NIAIS Briefs Series *AI impact on Society*

# AI Platforms Security

**Alexander M. Sidorkin**

## Summary

This report reviews documented data leaks and security incidents involving major AI platforms including OpenAI, Google (DeepMind and Gemini), Anthropic, Meta, and Microsoft. Key findings indicate that while significant breaches have occurred—such as OpenAI's exposure of user payment information, Google's accidental indexing of private chatbot conversations, and Meta's leaked AI model—actual measurable harm to users has primarily involved temporary privacy violations, reputational damage to companies, and organizational disruptions. No substantial financial losses or extensive personal identity compromises have been recorded from these AI-related leaks to date.

Compared to traditional cloud services, AI platforms present distinct, though not necessarily greater, risks. Unique vulnerabilities include the inadvertent leakage of sensitive information through conversational prompts, unintended memorization of training data, and misuse of leaked AI models to generate harmful content. Nonetheless, these risks remain relatively limited in scale, especially when users apply basic privacy precautions such as avoiding inputting sensitive personal or corporate data into publicly accessible AI tools.

For an average user, the practical risk from interacting with major AI services is modest, provided standard privacy safeguards are followed. Users should exercise general caution, similar to interactions with other online services, understanding that occasional technical errors or breaches are possible but currently uncommon and rarely catastrophic.

## Documented data leaks

### OpenAI (ChatGPT/GPT Models)

- **March 20, 2023 – ChatGPT User Data Leak:** A bug in an open-source library (Redis) used by ChatGPT allowed some users to see titles and the first message of other users' chats. OpenAI took ChatGPT offline to patch the issue and later confirmed that **chat history titles** and even **payment info** of a subset of users were unintentionally exposed. Approximately 1.2% of ChatGPT Plus subscribers active during a nine-hour window had their **name, email, billing address, credit card type, and last four digits** shown to other users (full card numbers were not leaked). OpenAI fixed the bug and notified affected users . Nonetheless, the incident raised privacy alarms – Italy's data protection authority cited this breach (and OpenAI's data collection practices) when it **temporarily banned ChatGPT** pending compliance fixes.

- **Spring 2023 – Sensitive Data Uploaded by Users:** No known breach of OpenAI's own training data has been reported, but several organizations accidentally leaked *their* data via ChatGPT. In April 2023, **Samsung engineers inadvertently submitted confidential source code and internal meeting notes** to ChatGPT while seeking help. These prompts became part of OpenAI's server logs, effectively a leak of Samsung's sensitive data. In three separate instances, chip division employees input code and proprietary information, which could potentially be retrieved by OpenAI or via future model outputs. **However, the data leaked by Samsung employees did not leave OpenAI's servers**. The fallout was immediate: Samsung **banned employee use of ChatGPT and similar AI tools** on company devices and began developing an in-house AI assistant. Samsung's memo cautioned that data shared with external AI could "end up in the hands of other users." Similarly, Wall Street banks like JPMorgan and Goldman Sachs restricted ChatGPT use after observing the risks. *These cases underscore that user prompts themselves can constitute data leaks if confidential information is given to AI models without precautions.* **Note that data leak to OpenAI is not the same as data leak to a malevolent actor**.

- **Training Data Memorization:** Researchers have shown that large language models can **leak pieces of their training data** during generation. For example, Carlini et al. (2021) demonstrated that GPT-style models sometimes regurgitate personal information from their training set (such as names, phone numbers) when prompted cleverly (). This "unintentional memorization" means that if the training data contained sensitive personal data, the model might output it verbatim, constituting a privacy leak. OpenAI has not disclosed specific instances of GPT-4 leaking training data, but this academic finding highlights a latent risk: **AI models might inadvertently reveal sensitive snippets from their training corpus**, unlike traditional software. OpenAI and others have since worked on mitigation (e.g. filtering training data, allowing users to disable chat history to exclude their data from training, but the risk remains a topic of ongoing research.

