



0 Questions U must Ask Every Vendor with Gen-AI

Description

Whether they are selling learning technology, Vendors are a clever bunch of folks. I saw this firsthand at LTUK. Those who have content/courses, or even learning systems and technology are noted very visibly in their system, platform content/courses, Gen-AI on the floor. Lots of attendees tried. Some who even have it; pitched it at vendor sessions focused on the information at hand, as though listened intently. Their eyeballs fully Svengali was in front.

the person just scraped the internet for their While observing one session, it was clear all, generative AI, aka Gen-AI, is all the rage, presentation. Still, eyes were on them. After sfi and darn it, if anyone is going to add it, us in the L& and Training communities. First, we must ask the right questions. Which for many so much a question, as a tñsi though, Gen-AI in action. ,ffo eyes are glazed ,wow listening intently and seeing

that who knows could be the sfl I get it. If done right can be spun into a magic elixir new. And there lies the problem. Although sfsimehcladream of turning anything into gold. turn X into Gold. Good thing he had that Gravity Newton actually thought he could validation as a backup.

Gen-AI in their system. This, of course, made Filtered was one of those vendors promoting there, as fast as my feet could go, and asked to me want to know more. I zoomed over system. I believe it is ChatGPT. One of the wen learn more about Gen-AI in their okay, components was curation plus, the other too. sulp recall, but it had a tñac this system to generate a learning pathway and it sounded great. I can use Gen-AI in

Then I saw it, and started to ask questions. Oh, boy!

- How do you count for Do you tell people they might get hallucination)The answer BTW, was no.(
- It looks like a search engine function why this was with one of the plus coming by who not just have a search)They noted that there were vendors function and wanted to add it to their systems.(were impressed by said search
- a chatbot and the assets are the ones I pay for)Answer is yes(
- each header, the assets underneath are the With the AI pathway, while it generates is yes(
- what can I do on the back end to make sure the Getting back to the hallucinations the content they see is accurate this was under the assumption that is free content(, is well from the net.)Answer We pull down from the net)which are vetted by us, that you can access. If they have identified several sources that sources we can add them.(In other words, client wants to add additional has scrape the entire internet rather selected sites that Filtered will go to. You still access free , whatever that entailed that the Gen-AI content, just selected

vendor limit such options, or tweak it in such a Now you may ask yourself why would a called Tokens. ChatGPT uses them. GPT4 uses way, as above. It has to do with something is the people behind the ChatGPT, how those vendors, make money. OpenAI LLM and several others they have their method of revenue. Any open-source and so forth. Then it is added to your system, you select can be modified, tweaked, and content)as it relates to our industry(I get how learning technology, authoring tool, will say that it into the technical of a free deal for any vendor out there but I Gen-AI to whatever you are doing, then an Open there. If you, Company X want to add 100% free, others are fee-based(. And those that Source LLM is the route to go)some are fee-based can get very expensive, which is why as it relates to our industry are will try to reduce the cost in any way possible. vendors

you need to understand at a high level about The Five Questions will come, but first, questions you will ask any vendor who has Gen-Tokens, because it is part of one of those AI.

Tokens

recommend just doing a search around Gen-AI and To go further down into Tokens, I a Tokens. For the basics, it works like this anytime the AI processes text it becomes

can be words or chunks of tokens. Azure OpenAI, for example,

way to think of tokens is pieces of words. Another

Multiple sites use the example of "reg rub mah". It looks like one word, but with Tokens, it counts as three words. And you are charged for each token.

is non-English, the token numbers increase. If the language is Spanish ("How are you?" contains 0 tokens) for 10 chars. The higher token-to-char ratio can make it more expensive to implement the API for languages other than English.

OpenAI provides a Tokenizer, which allows you to see the number of tokens based on the text. It works is with a white space prior to the word. Please note that with tokens, the way it itself. Example:

I tried it out, with the phrase "who use proprietary models without stating it is built on an open-source model, is a22." The number of tokens whether

what is the basis, or method for processing and Okay, I added some words/characters but ?see

to find a simple explanation for explaining this I did a deep dive around the net seeking in a given request and choose Azure OpenAI. total number of tokens processed output, and request parameters. The quantity of depends on the length of your input, your response latency and throughput for the tokens being processed will also affect) Azure OpenAI. seldom

types in a word, phrase, or whatever, and hits. In other words, the moment the end user begins. When the end user follows up, more the button for a response, the cost for tokens the numbers will rapidly increase. For example, an OpenAI token costs. As you can see, is reportedly, 700K a day. That is 700 hundred ChatGPT is free. The cost to do this thousand dollars a day.

varies. That is to say, if a vendor uses any of LLMs the cost per OpenAI Pricing for tokens the same as OpenAI, which most folks think is the same. Nor is the price for Stability.ai with ChatGPT and their latest version, GPT-4. will recognize

challenge for those of curious nature. Here is another Finding the pricing is another **pricing** models, but vendors would use the token more. I should note they offer two price what I use. The interesting piece around some of them say the subscription, which is

is fee-based, even before actual usage. sġAnepOmodels is that training your data

Other LLM Models

continue to be added at a faster rate than mostThere are a lot out there and they anticipated.

but these are based on other Gen-AI models,There are open-source AI that is 100% free, Thus, you might find tġsi thus while free, itsay, a foundational LLM fee-based one.

using any one of these)a way to avoid payingvendors who went 100% free open source token fees(. This is just a shortlist.

