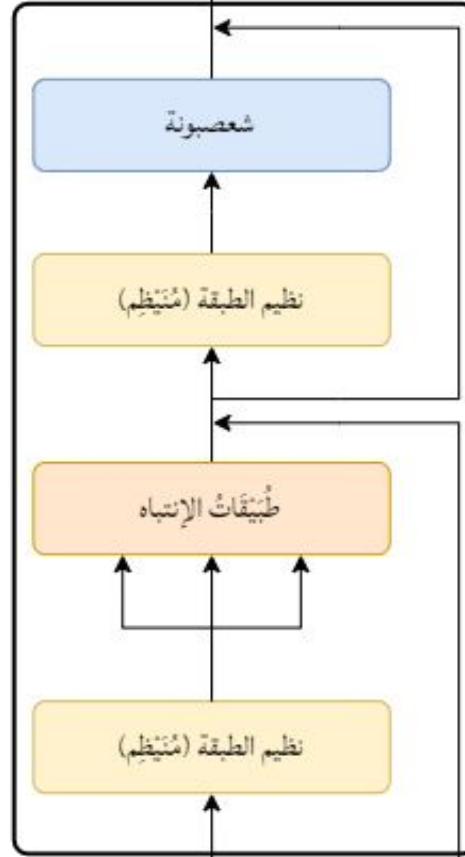
$$F_1 = \hat{Z}_1 W_1 \quad F_1 \in \mathbb{R}^{T \times d_{ff}}$$

$$F_2 = f(F_1)$$

$$Z_1 = Z_0 + \text{MultiHead}(Q, K, V)$$

$$\hat{Z}_1 = \text{LayerNorm}(Z_1) \quad \hat{Z}_1 \in \mathbb{R}^{T \times d_{model}}$$

$$\hat{Z}_0 = \text{LayerNorm}(Z_0) \quad \hat{Z}_0 \in \mathbb{R}^{T \times d_{model}}$$



$$Q = \hat{Z}_0 W_Q \quad K = \hat{Z}_0 W_K \quad V = \hat{Z}_0 W_V$$

$$Q, K, V \in \mathbb{R}^{T \times d_k}$$

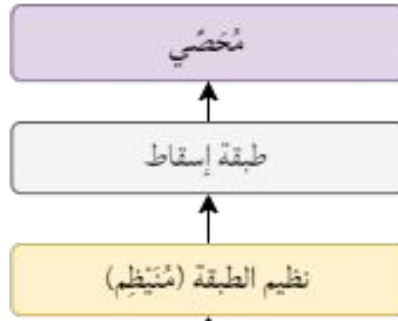
$$\text{Attention}(Q, K, V) = \text{Softmax} \left(\frac{QK^T}{\sqrt{d_k}} + M \right) V \quad M \in \mathbb{R}^{T \times T}$$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W_O$$

