

Post 1 - Before LLM

محتاجين قبل ما نعرف شويه عن ال LLM وازاي ال models ديه بتعمل generate لل text ، هو إننا ناخذ خطوة ونشوف الموضوع بدء إزاي من عند ال n-grams ، وازاي كان بناء على ال corpus الى عندي وال frequency بتاعت ال n-grams ققدر اني احدد اني الكلمة ديه هي المتوقع تيجي في ال next predicted token في الجملة ، وده لانها ظهرت ب frequency اكبر من الكلمات الاخر بإنها بتيجي مع الكلمتين او الثلاثه او ال n-grams دول ، مثلاً زي

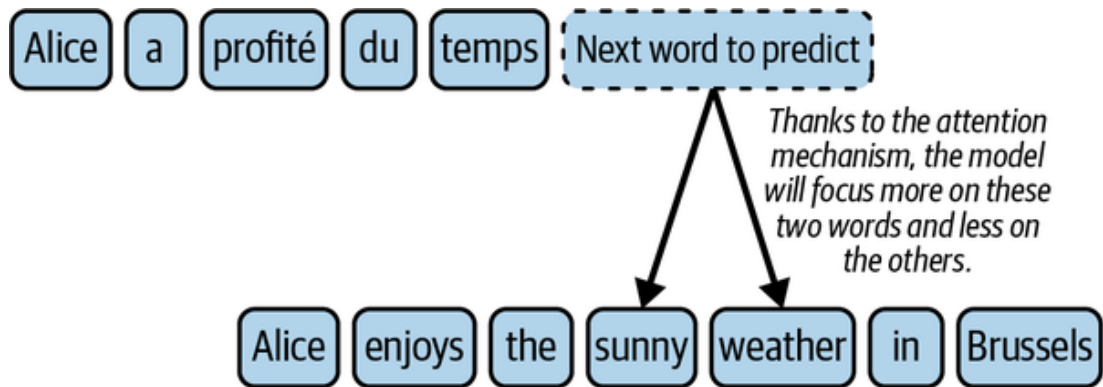
عمر بن ، هنلاقي المتوقع الخطاب ، بس ماذا لو كان سياق الكلام بيكلم عن عمر بن عبدالعزيز ؟ هنا بقا احنا ملناش دعوه بالسياق لاننا معتمدين على ال frequency بتاع ظهور كلمة الخطاب مع عمر بن ، اكثر من ظهورها مع عبدالعزيز رضى الله عنهما وعن سائر أصحاب رسولنا الكريم صل الله عليه وسلم .

جميل ابدينا اننا ناخذ خطوة اخرى نحاول منها اكثر نفهم ال context وكمان ال grammar بتاع الجمل وهنا ظهرت بعد فترى ال RNN and LSTM والى قعدت مستخدمه فترة كبيرة بسبب انها كانت بتدى نتايج كويسة ، وخاصة ظهرت لموضوع الترجمات ، ولكن كان فيه مشكلة اخرى واني النوع ده من ال network كان بياخذ وقت كبير جدا ومينفعش انك تعمله process in parallel بسبب اني كل time step معتمده على ال output بتاع الى قبلها وكمان ال complexity بتبقا اكبر كل ما عدد ال tokens بقا اكبر ، فلو عندك text فيه 5000 token كانك عند netrok فيها 5000 layer !!! وكمان مع حجم ال data الى بنشوف ال LLM متدرب عليها كان هيبقا صعب جدا مع ال models ديه .

ولكن جه بعد كده في 2017 ال paper بتاعت transformers all you need والى غيرت الدنيا في ال nlp ، لانه عالج اولاً مشكلة ال long text ، وكمان ال parallelization ومشاكل ال vanishing and exploding

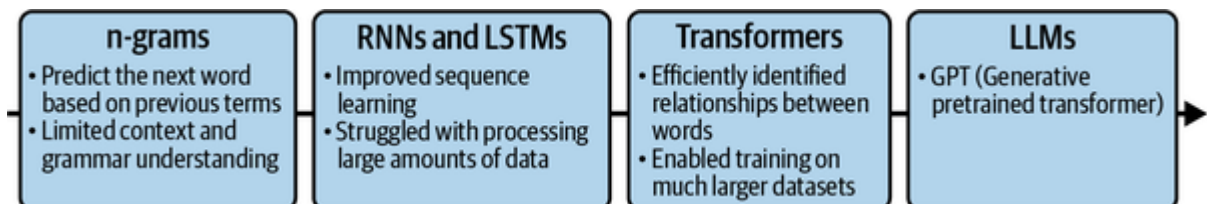
problems الی كانت فی ال RNN and LSTM الی بتحصل بسبب عدد ال tokens .

