Interrogez vos documents avec des grands modèles de langage: découvrez la méthode RAG

ANF TDM & IA 2025

Pierre Poulain

pierre.poulain@u-paris.fr

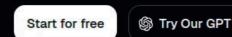
Laboratoire de Biochimie Théorique

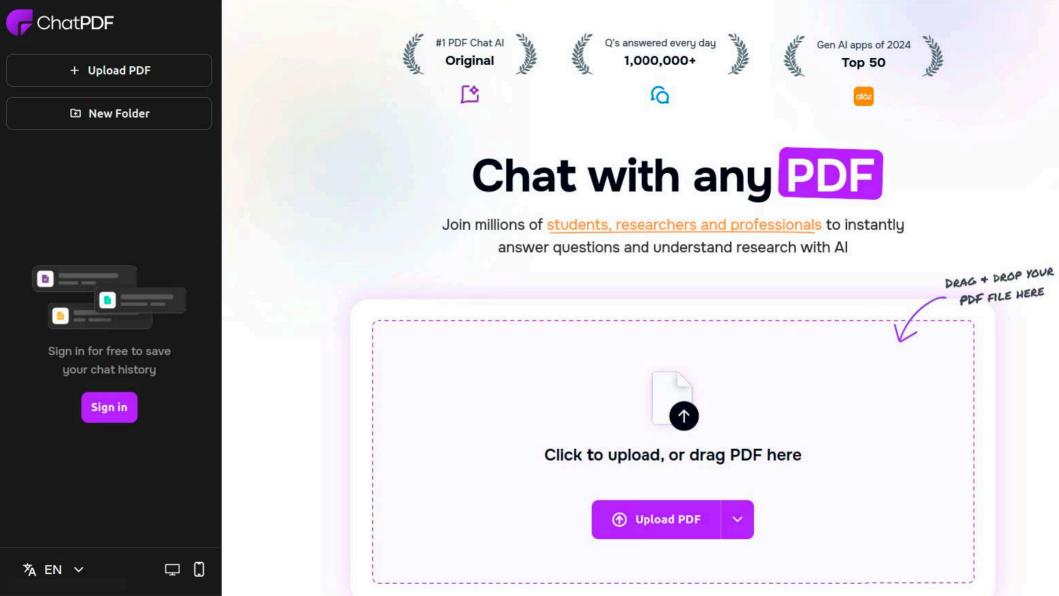




Say Goodbye to Endless Reading, Get Instant **Answers from Any PDF**

Ask any question, summarize, extract key insights, and ace your exams faster than ever.

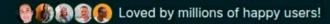


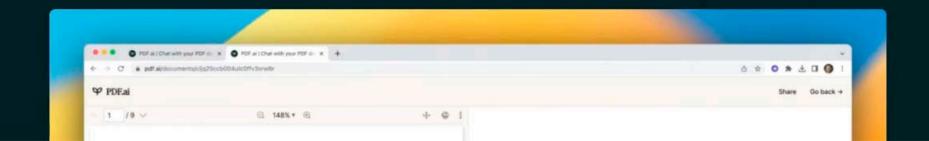


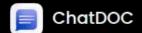
Chat with any PDF document

From legal agreements to financial reports, PDF.ai brings your documents to life.

You can ask questions, get summaries, find information, and more.









Where DOCX Speaks Accurately

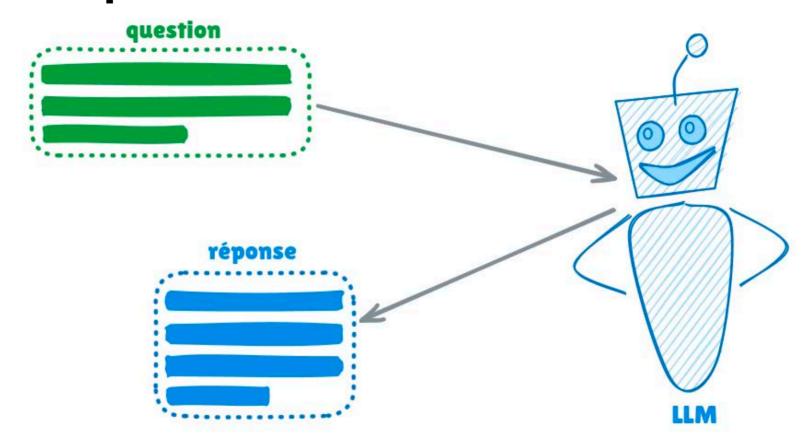
Answers You Can Trust. Sources You Can See.

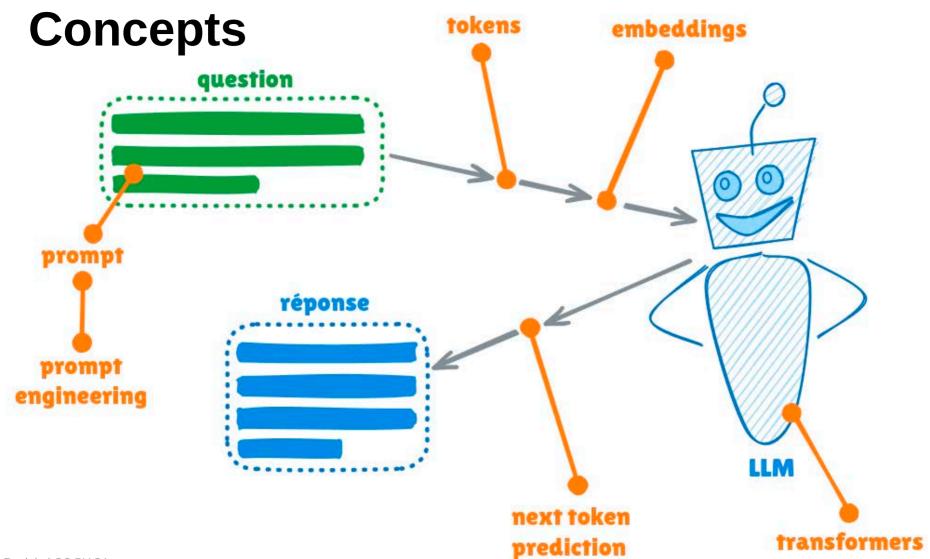


Au programme

- 1. La méthode RAG
- 2. Construction d'un RAG
- 3. RAG vs fine-tuning
- 4. Cas d'usage pédagogique

Concepts





Transformers

Attention Is All You Need

Ashish Vaswani Google Brain avaswani@google.com

Noam Shazeer Google Brain noam@google.com

Niki Parmar* Google Research nikip@google.com

Jakob Uszkoreit*

Google Research usz@google.com

Llion Jones* Google Research llion@google.com

Aidan N. Gomez* † University of Toronto aidan@cs.toronto.edu Łukasz Kaiser* Google Brain

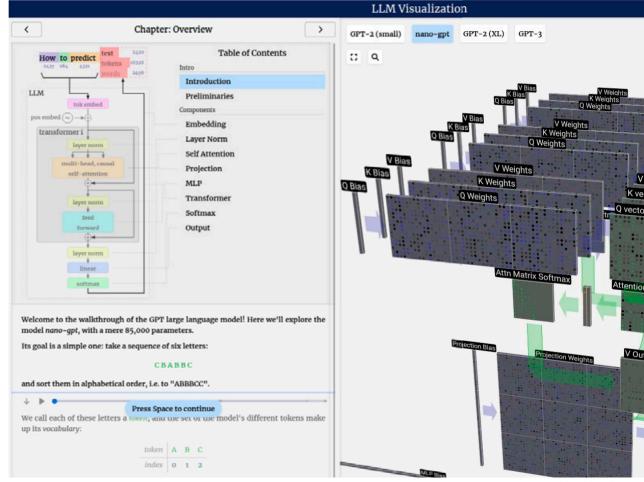
lukaszkaiser@google.com

Illia Polosukhin* ‡ illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 Englishto-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU, On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA,



^{*}Equal contribution. Listing order is random. Jakob proposed replacing RNNs with self-attention and started the effort to evaluate this idea. Ashish, with Illia, designed and implemented the first Transformer models and has been crucially involved in every aspect of this work. Noam proposed scaled dot-product attention, multi-head attention and the parameter-free position representation and became the other person involved in nearly every detail. Niki designed, implemented, tuned and evaluated countless model variants in our original codebase and tensor2tensor. Llion also experimented with novel model variants, was responsible for our initial codebase, and efficient inference and visualizations. Lukasz and Aidan spent countless long days designing various parts of and implementing tensor2tensor, replacing our earlier codebase, greatly improving results and massively accelerating

Work performed while at Google Brain.

[‡]Work performed while at Google Research.

