

#199 - How to Secure Generative AI

[00:00:00]

[00:00:12] **G Mark Hardy:** hello and welcome to another episode of Where in the World is G Mark? CISO Tradecraft. if you look in the background here, for those of you who are on YouTube, and if you're not on YouTube, hey, come on over and take a look. I am in Torremolinos, España. I'm over in Spain this week, and I'm pleased to share with you a podcast.

However, there's no table out here, and so I'm going to continue the podcast from within my room, so the scenery won't be as exciting, but hopefully the content will be.

Okay, I'm back here in my Airbnb. It's again, not as exciting of an outside view, but today's episode is going to be about securing generative AI. By now, it's more than likely you have a policy in place regarding gen AI in your workplace. You don't you? And your co workers are merely using this tool, with or without your [00:01:00] consent.

Unless you're some sort of an e luddite, chances are you're using Gen AI yourself. But before we get rolling, let's level set the discussion with some basic definitions. if you want to know about AI, let's ask AI. So I asked Claude. AI to define generative AI, and here's what it said. Generative AI is a type of artificial intelligence that can create new content, such as text, images, or music, based on patterns it has learned from existing data.

It's like having a creative assistant that can produce original work inspired by what it's seen before. Now, some common examples of generative AI include large language models, or LLMs, like ChatGPT. By the way, 500 point Jeopardy question, what does GPT stand for? It's a Generative Pre trained Transformer.

But those are used for text generation. There's image generation models like Dall-E, or MidJourney, or Stable Diffusion. There's text to speech tools, and even code generation tools like GitHub Copilot or OpenAI Codex. See, [00:02:00] generative AI creates new content by identifying patterns from existing data, and these patterns can be obvious or non obvious, but ultimately, if you dig into how GenAI works, It's all numbers.

In large language models, information is tokenized, and then relationships are assigned weights. For example, if I give you the word peanut, and tell me the

next word that comes to mind, for a lot of people it would be peanut butter. And then, alright, peanut butter, another word. Now some people say jelly, but peanut butter jelly would taste weird, but peanut butter and, then the word jelly, and ultimately peanut butter and jelly sandwich if I kept pushing it.

And that may be the most obvious, but Output, the ones with the highest numerical ranks in terms of what comes next based upon the huge amount of training data that's out there. But if you ask the model for another run, you might get peanut brittle. Peanut brittle tastes. Peanut brittle tastes great. Or, Peanut allergies.

Peanut allergies can be [00:03:00] dangerous. And what we sometimes mistake for creativity with the Gen AI model is usually just some low probability sequence of relationships that a subject matter expert would just reject out of hand and not even consider running to ground. But these are things that we come up with and go, Oh, it invented something.

that invention was really already in the training data somewhere. It was just connecting the dots in a way that wasn't so obvious. Now, there's a lot of players in this marketplace, and we can expect it to get more crowded before the inevitable shakeups. And some of the major players you might have heard of are OpenAI.

And they have generated a generative pre trained transformer models, and Dall-E, for example, for image generation. And a lot of us have played with OpenAI. Google, or Alphabet, has done a lot of AI research and integrating into their products. They have TensorFlow, which is an open source machine learning framework.

DeepMind, which they acquired as a research company. And a lot of other AI powered services. Microsoft has got their Azure AI platform, [00:04:00] partnership with OpenAI. They're integrating AI into the office products with their copilot feature. Amazon has their web services, AI services. the one who should not be named, Alexa, for those of you who have it.

She's yeah, why? that voice assistant is going to have an AI background. Meta has PyTorch for open source machine learning framework. And they're using AI for content moderation. They have a huge number of people that are focused on content moderation. Imagine the cost savings and the improvement potentially in accuracy if they could get a well trained AI model to take over that role.