- [Vicuna](#)
- [Stability.ai](#))They offer 100% freebies, and [DreamStudio](#) which is fee-based. Even only want to use the API(Their newest one is if you [StableLM](#).
- [GPT4All](#)
- [LAION](#)
- [Dolly-2](#)
- [Alpaca](#)site, you can reproduce it for less than \$600(But I slide it ġ)According to the into free here.
- [Koala](#))Built on LLaMA 13B(
- [Github](#)of a lot of 100% free open-source models(This link goes to ġ)A repository the Gen-AI ones.

models ġMost people will recognize OpenAIbut are likely unaware that they offer multiple edits images)Bing Create uses this(and Whisperincluding Dall-E2)which creates and and translates into hceepsġwhich is arecognition model that transcribes, identifies)OpenAI(ġsegaugnal multiple

is expected to announce the next version of Google [PaLM-2](#) for Bard with other LLM and not available elsewhere. To learn about whatmodels to follow. It is their own model PaLM was trained on)for those curious, read [here](#)(

ġcruos nepoġThese are all ġA Short ListLLM)Fee-Based(

- [LLaMA](#)there are numerous offshoots ġ)from Meta(ġFrom LLaMA, [view](#)
- [Azure OpenAI](#) For those using Azure ġ)Microsoft(
- [Amazon Titan](#) ġAvailable only to those who use AWS ġ)Amazon(They also offer [Bedrock](#),allows you to build from other LLM)for exampleyou have ġ which from that)foundational simply means it is theStablity.ai and want to build

above including OpenAI are the foundation of LLM which you need, and thus all of the foundational.

- **BERT** others from Google, for use with Google Cloud () and
- **Databricks** the AI Agents category which I believe is really) I place this more into coverage on how it works, etc for another time (. the power for our industry, and
- **Midjourney**) Image Creator (

want to pay for an LLM when freebies are out there? Why would someone

GPT-4 offers 1 trillion parameters. ChatGPT One of the key advantages is the parameters. Personally, it just comes down to what is yours next up with around 75 Billion. are in the business to make money. Thus, preference. Vendors who raise capital/funding you need to do a bit more than the freebies. by requiring it at some point,

the five questions, because one of them is) And What does this have to do with the most important(

1. Language Model do you use with your learning What Learning tool/content platform, offer two or ? . cte system/technology/authoring They may recommend, but I doubt at least right now that more, which experts in the LLM space is the way vendors are going.

model, which sounds amazing but it has to be A vendor might respond with a proprietary the unlimited funds to build 100% from scratch, built on something. They do not have a(the computing power to run gen-ai is massive high carbon footprint, b(to because due to heat, the cpus need water, and c(It just make any sense for the seed stay cool them to go that route.

they built it on some foundational model; made Thus the proprietary model to me, means would do, hence the open source part(, and said,, ada T changes to it) as anyone LLM are foundational(. If a rateirporp What you are seeking is the foundation) again, on Moodle, you might cringe, and yet there vendor said, hey we built are system originally system built on Moodle and then revamped it so are vendors who initially launched their knowing at the duoc much, and so on, that you even tell, when they debuted it. That said, you could explore via the net, what it is, least the name of the model, for those curious, the od how it works, and so forth. Maybe you care, but I would if a vendor said Dolly versus GPT-4. I would want the latter.

other day to show their system using Gen-AI went One vendor who reached out to me the

used. They told me that they used proprietary, beyond vague when I asked what LLM they on what that means and figured to corner them short, and big ones. Now, I have no idea I see the system) in early June (. Short and Big We are not talking about ? seno once for those super , yeH hummingbird feeders heredo you want a short one or a big one ? sdrimgnimmuh large

though no one will dare to ask what the is LLM, Yet, they are pushing this solution out as know what is a Learning Language Model. You need to tñod let alone that anyone would question here is to get the vendor to tell you know unless you are super curious, but the knowing it will aid you to follow up with what they used or are using. Because specifically which you must ask! a secondary question

what ?no 2. What data sets was your LLM trained You cannot skip this question and think Internet searches from 2022, ?rettam does it Because if they trained it on Wikipedia, and other learning models they selected that would be relevant. Take Kirkpatrick, 2021. Yes, it learns over time, but data sets ChatGPT for example. The data sets are from initially are relevant.

angle with their new authoring tool told me, One vendor who pushes the whole Gen-AI that were selected by experts they identified that the data sets included learning models And who are these ? seno or worked with. Which Unless they raised Donald ? strepxe the dead and pulled a Frankenstein with Gagne, not seeing an expert m Kirkpatrick from here. Oh, they also use a proprietary model.