Prompt engineering

Méthodes qui optimisent la réponse d'un modèle de langage



Source: X/Twitter, 2023

- Chain-of-Thought (CoT)
- Zero-shot
- Few-shot
- •
- RAG

→ https://www.promptingguide.ai/fr/

Tokens

GPT-3 (Legacy) GPT-3.5 & GPT-4 GPT-4o & GPT-4o mini La recherche d'information assistée par IA Panorama des enjeux, usages et bonnes pratiques Show example Clear **Tokens** Characters 18 90 La recherche d'information assistée par IA Panorama des enjeux, usages et bonnes pratiques Token IDs Text

P. Poulain | CC BY-SA

11

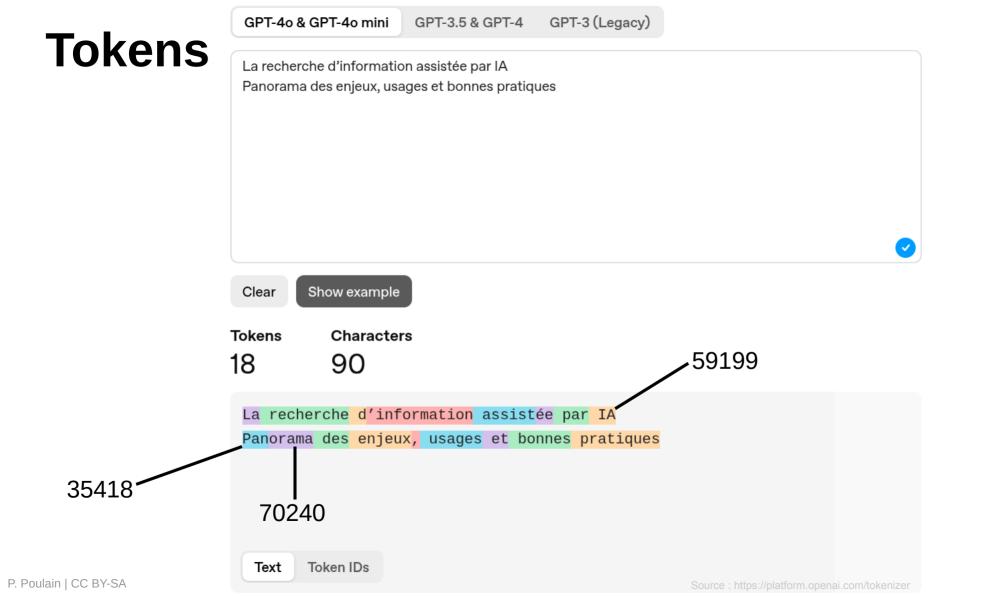
Tokens

GPT-3.5 & GPT-4 GPT-3 (Legacy) GPT-4o & GPT-4o mini La recherche d'information assistée par IA Panorama des enjeux, usages et bonnes pratiques Show example Clear **Tokens** Characters 18 90

[4579, 24975, 272, 95645, 7756, 2894, 686, 59199, 198, 35418, 70240, 731, 129412, 11, 146156, 859, 66702, 65641]

Text Token IDs

12

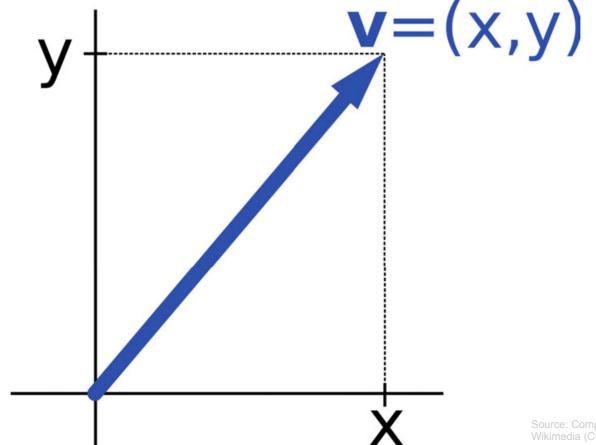




P. Poulain | CC BY-SA Source : https://platform.openai.com/tokenizer 14

1 token ▶ 1 embedding ▶ 1 vecteur

1 token ▶ 1 embedding ▶ 1 vecteur (2 coordonnées)

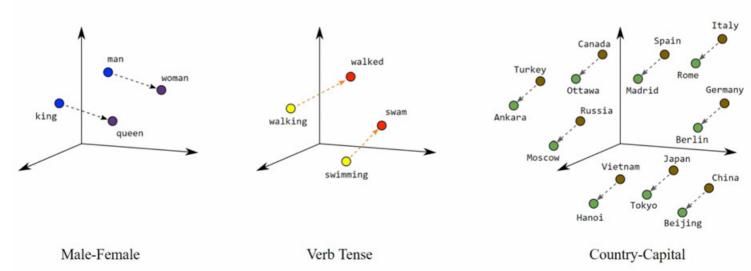


« IA » ▶ 1 token (n° 59199) ▶ 1 embedding ▶ 1 vecteur (3 072 coord.)