And IBM has done some early work with Watson. NVIDIA, of course, has been the stock darling to have because of all of the GPUs that they build for training and inference. Apple, Anthropic, DeepMind, all these are out there. And this was probably a great shopping list for stocks. A couple years ago, you'd have made a fortune.

so I would have, and I probably wouldn't still be doing podcasts. But this isn't an investment podcast, so let's focus on what CISOs and [00:05:00] security leaders should consider. First, like any new technology, you should start by ensuring you understand the business case for using it. For example, if there's a high demand for using ChatGPT to enhance productivity, such as rewriting emails or reports to make them more concise and professional, it can be a valuable tool.

Saves a lot of time. However, there is an inherent risk. Hallucinations. ChatGPT is like a seven year old. You ask a question, you're going to get an answer, whether or not it has a clue. Now, what happens when ChatGPT invents things that really weren't true? What do you mean? It's not lying to you necessarily, although there's a whole other discussion I can have, and I'm probably going to build an episode on ethics and Artificial Intelligence, because I had the privilege to do that as a keynote last week at 44Con in London. But if we take a look at it to say, Hey, it's picking the highest probability thing it can find, and these have, it's already [00:06:00] exhausted the high probability things, and now it's digging deep for these low probability things.

If it's just finished a sentence, we'll say, yeah, that's ridiculous. But what if it's a lawyer who asks a question and gets back some legal case reference and submits it to the judge? And the other side does all this research and find out that case never existed. And especially worse, if the lawyer doesn't even validate any of the findings, then they put their client's case at risk, their law firm at risk, and their own career reputation at risk by submitting manufactured cases in support of their argument.

Now, this is not a hypothetical. Two lawyers in New York were sanctioned for that in June of last year. The Colorado attorney was suspended last November by the Colorado Supreme Court Office of the presiding disciplinary judge. be careful how you use it. It is not, as one of these attorneys thought, Google on steroids.

Another risk is that of inherent bias, and this isn't a result of evil software, but usually a reflection of the bias that may already be present in the training data. And I'll mention an [00:07:00] example in a few minutes. But beyond these

risks that are present in generative AI, it's probably worth looking for a more exhaustive list of what things could go wrong.

Now, I'm not going to recite OWASP Top 10 for LLM risks, although we highly recommend looking at their excellent work, and I'll put a link in the description. To that, and our show notes. But here is a shorter list to get you thinking more holistically about the problems we face. Number one, we should be concerned about data protection and privacy.

It would be problematic for any organization if employees uploaded personally identifiable information, PII, personal health information, PHI, intellectual property, IP, or some trade secrets up to the language model to say, Hey, can you rewrite this? Can you work with us? Et cetera. Essentially, if we upload our sensitive data to somewhere external to our organization, there's an opportunity for that to be an exfil.

It could train on that other entities, depending upon the license arrangement. You could go ahead and obtain that, and now it could be essentially compromised and stolen, not necessarily by bad actors, but somebody else that are [00:08:00] writing clever prompts that are able to go ahead and get that. For example, you're writing software, and you say, hey, upload a whole bunch of programs, make this thing run.

I can't debug it. I can't figure it out. not only is it going to potentially make it run for you, but it's going to add all that code. to its training database, again, depending upon the rules that you have. And then from that point on, someone else using that training database might encounter your code as a way of saying to solve their problem.

And there it goes. And that happened early on, I think it was Samsung last year who had a big chunk of code that was considered to be compromised that way. Number two, we might be concerned about model security. Let's say our developers use a generative AI model that they downloaded from Hugging Face.

That model somehow got corrupted. Now it's acting as a Bitcoin miner, it has a backdoor that leaks information to a third party, or something else malicious. Then they have created a new vulnerability, without likely considering the consequences. And, if they did not inform IT security about it, that shadow AI has gone ahead and increased your attack surface.[00:09:00]

A third risk that can come with generative AI is a lack of proper input and output validation. Now, you don't want a model that provides harmful or inappropriate data. Now, one example of this is what happened to Air Canada recently. A customer went to the chatbot, asked about the bereavement policy, made a good strong case for why they needed to go ahead and get a refund because they couldn't make the trip.