the vendor uses a fee-based open source model, as in say GPT-4, will you see s fel 3. If in the coming years with your you the client ? tcartnoc an increase in fees

I said these models are not cheap due to the Well, if you think ? snekot Remember when then Bugs Bunny is right behind you, and the the vendor is going to eat that cost, is still a viable mode of transportation. Hindenburg

to cost a vendor. It is going to get pricey real just think how much money this is going then skirt it by going with a different layer quick, even if they use say part of the LLM and and to me, that is in pricing. Pricing has on top of it. Those costs have to go somewhere, can get \$60 per seat per user/year, at 2,500 always been an arbitrary model. If a vendor so) this is just an example(. I often tell the users, then they have no problem in doing story me \$48 per user/year) for a seat(, which I of a vendor who wanted to charge declined and price. The next day, it was down to around \$9 then declined on the next bucks. And they still made money on me.

charge it(, APIs) if they charge(, onboarding, Vendors can hide the price in setup) if they that contract you want to lock in that price for training, and so forth. Thus, when you sign increases of X percent that some vendors push as the length of your contract. No, price or living costs. None of that nonsense. Lock it though you have to eat the cost for inflation in. Because they will not eat that cost. That) referring to fee-based models, and I think the time to change the code and create your own think some 100% freebies because costing them money(. Oh, ditto on the authoring exclusive model with data sets is still getting for your e-learning that charges a fee tool, learning tech or whatever else you are

Gen-AI model they are using) within the company(ε. For those clients who already have a model for whatever data you want to bring in and? you can you API it into their Gen-AI I wonder how many vendors who have Gen-AI in This is going to be tricky here because tech or whatever has thought about it. But at their system or authoring tool or learning come up, and when it does the vendor is going to some point this question is going to have to know the answer.

because it is already starting to happen) not by I can see the inquiry around plugins, specific(. If client per se; but in general solutions not learning/training or e-learning other that exist for their LLM, which if they went I were a vendor I would start adding plugins the entire market, including 100% freebies, they 100% secret model that nobody else on though you can take your HRIS platform and ask are going to be in for a shock. The idea off, your this LLM, so don't even to add it to the something I would recommend doing. First to have an LLM. Secondly, it has to match what HRIS, HCM, ERP, Payroll, or whatever has vendor has) especially if the vendor is on AWS which the majority are, and they are the to have a solution like bedrock) using AWS as an using Titan(, OR the learning vendor has Y and it will work. Remember we are at a very example(whereas you can have X they can early stage here.

called data, as in what data is getting passed beyond serious Then there is this thing

around privacy, legality and who knows what which you would need to create concerns the vendor do with that data legal docs for those questions you have to ask what will be asked, responses, by who and so on.

its own LLM i.e. 100% free or fee-based (If your company is not at this stage of having this question. And if you have it, and think do it with X self then you can avoid many self respond back, uh, revert a hw learning not, because there are way too tiny steps here. Even by the end of 2024, I know out there, plus we are at extremely it goes, and if one vendor self would be leery. not forget that Gen-AI learns as trained their solution on data sets, to which experienced what a company did, when they was a then saw the LLM learned a new skill, that part or existed in the data sets) it on matching , DLOH language (, we should all be this had nothing to do with Q4, but just a surprise of what we are dealing with.

tech/authoring tool/etc deal with? no it can't. How will your system/learning

this as question 2 or 3, because it is they should have placed **BIG ELEPHANT** in the in their pitch. Hallucination exists in all rooms, and vendors are not so quick to respond who says, oh, we have a method to fix Gen-AI. You get rid of them today. So any vendor or a small number is selling you a magic elixir? devlo vni is Newton to have none

with any LLM and more importantly, Gen-AI. Hallucination is the biggest issue/problem And it occurs. NVidia just came out with a What they mean is fake or false information. solution they claim reduces hallucination,) [NeMo Guardrails](#) to reduce hallucinations but how well it will work over a period of time. it is truly unknown

When Fun Fact [Bard was tested by Google employees](#), one responded that it was a still went live with it. rail la cigolo htap Google

Bottom Line

out Gen-AI to me Here is my promise to you. going to be holding those vendors who push not just the ones listed above, but ones fire. I am going to ask those tough questions, the my findings. If that go further than that. report

keep you updated. I plan to create a massive As LLM and Gen-AI in general advance, I will you can access via the Learning Library on Gen-AI tools that you can use for directory and training. That is coming this summer. your learning

vendors who think, nobody is paying attention. And for

Worry,

Because we all are.

E-Learning 24/7

Category

1. authoring tools
2. Gen-AI
3. Generative AI
4. learning systems
5. learning technology
6. LLM
7. Tokens

Tags

1. e-learning
2. Gen-AI
3. Generative AI
4. Hallucination
5. Learning Language Model
6. learning systems
7. learning technology
8. LLM

Date Created

May 9, 2023

Author

diegoinstudiocity