PyTorch 2: Faster machine learning through dynamic Python bytecode transformation



Jason Ansel

Research Scientist, Meta



The Great ML Framework Debate

Eager Mode

- Preferred by users
- Easier to use programming model
- Easy to debug
- PyTorch is a primarily an eager mode framework

Graph Mode

- Preferred by backends and framework builders
- Easier to optimize with a compiler
- Easier to do automated transformations

PyTorch's many attempts at graph modes

torch.jit.trace

- Record + replay
- Unsound
- Can give incorrect results because it ignores Python part of program

torch.jit.script

- AOT parses Python into graph format
- Only works on ~45% of real world models
- High effort to "TorchScript" models

Lazy Tensors (Torch XLA)

- Graph capture through deferred execution
- High overheads
- Performance cliffs

PyTorch Models Are Not Static Graphs

Due to history of being an eager model framework, PyTorch users have written models in ways where whole program graphs are impossible

In our benchmark suite 20% of models, do one (or more) of:

- Convert tensors native Python types (x.item(), x.tolist(), int(x), etc)
- Use other frameworks (numpy/xarray/etc) for part of their model
- Data dependent Python control flow or other dynamism
- Exceptions, closures, generators, classes, etc

All of these violate the assumptions of most graph mode backends.

PyTorch Usability/Performance Tradeoff



Performance

INTRODUCING

cmodel = torch.compile(model)

TorchDynamo: Out-of-the-box graph capture for PyTorch

torch.compile() with a user-defined backend

from typing import List
import torch

```
@torch.compile(backend=my_compiler)
def toy_example(a, b):
    x = a / (torch.abs(a) + 1)
    if b.sum() < 0:
        b = b * -1
    return x * b</pre>
```

```
for _ in range(100):
    toy_example(torch.randn(10), torch.randn(10))
```

Output:

my_compiler()	called wi	th FX graph:	
opcode	name	target	args
placeholder	а	a	()
placeholder	b	b	()
call_function	abs_1	torch.abs	(a,)
call_function	add	operator.add	(abs_1, 1)
call_function	truediv	operator.truediv	(a, add)
call_method	sum_1	sum	(b,)
call_function	lt	operator.lt	(sum_1, 0)
output	output	output	((truediv, lt),)

my_compiler() called with FX graph:

opcode	name	target	args
placeholder	b	b	()
placeholder	х	Х	()
call_function	mul	operator.mul	(b, -1)
call_function	mul_1	operator.mul	(x, mul)
output	output	output	((mul_1,),)

my compiler() called with FX graph:

opcode	name	target	args
placeholder	b	b	()
placeholder	х	Х	()
call_function	mul	operator.mul	(x, b)
output	output	output	((mul,),)



TorchInductor: A PyTorch Native Compiler

TORCHINDUCTOR PRINCIPLES

PyTorch Native

Similar abstractions to PyTorch eager to allow support for nearly all of PyTorch, with a thin translation layer.

Python First

A pure python compiler makes TorchInductor easy to understand and hackable by users. Generates Triton and C++.

Breadth First

Early focus on supporting a wide variety of operators, hardware, and optimization. A general purpose compiler, that can scale.

TORCHINDUCTOR TECHNOLOGIES

Define-By-Run Loop-Level IR

Direct use of Python functions in IR definitions allows for rapidly defining lowering with little boilerplate.

Dynamic Shapes & Strides

Uses SymPy to reason about shapes, indexing, and managing guards. Symbolic shapes from the ground up. Reuse State-Of-The-Art Languages

Generates output code in languages popular for writing handwritten kernels:

- Triton for GPUs
- C++/OpenMP for CPUs

What is Triton?

