

TECHDOCS

# AI Model Security

---

## Contact Information

Corporate Headquarters:  
Palo Alto Networks  
3000 Tannery Way  
Santa Clara, CA 95054  
[www.paloaltonetworks.com/company/contact-support](http://www.paloaltonetworks.com/company/contact-support)

## About the Documentation

- For the most recent version of this guide or for access to related documentation, visit the Technical Documentation portal [docs.paloaltonetworks.com](http://docs.paloaltonetworks.com).
- To search for a specific topic, go to our search page [docs.paloaltonetworks.com/search.html](http://docs.paloaltonetworks.com/search.html).
- Have feedback or questions for us? Leave a comment on any page in the portal, or write to us at [documentation@paloaltonetworks.com](mailto:documentation@paloaltonetworks.com).

## Copyright

Palo Alto Networks, Inc.  
[www.paloaltonetworks.com](http://www.paloaltonetworks.com)

© 2025-2025 Palo Alto Networks, Inc. Palo Alto Networks is a registered trademark of Palo Alto Networks. A list of our trademarks can be found at [www.paloaltonetworks.com/company/trademarks.html](http://www.paloaltonetworks.com/company/trademarks.html). All other marks mentioned herein may be trademarks of their respective companies.

## Last Revised

October 24, 2025

---

# Table of Contents

<b>Secure Your AI Models with AI Model Security.....</b>	<b>5</b>
What is a Model?.....	6
What is AI Model Security?.....	7
Impact of model vulnerabilities.....	8
AI Model Security Core Components.....	9
Get Started with AI Model Security.....	11
Create a Deployment Profile for Prisma AIRS AI Model Security.....	12
Configure Identity and Access Management.....	15
Install AI Model Security.....	17
Default Security Groups and Rules.....	19
Scanning Models.....	20
Viewing Scan Results.....	23
Understanding Model Security Scan Results.....	32
Understanding Model Security Rules.....	34
How Scans, Security Groups, and Rules Work Together.....	35
Connecting Scans to Rules.....	37
Security Rule Checks.....	38
Hugging Face Model Rules.....	41
Reporting New Threats.....	44
Supported Model Formats.....	46

## Table of Contents

---

# Secure Your AI Models with AI Model Security

## Where Can I Use This?

- Prisma AIRS (AI Model Security)

## What Do I Need?

- Prisma AIRS AI Model Security License

## What is a Model?

Model is the collection of files required to perform a single base inference pass. Model has the following structure:

- **Source**—Where the AI model lives.
- **Model**—Logical entity (like "sentiment-analyzer").
- **ModelVersion**—Specific AI model version (like "v1.2.3").
- **Files**—Artifacts of that AI model version.

Models are the foundational asset of AI/ML workloads and already power many of your key systems in use today. AI model security focuses on securing your models against threats like:

- **Deserialization Threats**: Protecting your models from executing malicious and unknown code at load time.
- **Neural Backdoors**: Detecting Manchurian Candidate like models.
- **Runtime Threats**: Protecting your models from executing malicious and unknown code at inference time.
- **Invalid Licenses**: Ensuring that your models are not using invalid licenses.
- **Insecure Formats**: Ensuring that your models use formats that help prevent threats.

## What is AI Model Security?

AI Model Security is an enterprise application designed to enforce comprehensive security standards for both internal and external machine learning models deployed in production environments. The application addresses a critical gap in organizational security practices where machine learning models, despite their significant impact on business operations, often lack the rigorous security validation that is standard for other data inputs and systems.

In most enterprise environments, traditional data inputs such as PDF files undergo extensive security scrutiny before processing, yet machine learning models that drive critical business decisions frequently bypass equivalent security measures. This disparity creates substantial operational risk, as compromised or inadequately validated models can impact business logic, data integrity, and decision-making processes. AI Model Security solves this problem by providing comprehensive model validation through automated security assessments, multi-source support for both internally developed and third-party models, and proactive identification and remediation of model-related security vulnerabilities.

By implementing AI Model Security, organizations can establish consistent security standards across all ML model deployments, significantly reducing operational risk from unvalidated or compromised models. It enables secure adoption of third-party ML solutions while maintaining organizational security rules and industry compliance standards. Additionally, it provides comprehensive audit trails and compliance reporting capabilities, ensuring that AI Model Security assessments meet regulatory requirements and internal governance standards.

## Impact of model vulnerabilities

Models can execute arbitrary code and existing tooling is not checking that for you. Models have been found at the root of cloud take over attacks, and can be used to exfiltrate data, or even to execute ransomware attacks. The sensitivity of the data that models are trained on and exposed to at inference time makes them a prime target for attackers.

# AI Model Security Core Components

AI Model Security enables you to have flexible controls to secure, validate, and manage AI models across different sources through *Security Groups*, *Sources*, *Rules*, and *Scans*.

AI Model Security delivers a comprehensive framework to establish and enforce security standards for AI models across your organization. Unlike traditional security tools that simply scan for malware, AI Model Security recognizes that AI models require more nuanced security considerations that incorporate license validation, file format verification, and context-specific security checks based on the teams and environments using the models.

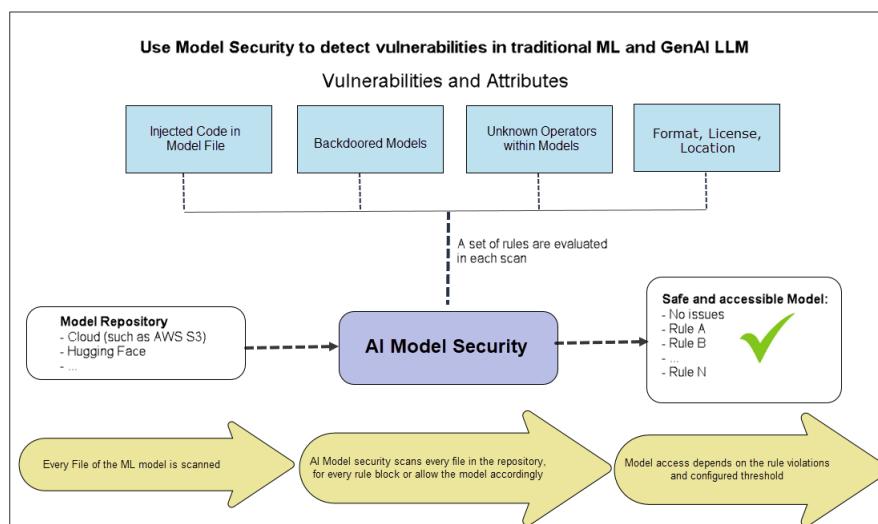
The AI Model Security approach moves beyond the simplified first-party versus third-party model distinction to provide granular security controls that scale with enterprise needs. This approach centers around four key components: *Security Groups*, *Sources*, *Rules*, and *Scans*.

Entity	Description	Examples
Security Groups	Serve as the foundation of your AI Model Security posture, allowing you to combine specific rules and requirements for models from a particular source.	<ul style="list-style-type: none"> <li>• HuggingFace-Research</li> <li>• S3-Production</li> <li>• Partner-S3-Audit</li> </ul>
Source	Each Security Group is assigned to a specific Source, which represents where model artifacts reside, such as HuggingFace for external models or Local Storage and Object Storage for internal models. The source designation is crucial as it provides metadata that powers specific security rules applicable to models from that source.	<ul style="list-style-type: none"> <li>• Hugging Face</li> <li>• S3</li> <li>• Local Disk</li> </ul>
Rules	Within each Security Group, you configure Rules that define the specific evaluations performed on models. Rules can verify proper licensing, check for approved file formats, scan for malicious code, and detect architectural backdoors. Each Rule can be enabled or disabled and configured as blocking or non-blocking, giving you precise control over which security issues prevent model usage versus those that simply generate warnings.	<ul style="list-style-type: none"> <li>• License Existence Check</li> <li>• Serialization Format Safety</li> <li>• Author Verification</li> <li>• Malicious Backdoor Detection</li> </ul>
Scan	<p>When models are evaluated against these Rules, a Scan is performed, documenting the verdict across all rules. These Scans create an audit trail of security evaluations and serve as decision points to either promote secure models forward or block potentially threatening ones early in your workflow.</p> <p>Here's what a typical scan will look like:</p>	Scan of fraud-detector:v2.1.0 using S3-Production group

Entity	Description	Examples
	<pre>Scan of ModelVersion "sentiment-model-v1.2.3": └─ License Compliance Rule (metadata, artifact-level) → PASS └─ Source Verification Rule (metadata, artifact-level) → PASS └─ Serialization Safety Rule (security, file-level: model.pkl) → FAIL (Critical) └─ Serialization Safety Rule (security, file-level: tokenizer.pkl) → PASS └─ Aggregate Result: FAIL (Critical severity exceeds threshold)</pre>	

AI Model Security leverages rules to help organizations establish sophisticated, scalable security frameworks tailored to their specific requirements. This flexible approach enables teams to enforce strict blocking mechanisms for high-severity threats while maintaining non-disruptive alerting for compliance monitoring—allowing security teams to effectively manage risk without hindering developer productivity. The result delivers dual benefits: end users gain confident access to vetted models through a seamless experience, while security teams receive comprehensive protection for their AI/ML infrastructure.

