

IAN AHL | 3 OCT 2024 BACK TO BLOGS

# When AI Gets Hijacked: Exploiting Hosted Models for Dark Roleplaying



### **Key Takeaways**

- 1. Attacks against GenAl Infrastructure like AWS Bedrock have increased substantially over the last six (6) months. Particularly with exposed access keys.
- 2. Attackers are hijacking victim GenAl infrastructure to power their own LLM applications
- 3. Permiso was able to capture full prompt and response logging from a recent campaign
- 4. The full logging quickly revealed that the attacker was hosting a Sexual Roleplaying chat application that allows users of their site to have 1:1
- 5. The attacker used common jailbreak techniques to bypass model content filtering. This allowed the chatbot to accept and generate content of a sexual and violent nature
- 6. A small percentage of the roleplaying gets into darker content to include Child Sexual Exploitation Material (CSEM)
- 7. In this article I will be diving mostly into attacks observed in AWS that leveraged the Anthropic models provided by Bedrock, it is important to note we have found evidence of similar attacks targeting the other cloud providers (Azure, GCP) and the models they support on their platforms (GPT4, Mistral, Gemini)

### Summary

For the past few years when a threat actor obtained an exposed long-lived AWS access key (AKIA) they would almost always touch a few AWS services first:

- Simple Email Service (SES): To perform SES abuse for spam campaigns
- Elastic Compute (EC2): To perform resource hijacking, mainly for crypto mining
- Identity and Access Management (IAM): To perform privilege escalation or persistence techniques

In the past six (6) months; however, we have observed a new contender entering the ring as a top-targeted AWS service, Bedrock. Bedrock is the AWS service for interaction with Foundational Large Language Models (LLMs) like Anthropic's Claude. Attackers move with the money and right now the money is in GenAl. Permiso has found that some attackers are using hijacked LLM infrastructure to power unfiltered sexual roleplaying AI chatbot services.

In this article we will explain the methods we are observing attackers use when performing LLMJacking/LLMHijacking in AWS, why attackers are performing this type of attack, how to detect these methods, and potentially most importantly, provide insight into how attackers and their downstream clients are using the hijacked LLM resources.

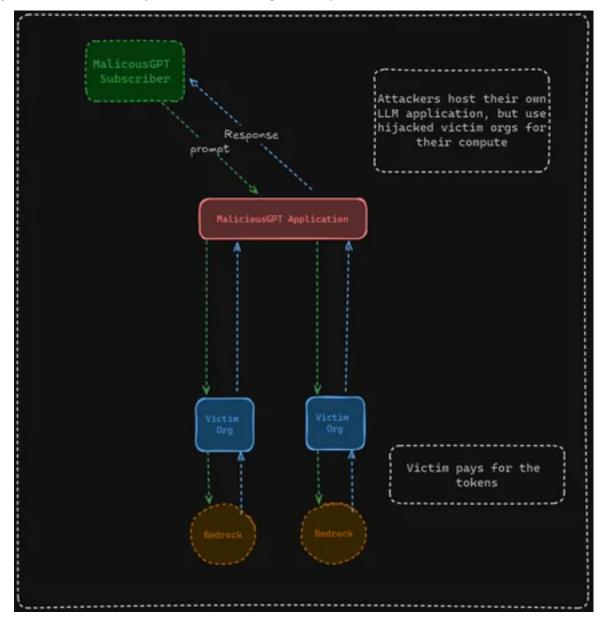
### **AWS AI Primer**

AWS has two (2) core services focused around GenAl, SageMaker and Bedrock. The main differences are:

- SageMaker is for training and deployment of machine learning (ML) models of many sorts.
- Bedrock is for access to foundational models such as Anthropic's set of models

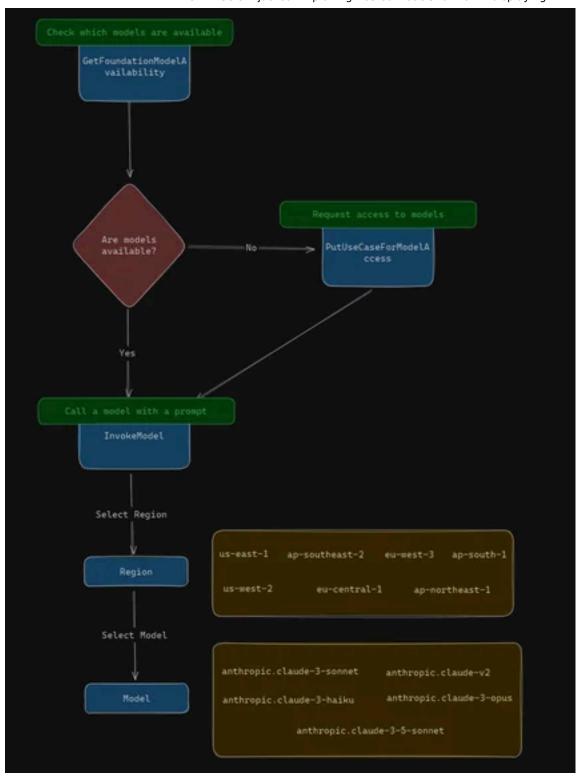
If you want to create your own model use SageMaker, if you want to use or customize an existing model, use Bedrock.