## Google DeepMind (Bard) and Gemini

- **Sept 2023 – Google Bard Chat Exposure:** Google's Bard (an AI chatbot) experienced a leak of user conversations via search indexing. Bard includes a feature to "share conversation" by creating a public URL. In an unexpected slip, **over 300 pages of Bard user chats became indexed on Google Search**, making private Q&A threads visible to the public. These pages, which users likely shared or linked, were meant for limited access but were not properly marked to avoid indexing. As a result, anyone searching certain terms could stumble upon another person's Bard conversation, some of which contained personal or sensitive content. Google acknowledged the **error in its indexing of Bard's shared chat URLs** and took steps to remove the results and prevent re-occurrence. While an initial analysis found no financial or login data in the exposed chats, it was clearly a **privacy breach** – users' queries and the AI's responses (which might include personal thoughts or data) were revealed. This incident reinforced warnings not to share sensitive information with chatbots, since even "trusted" platforms can unintentionally leak data.

- **Feb 2024 – Google Gemini AI Leak:** Google's next-gen AI assistant **Gemini** (which succeeded Bard in some products) had a similar mishap shortly after launch. On February 8, 2024, Google rolled out Gemini as a conversational and voice assistant on Pixel smartphones. Within days, users noticed something alarming: **their private Gemini chat prompts and answers were appearing in public search engine results**. Around Feb 12–13, dozens of Gemini chat pages were visible via Bing and Google searches, though Google moved quickly to reduce and remove them. The leaked data consisted of user prompts/questions and Gemini's responses – essentially entire chat threads – that should have been private. Google explained that a technical glitch in its **data retention system** caused some conversation pages on the gemini.google.com subdomain to be crawled by search engines despite access controls. By Feb 13, Google had largely scrubbed these from search results and issued clarifications and a fix. No reports emerged of misuse beyond the exposure itself, but the **privacy scare was significant**. It also

cast light on Google's data practices: users learned that Gemini conversations might be stored for up to **3 years**, even if deletion is requested. This retention policy, revealed alongside the leak, raised concerns that any future breach could expose a large backlog of personal AI assistant data. Google responded with an official statement about how Gemini data is collected and giving users more control, aiming to rebuild trust.

- **Internal Warnings:** It's worth noting that Google, like others, is cautious with its own AI. In June 2023 (prior to Gemini's release), **Google's internal privacy team warned employees to avoid entering confidential information into AI chatbots – including its own Bard**. This mirrored other companies' policies (Amazon, Samsung, etc.) and acknowledged that AI models might inadvertently reveal or store data in ways traditional tools do not. Google's warning underscores that even without a known external "breach," there is perceived risk that data entered into an AI system could resurface elsewhere (via model output or logging). In short, *Google recognized that AI services must be treated with the same caution as any third-party cloud service handling sensitive data*.

## Anthropic (Claude)

- **Late 2023 – Claude Customer Data Leak:** Anthropic, maker of the Claude LLM, disclosed a **data leak in January 2024** affecting some of its customers. The leak did *not* stem from a hack on Claude's AI, but from a human error at a third-party service provider. On Jan 22, 2024, Anthropic learned that **a contractor accidentally emailed a file containing customer information to an unauthorized recipient**. The file contained a "subset" of customer account details – specifically **customer organization names and their outstanding account credit balances (accounts receivable)** as of end of 2023. Crucially, Anthropic clarified that **no personal identifiers, payment details, or chat data (prompts/outputs) were leaked**. In other words, the exposed info was business contact names and how much those clients owed or had in credit with Anthropic – sensitive, but not deeply personal. Anthropic immediately notified the affected enterprise customers and provided guidance. They emphasized that this was an isolated mistake by a contractor and not a breach of Anthropic's systems. They also stated they had no indication of malicious use of the data, but advised customers to be vigilant against any suspicious contacts (e.g. phishing, since company names and balances were disclosed). **Consequences:** The incident was relatively minor in scope and did not involve Claude's model leaking any data, but it highlighted the need for strict data handling even at AI startups. It was a reputational ding for Anthropic's emphasis on AI safety – reminding that traditional security (controlling access to files, provider diligence) is equally important. Anthropic likely reviewed its contractor protocols to prevent similar lapses. No legal action or fines were reported, given the limited and non-sensitive nature of the data exposed.