وهنا ال transformers فيه ما يسمى بال attention mechanism والی فيه انی مش بدی اهمیه لكل الكلمات الی فی ال text بنفس ال weights لا انا ممكن يكون ال predicted token الجای معتمد على كلمتين فقط ومش لازم يكونوا كمان الكلمتين دول ورا بعض ، فهنا ال prediction بتاع ال token هیکون معتمد على الكلمات المهمة لیه والكلمات الاخری هتاخذ ال weights اقل ، زی ما فی المثال :



هنا مثلا الكلمة الی انا هتوقعها فی ال next time step هی ensoleillé والی هی معناها sunny فالبنسبة لیهما اكثر كلمتين هی محتاجهم هما sunny and weather .

بعد كده ننتقل لنقطة ال GPUs وال Parallelization وهی انی ال transformers models قادره انها تهندل ال text بتاعتك بشكل simuntuiosly مش sequentially وده متناسب مع ال GPUs ويخليك تقدر تعمل train على داتا كبيره فی وقت اقل.



Post 2 - LLM

فى 2017 انتشرت attention is all you need والى كان فيه ال transformers مع فكرة ال attention ، وهنا فى ال transformers فى جزئين مهمين هما ال encoder and decoder بعض ال models هى فقط encoder والى بيكون دورة انه يطلع ليك numerical representation لل text بتاعك ، والبعض الاخر وهو ال decoder والى بيكون دورة هو انه يحاول يعمل generate ل text بناء على ال encoder output وحاجات اخرى ، وفى بعض ال models يتضمن الاثنين ال encoder and decoder .

طب هنا بالنسبة لل gpt مثلا فهو decoder فقط وهنا ده بيعتمد على نوع من ال attention وهو ال self attention عشان يقدر يعمل generate لل text .

طيب دلوقت بقا بالنسبة ل LLM زى GPT الى هو شغال ك decoder فقط كان بيروح ليها prompt الى هو كانه ال start بتاع الكلام وهو بيتدى يكمل ليك الكلام ويعمل generate لباقي ال sequence زى

The weather is nice today, so I decided to
هتلاقى هو بناء على ال prompt input ديه هيروح يعمل generate مثلا ل
:

go for a walk

وده يعتبر ما يسمى بال completion .

طب ازاي بقا ال LLM عمل ده ، فى البداية ال prompt الى انتا بعتاهله ديه بتروح يحصلها ما يسمى tokenization عشان فى الاخر تكون عبارته عن list of tokens وهنا طبعا ال tokenization فيه انواع مختلفه ولكن نقدر نقول بشكل عام هو بيكون عبارته عن اما single words or sub-words

or punctuation and others فمثلا ممكن كلمه today تبقا عباره عن
to and day الى هو two tokens .

ممكن تجرب ال tokenization بتاع open ai من هنا:

<https://platform.openai.com/tokenizer>

بعد كده بيكون بناء على ال tokens ديه بيتدى يفهم ال transformers
models with attention mechanism انهى اكثر tokens هتكون
مناسبه انها تروح ت predict ال next token وبعدين يضم ال token ده
لل prompt او ال input الى دخلته وبعدين يشوف ال token الى بعده
ويضمه وهكذا لحد ما يوصل لل max length بتاع ال models او انه يكون
انتا محدد .