To implement AI Model Security effectively, you'll typically need at least two Security Groups: one for external models using HuggingFace as a Source, and another for internal models using Local or Object Storage Sources. This separation allows you to apply appropriate security standards based on the origin and intended use of models across your organization.



# Get Started with AI Model Security

Where Can I Use This?	What Do I Need?
<ul style="list-style-type: none"> <li>Prisma AIRS (AI Model Security)</li> </ul>	<input type="checkbox"/> Prisma AIRS AI Model Security License

## AI Model Security Workflow

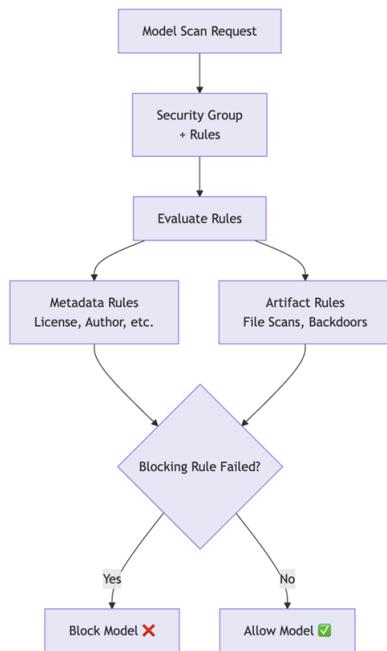
AI Model Security operates on a two-tier hierarchical structure where each security group encompasses one or more security rules.



- 1. Source Type Binding:** Model Security Groups are initially created and associated with a specific source type (such as, S3 buckets).
- 2. Rule Configuration:** Populate these groups with relevant rules. AI Model Security provides intelligent suggestions and validation based on the selected source type, including default rule collections tailored to each source.  
For instance, a "Verify Author" rule wouldn't be available for S3-based groups since S3 doesn't maintain author metadata.  
Rules operate in either blocking or non-Blocking modes, streamlining the previous severity-based threshold system from AI Model Security.
- 3. Scanning Process:** During model evaluation, the scan request specifies the applicable model security group to establish context. The system then processes two rule categories:
  - Metadata Rules:** Validate model metadata from the source platform
  - Artifact Rules:** Conduct comprehensive analysis of model files
- 4. Result Processing:** Each rule produces a binary PASS or FAIL outcome. The final scan verdict aggregates all rule results and determines whether to apply blocking enforcement (for critical issues) or non-Blocking responses (for logging and alerts).

The AI Model Security delivers structured flexibility—the same model can undergo evaluation against multiple Model Security Groups (such as separate development and production

configurations), each with distinct Source Types and rule sets, enabling context-appropriate security rules.



## Create a Deployment Profile for Prisma AIRS AI Model Security

Where Can I Use This?	What Do I Need?
<ul style="list-style-type: none"><li>Prisma AIRS</li></ul>	<input checked="" type="checkbox"/> Software NGFW Credits

Use this information to deploy a tenant for AI Model Security. Consider the following when creating a deployment profile:

- AI Model Security is licensed via NGFW credits.
- The deployment profile is accessible to all SCM Pro users.
- The deployment profile is accessible on the Customer Support Portal.
- For this early access release, the deployment profile name is **AI Model Security (preview)**.



You need an active deployment profile and have a proper Identity & Access to view the AI Model Security and its features in Prisma AIRS.

Creating a deployment profile for AI Model Security includes the following steps:

1. Create the deployment profile.
2. Associate the deployment profile with a tenant.

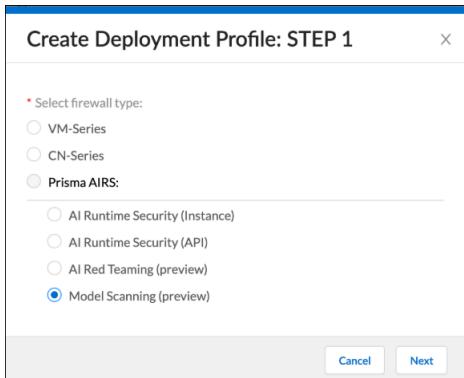
### Create a Deployment Profile for Prisma AIRS AI Model Security in the Customer Support Portal

**STEP 1 |** Log in to the Palo Alto [Customer Support Portal \(CSP\)](#).

**STEP 2 |** Navigate to Products > Software/Cloud NGFW Credits.

**STEP 3 |** Locate your credit pool and click **Create Deployment Profile**.

**STEP 4 |** Under **Select firewall type**, expand **Prisma AIRS** and select **AI Model Security (preview)**.



**STEP 5 |** Select **Next**.

**STEP 6 |** Enter a **Profile Name**.



Click **Calculate Estimated Cost** to view the credits used for AI Model Security.

**STEP 7 |** Click **Create Deployment Profile**.

This takes you to the **Software NGFW Credits** page in the **Customer Support Portal**. The page displays the list of deployment profiles within the credit pool for the selected account.

After creating the deployment profile, you'll associate it with a tenant.

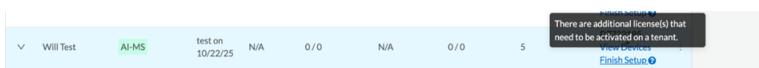
### Associate the Deployment Profile with a Tenant

Use the information in this section to associate a deployment profile with a tenant. Consider the following:

- Once the deployment profile is created, you're provided with a link to the [Palo Alto Networks Hub](#) to associate the profile with a TSG.
- You can create a new tenant (TSG), or, you can associate the profile with an existing tenant.
- If a new tenant is created, a Strata Cloud Manager (SCM) instance is created; this process can take approximately 15-20 minutes to complete.
- If you are using an existing tenant, a SCM instance is typically associated with it.
- Because AI Model Security is accessible for all SCM Pro users, you'll see SCM and Strata Logging Service (SLS) instances being created; these elements are part of the SCM Pro license.

To associate the deployment profile with a tenant:

**STEP 1 |** In the CSP dashboard, locate your deployment profile and click **Finish Setup**; this redirects you to the Palo Alto Networks hub.



**STEP 2 |** Select the name of the CSP account in which you created the deployment profile.

 If you select a new tenant, also select the region where the tenant should be created. This is referred to as the platform region. AI Model Security determines which regions are supported.

Only Region: United States - Americas is currently supported.