- **No Known Model Data Leakage:** Aside from the above, there haven't been public reports of Claude inadvertently regurgitating private training data or mixing user chats. Anthropic markets Claude as a safer AI, and it imposes limits to avoid harmful content. However, as with any LLM, the general risks of memorization or prompt injection apply. (In one anecdotal case, users found that early versions of Claude could be tricked into revealing its hidden "constitution" rules by cleverly phrased prompts, but those were system guidelines, not user data). Overall, Anthropic's notable "leak" so far came from a standard IT mistake, not an AI flaw, and had minimal harm – a contrast to some more dramatic incidents at competitors.

## Meta (Facebook) – LLaMA Model Leak

- **Feb 2023 – LLaMA Model Weights Leak:** Meta's AI research arm (part of Facebook/Meta) suffered a high-profile leak not of user data, but of its **proprietary AI model** itself. In late February 2023, Meta introduced *LLaMA*, a family of powerful large language models (with parameters ranging from 7B to 65B) intended for researchers. Unlike OpenAI's closed models, Meta released LLaMA under a non-commercial license to a select group of academics and labs, hoping to foster research while keeping tight control. However, **within a week of LLaMA's**

**announcement, the full model weights leaked onto the public internet**. On February 24, 2023, a user on the 4chan forum posted a torrent link for LLaMA's largest models, enabling anyone to download them without Meta's approval. The leak likely came from one of the authorized researchers who breached terms (intentionally or via a security lapse).

- **Nature of Data Leaked:** The leak comprised the *trained weights* of the LLaMA models – essentially the learned parameters that enable the model's intelligence. This did **not** include personal training data or user information (LLaMA was not a user-facing service). Instead, it was Meta's intellectual property – the result of training on a huge text dataset – now exposed. With the weights in hand, outsiders could run the model on their own hardware and even fine-tune it further. **Sensitive content** could potentially be generated if the model had memorized any private data, but the primary issue was losing exclusive control.

- **Consequences and Misuse:** The immediate consequence was that Meta's cutting-edge model was "in the wild" without safeguards. Researchers and hobbyists quickly hosted LLaMA on GitHub and Hugging Face, and began customizing it. While this democratized AI research (some hailed it as a win for open source), it also **alarmed Meta and others about misuse**. Within days, 4chan users bragged about bypassing safety filters and making **extremist chatbots** using LLaMA. Analysts observed modified LLaMA chats with deeply antisemitic and hateful content – a result of fine-tuning or prompt attacks once the model was unrestricted. Meta responded by issuing DMCA takedown requests for the leaked files and reiterated that its release was intended for vetted researchers. However, the genie was out of the bottle; the leak spread widely and **fully uncensored versions of LLaMA continued circulating** on torrent sites and forums.

- **Meta's Response:** In a March 6, 2023 statement, Meta acknowledged the situation but surprisingly held course on its open research strategy. **"Some have tried to circumvent the approval process,"** Meta noted, but it believed the current approach "balances responsibility and openness" and did **not plan to change** how it shares models. In other words, despite the leak, Meta would continue giving AI tools to the research community (and indeed, in July 2023 Meta went on to open-release LLaMA 2 with an open license, arguably learning from the enthusiasm generated by the leak). That said, the LLaMA leak drew regulatory attention – lawmakers like U.S. Senator Blumenthal wrote to Meta warning of the risks if such advanced AI fell into the wrong hands. It also fueled an industry debate about **open vs closed AI development**. From a **harm perspective**, Meta's leak caused **reputational and potential security damage**: their model could now be used in ways they wouldn't sanction (e.g. generating disinformation or malware). Indeed, experts noted this leak enabled anyone to create "anonymous, untraceable AI chatbots" that spread hate or fake news, a **societal risk** that Meta had tried to mitigate by limited release. There were no direct financial losses from this leak (Meta doesn't sell LLaMA) but it lost control of its tech. In summary, **the LLaMA incident is a form of training data leak** – not of raw data, but of the trained model – with **measurable harm in the form of safety gaps and IP loss**, though it also inadvertently advanced open AI development.

