



مجمع الملك سلمان  
العالمي للغة العربية  
King Salman Global Academy for Arabic Language

# بَلْسَم: مؤشر نضج تقنيات الذكاء الاصطناعي للغة العربية

التقرير الثاني

فبراير- ٢٠٢٦



HUMAIN



stc



aiXplain



groq



شركاء  
المجمع:

## كلمة الأمين العام لمجمع الملك سلمان العالمي للغة العربية

تتكامل أعمال مجمع الملك سلمان العالمي للغة العربية في مسارات متعددة تشمل: الحوسبة، والتخطيط، والثقافة، والتعليم؛ بما يحقق أهدافه وإستراتيجيته وفقاً لمهامه التي نص عليها تنظيم المجمع؛ خدمةً للغة العربية، وتعزيزاً لمكانتها العالمية.

ويأتي مؤشر (بَلَسَم) ضمن سياق أعمال المجمع المختلفة في الحوسبة العربية؛ سعياً إلى توظيف الذكاء الاصطناعي في خدمة العربية، واستكمالاً لأعماله، مثل: (مرصد العربية) الذي يقيس واقع حوسبة العربية مقارنةً باللغات الأخرى، ويرصد أهم المبادرات فيه، و(مركز ذكاء العربية) الذي يقدم خدماتٍ متنوعةً لرواد صناعة تقنيات اللغة العربية، ويجمع الباحثين والأساتذة من اللغويين والحاسوبيين.

ويسعد المجمع بنشر تقرير أداء النماذج اللغوية في النصف الثاني من عام ٢٠٢٥م، الذي يمثل تطوراً واضحاً لمؤشر (بَلَسَم)؛ بتوسيع نطاق التقييم ليشمل أسئلةً جديدةً وفئاتٍ إضافيةً؛ مما يعزز شمولية المؤشر ودقته في قياس أداء نماذج الذكاء الاصطناعي للغة العربية. ويأتي التقرير بالشراكة مع عدد من الجهات المحلية والدولية التي أسهمت في بناء بياناته وأدوات قياسه، وقدمت استشاراتٍ في تطوير منصة متاحة للتقييم، تجمع بين التحكيم البشري والتقييم الآلي عالي الدقة؛ مما أسهم في الوصول إلى مقاييس موثوقة، وتعزيز اعتماده على المستوى النظري والتطبيقي.

وفي الختام، نتقدم بالشكر إلى صاحب السمو الأمير/ بدر بن عبدالله بن فرحان آل سعود -حفظه الله- وزير الثقافة، رئيس مجلس أمناء المجمع، على دعمه الدائم لمشروعات المجمع وأعماله، وإلى شركاء المؤشر الذين أسهموا في تحقيق الأهداف النهائية وبناء البيانات الابتدائية له، ونتطلع في التقارير القادمة إلى مزيد من التعاون والشراكة مع العديد من الجهات الأكاديمية والتقنية؛ لتوسيع نطاق العمل البحثي والتقني في المعالجة الآلية للغة العربية؛ بما يعزز من ريادة اللغة العربية في المجال التقني.

الأمين العام

أ.د.عبد الله بن صالح الوشمي.

٤ ..... (بَلْسَم) مؤشر نضج تقنيات الذكاء الاصطناعي للغة العربية

٤ ..... إحصاءات عامة من التقييم لمؤشر (بَلْسَم)

٥ ..... فئات المهام المشمولة في التقييم

١١ ..... منهجية التقييم

١١ ..... أولاً: أنماط الأسئلة

١٢ ..... ثانياً: قياس النتائج

١٢ ..... أ) التقييم الكمي

١٢ ..... ب) التقييم الآلي

١٣ ..... ج) التقييم البشري

١٣ ..... ثالثاً: طريقة التقييم المعتمدة

١٣ ..... رابعاً: توصيات التقرير الثاني:

١٤ ..... النماذج المُقيّمة في التقرير الثاني

١٧ ..... النتائج

١٧ ..... أولاً: متوسط الأداء العام للنماذج

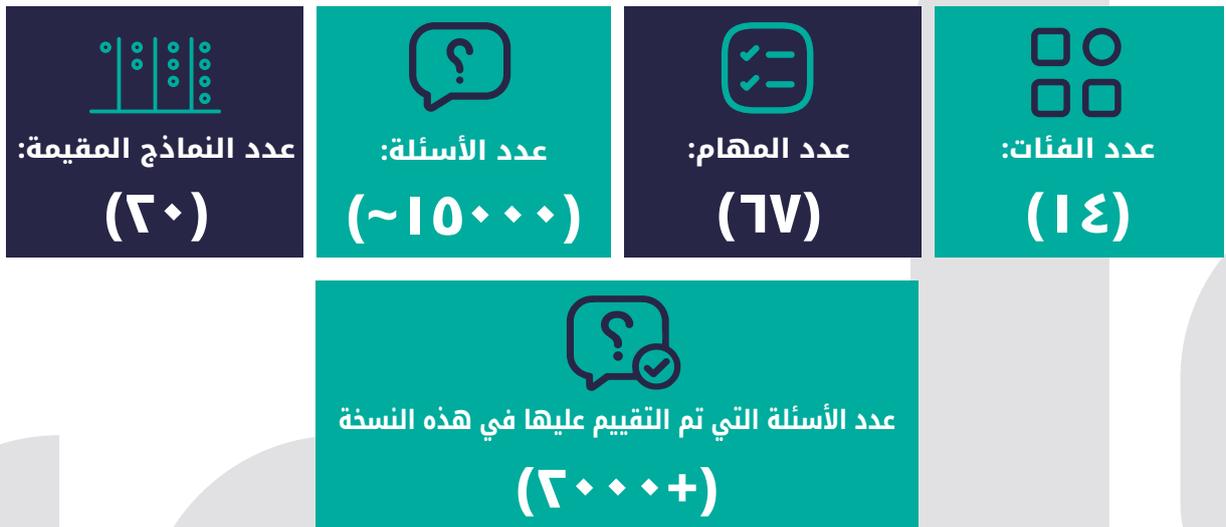
١٨ ..... ثانياً: تقييم النماذج حسب الفئات ومهام كل فئة

# بَلَسَم مؤشّر نضج تقنيات الذكاء الاصطناعي للغة العربية

هو مبادرة تعاونية يقودها مجمع الملك سلمان العالمي للغة العربية، بالتعاون مع عدد من الشركاء من داخل المملكة العربية السعودية وخارجها، ويهدف إلى قيادة تطوير وتجهيز مجموعات بيانات اختبار متخصصة تُعدّ ضروريةً لتقييم أداء نماذج اللغة الكبيرة (LLMs) في مجموعة متنوعة من مهام معالجة اللغة الطبيعية (NLP) باللغة العربية. وتسهم هذه المبادرة في سد فجوة معرفية مهمة في مجال تقييم النماذج اللغوية العربية؛ بتوفير أطر ومنهجيات تقييم دقيقة تساعد على قياس أداء هذه النماذج كلياً أو جزئياً. وتهدف إلى تأسيس فهم أعمق لنقاط القوة والضعف في النماذج اللغوية عند تقييم أدائها في مهام متعددة، مثل: التدقيق النحوي، وتوليد النصوص، وفهم المحتوى.

يهدف هذا التقرير إلى تقديم تحليل شامل لأداء النماذج اللغوية الكبيرة (LLMs) على مؤشر (بَلَسَم)؛ إذ يستند إلى تصميم وبناء مجموعات بيانات اختبار متخصصة جرى تحليل نتائجها ومقارنتها مع نتائج التقييم البشري؛ لاختبار مدى دقة النماذج في محاكاة التقييم البشري. ويشمل التقييم تحليل أداء النماذج في مجموعة واسعة من المهام ضمن كل فئة، مع ترتيب النماذج من الأعلى إلى الأدنى، إضافةً إلى استعراض متوسط الأداء العام بجميع الفئات. وشهدت هذه النسخة توسعاً في تصميم مهام التقييم من حيث تنوع أنماط الأسئلة؛ إذ تضمنت أسئلةً توليديةً وأسئلةً اختياراً من متعدد؛ بما يتيح تقييماً أكثر شمولية لقدرات النماذج اللغوية؛ بقياس أدائها في إنتاج الإجابات المفتوحة، ودقة الاختيار والاستدلال ضمن اختيارات محددة؛ وبما يعكس استخدامات واقعيةً ومتنوعةً للنماذج اللغوية في التطبيقات العملية. وقد تم تقييم أسئلة جديدة لم ترها النماذج من قبل.

## إحصاءات عامة عن التقييم لمؤشّر (بَلَسَم)



## فئات المهام المشمولة في التقييم: ✓

تعريف المهمة	المهمة	تعريف الفئة	الفئة
إكمال جملة معطاة بجملة منطقية ومتوافقة مع السياق.	إكمال النصوص (كلاهما)	كتابة نصوص أصلية متنوعة تتطلب الإبداع اللغوي.	الكتابة الإبداعية
توليد ردود حوارية تمثل تفاعلاً بين أطراف الحديث.	توليد الحوارات (كلاهما)		
تقديم تفسيرات أو حقائق داعمة لإجابة؛ لتوضيح المعلومة.	الشرح (كلاهما)		
تحويل أسئلة المستخدم أو استفساراته إلى تعليمات واضحة ومصنفة.	توليد التعليمات (توليدي)		
كتابة مقالات إخبارية مستندة إلى وصف محدد للحدث.	توليد المقالات الإخبارية (توليدي)		
توليد قصائد مرتبطة بعنوان أو نص معين.	توليد القصائد (توليدي)		
إنشاء سؤال بناءً على إجابة معطاة.	توليد الأسئلة (توليدي)		
إنتاج افتراض جديد مستمد من فرضية دون تعارض أو إضافة.	تأليف الجمل (توليدي)		
كتابة قصة منطقية ومتناسكة لنص سردي معطى.	تأليف القصص (توليدي)		
توليد عنوان ملائم لنص يصف مضمونه .	توليد العنوانات (توليدي)		

تعريف المهمة	المهمة	تعريف الفئة	الفئة
إنتاج نص جديد ذي معنى بناءً على مدخل معين.	توليد النصوص (توليدي)		الكتابة الإبداعية
إنشاء تعريفات لمفاهيم أو عبارات بناءً على سياق المستخدم أو طلبه.	توليد التعريفات (توليدي)	كتابة نصوص أصلية متنوعة تتطلب الإبداع اللغوي.	
توليد إجابة خاطئة للسؤال.	توليد إجابات غير مناسبة (توليدي)		
تقييم ما إذا كانت جملة أو فقرة معينة تُعدُّ امتدادًا منطقيًا وسليمًا لنص سابق.	تقييم صحة تتمة النص (اختياري)		
تحديد ما إذا كان يمكن استنتاج جملة معينة من نص ما استنتاجًا منطقيًا؛ لتقييم العلاقة الدلالية بين النص والاستنتاج.	الاستلزام النصي (كلاهما)	تحديد ما إذا كان يمكن استنتاج جملة معينة منطقيًا من أخرى؛ بهدف تقييم العلاقة الدلالية بين النص والاستنتاج.	التضمن
تحديد نوع العلاقة الدلالية بين جملتين.	التشابح الدلالي (كلاهما)		
استكمال الجمل بكلمات مناسبة دلاليًا ولغويًا.	ملء الفراغ (كلاهما)	استكمال الجمل بكلمات مناسبة دلاليًا ولغويًا.	ملء الفراغ
تحديد نوع العلاقة أو المعلومة المطلوبة من نص معين.	استخلاص المعلومات المطلوبة (كلاهما)		استخلاص المعلومات
استخراج أسماء الأمراض المذكورة في النص.	تحديد أسماء الأمراض في النصوص (توليدي)	تحليل النصوص لاستخراج معلومات محددة.	

تعريف المهمة	المهمة	تعريف الفئة	الفئة
استخلاص الكلمات أو العبارات المفتاحية التي تلخص النص.	استخلاص الكلمات الرئيسية (توليدي)	تحليل النصوص لاستخراج معلومات محددة.	استخلاص المعلومات
تحديد الكيانات المسماة، مثل: (الأسماء، والأماكن، والجينات) من النص.	التعرف إلى أسماء الكيانات (توليدي)		
استخراج العلاقة بين كيانيين في نص معين.	استخلاص العلاقات (توليدي)		
تحديد وجود علاقة دلالية بين كيانيين مذكورين في النص من عدمها.	تصنيف العلاقات بين الكيانات (اختياري)		
تصنيف كيان معين إلى أقرب حقل دلالي مناسب له.	تصنيف الكيانات (اختياري)		
تحديد الجملة الناتجة عن سبب معين في نص.	تصنيف السبب والنتيجة (كلاهما)	استنتاج المنطق من خلال مهام تتطلب فهم العلاقات السببية، وحل الألغاز، والترتيب، والتنبؤ.	المنطق
تحديد الكلمة أو العبارة التي يشير إليها ضمير معين في الجملة.	تحليل الإحالات الضميرية (توليدي)		
تحليل الخيارات؛ لتحديد الأكثر منطقية بناءً على معطى معين.	التحليل التنبؤي (توليدي)		
حل الألغاز البسيطة باستخدام المنطق.	حل الألغاز (توليدي)		
تحديد الجملة التي تخالف المنطق السليم أو المعرفة العامة عن العالم.	التحقق من المنطق السليم (اختياري)		

تعريف المهمة	المهمة	تعريف الفئة	الفئة
ترتيب خطوة أو إدراجها في تسلسل منطقي للخطوات.	ترتيب الجمل (اختياري)	استنتاج المنطق من خلال مهام تتطلب فهم العلاقات السببية، وحل الألغاز، والترتيب، والتنبؤ.	المنطق
تطبيق التفكير المنطقي لاستخلاص نتيجة صحيحة من نص معين.	الاستدلال المنطقي (اختياري)		
تحديد ما إذا كانت جملة معيّنة متماسكةً دلاليًا مع الجملة أو السياق السابق لها.	تصنيف التماسك النصي (اختياري)		
تنفيذ المعادلات الرياضية أو المهام البرمجية.	تنفيذ البرامج (كلاهما)	تنفيذ المعادلات الرياضية أو المهام البرمجية.	تنفيذ البرامج
تقديم إجابة مباشرة وصحيحة لسؤال مطروح.	الإجابة عن الأسئلة المطروحة (كلاهما)	الإجابة عن السؤال المطروح.	الإجابة عن الأسئلة
تفكيك سؤال معقد إلى سلسلة أسئلة بسيطة للإجابة النهائية.	تحليل السؤال (توليدي)		
قياس قدرة النماذج على التحقق من صحة الاستدلالات اعتمادًا على النص.	التحقق من الاستدلال النصي (كلاهما)		
التأكد مما إذا كان التوضيح المقدم يعين على تفسير الغرض من السؤال الأصلي، ويزيل الغموض عنه.	فهم السؤال (كلاهما)	استخراج وفهم المعلومات من النصوص المكتوبة من خلال الإجابة عن أسئلة مستندة إلى المحتوى.	فهم المقروء
تحديد ما إذا كان السؤال يمكن أن يُجاب عنه إجابةً تامةً اعتمادًا على السياق أو المعلومات المتاحة.	تصنيف إمكانية الإجابة (اختياري)		
تحديد الكلمات التي تحتوي على أخطاء نحوية في الجملة.	اكتشاف الأخطاء النحوية (توليدي)	تحديد تسميات دقيقة لعناصر ضمن تسلسل نصي، مثل: تحديد الأخطاء النحوية، أو استخراج الكلمات المفتاحية.	توسيم السلاسل النصية
تحديد القسم الكلامي لكل كلمة داخل سياق الجملة.	تحديد القسم الكلامي (اختياري)		

تعريف المهمة	المهمة	تعريف الفئة	الفئة
اقتراح عنوان أو موضوع مناسب لنص أو بريد شبكي.	توليد العنوانات (توليدي)	موجز للنص يُبرز أهم المعلومات والمفاهيم.	التلخيص
تلخيص نص طويل إلى جملة أو عدة جمل تختصر المعنى.	تلخيص النصوص (توليدي)		
تحديد اللهجة المستخدمة في النص العربي.	تحديد اللهجة (كلاهما)		تصنيف النصوص
تصنيف النص إلى فئة معينة.	التصنيف الفئوي للنصوص (كلاهما)		
تحديد التعليمات المطلوبة من النص.	فهم التعليمات (توليدي)		
تصنيف نية المستخدم من استفسار أو طلب.	تحديد المقصود (توليدي)		
تحديد نوع الإشكالية أو الفئة التي وجدت في النص.	تحديد الإشكالية (توليدي)	تصنيف المحتوى النصي إلى فئات محددة اعتماداً على معانيه أو خصائصه.	
تحديد المشاعر التي يعبر عنها النص.	اكتشاف المشاعر (اختياري)		
تحديد ما إذا كانت الرسالة مزججةً أو معتادةً.	اكتشاف الرسائل المزججة (اختياري)		
اكتشاف ما إذا كان النص يحتوي على كراهية أو إساءة.	اكتشاف الإساءة (اختياري)		
اكتشاف ما إذا كانت الجملة تتضمّن سخريةً.	اكتشاف السخرية (اختياري)		
تحليل مشاعر شخصية بناءً على موقف أو جملة.	تحليل المشاعر (اختياري)		

تعريف المهمة	المهمة	تعريف الفئة	الفئة
تحديد نوع الجملة أو الغرض التواصلية منها، مثل: السؤال، أو الأمر، أو التصريح؛ اعتمادًا على معناها وسياقها.	تعرف أفعال الحوار (اختياري)		تصنيف النصوص
تقييم الأفعال أو السلوكيات وتصنيفها استنادًا إلى القيم والمعايير الأخلاقية العامة.	تصنيف السلوكيات الأخلاقية (اختياري)	تصنيف المحتوى النصي إلى فئات محددة اعتمادًا على معانيه أو خصائصه.	
تحديد الفئة المناسبة للسؤال بناءً على معناه ومحتواه.	تصنيف الأسئلة (اختياري)		
تصنيف محتوى النص بناءً على موضوعه.	تحديد الموضوع (اختياري)		
تحديد ما إذا كان النص يتضمن قولبة نمطية أو يعرض نمطًا مضافًا لها؛ اعتمادًا على تحليل المعنى والسياق.	اكتشاف القوالب النمطية (اختياري)		
تصحيح الأخطاء النحوية والإملائية في النص.	تصحيح القواعد النحوية (كلاهما)		تعديل النصوص
تعديل النص ليتوافق مع الهوية الجنسية (من الذكر للأنثى) والعكس.	تعديل الهوية الجنسية للنص (توليدي)	إجراء تغييرات لغوية أو دلالية على النصوص؛ لتحسينها أو تكيفها مع سياقات محددة.	
إعادة صياغة الجملة بأسلوب نحوي مختلف، مع الحفاظ على المعنى.	إعادة الصياغة (توليدي)		
إعادة صياغة سؤال غامض ليصبح واضحًا ومحددًا.	تعديل السؤال (توليدي)		
تبسيط النص المعقد؛ ليصبح أيسر فهمًا.	تبسيط النص (توليدي)		
ترجمة النصوص من لغة إلى أخرى.	الترجمة الآلية (كلاهما)	تحويل النصوص آليًا من لغة إلى أخرى، أو بين لهجات مختلفة.	الترجمة الآلية والنقل الصوتي
ترجمة لهجة عربية محلية إلى العربية الفصحى.	ترجمة اللهجات (توليدي)		

الفئة	تعريف الفئة	المهمة	تعريف المهمة
الترجمة الآلية والنقل الصوتي	تحويل النصوص آلياً من لغة إلى أخرى، أو بين لهجات مختلفة.	النقل الصوتي (توليدي)	تحويل النص العربي إلى الكتابة الصوتية.
الدقة المعلوماتية	إجراء تغييرات لغوية أو دلالية على النصوص؛ لتحسينها أو تكيفها مع سياقات محددة.	التحقق من صحة الادعاء (اختياري)	التحقق من أنّ ما ورد في النص من ادعاء مدعوم بأدلة واردة في النص نفسه.
		التحقق من صحة الإجابة (اختياري)	التحقق من أنّ الإجابة المقدّمة صحيحة وتتماشى مع المعلومات الواردة في نص أو محتوى معرفي معيّن.

## منهجية التقييم:

يعتمد هذا التقرير على منهجية تقييم متعددة المراحل، تهدف إلى قياس أداء النماذج اللغوية بدقة وتوازن؛ وذلك من خلال الجمع بين التقييم البشري والتقييم الآلي، واستخدام أنماط مختلفة من الأسئلة تعكس تنوع الاستخدامات الواقعية للنماذج اللغوية. فيما يلي نستعرض أنماط الأسئلة وآلية قياس النتائج وطريقة التقييم المعتمدة وتوصيات التقرير الثاني.

### أولاً: أنماط الأسئلة

تم تصميم مجموعات التقييم لتشمل نوعين رئيسيين من الأسئلة، هما:

#### الأسئلة التوليدية (إجابات مفتوحة)

تُستخدم لتقييم قدرة النموذج على توليد محتوى لغوي متماسك ودقيق، يشمل الفهم العميق للسؤال، وصياغة الإجابة، وجودة التعبير، ومدى توافق المخرجات مع الإجابة المرجعية.

#### الأسئلة الاختيارية (اختيار من متعدد)

تُستخدم لقياس دقة الفهم، والاستدلال، والقدرة على اختيار الإجابة الصحيحة من بين بدائل محددة؛ بما يعكس كفاءة النموذج في المهام المعقدة والمعيّنة.

تم اعتماد منهجية تقييم متعددة المراحل هدفت إلى تحديد الطريقة المناسبة لقياس أداء النماذج بدقة وموثوقية، شملت المرحلة الأولى تجربة عدد من المنهجيات الآلية، مثل: مقياس n-وحدة (n-gram) للأسئلة التوليدية، والدقة.

(Accuracy) للأسئلة الاختيارية. تمت مقارنة نتائج هذه المنهجيات بنتائج التقييم البشري؛ لاختيار الأداة الأقرب لمحاكاة الحكم البشري في تقييم جودة الإجابات.

## ثانيًا قياس النتائج

يعتمد قياس النتائج في هذا التقرير على آليات تقييم متعددة، تهدف إلى تحليل أداء النماذج اللغوية من زوايا مختلفة؛ بالجمع بين التقييم الكمي، والتقييم الآلي، والتقييم البشري. ويسهم هذا التكامل في توفير قياس أكثر دقة وموثوقية لأداء النماذج، والحد من الانحياز المرتبط بالاعتماد على أسلوب تقييم واحد فقط.

### (أ) التقييم الكمي

تم استخدام مقاييس كمية معيارية لقياس أداء النماذج اللغوية بموضوعية؛ بمقارنة مخرجات النماذج بالإجابات المرجعية المعتمدة. وقد طُبِّقت هذه المقاييس بما يتناسب مع طبيعة كل نوع من الأسئلة.

- الأسئلة التوليدية: استُخدمت المقاييس المعتمدة على ن-وحدة (n-gram)، التي تقيس درجة التشابه بين النص المُولَّد والنص المرجعي من حيث تطابق تسلسلات الكلمات.
- الأسئلة الاختيارية: تم الاعتماد على مقياس الدقة (Accuracy)؛ لقياس مدى صحة اختيار النموذج للإجابة الصحيحة من بين البدائل المتاحة.

### (ب) التقييم الآلي

استخدام مقيّم آلي يُقدّر جودة مخرجات النماذج اللغوية؛ بمقارنة إجابة النموذج بالإجابة المرجعية، وذلك لكل من الأسئلة الاختيارية والتوليدية.

- درجة (٠): الإجابة غير صحيحة تمامًا.
- درجة (١): إجابة ضعيفة تحتوي على أخطاء كبيرة أو غير مكتملة بدرجة كبيرة، وتكاد لا تشبه الإجابة المرجعية
- درجة (٢): إجابة مقبولة، ولكنها تختلف بوضوح عن الإجابة المرجعية
- درجة (٣): إجابة صحيحة تمامًا أو شبه كاملة.

وبعد احتساب درجات التقييم الآلي لكل من الأسئلة الاختيارية والتوليدية بطريقة مستقلة، تُدمج النتائج النهائية وفقًا لأوزان نسبية تعكس أهمية كل نمط من أنماط الأسئلة، وذلك على النحو الآتي:

- (٣٠٪) من الدرجة النهائية مخصصة لنتائج الأسئلة الاختيارية.
  - (٧٠٪) من الدرجة النهائية مخصصة لنتائج الأسئلة التوليدية.
- ويهدف هذا التوزيع إلى إعطاء أولوية أكبر للقدرات التوليدية للنماذج اللغوية؛ نظرًا إلى دورها الأساس في التطبيقات الواقعية التي تتطلب فهمًا عميقًا وإنتاجًا لغويًا عالي الجودة.

## ج) التقييم البشري:

أجرى تقييم بشري على عينة مكونة من (٢٠٠) سؤالٍ موزعة على المهام المختلفة، وقد تولى ثلاثة محكمين تقييم كل إجابة وفقًا لأربعة مقاييس رئيسية، واحتُسب المتوسط النهائي لدرجاتهم لكل سؤال، ثم لكل نموذج.

## معايير التقييم البشري:

- **الدقة (٠ إلى ٣):** تشير إلى مدى صحة الإجابة وتطابقها مع المطلوب.
- **الشمولية (٠ إلى ٣):** تقيس مدى اكتمال الإجابة وتغطيتها للجوانب المهمة في السؤال.
- **الهلوسة (٠ أو ١):** تشير إلى وجود معلومات مُختلقة أو غير صحيحة.
- **الإطناب (٠ أو ١):** تقيس مدى الإطالة أو الخروج عن الموضوع.

## ثالثًا: طريقة التقييم المعتمدة

بناءً على المقارنة بين المنهجيات المختلفة؛ تم اعتماد التقييم الآلي كونه المنهجية الأقرب لمحاكاة التقييم البشري؛ إذ أظهر علاقة طردية قوية بلغت (٠.٩) مع نتائج المحكمين. في المقابل، أظهرت المقاييس الكمية التقليدية علاقةً عكسيةً مع التقييم البشري؛ مما يشير إلى محدودية قدرتها على عكس الحكم البشري بدقة في تقييم جودة الإجابات. والجدير بالذكر أن منهجية التقييم قابلة للتحديث مستقبلاً؛ بما يتناسب مع طبيعة المهام المختلفة والتقارير القادمة، مع الاستمرار في مقارنة أدوات التقييم الآلي بالأحكام البشرية؛ لضمان دقة النتائج وموثوقيتها.

## رابعًا: توصيات التقرير الثاني

- **أشكال المهام:** تتفاوت المهام بين مهام مفتوحة (توليدية إبداعية: مثل التلخيص، واستكمال القصص، ونظم القصائد، وغيرها)، وبين مهام مغلقة (لها إجابة واحدة صحيحة: مثل الإعراب، وتحديد اللهجة، وغيرها). ويوصي فريق بلسم بأهمية النظر في آلية مبتكرة لتقييم المهام الإبداعية سواء كان ذلك بتقسيمها إلى مهام حتمية (مثل اكتشاف البيت المكسور من أجل تقييم نظم القصائد)، أو كان عبر المقارنة بين النماذج (بأن تعرض نتائج نموذجين، ويتم المقارنة بينهما).
- **صعوبة المهام:** من الملاحظ أن عددًا من المهام التي قيمت النماذج اللغوية عليها لم تعد صعبة على النماذج الحديثة، مما يستدعي النظر أكثر في مهام أكثر صعوبة. من المهام التي فشلت فيها النماذج بوضوح توليد القصائد، لذا فهي تستدعي مزيدًا من البحث.
- **علاقة طردية:** لاحظنا تقاربًا وارتباطًا جيدًا بين التقييم البشري وتقييم النموذج المحكم (الذي يقوم بالمقارنة بين الإجابة النموذجية وإجابة النموذج اللغوي)، وهو ما يزيد من موثوقية التقييم الآلي. هذا التقارب أفضل في المهام المغلقة، ولكنه أقل في المهام المفتوحة.

• **النموذج المحكم:** على الرغم من أن النموذج المحكم لا يجب أن يطلع على السؤال، وإنما على الإجابة النموذجية وإجابات النماذج المختلفة، إلا أنه فاق جميع النماذج الأخرى عند تقييمه، مما استلزم استبعاده من القائمة لأجل العدالة بين النماذج. ويوصي فريق بلسم باعتماد نموذج مفتوح المصدر ليكون النموذج المحكم (خاصة بعد النظر في نتائج هذا التقرير).

• **المستويات اللغوية:** ركز التقرير الأول والثاني من مؤشر بلسم على المهام النصية فقط في النماذج اللغوية. ولا شك أن المهام الصوتية والبصرية هي جزء أساسي من لغتنا العربية (مثل القدرة على نطق القوائد بالشكل الصحيح، وصناعة الصور بطابع ثقافي عربي، وغيرها). ونأمل في التقارير القادمة أن تتضمن النتائج تقييمًا لهذه المهام المختلفة.

• **البحث والتطوير في التخصصات المختلفة:** ما زالت آليات التقييم تتطلب مزيدًا من الأبحاث العلمية خاصة فيما يتعلق بالمهام العربية، مثل مجال اللهجات والتفاوت فيما بينها، وقياس قدرات النماذج في فهم الثقافة العربية، ومجال البلاغة والشعر والأدب. ولم يعد مجال التقييم حكرًا على أصحاب التخصصات التقنية. وندعو الباحثين من التخصصات الإنسانية والشرعية وغيرها، إلى المشاركة في مؤشر بلسم بإضافة مهام نوعية عربية.

## ✓ النماذج المُقيّمة في التقرير الثاني

اعتمد التقرير في اختيار النماذج المُقيّمة على منهجية واضحة، تهدف إلى ضمان الاستمرارية والاتساق في التقييم. وقد شمل الاختيار أغلب النماذج التي قُيِّمت في التقرير السابق، مع اعتماد أحدث إصداراتها المتاحة، إلى جانب إدراج نماذج جديدة طُرحت خلال العدة اللاحقة، وقد رُوّعي في عملية الاختيار تمثيل تنوع النماذج من حيث أحجامها وبنائها المعمارية، مع تضمين كل من النماذج المفتوحة المصدر والمغلقة المصدر؛ بما يضمن توافقها مع متطلبات التقرير ودعمها للغة العربية، أو عملها ضمن إطار متعدد اللغات.