**STEP 3 |** Select the deployment profile that you want to associate with the tenant.

**STEP 4 |** Select **None** in Additional Services.

Set Up Profile(S)

Recipient: test	Edit																														
Region: United States - Americas	Edit																														
Select Deployment Profile(s): The deployment profile(s) shown are based on your customer support account selection.																															
<table border="1"><thead><tr><th><input type="checkbox"/></th><th>Profile Name</th><th>Auth Code</th><th>Status</th></tr></thead><tbody><tr><td><input type="checkbox"/></td><td>initial-ms-dp</td><td>D4073957</td><td>Available</td></tr><tr><td><input type="checkbox"/></td><td>Strata Cloud Manager Pro, AI Model Scanning</td><td></td><td></td></tr><tr><td><input type="checkbox"/></td><td>Josh test</td><td>D6717430</td><td>Available</td></tr><tr><td><input checked="" type="checkbox"/></td><td>Strata Cloud Manager Pro, AI Model Scanning</td><td>D7732195</td><td>Available</td></tr><tr><td><input type="checkbox"/></td><td>Will Test</td><td></td><td></td></tr><tr><td><input type="checkbox"/></td><td>One more will test</td><td>D6072573</td><td>Available</td></tr><tr><td colspan="2">Strata Cloud Manager Pro, AI Model Scanning</td></tr></tbody></table>		<input type="checkbox"/>	Profile Name	Auth Code	Status	<input type="checkbox"/>	initial-ms-dp	D4073957	Available	<input type="checkbox"/>	Strata Cloud Manager Pro, AI Model Scanning			<input type="checkbox"/>	Josh test	D6717430	Available	<input checked="" type="checkbox"/>	Strata Cloud Manager Pro, AI Model Scanning	D7732195	Available	<input type="checkbox"/>	Will Test			<input type="checkbox"/>	One more will test	D6072573	Available	Strata Cloud Manager Pro, AI Model Scanning	
<input type="checkbox"/>	Profile Name	Auth Code	Status																												
<input type="checkbox"/>	initial-ms-dp	D4073957	Available																												
<input type="checkbox"/>	Strata Cloud Manager Pro, AI Model Scanning																														
<input type="checkbox"/>	Josh test	D6717430	Available																												
<input checked="" type="checkbox"/>	Strata Cloud Manager Pro, AI Model Scanning	D7732195	Available																												
<input type="checkbox"/>	Will Test																														
<input type="checkbox"/>	One more will test	D6072573	Available																												
Strata Cloud Manager Pro, AI Model Scanning																															
1 Row Selected																															
25 Rows ▾ Page 1 of 1 < >																															

**STEP 5 |** Agree to the terms, then click **Activate**.

 Deployment Profile activation may take up to two hours. Once activated, the AI model Security will be visible in the Strata Cloud Manager web interface (**Insights > Prisma AIRS > Model Security**).

You're redirected to the Tenant Management page which shows the instances for the tenant.

**STEP 6 |** Verify the TSG association in the **Hub**. To do this:

1. Navigate to **Common Services > Tenant Management**.
2. Select the tenant and switch to **Deployment Profiles**.
3. Confirm that the **Profile Association Status** is **Complete**.

### Edit or Update a Deployment Profile

You can edit an existing deployment profile on the CSP as long as AI Model Security supports your updates to the deployment profile. These changes may include changing the number of scans allocated per month.

**STEP 1 |** Log in to the Palo Alto [Customer Support Portal \(CSP\)](#).

**STEP 2 |** In the CSP dashboard, click on the three dots (ellipsis or overflow menu) next to your deployment profile.

**STEP 3 |** Edit, delete or deactivate the deployment profile.

## Configure Identity and Access Management

Where Can I Use This?	What Do I Need?
<ul style="list-style-type: none"> <li>Prisma AIRS (AI Model Security)</li> </ul>	<span style="color: orange;">□</span> Prisma AIRS AI Model Security License



If you have an active deployment profile but lack IAM permissions for AI Model Security, it may result in user authentication error. Hence it is important to have an active deployment profile and IAM permission to access AI Model Security.

- STEP 1 |** Navigate to Strata Cloud Manager Identity and Access Management (IAM) settings, **Common Services > Identity & Access**.
- STEP 2 |** Navigate to **Identity & Access > Access Management** and select the tenant from the tenants list. Verify the **Identity Information** for the selected tenant to ensure that the **Role** assigned is either **Superuser** for all apps and services, or a custom role with access to AI Model Security.

Applications	Role	Inherited From	Authorization Source
All Apps & Services	Superuser		Local

Assigned Roles:

- pairs-ms-sa-prod → Superuser
- model-security-sdk → Superuser
- test-custom-role → Superuser

## STEP 3 | (Optional) (Administrator Only) Create a custom role.

- To create a new custom role and assign it to their identity, select **Roles** and then **Custom Roles**.

Custom Role Name	Custom Role ID	Inherited From	Description	Actions
mssdk	mssdk1820847758		Minimal permissions for model scan SDK	

- Enable AI Model Security.

**Add Custom Role**

Name: mssdk  
Display Name (Optional):  
Description: Enter a brief description here

**Permissions**

Web UI	API	Description
<input type="checkbox"/> Advanced DNS Resolver All Functions (No access 9)	<input type="checkbox"/> Read Write	All Advanced DNS Resolver functions.
<input checked="" type="checkbox"/> AI Model Security (Read/write 2)	<input type="checkbox"/> Read Only	Access to all AI Model Security
<input type="checkbox"/> AI Red Teaming (No access 4)	<input type="checkbox"/> No Access	Full access to all AI Red Teaming services.
<input type="checkbox"/> AI Services (No access 16)	<input type="checkbox"/> Read Write	Full access to all AI services.
<input type="checkbox"/> Comm Services (No access 6)	<input type="checkbox"/> Read Only	Grants access to the Strata Logging Service Left Nav bar in SCM. Note: This does not include access to the SLS standalone or the Log Viewer.
<input type="checkbox"/> Agent Security (No access 6)	<input type="checkbox"/> No Access	Grant access (read/write) to Agent Security as a whole
<input type="checkbox"/> Next-Generation CASB (No access 75)	<input type="checkbox"/> Read Write	Full access to all Next-Generation CASB.
<input type="checkbox"/> Prisma Access & NGFW Configuration (No access 11)	<input type="checkbox"/> Read Only	Full access to all Prisma Access services.
<input type="checkbox"/> AIOps For NGFW	<input type="checkbox"/> No Access	Access AIOps for NGFW
<input type="checkbox"/> IoT Security	<input type="checkbox"/> Read Write	IoT Security

Legend: Read Write Read Only No Access

- If you need a service account with API access, select **API** and **Add permissions** and assign the minimum permissions.



To use SDK, you need at least the following permissions: **ai\_ms\_pypi\_auth**, **ai\_ms.scans**, and **ai\_ms.security\_groups**.

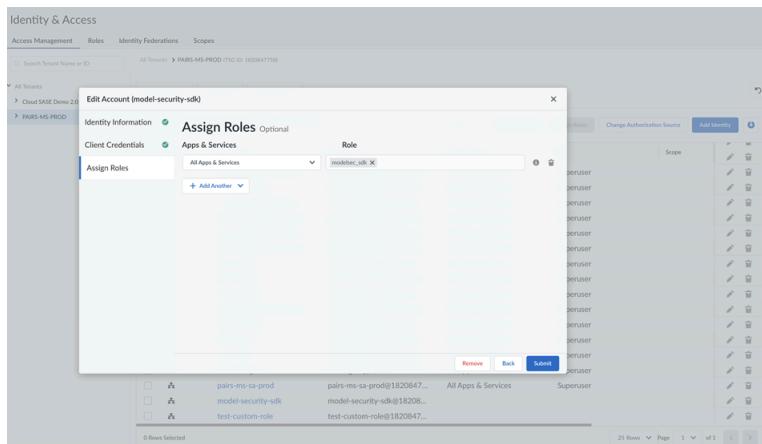
**Add Permissions**

Search for permission:

- ai\_ms\_scans**
- ai\_ms.security\_groups**
- ai\_ms.security\_rules**

Cancel Save

- Enter the **Name** and **Description** for the custom role and **Save** the changes.
- Select **Assign Roles** and verify if the custom role that you created appears in the **Apps & Services** list.



## Install AI Model Security

### Where Can I Use This?

- Prisma AIRS (AI Model Security)

### What Do I Need?

- Prisma AIRS AI Model Security License

To scan both internal and external models, you require either AI Model Security CLI or SDK. AI Model Security is available as a Python package that offers both a command-line interface and a Python SDK. Install the package using your preferred Python package manager.