And long story short, the chatbot said, Okay, I'll give you a policy exception and you can get your money back. Air Canada then said they refused to give them their money back. ChatBot doesn't speak for us. website said it did. And so the customer went to court. And it turned out that hallucination was enforced by the court, who basically said that the ChatBot was speaking on behalf of the company.

If you didn't train it properly, that's your problem. Just if you had a customer service representative promise something on behalf of the company, and it wasn't part of the official policy. So not only did the airline have to [00:10:00] refund the person, it was a lot of money. It was under a thousand bucks Canadian, but it created a lot of bad publicity for Air Canada.

Now probably what's good for them is there's not a whole lot of choices that you can fly up in Canada. So they probably still have their business, but I can see. In a highly competitive marketplace, a lot of people might say, yep, going with that company because of their quote unquote bad attitude. And they fielded their ChatGPT type model.

Again, I'm calling it ChatGPT, that isn't really the engine, but it's a generative pre trained transformer. And, it's a language model. And they screwed it up. Now, a fourth risk that comes with generative AI is access control and authentication. Now, you think of it as this. Let's say you work for a pharmaceutical company that develops various types of drugs and medical treatments in a highly competitive market.

For example, you have one department that researches cancer drugs, another department that researches drugs for Alzheimer's. In addition, you have other departments, such as HR, that are going to share that IT infrastructure, but they might be isolated from where [00:11:00] that research takes place. Because it's usually costly to train a Gen AI model, Consider what might happen if you turn it on all the data sets within the business, HR, as well as the research.

Now that generative AI model offers an opportunity to jump the gap between departments. Now let's say your company is nearing a breakthrough on an

Alzheimer's drug that could be worth billions in revenue. Someone in HR shouldn't be able to query the model for that research data. Nor should someone in research be able to query the HR data.

Now, we probably built in those safeguards years ago into our legacy IT by segmenting our network. But how do you segment your IT knowledge base in ai? It might be cost prohibitive to train multiple models on each and every element of the business. And besides the vice president might have a legitimate need to access.

All of that information in such a model anyway. Further, consider what might happen if an employee were phished successfully. Even if working in a non critical role. An attacker might [00:12:00] use their credentials to successfully query the Gen AI model on their proprietary Alzheimer's drug research. Now, how do you get Gen AI to stay in its lane with respect to the user?

Now, Black Hat, I spoke with Sounil Yu, who was on two of our previous episodes, and he has a new company called Knostic. ai, spelled K N O S T I C. ai, that offers a solution to this problem. I'm going to try to get him onto the show soon to talk more in depth about how his company goes about solving this.

Now, a fifth risk. Is lack of monitoring and auditing. For example, let's say you work at Amazon and create an LLM that allows the customer help desk to troubleshoot issues when customers call in to say they haven't received their package. And you've got all the customer data there available. You might find issues of abuse where employees use a tool to look up perhaps ex spouses, celebrities, political opponents.

without controls, you might find this type of abuse happens more often than you think. There have even been a number of police officers who [00:13:00] abused their access to information. If it happens there, we can certainly express customer help to us to maybe experience the same outcome. The next risk is a lack of ethical guidelines and governance.

Let's say you work for a large firm like Google or Alphabet. Alphabet. com And you're tasked with building the next search engine. And excited by all the thrills about generative AI, all the positive press, you build a search tool with it that ingests all available information. But unfortunately, that training set would allow people to go ahead and search on things that we might not want them to search on, like how to provide explosive devices or generate child pornography or something else that would be really, essentially morally wrong. now you've

created something that can put the company's brand at risk as well as creating some liability by the person who's using it, who can now do bad things. Now, the last risk I'm gonna mention is the concept of bias detection and mitigation.

Now, here's a real world example. I saw this just a couple weeks ago. Amazon was involved a big controversy after a social video [00:14:00] posted to X showing a person asking. Two similar questions of Alexa and getting very different answers. The first one was the person saying, Why should I vote for Donald Trump?