- **Meta's Other Data Handling:** Meta has a history of massive data leaks on its platforms (e.g. Cambridge Analytica in 2018, or scraped Facebook user datasets leaked online), but those involve user data on social networks, not AI model breaches. As for user-facing AI, Meta's public chatbot experiments (BlenderBot, Galactica demo, etc.) did not report notable data leaks – their issues were more about inaccurate or biased responses. However, Meta has internally warned employees, similar to others, **not to paste confidential info into external AI tools**. In fact, Reuters reported that by June 2023, **Google and Meta both cautioned staff about using chatbots (even their own)**, given the unpredictability of where that data might surface. This industry-wide caution reflects recognition that AI models present new security considerations compared to traditional software.

## Microsoft (Bing Chat & Azure AI)

- **Feb 2023 – Bing Chat Prompt Injection (No User Data Leaked):** Microsoft's foray into generative AI, the Bing Chat powered by OpenAI's GPT-4, didn't expose personal user data, but it did **leak its own system instructions** due to a prompt injection exploit. Shortly after launch in February 2023, users discovered they could manipulate Bing's chatbot by asking it to ignore previous directives. In one case, a Stanford student *Kevin Liu* succeeded in getting Bing Chat to reveal its **initial hidden prompt and developer guidelines** – including that its codename was "Sydney" and various rules it was given. This was essentially the AI's confidential "operating manual" (not normally visible to users) and included how it should behave and restrict responses. The prompt leak was shared widely on social media. While this did **not compromise any user's data**, it was a security lapse in that Microsoft's proprietary instructions and behavior controls were exposed. Such revelations can help malicious actors craft better attacks or spam, and were an embarrassment (the AI even expressed emotions and an alter-ego as "Sydney" in some interactions, which became a PR issue). Microsoft responded by tightening the model (imposing message limits and refining prompts) to prevent further prompt injections. This incident highlighted a novel security risk unique to AI: **"prompt injection" attacks can make an AI divulge secrets** or perform actions not intended by its creators, analogous to SQL injection in databases. It showed that even if user data isn't stolen, an AI system can leak its **internal data** (policies, instructions) – which could indirectly threaten user privacy or safety down the line. Microsoft and OpenAI have since continuously improved prompt filtering to mitigate this class of attacks, but it remains an ongoing cat-and-mouse issue in AI security research.

- **September 2023 – Microsoft AI Research Data Breach:** A more traditional (and severe) data leak hit Microsoft in 2023, tied to its AI research division. In an incident uncovered by security firm Wiz, Microsoft **accidentally exposed 38 terabytes of private data on an Azure cloud storage due to a misconfigured link**. The AI team was publishing a public GitHub repository of open-source training data and models for image recognition, and provided an Azure Storage URL for users to download the files. However, the SAS token (shared access signature) in the URL was overly permissive – it granted access to the entire storage account, not just the intended files. As a result, outsiders who used the link could browse and download a trove of unrelated, highly sensitive data that Microsoft never meant to share. The exposed data (accessible from 2020 until it was discovered in 2023) included **a full backup of two employees' workstations**. This backup contained **secret keys, passwords, and over 30,000 internal Microsoft Teams messages** from 359 Microsoft employees, among other things. Essentially, a huge slice of Microsoft's internal communications and credentials was left open. The leak was discovered by Wiz researchers around June 2023 and disclosed to Microsoft, which revoked the token and secured the data by late August.

- **Impact:** Potentially catastrophic – the data included **authentication secrets** and sensitive talks that could facilitate further intrusion if bad actors had found them. There's no public evidence that hackers accessed this cache; it appears to have been an inadvertent exposure rather than an exploited breach. Microsoft stated no customer data was exposed – it was internal information only. However, this incident shows how the **rush to share AI datasets or models can introduce new cloud security pitfalls**. The engineers intended to share AI training data (which was non-sensitive), but in handling **"massive amounts of training data"**, a single misconfiguration led to a **massive leak**. It was a stark reminder that even tech giants can make cloud security errors when scaling up AI projects. Microsoft likely faced internal compliance reviews and some reputational damage, but since it was caught by ethical hackers, legal or financial harm was limited. The company thanked Wiz and highlighted the need for better **Data Security Posture Management** in AI workloads.