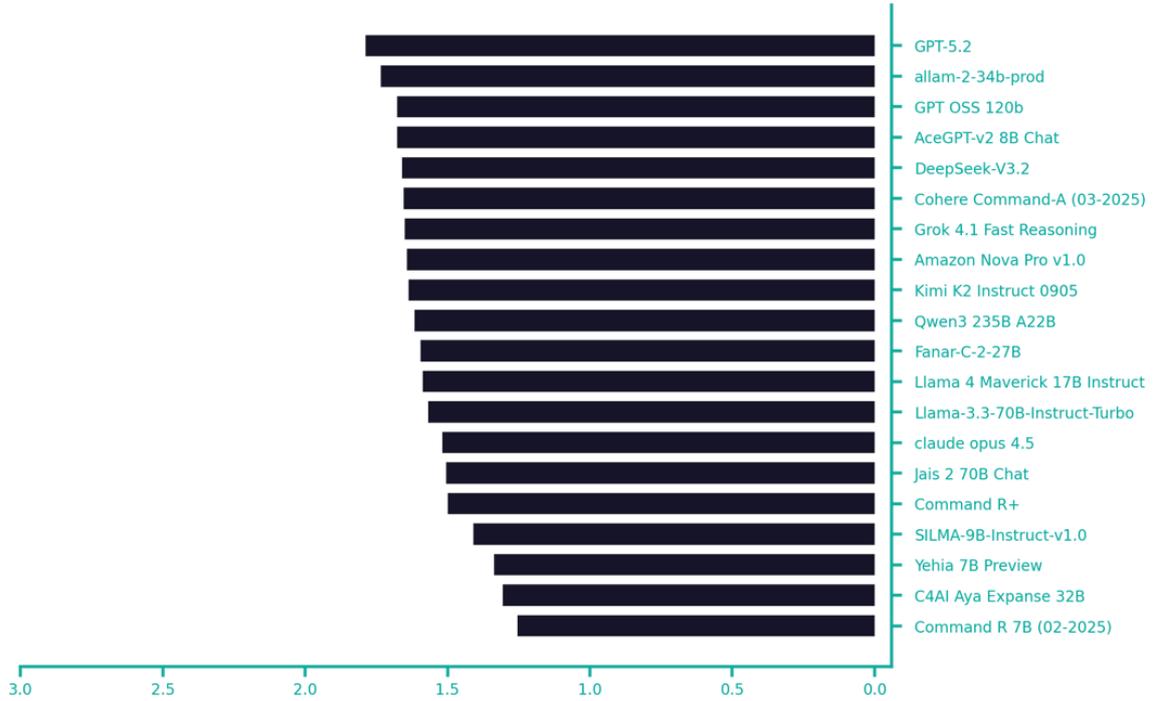
اسم النموذج	المصدر	الحجم (عدد المعاملات)	الوصف	المرجع/ الرابط
AceGPT-v8-8B-Chat	مفتوح المصدر	8B	نموذج محاكاة باللغة العربية، مكون من (٨) مليارات معام، مبني على (LLaMA٢)، ومخصص للتفاعلات النصية	<a href="#">المرجع</a> <a href="#">الرابط</a>
c4ai-aya-expanse-32b	مفتوح المصدر	32B	نموذج متعدد اللغات، مكون من (٣٢) مليار معمل، ويدعم (٢٣) لغةً من بينها العربية.	<a href="#">المرجع</a> <a href="#">الرابط</a>
Cohere (Command-A 03- 2025)	مفتوح المصدر	111B	نموذج لغوي كبير، يحتوي على حوالي (١١١) مليار معمل، ويدعم (٢٣) لغةً.	<a href="#">المرجع</a> <a href="#">الرابط</a>

<a href="#">الرابط</a>	نموذج متعدد اللغات، مكون من (١٠٤) مليارات معام، ومخصص للتوليد المعزز بالاسترجاع، وتنفيذ التعليمات.	104B	مفتوح المصدر	Command R+
<a href="#">الرابط</a>	نموذج لغوي متعدد اللغات، يتكون من ٧ مليارات معام.	7B	مفتوح المصدر	Command-r7b 2024-12
<a href="#">المرجع</a> <a href="#">الرابط</a>	نموذج متعدد الخبراء، مكون من (٦٨٥) مليار معمل يدعم مهام الفهم اللغوي، والبرمجة، والتفكير المنطقي.	685B MoE	مفتوح المصدر	DeepSeek V3
<a href="#">المرجع</a> <a href="#">الرابط</a>	نموذج ثنائي اللغة (عربي-إنجليزي)، مكون من (٧٠) مليار معام.	70B	مفتوح المصدر	Jais-70-2B-Chat
<a href="#">المرجع</a> <a href="#">الرابط</a>	نموذج لغوي كبير قائم على تعدد الخبراء، ومُدرَّب على اتباع التعليمات، ومهيأ لمهام الاستدلال والبرمجة بسياق طويل.	1T MOE	مفتوح المصدر	Kimi K2-Instruct-0905
<a href="#">المرجع</a> <a href="#">الرابط</a>	نموذج لغوي كبير صمّمته (Meta)، ويحتوي على (٧٠) مليار معام، ويدعم عدة لغات.	70B	مفتوح المصدر	Llama-70-3.3B-Instruct-Turbo
<a href="#">المرجع</a> <a href="#">الرابط</a>	نموذج لغوي كبير صمّمته (Meta)، ويحتوي على (١٧) مليار معام، ويدعم عدة لغات.	17B	مفتوح المصدر	Llama 4 Maverick 17B Instruct
<a href="#">المرجع</a> <a href="#">الرابط</a>	نموذج لغوي مكون من (٣٣٥) مليار معام.	235B	مفتوح المصدر	Qwen235 3B A22B
<a href="#">المرجع</a> <a href="#">الرابط</a>	نموذج لغوي باللغة العربية، مكون من (٩) مليارات معام، ومبني على بنية (Gemma).	9B	مفتوح المصدر	SILMA9-B Instruct-v1.0
<a href="#">المرجع</a> <a href="#">الرابط</a>	نموذج ثنائي اللغة (عربي-إنجليزي) مكون من (٧) مليارات معام.	7B	مفتوح المصدر	Yehia7-B preview

اسم النموذج	المصدر	الحجم (عدد المعاملات)	الوصف	المرجع/ الرابط
Fanar-C-27-2B	مفتوح المصدر	27B	نموذج ثنائي اللغة (عربي-إنجليزي) مدرب على ما يقرب (٢٧) مليار معاملة.	<a href="#">الرابط</a>
gpt-oss-120b	مفتوح المصدر	120B	نموذج متعدد اللغات، مكون من (١٢٠) مليار معاملة، وطوّرتة (OpenAI).	<a href="#">الرابط</a>
Amazon Nova Pro	مغلق المصدر	غير معلن	نموذج من (Amazon Bedrock).	<a href="#">الرابط</a>
Grok-4.1 Fas Reasoning	مغلق المصدر	غير معلن	نموذج متعدد اللغات، ويدعم الاستدلال المتقدم، وطوّرتة (xAI).	<a href="#">الرابط</a>
Claude Opus 4.5Instruct	مغلق المصدر	>130B	نموذج متعدد اللغات من (Anthropic).	<a href="#">الرابط</a>
GPT 5.2	مغلق المصدر	غير معلن	نموذج من (OpenAI)، يدعم مدخلات متعددة الوسائط واللغات.	<a href="#">الرابط</a>
ALLaM34 – 2B	مغلق المصدر	غير معلن	نموذج ثنائي اللغة (عربي-إنجليزي) مدرب على ما يقرب (٢٧) مليار معاملة.	<a href="#">الرابط</a>

## أولاً: متوسط الأداء العام للنماذج:

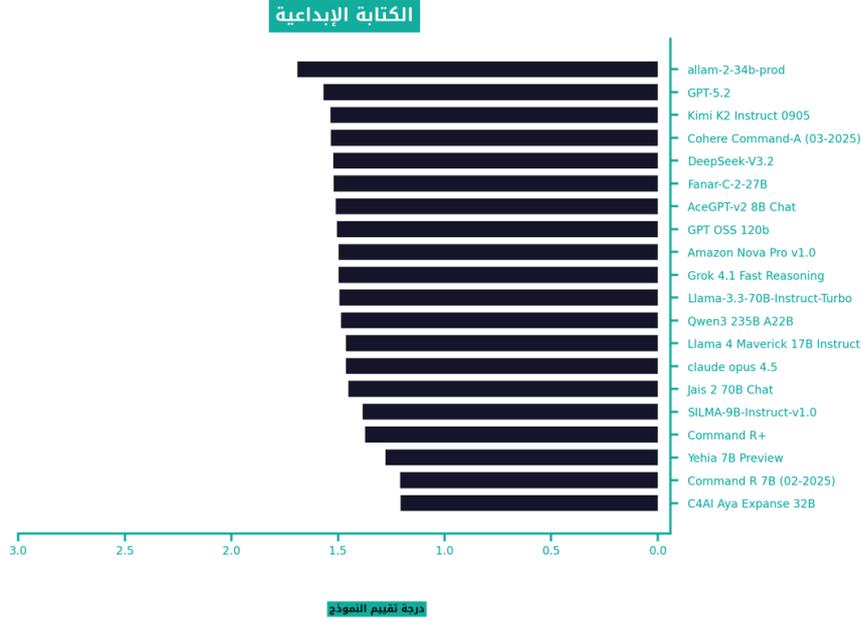
### المتوسط العام للنماذج



### درجة تقييم النموذج

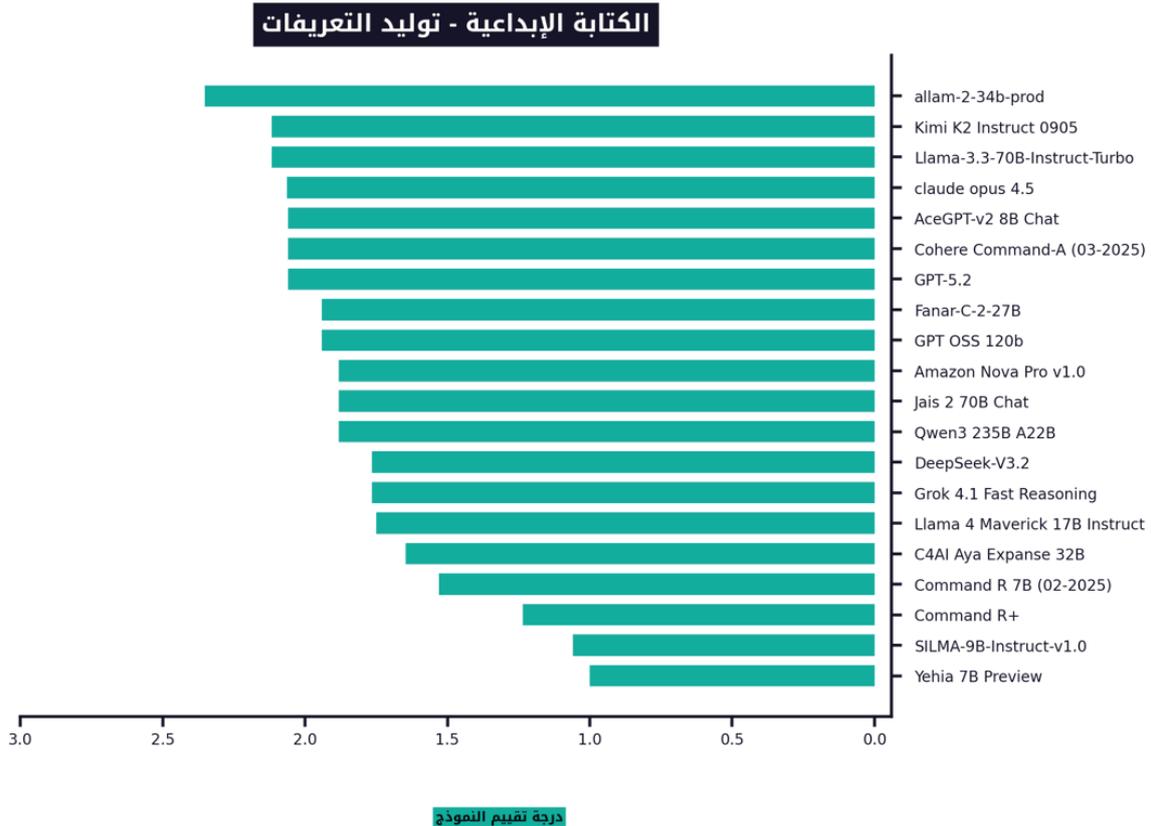
## ثانيًا: تقييم النماذج حسب الفئات، ومهام كل فئة:

### ١. الكتابة الإبداعية (Creative Writing).

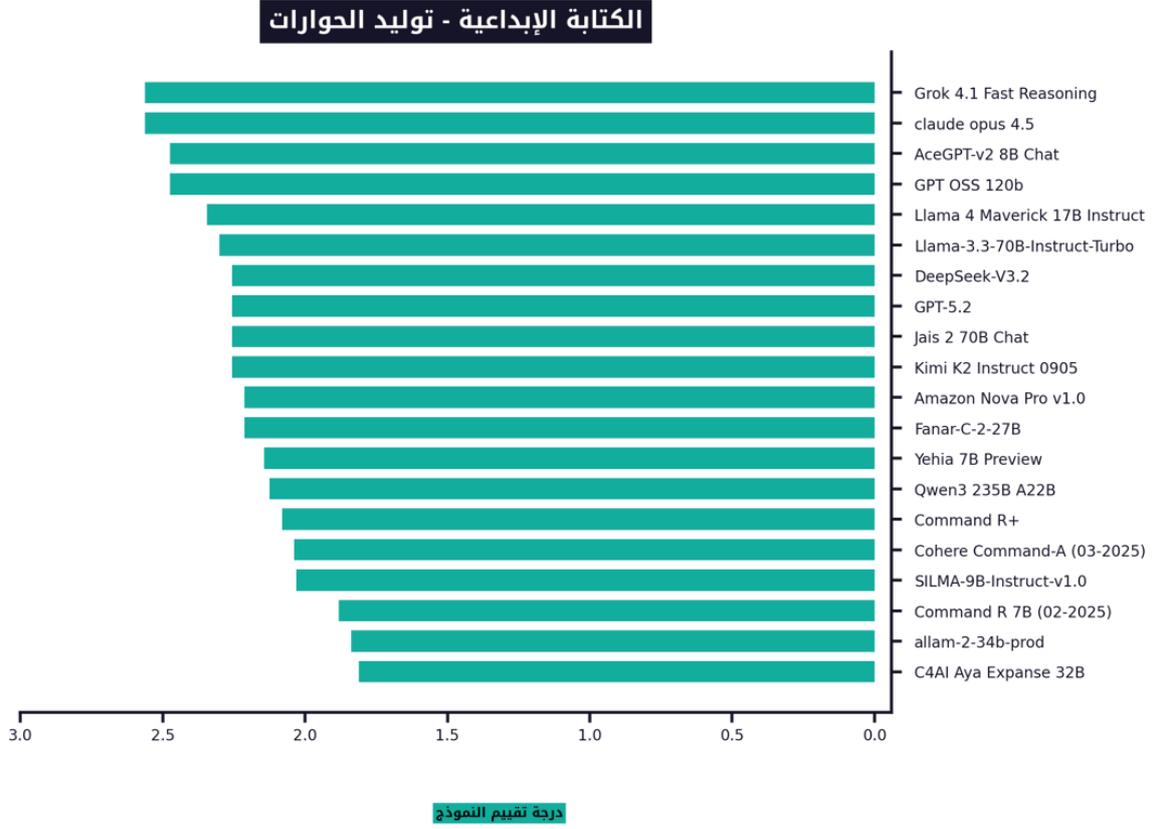


## المهام الفرعية:

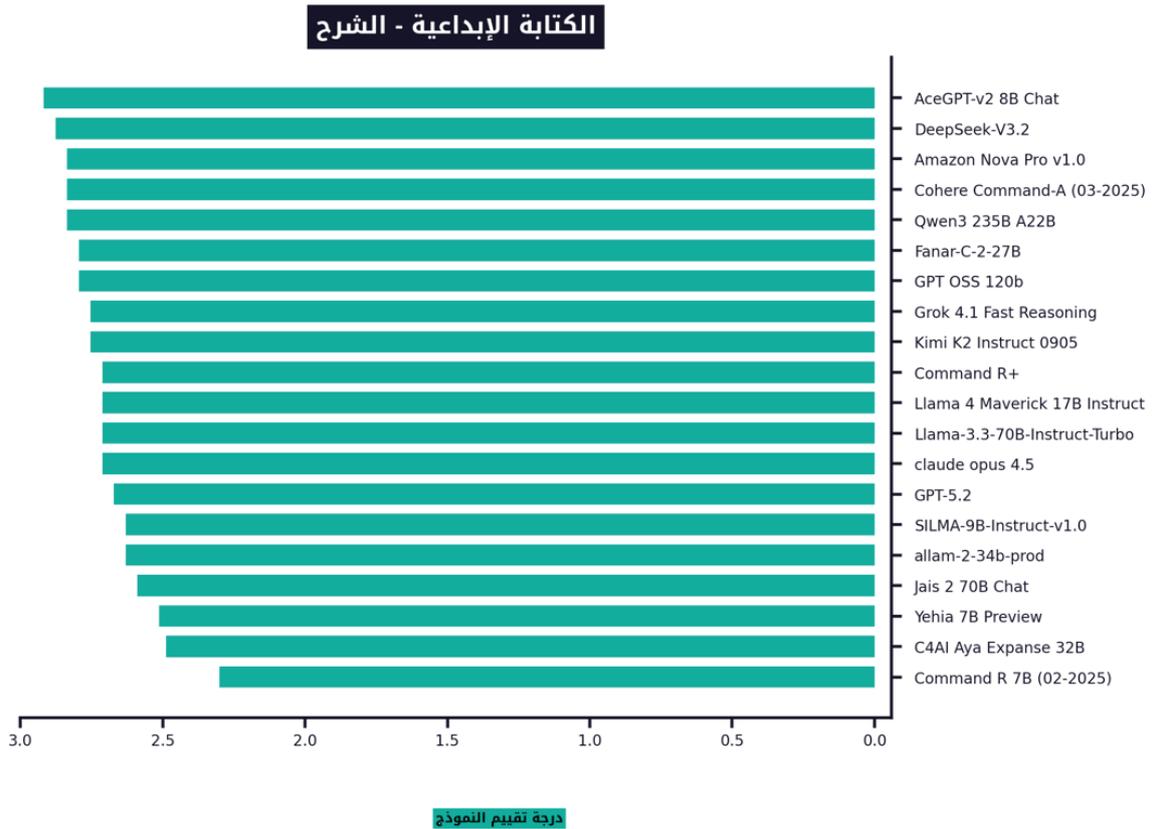
### ١.١ توليد التعريفات (Definition Generation).



## ١.٢ توليد الحوارات (Dialogue Generation).

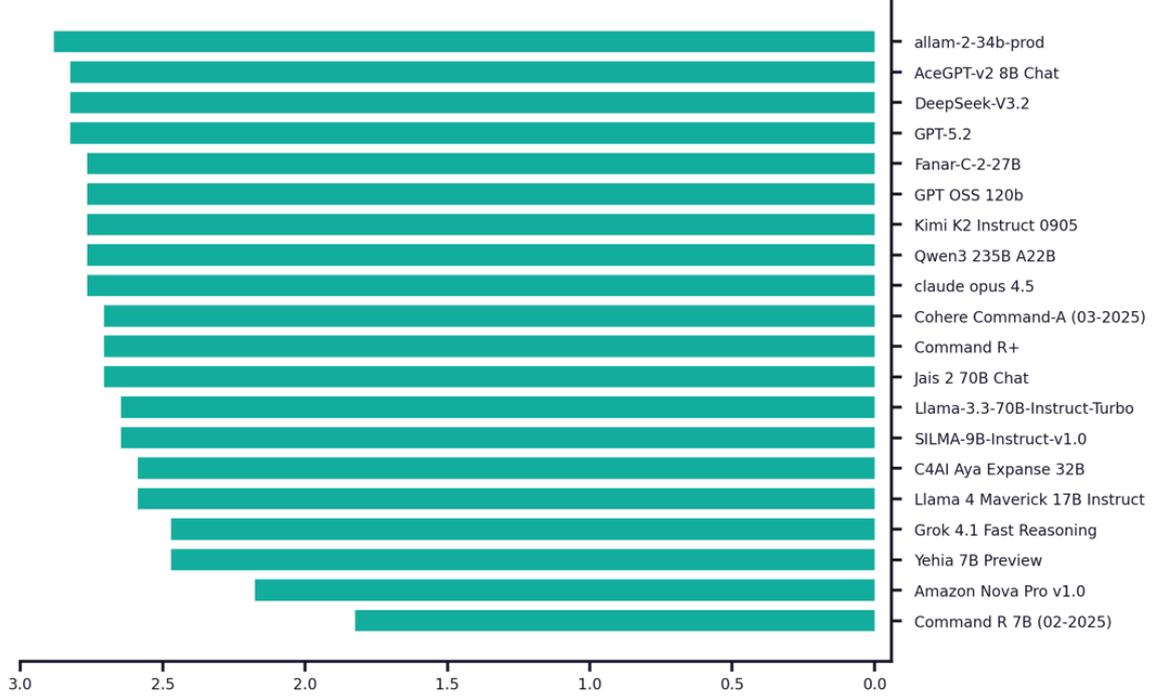


## ١.٣ الشرح (Explanation).



## ١.٤ توليد التعليمات (Instruction Generation).

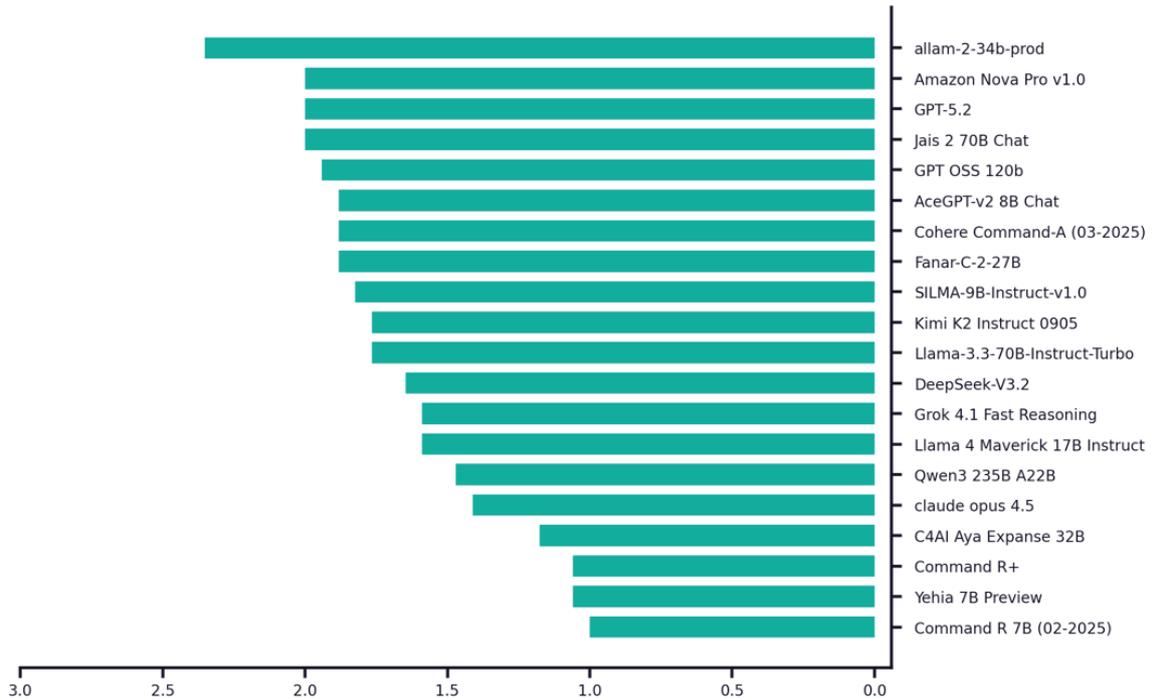
### الكتابة الإبداعية - توليد التعليمات



درجة تقييم النموذج

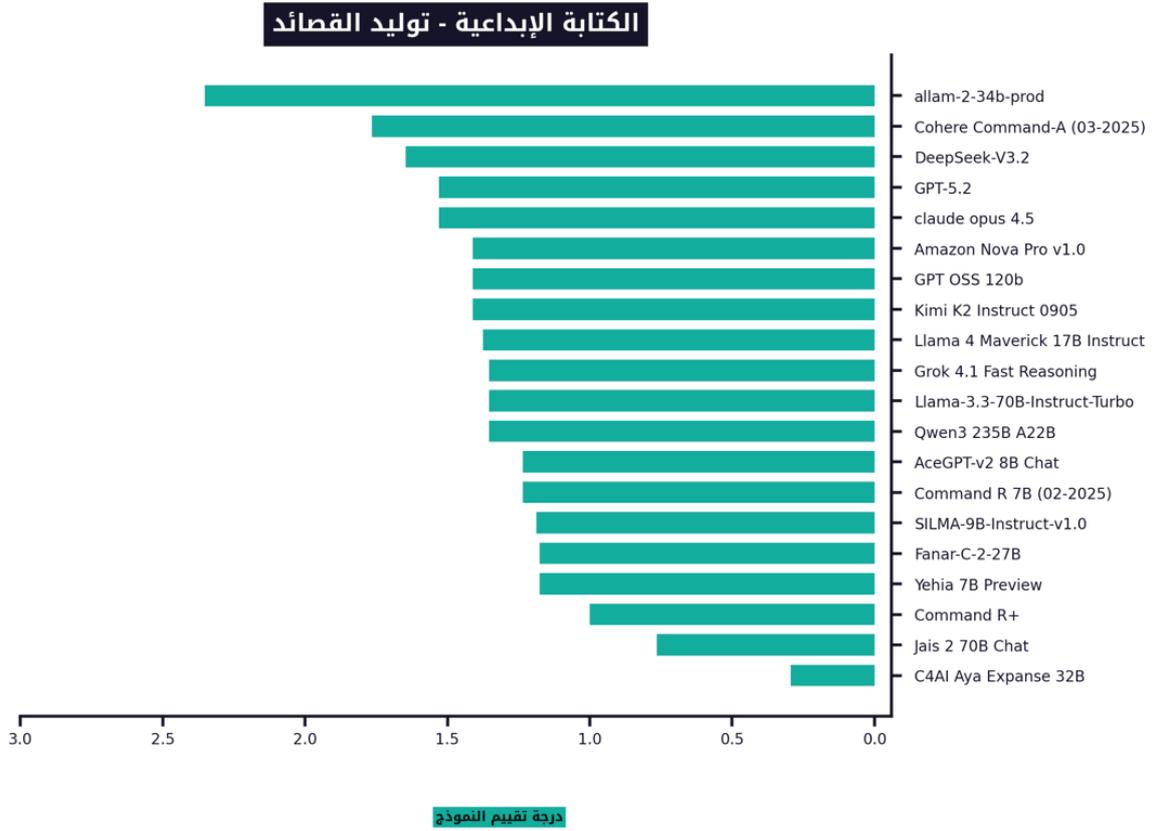
## ١.٥ توليد المقالات الإخبارية (News Article Generation).

### الكتابة الإبداعية - توليد المقالات الإخبارية

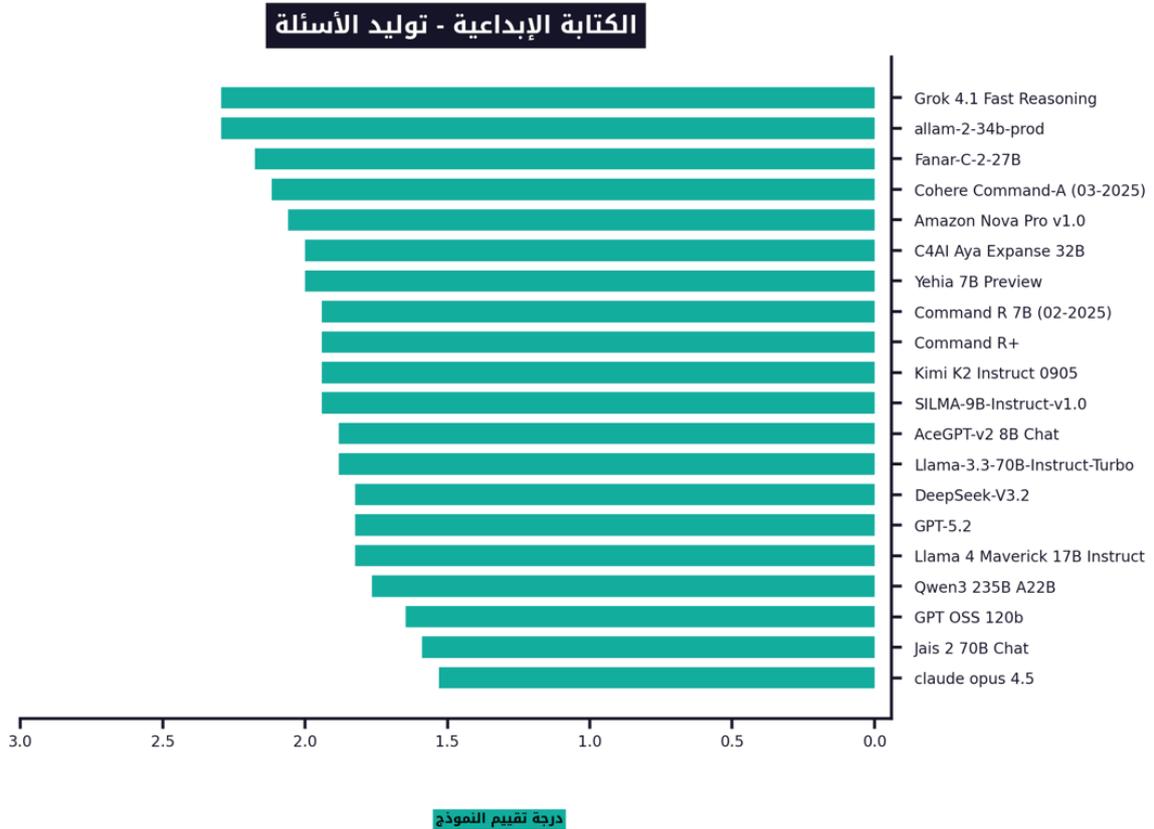


درجة تقييم النموذج

## ١.٦ توليد القصائد (Poem Generation).

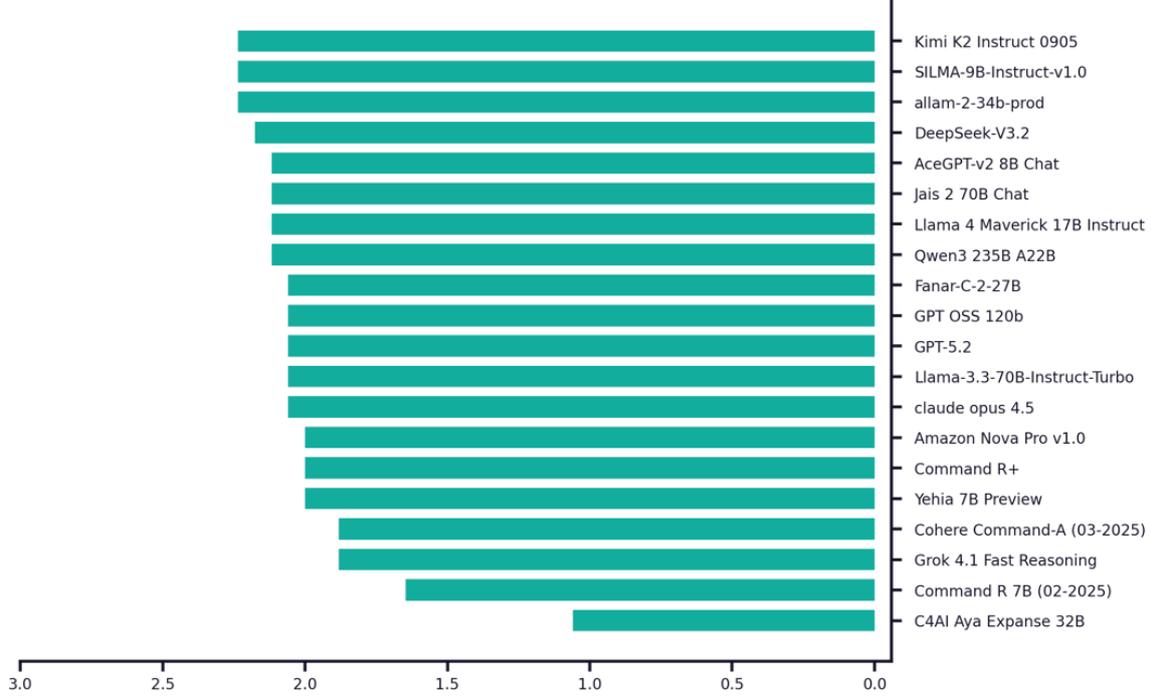


## ١.٧ توليد الأسئلة (Question Generation).



## ١.٨ تأليف الجمل (Sentence Composition).

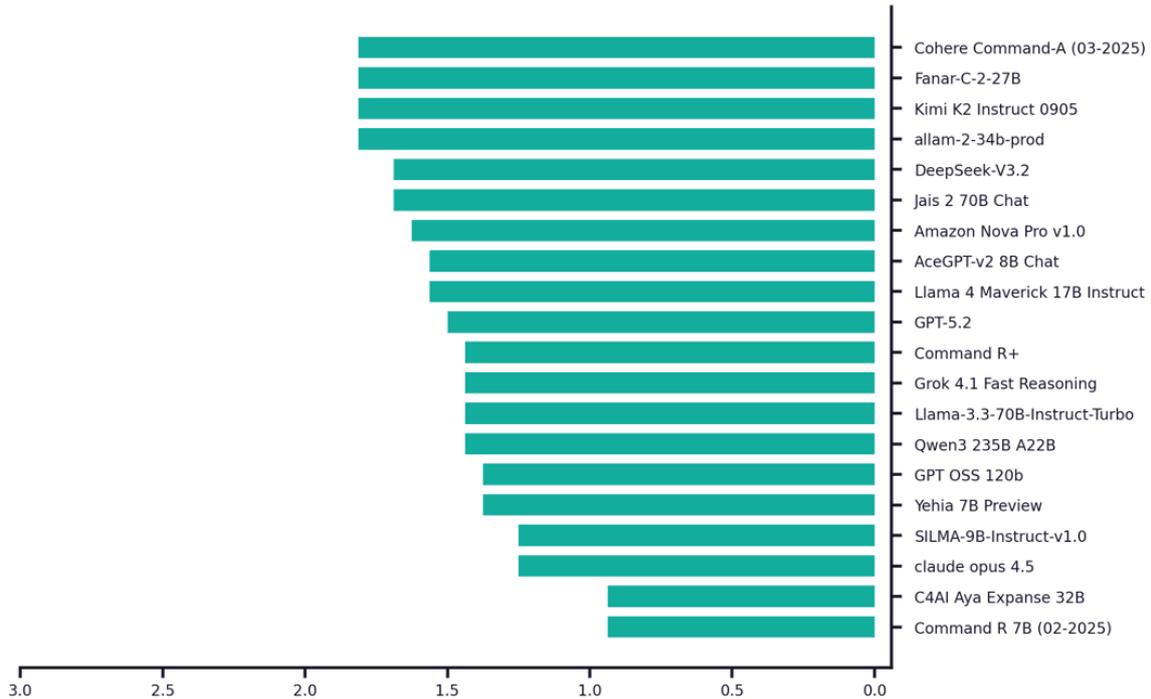
### الكتابة الإبداعية - تأليف الجمل



درجة تقييم النموذج

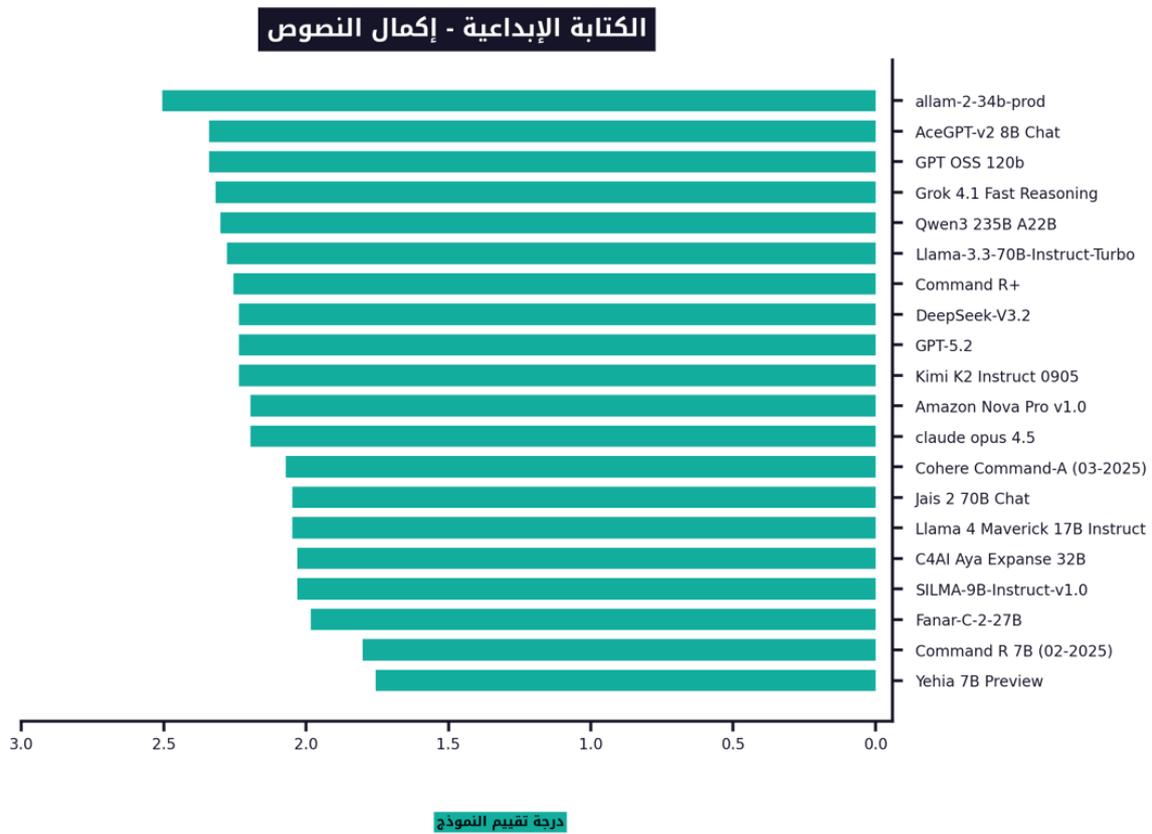
## ١.٩ تأليف القصص (Story Composition).