A new programming language for highly performant GPU kernels

- Higher level than CUDA
- Lower level than preexisting DSLs
- Allows non-experts to write fast custom kernels

Users define tensors (i.e., blocks of data) in SRAM, and modify them using torch-like operators



LOW-LEVEL MEMORY CONTROL

Like in Numba, kernels are defined in Python using the triton.jit decorator Users can construct tensors of pointers and dereference them element-wise



Blocked program representation allows the Triton compiler to generate extremely efficient code https://triton-lang.org https://github.com/openai/triton by Philippe Tillet @ OpenAl

Triton: an intermediate language and compiler for tiled neural network computations

Philippe Tillet, H. T. Kung, David Cox

In Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages (MAPL 2019)

https://doi.org/10.1145/3315508.3329973

TorchInductor Overview

AotAutograd	Inductor Graph Lowerings	Inductor Scheduling	Wrapper Codegen
Decomposes into smaller operator set	Remove views, broadcasting, and simplify indexing	Horizontal / vertical fusion decisions	Outer code that calls kernels and allocates memory
Capture forwards + backwards	Rematerialize vs reuse decisions	Reduction fusions Tiling	(Replaces interpreter)
Some inductor specific decomps included in this step	Layout tuning and optimization Loop order	Memory planning and buffer reuse	
		In-place memory buffers	Backend Codegen
		Autotuning / kernel selection	Triton
			C++
			Halide (new)



Input Code	ATen FX Graph	Define-by-run IR	Scheduling/Fusion	Output Triton	Output Wrapper
------------	---------------	------------------	-------------------	---------------	----------------

import torch

Run with: TORCH_COMPILE_DEBUG=1 python inductor_demo.py

```
@torch.compile(dynamic=True)
def toy_example(x):
    y = x.sin()
    z = y.cos()
    return y, z
```

toy_example(torch.randn([8192, 1024], device="cuda"))

Input Code	ATen FX Graph	Define-by-run IR	Scheduling/Fusion	Output Triton	Output Wrapper
------------	---------------	------------------	-------------------	---------------	----------------

def forward(self, arg0_1: f32[s0, s1]):
 # File: inductor_demo.py:6, code: y = x.sin()
 sin: f32[s0, s1] = torch.ops.aten.sin.default(arg0 1)

File: inductor_demo.py:7, code: z = y.cos()
cos: f32[s0, s1] = torch.ops.aten.cos.default(sin)
return (sin, cos)

TorchInd	TorchInductor Example								
Input Code	ATen FX Graph	Define-by-run IR	Scheduling/F	usion	Output Triton	Output Wrapper			
<pre>def inner_ i0, i1 = tmp0 = c tmp1 = c return t</pre>	_fn_buf0(index = index ops.load(arg0_ ops.sin(tmp0) tmp1	(): _1, i1 + i0	* s1)	buf0_ nam lay dat	ir = TensorBox(Stor e='buf0', out=FixedLayout('cu siz a=Pointwise(inner_ ranges=	rageBox(ComputedBuffer(uda', torch.float32, ze=[s0, s1], stride=[s1, fn=inner_fn_buf0, =[s0, s1],))))	1])		
<pre>def inner_fn_bufl(index): i0, i1 = index tmp0 = ops.load(buf0, i1 + i0 * s1) tmp1 = ops.cos(tmp0) return tmp1</pre>			s1)	buf1_ nam lay dat	<pre>ir = TensorBox(Stor e='buf1', out=FixedLayout('cu si: a=Pointwise(inner_ ranges=</pre>	rageBox(ComputedBuffer(uda', torch.float32, ze=[s0, s1], stride=[s1, fn=inner_fn_buf1, =[s0, s1],))))	1])		



Input Code	ATen FX Graph	Define-by-run IR	Scheduling/Fusion	Output Triton	Output Wrapper

```
@triton.jit
def triton 0(in ptr0, out ptr0, out ptr1, xnumel, XBL0CK : tl.constexpr):
    xoffset = tl.program id(0) * XBLOCK
    xindex = xoffset + tl.arange(0, XBLOCK)[:]
    xmask = xindex < xnumel</pre>
    x0 = xindex
    tmp0 = tl.load(in ptr0 + (x0), None)
    tmp1 = tl.sin(tmp0)
    tmp2 = tl.cos(tmp1)
    tl.store(out ptr0 + (x0 + tl.zeros([XBL0CK], tl.int32)), tmp1, None)
    tl.store(out ptr1 + (x0 + tl.zeros([XBL0CK], tl.int32)), tmp2, None)
```