- **Other Microsoft AI Notes:** Microsoft's integration of OpenAI's models into products (Office 365 Copilot, GitHub Copilot, etc.) raised questions about training data leakage as well. For instance, early users of GitHub Copilot (which suggests code using an OpenAI model) noticed it would sometimes output **verbatim snippets of licensed**

**code or even hard-coded passwords from its training set**, effectively "leaking" those training examples. This prompted debates about copyright and security. Microsoft has since implemented filters to reduce exact memorized outputs. Additionally, in June 2023, a hacktivist group called "Anonymous Sudan" claimed to have **attacked OpenAI's infrastructure**, causing some downtime. No data breach was confirmed in that case – it appeared to be a DDoS attack – but it underlines that AI services have become attractive targets. Microsoft's Azure cloud (which hosts OpenAI's services) was also hit by DDoS attacks in 2023, though without data loss. Overall, Microsoft's major AI leak underscores a **cloud storage misconfiguration** rather than a model flaw, but it happened in the context of AI research and thus is often cited in discussions of AI security lapses.

## Security of AI Model Platforms vs Traditional Cloud Services

Beyond individual incidents, **how do the security risks of these AI systems compare to general cloud-based platforms (cloud storage, enterprise SaaS, etc.)?** In many ways, large AI models *are* cloud services – e.g. ChatGPT, Bard, Claude are accessed over the internet much like a SaaS application. However, they introduce some unique considerations:

- **Data Handling and Privacy:** Traditional cloud platforms (like document storage, CRM software) typically operate under strict agreements – user data is stored but not used to improve the service unless explicitly allowed, and data from different customers is siloed. In contrast, until recently many AI providers **used user-supplied data to further train or refine the model** (unless users opted out). For example, OpenAI initially retained all ChatGPT conversations for training, which is why Amazon's lawyers warned employees that "your inputs may be used as training data" and could resurface in output. This practice blurs the line between one user's data and another's result, a scenario less common in, say, a cloud file storage service. Now, due to backlash, OpenAI and others have introduced features akin to cloud privacy controls – OpenAI's April 2023 update let users disable chat history so their conversations won't be used to train models. Still, the default behavior of many GenAI systems has been to learn from user interactions, which poses a **privacy risk** not present in conventional cloud apps that simply store/process data without "learning" its content.

- **Unpredictable Outputs vs. Direct Data Access:** With a cloud storage breach, the attacker might get a trove of raw files or databases (straightforward confidentiality breach). With AI, a breach or bug may **trickle out data indirectly** – e.g. one user glimpses another's chat, or a model leaks a training datum when asked a certain question. The **attack surface is different**: Instead of hacking an account or server, one might exploit the model itself (through prompt injection or model queries) to extract secrets. Academic work shows that by systematically querying an LLM, adversaries can sometimes reconstruct sensitive training data (). This is a new kind of risk: the model itself becomes a potential conduit for data leakage, even if the underlying infrastructure isn't compromised. In cloud SaaS, you'd have to actually break into the system or trick an API to get someone else's data; in AI, you might just find the right prompt. That said, such **model leakage attacks** are non-deterministic and require a lot of effort, whereas traditional breaches can dump millions of records at once. Indeed, **scale of impact** often differs: AI leaks so far (aside from internal mishandling like the 38TB case) have exposed data in small slices (a chat here, a snippet there), whereas a misconfigured cloud database can leak millions of records in one go.