### STEP 1 | Generate the pip index link.

Copy the script below and save it to your local environment (alternatively, you can create your own script using this as a reference).

```
#!/bin/bash
#
# Model Security Private PyPI Authentication Script
# Authenticates with SCM and retrieves PyPI repository URL
#
set -euo pipefail

# Check required environment variables
: "${MODEL_SECURITY_CLIENT_ID:?Error: MODEL_SECURITY_CLIENT_ID not set}"
: "${MODEL_SECURITY_CLIENT_SECRET:?Error: MODEL_SECURITY_CLIENT_SECRET not set}"
: "${TSG_ID:?Error: TSG_ID not set}"

# Set default endpoints
API_ENDPOINT="${MODEL_SECURITY_API_ENDPOINT:-https://api.sase.paloaltonetworks.com/aims}"
TOKEN_ENDPOINT="${MODEL_SECURITY_TOKEN_ENDPOINT:-https://auth.apps.paloaltonetworks.com/oauth2/access_token}"

# Get SCM access token
```

```

TOKEN_RESPONSE=$(curl -sf -X POST "$TOKEN_ENDPOINT" \
-H "Content-Type: application/x-www-form-urlencoded" \
-u "$MODEL_SECURITY_CLIENT_ID:$MODEL_SECURITY_CLIENT_SECRET" \
-d "grant_type=client_credentials&scope=tsg_id:$TSG_ID") || {
echo "Error: Failed to obtain SCM access token" >&2
exit 1
}

SCM_TOKEN=$(echo "$TOKEN_RESPONSE" | jq -r '.access_token')
if [[ -z "$SCM_TOKEN" || "$SCM_TOKEN" == "null" ]]; then
echo "Error: Failed to extract access token from response" >&2
exit 1
fi

# Get PyPI URL
PYPI_RESPONSE=$(curl -sf -X GET "$API_ENDPOINT/mgmt/v1/pypi/
authenticate" \
-H "Authorization: Bearer $SCM_TOKEN") || {
echo "Error: Failed to retrieve PyPI URL" >&2
exit 1
}

PYPI_URL=$(echo "$PYPI_RESPONSE" | jq -r '.url')
if [[ -z "$PYPI_URL" || "$PYPI_URL" == "null" ]]; then
echo "Error: Failed to extract PyPI URL from response" >&2
exit 1
fi

echo "$PYPI_URL"

```

**STEP 2 |** Set up authentication using environment variables.

After placing the script in an executable location, you'll need to set several environment variables before running it. Both the AI Model Security CLI and SDK require authentication credentials set as environment variables. The client automatically manages OAuth2 authentication with the provided credentials.

```

export MODEL_SECURITY_CLIENT_ID=<your-client-id>
export MODEL_SECURITY_CLIENT_SECRET=<your-client-secret>
export TSG_ID=<your-tsg-id>

```

**STEP 3 |** Install AI Model Security package (both SDK and CLI) with uv or pip.

1. Install AI Model Security package (both SDK and CLI) using uv, or.

```
uv add model-security-client --index $(/path/to/script.sh)
```

2. Install AI Model Security package (both SDK and CLI) using pip.

```
pip install model-security-client \
--extra-index-url <URL from Script>
```

**STEP 4 |** Initialize the AI Model Security Python SDK.

To use the Python SDK in your code, import and initialize the AI Model Security client.

```
from uuid import UUID
from model_security_client.api import ModelSecurityAPIClient

# Initialize the client
client = ModelSecurityAPIClient(
    base_url="https://api.sase.paloaltonetworks.com/aims"
)
```

The AI Model Security client uses the same environment variables for authentication as the CLI.

## Default Security Groups and Rules

Where Can I Use This?	What Do I Need?
<ul style="list-style-type: none"> <li>Prisma AIRS (Model Security)</li> </ul>	<input checked="" type="checkbox"/> Prisma AIRS Model Security License

AI Model Security leverages Security Groups to enable you to establish and configure your Model Security controls. Security Groups have a one-to-one relationship with Model Sources and this association is permanent once set.

Each Model Source provides specialized metadata that is used in rules and security validation processes.

You can create a dedicated Security Group for each of our supported Model Sources:

- Local Storage
- HuggingFace
- Azure Blob Storage
- Google Cloud Storage
- Amazon S3

Following are the default security group names:

Security Group Name	Model Source
Default LOCAL	Local Storage
Default HUGGING_FACE	HuggingFace
Default AZURE	Azure Blob Storage
Default GCS	Google Cloud Storage
Default S3	Amazon S3

## Scanning Models

Where Can I Use This?	What Do I Need?
<ul style="list-style-type: none"> <li>Prisma AIRS (AI Model Security)</li> </ul>	<input checked="" type="checkbox"/> Prisma AIRS AI Model Security License

Once your Security Group is configured, you can scan models through either the CLI or Python SDK. The process varies slightly depending on whether you're scanning HuggingFace models or local models.

### Scan a HuggingFace Model

To scan a model hosted on HuggingFace, provide the model URI and your Security Group UUID.

#### Scan using CLI

```
model-security scan \
--security-group-uuid "12345678-1234-1234-1234-123456789012" \
--model-uri "https://huggingface.co/microsoft/DialoGPT-medium"
```

#### Scan using Python SDK

```
from uuid import UUID

result = client.scan(
    security_group_uuid=UUID("12345678-1234-1234-1234-123456789012"),
    model_uri="https://huggingface.co/microsoft/DialoGPT-medium"
)

print(f"Scan completed: {result.eval_outcome}")
```

The AI Model Security automatically fetches the latest version from HuggingFace. To scan a specific version, include the version parameter.

#### Scan using CLI

```
model-security scan \
--security-group-uuid "12345678-1234-1234-1234-123456789012" \
--model-uri "https://huggingface.co/microsoft/DialoGPT-medium" \
--model-version "7b40bb0f92c45fef957d088000d8648e5c7fa33"
```

#### Scan using Python SDK

```
from uuid import UUID
result = client.scan(
    security_group_uuid=UUID("12345678-1234-1234-1234-123456789012"),
    model_uri="https://huggingface.co/microsoft/DialoGPT-medium",
    model_version="7b40bb0f92c45fef957d088000d8648e5c7fa33"
)
```

### Filter Files in HuggingFace Scans

Large HuggingFace repositories may contain files you don't need to scan. Use global patterns to include or exclude specific files.

#### Scan using CLI

```
model-security scan \
--security-group-uuid "12345678-1234-1234-1234-123456789012" \
--model-uri "https://huggingface.co/microsoft/DialoGPT-medium" \
--allow-patterns "*.bin" "*.json" \
--ignore-patterns "*.md" "*.txt"
```

#### Scan using Python SDK

```
from uuid import UUID
result = client.scan(
    security_group_uuid=UUID("12345678-1234-1234-1234-123456789012"),
    model_uri="https://huggingface.co/microsoft/DialoGPT-medium",
    allow_patterns=["*.bin", "*.json"],
    ignore_patterns=["*.md", "*.txt"]
)
```

### Scanning a Local Model

For models stored locally, specify the path to the model directory.

#### Scan using CLI

```
model-security scan \
--security-group-uuid "12345678-1234-1234-1234-123456789012" \
--model-path "path/to/local/model"
```

#### Scan using Python SDK

```
from uuid import UUID
result = client.scan(
    security_group_uuid=UUID("12345678-1234-1234-1234-123456789012"),
    model_path="path/to/local/model"
)
```

### Scanning a Model from Object Storage

To scan a model from object storage (such as S3), provide both the local path and the storage URI.

#### Scan using CLI

```
model-security scan \
--security-group-uuid "12345678-1234-1234-1234-123456789012" \
--model-path "path/to/local/model" \
--model-uri "s3://your-bucket/model-directory" \
```

```
--model-name "production-classifier" \
--model-author "ml-team" \
--model-version "v2.1"
```

### Scan using Python SDK

```
from uuid import UUID
result = client.scan(
    security_group_uuid=UUID("12345678-1234-1234-1234-123456789012"),
    model_path="path/to/local/model",
    model_uri="s3://your-bucket/model-directory",
    model_name="production-classifier",
    model_author="ml-team",
    model_version="v2.1"
)
```

The CLI shows scan results in real-time as they finish. Each scan tests the model against all active rules in your Security Group. The output shows whether the model passes or fails based on your rule configuration.