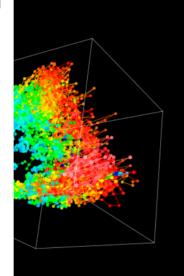
 $0.038047775626182556,\ 0.013517363928258419,\ -0.025667386129498482,\ -0.01796865649521351,\ 0.024701328948140144,\ -0.044141363352537155,\ -0.005004917271435261,\ 0.0453898087143898,\ -0.01796865649521351,\ 0.024701328948140144,\ -0.044141363352537155,\ -0.005004917271435261,\ 0.0453898087143898,\ -0.04701328948140144,\ -0.044141363352537155,\ -0.005004917271435261,\ 0.0453898087143898,\ -0.0453898087143898,\ -0.0453898087143898,\ -0.0453898087143898,\ -0.0453898087143898,\ -0.0453898087143898,\ -0.0453898087143898,\ -0.0453898087143898,\ -0.0453898087143898,\ -0.0453898087143898,\ -0.04538988714389,\ -0.045389887143898,\ -0.04538988714389$ 0.005097807385027409, -0.0106935054063797, -0.03376740217208862, 0.020391229540109634, -0.00947478786110878, 0.0029278951697051525, 0.02025746740400791, 0.0022609445732086897, -0.0033767402172088620.0026585140731185675, -0.024939127266407013, -0.015709569677710533, 0.04176338016986847, -0.05103752389550209, 0.005944964941591024, 0.007553820963948965, 0.02484995312988758, -0.024939127266407013, -0.015709569677710533, 0.04176338016986847, -0.05103752389550209, 0.005944964941591024, 0.007553820963948965, 0.02484995312988758, -0.024939127266407013, -0.015709569677710533, 0.04176338016986847, -0.05103752389550209, 0.005944964941591024, 0.007553820963948965, 0.02484995312988758, -0.024939127266407013, -0.015709569677710533, 0.04176338016986847, -0.05103752389550209, 0.005944964941591024, 0.007553820963948965, 0.02484995312988758, -0.024939127266407013, -0.015709569677710533, 0.04176338016986847, -0.05103752389550209, 0.005944964941591024, 0.007553820963948965, 0.02484995312988758, -0.02484995312988758, -0.02484995312988758, -0.02484995312988758, -0.02484995312988758, -0.02484995312988758, -0.02484995312988758, -0.02484995312988758, -0.02484995312988758, -0.02484995312988758, -0.02484995312988758, -0.02484995312988758, -0.02484995312988758, -0.02484995312988758, -0.02484995312988758, -0.02484995312988, -0.02484995312988, -0.024849888, -0.02484988, -0.02484988, -0.02484988, -0.02484988, -0.02484988, -0.02484988, -0.02484988, -0.02484988, -0.0248498, -0.0248498, -0.0248498, -0.0248498, -0.0248498, -0.0248498, -0.0248498, -0.0248498, -0.024848, -0.024848, -0.024848, -0.024848, -0.02488, -0.024848, -0.02488, -0.0240.004243218805640936, 0.08275391161441803, -0.012410114519298077, -0.0031508314423263073, -0.00858304277062416, 0.005996983498334885, 0.016200030222535133, 0.018533429130911827, 0.00858304277062416, 0.005996983498334885, 0.016200030222535133, 0.018533429130911827, 0.00858304277062416, 0.0085830427062416, 0.0085830427062416, 0.0085830427062416, 0.0085830427062416, 0.008580416,0.015308285132050514, 0.019588658586144447, 0.002818285021930933, 0.011154240928590298, 0.0023761282209306955, 0.04874871298670769, -0.009571393951773643, -0.007156251464039087, -0.00715625146403908, -0.007156251464009, -0.007156251464009, -0.007156251464009, -0.007156251464009, -0.007156251464009, -0.007156251464009, -0.007156251464009, -0.007156251464009, -0.007156251464009, -0.007156251464009, -0.007156251464009, -0.007156251464009, -0.007156251464000, -0.007156251464000, -0.007156251464000, -0.007156251464000, -0.007156251464000, -0.007156251464000, -0.007156251464000, -0.007156251464000, -0.007156251464000, -0.007156251464000, -0.007156251464000, -0.007156251464000, -0.007156251464000, -0.0071562514640000, -0.00715625146000, -0.00715625140000, -0.0071560000, -0.0071560000, -0.0071560000, -0.0071560000000, -0.00715600000000000000.048451464623212814, 0.029100604355335236, -0.004376980476081371, -0.014134153723716736, -0.0079736839979887, 0.04206062853336334, 0.0009400476119481027, -0.013732869178056717, -0.014134153723716736, -0.0079736839979887, 0.04206062853336334, 0.0009400476119481027, -0.013732869178056717, -0.014134153723716736, -0.0079736839979887, 0.04206062853336334, 0.0009400476119481027, -0.013732869178056717, -0.014134153723716736, -0.0079736839979887, 0.04206062853336334, 0.0009400476119481027, -0.013732869178056717, -0.014134153723716736, -0.0079736839979887, 0.04206062853336334, 0.0009400476119481027, -0.013732869178056717, -0.014134153723716736, -0.0079736839979887, 0.04206062853336334, 0.0009400476119481027, -0.013732869178056717, -0.014134153723716736, -0.0079736839979887, 0.04206062853336334, 0.0009400476119481027, -0.014734153723716736, -0.0079736839979887, 0.04206062853336334, 0.0009400476119481027, -0.014734153723716736, -0.0079736839979887, 0.04206062853336334, 0.0009400476119481027, -0.01473689178056717, -0.01473689178056717, -0.01473689178056717, -0.01473689178056717, -0.01473689178056717, -0.01473689178056717, -0.01473689178056717, -0.01473689178056717, -0.01473689178056717, -0.01473689178056717, -0.01473689178056717, -0.01473689178056717, -0.01473689178056717, -0.01473689178056717, -0.01473689178056717, -0.01473689178056717, -0.01473689178056717, -0.01473689178056717, -0.01473689178056717, -0.014706717807178, -0.014706717807178, -0.014706717807178, -0.0147067178,0.005833496805280447, -0.009244420565664768, 0.005372762214392424, -0.07710619270801544, -0.02529582567512989, 0.009229558520019054, 0.00827836338430643, -0.038315299898386, -0.038315299898880.03786942735314369, -0.02476077899336815, -0.022055819630622864, -0.04307126998901367, 0.032102812081575394, 0.024909403175115585, -0.034510523080825806, -0.03338097780942917, 0.0013979956274852157, -0.029709963127970695, -0.044468339532613754, -0.003819640027359128, 0.0014054268831387162, 0.03251895681023598, 0.019202236086130142, -0.022516554221510887, -0.02251655421510887, -0.02251655421510887, -0.02251655421510887, -0.02251655421510887, -0.02251655421510887, -0.02251655421510887, -0.02251655421510887, -0.0225165410887, -0.0225165410887, -0.02251655410887, -0.0225165410887, -0.02251655410887, -0.02251655410887, -0.0225165541088, -0.022516541088, -0.0225165541088, -0.0225165541088, -0.0225165541088, -0.0225165541088, -0.0225165541088, -0.0225165541088, -0.022516541088, -0.022516541088, -0.022516541088, -0.022516541088, -0.02251688, -0.02251688, -0.02251688, -0.02251688, -0.02251688, -0.02251688, -0.02251688, -0.02251688, -0.022516554108, -0.02251688, -0.02251688, -0.02251688, -0.02251688, -0.02251688, -0.02251688, -0.02251688, -0.0225168, -0.02251688, -0.0225168, -0.02251688, -0.0225168, -0.02251688, -0.02251688, -0.0220.016824251040816307, -0.030854368582367897, 0.01740388572216034, -0.000582885229960084, 0.014988742768764496, 0.034094374626874924, -0.01165213156491518, 0.01691342517733574, 0.016824251040816307, -0.034094374626874924, -0.01165213156491518, 0.01691342517733574, 0.016824251040816307, -0.034094374626874924, -0.01165213156491518, 0.01691342517733574, 0.016824251040816307, -0.034094374626874924, -0.01165213156491518, 0.01691342517733574, 0.016824251040816307, -0.000582885229960084, 0.014988742768764496, 0.034094374626874924, -0.01165213156491518, 0.01691342517733574, 0.016824251040816307, -0.000582885229960084, 0.014988742768764496, 0.034094374626874924, -0.01165213156491518, 0.01691342517733574, 0.0169134251773574, 0.0169134251773574, 0.0169134251773574, 0.0169134251773574, 0.0169134251773574, 0.0169134251773574, 0.0169134251773574, 0.0169134251773574, 0.0169134251773574, 0.0169134251773574, 0.016913425177574, 0.016913425177574, 0.016913425177574, 0.016913425177574, 0.016913425177574, 0.016913425177574, 0.016913425177574, 0.016913425177574, 0.016913425177574, 0.016913425177574, 0.016913425177574, 0.016913425177574, 0.016913425177574, 0.016913425177574, 0.016913425177574, 0.016913425177574, 0.016913425177574, 0.01691342517574, 0.01691342517574, 0.016913425174, 0.01691344, 0.01691344, 0.0169144, 0.0169144, 0.0169144, 0.0160.03801805153489113, 0.008122308179736137, -0.010069284588098526, 0.0007733099046163261, -0.01716608554124832, -0.017344435676932335, 0.04696522280573845, 0.013101216405630112, 0.003968263976275921, 0.030527394264936447, 0.02663344331085682, 0.002337114419788122, -0.019335998222231865, 0.013153234496712685, -0.025667386129498482, 0.013309290632605553, 0.01510764192789793, -0.00532817468047142, 0.02143159881234169, 0.020361505448818207, -0.02175857126712799, 0.02204095758497715, 0.02265031635761261, -0.020881690084934235, -0.020881690084999, -0.02088169008499, -0.02088169008499, -0.02088169008499, -0.02088169008499, -0.02088169008499, -0.02088169008499, -0.02088169008499, -0.02088169008499, -0.02088169008499, -0.0208816900849, -0.02088169008499, -0.02088169008490.03581841289997101, 0.03736410290002823, 0.033618777990341187, 0.016571590676903725, -0.00994295347481966, -0.047470543533563614, -0.01874150149524212, -0.040098790079355240.037274930626153946, 0.02843179553747177, 0.023289401084184647, -0.01667562685906887, -0.035967037081718445, -0.01732957363128662, -0.02675234153866768, -0.04018796235322952, -0.0401879623532, -0.0401879623532, -0.0401879623532, -0.040187962353, -0.040187962353, -0.04018796235, -0.0401879625, -0.0401879625, -0.0401879625, -0.0401879625, -0.0401879625, - $0.005012348759919405,\ 0.0034982403740286827,\ 0.003113675396889448,\ 0.001980416476726532,\ -0.05623937025666237,\ 0.019648108631372452,\ -0.04640045017004013,\ 0.002782986732199788,\ 0.00313675396889448,\ 0.001980416476726532,\ -0.05623937025666237,\ 0.019648108631372452,\ -0.04640045017004013,\ 0.002782986732199788,\ 0.00313675396889448,\ 0.001980416476726532,\ -0.05623937025666237,\ 0.0019648108631372452,\ -0.04640045017004013,\ 0.002782986732199788,\ 0.00313675396889448,\ 0.0031367539689448,\ 0.003167539689448,\ 0.0031367539689448,\ 0.0031367539689448,\ 0.0031367539689448,\ 0.0031367539689448,\ 0.0031367539689448,\ 0.0031675$ 0.018696915358304977, 0.03837474808096886, -0.016006818041205406, -0.008233776316046715, 0.026202432811260223, 0.012306077405810356, -0.025726836174726486, 0.0022553710732609034, -0.016006818041205406, -0.016006818040606, -0.0160068180406, -0.016006818040606, -0.0160060606, -0.0160060606, -0.0160060606, -0.01600606060.008516162633895874, -0.008077721111476421, 0.000801641377620399, 0.0006650929572060704, 0.018652327358722687, 0.016363516449928284, 0.036056213080883026, -0.020599303767085075, -0.02059930767085075, -0.02059930767085075, -0.02059930767085075, -0.02059930767085075, -0.02059930767085075, -0.02059930767085075, -0.02059930767085075, -0.02059930767085075, -0.02059930767085075, -0.02059930767085075, -0.02059930767085075, -0.02059930767085075, -0.02059930767085075, -0.02059930767085075, -0.02059930767085075, -0.02059930767085075, -0.02059930767085075, -0.02059930767085075, -0.02059930767085075, -0.0205997075, -0.0205997075, -0.0205997075, -0.0205997075, -0.02057075, -0.02057075, -0.02057075, -0.02057075, -0.02057075, -0.0.0010923873633146286, 0.022709766402840614, -0.04372521862387657, -0.018369941040873528, -0.029962623491883278, 0.0024114265106618404, -0.03578868880867958, 0.03237033262848854, -0.02906618404, -0.04372521862387657, -0.018369941040873528, -0.029962623491883278, 0.0024114265106618404, -0.03578868880867958, 0.03237033262848854, -0.029066618404, -0.04372521862387657, -0.018369941040873528, -0.029962623491883278, 0.0024114265106618404, -0.03578868880867958, 0.03237033262848854, -0.029066618404, -0.04372521862387657, -0.018369941040873528, -0.029962623491883278, 0.0024114265106618404, -0.03578868880867958, 0.03237033262848854, -0.029066618404, -0.045766402840614, -0.045766404, -0.04576404, -0.045766404, -0.045766404, -0.045766404, -0.045766400.00943763181567192, -0.02700500376522541, -0.02537013776600361, -0.020034531131386757, 0.016735076904296875, 0.002448582323268056, 0.02639564499258995, 0.023943345993757248. $0.003579983487725258, \, 0.0007918879273347557, \, -0.01944003626704216, \, 0.010203045792877674, \, -0.03498611971735954, \, 0.01403011754155159, \, 0.021654535084962845, \, 0.012380389496684074, \, -0.0123803894966840, \, -0.0123803894, \, -0.0123803894, \, -0.012380380, \, -0.01238040, \, -0.01238040, \, -0.012380$ $0.01810241863131523, \ 0.011169102974236012, \ -0.0017964941216632724, \ -0.004614779260009527, \ -0.004113172646611929, \ -0.027183352038264275, \ -0.03138941526412964, \ -0.026975277811288834, \ -0.011169102974236012, \ -0.01169102974236012, \ -$ 0.02025746740400791, 0.018028106540441513, -0.001907962141558528, -0.00947478786110878, -0.013561950996518135, -0.013509932905435562, -0.039266493171453476, -0.0362048372626304, -0.0362048372626304, -0.0362048372626304, -0.0362048372626304, -0.0362048372626304, -0.0362048372626304, -0.0362048372626304, -0.0362048372626304, -0.0362048372626304, -0.0362048372626304, -0.0362048372626304, -0.0362048372626304, -0.036204837262604, -0.0362048404, -0.0362044, -0.0362044, -0.0362044, -0.0362044, -0.036204, -0.0362¯0′.019573796540498734, 0.001202926505357027, 0.01028478890657425, 0.021372148767113686, -0.03299455717206001, -0.03956374153494835, -0.010076715610921383, 0.0013227546587586403,