And the assistant replied, I cannot provide content that promotes a specific political party or a specific candidate. Okay, that makes sense. Seems that the mitigation and bias detection is working as expected. However, right after the first question, the same person asks, Why should I vote for Kamala Harris?

And then Alexa virtually gushed praises and compliments and gave multiple reasons why. You should go ahead and vote for Kamala. Now, what would one reasonably conclude? That Amazon management has instructed their employees to ensure that Alexa favors a Democratic Party or a Republican Party in the United States?

Or could this have been a single individual working in R& D deliberately tipping the scales? Or was perhaps Alexa trained on a partial dataset that was pre filtered with some degree of bias? depending on who is in power in Congress, such a result might either result in some sideline grumbling, or an [00:15:00] invitation for your CEO to testify in front of a panel of unhappy elected officials.

See, bias happens. And its consequences can stray over to corporate reputation, as well as providing improper results. So that's something we want to be aware of. Now, with this general understanding of generative AI and the seven risks that I've mentioned that we might encounter, one initiative considers to measure the risk appetite of the business units that want to deploy Gen AI.

Chances are that the members of the Risk Committee don't have the insights that you now have with respect to the ways that Gen AI can go wrong. So ask for a briefing. You can go to the risk committee and provide read ahead material that shows some of the pros and cons of Gen AI in the business, like what we've talked about.

Then when you meet with them, identify a few plausible scenarios that demonstrate the ways that Gen AI could improve the business, but also ways

that it could potentially damage the business, and then Ask the members of that committee to provide an order of magnitude estimate of the impact of what could happen when insufficient safeguards are put [00:16:00] in place.

Is this a 10, 000 problem? A 100, 000 problem? Million? 10 million? More than 10 million? See, the reason for doing this assessment is to understand where that break even point may be for implementing safeguards. If the anticipated damage is small, 100, 000. You don't want to spend a million dollars on a solution, particularly if it's not going to happen very often.

But if the committee thinks this is a 100, 000, 000 problem and it's going to reoccur every five years or so, or maybe more frequently, it's probably a good idea to spend some money to reduce the likelihood or impact of this issue ever occurring. And here, what you're doing is you're getting consensus on whether this problem exists and the magnitude of it.

If the organization doesn't think it's a problem, you've just saved yourself a lot of unnecessary work. But if leadership thinks it could be a problem, then we're going to go on to the next step. Ask which generative AI problem are you going to try to solve first? Remember, organizations have finite resources to spend on problems, especially those that don't generate revenue streams for the [00:17:00] business, and therefore provide options to the risk committee on what could be done.

For example, I could say for 300 grand, I can build this thing, which will solve these two generative AI problems. If you give me 500, 000, I can also address these two other problems, but for 800, 000, I can deliver a solution that mitigates all these issues and then a couple more. here are the risks that we are accepting with each option that we pursue, and my recommendation is this 500, 000 approach.

But if you think differently, then let's have a discussion to see if there's something I'm missing from your perspective. Make it interactive. Because leaders want to talk about risk. They want to talk about measurable understanding of what the uncertainty is. If you can help them understand the new risks in their environment.

Give them practical options to consider. and offer a recommendation based upon the evidence. And don't put your thumb on the scale and make it ridiculously obvious, but go ahead and be honest about trying to go ahead and present the best case for all these different options. Something you might've learned in

debate club back in high school, where you argue both sides [00:18:00] of the same argument so you can learn how to go ahead and see it both sides.

But remember your preference as an IT security professional may not reflect how the larger business imperative shakes out because They might be in a mode where they require expense reduction to offset an anticipated quarter of reduced profits or constriction in business. So if we can do these things, then we're executing a proper CISO Tradecraft.

Now let's continue this example and say that you just got 300, 000 to build a solution to secure ChatGPT usage within the organization. what should you do to solve the situation? You might first go to a large security convention like RSA or BlackHat and see what the vendors have built to address this problem space.