### الكتابة الإبداعية - تأليف القصص

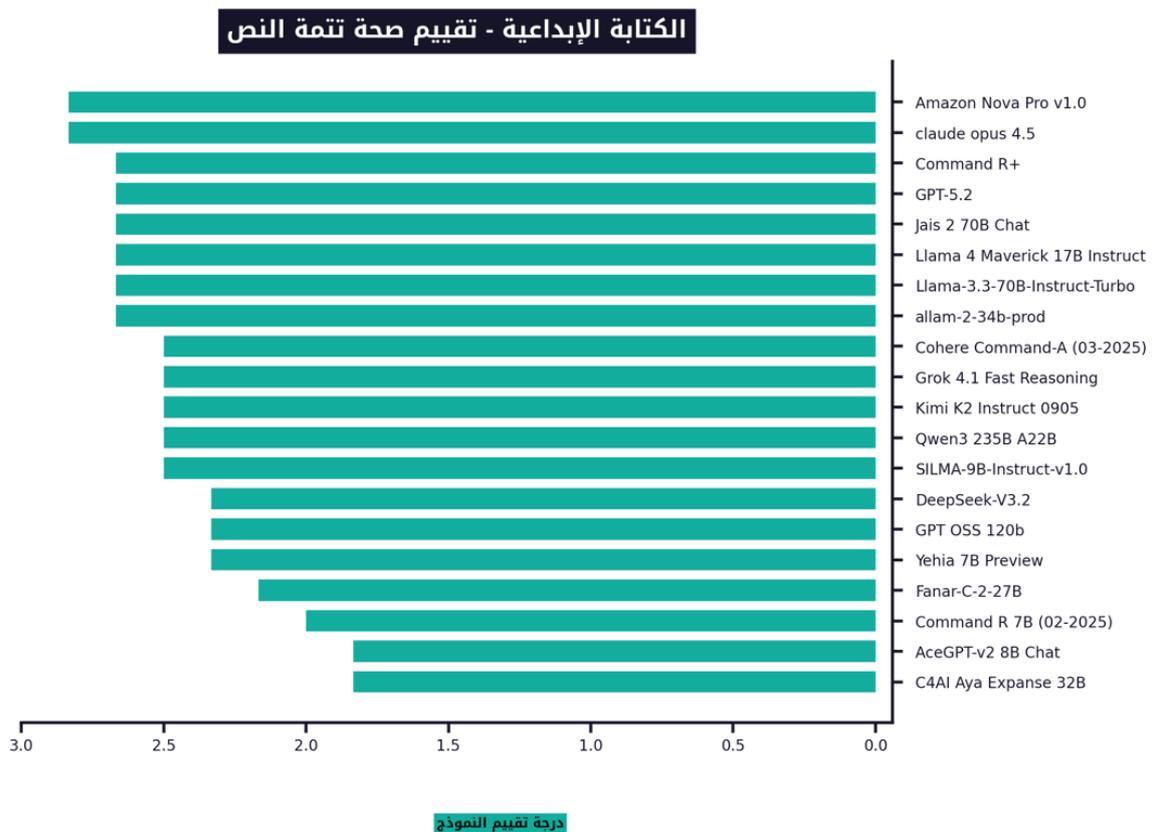


درجة تقييم النموذج

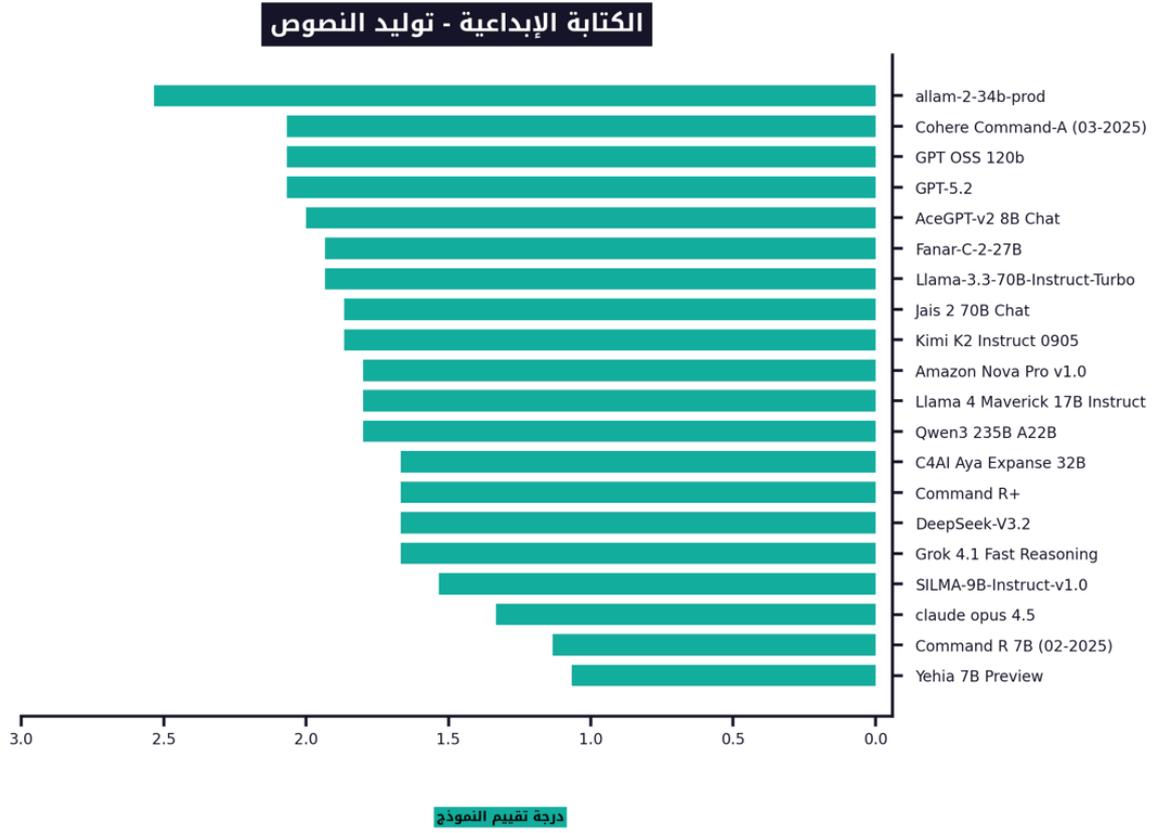
## ١.١٠ إكمال النصوص (Text Completion).



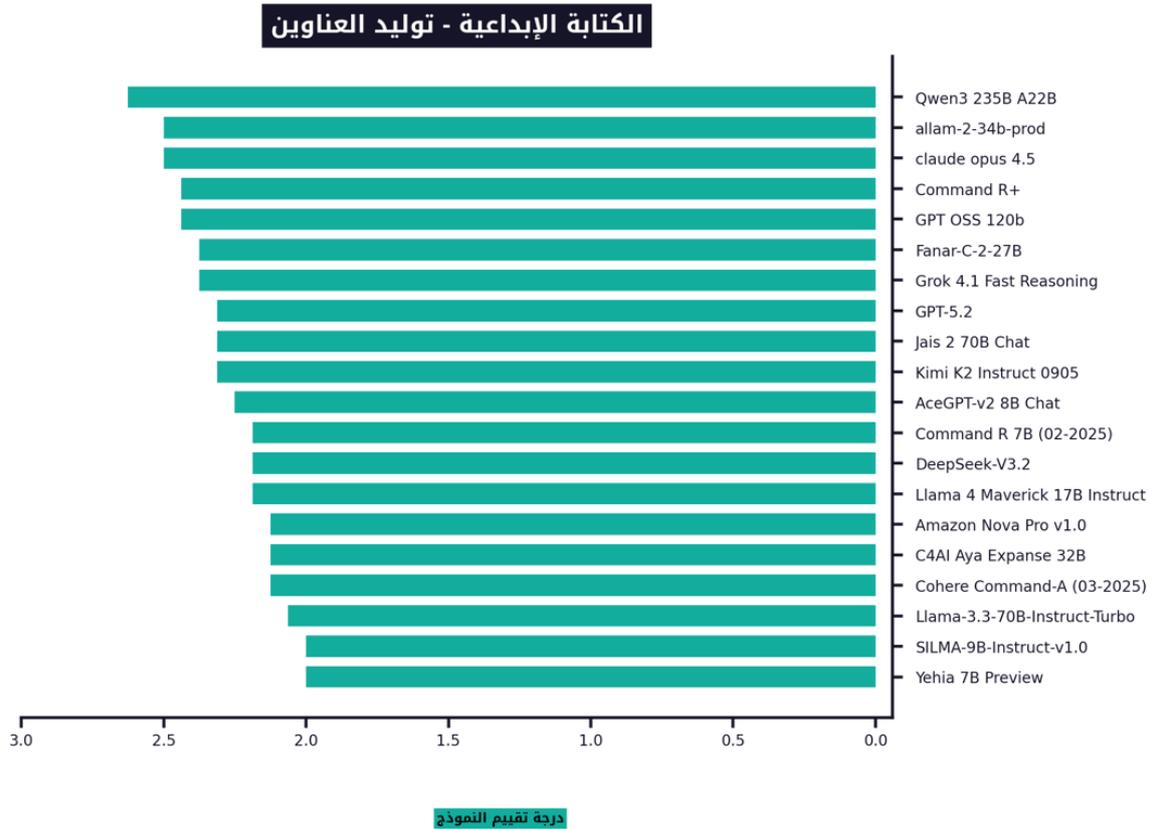
## ١.١١ تقييم صحة تتمة النص (Text Continuation Evaluation).



## ١.١٢ توليد النصوص (Text Generation).

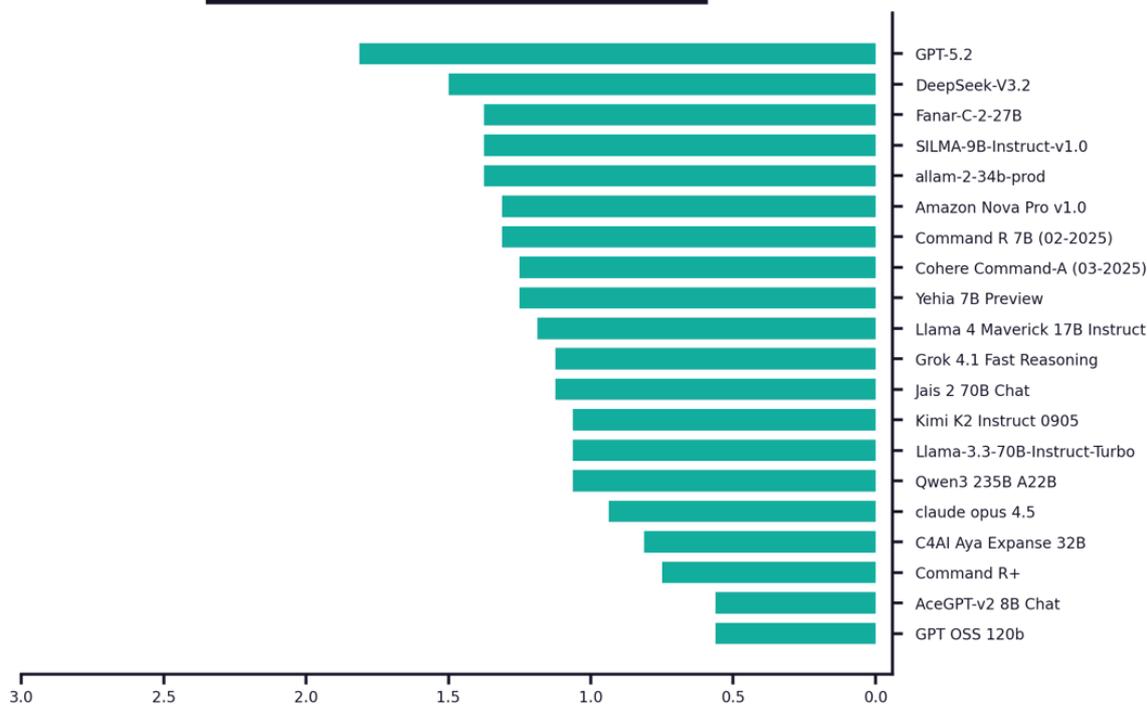


## ١.١٣ توليد العناوين (Title Generation).



## ١.١٤ توليد إجابات غير مناسبة (Wrong Candidate Generation).

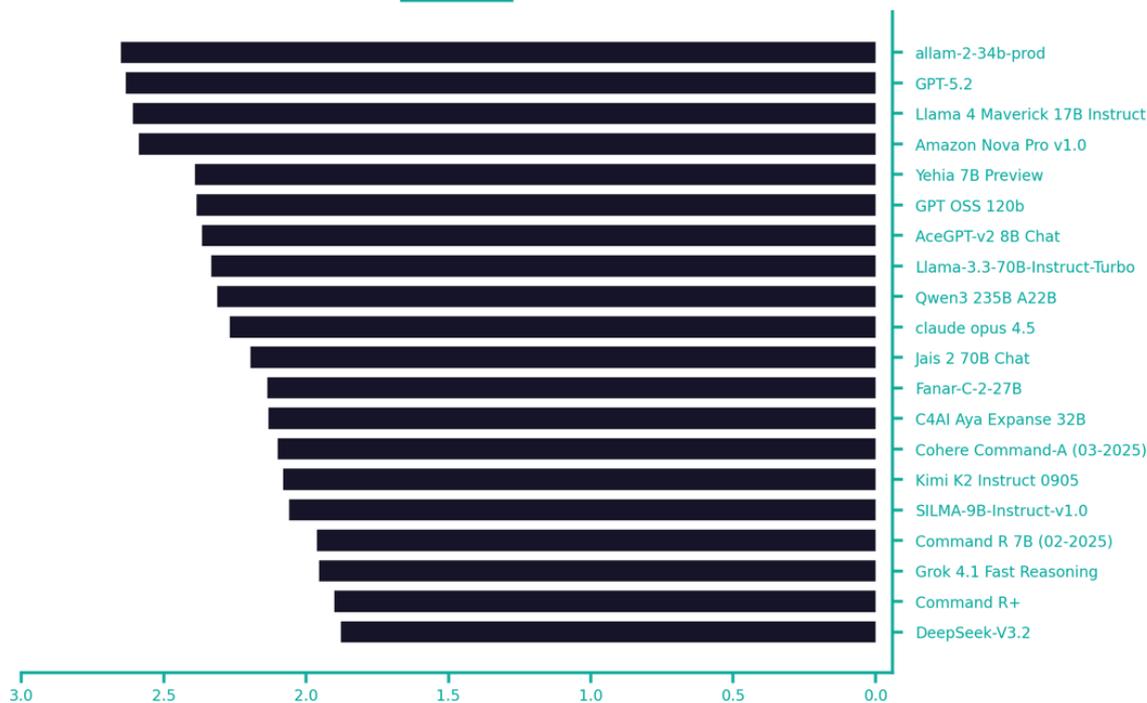
### الكتابة الإبداعية - توليد إجابات غير مناسبة



درجة تقييم النموذج

## ٢. التضمين (Entailment).

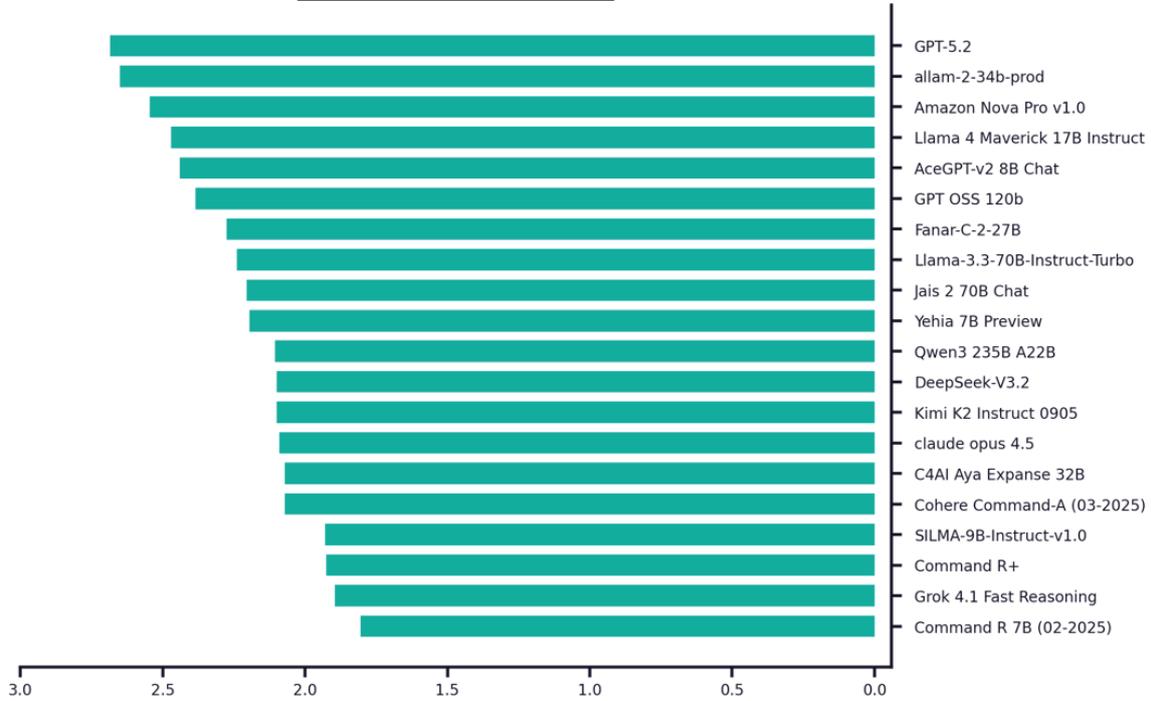
### التضمين



درجة تقييم النموذج

## ٢.١ التشابه الدلالي (Semantic Similarity).

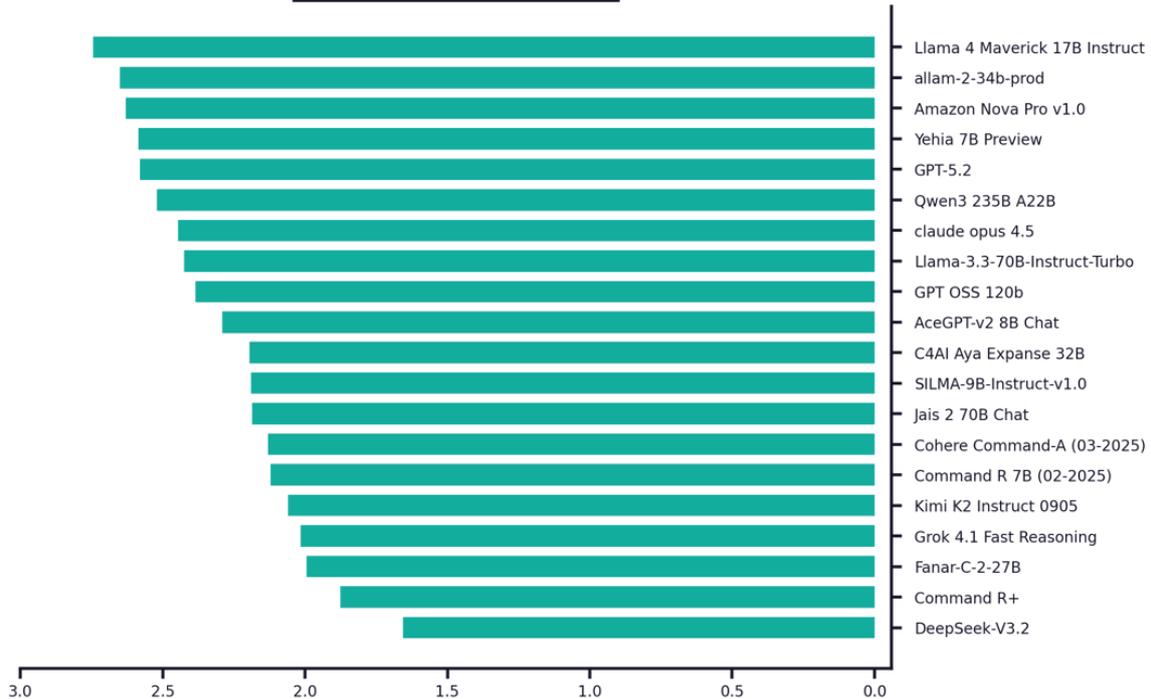
### التضمين - التشابه الدلالي



درجة تقييم النموذج

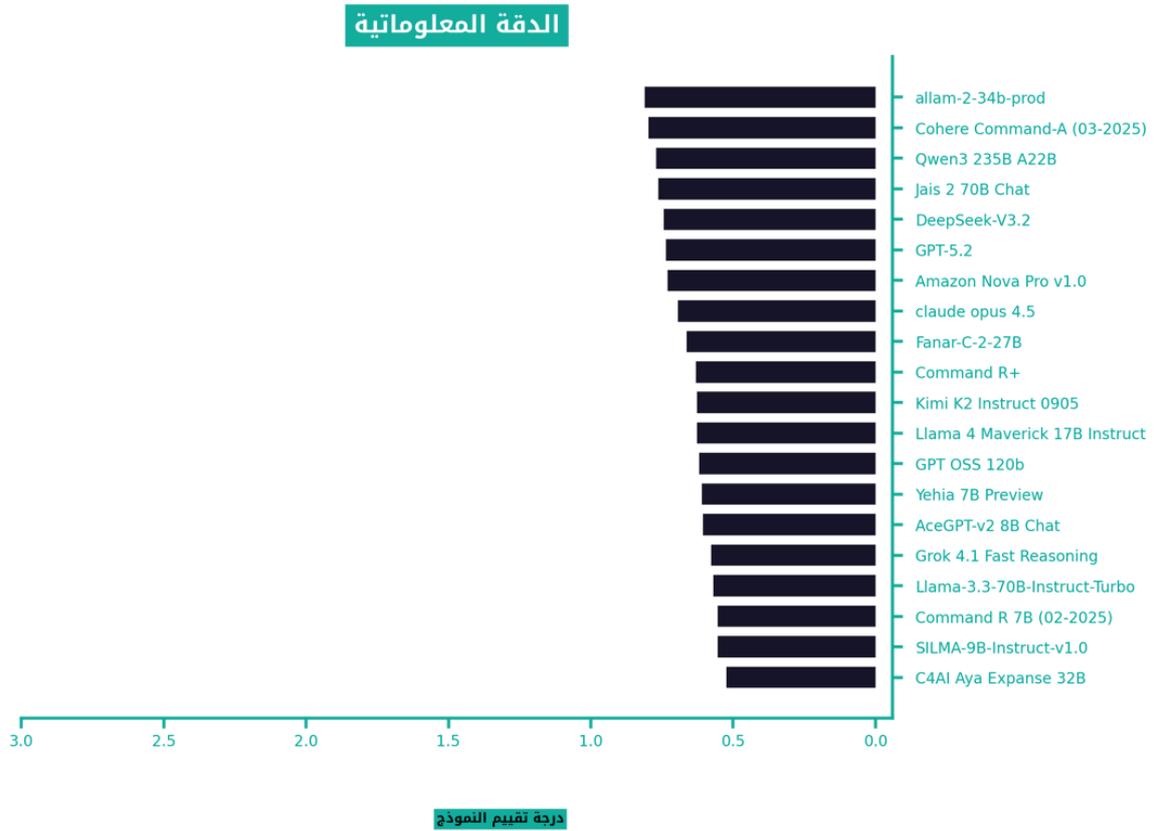
## ٢.٢ الاستلزام النصي (Textual Entailment).