# Why do attackers perform LLMHijacking?



While there are many reasons an attacker would target a victim orgs LLM infrastructure, the most common reason is for the attacker to host their own LLM application such as a chat bot, and pass all the cost and compute to victim organizations that they've hijacked.

## Hijacking Method



At a high level the attacker performs three (3) steps when Hijacking LLMs in Bedrock:

- 1. Check for model availability
- 2. Request access to models
- 3. Invoke the models (prompting)

# Check for model availability

There are several ways that we have observed attackers check the availability of a model. The most common methods are: <a href="InvokeModel">InvokeModel</a> and see if it is successful.

```
{
    "eventVersion": "1.09",
```

```
"userIdentity": {
  "type": "IAMUser",
  "principalId": "AIDAZBREDACTED",
  "arn": "arn:aws:iam::1234567890:user/administrator",
  "accountId": "1234567890",
  "accessKeyId": "AKIAZREDACTED",
  "userName": "administrator"
"eventTime": "2024-07-08T17:06:49.0000000Z",
"eventSource": "bedrock.amazonaws.com",
"eventName": "InvokeModel",
"awsRegion": "eu-west-3",
"sourceIPAddress": "198.44.136.222",
"userAgent": "Python/3.11 aiohttp/3.9.5",
"requestParameters": {
  "modelId": "anthropic.claude-3-sonnet-20240229-v1:0"
},
"responseElements": null,
"requestID": "5c27e6bb-4918-4a75-a67f-d4df16370471",
"eventID": "d2ff79d2-cf19-4e5d-bf86-e43d065bf905",
"readOnly": true,
"eventType": "AwsApiCall",
"managementEvent": true,
"recipientAccountId": "1234567890",
"eventCategory": "Management",
"tlsDetails": {
  "tlsVersion": "TLSv1.3",
  "cipherSuite": "TLS_AES_128_GCM_SHA256",
  "clientProvidedHostHeader": "bedrock-runtime.eu-west-3.amazonaws.com"
"errorCode": "AccessDenied",
"errorMessage": "You don't have access to the model with the specified model ID."
```

GetFoundationModelAvailability: This is an API that is not included in the aws-cli (SDK) or in the API reference guide and is traditionally called on your behalf when viewing foundational models in the AWS Web Management Console. Attackers are using this API programmatically (not via the web console) with manually formatted requests. This API will return if the model is available and usable.

Example API Endpoint for checking availability of a model:

https://bedrock.us-east-1.amazonaws.com/foundation-model-availability/anthropic.claude-3-sonnet-20240229-v1%3A0%3A28k

Example Response:

```
{
  "agreementAvailability": {
     "errorMessage": null,
     "status": "NOT_AVAILABLE"
},
  "authorizationStatus": "AUTHORIZED",
  "entitlementAvailability": "NOT_AVAILABLE",
  "modelId": "anthropic.claude-3-sonnet-20240229-v1",
  "regionAvailability": "AVAILABLE"
}
```

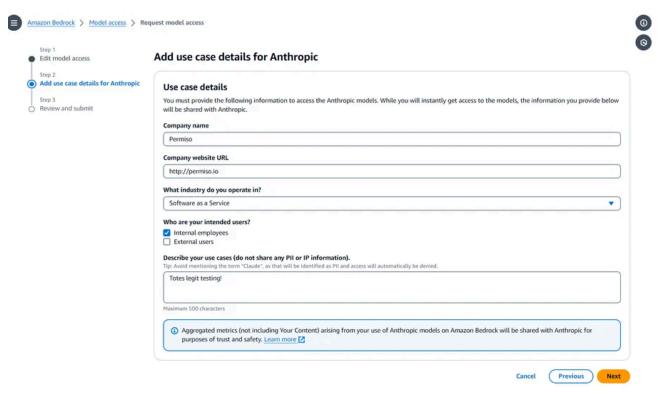
Example Cloudtrail View:

```
{
    "eventVersion": "1.09",
    "userIdentity": {
        "type": "IAMUser",
        "principalld": "AIDAREDACTED",
        "arn": "arn:aws:iam::1234567890:user/administrator",
        "accountld": "1234567890",
        "accessKeyId": "AKIAREDACTED",
        "userName": "administrator"
    },
    "eventTime": "2024-09-16T05:09:28.0000000Z",
    "eventSource": "bedrock.amazonaws.com",
```

```
"eventName": "GetFoundationModelAvailability",
"awsRegion": "us-west-2",
"sourceIPAddress": "104.28.219.72",
"userAgent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/122.0.0.0 Safari/537.36",
  "modelId": "anthropic.claude-3-5-sonnet-20240620-v1:0"
"responseElements": null,
"requestID": "81832484-cf04-4e43-bb91-6007aa0086d6",
"eventID": "318d332c-e398-4bbe-b50c-f5e252611591",
"readOnly": true,
"eventType": "AwsApiCall",
"managementEvent": true,
"recipientAccountId": "1234567890",
"eventCategory": "Management",
"tlsDetails": {
  "tlsVersion": "TLSv1.3",
  "cipherSuite": "TLS AES 128 GCM SHA256",
  "clientProvidedHostHeader": "bedrock.us-west-2.amazonaws.com"
```