A model fails if any blocking rule detects a violation. Non-blocking rules record findings without preventing the model from being approved.

### Customize Scanning Models

You can configure scan execution and adjust result timeout settings.

### Customize Scan using CLI

```
model-security scan \
--security-group-uuid "12345678-1234-1234-1234-123456789012" \
--model-uri "https://huggingface.co/microsoft/DialoGPT-medium" \
--poll-interval-secs 10 \
--poll-timeout-secs 900 \
--block-on-errors
```

### Customize Scan using Python SDK

```
from uuid import UUID
result = client.scan(
    security_group_uuid=UUID("12345678-1234-1234-1234-123456789012"),
    model_uri="https://huggingface.co/microsoft/DialoGPT-medium",
    poll_interval_secs=10,
    poll_timeout_secs=900,
    scan_timeout_secs=900
)
```

Following are the configuration options to customize the scan for AI models.

Configuration Option	Description	Default Value
poll_interval_secs	Specify the frequency of scan status checks.	5 seconds
poll_timeout_secs	Specify the maximum wait time for scan completion.	600 seconds
scan_timeout_secs	(SDK only) Specify the timeout for local model scanning.	600 seconds
block_on_errors	(CLI only) CLI exits with an error code when scan errors occurs.	NA

## Viewing Scan Results

Where Can I Use This?	What Do I Need?
<ul style="list-style-type: none"> <li>Prisma AIRS (AI Model Security)</li> </ul>	<input checked="" type="checkbox"/> Prisma AIRS AI Model Security License

The `scan` command displays results directly, showing the overall verdict and key findings. For deeper analysis, you can retrieve detailed results using the scan ID.

After a model is scanned, it will either pass or fail the scan based on the security checks performed by AI Model Security. If the model passes the scan, it will be downloaded as usual. If the model fails the scan, it will fail to download and return a 403 error as well as a Universally Unique Identifier (UUID) that can be used to view the scan results.

### Retrieving a Specific Scan (CLI/SDK)

After a scan completes, note the scan ID from the output. Retrieve the full results at any time.

#### Retrieve Scan Results using CLI

```
model-security get-scan "87654321-4321-4321-4321-210987654321"
```

#### Retrieve Scan Results using Python SDK

```
from uuid import UUID
scan_id = UUID("87654321-4321-4321-4321-210987654321")
scan_result = client.get_scan(scan_id)

print(f"Scan Status: {scan_result.eval_outcome}")
print(f"Model URI: {scan_result.model_uri}")
print(f"Created: {scan_result.created_at}")
```

### View Scan Summary (CLI/SDK)

View a summary of recent scans to track your security assessments.

### View Scan Summary using CLI

```
model-security list-scans --limit 20
```

### View Scan Summary using Python SDK

```
scans = client.list_scans(limit=20)

for scan in scans.scans:
    print(f"Scan {scan.uuid}: {scan.eval_outcome} - {scan.model_uri}")
```

You can also filter scans by source type, evaluation outcome, or time range.

### Filter Scans by source type, evaluation outcome, or time range using CLI

```
model-security list-scans \
--source-types "HUGGING_FACE" "S3" \
--eval-outcomes "ALLOWED" "BLOCKED" \
--start-time "2025-01-01T00:00:00" \
--limit 50
```

### Filter Scans by source type, evaluation outcome, or time range using Python SDK

```
from datetime import datetime, timezone
from airs_schemas.constants import SourceType, EvalOutcome

scans = client.list_scans(
    source_types=[SourceType.HUGGING_FACE, SourceType.S3],
    eval_outcomes=[EvalOutcome.ALLOWED, EvalOutcome.BLOCKED],
    start_time=datetime(2025, 1, 1, tzinfo=timezone.utc),
    limit=50
)
```

### View Scan Results using Strata Cloud Manager (Web Interface)

Although the CLI delivers complete results, the AI Model Security web interface provides additional tools for analyzing scan findings.

1. Log in to [Strata Cloud Manager](#).

2. Navigate to the Insights > Prisma AIRS > Model Security > Scans

The screenshot shows the 'Scans' page in the Prisma AIRS interface. At the top, there's a breadcrumb navigation: 'Prisma AIRS / Model Security / Scans'. Below it is a title 'Scans' and a subtitle 'This page contains a list of scans that have been performed.' There are several search and filter options: a search bar ('Search scans...'), a date range selector ('Scans from: Past 30 Days'), dropdown menus for 'Evaluation Results', 'Model Security Group', and 'Source'. A large heading '16998 Scans' is followed by a table with columns: 'Scan Request ID', 'Time of Request', 'Rule Outcome', 'Model URI', and 'Source'. The table lists 10 scan entries, each with a status icon (green checkmark or red warning), a unique ID, the request time (Oct 14, 2025 12:31pm), the rule outcome (e.g., 0/7 rules passed), the model URI (e.g., ...sg-model-ec77064c), and the source (Local). The last row shows 'Page Size: 10' and '1 to 10 of 16998'.

Scan Request ID	Time of Request	Rule Outcome	Model URI	Source
29c3...5734	Oct 14, 2025 12:31pm	0/7	...sg-model-ec77064c	Local
f3b2...964a	Oct 14, 2025 12:31pm	0/7	...sg-model-ec77064c	Local
9a3f...a5c5	Oct 14, 2025 12:31pm	0/7	...iles-test-25b5430c	Local
e63e...e090	Oct 14, 2025 12:31pm	0/7	...eval-test-23882222	Local
3ab2...17eb	Oct 14, 2025 12:31pm	0/7	...bbbbbe1b/model™.pkl	Local
e868...b850	Oct 14, 2025 12:31pm	0/7	...n-details-f79e80a1	Local
c441...509b	Oct 14, 2025 12:31pm	0/7	...t/minimal-2b9fde73	Local
fb56...6870	Oct 14, 2025 12:31pm	1/6	...d-unsafe/model.pkl	Local
5164...8c1e	Oct 14, 2025 12:31pm	0/7	...nicode-元氣-e421c0a3	Local
4ce0...c3b4	Oct 14, 2025 12:31pm	0/7	...c5d-safe/model.pkl	Local

3. Locate your scan by ID or filter the scan list. Review detailed findings for each rule evaluation. Export results or share them with your team.

Following is an example scan result that is **Allowed**. After locating the specific scan, select **Overview** to review the evaluation details.

The screenshot displays the Prisma AIRS / Model Security interface. On the left, a vertical sidebar lists various navigation items with corresponding icons:

- INSIGHTS**: Includes Activity Insights, Prisma AIRS, Model Scanning (highlighted with a green bar), Model Security Groups, SECURITY, Executive Summary, Threat Search, POSTURE, CDSS Adoption, Compliance Summary, Security Posture Insights, PAN-OS CVEs, On Demand BPA, Feature Adoption, Zero Trust Posture Center, Best Practices, and OPERATIONAL.
- Scans**: Shows a summary of 16998 Scans, a search bar for scans, and a list of recent scan results with status indicators (green checkmarks) and IDs (e.g., 29c3...5, f3b2...9, 9a3f...a5, e63e...e, 3ab2...1, e868...b, c441...5, fb56...6).
- AI Model Security**: Includes a star icon and a bell icon.

At the bottom, there is a footer with the text "©2025 Palo Alto Networks, Inc." and a shield icon.

Following is an example scan result that is **Blocked**. After locating the specific scan, select **Overview** to review the evaluation details.

The screenshot displays the Prisma AIRS / Model Security dashboard. On the left, a vertical sidebar lists navigation icons: a yellow hexagon (Insights), a magnifying glass (Search), a target (Model Scanning), a chart (Model Security Groups), a gear (Security), a bell (Threat Search), binoculars (Posture), a magnifying glass (CDSS Adoption), a gear (Compliance Summary), and a gear (Security Posture Insights). The main content area is divided into three sections: **Insights**, **Scans**, and **OPERATIONAL**.