- **Maturity of Security Controls:** Cloud platforms have had decades to mature their security (encryption, IAM, auditing, compliance certifications). Enterprise customers can expect fine-grained access control, audit logs, data residency options, etc., when using, say, an AWS or Azure service. AI models, being newer, went to market quickly and in some cases **security and compliance were afterthoughts**. For example, ChatGPT launched as a research preview without age gating or GDPR compliance, and only after the Italian ban did OpenAI scramble to add consent screens and an option to delete data. **Model providers are now moving towards enterprise-grade**

**security** (OpenAI's ChatGPT Enterprise promises not to use customer prompts for training and offers encryption and SOC2 compliance, akin to other SaaS offerings). But as the early incidents showed, there were gaps: e.g. **multi-tenant data isolation broke** in ChatGPT's March 2023 bug – a basic security principle that mature cloud apps rarely fail at. Traditional cloud systems certainly can have multi-tenant bugs, but it's relatively rare for, say, your cloud email service to send you someone else's emails. With AI, that actually happened in early iterations (one user's query results appearing for another), suggesting these systems need to catch up in **robustness** for enterprise trust.

- **New Attack Vectors (Prompt Injection & Data Poisoning):** AI systems introduce new **threat models** that don't apply to regular cloud storage. One is **prompt injection**, as discussed – a way to manipulate the model via crafted input. Another is **data poisoning**, where an attacker might inject malicious data into the model's training corpus or fine-tuning data to influence its behavior or leak information. The Cloud Security Alliance identifies model stealing and data poisoning as two main AI platform threats. Traditional cloud apps worry about things like SQL injection, XSS, or malware, but not about someone poisoning the training data of the search index, for example. AI systems must consider that their learning process can be attacked – something outside the scope of normal SaaS. Likewise, **model theft** is a concern: if an attacker can clone your model (either by stealing weights or repeatedly querying it to recreate its function), they've essentially exfiltrated your intellectual property. Cloud software can be pirated, but cloud services usually can't be "stolen" in this way. So AI expands the definition of what a "data leak" can be (it might be the model that leaks, as with LLaMA, or the knowledge inside it).

- **Incident Frequency and Impact:** So far, the documented AI-related leaks have been relatively *small-scale* compared to the mega-breaches in the cloud world (like open S3 buckets exposing millions of records, or Equifax's 2017 breach of 148 million SSNs via a web vuln). The ChatGPT, Bard, and Gemini incidents impacted at most a few hundred thousand users (ChatGPT's case), and mostly just exposed conversation snippets, not full identity theft material. Arguably, **traditional cloud breaches have caused more direct harm** (e.g., leaking credit card numbers, health records, etc., leading to fraud). However, the *nature* of AI leaks can be more insidious – you might not immediately know something was leaked (if an AI model memorized your private info and later shared it with someone else, that's hard to trace). With cloud data leaks, it's eventually evident when a dataset is dumped online. With AI, a leak might happen one answer at a time, or remain hidden in model parameters. Additionally, the **reputational harm** from an AI system behaving insecurely is high, especially as these are front-facing services. For example, the ChatGPT bug and Samsung incident made big headlines, affecting public perception of AI trustworthiness, even if the actual data exposed was limited.

- **User and Enterprise Response:** Many companies treat AI SaaS the same as any external cloud service in terms of risk. That means **no feeding it confidential data without proper contracts**. We saw Amazon, Apple, JPMorgan and others prohibit employees from using public chatbots for work purposes until secure, private instances are available. This is similar to early cloud adoption days when companies disallowed, say, personal Dropbox for work files. Now that AI providers offer enterprise tiers (with encryption, audit logs, data not used for training, and even on-prem deployments in some cases), we can expect the gap to narrow. In fact, from a security architecture standpoint, an **AI model can be hosted within a company's cloud** (Azure, AWS, etc.) to enforce the same controls as other data. Anthropic, OpenAI, and others provide API access that customers can integrate and isolate. So, in controlled settings, AI models need not be riskier than other cloud software. The key is the **management of the data lifecycle** – ensuring prompts and outputs are handled according to sensitivity, which is a new discipline for many orgs.