$$W_O \in \mathbb{R}^{(h \times d_k) \times d_{model}}$$

بُعْدَة الجامعة

$$\mathcal{L}(\hat{y}, y) = \text{CrossEntropy}(\hat{y}, y)$$



$$Z_2 = F_2 W_2 \quad Z_2 \in \mathbb{R}^{T \times d_{model}}$$

$$Z_3 = Z_1 + Z_2 \quad Z_3 \in \mathbb{R}^{T \times d_{model}}$$

$$Y = \text{Softmax}(Z_3 W_Y) \quad Y \in \mathbb{R}^{T \times d_{vocab}}$$

مبرمجة

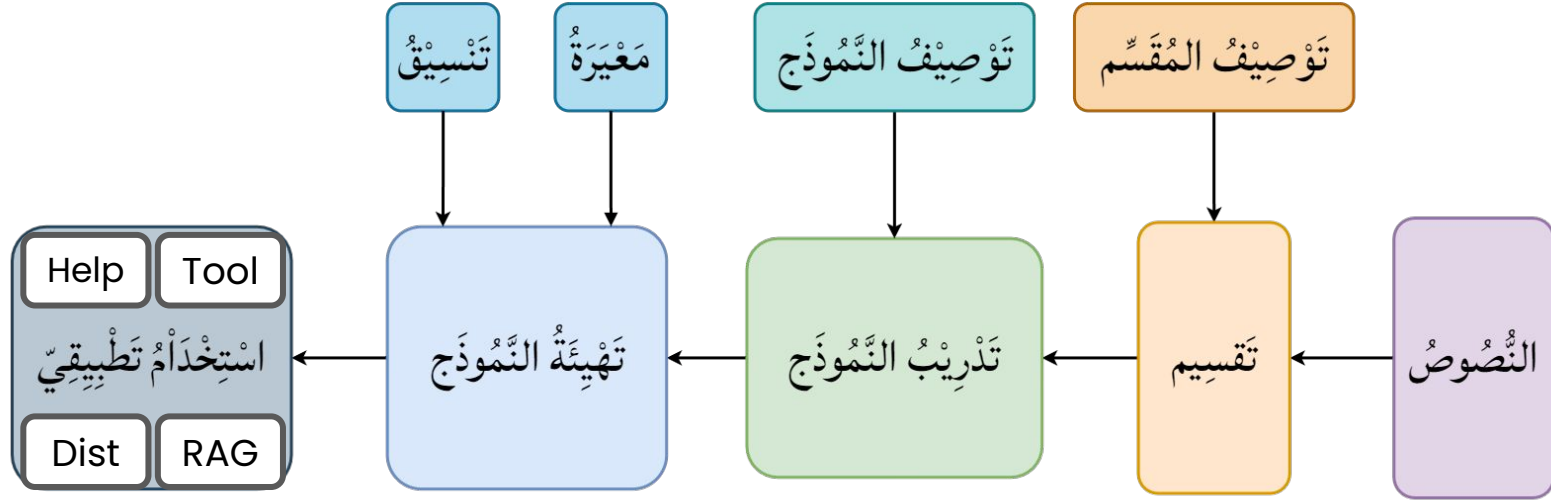
```
1 import numpy as np
2
3 def gelu(x):
4     return 0.5 * x * (1 + np.tanh(np.sqrt(2 / np.pi) * (x + 0.044715 * x**3)))
5
6 def softmax(x):
7     exp_x = np.exp(x - np.max(x, axis=-1, keepdims=True))
8     return exp_x / np.sum(exp_x, axis=-1, keepdims=True)
9
10 def layer_norm(x, g, b, eps: float = 1e-5):
11     mean = np.mean(x, axis=-1, keepdims=True)
12     variance = np.var(x, axis=-1, keepdims=True)
13     return g * (x - mean) / np.sqrt(variance + eps) + b
14
15 def linear(x, w, b):
16     return x @ w + b
17
18 def ffn(x, c_fc, c_proj):
19     return linear(gelu(linear(x, **c_fc)), **c_proj)
20
21 def attention(q, k, v, mask):
22     return softmax(q @ k.T / np.sqrt(q.shape[-1]) + mask) @ v
23
24 def mha(x, c_attn, c_proj, n_head):
25     x = linear(x, **c_attn)
26     qkv_heads = list(map(lambda x: np.split(x, n_head, axis=-1), np.split(x, 3, axis=-1)))
27     causal_mask = (1 - np.tri(x.shape[0])) * -1e10
28     out_heads = [attention(q, k, v, causal_mask) for q, k, v in zip(*qkv_heads)]
29     x = linear(np.hstack(out_heads), **c_proj)
30     return x
31
32 def transformer_block(x, mlp, attn, ln_1, ln_2, n_head):
33     x = x + mha(layer_norm(x, **ln_1), **attn, n_head=n_head)
34     x = x + ffn(layer_norm(x, **ln_2), **mlp)
35     return x
36
37 def gpt2(inputs, wte, wpe, blocks, ln_f, n_head):
38     x = wte[inputs] + wpe[range(len(inputs))]
39     for block in blocks:
40         x = transformer_block(x, **block, n_head=n_head)
41     return layer_norm(x, **ln_f) @ wte.T
42
43 def generate(inputs, params, n_head, n_tokens_to_generate):
44     from tqdm import tqdm
45     for _ in tqdm(range(n_tokens_to_generate), "generating"):
46         logits = gpt2(inputs, **params, n_head=n_head)
47         next_id = np.argmax(logits[-1])
48         inputs = np.append(inputs, [next_id])
49     return list(inputs[len(inputs) - n_tokens_to_generate :])
50
51
```


معيّرة النماذج

هي عملية تدريب عاديّة ولكنّ المُدخلات (المدوّنة) تكون بصيغة مُهيكلّة ذات بُنية محدّدة ليتعلّمها النموذج بإحاطة المُدخلات بأجزاء لغوية مخصّصة يتبع النموذج أنماطها بعد التعلم