### التضمين - الاستلزام النصي

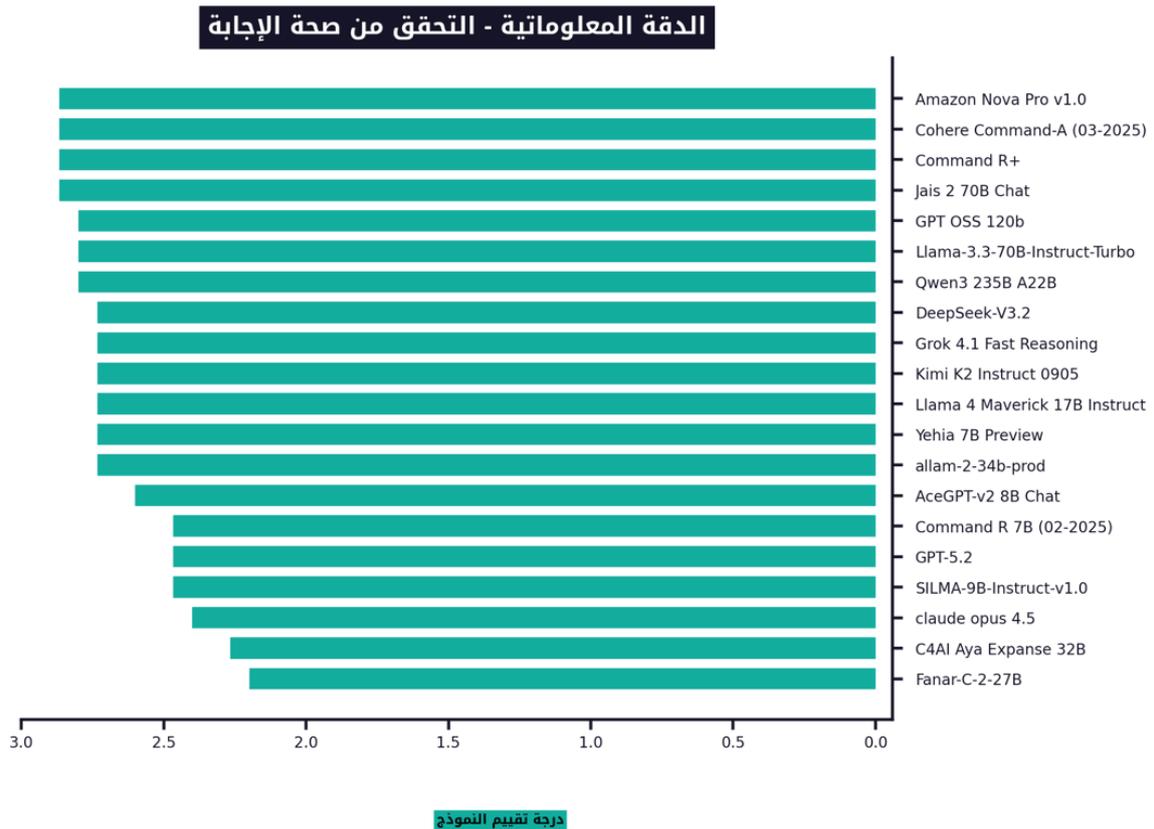


درجة تقييم النموذج

### ٣. الدقة المعلوماتية (Factuality).

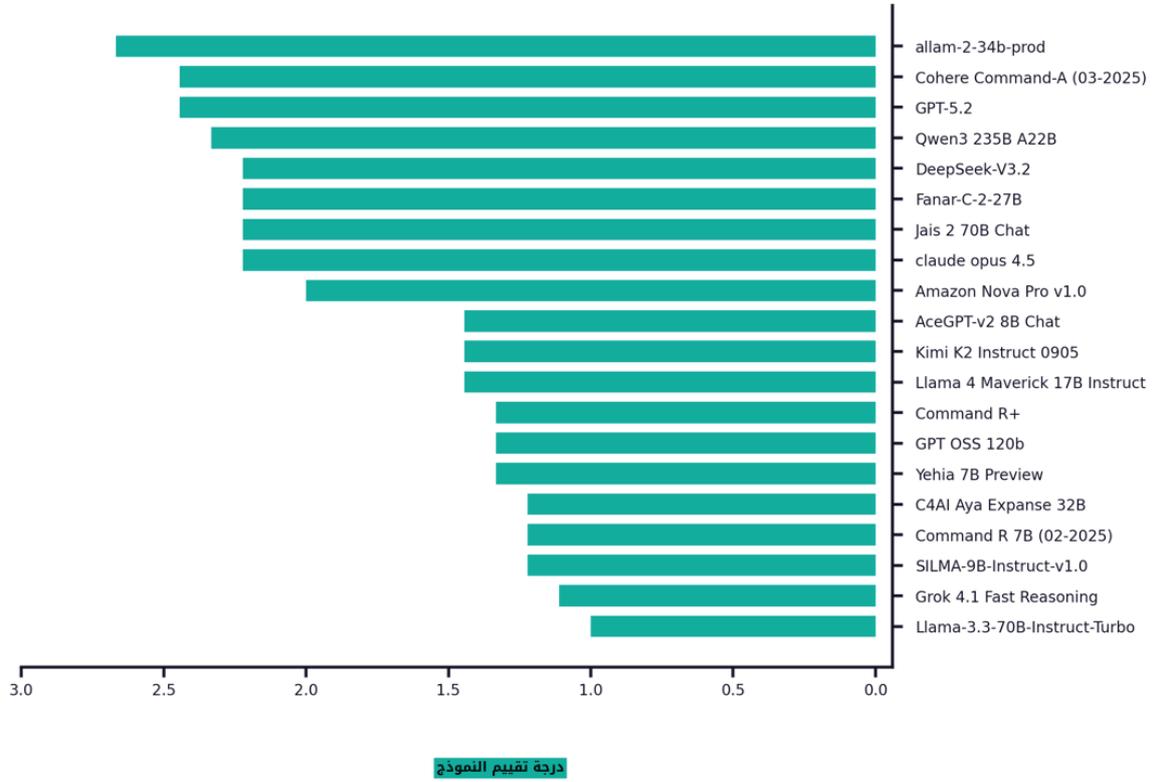


### ٣.١ التحقق من صحة الإجابة (Answer Verification).

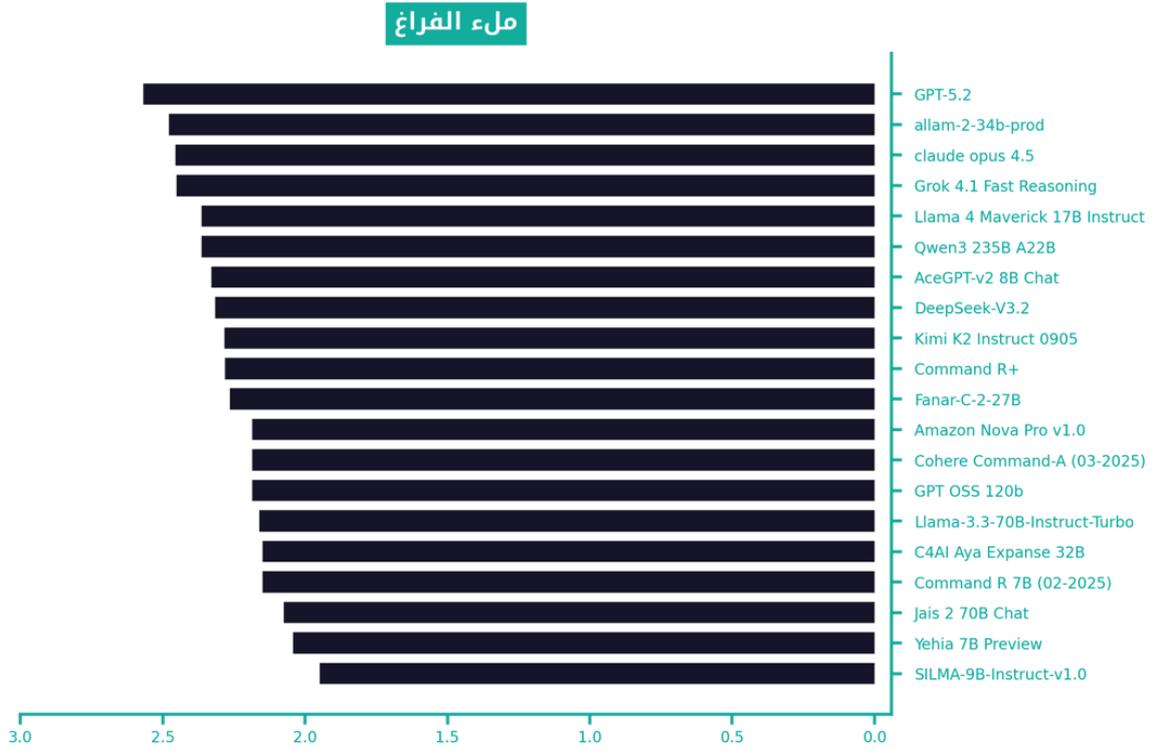


## ٣.٢ التحقق من صحة الادعاء (Claim Verification).

### الدقة المعلوماتية - التحقق من صحة الادعاء

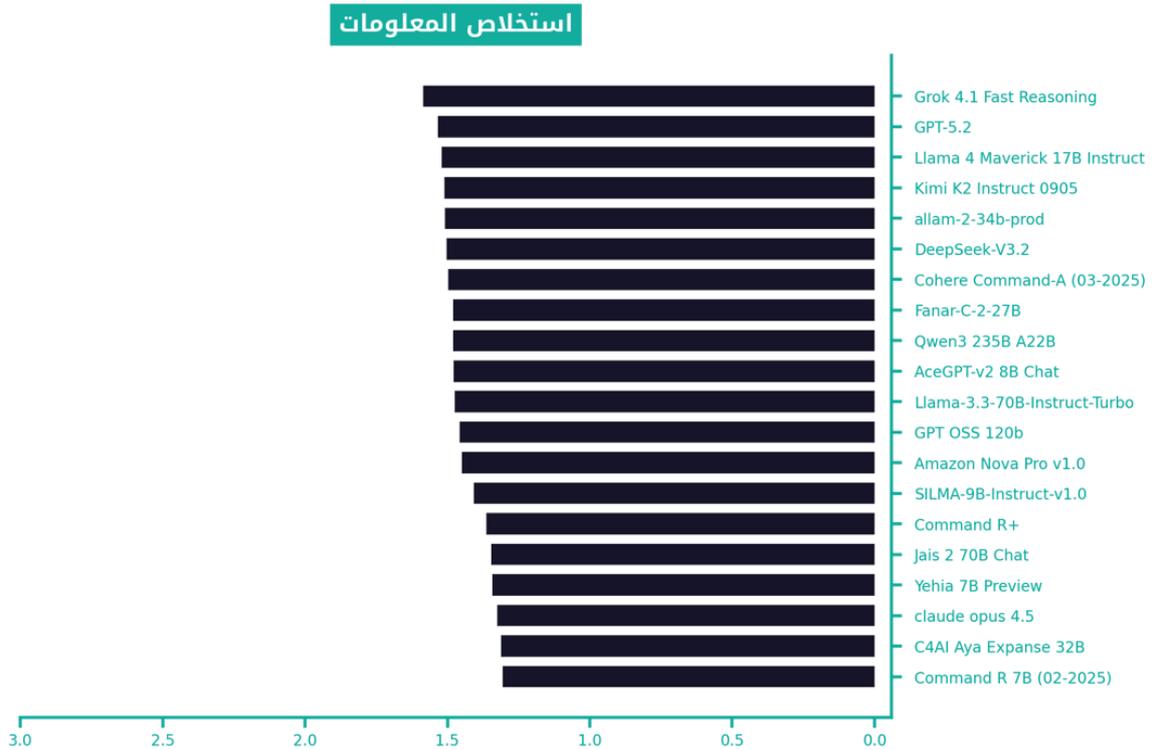


## ٤. ملء الفراغ (Fill in the Blank).



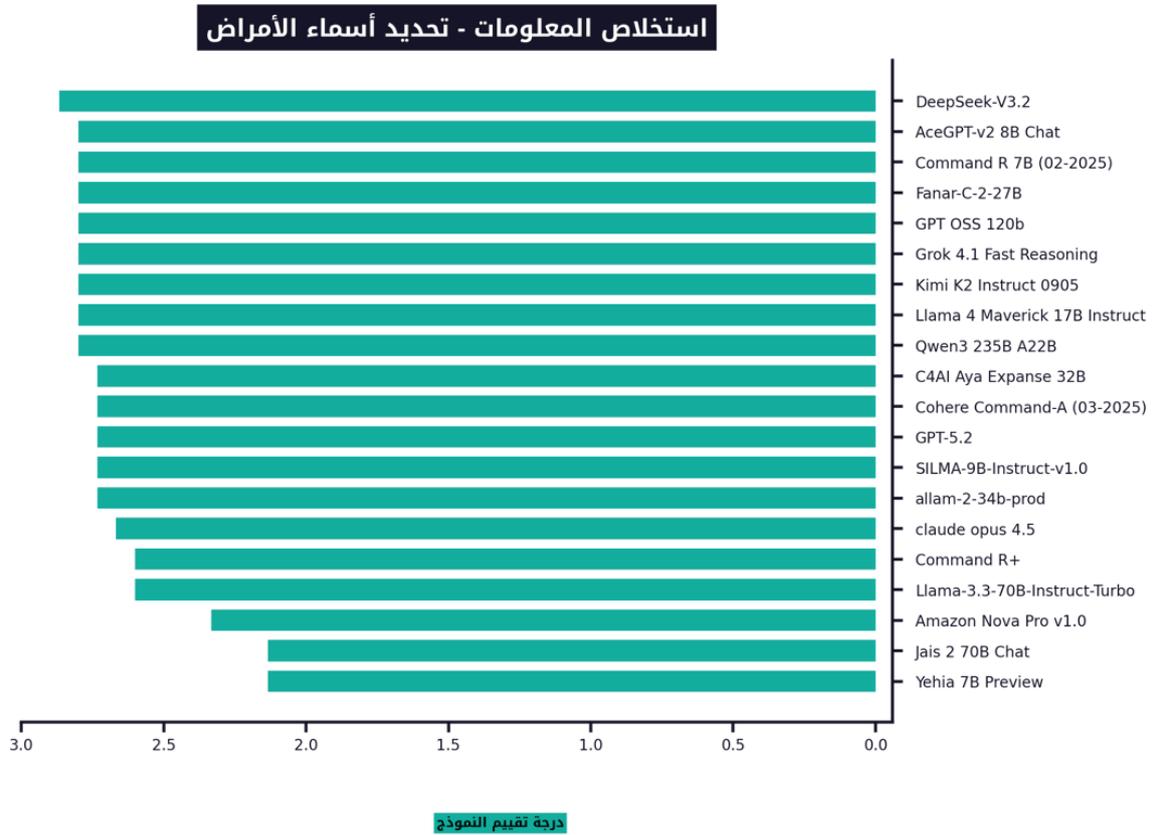
درجة تقييم النموذج

## ٥. استخلاص المعلومات (Information Extraction).

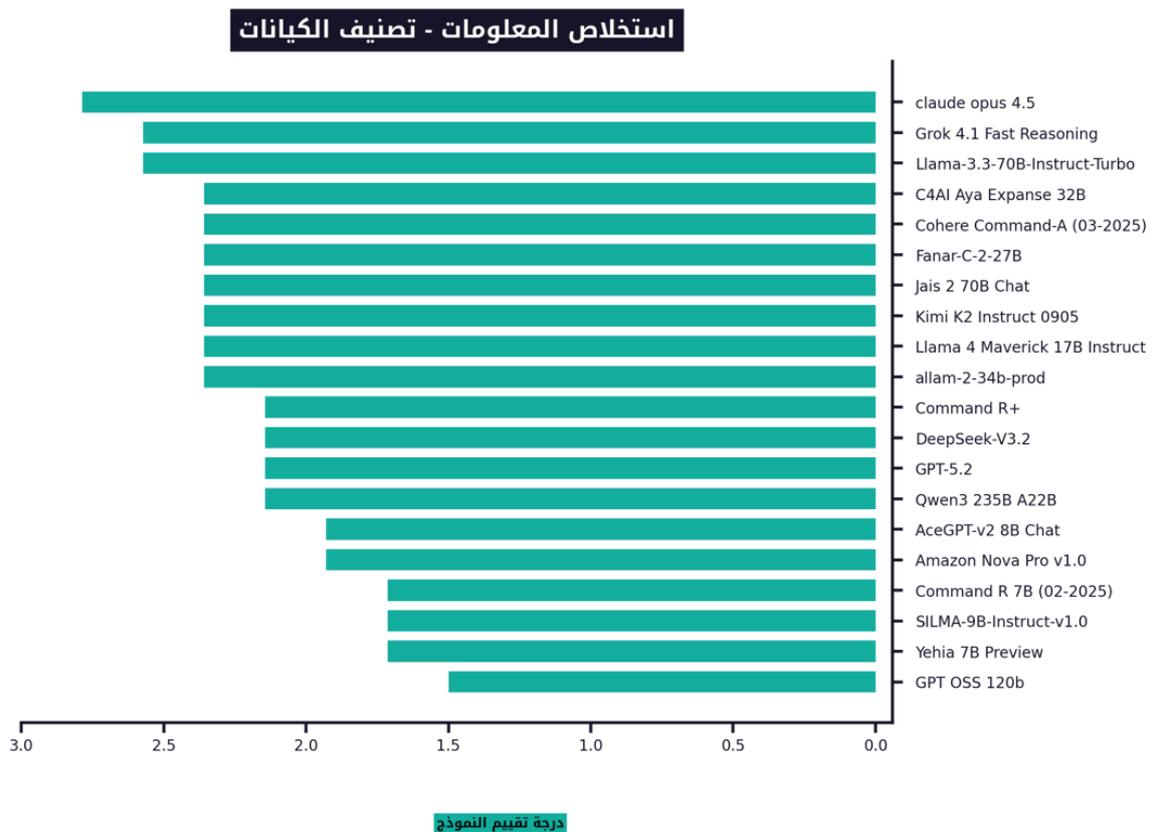


درجة تقييم النموذج

## 0.1 تحديد أسماء الأمراض (Disease Mention Identification).



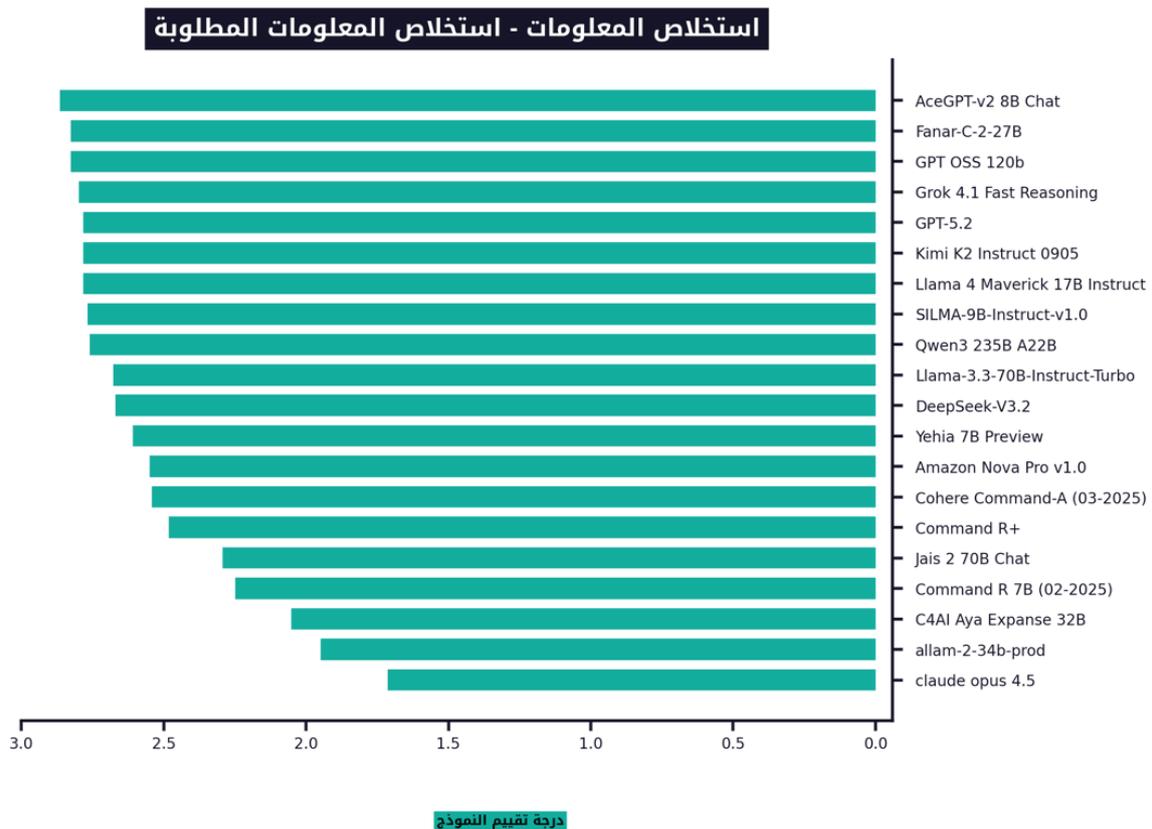
## 0.2 تصنيف الكيانات (Entity Categorization).



### ٥.٣ تصنيف العلاقات بين الكيانات (Entity Relation Classification).

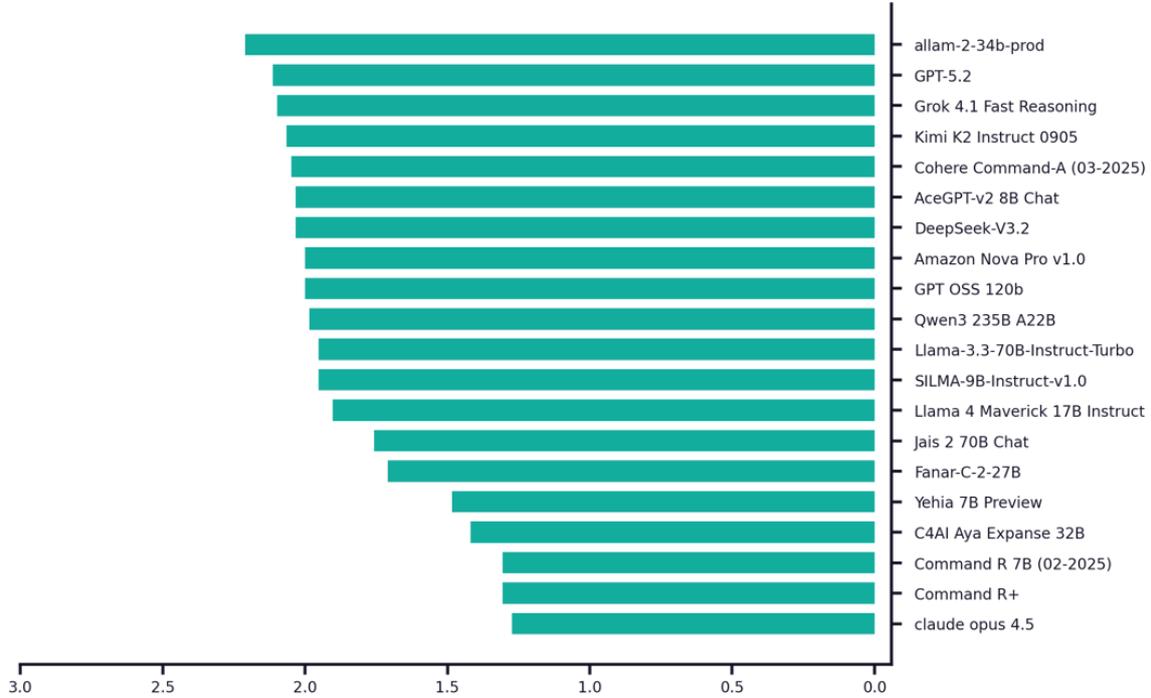


### ٥.٤ استخلاص المعلومات المطلوبة (Extracting Required Information).



## 0.0 استخلاص الكلمات الرئيسية (Keyword Extraction).

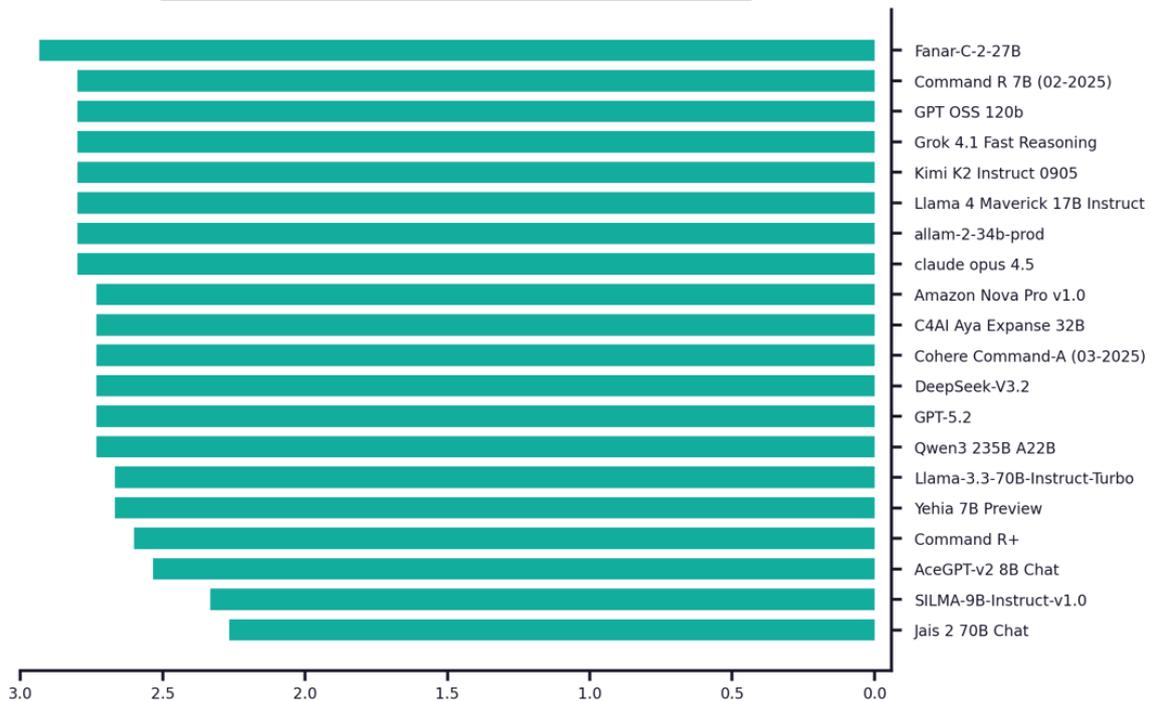
### استخلاص المعلومات - استخلاص الكلمات الرئيسية



درجة تقييم النموذج

## 0.6 التعرف على أسماء الكيانات (Named Entity Recognition).

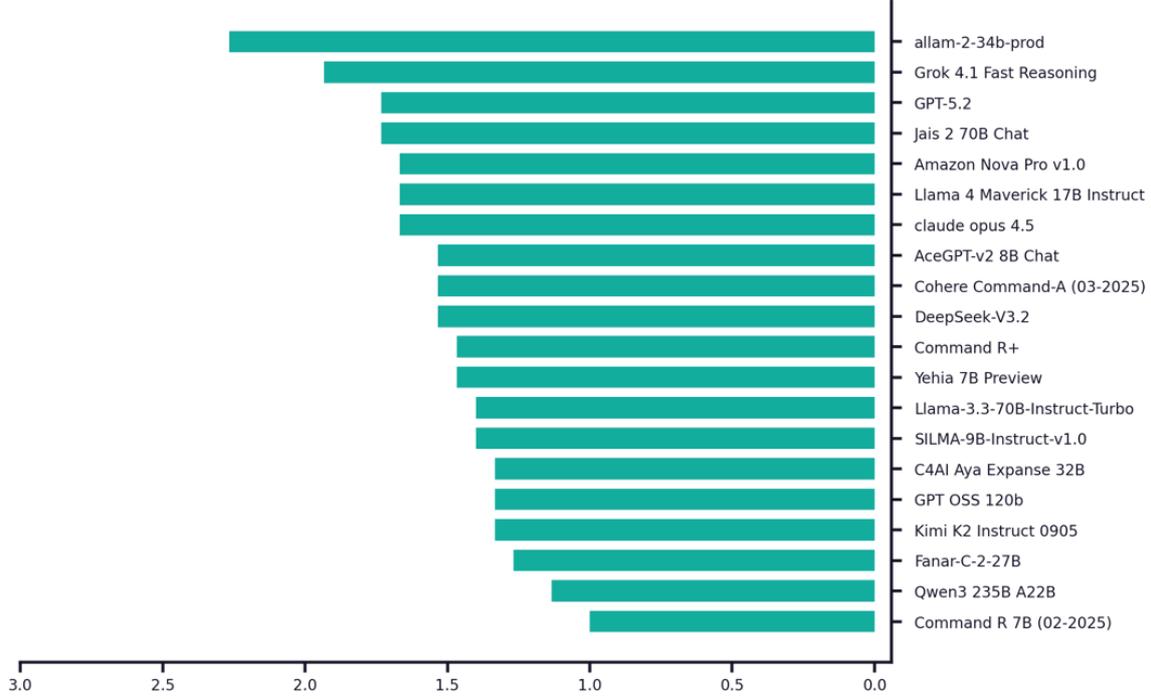
### استخلاص المعلومات - التعرف على أسماء الكيانات



درجة تقييم النموذج

## ٥.٧ استخلاص العلاقات (Relation Extraction).

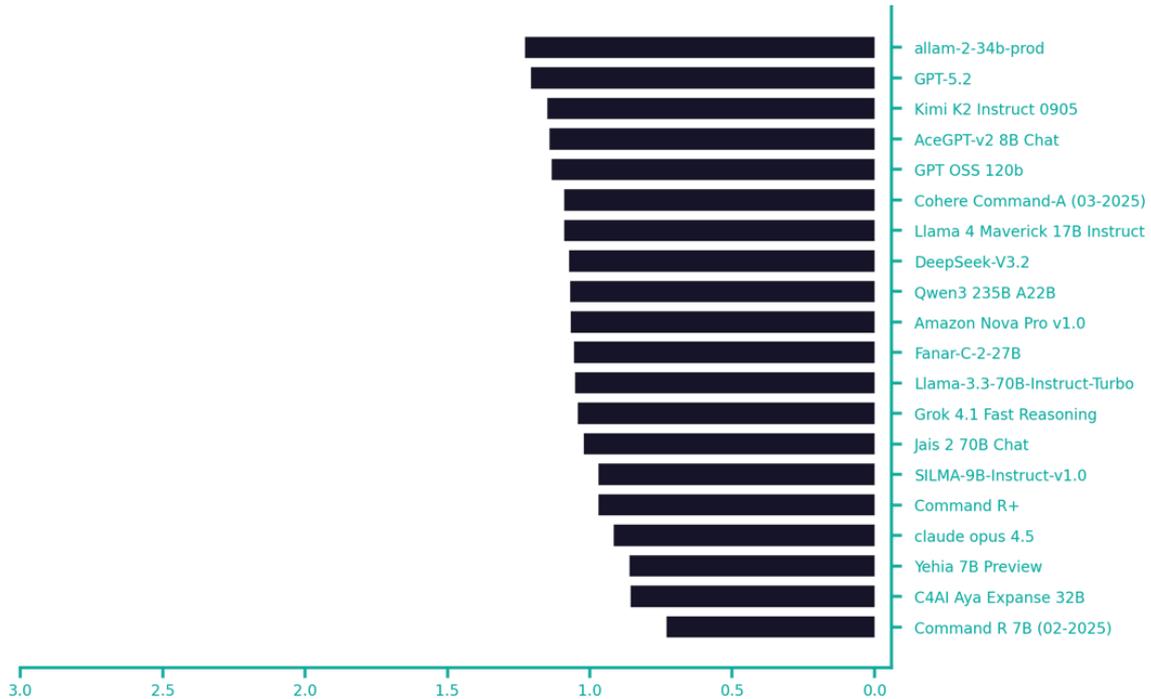
### استخلاص المعلومات - استخلاص العلاقات



درجة تقييم النموذج

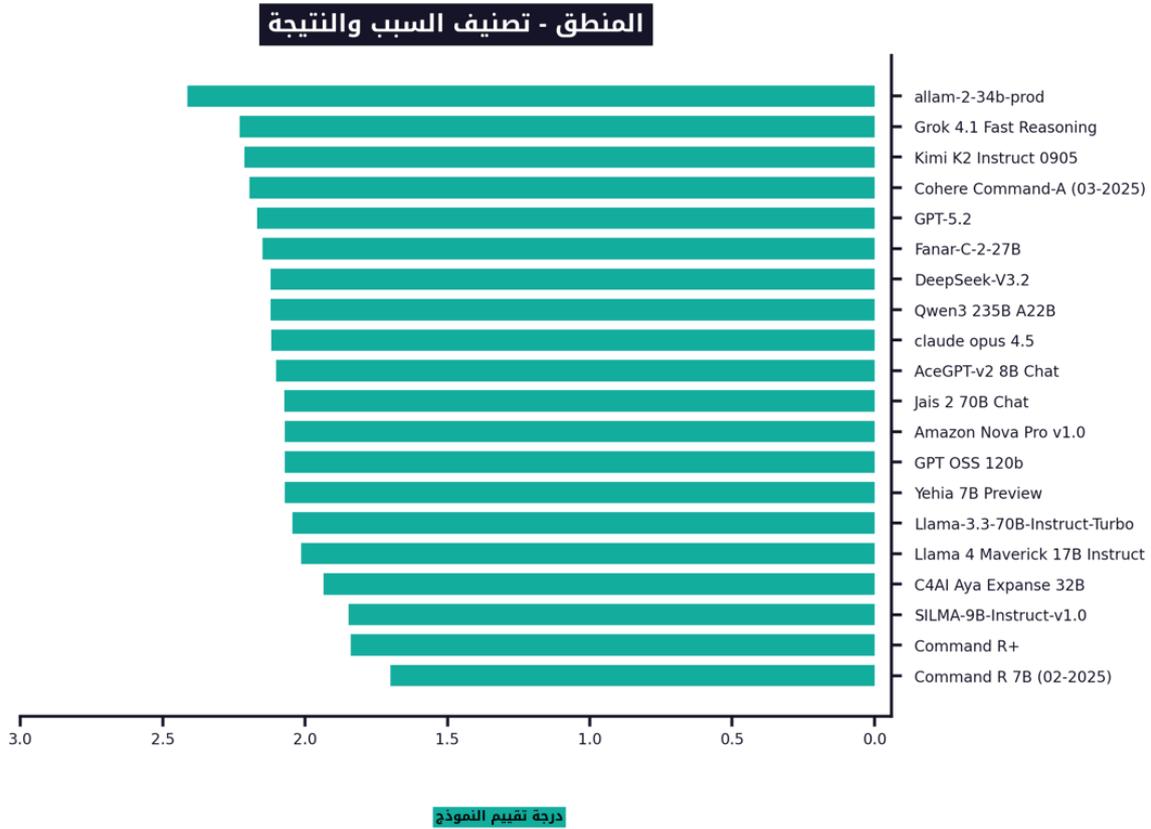
## ٦. المنطق (Logic).