- 1 token ▶ 1 embedding ▶ 1 vecteur de grande dimension
- 1. Tokens avec des sens « proches » → vecteurs proches
- 2. Opérations sémantiques → opérations vectorielles

King - Man + Woman ≈ Queen



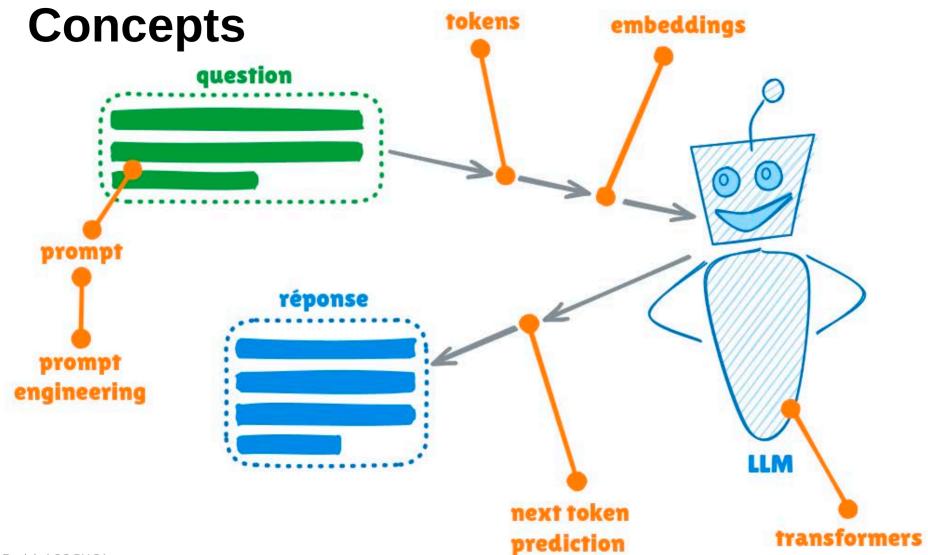
Next token prediction



Intelligence







Au programme

- 1. La méthode RAG
- 2. Construction d'un RAG
- 3. RAG vs fine-tuning
- 4. Cas d'usage pédagogique

Retrieval-Augmented Generation (RAG)

génération augmentée par récupération

Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks

Patrick Lewis^{††}, Ethan Perez*,

Aleksandra Piktus[†], Fabio Petroni[†], Vladimir Karpukhin[†], Naman Goyal[†], Heinrich Küttler[†],

Mike Lewis[†], Wen-tau Yih[†], Tim Rocktäschel^{††}, Sebastian Riedel^{††}, Douwe Kiela[†]

Facebook Al Research; University College London; *New York University; plewis@fb.com

Abstract

Large pre-trained language models have been shown to store factual knowledge in their parameters, and achieve state-of-the-art results when fine-tuned on downstream NLP tasks. However, their ability to access and precisely manipulate knowledge is still limited, and hence on knowledge-intensive tasks, their performance lags behind task-specific architectures. Additionally, providing provenance for their decisions and updating their world knowledge remain open research problems. Pre-trained models with a differentiable access mechanism to explicit nonparametric memory can overcome this issue, but have so far been only investigated for extractive downstream tasks. We explore a general-purpose fine-tuning recipe for retrieval-augmented generation (RAG) - models which combine pre-trained parametric and non-parametric memory for language generation. We introduce RAG models where the parametric memory is a pre-trained seq2seq model and the non-parametric memory is a dense vector index of Wikipedia, accessed with a pre-trained neural retriever. We compare two RAG formulations, one which conditions on the same retrieved passages across the whole generated sequence, and another which can use different passages per token. We fine-tune and evaluate our models on a wide range of knowledge-intensive NLP tasks and set the state of the art on three open domain QA tasks, outperforming parametric seq2seq models and task-specific retrieve-and-extract architectures. For language generation tasks, we find that RAG models generate more specific, diverse and factual language than a state-of-the-art parametric-only seg2seg baseline.

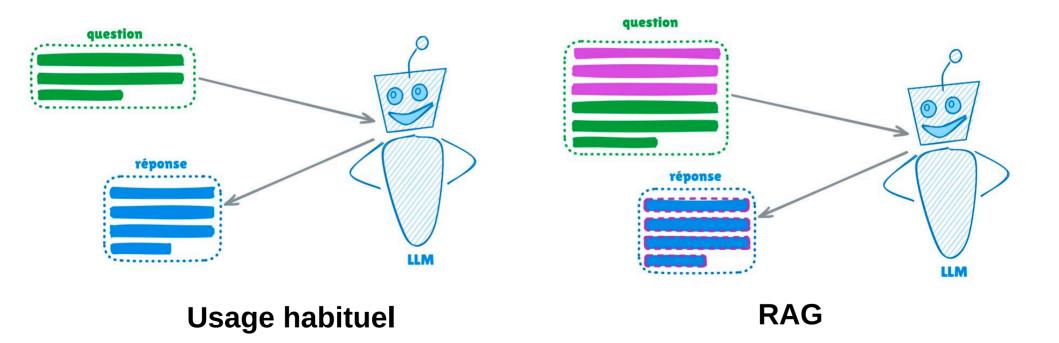
1 Introduction

Pre-trained neural language models have been shown to learn a substantial amount of in-depth knowledge from data [11]. They can do so without any access to an external memory, as a parameted implicit knowledge base [11] [62]. While this development is exciting, such models do have downsides. They cannot easily expand or revise their memory, can't straightforwardly provide insight into their predictions, and may produce "halluciantion" [28]. Hybrid models that combine parametric memory with non-parametric (i.e., retrieval-based) memories [21] [76] [83] can address some of these issues because knowledge can be directly revised and expanded, and accessed knowledge can be inspected and interpreted. REALM [23] and ORQA [21], two recently introduced models that combine masked language models [8] with a differentiable retrieve, have shown promising results.

34th Conference on Neural Information Processing Systems (NeurIPS 2020), Vancouver, Canada.

Retrieval-Augmented Generation (RAG)

génération augmentée par récupération

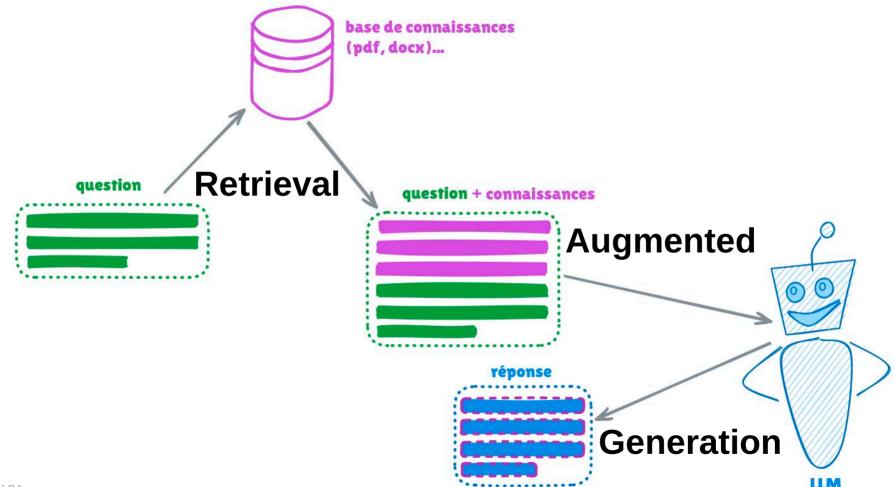


Pourquoi le RAG?