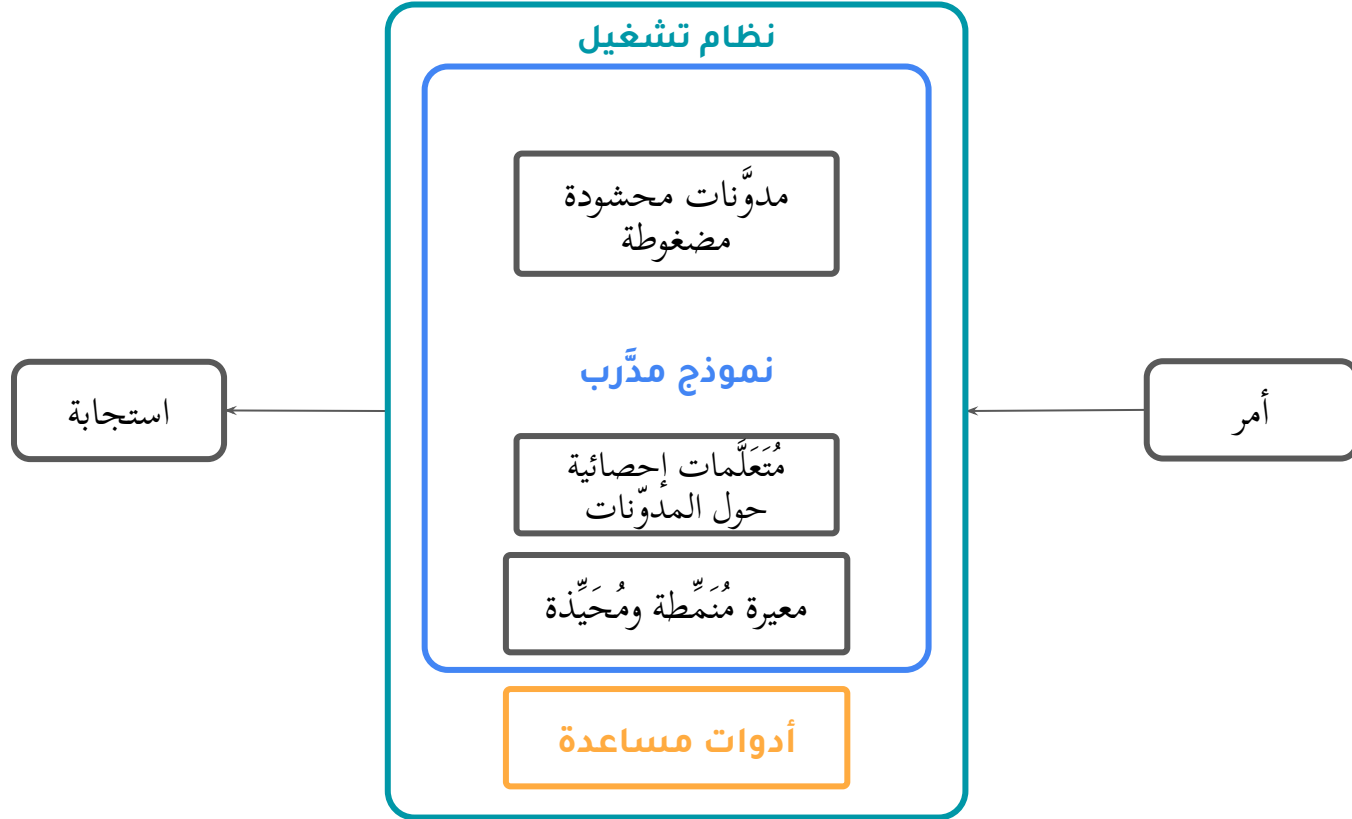
سؤال		جواب	[مستخدم] ... [/مستخدم] [نموذج] ... [/نموذج]	
شخصية	سياق	سؤال	جواب	
[مستخدم] [سياق] ... [/سياق] [سؤال] ... [/سؤال] [/مستخدم] [نموذج] ... [/نموذج]				

عَوْدًا عَلَى بَدْءِ



رَأَيْتُ النَّاسَ غَالِبَهُمْ رُسُومِي... وَلَيْسَ بِقَارِيٍّ يَفْرِي السُّطُورَا

نظام تشغيل جديد



خاتمة وتوصيات

عرضنا (بل إن شئت قُلْ ضغطنا وحشرنا) في هذه الورقة عُلالةً مُختصرة وإشارات مُقتصرة حول نماذج المحولات التوليدية العربية وما يتعلق بها من مُدُونات وخُرُزمات. فقد شرحنا مركبات النموذج ومن ثَمَّ أشهر النماذج العربية ثم ختمت بتطبيقات النماذج اللغوية. وأختم بأمر يجول بخاطري حول علاقة العربية بهذا المجال من مُنطلقين: الأول حول ضرورة الإهتمام بها ابتداءً قبل التسارع إلى إدخالها في النماذج اللغوية أو غيرها من التطبيقات الحاسوبية فإنه لا يُتصور أن تكتب الأبحاث وتُنشر المقالات بلغة غيرها والمُتَكَلِّم عنه هي. والثاني حول المجالات المُهتَم بها وأولويَّتها في خدمة العربية، فإن الناظر يُدرك بلا شك أن غالب العلوم الآن منشورة بالكلزية ولا بُدَّ أن تُعرَّب وتُعدَّل المناهج كذلك فأظن أنه ينبغي الإهتمام بدءًا بالترجمة وتكثير بحوثها وتطويرها حتى نقدر على تقديم هذه العلوم بأصولها وأُسُسها بلسان عربي ثم ننطلق بعد ذلك لضغط هذه المعارف وصهرها في النماذج اللغوية. وأخيرا أختم بنصيحة استللتها من كلام تشومسكي لَمَّا سُئل عن هذه النماذج فقال "أنها لم تُقدِّم شيئا في المجال العلمي وإنما هي أدوات مفيدة هندسيا وتطبيقيا" [74] أقدمُها للباحث لألا يترك ما بين يديه من بحوث نظرية ومسائل علمية ويَهْرَع لهذه المُندَلقات يَكُتِب حولها بل فليَلْزَم كراسته وقلَمه ويَحْتِمْ فإنما هي زوبعةٌ في فنيجان وجعجعة بلا إطحان إذ:

(من راقب الناس لم يظفر بحاجته... وفاز بالطيبات الفاتك اللهج)²⁸

و (يا باري القوس بريا ليس يُحسُّنه ... لا تظلم القوس أعط القوس باريها)²⁹

للتواصل:

<https://www.linkedin.com/in/serry-sibae>

<https://github.com/serrysibae/>

serrytowork@gmail.com

المراجع

1. ورقة المحولات التوليدية العربية بنظرة تقنية رياضية بحثية مُختصرة، سري بن تيسير السباعي 2024.
https://drive.google.com/file/d/1rP8jpJOVM82OISqIT0Kq0PUmsoab_bEU/view?usp=sharing
2. كتاب "كيف يعمل متمم" (مقدمة ممتازة)
<https://writings.stephenwolfram.com/2023/02/what-is-chatgpt-doing-and-why-does-it-work>
3. ورقة المحولات الأصلية Attention is all you need
https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf
4. دراسة مسحية عن النماذج العربية
<https://arxiv.org/abs/2410.20238>
5. برمجة "متمم" بطريقة مختصرة
https://github.com/jaymody/picoGPT/blob/a750c145ba4d09d5764806a6c78c71ffa88e64/gpt2_pico.py#L3-L58
6. موقع Huggingface لتجربة النماذج والمدونات
<https://huggingface.co>
7. مختصر تعريفى ببعض النماذج العربية المشهورة
<https://docs.google.com/document/d/1pNt4siq6cbTeeT7CEBXJER8yxBOQPxd5-rUbDLVGhHQ/edit?usp=sharing>
8. معجم سدايا للذكاء الاصطناعي
<https://sdaia.gov.sa/ar/MediaCenter/KnowledgeCenter/ResearchLibrary/SDAIAPublications15.pdf>

المراجع

1. مقالات سري السباعي على كراميل (فيها مقدمة عن المجال وكذا مقالتان أخريان)

<https://caramellaapp.com/serrysibae>

2. كتاب الدكتور نزار حبش "مقدمة إلى معالجة اللغة العربية"

3. كتاب جرفاسكي [/https://web.stanford.edu/~jurafsky/slp3](https://web.stanford.edu/~jurafsky/slp3)

4. مقرر ستانفورد CS124 وهناك مرئيات له [/https://web.stanford.edu/class/cs124](https://web.stanford.edu/class/cs124)

5. معاجم الذكاء الصناعي المذكورة: معجم طعيمة

<https://dlarabic.com/%d9%83%d8%aa%d8%a7%d8%a8-%d9%85%d8%b9%d8%ac%d9%85-%d9%85%d8%b5%d8%b7%d9%84%d8%ad%d8%a7%d8%aa-%d8%a7%d9%84%d8%aa%d8%b9%d9%84%d9%85-%d8%a7>

[5%d8%b7%d9%84%d8%ad%d8%a7%d8%aa-%d8%a7%d9%84%d8%aa%d8%b9%d9%84%d9%85-%d8%a7](https://dlarabic.com/%d9%83%d8%aa%d8%a7%d8%a8-%d9%85%d8%b9%d8%ac%d9%85-%d9%85%d8%b5%d8%b7%d9%84%d8%ad%d8%a7%d8%aa-%d8%a7%d9%84%d8%aa%d8%b9%d9%84%d9%85-%d8%a7)

[/ %d9%84%d8%a2%d9%84%d9%8a-%d9%88%d8%a7%d9%84%d8%aa%d8%b9%d9%84](https://dlarabic.com/%d9%83%d8%aa%d8%a7%d8%a8-%d9%85%d8%b9%d8%ac%d9%85-%d9%85%d8%b5%d8%b7%d9%84%d8%ad%d8%a7%d8%aa-%d8%a7%d9%84%d8%aa%d8%b9%d9%84%d9%85-%d8%a7)

6. محاضرات كراثي <https://www.youtube.com/playlist?list=PLAqhIrjkbxWl23v9cThsA9GvCAUhRvKZ>