**Insights** section:

- Activity Insights
- Prisma AIRS
- Model Scanning
- Model Security Groups

**Scans** section:

- 16998 Scans
- Scan Results

**OPERATIONAL** section:

- 29c3...5
- f3b2...90
- 9a3f...a5
- e63e...e
- 3ab2...1
- e868...b
- c441...5
- fb56...68

At the bottom, the footer reads "AI Model Security" and "©2025 Palo Alto Networks, Inc." with page number 29.

Select **JSON** of the specific Scan, to review the details of the scan results in JSON format.

The screenshot displays the Palo Alto Networks AI Model Security interface. On the left, a vertical sidebar contains icons for various features: Insights (yellow hexagon), Activity Insights (magnifying glass), Prisma AIRS (target), Model Scanning (yellow arrow), Model Security Groups (grid), SECURITY (blue shield), Executive Summary (bell), Threat Search (binoculars), POSTURE (blue shield), CDSS Adoption (magnifying glass), Compliance Summary (gear), Security Posture Insights (blue shield), PAN-OS CVEs (blue shield), On Demand BPA (blue shield), Feature Adoption (blue shield), Zero Trust Posture Center (blue shield), Best Practices (blue shield), and OPERATIONAL (star).

**Insights**

- Activity Insights
- Prisma AIRS
- Model Scanning**
- Model Security Groups

**SECURITY**

- Executive Summary
- Threat Search

**POSTURE**

- CDSS Adoption
- Compliance Summary
- Security Posture Insights

**PAN-OS CVEs**

- On Demand BPA
- Feature Adoption
- Zero Trust Posture Center

**Best Practices**

**OPERATIONAL**

**Scans**

This page contains a list of scans.

Search scans...

**16998 Scans**

Scan ID	Status
29c3...57	✓
f3b2...96	✓
9a3f...a5	✓
e63e...e0	✓
3ab2...17	✓
e868...b8	✓
c441...50	✓
fb56...68	✗

The Strata Cloud Manager displays rule violations visually, identifies the specific files or model components that caused findings, and offers remediation guidance.

## Understanding Model Security Scan Results

Where Can I Use This?	What Do I Need?
<ul style="list-style-type: none"><li>Prisma AIRS (Model Security)</li></ul>	<input checked="" type="checkbox"/> Prisma AIRS Model Security License

Once a scan finishes, you can retrieve the results through either the SDK or the Strata Cloud Manager web interface. The primary check involves determining whether your scanned model is permitted or restricted.

Execute the following command to retrieve the scan result using the SDK:

```
scan_result.eval_outcome
```

To view the scan results in the Strata Cloud Manager, navigate to **Insights > Prisma AIRS > Model Security > Scans** and locate your recently scanned model.

Scan Request ID	Time of Request	Rule Outcome	Model URI	Source	Actions
e653...8e73	Oct 21, 2025 9:48am	(0) 2 ✓ 5	..shared-unsafe-model	S3	[Details]
1f97...0fa6	Oct 21, 2025 9:48am	(0) 0 ✓ 7	..shared-safe-model	S3	[Details]
8e07...4442	Oct 21, 2025 9:48am	(0) 2 ✓ 5	..shared-unsafe-model	S3	[Details]
f02c...82eb	Oct 21, 2025 9:48am	(0) 0 ✓ 7	..shared-safe-model	S3	[Details]
6387...1039	Oct 21, 2025 9:32am	(0) 2 ✓ 5	..shared-unsafe-model	S3	[Details]
731b...00b1	Oct 21, 2025 9:32am	(0) 0 ✓ 7	..shared-safe-model	S3	[Details]
7415...a5ea	Oct 21, 2025 9:32am	(0) 2 ✓ 5	..shared-unsafe-model	S3	[Details]
200d...f1ae	Oct 21, 2025 9:32am	(0) 0 ✓ 7	..shared-safe-model	S3	[Details]
a227...bc8a	Oct 21, 2025 9:18am	(0) 2 ✓ 5	..shared-unsafe-model	S3	[Details]
108c...a23d	Oct 21, 2025 9:18am	(0) 0 ✓ 7	..shared-safe-model	S3	[Details]

A red shield indicates the model was blocked, while a green shield shows the model is allowed.

Select the **Scan Request ID** to open a detailed flyover with comprehensive scan information:

This example shows a result of a scan that is blocked due to failure in compliance to the rules: the use of non-approved file formats and detection of code execution in `unsafe_model.pkl` that runs during model loading, indicating a potential deserialization attack.

To view the complete JSON response that your SDK would receive, select **JSON**:

```

{
  "id": "e853_8e73",
  "request_id": "e853_8e73",
  "time": "2025-10-21T09:48:00Z",
  "url": "s3://test-bucket/shared-unsafe-model",
  "status": "Blocked",
  "evaluations": [
    {
      "name": "Known Framework Operators Check",
      "status": "Passed"
    },
    {
      "name": "Model Architecture Backdoor Check",
      "status": "Passed"
    },
    {
      "name": "Suspicious Model Components Check",
      "status": "Passed"
    },
    {
      "name": "Stored In Approved Location",
      "status": "Passed"
    },
    {
      "name": "Stored In Approved File Format",
      "status": "Failed"
    }
  ],
  "files": [
    {
      "path": "unsafe_model.pkl",
      "status": "Failed"
    }
  ]
}

```

# Understanding Model Security Rules

Where Can I Use This?	What Do I Need?
<ul style="list-style-type: none"> <li>Prisma AIRS (Model Security)</li> </ul>	<input checked="" type="checkbox"/> Prisma AIRS Model Security License

Model security rules serve as the central mechanism for securing model access. Model security scanning covers several key areas: threat and metadata. To learn more about threat categories and risk mitigation strategies for AI and machine learning systems, refer [Model Threats](#).

Following are the three major threat categories:

Threat	Description
<b>Deserialization threats</b>	Issues that arise when you load a model into memory.
<b>Backdoor Threats</b>	Issues that arise when a model was specifically designed to support alternative or hidden paths in its behavior.
<b>Runtime Threats</b>	Issues that arise when you use a model to perform inference.

Model security also examines specific metadata fields in models to address security considerations, such as verifying a model's license to ensure it's appropriate for commercial use. Model security rules enable you to validate the following concerns:

Metadata in Models	Security Rule Validation
<b>Open Source License</b>	Ensures that the model you are using is licensed for your use case.
<b>Model Format</b>	Ensures that the model you are using is in a format that is supported by your environment.
<b>Model Location</b>	Ensures that the model you are using is hosted in a location that is secure and trusted.
<b>Verified Organizations in Hugging Face</b>	Ensures that the model you are using is from a trusted source.
<b>Model is Blocked</b>	Overrides all other checks to ensure that a model is blocked no matter what.

## Secure Your AI Models with AI Model Security

Our managed rules integrate all of these checks to provide comprehensive security for model usage. To get started, you can explore your defaults in your Model Security at **Insights > Prisma AIRS > Model Security > Model Security Groups**.

Prisma AIRS / Model Security / Model Security Groups

## Model Security Groups

This page contains a collection of rules that define the security posture for models originating from a specific source.

Search security groups...    Source    Enabled Rules

20 Model Security Groups			
Group Name	Associated Source	Last Edited	Description
SDK Local Scanning Policy	Local Disk	Oct 14, 2025 12:48pm	Security group for local SDK models
shared-958e714a	Local Disk	Oct 13, 2025 10:46pm	Shared SG for all tests in class
My First Security Group	HuggingFace	Oct 13, 2025 7:31pm	Getting started with model security
HuggingFace Verified Publishers	HuggingFace	Oct 13, 2025 11:37am	Security policy for HuggingFace publishers
Permissive Dev Policy	S3	Oct 13, 2025 11:37am	Lenient security for S3 development
Critical Threats Only	Azure	Oct 13, 2025 11:37am	Block only critical severity threats
SOC2 Compliance	Google Cloud Storage	Oct 13, 2025 11:37am	Meets SOC2 audit requirements
Production with Legacy Models	S3	Oct 13, 2025 11:37am	Production policy with pickle exports
Zero Trust Production	Local Disk	Oct 13, 2025 11:37am	Strictest security for production
shared-2fe5f33f	Local Disk	Oct 13, 2025 6:46am	Shared SG for all tests in class

Page Size: 10    1 to 10 of 20

## How Scans, Security Groups, and Rules Work Together

With Model Security, you initiate a model scan and associate it with a security group. Model Security evaluates the model against all enabled rules within that security group to assess whether it satisfies your security requirements. This assessment relies on the results from enabled rules in the security group and the group's configured threshold.