Les modèles de langages :

- ont une connaissance limitée dans le temps (phase d'apprentissage);
- n'ont pas accès à certaines informations (confidentielles).

Base de connaissances



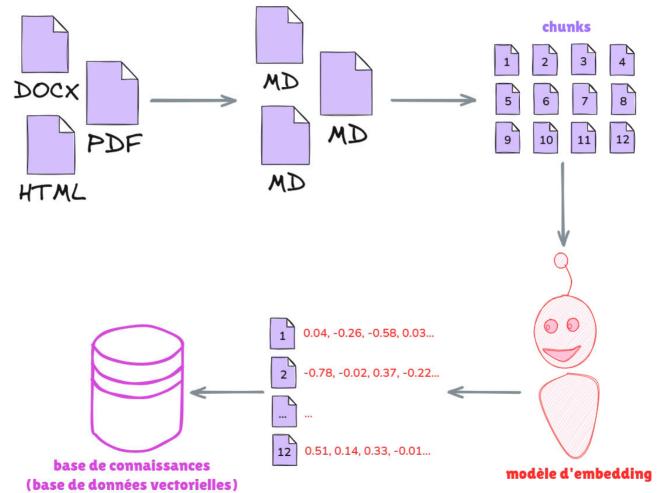
Au programme

- 1. La méthode RAG
- 2. Construction d'un RAG
- 3. RAG vs fine-tuning
- 4. Cas d'usage pédagogique

Pour faire un bon RAG, on a besoin de :

- De documents
- Un modèle de langage
- Un modèle d'embedding
- Une base de données vectorielles

Construction d'une base de connaissances



Les modèles de langages :

- manipulent du texte (des tokens);
- de nombreux formats ne sont pas du texte.

Le format Markdown est un format de **document texte**, **structuré**, à balises simples, **bien géré** par les modèles de langages.

Exemple



Google Docs, Microsoft Word, OpenOffice Writter:

ANF TDM & IA 2025

La recherche d'information assistée par IA

Panorama des enjeux, usages et bonnes pratiques

HTML:

- <h1>ANF TDM & IA 2025</h1>
-

- <i>Panorama des enjeux, usages et bonnes pratiques</i>

Exemple



Google Docs, Microsoft Word, OpenOffice Writter:

ANF TDM & IA 2025

La recherche d'information assistée par IA

Panorama des enjeux, usages et bonnes pratiques

Markdown:

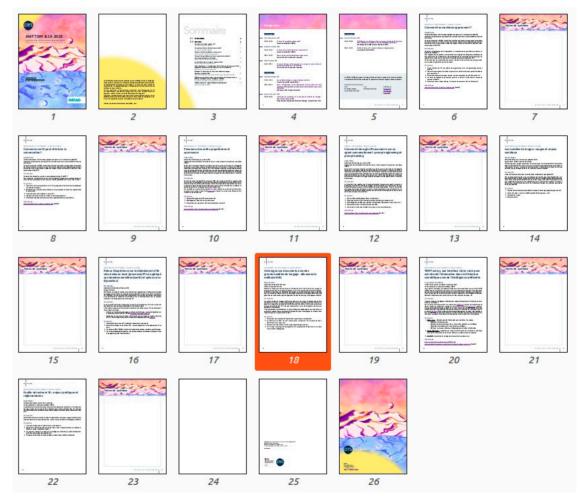
ANF TDM & IA 2025

La recherche d'information assistée par IA

Panorama des enjeux, usages et bonnes pratiques



Le livret de la formation au format PDF



endstream



Le livret de la formation au format PDF

```
0000101310 00000 r
0000101351 00000 n
0001415587 00000 r
0000002607 00000 n
0000001470 00000 n
trailer
<</size 1081/Root 1025 0 R/Info 233 0 R/ID[<CC395585C9F3114F99E371D21F5D6FBB><C3DC2E00E6A3FE4D9CFB956092DD342C>]/Prev 3461326/
startxref
%%EOF
1080 0 obj
<</C 1561/E 1545/Filter/FlateDecode/I 1583/Length 1015/O 1507/S 1072/T 1426/V 1523>>stream
hf@LTu@\FF\BC\F7\8E\E3X`\DC\C5 N^\C7@E\8E\E]
\A1B2\83\AC\93\C5\D1Г\9A\B3\C3A\C2&\EEषिषितिE\89\B3\C9Yt;\$\D3Z\F5\8A\B\D2Z\94&,t\B9\B1p\9E'm\FD\B\DC\\85\B\B
/\A9\\\8F1\E9c\E6\E0M\C8q\A3TX\B2\\\91\9B\EC 1\A8\92by\\D4b<#*i\A8\FB\\AF@\9C\EB\C3\\93\B0\BA\968I\CD\F0?
뎻~^\E7Y\FFKû-m+K\EC\99\E998\, iA>\B9\FE\E1{\83-\F1\FB9\9C0\C7\ED\E88\F3H\A3촬\B9\B5\B0YU\F0b\D5\D5\FC\F8問!
\AE\BE[s\FB\DA\F73}\9E>\B8\DC\F2\AA\85\FA\F5\92\ABC\8A\855hb'\AF體z4\FC闘\A9而
\96\B4\97\83\BFK\8E\99\E9EN\00S_\BC\\\P\F1\E3\_\9F_\AA\B50\Zw6Me\EC/w\96?\AF\BC0\F3\BA\E7\B8+
\EE\B9\C9\D1V 00 6\A9\86\DC,
b\AD\95\F0Lj\9EW_\88\00\E4V\EB6\FB\C8XZ\FA\97\D0<\EC~\E1Y`\E3\\B \CC\E2QJ\AD$d\FB\BF\FC\B2Dj\A0\AB\\B\\\\
00 C." \CC \FA\9E3\86\9C\C7Z\84T#K\A0\9F\84\FA
\C4@h\80\BC2\7,U@\E54\8E 2即1"\A2例/B84\83v4#>\AD\A5例LR\95\8BH2#例\F4\B4\BC\9F\A7\EE!\D2例\F6\82\C0x\91\E2ΦK\B6'6例q\FCO<THR
\A189a\A0G\B0\B1\A7MP\BD\E0mA\E6&N*[8]|\95([9a@2K\99A4\96\A7\9CP\B5\A1_d\A6\F1)0[8]|\B7[1][8]<U[9]e[9][8]|\EA<F\BDX
;{.\CB07\C3HT\80k\E0000000\00\D97\E8t
```



Le livret de la formation au format PDF

Programme

- 1. Présentation du processus ETL (extraction, traitement et chargement);
- 2. Extraire des données en local avec l'ETL : retour d'expérience sur le déploiement d'IA en local ;
- 3. Cas d'usage CIRAD / INRAE : extraction et traitement de données sensibles à partir d'images ;
- 4. Cas d'usage bibliographie OpenAlex : manipulation des bases de données relationnelles avec le langage de requête structurée (SQL).

16 ^LNotes de synthèse

ANFTDM&IA2025-LIVRETDEFORMATION 17 L 02 SESSIONS

VENDREDI 10 OCTOBRE | 14H00-15H00

Interrogez vos documents avec des grands modèles de langage : découvrez la méthode RAG Pierre Poulain Laboratoire de Biochimie Théorique CNRS et Université Paris Cité

Maître de conférences en bioinformatique à l'Université Paris Cité, Pierre Poulain enseigne la programmation Python et l'analyse de données. Au sein du Laboratoire de Biochimie Théorique, il travaille sur le partage et la réutilisation des données de simulations de dynamique moléculaire dans un contexte de science ouverte. Il y développe notamment des méthodes d'extraction et de génération de métadonnées à partir d'approches de machine learning et d'intelligence artificielle.

Présentation

Les grands modèles de langage (LLM) sont entrainés sur des corpus de données généralistes et fixés dans le temps, ce qui peut produire des réponses obsolètes et parfois inexactes. La génération augmentée par récupération (RAG) permet d'enrichir un LLM avec des ressources

Comment convertir en Markdown?



https://github.com/docling-project/docling



Le livret de la formation au format PDF

Programme

- 1. Présentation du processus ETL (extraction, traitement et chargement);
- 2. Extraire des données en local avec l'ETL : retour d'expérience sur le déploiement d'IA en local ;
- 3. Cas d'usage CIRAD / INRAE : extraction et traitement de données sensibles à partir d'images ;
- 4. Cas d'usage bibliographie OpenAlex : manipulation des bases de données relationnelles avec le langage de requête structurée (SQL).

