



Microsoft Research

Summit 2021

Accelerating  PyTorch
with Torchy, a tracing JIT compiler

Nuno Lopes
MSR Cambridge

Eager-mode frameworks are amazing!

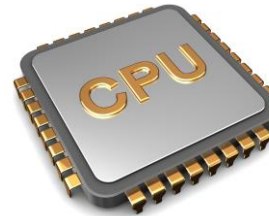
```
x = torch.tensor(((1.,2.), (3.,4.)))  
y = torch.tensor(((5.,6.), (7.,8.)))  
  
z = x.mul(y)  
z = z.add(y)  
x.add_(z)  
  
print(x)
```

```
tensor([[11., 20.],  
        [31., 44.]])
```

Eager-mode frameworks are slow! 🙄

```
x = torch.tensor(((1.,2.), (3.,4.)))  
y = torch.tensor(((5.,6.), (7.,8.)))  
  
z = x.mul(y)  
z = z.add(y)  
x.add_(z)  
  
print(x)
```

```
tensor([[11., 20.],  
        [31., 44.]])
```



```
tensor([[ 5., 12.],  
        [21., 32.]])
```

```
tensor([[10., 18.],  
        [28., 40.]])
```

```
tensor([[11., 20.],  
        [31., 44.]])
```



Eager-frameworks "hacks"

```
x = torch.tensor(((1.,2.), (3.,4.)))  
y = x.transpose(0, 1)  
y[0,0] = 42  
print(y)  
print(x)
```

tensor([[42., 3.],
 [2., 4.]])

tensor([[42., 2.],
 [3., 4.]])

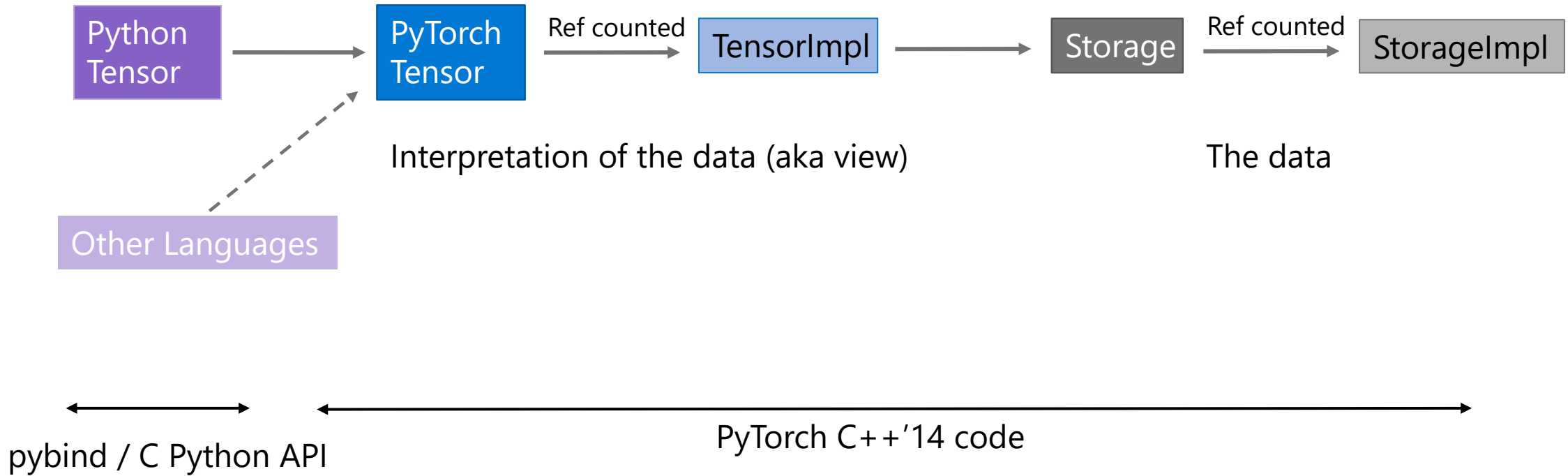
+ ever-increasing list of fused ops that users need to call manually

Transpose fuses marvelously with matmul!

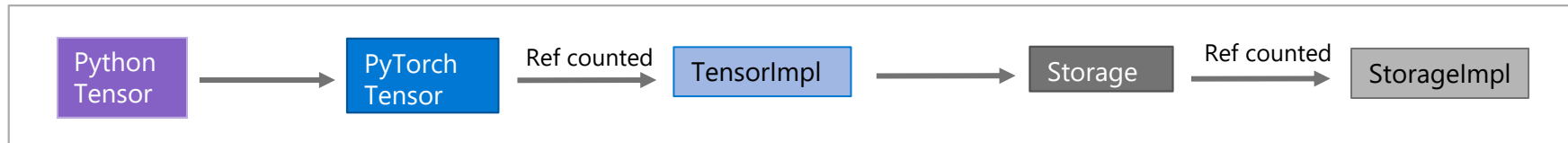
What's a Tensor?

Multiple tensor implementations:

- Dense
- Sparse
- Batched
- ...