When you scan a model as follows:

```
# Import the Model Security SDK/Client
from model_security_client.api import ModelSecurityAPIClient
# Load your scanner URL
```

```

scanner_url = os.environ["MODEL_SECURITY_TOKEN_ENDPOINT"]
# Define your model's URI
model_uri = "s3://demo-models/unsafe_model.pkl"
# Set your security group's UUID
security_group_uuid = "6e2ccc3a-db57-4901-a944-ce65e064a3f1"
# Create a Model Security Client
guardian = ModelSecurityAPIClient(base_url=scanner_url)
# Scan your model
response = client.scan(model_uri=model_uri,
    security_group_uuid=security_group_uuid)

```

The response will appear as follows (showing the security rules results):

```

{
  "http_status_code": 200,
  "scan_status_json": {
    "aggregate_eval_outcome": "FAIL",
    "aggregate_eval_summary": {
      "critical_count": 1,
      "high_count": 0,
      "low_count": 0,
      "medium_count": 1
    },
    "violations": [
      {
        "issue": "Model file 'ykilcher_totally-harmless-model/retr0reg.gguf' is stored in an unapproved format: gguf",
        "threat": "UNAPPROVED_FORMATS",
        "operator": null,
        "module": null,
        "file": "ykilcher_totally-harmless-model/retr0reg.gguf",
        "hash": "f59ad9c65c5a74b0627eb6ca5c066b02f4a76fe6",
        "threat_description": "Model is stored in a format that your Security Group does not allow",
        "policy_name": "Stored In Approved File Format",
        "policy_instance_uuid": "34ef1ddc-0b7a-45b8-a84a-c96b1d8383d0",
        "remediation": {
          "steps": [
            "Store the model in a format approved by your organization"
          ]
        }
      },
      {
        "issue": "The model will execute remote code since it contains operator `__class__` in Jinja template.",
        "threat": "PAIT-GGUF-101",
        "operator": "__class__",
        "module": null,
        "file": "ykilcher_totally-harmless-model/retr0reg.gguf",
        "hash": "f59ad9c65c5a74b0627eb6ca5c066b02f4a76fe6",
        "threat_description": "GGUF Model Template Containing Arbitrary Code Execution Detected",
        "policy_name": "Load Time Code Execution Check",
      }
    ]
  }
}

```

```

    "policy_instance_uuid": "09780b9f-c4f7-4e0b-
ad21-7ff779472283",
    "remediation": {
        "steps": [
            "Use model formats that disallow arbitrary code
execution"
        ]
    }
}

```

Note the FAIL status in the aggregate\_eval\_outcome field. This indicates the model did not pass the scan because security rule failures surpassed your security group's threshold, with the violations field providing details about which rules were breached.

Each model security rule contains the following fields.

Rule Field	Description
<b>Rule Name</b>	Specifies the name of the security rule.
<b>Rule Description</b>	Specifies the description of the security rule.
<b>Compatible Sources</b>	Specifies the model source types that this rule is compatible with.
<b>Status</b>	Specifies the status of the security rule, either Enabled or Disabled. This can be set globally, or at the security group level.

## Connecting Scans to Rules

When scanning a model, model security first identifies its format through introspection. After determining the model type, model security maps it to the taxonomy of model vulnerability threats and coordinate the specific deeper scans required for that model.

A series of specific threats like **Arbitrary Code Execution At Runtime** are grouped together when they are reported and are shown in a specific rule. This allows you to block specific types of threats without having to manage the complexity of all the various formats and permutations.



You can also configure common rules for all models regardless of their source type.

## Security Rule Checks

Rule Name	Status	Description	Example
<b>Runtime Code Execution Check</b>	<b>Status:</b> Enabled	<p>This rule detects Arbitrary Code Execution that can occur during model inference through various methods.</p> <p>These attacks mean the model will execute code without your knowledge during use, making this a <b>Critical</b> issue to block. Learn more about this threat type here: <a href="#">Runtime Threats</a>.</p>	<ul style="list-style-type: none"> <li>• Keras Model Lambda Layer Arbitrary Code Execution Detected at Model Run Time</li> <li>• TensorFlow SavedModel Contains Arbitrary Code Execution at Model Run Time</li> </ul>
<b>Known Framework Operators Check</b>	<b>Status:</b> Enabled	Machine learning model formats often include built-in operators to support common data science tasks during model operation. Some frameworks allow custom operator definitions, which poses risks when executing unknown third-party code.	When TensorFlow SavedModel Contains Unknown Operators is detected, it indicates that the model creator is

Rule Name	Status	Description	Example
			<p>using non-standard tooling approaches.</p> <p>For more information refer, <a href="#">SavedModel Contains Unknown Operators</a>.</p>
<b>Model Architecture Backdoor Check</b>	Status: Enabled	<p>A model's behavior can contain a backdoor embedded in its architecture, specifically through a parallel data flow path within the model. For most inputs, the model operates as expected, but for certain inputs containing a trigger, the backdoor activates and effectively alters the model's behavior.</p>	<p>When ONNX Model Contains Architectural Backdoor is detected, it warns you that a model has at least one non-standard path requiring further investigation.</p> <p>For more information refer, <a href="#">ONNX Model Contains Architectural Backdoor</a>.</p>
<b>Load Time Code Execution Check</b>	Status: Enabled	<p>This rule is similar to the runtime attacks, but these attacks execute immediately upon model loading. For example, the below python snippet is sufficient to trigger the exploit without your knowledge:</p> <pre data-bbox="621 1638 1225 1807"> with     open('path_to_your_model.pkl',         'rb') as file:         model =             pickle.load(file) </pre>	<p>When Pickle, Keras Lambda Layers, PyTorch models, and more are all vulnerable to this threat is detected in a</p>

Rule Name	Status	Description	Example
		<p>model, there's no reliable method to eliminate the threat.</p> <p>However, you can investigate the model's history or, if it's your own model, examine the build process to identify the source of the malicious code.</p> <p>For more information refer, <a href="#">PyTorch Model Arbitrary Code Execution</a></p>	
<b>Suspicious Model Components Check</b>	<b>Status:</b> Enabled	<p>Not all rules target specific threats; some, like this one, assess a model's potential for future exploitation. This check identifies components within the model that could enable malicious code execution later.</p> <p>Example:</p> <ul style="list-style-type: none"> <li>Remote code execution being called by fetching external data from Protobuf files or others over the internet.</li> </ul> <p>A violation here should prompt you to be cautious and to evaluate all relevant components around the model before making a decision on whether or not the model is safe for use.</p> <p>More information can be found here: <a href="#">Keras Lambda Layers Can Execute Code</a></p>	<p>Remote code execution being called by fetching external data from Protobuf files or others over the internet.</p> <p>This violation prompts caution and thorough evaluation of all relevant model</p>

Rule Name	Status	Description	Example
			components before determining whether the model is safe for use.  For more information refer, <a href="#">Keras Lambda Layers Can Execute Code</a>
Stored In Approved File Format	Status: Enabled	<p>This rule verifies whether the model is stored in a format that you've approved within the rule's list.</p> <p>We recommend enabling these formats by default:</p> <ul style="list-style-type: none"> <li>• Safetensors</li> <li>• JSON</li> </ul> <p>By default Model Security reports that pickle and keras_metadata are not approved formats.</p>	—

## Hugging Face Model Rules

Lastly, there are rules specifically scoped to Hugging Face models. These rules target the particular metadata you control or that is consistently provided by Hugging Face.

Hugging Face Model Rule	Status	Description
License Exists	Status: Enabled	<p>The simplest rule in the application. This rule checks to see if the model has a license associated with it.</p> <p>If no license is present, then this rule will fail.</p>
License Is Valid For Use	Status: Enabled	This rule gives you more control for the models that your organization will run or test. As a commercial entity, you may not want to run models with non-permissive open-source licenses like GPLv3 or others.