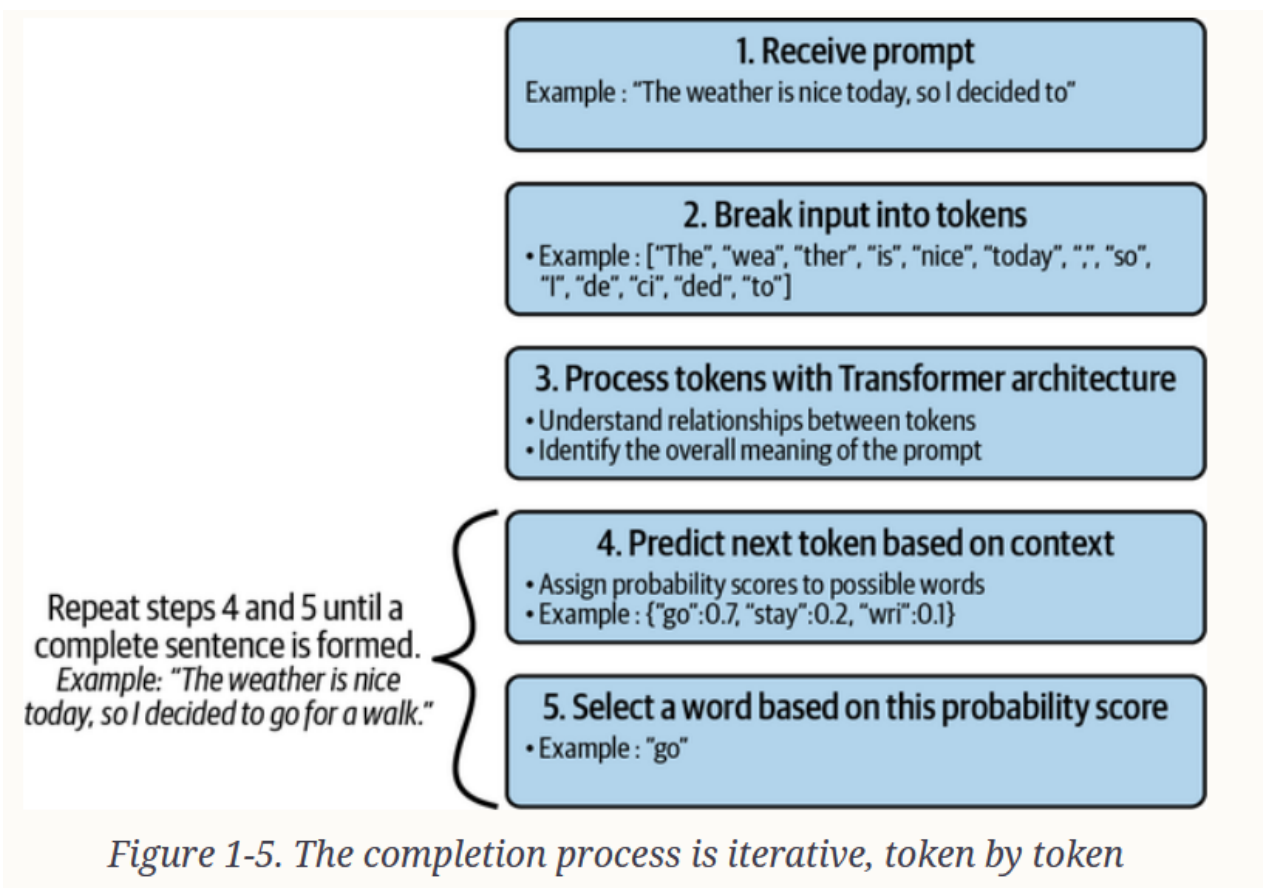


Figure 1-5. The completion process is iterative, token by token

Post 3 - History GPT-1 to GPT-4

قبل ما نكلم على GPT-1 محتاجين نعرف انى كان معظم ال NLP tasks معتمده بشكل او باخر على ال labeled data ، عشان تقدر تعمل train ل model على ال task الى انتا عايزه زى ال sentiment classification ، وده بيستهلك وقت كبير جدا وفلوس كتيره عشان تقدر تبني dataset معتمده على ال manual classification الى من خلاله هتقدر انك تبني ال model الى انتا عايزه .

فجه هنا بعد سنه من ال attention is all you need ، انتشرت ال paper الخاصه ب GPT-1 من open ai ، والى كان فيها خطوة قبل عمليه ال nlp model الى هو مبنى على ال task بتاعك ، وهو إنك الأول ممكن تبني pretrained model يقدر يتوقع ال next token من داتا كبيرة جدا ، وبعدين تاخذ ال model ده تعمله fine tune على ال nlp task الى عندك فى ال classification .

بالطريقه ديه ال model الاول ده عرف معلومات كتيره عن اللغة عرف ال context بتاع اللغة وال grammer بتاع اللغة وغيره وعلى حسب حجم الداتا برضه ال model بيقدر يكون افضل معاك فى ال downstream task ال عندك ، وبكده تقدر تاخذ ال model ده مش بس ل sentiment problem لا تقدر تاخذه لاكثر من task .