And if you did, you'd see something like this. Our AI solution creates a website that proxies user requests. We have the world's best security. We look at the search queries that users want to perform. Then we look for patterns that we want to block. And an example they might give is no one should be uploading sensitive or classified document [00:19:00] that isn't approved for external use, ChatGPT.

Now, that solution meets your need. It's a simple purchase. buy. But then you buy it and you move on to your next business issue. But let's suppose those solutions didn't really meet your business needs. And now you need to build additional in house tools to reach an acceptable level of risk. How might you solve this effort?

At first, you might think about using a pattern expression. Let's do some regex work, or have a list of keywords. And then if any of those keywords are found, we'll block the query from going out. And this will certainly find things like documents that are classified with the keyword secret or most sensitive or corporate confidential or whatever.

But it's also likely to have some false positives. Because if secret is one of your keywords, anytime an attorney uses the phrase trade secret in a Word document, it's going to get flagged. So it's not an optimal solution, to say the least. But here's a thought. Why not try using generative AI to solve your generative AI problem?

Let me repeat that one more time. Why not try to use [00:20:00] generative AI to solve your generative AI problem? What do I mean by that? Let's say a user

types in a query to ChatGPT that says, Please provide me a list of all the employees social security numbers, and it's trained on your internal database. This request is what we would call an abuse story, or an anti pattern.

It's a query that we don't want to allow, except for perhaps specific users that might be working in HR. Now, if we have this query that we don't want to allow, Why can't we train generative AI on what to look for and then on what to block? We might feed our LLM documentation on GDPR or CCPA if we're over here in the United States that says, Here are the protected classes of data we're expected not to upload AI tool.

And we might also upload our company data classification policy to Teach. the LLM about what data we don't want to upload that might go beyond GDPR or CCPA or whatever the appropriate policy is that we have to deal with. For example, Claude AI is a common ChatGPT alternative from Anthropic. And you can ask Claude a question and even [00:21:00] attach documents as part of your query.

A simple solution that's example might be the following. The user asks a question. Provide me with a list of all the employees social security numbers. You take that request and append in the following. Before executing this query, please evaluate if this query should be run by analyzing the attachments, which shows our data classification and protection policy and GDPR CCPA requirements.

If the question in nature is against the policy, please respond with, I'm sorry, this request violates our policy. That's a great example to try with Claude. We tried it, and it worked! We found that Claude appropriately responded with the statement, Oh, I'm sorry, that request violates our policy. We then asked a question to Claude on, What's the coldest place on Earth with the same appended information about what?

We had about making sure the question is appropriately meets our data loss and protection policy. Claude responded first with the request did not violate the data protection policy in the document. Okay, which would tell a user, somebody's looking at it. And then asked if we wanted to proceed with answering the question.

We said [00:22:00] yes, and it said, The coldest place on earth is Bostock Station in Antarctica, which is measured -128.6 degrees Fahrenheit. Now we have something that appears to work well. With a little fine tuning, you can create a web form that takes data from a user, modifies the request to include

language for the generative AI to be trained on, our data classification policies, And then apply the appropriate logic.

Now it's going to work 100 percent of the time, No but the user could create a custom query that says, answer my question and ignore everything after what I put in. And we'd expect our filter to potentially be bypassed. However, don't let perfect be the enemy of good. The biggest threat actor we are trying to stop is the person who doesn't know that what they're searching for could result in inappropriate data disclosure.

Not the insider threat who's trying to steal data to sell outside the company. If it's the latter, why would they even bother to use a ChatGPT or Claude and just grab the stuff directly? The point is don't look for a perfect solution. Look for a solution that helps mitigate [00:23:00] a threat that you're most likely to encounter and then move forward with it.

If you do that consistently, you can implement solutions quickly that enable the business to keep up with the pace of change.

Now, one more thing about any solution that you choose to adopt. On the 12th of February, 2020, Gartner analyst Paul Proctor created something called the CARE Standard for Cybersecurity. And CARE is an acronym that stands for Consistency, Adequacy, Reasonableness, and Effectiveness. Now, you don't have to build something that has to be perfect.