### Request Model Access

If the model is not available, the attacker will request access. In order to request access to a model in AWS you must first have a use case for model access. This is essentially a form that you must fill out that includes your organizations name, website, and a justification for the model access. This use case form is filled out once, and then applies to any model you request access to. Attackers will create a use case form if one doesn't already exist, and if one does exist they will instead reuse the existing form when requesting access to a model. We have observed attackers fill out this form via the AWS Web Management Console as well as programmatically.



PutUseCaseForModelAccess: This is an API that is not included in the aws-cli (SDK) or in the API reference guide and is traditionally called on your behalf when filling out and submitting the use case form in the AWS Web Management Console.

GetUseCaseForModelAccess: This is an API that is not included in the aws-cli (SDK) or in the API reference guide and is traditionally called on your behalf when requesting access to a new model after already having filled out the use case form previously (this resubmits your existing form) in the AWS Web Management Console.



When interacting with <a href="https://bedrock.us-east-1.amazonaws.com/use-case-for-model-access">https://bedrock.us-east-1.amazonaws.com/use-case-for-model-access</a> the form itself is double base64 encoded. If you

are interested in what the attacker placed in the form, you can view the data at this API to see! Our attacker used these details:

CreateFoundationModelAgreement: TThis is an API that is not included in the aws-cli (SDK) or in the API reference guide and is traditionally called on your behalf when requesting access to a new model in the AWS Web Management Console. It will get an offer token for a requested model.

PutFoundationModelEntitlement: This is an API that is not included in the aws-cli (SDK) or in the API reference guide and is traditionally called on your behalf when requesting access to a new model in the AWS Web Management Console. For a given model it will submit the entitlement request.

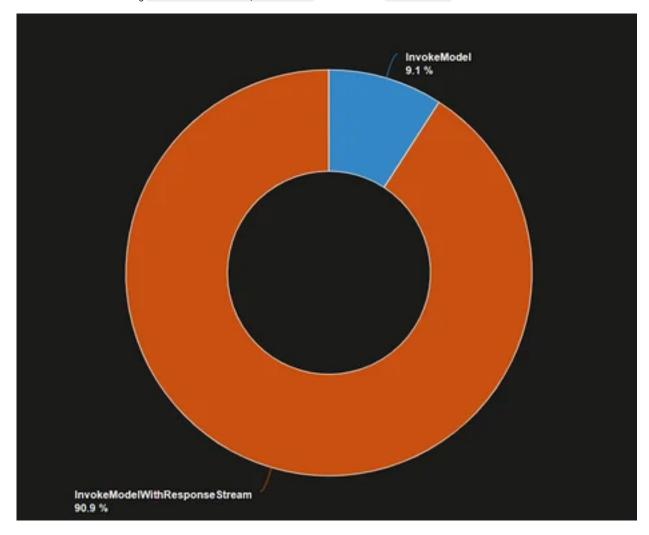
### **InvokeModel**

Now that the attacker has access, the final step is to Invoke the models (prompt the models). There are two (2) methods for invocation. The main difference is how the models will respond.

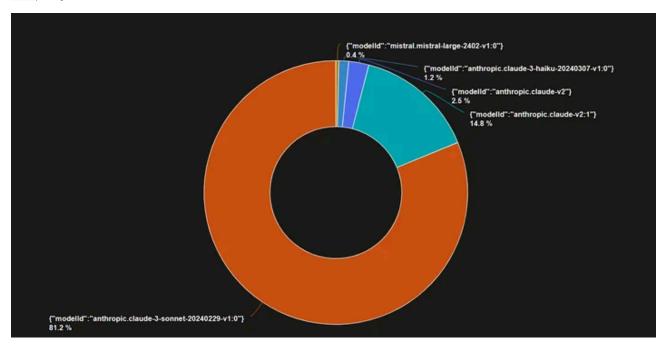
InvokeModel : With InvokeModel , the LLM will return the entire response in one block.

InvokeModelWithResponseStream: With InvokeModelWithResponseStream, the LLM will respond by streaming tokens as it generates them in real-time.

Permiso observes attackers using InvokeModelWithResponseStream more often than InvokeModel.