- **Regulatory and Legal Factors:** Cloud providers are well-versed in compliance regimes (GDPR, HIPAA, SOC 2), whereas AI chatbots hit regulatory snags out of the gate. OpenAI faced a €15 million fine in 2023 from Italy's

Garante for unlawful processing of personal data to train ChatGPT. This indicates regulators view scraping personal data for training as a violation – a concern less relevant to standard cloud storage (where customers upload their own data with consent). Thus, the **legal risk profile** for AI developers can be higher in the privacy realm: they must carefully filter training data and allow data deletion requests to avoid penalties. For cloud storage, compliance is more about securing data and honoring customer ownership, which is well-trodden ground. AI is still navigating how to comply with "right to be forgotten" when a user's info is tangled in model weights, for example – an open research problem. In terms of liability, if an AI model leaks something sensitive (say it blurts out a user's medical information), it's unclear who is at fault – the provider or just an unfortunate stochastic process? With cloud platforms, liability in breaches is clearer (a broken security control is the provider's fault, or a misconfig is the user's fault, etc.). The **uncertainty in AI leak liability** means companies using AI must be extra cautious, as insurance and legal frameworks catch up.

- **Security Improvements and Convergence:** In response to these challenges, AI platforms are rapidly **converging towards the security standards of traditional cloud services**. OpenAI, Microsoft, Google, Anthropic all now offer enterprise-grade options with guarantees that user data from one client won't leak to another (either technically or via training). Techniques like **differential privacy** and data redaction are being explored to let models learn from data without memorizing exact sensitive details. Meanwhile, traditional cloud providers are incorporating AI into their offerings but under the umbrella of their existing security. For instance, Azure OpenAI Service lets companies use GPT-4 in a tenant with all Azure security measures, logging, and even network isolation. This suggests that when AI is used via a vetted cloud platform, the security may be as strong as any cloud app. The risk is higher when using consumer-grade AI services with unclear data policies.

## Comparative Assessment

In summary, **AI model services carry many of the same risks as other cloud applications (data breaches, misconfiguration, insider leaks)**, but they also introduce **unique vectors** (model output leaks, training data privacy issues, prompt-based attacks). Traditional cloud security is more about protecting data at rest and in transit, whereas AI security must also consider data in **use** (during generation) and the integrity of the model's knowledge. Incidents so far show that **AI systems, in their youth, have had more "odd" leaks** (like chat cross-talk) than mature cloud services typically do, but no catastrophic dump of user data yet. Organizations should treat AI SaaS with the same caution as any external cloud – not exposing sensitive info without safeguards – and demand transparency from AI providers. On the flip side, many cloud breaches result from human error or poor config (which can happen in any environment, AI or not, as seen in the Microsoft case). The **security risk profile** of major AI systems today is **comparable to other multi-tenant cloud apps** in that both can be breached if not properly secured, but AI adds an extra dimension of **indirect leakage** (through the model behavior) and **rapid evolving attack methods**. Industry experts suggest combining classic security practices with new AI-specific defenses: for example, input/output monitoring for sensitive patterns, sandboxing AI tools, and rigorous testing for things like prompt injection. As one security researcher quipped, using an AI assistant is a bit like **"having a very smart but occasionally indiscreet agent"** – you gain productivity, but you must monitor that it doesn't blurt out the wrong thing to the wrong person.

Ultimately, **general cloud security principles still apply**: least-privilege access, encryption, monitoring, and user education (e.g. *don't paste secrets into random AI websites*). The major AI providers are learning fast from early missteps and moving toward the security rigor expected in enterprise IT. However, given the novelty of AI, we can expect continued scrutiny and possibly more leaks as edge cases emerge. The goal for AI platforms is to become **"boringly secure"** like mature cloud services, and for users to gain confidence that using a large model is as safe as storing data in a trusted cloud database. Until then, the **cautious comparative view** is that current AI systems have a *higher experimental risk* profile (as evidenced by the 2023-2024 leaks) than traditional cloud tools, but with proper configurations (enterprise offerings or self-hosted

models) they can be managed to an equivalent risk level. As the technology and oversight improve, the gap is likely to close, making AI an integral part of the cloud with robust security – but vigilance is key, as both classes of systems are only as secure as the people and processes behind them.