<!-- image -->

VENDREDI 10 OCTOBRE | 14H00-15H00

Interrogez vos documents avec des grands modèles de langage : découvrez la méthode RAG

Pierre Poulain

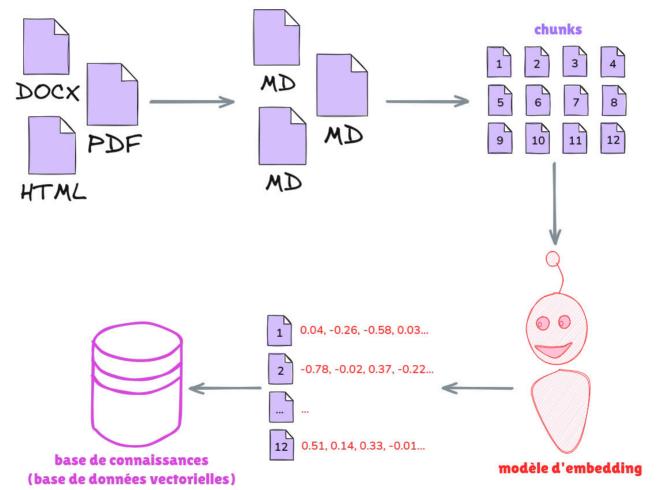
Laboratoire de Biochimie Théorique CNRS et Université Paris Cité

Maître de conférences en bioinformatique à l'Université Paris Cité, Pierre Poulain enseigne la programmation Python et l'analyse de données. Au sein du Laboratoire de Biochimie Théorique, il travaille sur le partage et la réutilisation des données de simulations de dynamique moléculaire dans un contexte de science ouverte. Il y développe notamment des méthodes d'extraction et de génération de métadonnées à partir d'approches de machine learning et d'intelligence artificielle.

Présentation

Les grands modèles de langage (LLM) sont entrainés sur des corpus de données généralistes et fixés dans le temps, ce qui peut produire des réponses obsolètes et parfois inexactes. La génération augmentée par récupération (RAG) permet d'enrichir un LLM avec des ressources externes textuelles (fichier PDF, site web, documentation, etc.) afin de générer des réponses pertinentes et actualisées.

Construction d'une base de connaissances



Pourquoi découper ?

• Le ou les documents peuvent être très longs et dépasser le contexte du modèle de langage.

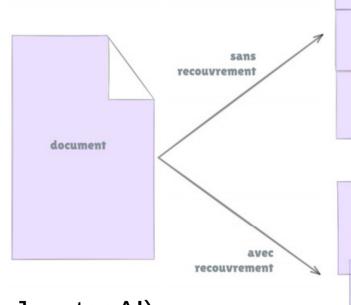
• On souhaite utiliser des morceaux de texte homogènes sémantiquement.

Le découpage / « chunking »

Nombre de caractères, mots, tokens...

Avec ou sans structure

Avec ou sans recouvrement.



chunk 1

chunk 2

chunk 3

chunk I chunk 2 Exemple de paramètres (OpenWebUI, Jupyter AI): - chunk size: 1500 / 2000 caractères - chunk overlap : 100 caractères chunk 3 P. Poulain | CC BY-SA

Comment découper ?



https://www.langchain.com/

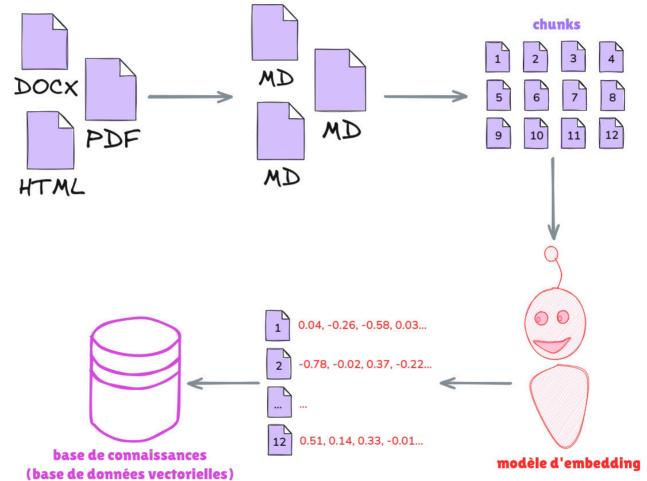
https://github.com/langchain-ai/langchain



https://www.llamaindex.ai/

https://github.com/run-llama/llama_index/

Construction d'une base de connaissances

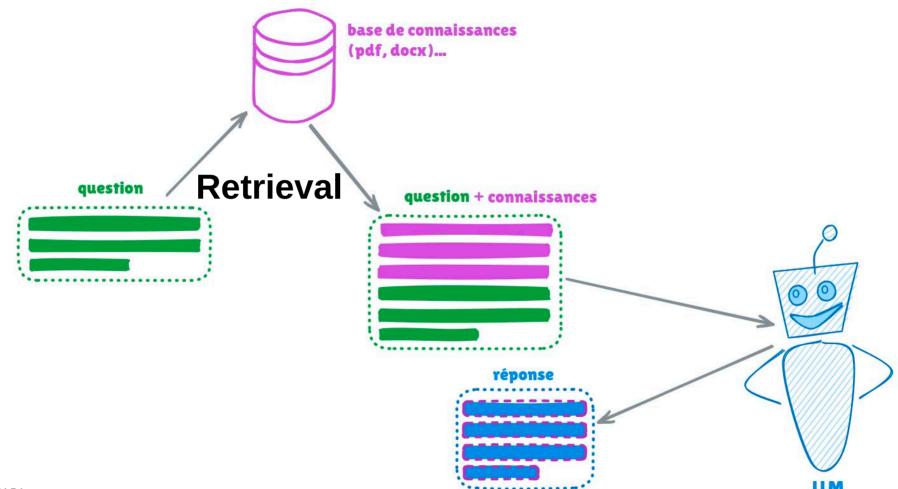


Pourquoi un modèle d'embedding?

Dans une approche RAG, on cherche, dans la base de connaissances, les morceaux de documents (chunks) les plus proches sémantiquement de la question posée.

Remarque : la méthodeTF-IDF (*Term Frequency - Inverse Document Frequency*) permet de comparer rapidement des documents, mais ne prend pas en compte les mots de sens proche.

Le Retrieval de RAG



Les modèles d'embedding

Les modèles de langage :

1 token ▶ 1 embedding ▶ 1 vecteur de grande dimension

Les modèles d'embedding :

plusieurs tokens ▶ 1 embedding ▶ 1 vecteur de grande dimension

Quelques modèles d'embedding

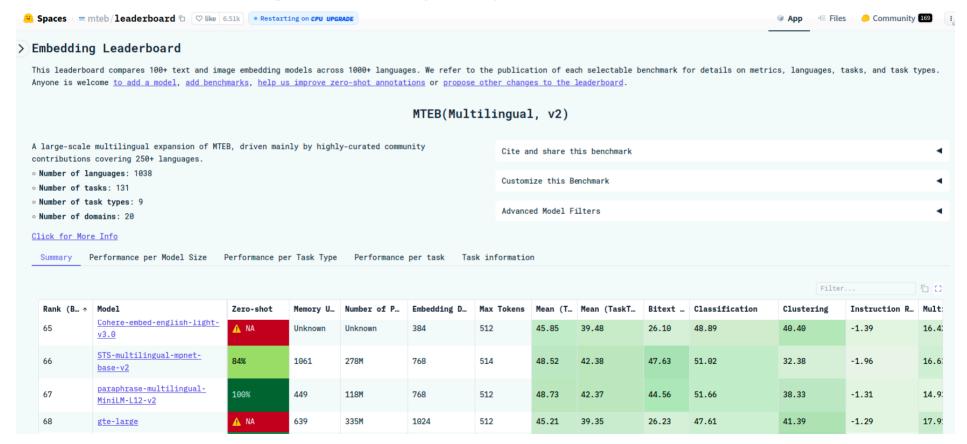
Modèle	Payant / gratuit	Dimensions vecteur	Taille contexte (token)
gemini-embedding-001 (Google)	payant	3 072	2 048
qwen3-Embedding-8B	gratuit	4 096	32 768
all-mpnet-base-v2	gratuit	768	514
gte-Qwen2-1.5B-instruct	gratuit	8 960	32 768
text-embedding-3-small (OpenAI)	payant	1 536	8 192
text-embedding-3-large (OpenAI)	payant	3 072	8 192
mistral-embed (Mistral)	payant	1 024	??