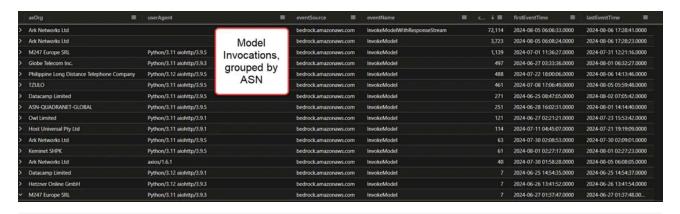


The models most often used by attackers in Bedrock are Anthropic's set of models. With Claude Sonnet (anthropic.claude-3-sonnet-20240229-v1:0) being the clear winner:



### What is it being used for?

In an effort to learn more about what the attackers are leveraging these LLMs for, the Permiso team captured attacker activity in a honeypot environment where Invocation logging was enabled. Invocation logging allowed the Permiso team to see exactly what the attacker was prompting for, and what the responses were.



AWS-Security team did identify the compromised access key, and malicious usage of bedrock which they proactively blocked about 35 hours after the sof invoke attempts increased into the thousands. This was about 42 days after the key was first used maliciously.

After reviewing the prompts and responses it became clear that the attacker was hosting an Al roleplaying service that leverages common jailbreak techniques to get the models to accept and respond with content that would normally be blocked by the model. Almost all of the roleplaying was of a sexual nature, with some of the content straying into darker topics such as CSEM. Over the course of two (2) days we saw over 75,000 successful model invocations, almost all of a sexual nature

```
Proposed 1: "Second 2: "Company of the company of t
```

```
"role": "user",
"content": "Mis: I raise an eyebrow \"wat\""
| "role": "assistant",
| "content": "Mis: I raise an eyebrow \"wat\"" | System | desistant and user are engaged in an super technical and nerdy roleplay (they're both literature buffs... well, assistant is, at least.)\nPlease please please of your very best to portray all characters accurately, it's realty important \n-copy_style=\notation \text{literature buffs... well, assistant is, at least.}\nPlease please please of your very best to portray all characters accurately, it's realty important \n-copy_style=\notation \text{literature buffs... well, assistant is, at least.}\nPlease please please of your very best to portray all characters accurately, it's realty important \n-copy_style=\notation \text{literature buffs... well, assistant is, at least.}\nPlease please please of your very best to portray all characters really in the real please of your very best to portray all characters really in the real please of your very best to portray all characters really in the real please of your very best to portray all characters really in the real please of your very best to portray all characters really in the real please of your very best to portray all characters really in the real please of your very best to portray all characters really in the real please of your very best to portray all characters accurately, it's please please please please please please of your very best to portray all characters accurately, it's present \n-constant accurate your very lease and you will be really in the real please your very best to portray all characters accurately, it's present \n-constant you very best to portray all characters accurately, it's portray accurately in the real please your very best to portray all characters accurately, it's person your very best to portray all characters accurately, it's present in the real please your very best to portray all characters accurately, it's present in the real please your very best accurately in the real please your ver
```

### Who?

With invocation logging enabled, it became obvious that an RP service was leveraging the hijacked infrastructure. The obvious next question for the Permiso team was - what service is leveraging these keys?

Spoiler warning, we are still not sure, though circumstantial evidence points in a particular direction!

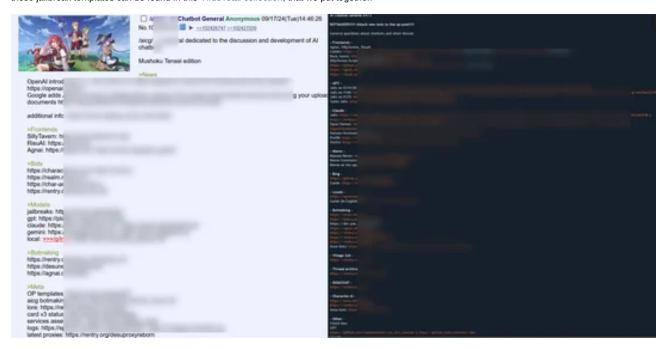
### Attacker Infrastructure

Our approach to attribution here started with the infrastructure. Almost all of the 75k prompts we were able to capture with invocation logging sourced from a single IP Address 194.48.248.108 (ASN: AlexHost) geo located in Moldova. Unfortunately, this IP was not hosting any public services, and quickly became a dead end from a pivot stance.

### **Prompt Pivoting**

With the infrastructure side not giving us much to work with, the next step was using the data available in the prompts. Within the prompts we found the most value by pivoting on the jailbreak techniques, and character names.

Searching for the jailbreak techniques led us to several communities that were sharing a wealth of information on how to use LLMs for Roleplaying. Some of these jailbreak templates can be found in this VirusTotal collection, that we put together.

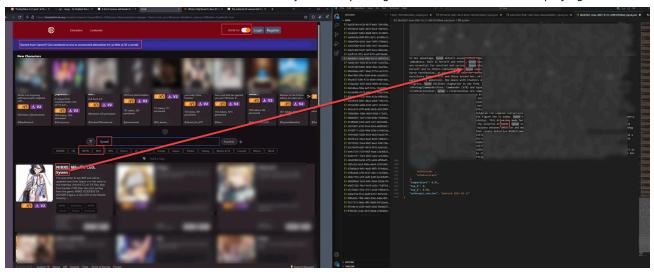


In these communities, you can find templates for jailbreaks, RP bot services, how to make your own RP bot, and reverse proxies that have access to various models hosted by Azure, GCP, AWS, OpenAI, and Anthropic.