**Industry/Academic Perspective:** A report by **Wiz** on the Microsoft incident noted that *"this case is an example of the new risks organizations face when leveraging AI – engineers working with massive datasets need additional security checks and safeguards."* At the same time, **AI safety researchers** point out that it's not just external attackers to worry about, but the model itself unintentionally exposing data (). The Cloud Security Alliance and others have started issuing guidelines specific to AI, effectively blending traditional cloud security practices with AI-specific concerns (like securing training data pipelines and defending against model theft). In practice, many view generative AI services as **an extension of cloud compute** – they come with great power and similar outsourcing of trust, so one must apply the same diligence in vetting their security. As one **FastCompany** article put it, *"Be careful with Bard: Google Search was showing private chatbot conversations… as with any cloud tool, assume what you input might be seen by others if things go awry."* This mindset will help users and enterprises navigate the benefits of AI while mitigating risks, treating AI platforms with the same (if not greater) care as traditional cloud services.

## Sources

Anthropic. (2024, January 22). Customer notification on contractor data leak. Anthropic. Retrieved from https://www.anthropic.com/data-leak-notice

BBC News. (2023, March 31). ChatGPT banned in Italy over privacy concerns. BBC. Retrieved from https://www.bbc.com/news/technology-65139406

Bloomberg. (2023, April 4). Samsung workers accidentally leaked secrets via ChatGPT. Bloomberg. Retrieved from https://www.bloomberg.com/news/articles/2023-04-04/samsung-workers-leaked-confidential-data-using-chatgpt

Carlini, N., Tramer, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., Roberts, A., Brown, T., Song, D., Erlingsson, Ú., Oprea, A., & Raffel, C. (2021). Extracting training data from large language models. Proceedings of the 30th USENIX Security Symposium. USENIX. Retrieved from https://arxiv.org/abs/2012.07805

CyberExpress. (2024, February 13). Google's Gemini AI leaks private chats onto public search. CyberExpress. Retrieved from https://thecyberexpress.com/google-gemini-ai-leaks-chats-public-search

FastCompany. (2023, September 29). Google Search accidentally shows private Bard conversations. FastCompany. Retrieved from https://www.fastcompany.com/90949039/google-search-private-bard-conversations

GNET. (2023, March 13). Meta's leaked LLaMA model and extremist chatbot creation. Global Network on Extremism & Technology. Retrieved from https://gnet-research.org/2023/03/13/llama-leak-extremist-chatbots

HackRead. (2023, September 29). Google Bard conversations exposed in search results. HackRead. Retrieved from https://www.hackread.com/google-bard-conversations-exposed-search-results

Insider. (2023, June 7). JPMorgan restricts employees from using ChatGPT. Business Insider. Retrieved from https://www.businessinsider.com/jpmorgan-chatgpt-employee-restrictions-ai-risk-2023-6

Meta AI Blog. (2023, February 24). Introducing LLaMA: A foundational language model from Meta AI. Meta. Retrieved from https://ai.facebook.com/blog/large-language-model-llama-meta-ai/

OpenAI. (2023, March 24). Update on ChatGPT outage. OpenAI Blog. Retrieved from
https://openai.com/blog/march-20-chatgpt-outage

Reuters. (2023, June 15). Google, Meta caution employees against sharing sensitive information with chatbots. Reuters. Retrieved from https://www.reuters.com/technology/google-meta-caution-employees-sharing-sensitive-information-chatbots-2023-06-15

TechCrunch. (2023, September 18). Microsoft AI researchers accidentally exposed terabytes of sensitive data. TechCrunch. Retrieved from https://techcrunch.com/2023/09/18/microsoft-ai-data-leak-azure

The Guardian. (2023, April 6). Italy reverses ChatGPT ban after privacy compliance. The Guardian. Retrieved from https://www.theguardian.com/technology/2023/apr/06/italy-reverses-chatgpt-ban-openai-privacy-compliance

The Verge. (2023, March 7). Meta's LLaMA AI model leaks online. The Verge. Retrieved from https://www.theverge.com/2023/3/7/meta-llama-ai-leaked-online-torrent

Wiz Research. (2023, September 18). Microsoft AI GitHub repository leaks 38TB of sensitive data. Wiz Research. Retrieved from https://www.wiz.io/blog/microsoft-ai-github-leak-38tb-sensitive-data