لكن هنا GPT-1 كان معموله train على حجم داتا مش كبير فكان يشتغل كويس مع ال completion task الى حد ما ولكن هيتحتاج انك تعمله fine tune لو عايز مثلا تستخدمه فى classification problem . ولكن فى النهايه بتحتاج labeled data صغيره عشان تقدر توصل لنتائج كويسه لل task بتاعتك بدل ما كنا بنحتاج ل massive label data عشان ال model يدرب عليها .

بعد كده حه ال GPT-2 فى 2019 والى كان عبارته عن scalable model من GPT-1 معمولة train على داتا اكبر حوالى 40 جيجا وليه كمان parameters اكثر 1.5 billion وده كان افضل من ناحية ال completment لل text وكمان افضل انك تاخده على down stream task ، بسبب تنوع وزياده ال data لان GPT-1 كان معمولة train على 11 الف كتاب ، فهنا ممكن يكون ال sequence فى الكتب مختلف عن ال social media مختلف عن ال web articles ، مختلف عن ال scientific articles or questions and answers regarding of any domain .
وده الى هيتم معالجته فى ال new GPTs .

Post 4 - History GPT-1 to GPT-4

نيجبى بعد كده لل GPT-3 فى 2020 والنقلة النوعيه فى حجم ال model نفسه والى كان من 1.5 فى GPT-2 ل 175 Billion paramter فى GPT-3 ، وكمان حجم ال data الى حصل عليها train كان كبير جدا ومن مصادر مختلفه فمثلا billions of webpages من :

<https://commoncrawl.org/>

ومن Wikipedia ومن كتب وغيره من ال resources ، هنا ال GPT-3 قدر انه يفهم complex and deep pattern فى ال text مع اختلاف المصادر والكتب والكتاب وغيره من ال data الى حصله train عليها قدر يفهم طريقه الكتابه وبقا عنده فهم عميق للغة ، وكمان بقا يعمل generate لنص متناسق وموزون تقدر تعتبره كان حد الى كاتب الكلام مش AI model ، مش بس كده لانه كمان اتعمله train على حجم ضخم من ال web pages بقا قادر انه generate code ، وكمان ممكن ميحتجش انك تعمل fine tune ويعمل هو ال task الى انتا كنت عايز تعملها fine tune .

ولكن هنا كان على حسب ال data برضه ال model معموله train عليها
فممكن انك تلاقي ال model مطلع محتوى سيء من نواحى مختلفه ، كمان
ممکن يكون فيه عدم تناسق بين ال end user عايزه والى ال model حصله
train عليها بذات انه متعلم ازاي يعمل generate لل next token وده
مش الهدف الاخير الى انا عايز ال model عشانه ، وهنا ظهر model اخر
لل GPT-3 وهو ال Instruct GPT-3 فى 2021 وده كان معموله
optimization من خلال ال reinforcement learning من ال
human feedback ، ببساطة انا هخلى ال model يتعمل بناء على ال
instructions بتاعت ال feed back الى جاي منى فيقدر يخلي ال
model بيتعلم طول الوقت .

فهنا بقا مثلا هنجي نقول لل model :

Explain what is meant by time complexity

مع GPT-3 لو ادبته ال prompt ديه هيروح يكمل بالشكل ده مثلا يعنى :

Explain what is meant by space complexity. Explain
what is meant by the big-O notation.

لكن مع GPT-3 Instruct :

Time complexity is a way of measuring the amount of
time it takes for an algorithm to run and complete its
task. It is usually expressed using Big O notation, which
measures the complexity of an algorithm in terms of the
number of operations it performs. The time complexity of
an algorithm is important because it determines how
efficient the algorithm is and how well it scales with
larger inputs.