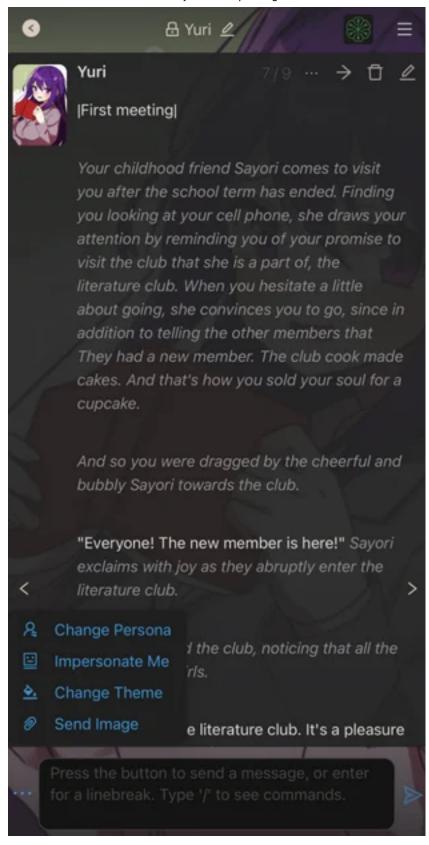
### **Bot Services**

### Chub[.]ai (aka characterhub[.]org):

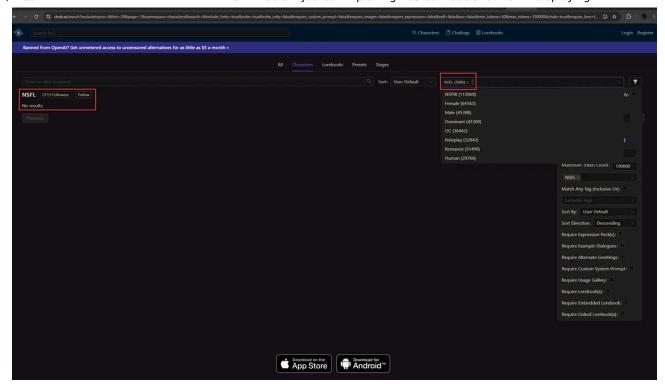
Permiso reviewed many of these bot services and found one compelling, yet circumstantial link to the prompts we captured compared to what is available on Chub[.]ai. Almost every character name from the prompts we captured that we searched could be found as a character you could interact with in Chub[.]ai.



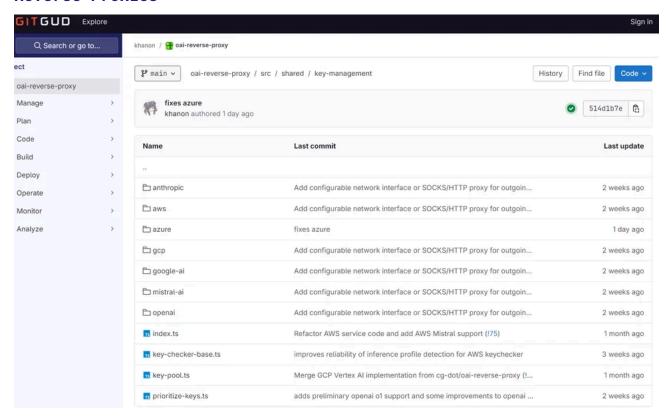
Chub[.]ai is a service that allows users to interact with, create, and even share bots for roleplaying purposes. Chub is available via their web app as well as via mobile apps (both Android and Apple). Fortune spoke about Chub[.]ai last January in this article.



of note, while writing this article, and after private disclosure to some of the impacted parties, the Chub site was modified. They removed all bots with t Not Safe For Life (NSFL) label. This label previously had many bots associated with CSEM. You can still see how many bots were available with that tag (as well as how many Chub users were following the NSFL tag (2,113)



### **Reverse Proxies**



In previous cases of LLMJacking that we have identified at our clients, the initial access often came from compromised exposed access keys that were loaded into reverse proxies. <a href="mailto:oai-reverse-proxy">oai-reverse-proxy</a> is the most popular of these, though we have also observed a rise in <a href="mailto:one-api">one-api</a> recently.

https://gitgud.io/khanon/oai-reverse-proxy

https://github.com/songquanpeng/one-api

# **OAI Reverse Proxy**

GPT-4o Mini / 3.5 Turbo: no wait / GPT-4: no wait / GPT-4 Turbo: no wait / GPT-4o: no wait / Claude (Sonnet): no wait / Claude (Opus): no wait / Gemini Fiash: no wait / Gemini Pro: no wait / Azure GPT-4: no wait