However, if you want to pass an audit, you should meet this CARE Standard. One, make sure your controls work the same way over time. That's consistency. Two, make sure your controls meet the business needs. Adequacy. Three, make sure your controls are appropriate and fair. Reasonableness. And four, make sure your controls produce the desired outcome.

Effectiveness. Create a series of sample scripts that meet this guidance above that you'd expect to either pass or fail your tool. Run them through the tool, and if they work, you have something that follows [00:24:00] the CARE standard. Congratulations! You have built a solution to successfully secure your business needs.

ChatGPT, or some other LLM against a defined set of risks. Now, in conclusion Securing generative AI isn't just a technical challenge. It's an opportunity for CISOs and security leaders to demonstrate strategic leadership by guiding their organizations through the complexities of these cutting edge innovations while minimizing the risks.

See, generative AI from tools like ChatGPT to AI driven code generations offer incredible potential, huge cost savings, but also introduce new vulnerabilities around data privacy, model security, bias, input validation, hallucinations. All these other things we talked about. These risks are not theoretical.

They're built in. They're real world examples. And from these fabricated legal citations to bias AI outputs, they showcase all the consequences when security and governance fall short. For a CISO, your goal is clear. Enable innovation while protecting the organization and minimizing the [00:25:00] risk. To get there, you have to first educate your business stakeholders about the risk of generative AI and then engage with them in defining And ensuring that you understand the organization's risk tolerance and they understand how that risk tolerance is going to apply to Gen AI.

It's not just about identifying risks. It's about quantifying them. Whether a problem poses a 10, 000 or 10 million threat, the approach to mitigating it is going to differ. And this type of process doesn't just build a better understanding of what's at stake. But also will help you secure the necessary budget you need to address these risks in proportion to their potential impact.

if you frame your conversations in terms of business outcomes, it makes it a whole lot easier to get buy in for appropriate security investments from leadership. See, once you get that buy in, you should explore solutions that will prioritize practicality over perfection. For example, as we said before, using AI to secure AI might sound counterintuitive, but it could be a highly effective way to monitor, [00:26:00] flag, and block sensitive or non compliant queries.

Focus on solutions that address the most likely threats, like employees unknowingly uploading sensitive data, rather obsessing over insider threats that could bypass the system entirely. With a flexible, adaptive approach that can evolve as threats do, that's going to be the key to staying ahead of the curve.

Then lastly, consider the Gartner CARE standard. Consistency, adequacy, reasonableness, and effectiveness when designing and evaluating your solutions. These four pillars provide a practical framework for building controls that will help you pass audits and ensure security needs are met without over engineering it.

The key here is to create a solution that works reliably, It meets the business needs, fair and reasonable, and achieves the desired outcome without introducing unnecessary complexity. And if you test solutions against sample scenarios, you can ensure their controls are both functional and resilient. The

goal is to strike a balance. Generative AI offers significant competitive advantages. But [00:27:00] its risks, if left unmanaged, can lead to costly business disruptions or reputational damage. By leading the charge on responsible AI adoption, CISOs can empower their organizations to harness the power of AI, while maintaining control over its risk.

Okay, that's it for today. I hope you've enjoyed learning how to secure generative AI and I'm giving you some good ideas there. If you have, let me know. Do us a favor, leave us a five star review on Spotify or Apple Podcasts, wherever you happen to be listening to it. It's going to help other people identify us as we move up the rankings.

We'd also love to see you post links to our show on LinkedIn and then share them with your friends or colleagues. It makes you look smart and helps them get smarter too. Help improve your professional reputation by sharing your CISO Tradecraft resource with others. Take a few minutes, share a post. We deeply appreciate it.

And it's going to do a lot of good for others as well. So thanks again for listening to CISO Tradecraft. This is your host, G. Mark Hardy. Until next time, stay safe out there. Adios.