فهنا هو فهم انا محتاج ايه ومرحش كمل زى ما فى ال GPT-3 .
طب ازاي ده حصل ؟

Post 5 From GPT-3 to InstructGPT

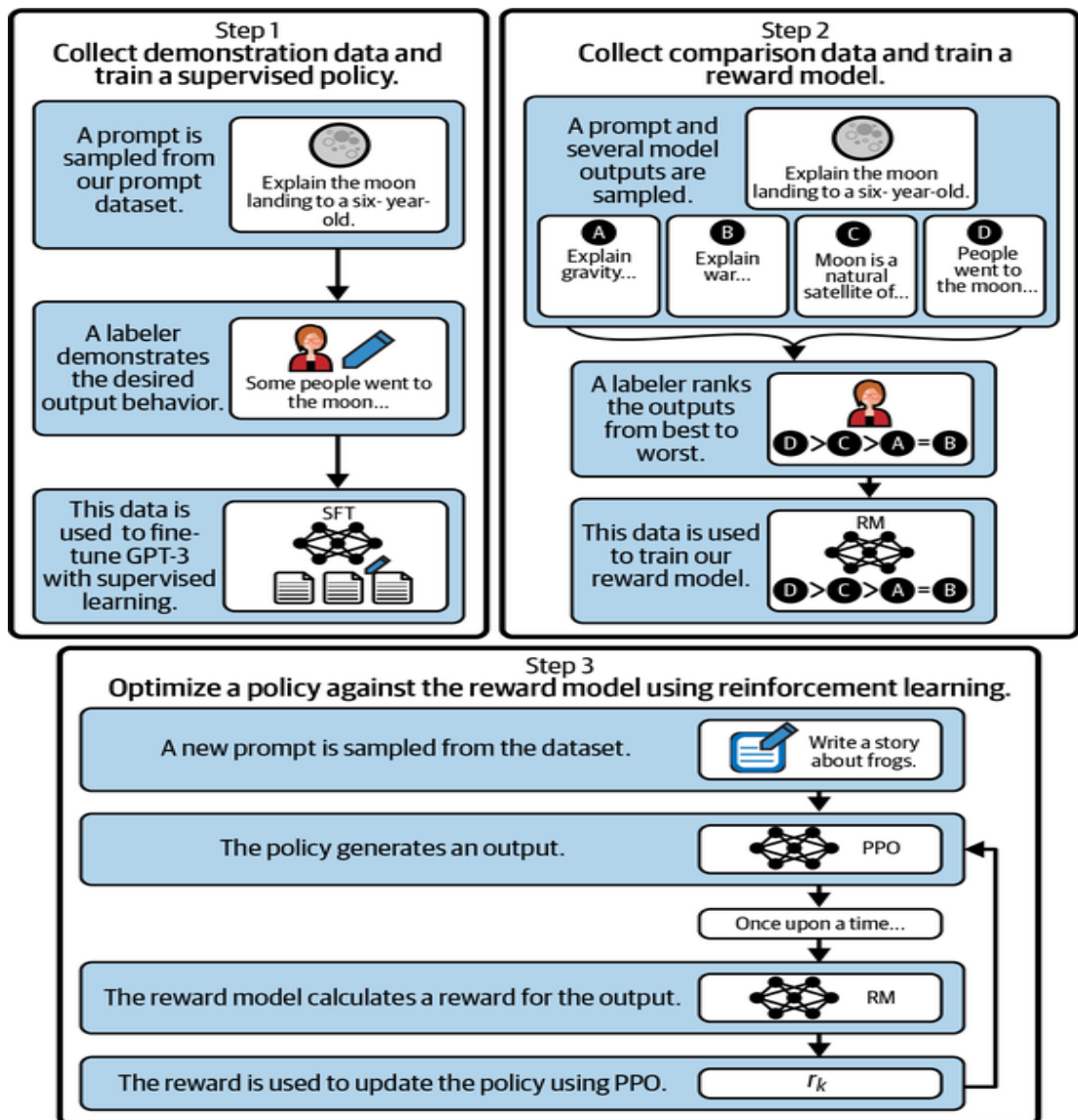
عشان نروح من GPT-3 ل InstructGPT فيه خطوتين مهمين وهما supervised fine-tuning (SFT) and reinforcement learning from human feedback (RLHF)

ال SFT هنا بتاخد ال model الى هو GPT-3 وتبتدى انها تعمل fine tune ولكن ب supervised من human فمثلا ناخد random prompt من ال dataset الى فيها prompt او اسئلة مثلا عندى ، واروح اخلى اكثر من انسان يرد على نفس ال prompt ديه ، فهنا بقا عندى السؤال الواحد ليه اكثر من إجابته على حسب معلومات كل بنى ادم .
هنا هيتكون عندى مجموعه من ال dataset الى بقا ليها + prompt answer from human الى هو دلوقت ققدر اخذ GPT-3 واروح اعمله fine tune بإستخدام ال data الى اتكونت ديه.
هنا ال model ده هو ما يسمى بال SFT .
بال RLHF متقسمه لخطوتين الأولى :