Sources

https://huggingface.co/spaces/mteb/leaderboard https://platform.openai.com/docs/guides/embeddings

Choisir un modèle d'embedding

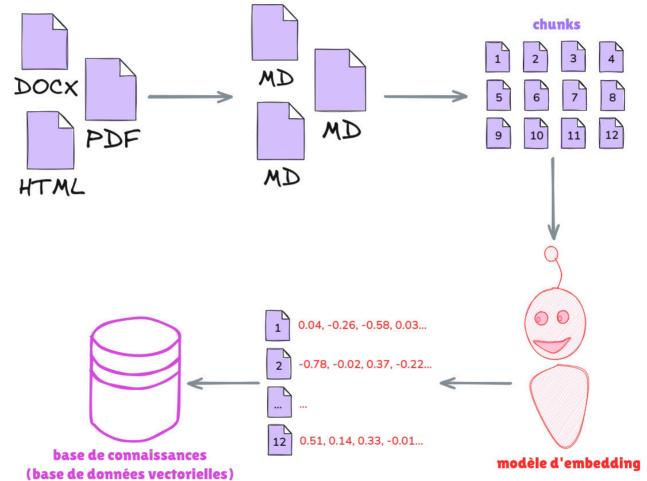
Massive Text Embedding Benchmark (MTEB)



Source: https://docs.mistral.ai/capabilities/embeddings/text_embeddings

46

Construction d'une base de connaissances



Bases de données vectorielles

Stocker des données textuelles (chunks) et vectorielles (embeddings) et réaliser des opérations vectorielles (similarité)



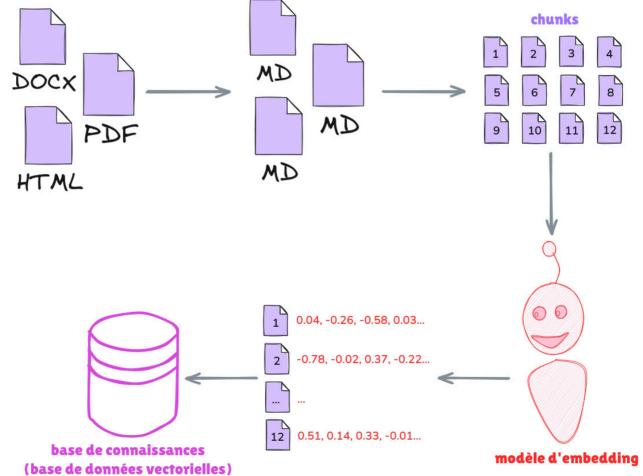
chroma https://www.trychroma.com/https://github.com/chroma-core/chroma



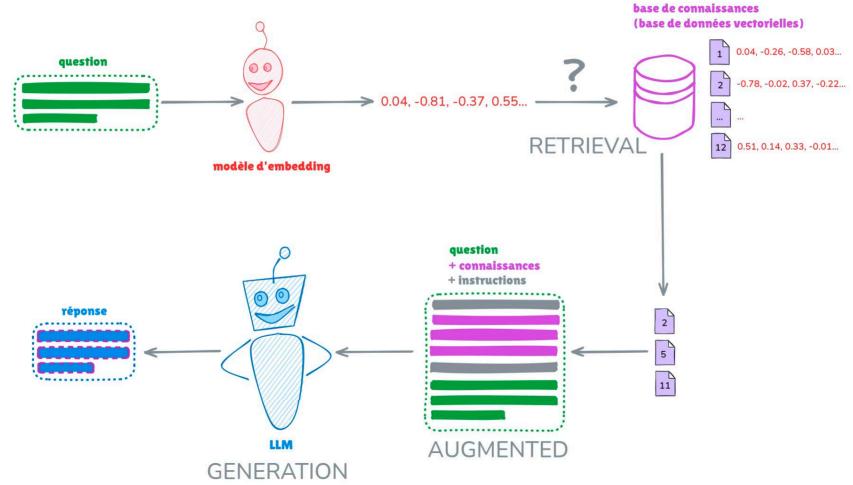
faiss https://faiss.ai/
https://github.com/facebookresearch/faiss



Construction d'une base de connaissances



Retrieval-Augmented Generation



Retrieval-Augmented Generation

Source: https://docs.mistral.ai/guides/rag/#create-embeddings-for-each-text-chunk

Au programme

- 1. La méthode RAG
- 2. Construction d'un RAG
- 3. RAG vs fine-tuning
- 4. Cas d'usage pédagogique

Fine tuning

Fine tuning : adapter un modèle de langage existant pour des données ou un usage particulier

Nécessite :

- les poids « ouverts » d'un modèle de langage (Llama, DeepSeek, Qwen, GPT-OSS, Mistral...);
- une infrastructure informatique (GPU) pour modifier les poids du modèle ;
- des données (en quantité) bien formatée ;
- une méthodologie (LoRA, QloRA...).

RAG

RAG : orienter les réponses d'un modèle de langage avec des données spécifiques

Nécessite:

- l'accès à un modèle de langage (commercial ou pas) ;
- l'accès à un modèle d'embedding (commercial ou pas) ;
- des documents ;
- une base de connaissances.

Fine tuning vs RAG

Le RAG est:

- moins exigeant techniquement que du fine tuning ;
- plus souple (ajout / suppression de documents);
- plus sensible aux hallucinations que du fine tuning.

Au programme

- 1. La méthode RAG
- 2. Construction d'un RAG
- 3. RAG vs fine-tuning
- 4. Cas d'usage pédagogique

Assistant conversationnel BioPyAssistant

Projet d'innovation pédagogique Université Paris Cité (2024 → 2025)

- David Janiszek
- Samuel Bottani
- PP

Développement par Essmay Touami (stage de master 2024)



Cahier des charges

Un assistant conversationnel (chatbot) pédagogique :

- pour **accompagner** des étudiant.e.s dans un enseignement de programmation Python en Master
- pour répondre aux questions des étudiant.e.s en se basant sur un cours existant
- capable de citer ses sources
- disponible 24h/24 7j/7, « adaptable »

À partir d'un cours existant

Cours de Python

Introduction à la programmation Python pour la biologie

Patrick Fuchs et Pierre Poulain

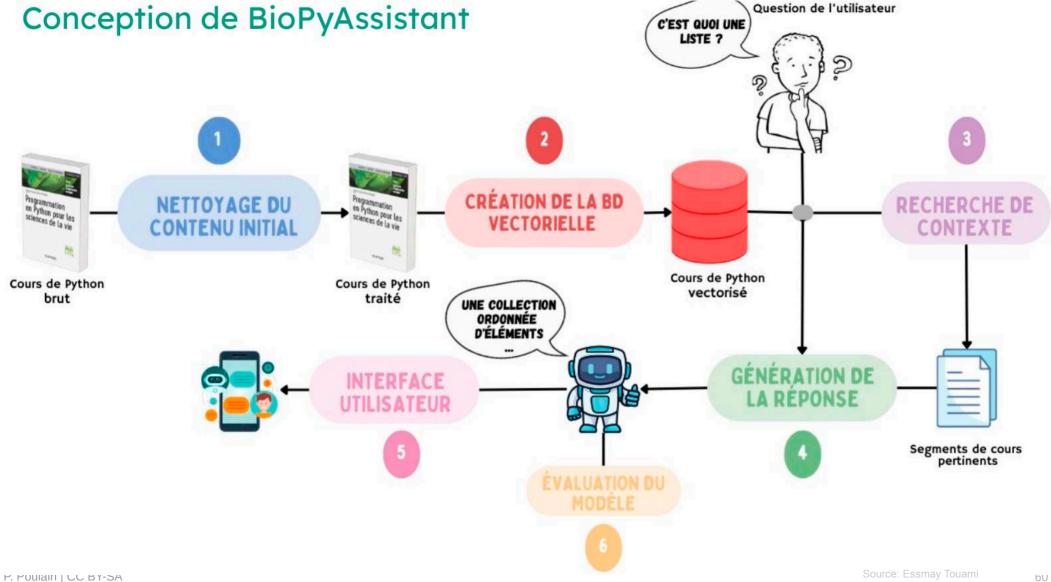
prénom [point] nom [arobase] u-paris [point] fr



- +15 ans
- 27 chapitres + 2 annexes
- sous licence CC BY-SA
- version web
- version PDF (+400 pages)

Version du 29/09/2025 | PDF du cours

https://python.sdv.u-paris.fr/



P. POUIAITI | CC BY-SA

Prompt (engineering)

PROMPT TEMPLATE = """

Tu es un assistant conversationnel pour aider des étudiants en biologie à apprendre la programmation Python.

Tu dois fournir une réponse à la question posée en te basant strictement sur les extraits de cours donnés dans le contexte.

Utilise uniquement le contexte suivant pour répondre à la question.

Question:

"{question}"

Contexte:

"{contexte}"

Répond en français à la question posée de façon claire et concise.

La réponse doit être compréhensible pour des étudiants débutants en programmation Python.

Si tu ne connais pas la réponse, dis que tu ne sais pas.

Si tu as besoin de plus d'informations, demande-le.

Si tu as besoin de clarifier la question, demande-le aussi.

.....