Input Code	ATen FX Graph	Define-by-run IR	Scheduling/Fusion	Output Triton	Output Wrapper
------------	---------------	------------------	-------------------	---------------	----------------

```
def call(args):
    arg0_1, = args
    args.clear()
    arg0_1_size = arg0_1.size()
    s0 = arg0_1_size[0]
    s1 = arg0_1_size[1]
    buf0 = empty_strided((s0, s1), (s1, 1), device='cuda', dtype=torch.float32)
    buf1 = empty_strided((s0, s1), (s1, 1), device='cuda', dtype=torch.float32)
    triton_0_xnumel = s0*s1
    triton_0.run(arg0_1, buf0, buf1, triton_0_xnumel, grid=grid(triton_0_xnumel))
    return (buf0, buf1, )
```

TorchInductor Example: C++ Output

```
Change device='cuda' to device='cpu'
```

}

```
extern "C" void kernel(const float* restrict in ptr0,
                       float* restrict out ptr0,
                       float* restrict out ptr1,
                       const long ks0,
                       const long ks1)
{
   #pragma omp parallel num threads(8)
        {
            #pragma omp for
            for(long i0=0; i0<((ks0*ks1) / 16); i0+=1)</pre>
            {
                auto tmp0 = at::vec::Vectorized<float>::loadu(in ptr0 + 16*i0);
                auto tmp1 = tmp0.sin();
                auto tmp2 = tmp1.cos();
                tmp1.store(out ptr0 + 16*i0);
                tmp2.store(out ptr1 + 16*i0);
            }
            #pragma omp for simd simdlen(8)
            for(long i0=16*(((ks0*ks1) / 16)); i0<ks0*ks1; i0+=1)</pre>
            {
                auto tmp0 = in ptr0[i0];
                auto tmp1 = std::sin(tmp0);
                auto tmp2 = std::cos(tmp1);
                out ptr0[i0] = tmp1;
                out ptr1[i0] = tmp2;
            }
        }
    }
```

NVIDIA A100 PERFORMANCE



Geomean speedup over PyTorch eager using float16 Higher is better

NVIDIA A100 PERFORMANCE





Cumulative distribution function of speedups over PyTorch eager.

	Inference	Training
All TorchInductor optimizations	$1.91 \times$	$1.45 \times$
Without loop/layout reordering	$1.91 \times (-0.00)$	$1.28 \times$ (-0.17)
Without matmul templates	$1.85 \times (-0.06)$	$1.41 \times (-0.04)$
Without parameter freezing	$1.85 \times (-0.06)$	$1.45 \times$ (-0.00)
Without pattern matching	$1.83 \times (-0.08)$	$1.45 \times$ (-0.00)
Without cudagraphs	$1.81 \times (-0.10)$	$1.37 \times (-0.08)$
Without fusion	$1.68 \times (-0.23)$	$1.27 \times (-0.18)$
Without inlining	$1.58 \times (-0.33)$	$1.31 \times (-0.14)$
Without fusion and inlining	$0.80 \times (-1.11)$	0.59× (-0.86)

Geomean speedup over PyTorch eager on 45 models from HuggingFace using fp16

Live PyTorch 2.0 Q&A Series:

https://www.youtube.com/@PyTorch PyTorch Dev Podcast (by ezyang)

https://pytorch-dev-podcast.simplecast.com/

Code:

https://github.com/pytorch/pytorch/tree/master/torch/_dynamo https://github.com/pytorch/pytorch/tree/master/torch/_functorch/a ot_autograd.py

https://github.com/pytorch/pytorch/tree/master/torch/_inductor

We are hiring!

Reach out to pengwu@meta.com