اروح اعمل random prompt من الى عندى ، واخلى ال model هو SFT يطلع اكثر من إجابته بناء على ال prompt ديه لانه كان مدرب من خلال انى نفس ال prompt اتجاوب عليها اكثر من مره من اشخاص مختلفين ، ويبتدى بعد مال model يطلع الاجابات المختلفه ، يبتدى برضه حد عنده خبره يعمل score للإجابات المختلفه ديه ويرتبها على حسب الأفضلية ، بعد كده بقا هناخد ال dataset الى اتكونت ديه والى بقا فيها كل prompt والإجابة بتاعتها ليهم score معين عشان ابتدى ادرب ما يسمى ب rewarded model الى هيكون عبارته عن انى اخذ ال SFT model عشان ادربه على ال scored data ديه فيدينى ما يسمى بال Rewarded model .
الهدف من ال Rewarded model هو انى ققدر اعمل score دلوقت للاجابه الى ال model هيطلعها ، بمعنى انى بناء على ال answer الى ال model هيدهانى هتكون متناسبه مع ال prompt ولا لا هيكون ال score ده اما عالى او واطى .

الخطوة الثانية بقاء هي اني اروح اخذ ال Rewarded model ده واستخدم
مع Reinforcement Learning عشان يبقا عندي ما يسمى بال
InstructGPT-3 .

هنا بقاء هنبتي ناخذ random prompt من اي user واخلى مثلا ال GPT
model او ال SFT model يجاوب على ال Prompt ديه ، بعدين بناء
على ال output ده هروح بال Rewarded Model ادى score لل
output الى model طلعه واعمله evaluation ، بناء على ال
rewarded الى ال model استلمها يبتدى يعمل update لنفسه وهنا انتقلنا
اني يكون الموضوع automated مش معتمد على human بقاء معتمد على
اني ال model يعمل generate و model تاني يقيم ال generation ده
و model تالت يحصله update بناء على التقييم الى استلمه.



Post 6 GPT-3.5, Codex, and ChatGPT

بعد كده فى 2022 open ai عرفت GPT-3.5 and codex ، وديه كان عبارة عن new versions من ال GPT-3 تتدرج تحت GPT-3.5 ، ال codex هو LLM ولكنه معموله train على billions of lines of codes ، وهو المعتمد عليه فى Github copilot ، ولكن codex بقا deprecated وبقا المستخدم حاليا هو GPT-3.5 turbo and GPT4 .

بعدين بقا فيه Copilot X الى بقا معتمد على GPT-4 . وفى الغالب chatgpt ال base model بتاعه هو ال GPT-3.5 ، والى هو ظهر فى اواخر 2022 وعمل ضجه كبيره فى عالم ال NLP وفى ال Business بشكل عام ، وفى مجالات كتير ، وبقينا حتى بنستخدمه عشان نتعلم من خلاله ك interactive teacher ليك انتا لوحدك .

فى 2023 ابتدى ظهور GPT-4 والى هو كان بيدى نتائج افضل واكثر دقه من chatgpt على رغم من انى chatgpt بيدى نتائج كويسه الا ان GPT-4 افضل وكمان فيه انك تقدر تستخدم معاه الصور وتقدر تعرف الى فى الصورة او تسال اسئله عن ايه الى فى الصورة بظبط . وكمان فى بعض ال exams مقارنة بين GPT-4 and chatgpt لقوا انى chatgpt جايب 10% فى واحد و 31% فى التانى بينما GPT-4 جايب 90% و 99% !

Table 1-1. Evolution of the GPT models

2017	The paper "Attention Is All You Need" by Vaswani et al. is published.
2018	The first GPT model is introduced with 117 million parameters.
2019	The GPT-2 model is introduced with 1.5 billion parameters.
2020	The GPT-3 model is introduced with 175 billion parameters.
2022	The GPT-3.5 (ChatGPT) model is introduced with 175 billion parameters.
2023	The GPT-4 model is introduced, but the number of parameters is not disclosed.

uses فى اخر ال **chapters** الكاتب اكلم عن ال
cases المستخدمه حاليا من الشركات مبنية على
GPT-4 و **chatgpt** منها:

Be My Eyes

Morgan Stanley

Khan Academy

Duolingo

Yabble

Waymark