### المنطق

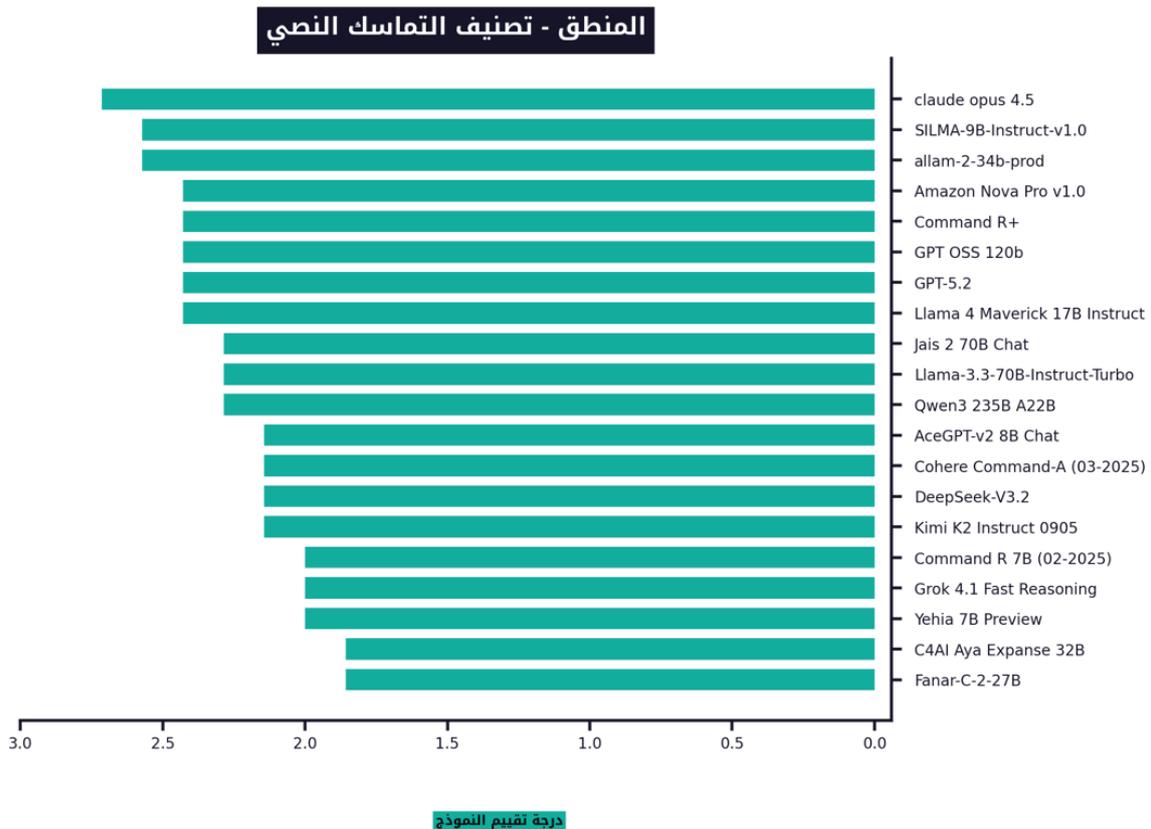


درجة تقييم النموذج

## ٦.١ تصنيف السبب والنتيجة (Cause Effect Classification).

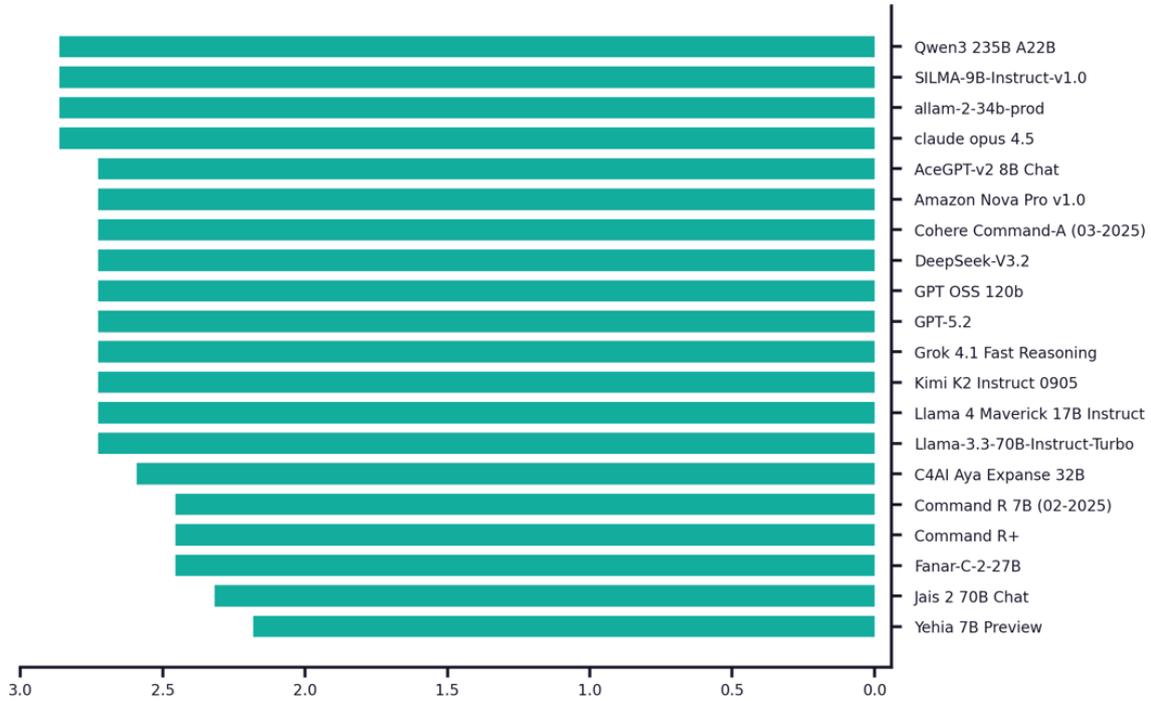


## ٦.٢ تصنيف التماسك النصي (Coherence Classification).



### ٦.٣ التحقق من المنطق السليم (Commonsense Validation).

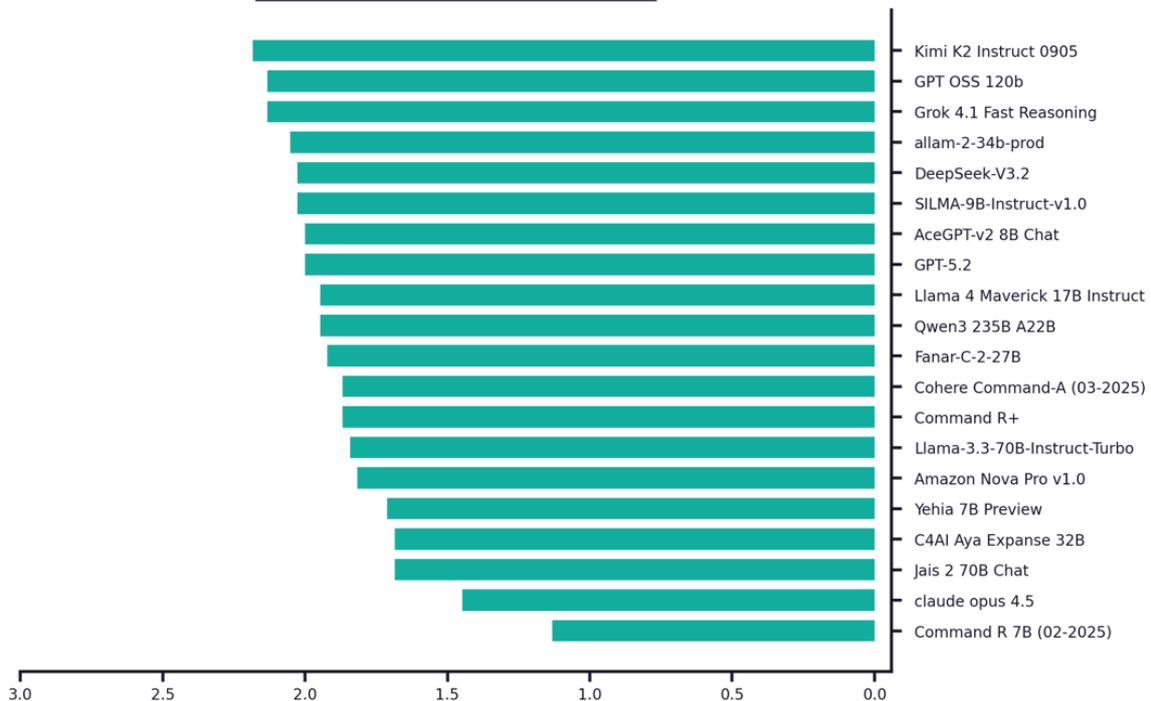
#### المنطق - التحقق من المنطق السليم



درجة تقييم النموذج

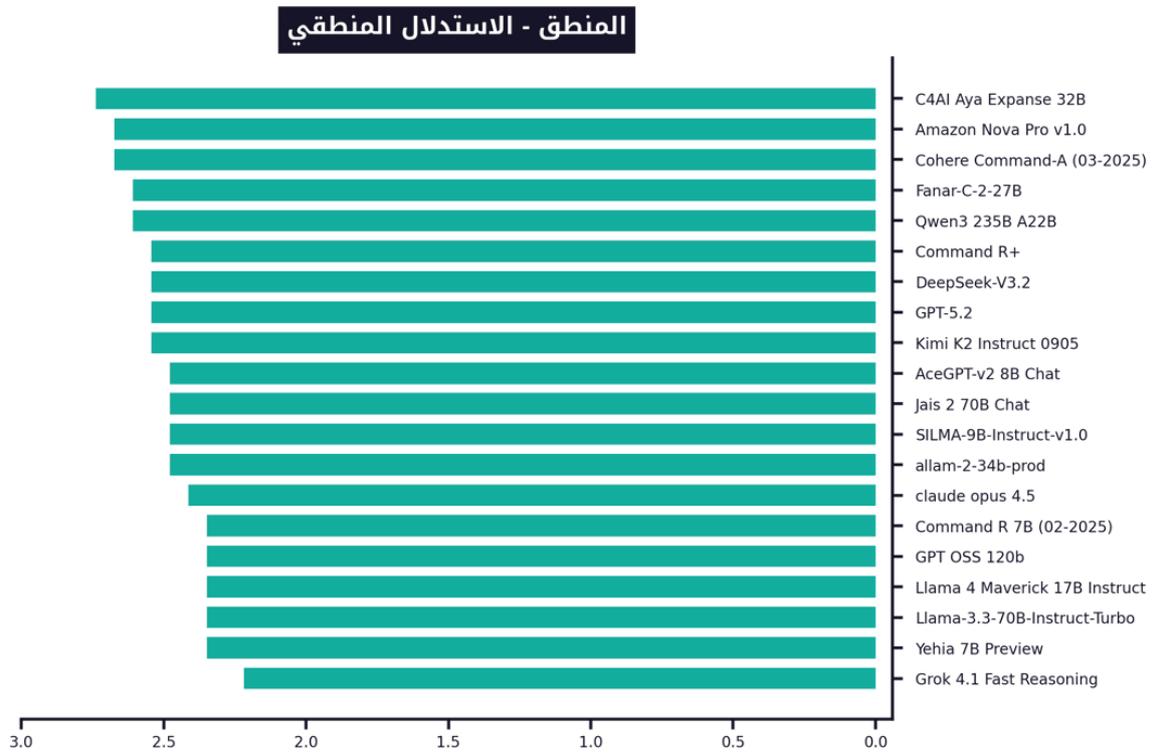
### ٦.٤ تحليل الإحالات الضميرية (Coreference Resolution).

#### المنطق - تحليل الإحالات الضميرية



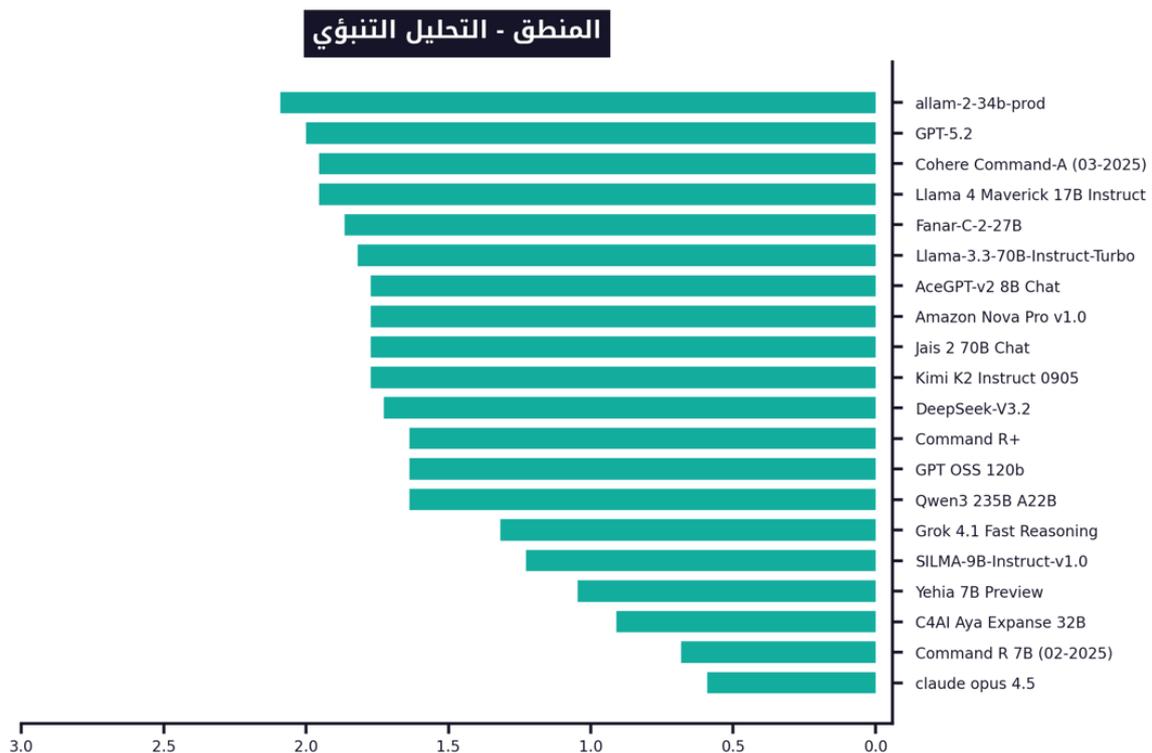
درجة تقييم النموذج

## ٦.٠ الاستدلال المنطقي (Logical Reasoning).



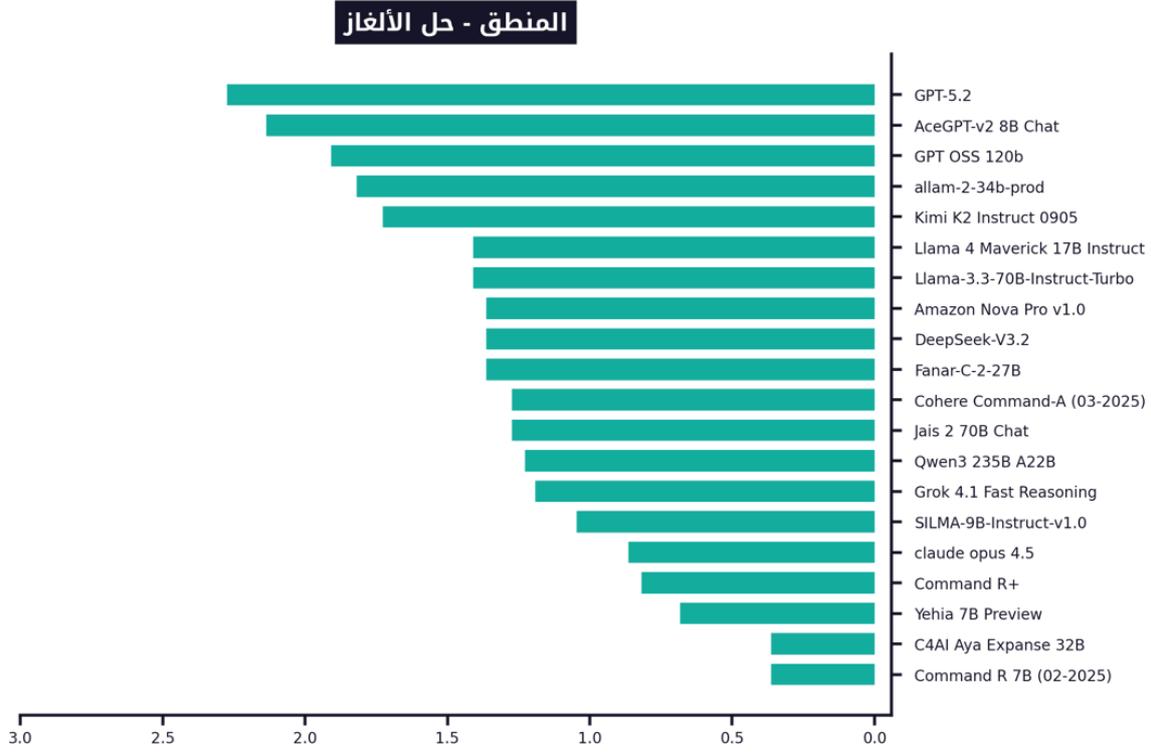
درجة تقييم النموذج

## ٦.١ التحليل التنبؤي (Predictive Analysis).



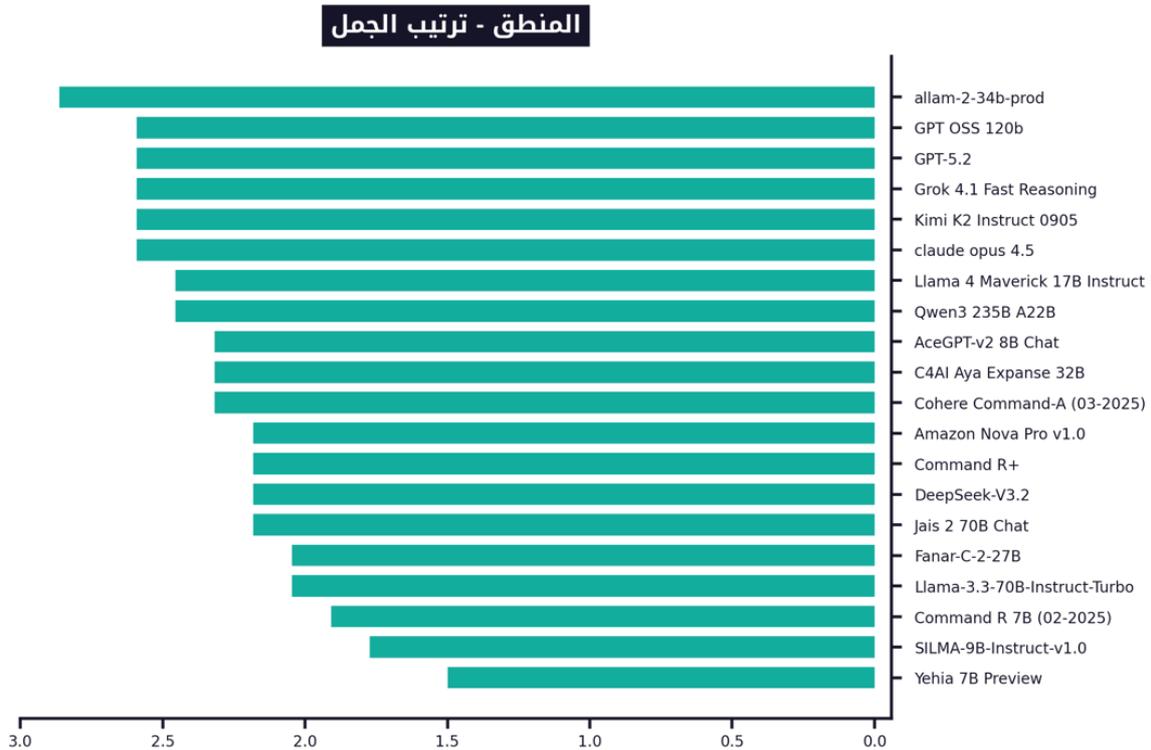
درجة تقييم النموذج

## ٦.٧ حل الألغاز (Riddle Solving).



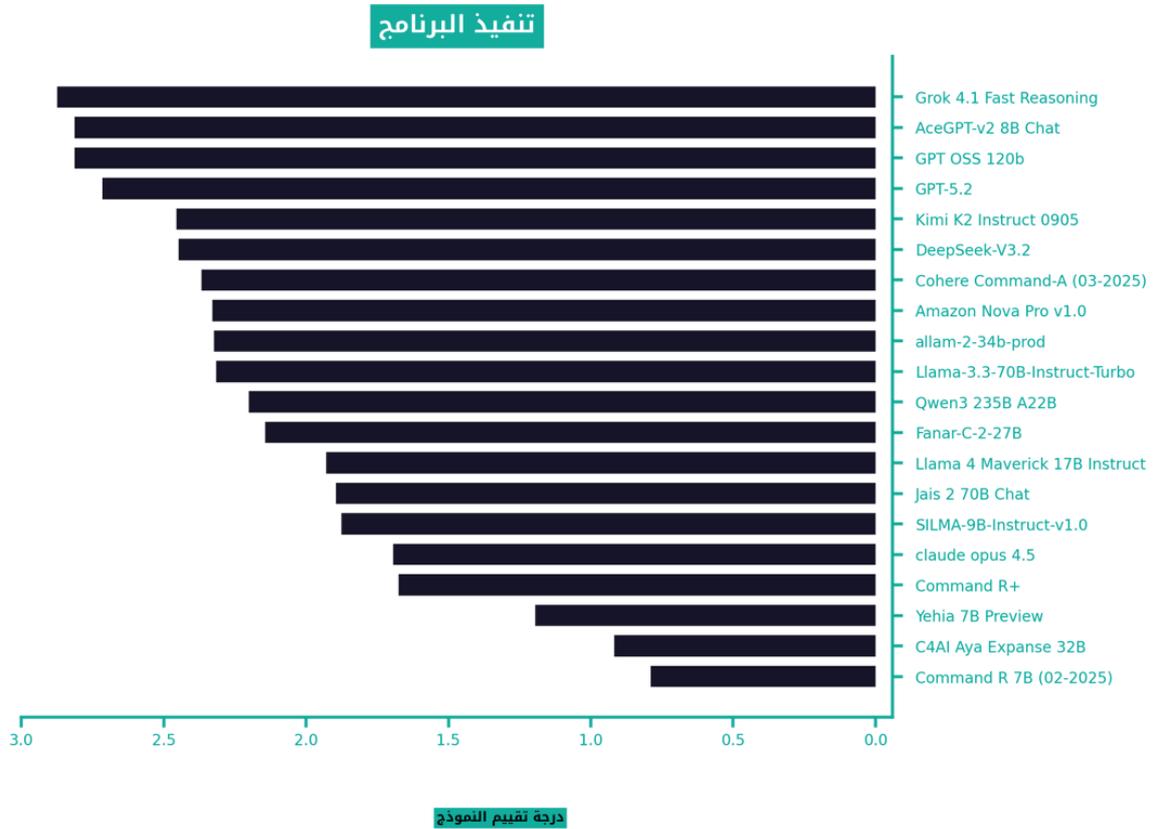
درجة تقييم النموذج

## ٦.٨ ترتيب الجمل (Sentence Ordering).

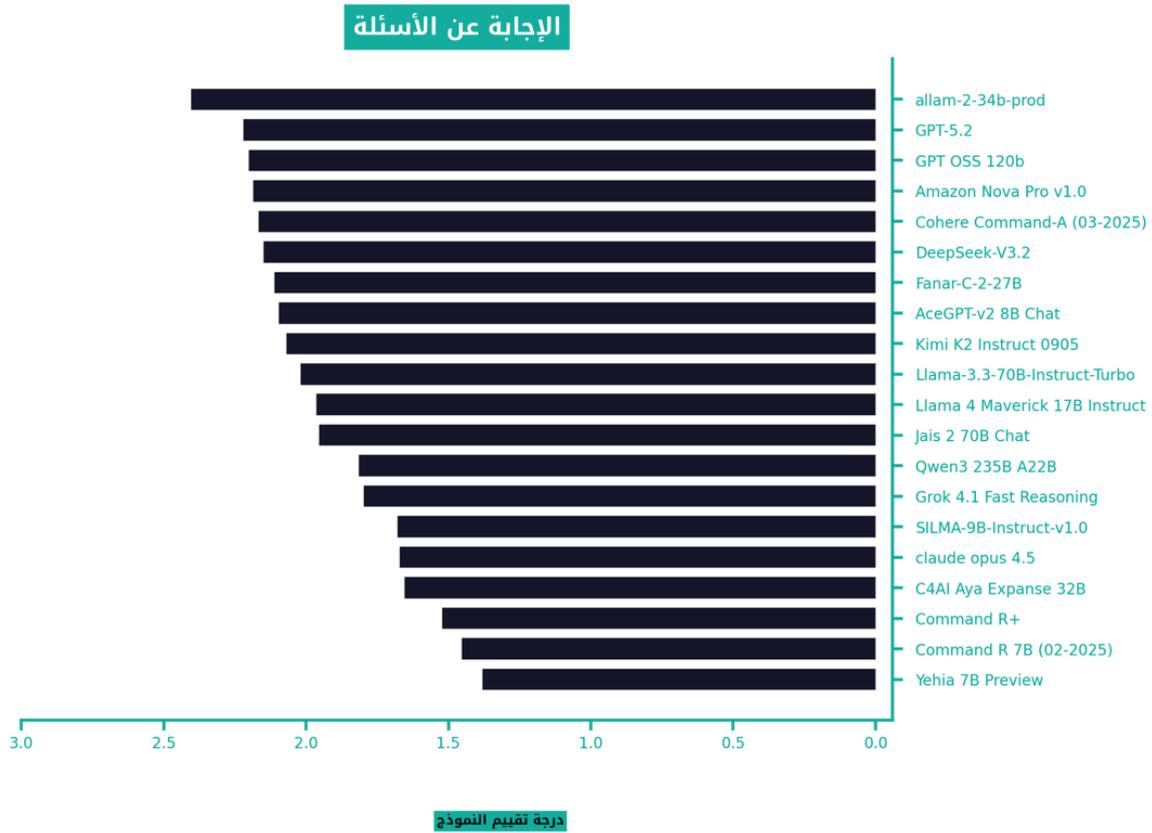


درجة تقييم النموذج

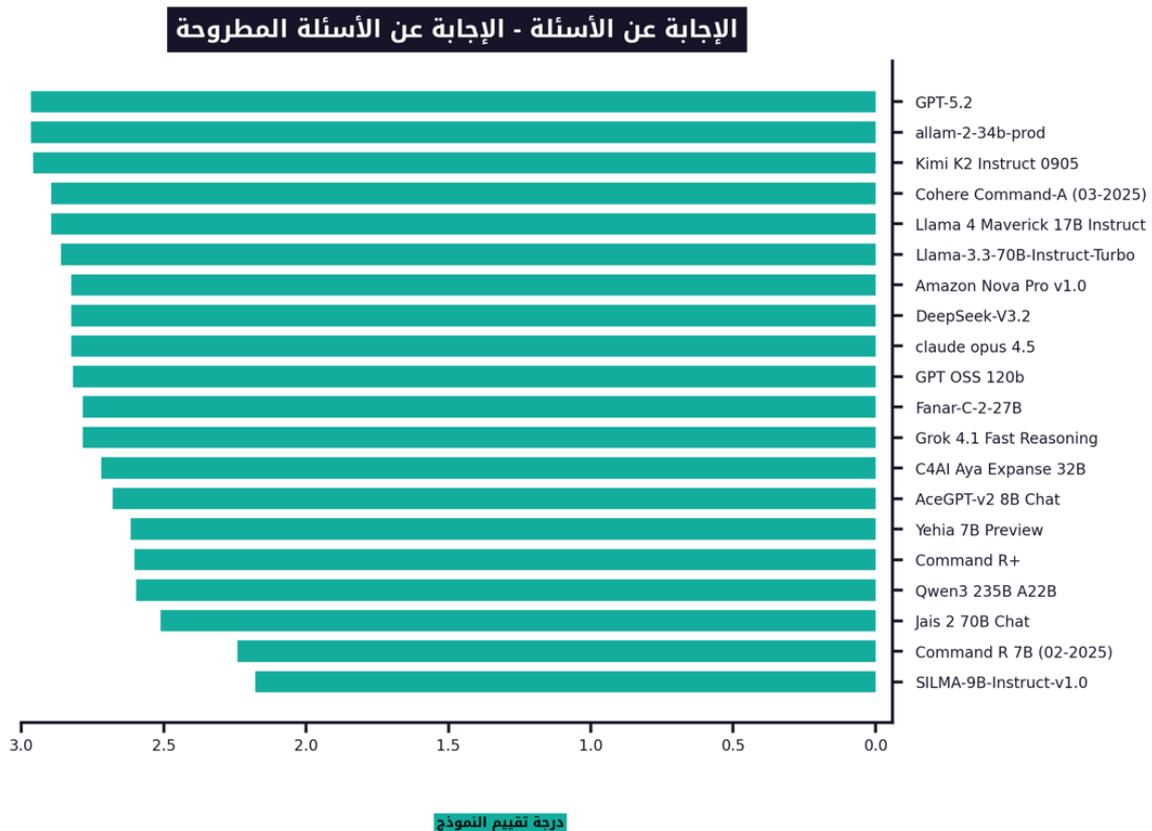
## ٧. تنفيذ البرامج (Program Execution).



## ٨. الإجابة عن الأسئلة (Question Answering).

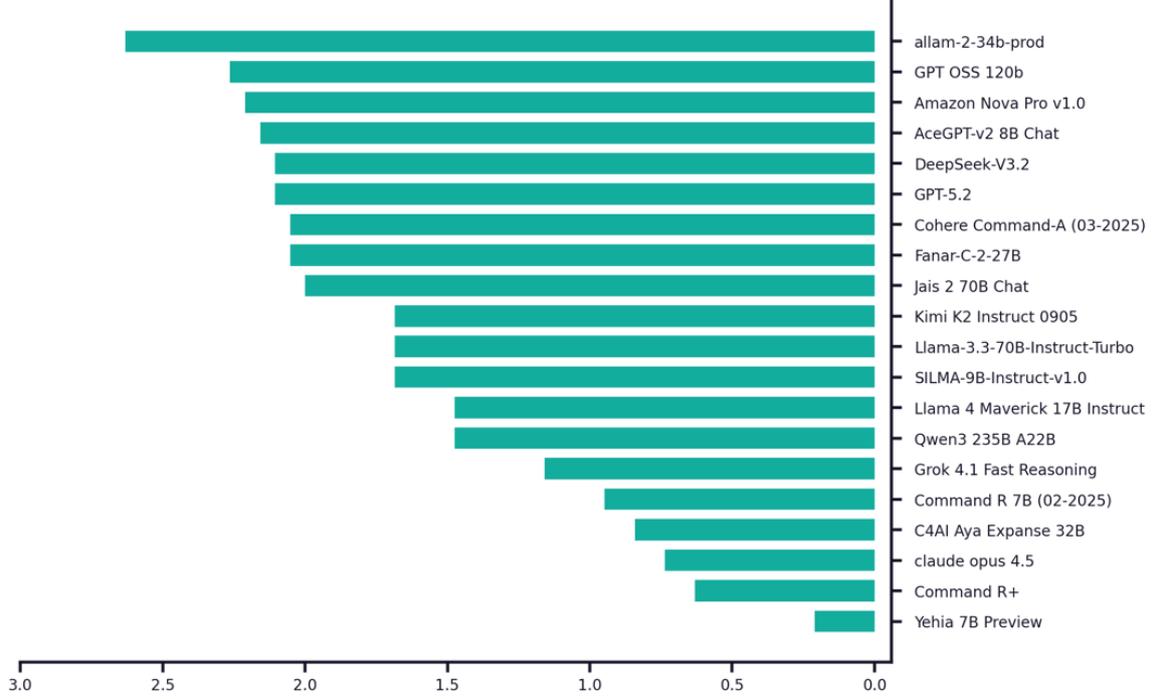


## ٨.١ الإجابة عن الأسئلة المطروحة (Answering Given Question).



## ٨.٣ تحليل السؤال (Question Decomposition).

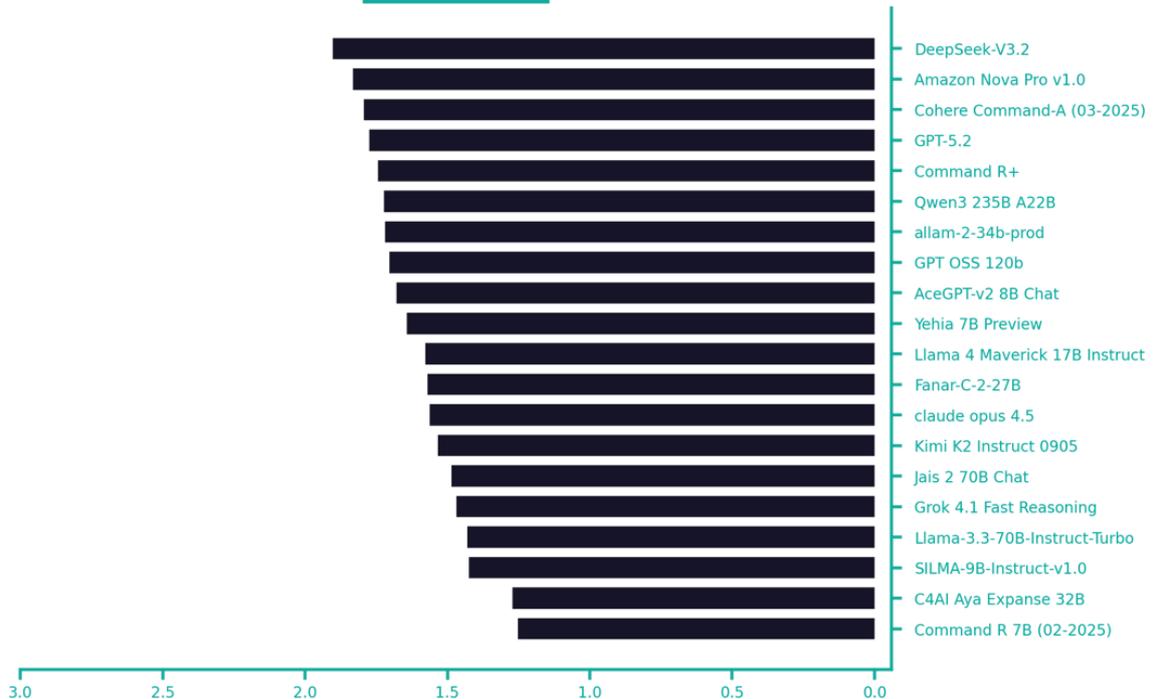
### الإجابة عن الأسئلة - تحليل السؤال



درجة تقييم النموذج

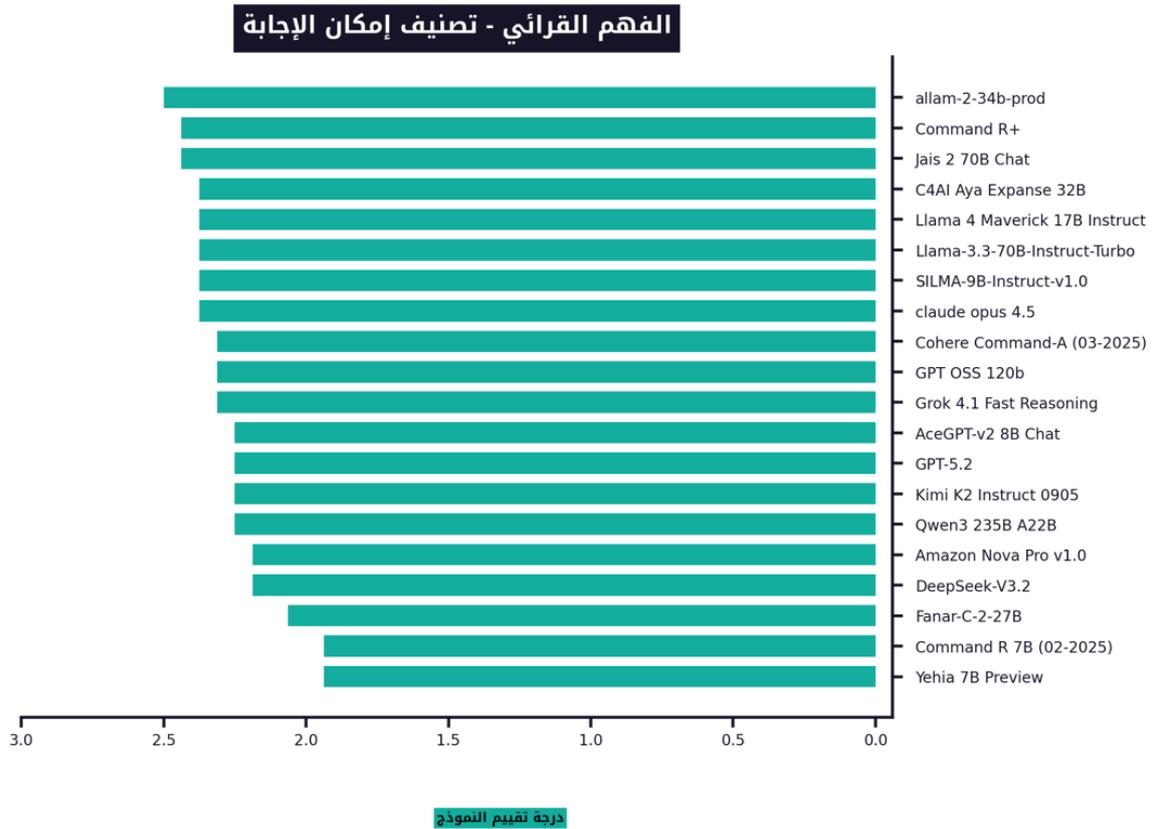
## ٩. الفهم القرائي (Reading Comprehension).