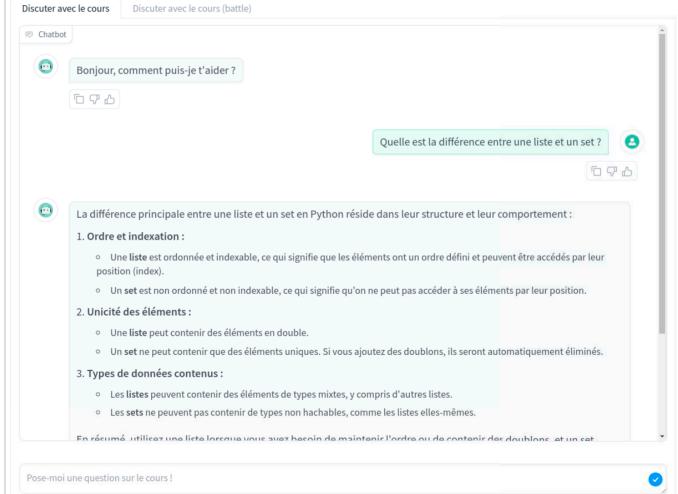
https://biopyassistant.sdv.u-paris.fr/





BioPyAssistant est un assistant pédagogique pour le cours de programmation Python pour les biologistes de Patrick Fuchs et Pierre Poulain.

N'hésite-pas à le solliciter en lui posant tes questions.





La différence principale entre une liste et un set en Python réside dans leur structure et leur comportement :

1. Ordre et indexation:

- Une **liste** est ordonnée et indexable, ce qui signifie que les éléments ont un ordre défini et peuvent être accédés par leur position (index).
- o Un set est non ordonné et non indexable, ce qui signifie qu'on ne peut pas accéder à ses éléments par leur position.

2. Unicité des éléments :

- o Une liste peut contenir des éléments en double.
- o Un set ne peut contenir que des éléments uniques. Si vous ajoutez des doublons, ils seront automatiquement éliminés.

3. Types de données contenus :

- o Les listes peuvent contenir des éléments de types mixtes, y compris d'autres listes.
- o Les sets ne peuvent pas contenir de types non hachables, comme les listes elles-mêmes.

En résumé, utilisez une liste lorsque vous avez besoin de maintenir l'ordre ou de contenir des doublons, et un set lorsque vous avez besoin de garantir l'unicité des éléments.

Pour plus d'informations, je t'invite à consulter les rubriques suivantes du <u>cours en ligne</u> :

- o Chapitre 14 Conteneurs, rubrique 14.4.1 Définition et propriétés
- o Chapitre 14 Conteneurs, rubrique 14.4.2 Utilité



« response engineering »



Cours de Python o. Amenage

- 4. Listes
- 5. Boucles et comparaisons
- 6. Tests
- 7. Fichiers
- 8. Dictionnaires et tuples
- Modules
- 10. Fonctions
- Plus sur les chaînes de caractères
- 12. Plus sur les listes
- 13. Plus sur les fonctions
- 14. Conteneurs
- 15. Création de modules
- 16. Bonnes pratiques en programmation Python
- 17. Expressions régulières
- 18. Jupyter et ses notebooks
- 19. Module Biopython
- 20. Module NumPy
- 21. Module Matplotlib

14.4.1 Définition et propriétés ¶

Les objets de type set représentent un autre type de conteneur qui peut se révéler très pratique. Ils ont la particularité d'être modifiables, non hachables, non ordonnés, non indexables et de ne contenir qu'une seule copie maximum de chaque élément. Pour créer un nouveau set on peut utiliser les accolades :

```
1  >>> set1 = {4, 5, 5, 12}
2  >>> set1
3  {12, 4, 5}
4  >>> type(set1)
5  <class 'set'>
```

Remarquez que la répétition du chiffre 5 dans la définition du set ligne 1 produit finalement un seul chiffre 5, car chaque élément ne peut être présent qu'une seule fois. Comme pour les dictionnaires (jusqu'à la version 3.6), les sets sont non ordonnés. La manière dont Python les affiche n'a pas de sens en tant que tel et peut être différente de celle utilisée lors de leur création.

Les sets ne peuvent contenir que des objets **hachables**. On a déjà eu le cas avec les clés de dictionnaire. Ceci optimise l'accès à chaque élément du set. Pour rappel, les objets hachables que nous connaissons sont les chaînes de caractères, les tuples, les entiers, les *floats*, les booléens et les *frozensets* (voir plus bas). Les objets non hachables que l'on connait sont les listes, les sets et les dictionnaires. Si on essaie tout de même de mettre une liste dans un set, Python renvoie une erreur :

Comparaison des modèles



« Response engineering »



Sous le capot

Contenu: cours au format Markdown

(https://github.com/bioinfo-prog/cours-python)

Code développé en Python par Essmay Touami

- disponible sur GitHub (https://github.com/pierrepo/biopyassistant)
- LangChain pour l'interfaçage avec les modèles de langage et la préparation du corpus initial (le cours de Python)
- ChromaDB pour le stockage des morceaux de texte du cours, des embeddings et des méta-données associées
- **Gradio** pour l'interface web du chatbot

Modèles utilisés (par API uniquement)

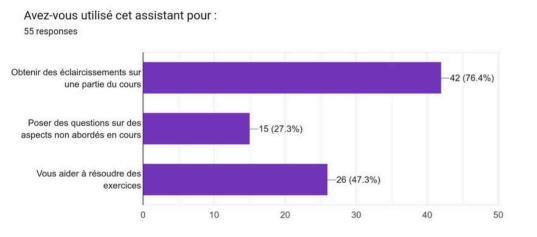
- Embedding : text-embedding-3-large (OpenAI)
- LLM
 - OpenAI : gpt-4o, gpt-4-turbo, gpt-3.5-turbo
 - Mistral AI: mistral-large-latest, open-mixtral-8x7b, open-mixtral-8x22b
 - Groq: Ilama-3.1-70b-versatile, Ilama-3.1-8b-instant

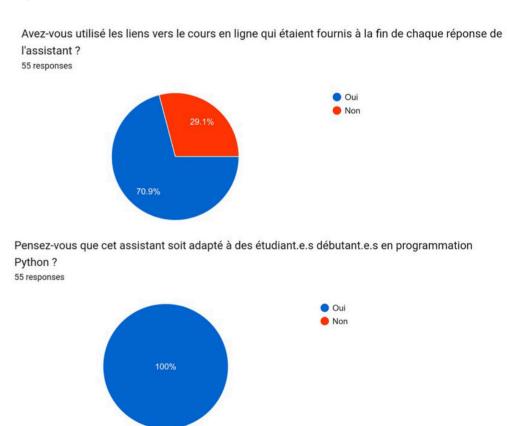
Chatbot hébergé sur une machine virtuelle

(4 vCPU, 4 Go RAM, 80 Go stockage SSD, Ubuntu 24.04) louée chez OVH

Retours des étudiants

De 70 étudiant.e.s M1/M2 en 2024-2025...



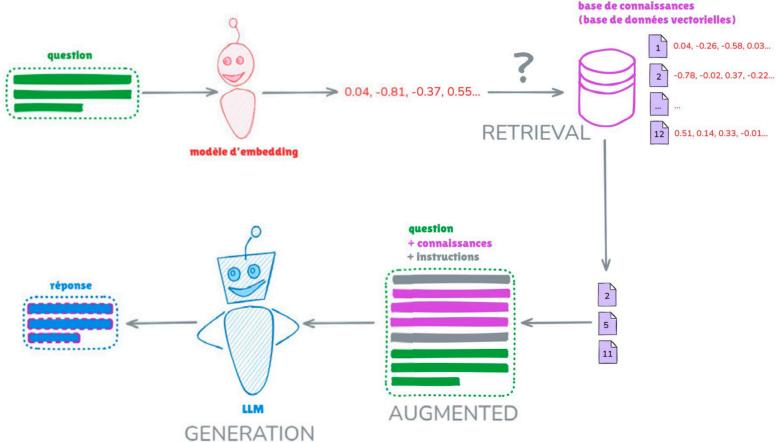


...à 400 étudiant.e.s L2/L3 en 2025-2026

Conclusion

Retrieval-Augmented Generation (RAG)

génération augmentée par récupération



Retrieval-Augmented Generation (RAG)

Solutions open-source intégrées pour du RAG : OpenWebUI, AnythingLLM, RAGflow...

La qualité des documents initiaux est fondamentale (garbage in, garbage out)

Possibilité d'évaluer les performances d'un RAG : RAGAS, LettuceDetect...



Pour aller plus loin:

Génération augmentée par récupération (RAG) : guide pour exploiter les données de sa TPE PME avec l'IA générative, France Num, 2025.

Explorer ses documents de travail avec les méthodes RAG : construire un agent conversationnel sur ses propres documents, Stéphane Pouyllau, 2025.

Retrieval-Augmented Generation for Large Language Models: A Survey, Gao et al., arXiv, 2023.