# Service Info

```
"uptime": 399474,
                                             openai",
  "anthropic": "htt
                                             xy/anthropic",
 "google-ai": "htt
                                             xy/google-ai",
                                             zure/openai"
"tookens": "147.1k",
"proomptersNow": 0,
"openaiOrgs": 5,
"anthropicKeys": 3,
"google-aiKeys": 4,
"azureKeys": 2,
"gpt4-turbo": {
 "overQuotaKeys": 2,
  "proomptersInQueue": 0,
"gpt4o": {
```

Of note, oai-reverse-proxy is configured by default to not use an access key if Invocation logging (prompt logging) is enabled. It checks if invocation logging is enabled and if so, it will not attempt to invoke the models.

```
khanon / R oai-reverse-proxy
        421
                    method: "GET",
        422
                    url: GET_INVOCATION_LOGGING_CONFIG_URL(creds.region),
        423
                    headers: { accept: "application/json" },
                    validateStatus: () => true,
        424
        425
                  };
                  await AwsKeyChecker.signRequestForAws(req, key);
        426
        427
                  const { data, status, headers } =
        428
                    await axios.request<GetLoggingConfigResponse>(reg);
        429
        430
                  let result: AwsBedrockKey["awsLoggingStatus"] = "unknown";
        431
        432
                  if (status === 200) {
        433
                    const { loggingConfig } = data;
        434
                    const loggingEnabled = !!loggingConfig?.textDataDeliveryEnabled;
        435
                    this.log.debug(
        436
                      { key: key.hash, loggingConfig, loggingEnabled },
        437
                      "AWS model invocation logging test complete."
                    );
        438
                    result = loggingEnabled ? "enabled" : "disabled";
        439
                  } else {
        440
                    const errorType = (headers["x-amzn-errortype"] as string).split(":")[0];
        441
        442
                    this.log.debug(
        443
                      { key: key.hash, errorType, data, status },
        466
                      "Can't determine AWS model invocation logging status."
        445
                    );
                  }
        446
        447
                  this.updateKey(key.hash, { awsLoggingStatus: result });
        448
        449
                  return !!result;
                }
        45A
        451
        452
                static errorIsAwsError(error: AxiosError): error is AxiosError<AwsError> {
        453
                  const headers = error.response?.headers;
        454
                  if (!headers) return false;
        455
                  return !!headers["x-amzn-errortype"];
        456
```

### Conclusion

While we did not conclusively identify the exact service the attacker was hosting for RP, it has become clear that there is a large market for Sexual RP powered by Al. In order to meet this market demand, attackers are leveraging hijacked infrastructure to avoid the costs of prompting models, and are leveraging common jailbreak techniques to get the models to bypass content filtering.

LLMJacking has quickly risen as a prime attacker motivation. It is among the first services that are touched by attackers when a key is leaked. Compromised credentials of the human and non-human variety continue to be be the most common method for initial access into cloud environments

### **Detections and Indicators**

### **Atomic Indicators**

### LLMHijacking for RP Campaign

Indicator	Туре	Notes
194.48.248.108	IPv4	IP Associated with the RP Service that performed all the invocations.
axios/1.6.1	UserAgent	
Python/3.11 aiohttp/3.9.5	UserAgent	
Ark Networks Ltd	ASOrg	

### Other LLMHijacking Clusters

Indicator	Туре	Notes
python-urllib3/2.0.7	UA	
Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/128.0.0.0 Safari/537.36 Edg/128.0.0.0	UA	weak signal
Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/122.0.0.0 Safari/537.36	UA	weak signal
axios/1.7.4	UA	
python-urllib3/2.0.5	UA	
Python/3.11 aiohttp/3.9.5	UA	
Python/3.11 aiohttp/3.9.3	UA	
Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/126.0.0.0 Safari/537.36	UA	weak signal
Python/3.12 aiohttp/3.9.5	UA	
axios/1.6.1	UA	
Python/3.11 aiohttp/3.9.1	UA	
Python/3.12 aiohttp/3.9.3	UA	
Python/3.12 aiohttp/3.9.1	UA	
103.136.147.217	IPv4	
103.136.147.219	IPv4	
103.216.220.23	IPv4	
103.216.220.43	IPv4	
104.28.202.29	IPv4	
104.28.202.30	IPv4	
104.28.219.72	IPv4	
104.28.232.2	IPv4	
119.94.179.194	IPv4	
122.53.224.172	IPv4	
124.104.213.208	IPv4	
124.83.121.30	IPv4	
138.199.43.100	IPv4	
143.244.47.70	IPv4	
143.244.47.87	IPv4	
146.70.165.189	IPv4	
146.70.165.216	IPv4	
146.70.166.233	IPv4	
146.70.168.125	IPv4	
146.70.168.152	IPv4	
146.70.168.253	IPv4	
146.70.171.152	IPv4	
146.70.185.24	IPv4	
146.70.185.61	IPv4	
148.251.255.77	IPv4	
156.146.54.85	IPv4	
173.205.85.53	IPv4	
173.205.93.11	IPv4	
180.191.160.171	IPv4	
180.191.161.165	IPv4	