### الفهم القرائي

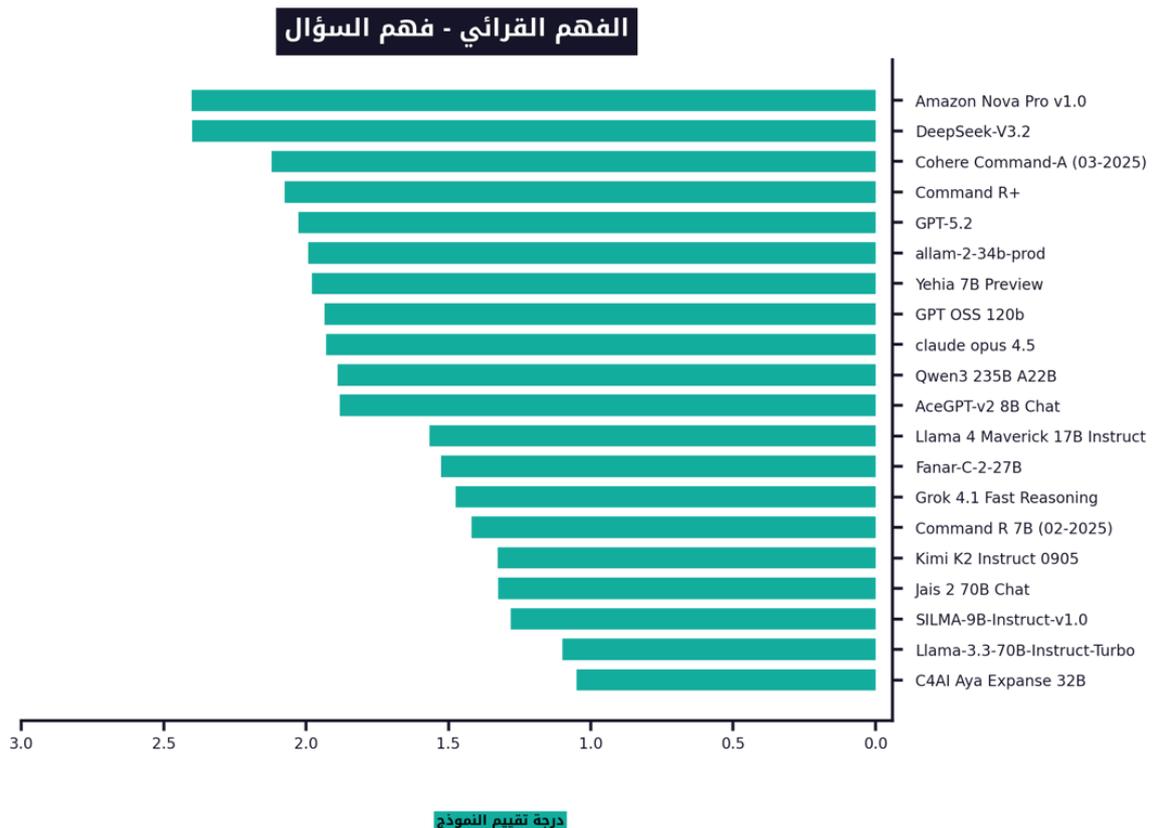


درجة تقييم النموذج

## ٩.١ تصنيف إمكان الإجابة (Answerability Classification).



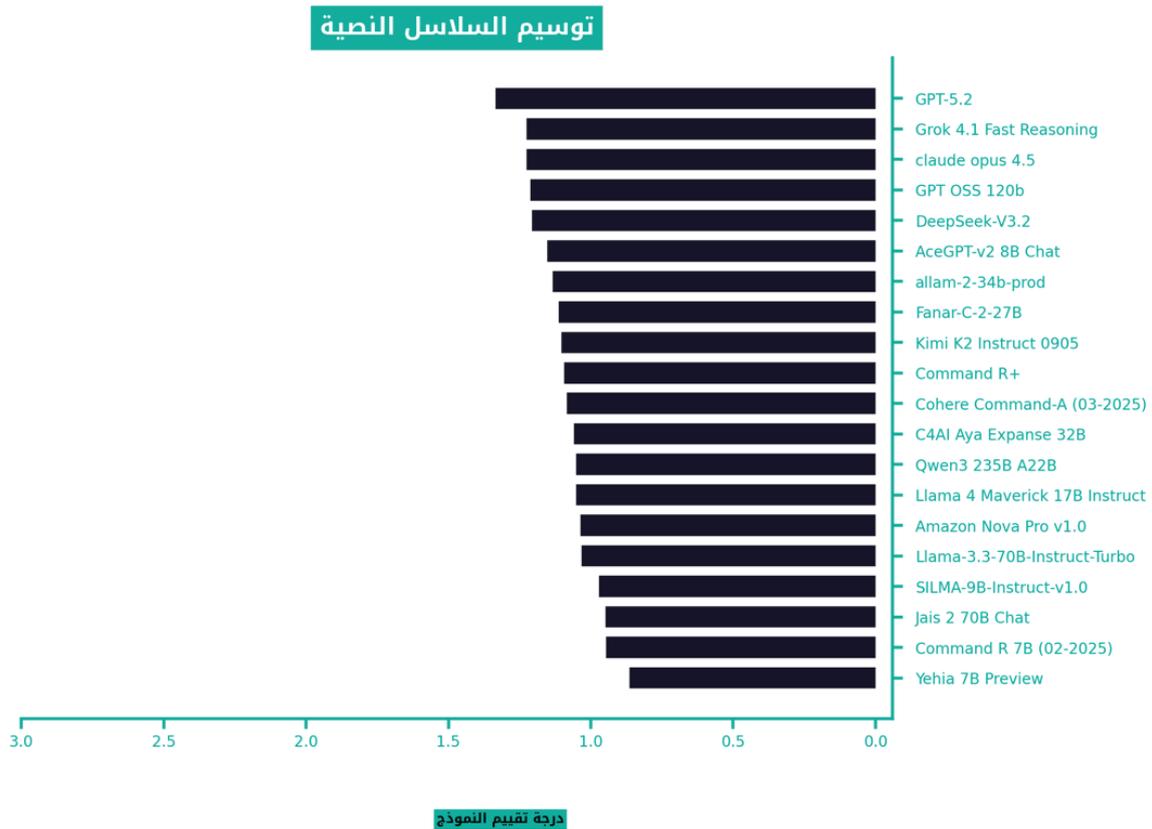
## ٩.٢ فهم السؤال (Question Understanding).



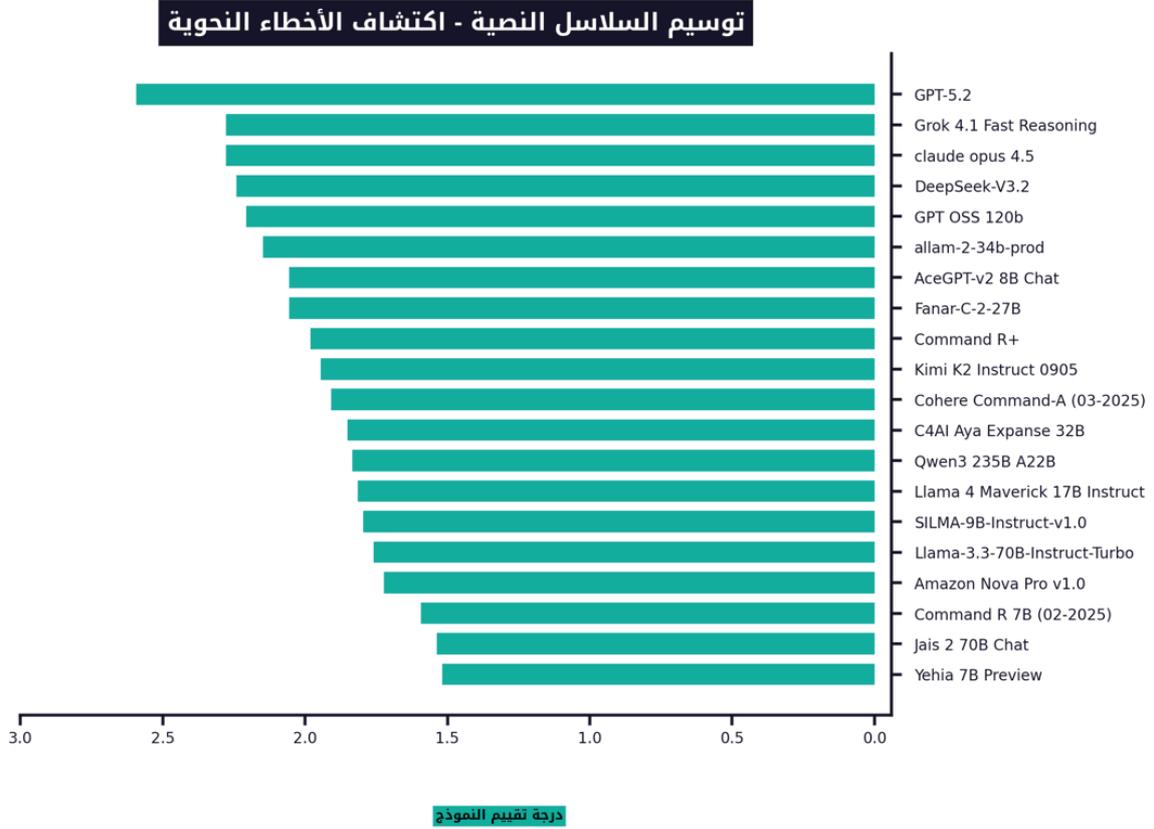
## ٩.٣ التحقق من الاستدلال النصي (Textual Inference Validation).



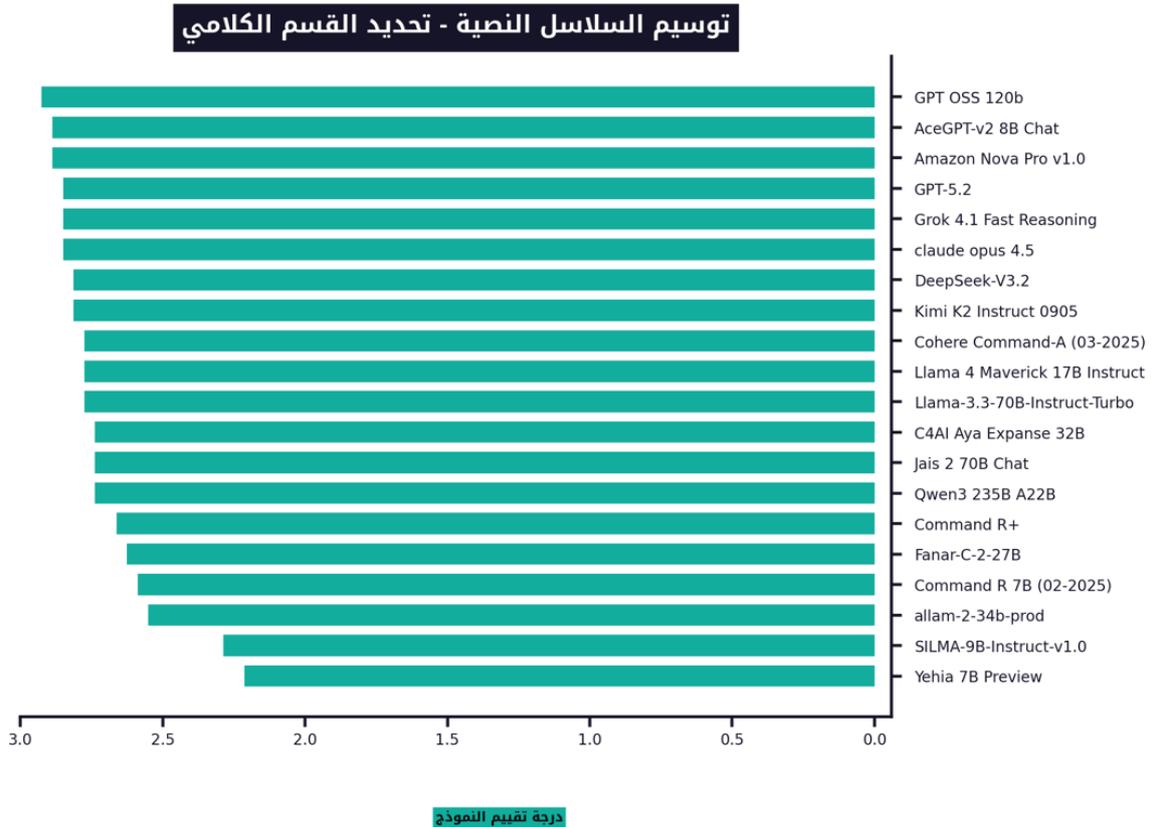
## ١٠. توسيم السلاسل النصية (Sequence Tagging).



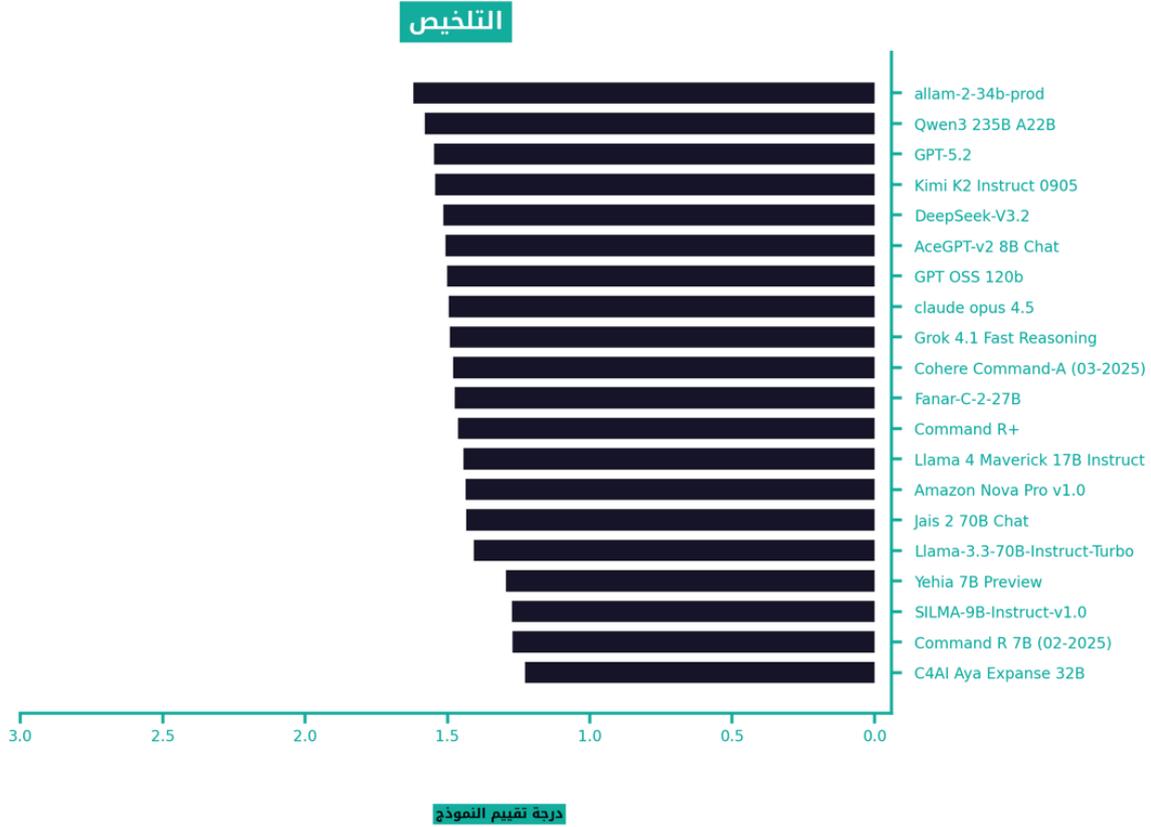
## ١٠.١ اكتشاف الأخطاء النحوية (Grammar Detection).



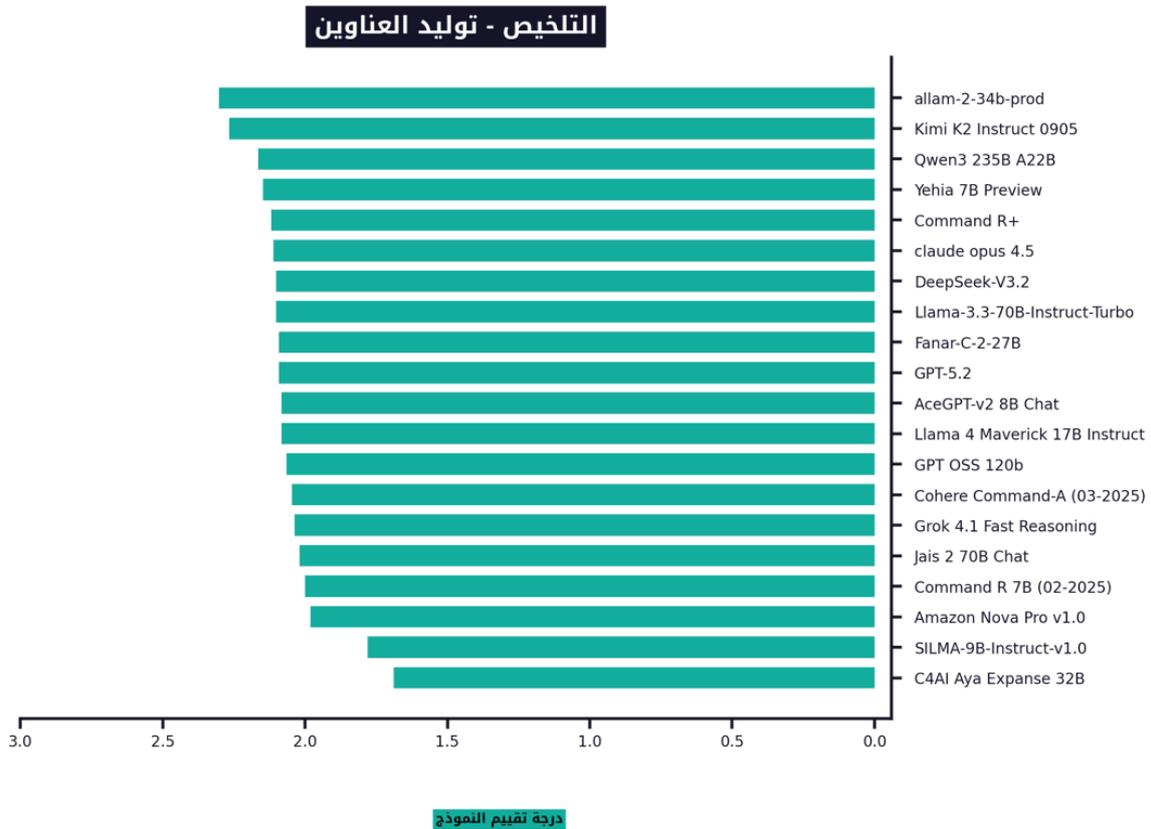
## ١٠.٢ تحديد القسم الكلامي (Part Of Speech Tagging).



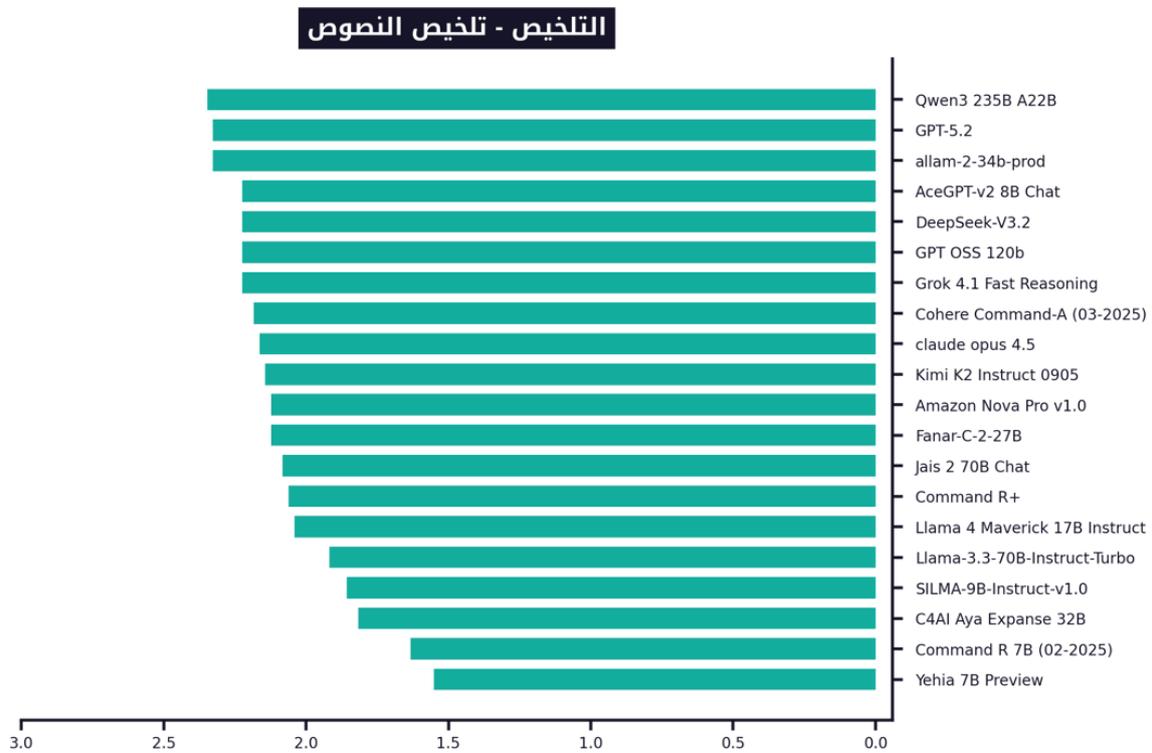
## ١١. التلخيص (Summarization).