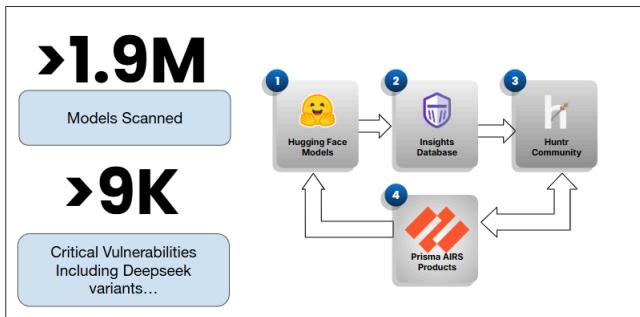
Hugging Face Model Rule	Status	Description
		<p>Adding licenses to the rule will expand the list of licenses that are allowed for use.</p> <p>A reasonable set of defaults are:</p> <ul style="list-style-type: none"> <li>• apache-2.0</li> <li>• mit</li> <li>• bsd-3.0</li> </ul> <p>You can find a full list of license options here: <a href="#">Hugging Face Licenses</a></p> <p>Note we use the <code>License identifier</code> field to map to the license in the model metadata.</p>
Model Is Blocked	<b>Status: Enabled</b>	<p>Sometimes you just cannot tolerate a model for whatever reason. This rule allows you to block a specific model from being used.</p> <p>The format for blocking a model from Hugging Face relies on the Organization and Model Name.</p> <p>For example <code>opendiffusion/sentiment - check</code> can be entered and it would block the model <code>sentiment - check</code> from <code>opendiffusion</code>.</p> <p> When a model is both blocked and allowed simultaneously, the block takes precedence.</p>
Organization Verified By Hugging Face	<b>Status: Enabled</b>	<p>Hugging Face is an excellent site for the latest models from all over the world, giving you access to cutting edge research.</p> <p>You may want to restrict organizations from providing unverified models from Hugging Face. This prevents accidentally running models from deceptively similar sources like <code>facebook - llama</code> instead of the legitimate <code>meta - llama</code> (where legitimate <code>meta - llama</code> is the correct model).</p> <p>This rule simply checks that Hugging Face has verified the organization, if that passes</p>

Hugging Face Model Rule	Status	Description
Organization Is Blocked	Status: Enabled	<p>the models from the organization will pass this check.</p> <p>If you find a particular organization that just delivers problematic models or for any other reason, you'd like to block them, this is your rule.</p> <p>Enter the organization name into the rule and all models from that organization will be blocked.</p> <p>For example <code>facebook-llama</code> would block ALL of the models provided by that organization.</p>

# Reporting New Threats

Where Can I Use This?	What Do I Need?
<ul style="list-style-type: none"> <li>Prisma AIRS (Model Security)</li> </ul>	<input type="checkbox"/> Prisma AIRS Model Security License

Our comprehensive vulnerability detection begins with the industry's most extensive dataset of model security issues. This intelligence is driven by our Huntr community—a network of over 16,000 security researchers who continuously discover novel attack vectors through our bug bounty program and strategic partnership with [Hugging Face](#).



You can explore our most recent discoveries in the [Insights DB](#), where we publish vulnerability assessments from scanning every public model in the Hugging Face repository. The platform includes mechanisms for security researchers and ML practitioners to challenge or validate our findings.

InsightsDB classifies and [lists the models](#) as **Safe**, **Unsafe**, and **Suspicious**. You can dispute a finding in [Insights DB](#) by selecting the specific model from the list and **Report an issue**:

- **Report your finding** if you've found a new threat.
- **Report an incorrect threat** if you disagree with our findings

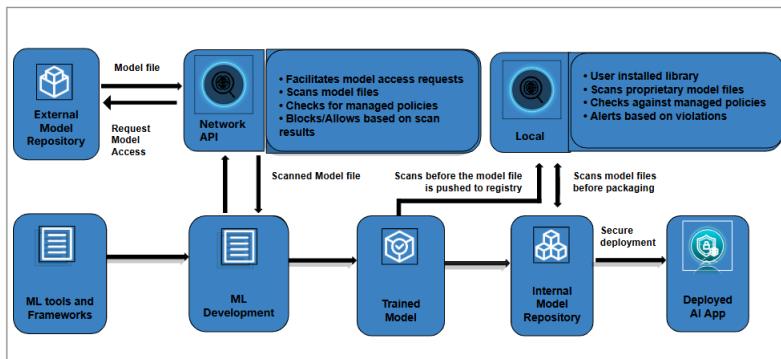
This continuous intelligence stream enables us to identify a broad spectrum of vulnerabilities in your models, including many AI-specific security issues that fall outside the scope of traditional software vulnerabilities—threats that typically won't appear in the National Vulnerability Database (NVD) or conventional security feeds.

Model Security operates as a centralized repository of models from third-party sources like Hugging Face.

Additionally, we provide an on-premises scanning solution that deploys directly within your infrastructure, enabling you to assess your proprietary models locally without data transmission to our systems—ensuring complete privacy and security of your assets.

## Secure Your AI Models with AI Model Security

---



# Supported Model Formats

Where Can I Use This?	What Do I Need?
<ul style="list-style-type: none"> <li>Prisma AIRS (Model Security)</li> </ul>	<input checked="" type="checkbox"/> Prisma AIRS Model Security License

AI Model Security checks are supported on the following formats:

1. **CNTK Models:** Models saved in Microsoft Cognitive Toolkit format.
2. **Flax Models:** Models created with Flax, a neural network library for JAX.
3. **GGUF Models:** General-purpose model format using GGUF.
4. **Keras Models:**
  - **Legacy Keras Models:** Older Keras models, often saved with HDF5.
  - **Keras 3 Models:** Newer Keras models using the latest Keras version.
  - **Keras Pickle Models:** Keras models saved with Python pickle.
  - **Keras H5 Models:** Models in HDF5 format, compatible with legacy and newer versions.
  - **Keras Weights:** Separate files storing only model weights.
  - **Keras JSON:** Models saved in JSON format for architecture storage.
  - **Keras Metadata:** Auxiliary files that store metadata for Keras models.
5. **KModel:** KModel files specific to Keras.
6. **LightGBM Models:** Gradient boosting models using LightGBM.
7. **MS Lite Models:** Microsoft Lite format for lightweight models.
8. **MXNet Models:** Models saved in Apache MXNet format.
9. **Numpy Models:**
  - **Numpy Array Files:** Arrays saved in .npy format.
  - **Numpy Zip Files:** Arrays compressed in .npz format.
  - **Numpy Pickle Files:** Arrays serialized with pickle.
10. **OM Models:** Models in Huawei Ascend's OM format.
11. **ONNX Models:** Models saved in Open Neural Network Exchange format.
12. **OpenVINO Models:**
  - **OpenVINO Binary Files:** Compiled binary files for OpenVINO.
  - **OpenVINO XML Files:** XML files storing OpenVINO model metadata.
13. **Pickle Files:** Models serialized using Python's pickle.
14. **PyTorch Models:**
  - **Various PyTorch Versions:** Models saved with different PyTorch versions.
  - **TorchScript:** PyTorch's format for serializing models.
  - **PyTorch Archives:** Archived files containing serialized models.

**15.RKNN Models:** Models saved in Rockchip Neural Network (RKNN) format.

**16.Safetensors:**

- **Safetensors Models:** Models saved using safetensors format for secure tensor storage.
- **Safetensors Index:** Index files for safetensors.

**17.SKLearn Models:** Scikit-learn models serialized for deployment.

**18.TensorRT Models:** NVIDIA's TensorRT models optimized for inference.

**19.TensorFlow Models:**

- **SavedModel:** TensorFlow's standard saved model format.
- **TFHub:** Models from TensorFlow Hub.
- **MetaGraph:** TensorFlow's MetaGraph format for exporting graphs.
- **TFLite:** Lightweight format for mobile and embedded devices.
- **TFJS:** TensorFlow.js format for models running in the browser.

**20.Torch Models:** General format for PyTorch models.

**21JSON Files:** JSON-based configurations or model descriptions.