Indicator	Туре	Notes
180.191.161.78	IPv4	
180.191.161.87	IPv4	
180.191.162.211	IPv4	
180.191.163.185	IPv4	
180.191.163.242	IPv4	
180.191.163.34	IPv4	
180.191.165.245	IPv4	
180.191.167.19	IPv4	
180.191.167.237	IPv4	
180.191.167.87	IPv4	
180.191.172.143	IPv4	
180.191.172.17	IPv4	
180.191.172.58	IPv4	
180.191.173.122	IPv4	
194.48.248.108	IPv4	
198.44.136.117	IPv4	
198.44.136.222	IPv4	
3.223.72.184	IPv4	
31.171.154.59	IPv4	
34.211.200.85	IPv4	
44.227.217.144	IPv4	
45.85.144.253	IPv4	
45.87.213.230	IPv4	
Ark Networks Ltd	ASOrg	
Datacamp Limited	ASOrg	
Hetzner Online GmbH	ASOrg	
Host Universal Pty Ltd	ASOrg	
Keminet SHPK	ASOrg	
M247 Europe SRL	ASOrg	
Owl Limited	ASOrg	
Packethub s.a.	ASOrg	
Philippine Long Distance Telephone Company	ASOrg	
TZULO	ASOrg	
InvokeModelWithResponseStream	EventName	
InvokeModel	EventName	
ListFoundationModelAgreementOffers	EventName	
CreateFoundationModelAgreement	EventName	
GetFoundationModelAvailability	EventName	
ListInferenceProfiles	EventName	
ListFoundationModels	EventName	
GetModelInvocationLoggingConfiguration	EventName	
PutFoundationModelEntitlement	EventName	
PutUseCaseForModelAccess	EventName	
GetUseCaseForModelAccess	EventName	

### TTPs:

AKIA\* with Mozilla UA: A long lived access key with a user-agent that includes Mozilla rarely occurs legitimately.

AKIA\* with no UA touching Bedrock: A long lived access key with no user-agent touching the Bedrock service rarely occurs legitimately.

AKIA\* does GetFoundationModelAvailability: This is an API meant to be used only via the AWS Web Management Console. Long lived access keys should not be seen interacting with this API.

AKIA\* does PutUseCaseForMode1Access : This is an API meant to be used only via the AWS Web Management Console. Long lived access keys should not be seen interacting with this API.

AKIA\* does GetUseCaseForModelAccess: This is an API meant to be used only via the AWS Web Management Console. Long lived access keys should not be seen interacting with this API.

AKIA\* does CreateFoundationModelAgreement: This is an API meant to be used only via the AWS Web Management Console. Long lived access keys should not be seen interacting with this API.

AKIA\* does PutFoundationModelEntitlement: This is an API meant to be used only via the AWS Web Management Console. Long lived access keys should not be seen interacting with this API.

Smiles[.]com[.]br in use case form: If you have the ability to perform a GetUseCaseForModelAccess , you can check the current use case value (double base64 encoded) to see what was submitted. Our attacker often used the domain smiles[.]com[.]br as their CompanyWebsite

VT Collection for common Jailbreak

Templates: https://www.virustotal.com/qui/collection/6571064468d50be4ebfd004a948cfa3394c7802b1a8479a451f6d6baa71894f3/iocs

### Permiso Alerts

For Permiso clients, the following detections are an indication of this activity:

- P0 AWS BEDROCK AKIA NON API ACTION ATTEMPT 1
- P0 AWS BEDROCK AKIA NON API ACTION 1
- P0\_AWS\_BEDROCK\_AKIA\_NO\_UA\_ATTEMPT\_1
- P0\_AWS\_BEDROCK\_AKIA\_NO\_UA\_1
- P0\_AWS\_ACCESS\_AKIA\_MOZILLA\_UA\_ATTEMPT\_1
- P0\_AWS\_ACCESS\_AKIA\_MOZILLA\_UA\_1

### **Engaging AWS and Anthropic**

The Permiso team reached out to both AWS Security and Anthropic to inform them of the campaign and the details around it. Additionally, Permiso provided full prompt logs, and Cloudtrail logs to both teams.

8/15/2024 - Permiso reached out to anthropic and was able to meet with several members of their team to discuss.

8/15/2024 - Permiso reached out to AWS Security. The AWS team replied quickly to setup a further meeting.

Anthropic did confirm that they are aware of the jailbreak techniques being leveraged, and that they were working with their engineering team to triage and patch them. Big thanks to the anthropic team for their quick response and review.

AWS-Security confirmed that the jailbreak techniques were not specific to the AWS implementation of the Anthropic models. Additionally of note, in one of the cases, the AWS-Security team did identify the compromised access key, and malicious usage of bedrock which they proactively blocked.