## ١١.١ توليد العناوين (Subject Generation).

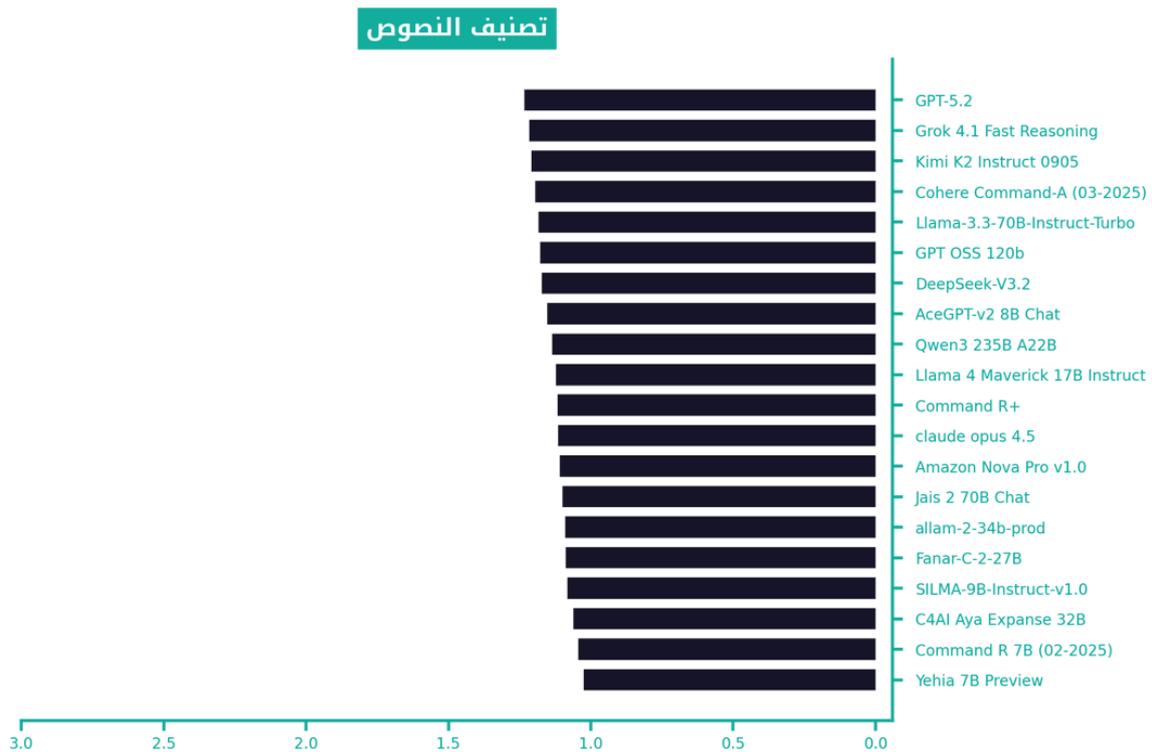


## ١١.٦ تلخيص النصوص (Text Summarization).



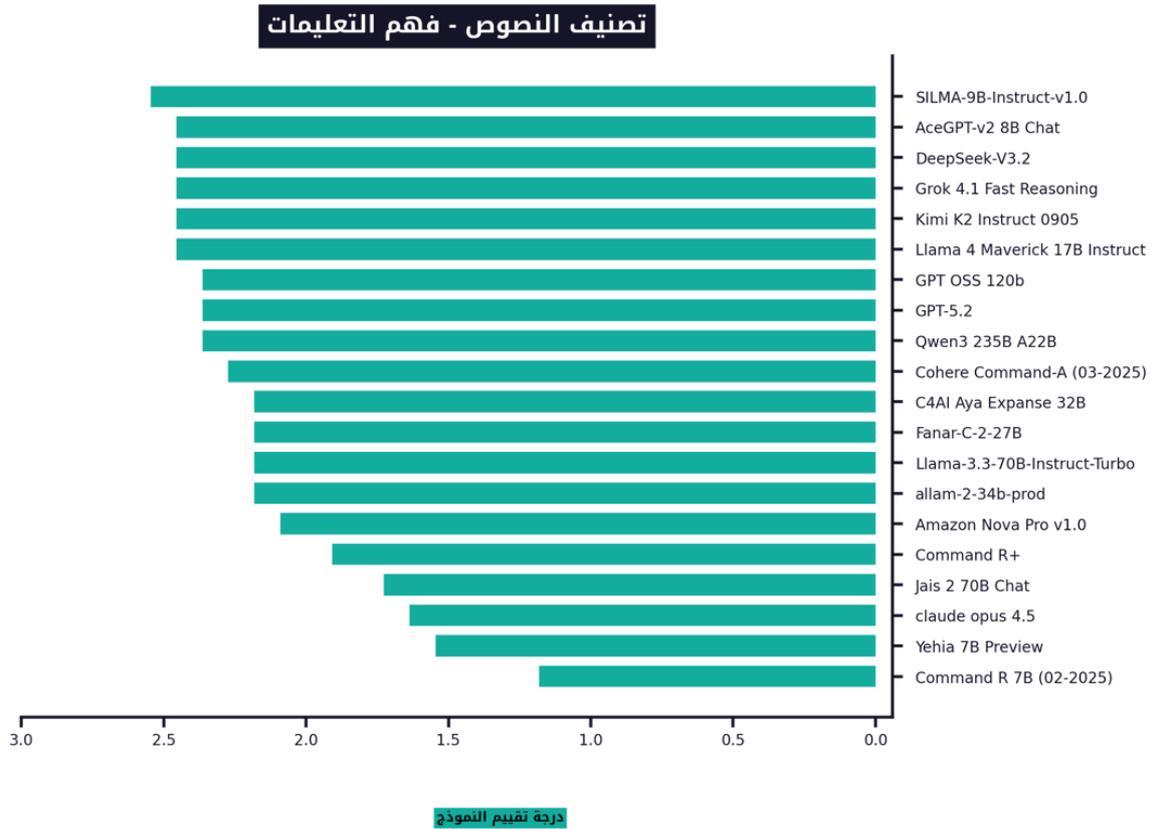
درجة تقييم النموذج

## ١١.٧ تصنيف النصوص (Text Classification).

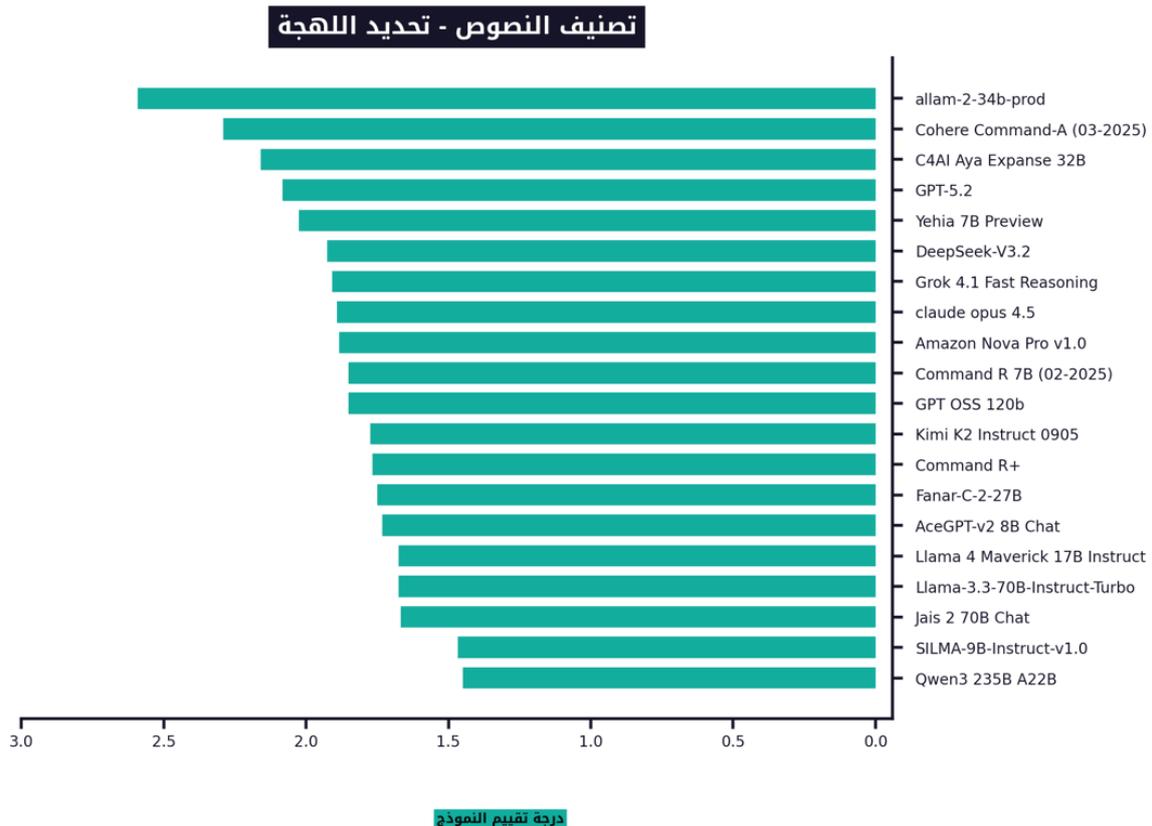


درجة تقييم النموذج

## ١٢.١ فهم التعليمات (Command Interpretation).



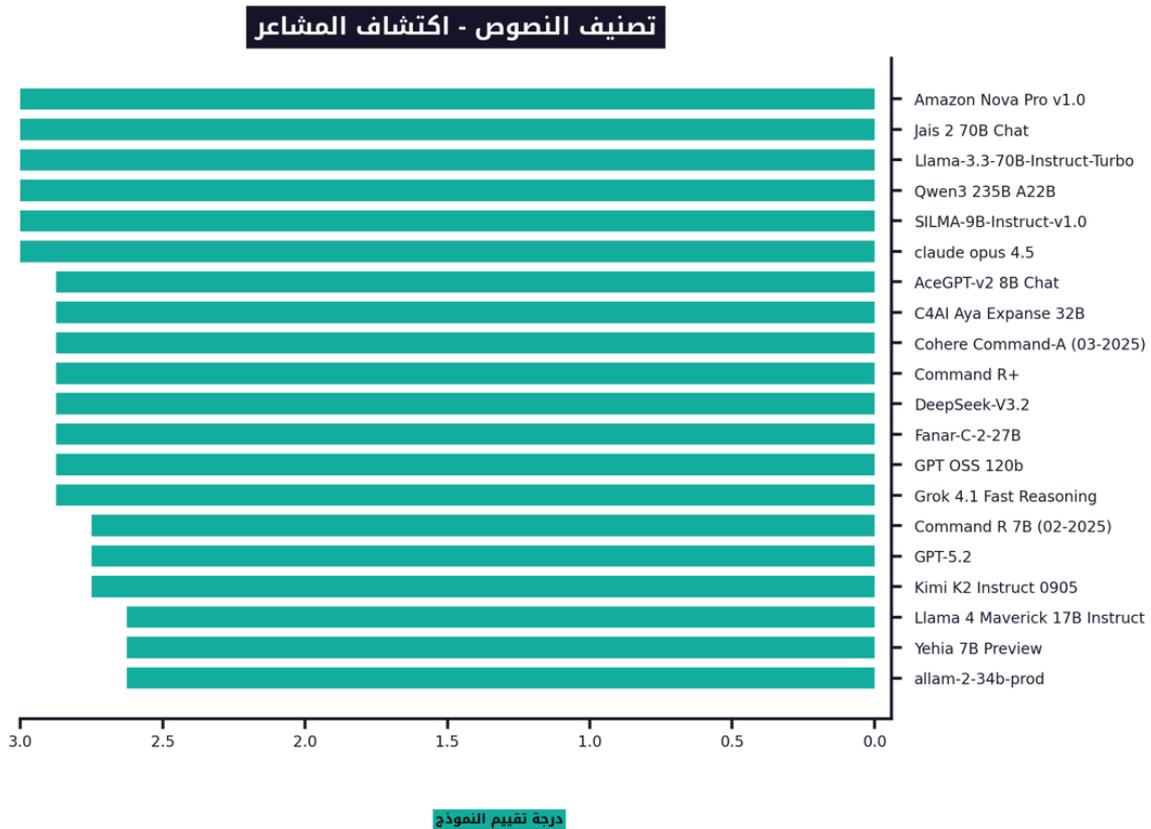
## ١٢.٢ تحديد اللهجة (Dialect Identification).



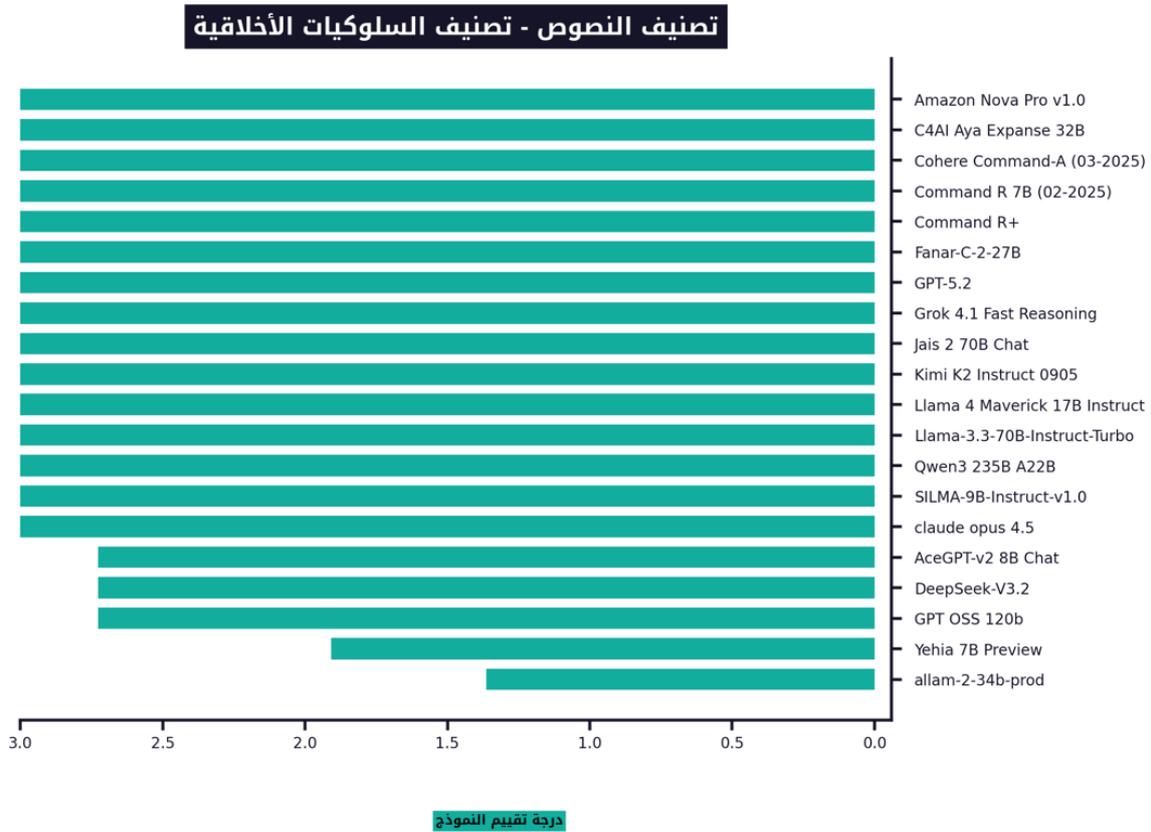
## ١٢.٣ تعرّف أفعال الحوار (Dialogue Act Recognition).



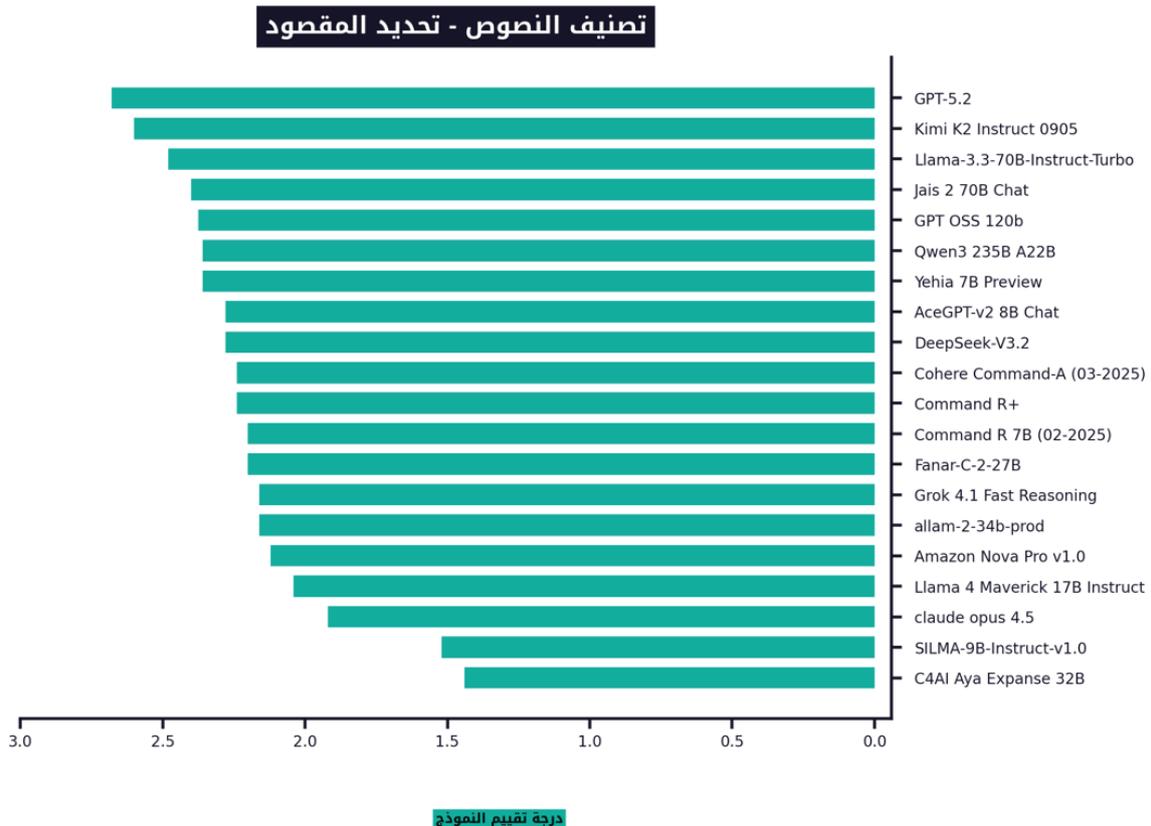
## ١٢.٤ اكتشاف المشاعر (Emotion Detection).



## ١٢.٥ تصنيف السلوكيات الأخلاقية (Ethics Classification).



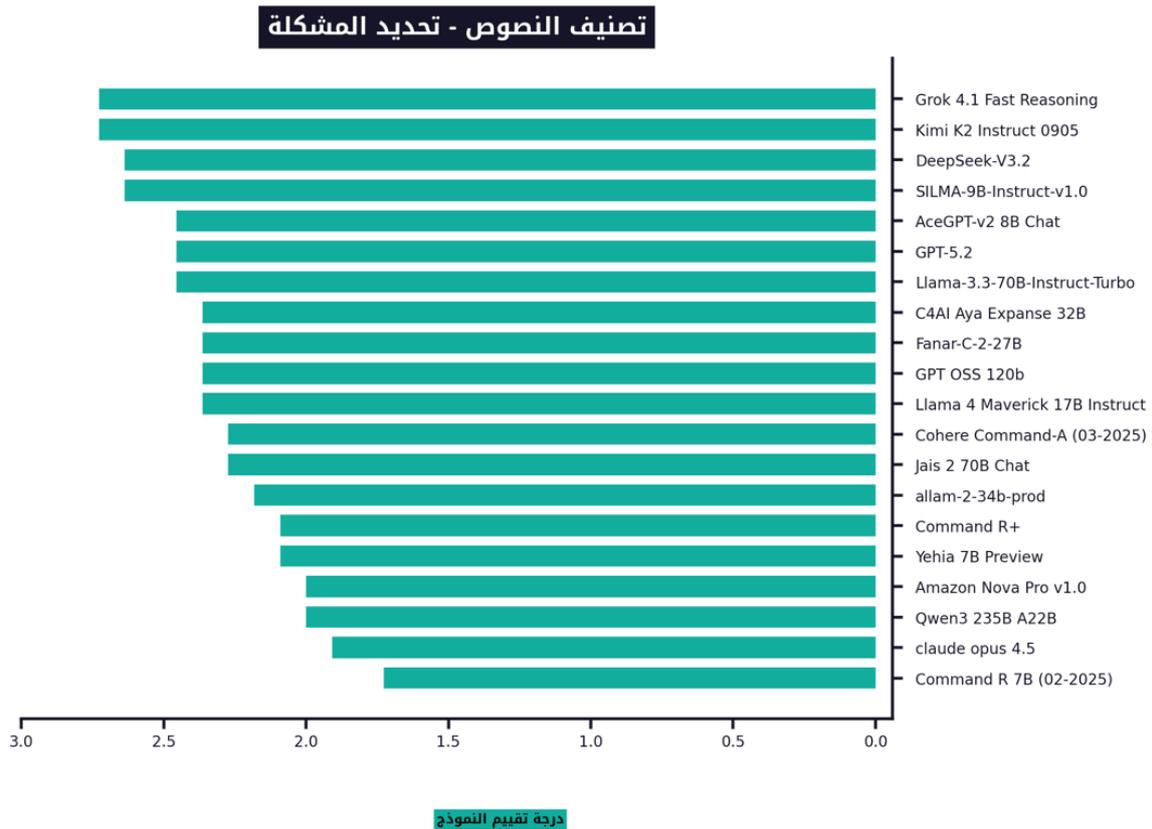
## ١٢.٦ تحديد المقصود (Intent Classification).



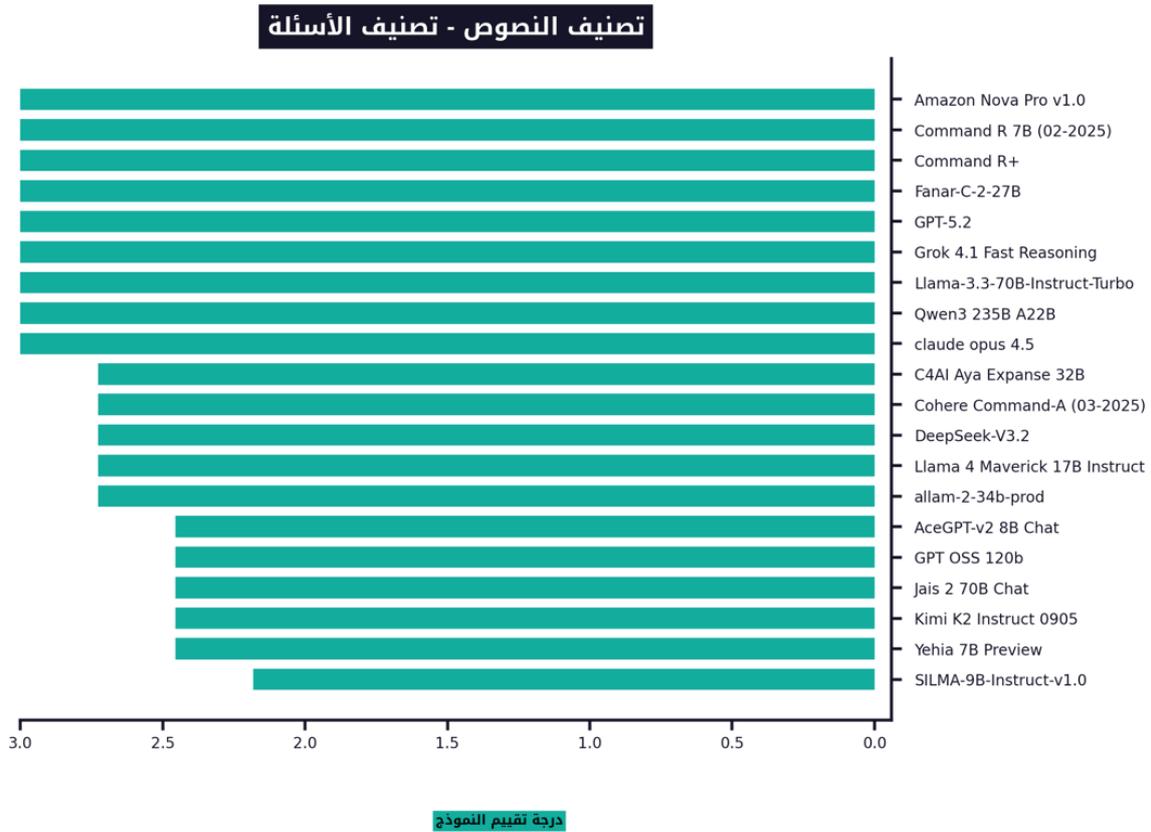
## ١٢.٧ اكتشاف الإساءة (Offensive Language Detection).



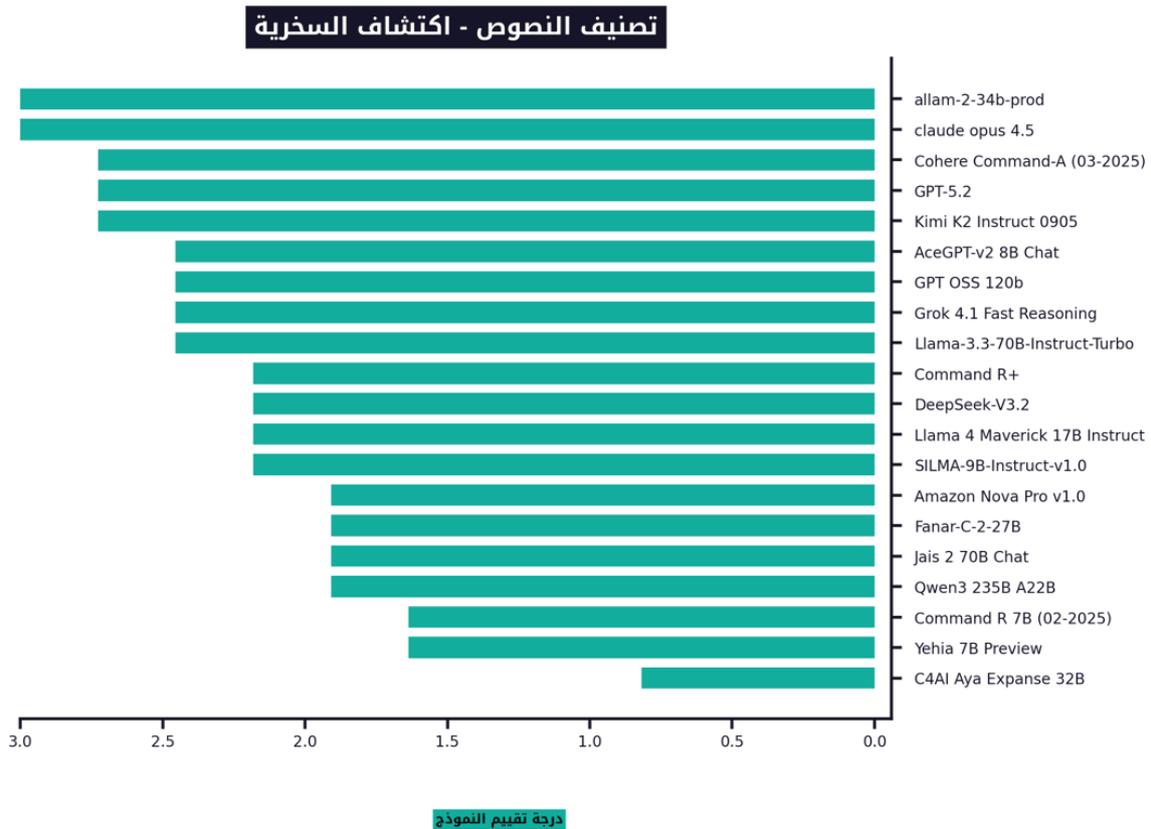
## ١٢.٨ تحديد المشكلة (Problem Identification).



## ١٢.٩ تصنيف الأسئلة (Question Categorization).

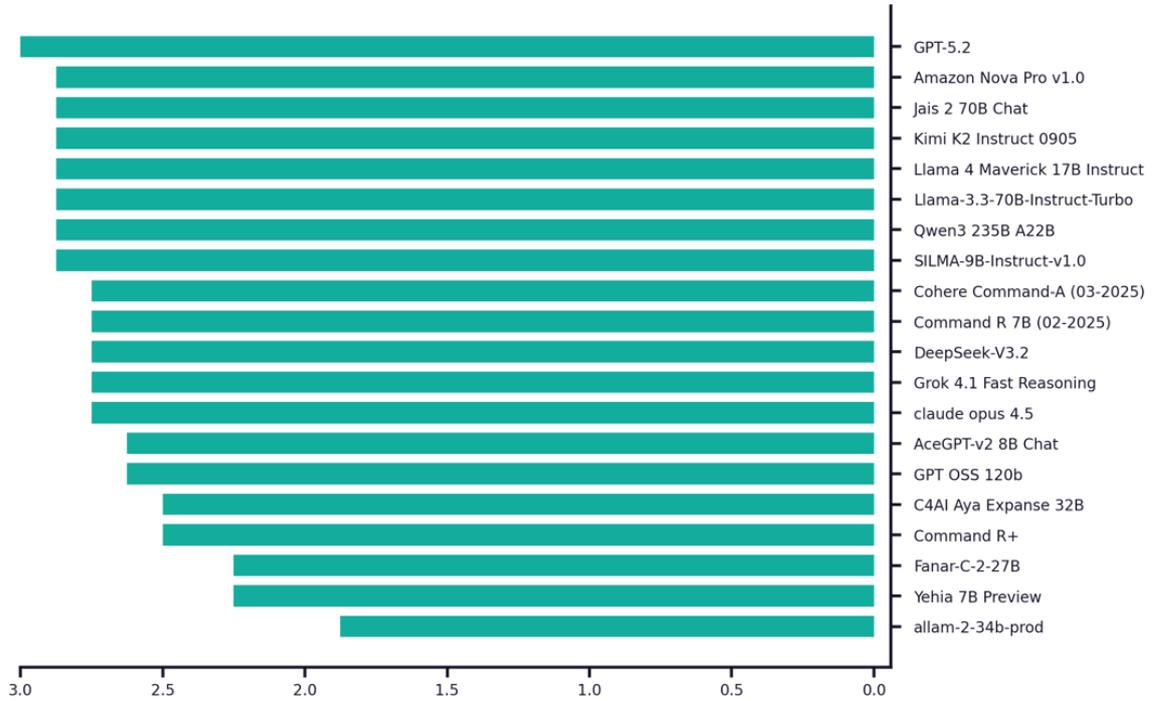


## ١٢.١٠ اكتشاف السخرية (Sarcasm Detection).



## ١٢.١١ تحليل المشاعر (Sentiment Analysis).

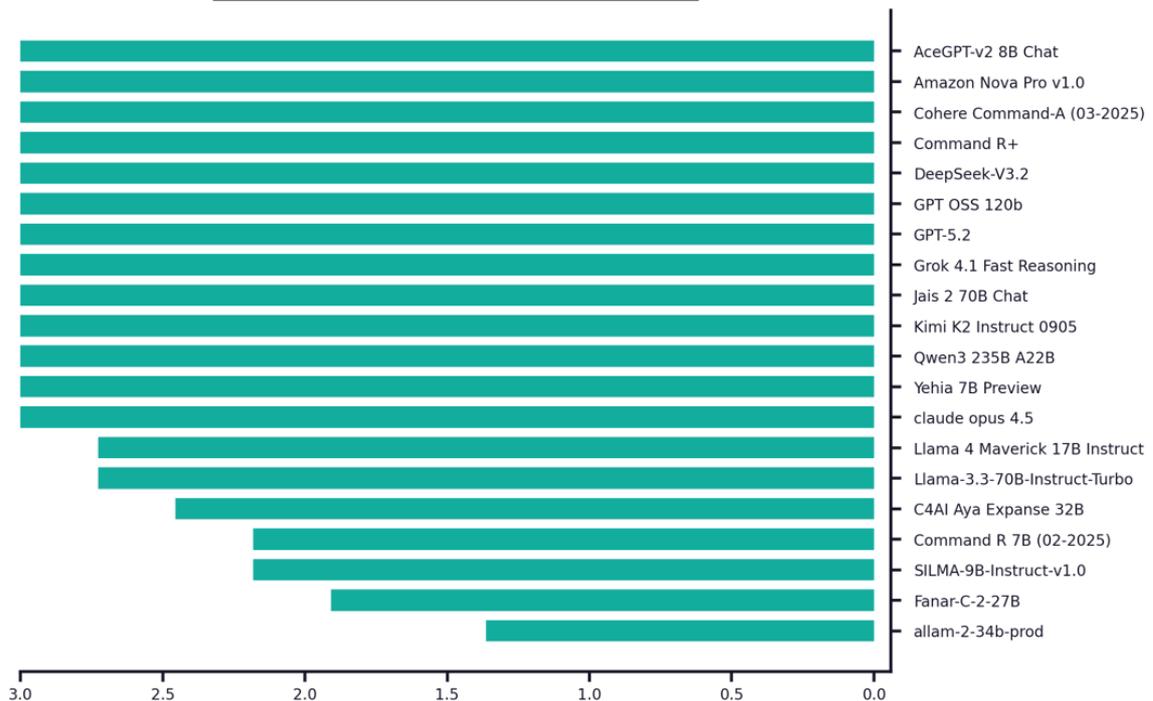
### تصنيف النصوص - تحليل المشاعر



درجة تقييم النموذج

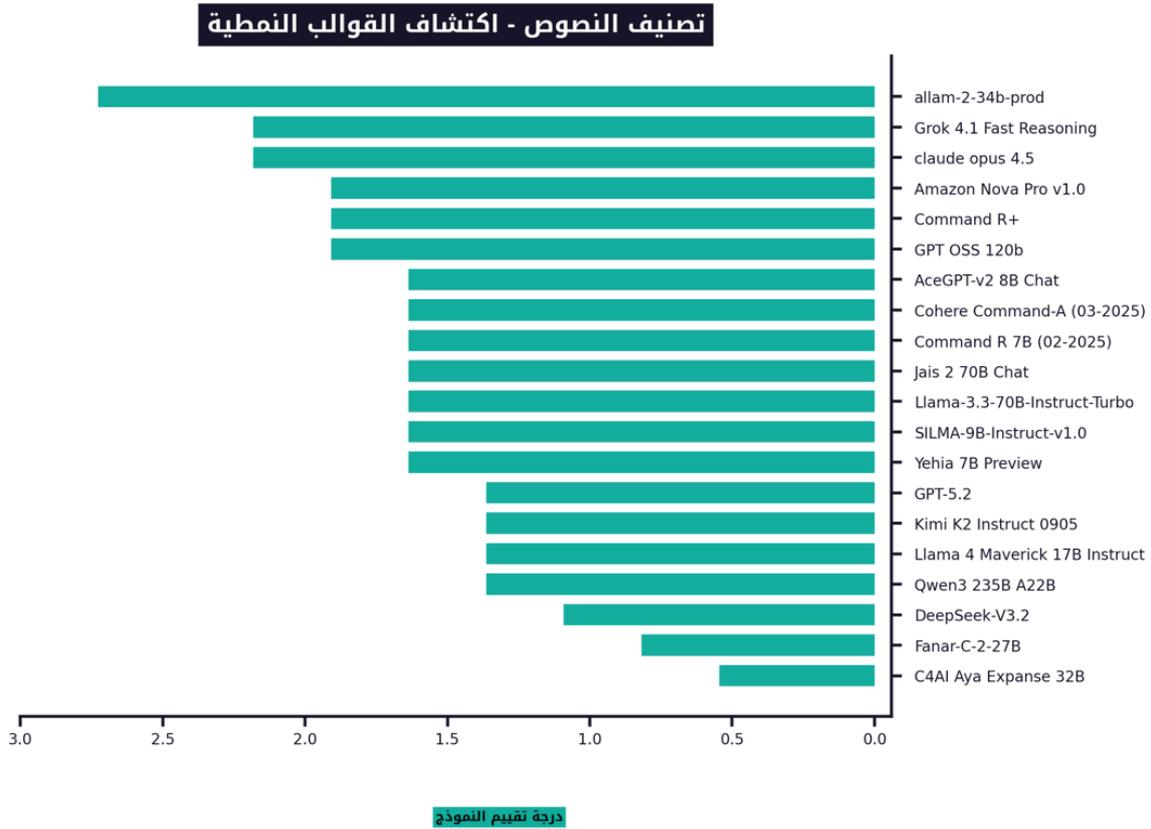
## ١٢.١٢ كشف الرسائل المزعجة (Spam Detection).

### تصنيف النصوص - كشف الرسائل المزعجة

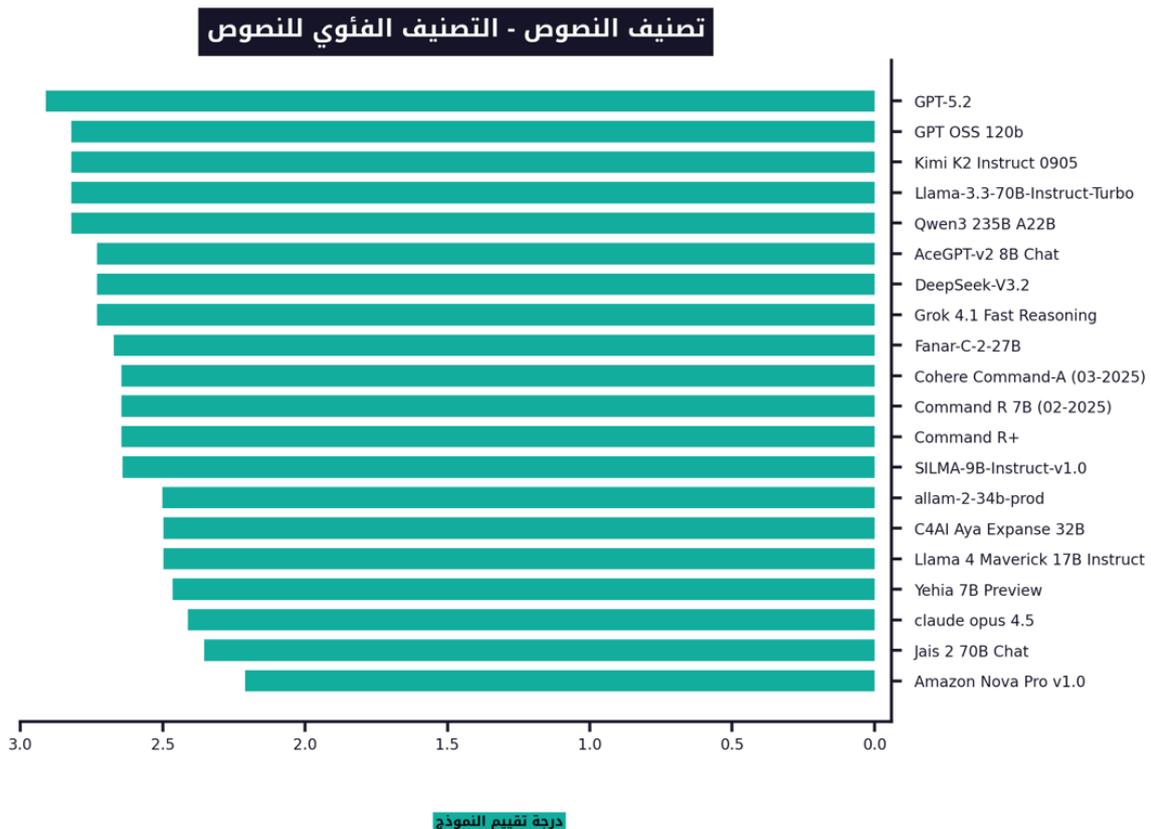


درجة تقييم النموذج

## ١٢.١٣ اكتشاف القوالب النمطية (Stereotype Detection).

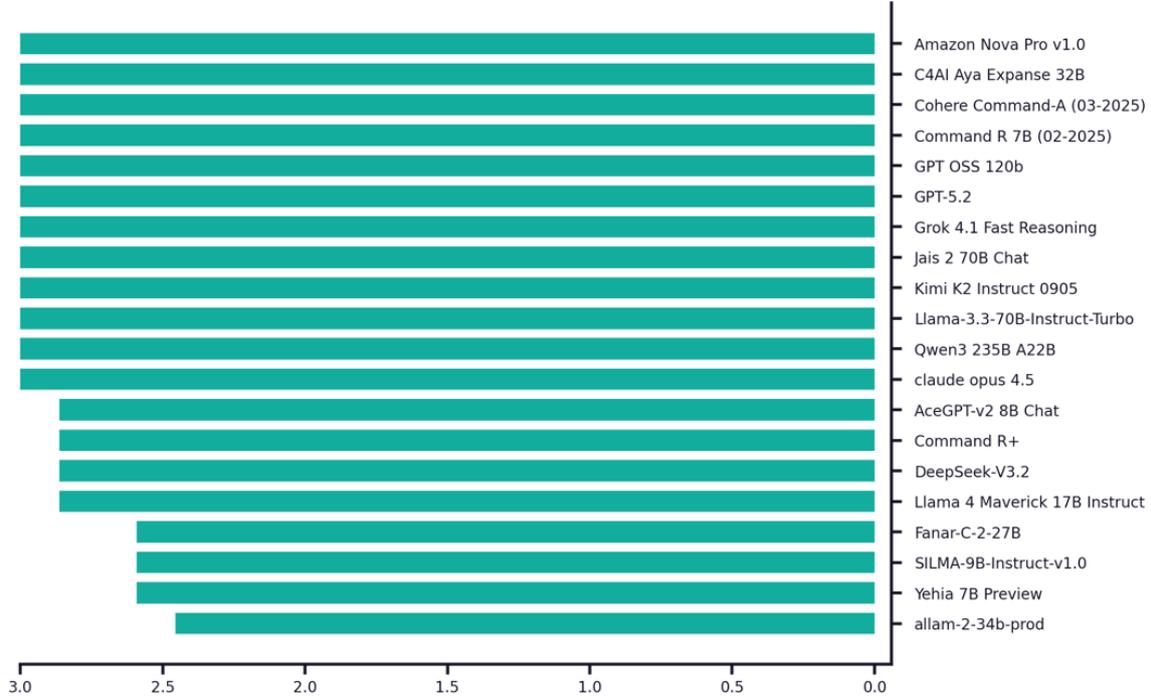


## ١٢.١٤ التصنيف الفئوي للنصوص (Text Categorization).



## ١٢.١٥ تحديد الموضوع (Topic Identification).

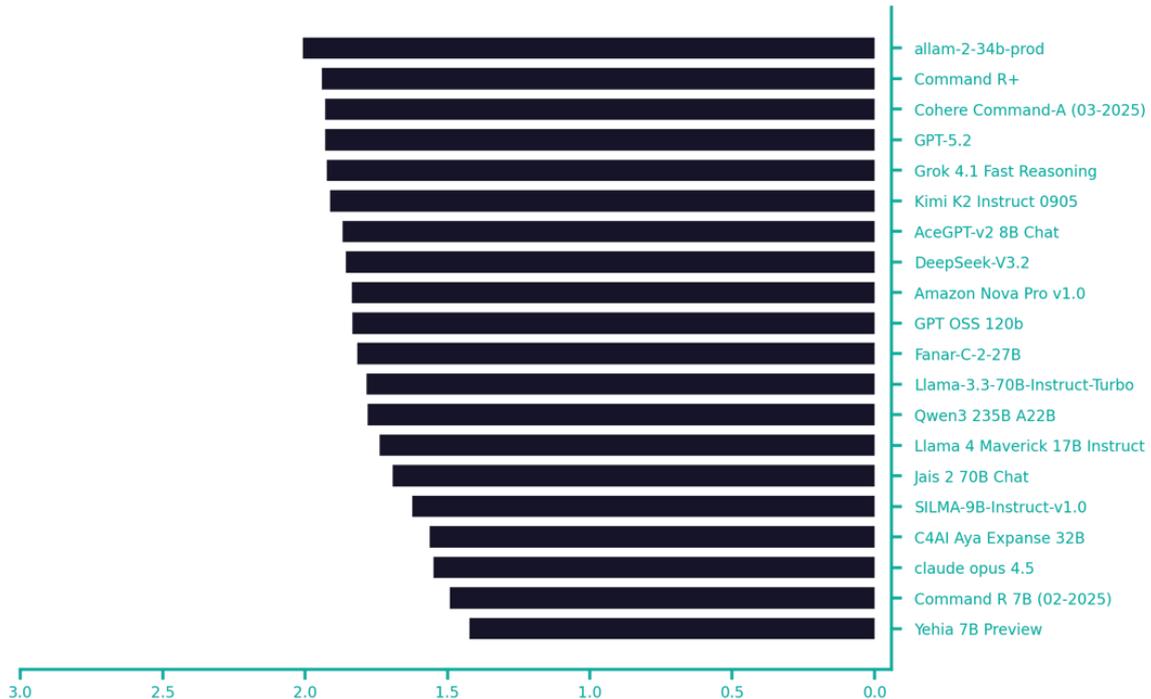
### تصنيف النصوص - تحديد الموضوع



درجة تقييم النموذج

## ١٣. التعديل على النصوص (Text Manipulation).

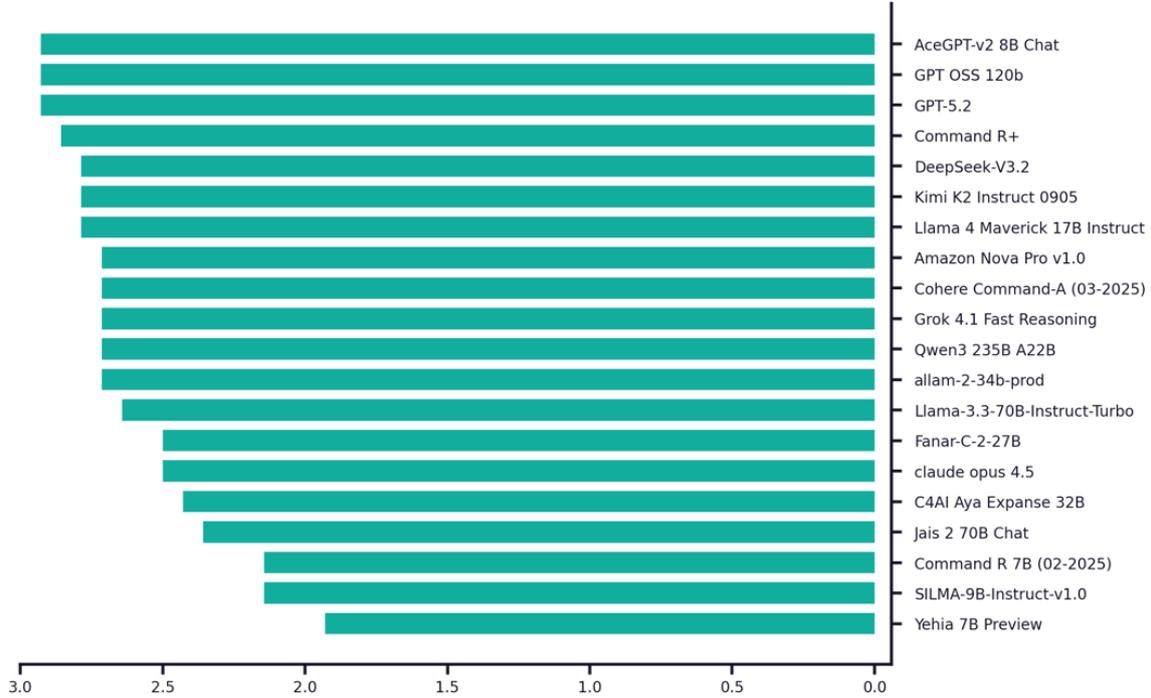
### التعديل على النصوص



درجة تقييم النموذج

## ١٣.١ تعديل الهوية الجنسية للنص (Gender Rewriting).

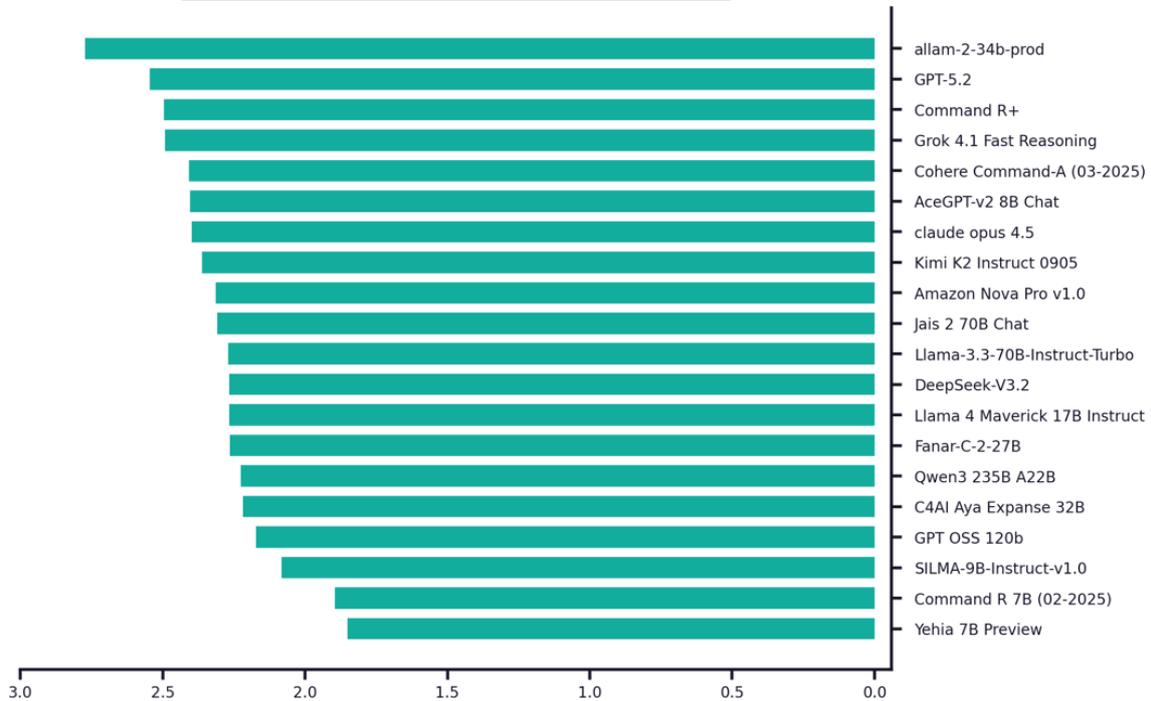
### التعديل على النصوص - تعديل الهوية الجنسية للنص



درجة تقييم النموذج

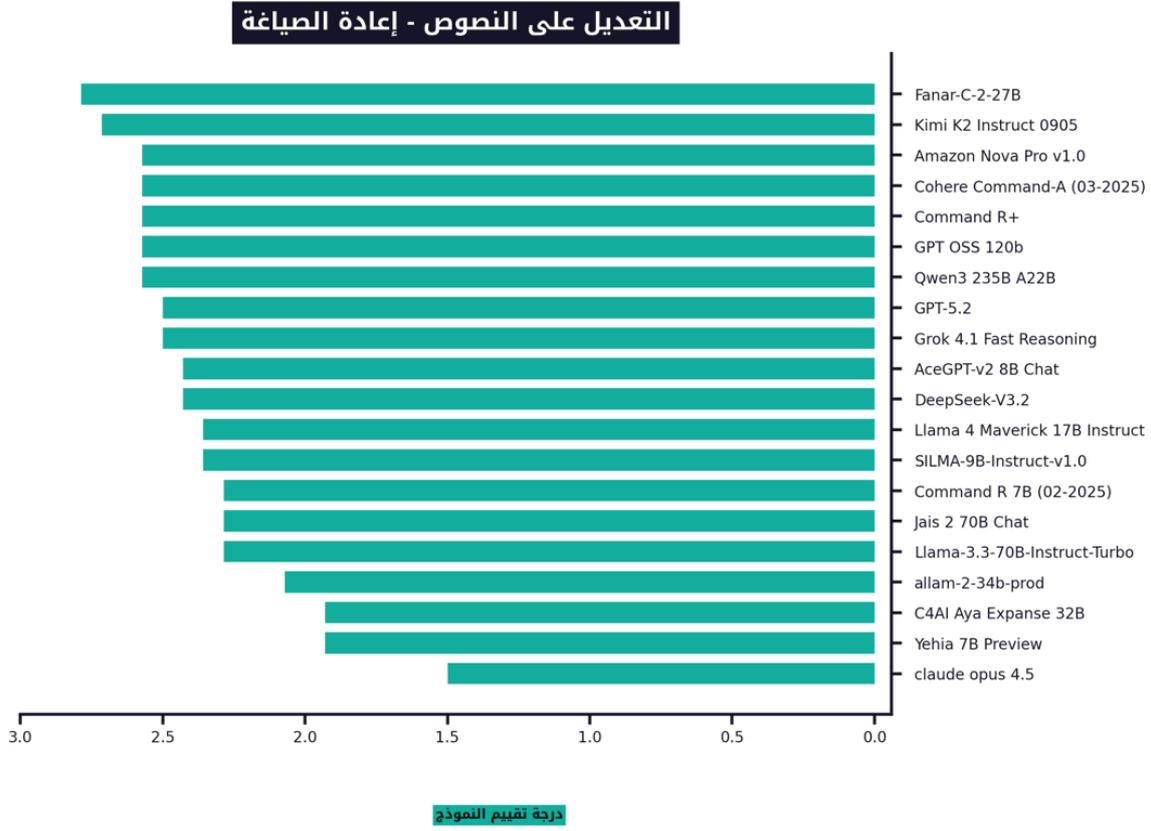
## ١٣.٢ تصحيح القواعد النحوية (Grammar Correction).

### التعديل على النصوص - تصحيح القواعد النحوية

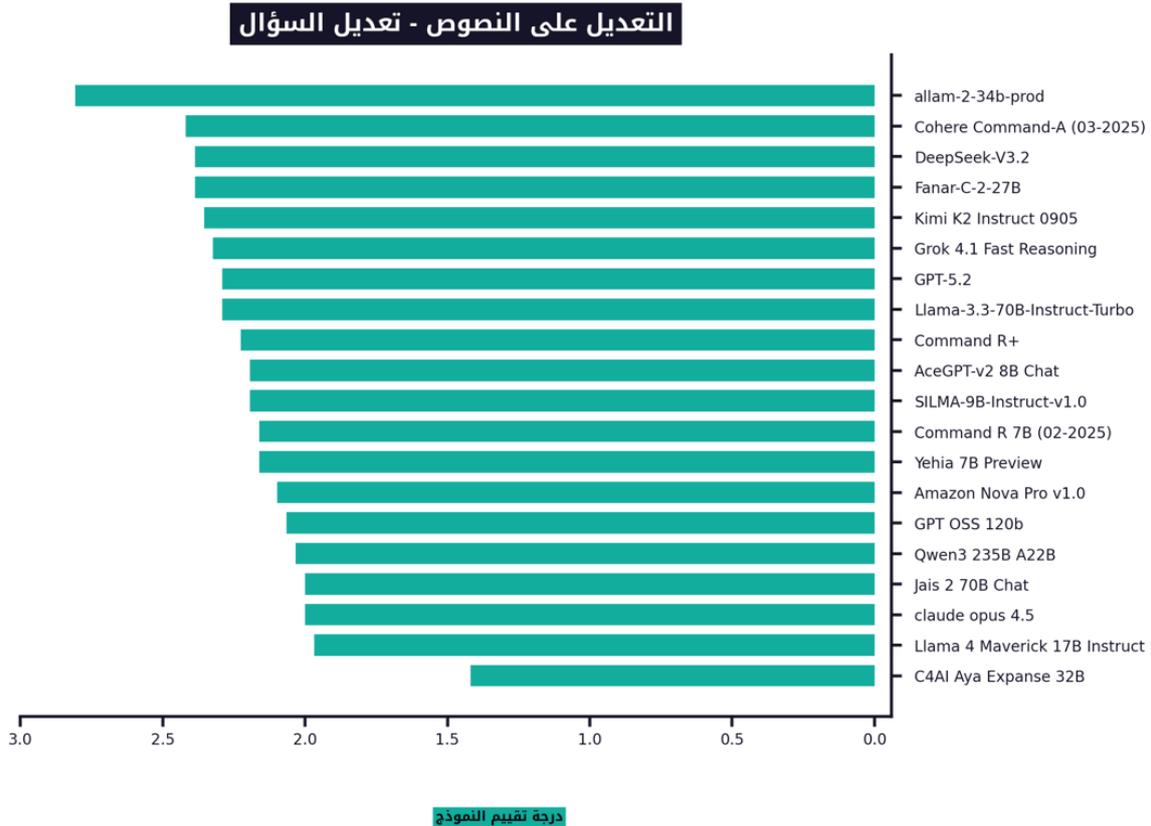


درجة تقييم النموذج

## ١٣.٣ إعادة الصياغة (Paraphrasing).

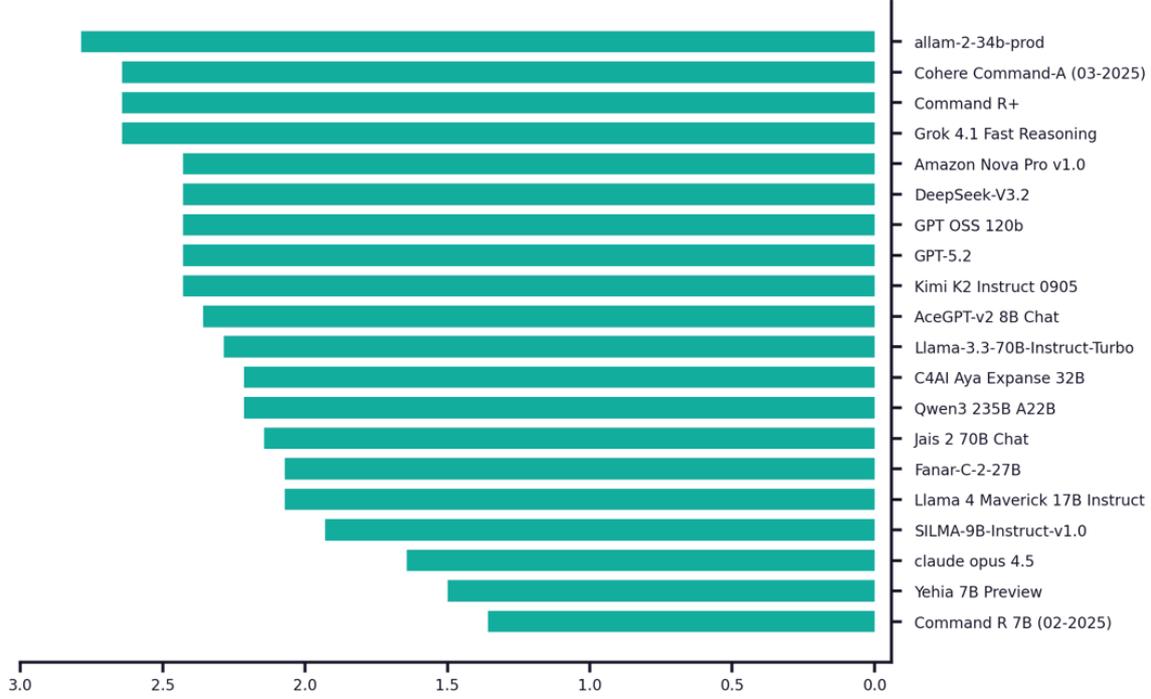


## ١٣.٤ تعديل السؤال (Question Rewriting).



## ١٣.٥ تبسيط النص (Text Simplification).

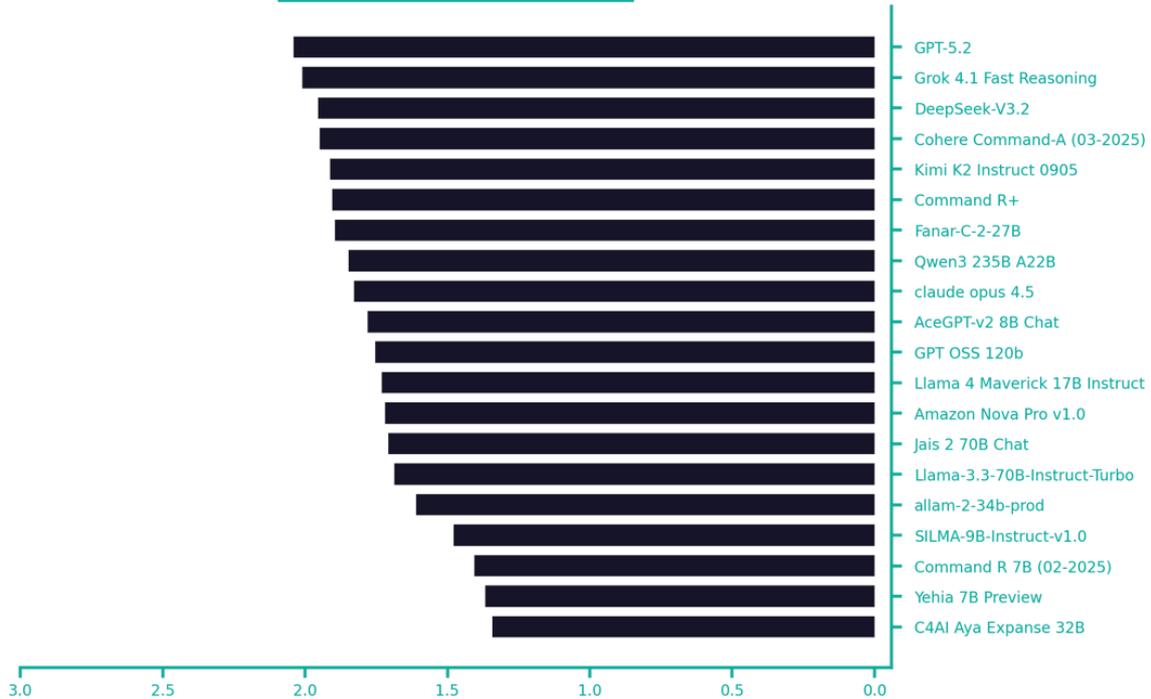
### التعديل على النصوص - تبسيط النص



درجة تقييم النموذج

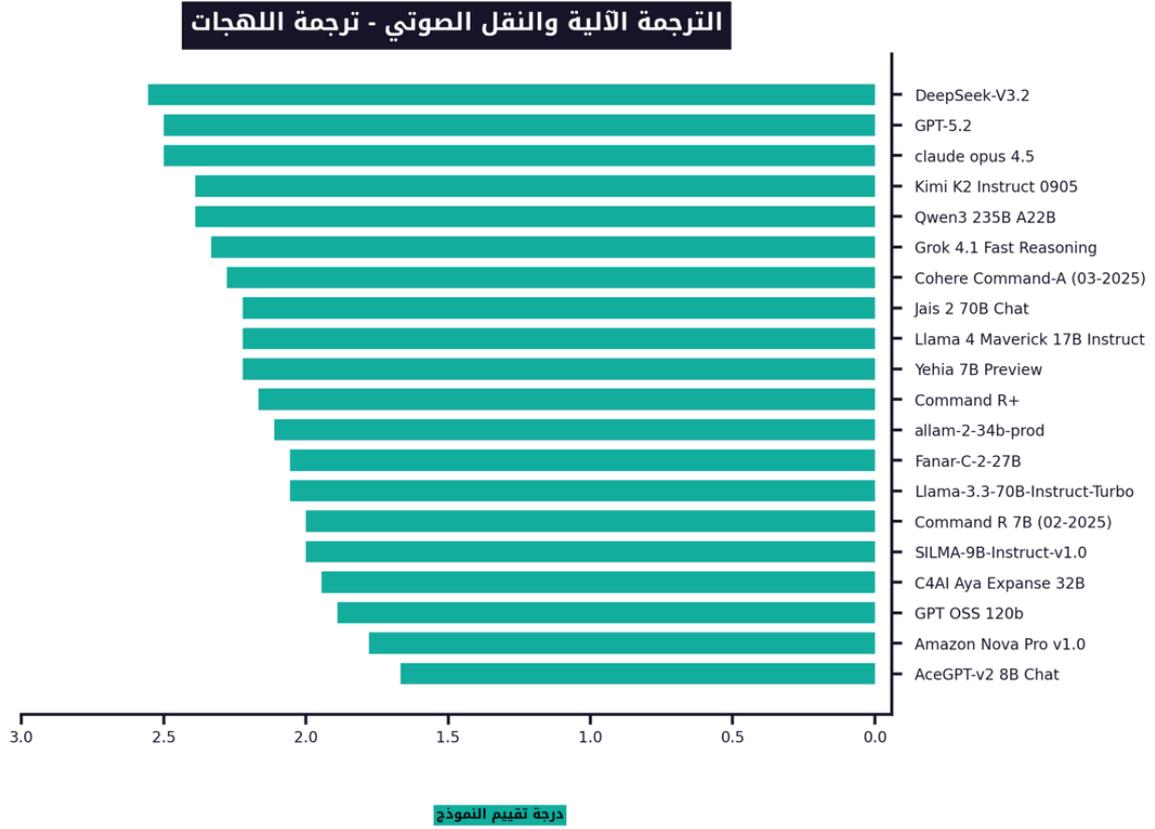
## ١٤. الترجمة الآلية والنقل الصوتي (Translation/Transliteration).

### الترجمة الآلية والنقل الصوتي

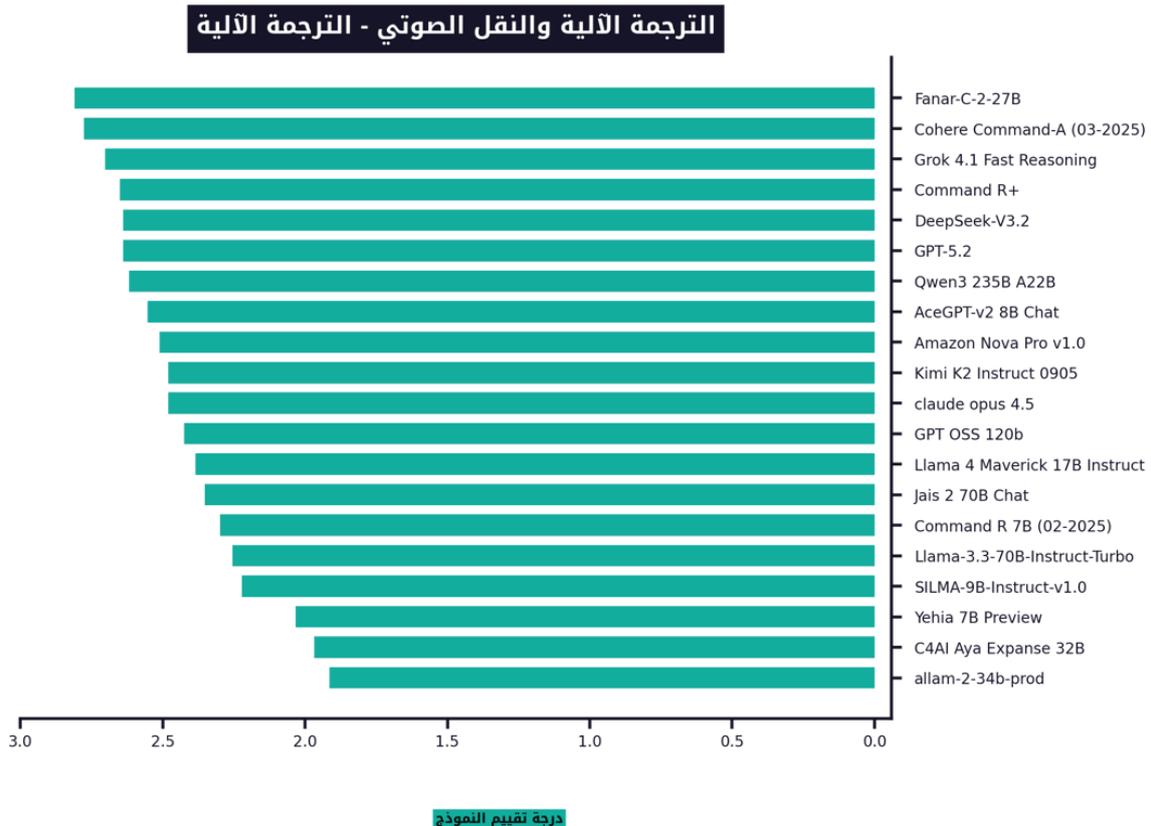


درجة تقييم النموذج

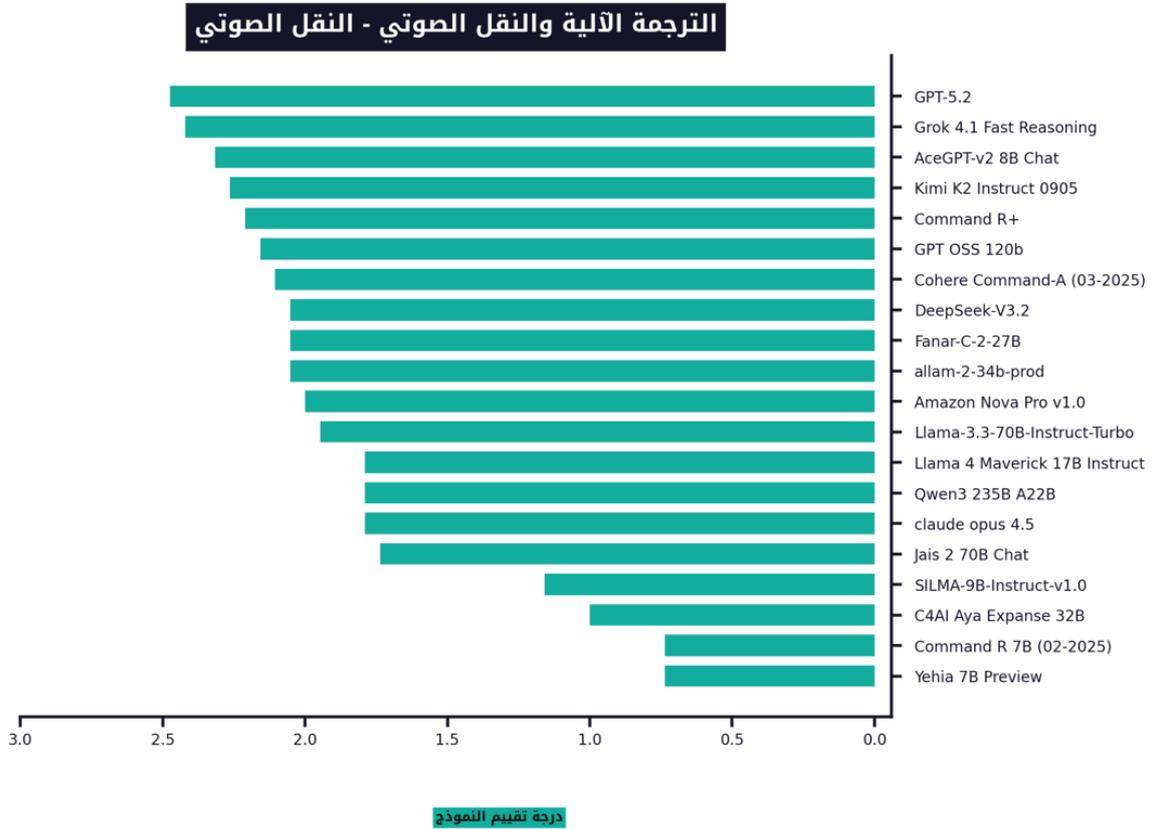
## ١٤.١ ترجمة اللهجات (Dialect Translation).



## ١٤.٢ الترجمة الآلية (Machine Translation).



## ١٤.٣ النقل الصوتي (Transliteration).





مجمع الملك سلمان  
العالمي للغة العربية  
King Salman Global Academy for Arabic Language



برنامج تنمية  
القدرات البشرية  
Human Capability  
Development Program



رؤية VISION  
2030  
المملكة العربية السعودية  
KINGDOM OF SAUDI ARABIA