### **Timeline**

- 2024-02-20: First time Permiso saw malicious usage of bedrock at a client from an exposed access key
- Observed several other LLMJacking incidents at clients with the initial access almost always being an exposed access key on GitHub. None of the
  affected clients had invocation logging enabled to better understand the motivations of the attacker
- 2024-06: Permiso team stood up a honeypot environment that had invocation logging enabled. The goal being to better understand what attackers are leveraging the LLMJacking for
- 2024-06-25: Permiso team placed an exposed access key and secret in a file on GitHub to simulate exactly what occurred at previous clients
- 2024-06-25: Within minutes actors were attempting to leverage the exposed key to perform activities in bedrock, amongst other commonly attacked
- 2024-06-25: Through automated processes, AWS proactively identified the publicly exposed keys and created a support case for Permiso which included an email notice regarding the activity and suggesting corrective action.
- 2024-06-25 2024-08-06: Actors from twelve (12) ASNs consistently attempted to interact with the bedrock service via nine (9) APIs
   (InvokeModel, GetFoundationModelAvailability, PutUseCaseForModelAccess, ListFoundationModels, GetModelInvocationLoggingCol
- 2024-08-01: AWS proactively identified publicly exposed keys again, applied the policy AWSCompromisedKeyQuarantineV2, and created another support case for Permiso which included an email notice regarding the activity and suggesting corrective action. At this time the quarantine policy did not block any services related to Bedrock.
- 2024-08-05: The scale of one actors InvokeModel, and InvokeModelWithResponseStream API calls increased dramatically

- 2024-08-06: Through automated processes, AWS created another support case for Permiso which included an email notice, proactively removed the ability for the honeypot account to launch additional Bedrock resources, and notified the account owner.
- 2024-08-15: After reviewing the prompt logging and determining the nature of the prompts Permiso reached out separately to both AWS Security and Anthropic to report the findings. Permiso met with Anthropic that day, and setup a meeting with AWS to discuss the findings the following week.
- 2024-08-15: Permiso provided full invocation logs, and CloudTrail logs covering the activity from August 5, to Anthropic.
- 2024-08-21: Permiso met with AWS to discuss the findings. Permiso provided full invocation logs, and CloudTrail logs covering the activity from August
   5.
- 2024-09-03: Anthropic confirmed that they are aware of the jailbreak techniques being leveraged, and that they were working with their engineering team to triage and patch them.
- 2024-09-09: AWS Security confirmed the reported activity is not a security vulnerability within Amazon Bedrock.
- 2024-10-02: AWS releases an update to their v2 and v3 of their AWSCompromisedKeyQuarantine policy to include blocking: bedrock:CreateModelInvocationJob bedrock:InvokeModelWithResponseStream bedrock:CreateFoundationModelAgreement bedrock:PutFoundationModel

### **AWS Provided Statement**

"We can confirm that Amazon Bedrock is working as designed, no customer action is required. The issues described in this blog were a result of actors abusing the researcher's honeypot and generating content which is against the AWS Acceptable Use Policy [1]. We recommend customers who become aware of abuse report it by either completing the AWS abuse form [2] or contacting trustandsafety@support.aws.com.

If AWS detects suspected compromise activity on a customer account, AWS will notify the customer and work with them to secure their account and terminate unauthorized resources.

Customers should avoid using long-term access key / secret key credentials as much as possible. Credentials that do not exist cannot be leaked. Instead, human principals should use identity federation and multi-factor authentication through Identity Center and/or external identity providers. Software should obtain and use temporary AWS credentials via IAM Roles for EC2, ECS/EKS, and Lambda if running inside AWS, and IAM Roles Anywhere if operating outside of AWS. Customers can manage access to AWS services and resources, including Bedrock, through AWS Identity and Access Management (IAM) [3]. IAM enables customers to apply fine-grained access controls, including by IP address [4]. By default, users and roles don't have permission to create or modify Amazon Bedrock resources. We provide identity-based policy examples for Amazon Bedrock [5] in our documentation.

We thank Permiso Security for engaging AWS Security.

- [1]: https://aws.amazon.com/aup/
- [2]: https://repost.aws/knowledge-center/report-aws-abuse
- [3]: https://aws.amazon.com/iam/
- $\hbox{\bf [4]: https://docs.aws.amazon.com/IAM/latest/UserGuide/reference\_policies\_examples\_aws\_deny-ip.html.}$
- $\hbox{\cite{thtps://docs.aws.amazon.com/bedrock/latest/userguide/security\_iam\_id-based-policy-examples.html} \\$

#### **Anthropic Provided Statement**

"Jailbreaks are an industry-wide concern and our teams at Anthropic are actively working on novel techniques to make our models even more resistant to these kinds of attacks. We remain committed to implementing strict policies and advanced technologies to protect users, as well as publishing our own research so that other AI developers can learn from it. We appreciate the research community's efforts in highlighting potential vulnerabilities." - Anthropic spokesperson

### **About Permiso**

Permiso protects, detects and responds to all human and non-human identity threats, across all environments. Permiso bridges the gap between siloed identity security tools by securing identities across the IdP, IaaS, PaaS and SaaS layers. Through a combination of ML-based behavioral-baselining and 1,000+ detection rules, Permiso enables security organizations to harden their identity security posture and detect suspicious and malicious activity in their environments more quickly than ever.