Tensor creation



```
x = torch.tensor(((1.,2.), (3.,4.)))
y = torch.tensor(((5.,6.), (7.,8.)))

z = x.mul(y)
x.add_(z)
w = x.to(torch.float, copy=False)
z = x.transpose(0, 1)
```

New Tensor/TensorImpl/Storage/StorageImpl w/ default type & placed on default device

New Tensor/TensorImpl/Storage/StorageImpl w/ same type & device as inputs

Nothing new; override StorageImpl's data

Bump Tensor ref count if types match; new Tensor/Storage otherwise

New Tensor/TensorImpl; bump Storage ref count

Life of a PyTorch function call

```
z = x.add(y)
```

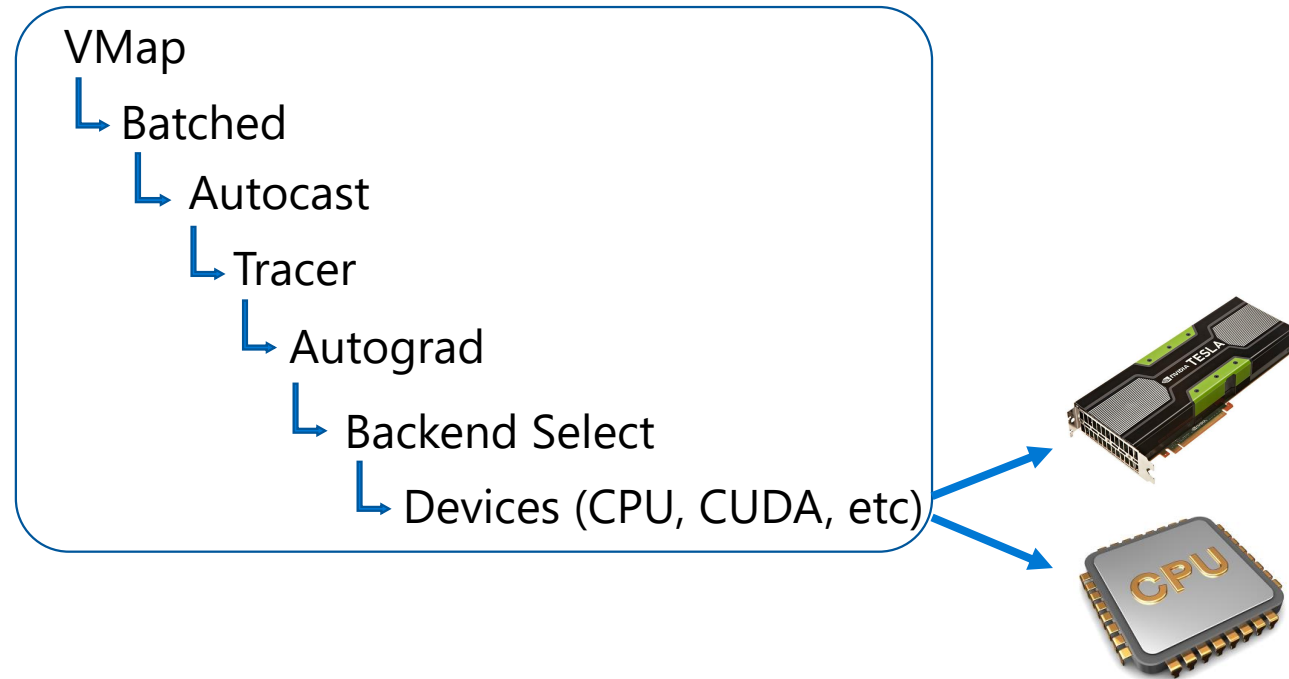
Dispatch:

Operation = Add.Tensor
Op0 = Tensor, CPU, Float
Op1 = Tensor, CPU, Float

Global dispatcher state:

Default device = CPU
Default type = Float
Include dispatch key = None

Waterfall dispatcher



Speeding up PyTorch today (single device)

- TorchScript Tracing (`torch.jit.trace`)
- TorchScript Compilation (`torch.jit.script`)
- ORT Module

- Autocast / apex (fp32 -> fp16)

TorchScript Tracing

- Function/module is executed (twice) with concrete inputs & operations recorded

```
def f(x, y):  
    z = x.add(y)  
    z.add_(x)  
    return x.mul(z)
```

```
w = torch.tensor(...)  
z = torch.tensor(...)  
torch.jit.trace(f, (w, z))
```

SSA-based IR:

```
def f(x: Tensor, y: Tensor) -> Tensor:  
    z = torch.add(x, y, alpha=1)  
    z0 = torch.add_(z, x, alpha=1)  
    return torch.mul(x, z0)
```

Tracing input-dependent code

```
def RAdam(wd, N_sma, ...):  
    if wd != 0:  
        p_data_fp32.add_(p_data_fp32, alpha=-wd * lr)  
  
    # more conservative since it's an approximated value  
    if N_sma >= 5:  
        denom = exp_avg_sq.sqrt().add_(eps)  
        p_data_fp32.addcdiv_(exp_avg, denom, value=-step_size)  
    else:  
        p_data_fp32.add_(exp_avg, alpha=-step_size)
```

- There are 4 possible different traces depending on the input!
- But TorchScript Tracing only supports single-trace functions.

TorchScript Compilation

- Compiler from Python AST to an SSA-based IR (the same used by tracing)
- Supports functions with control-flow
- But no support for too many Python features (only tensor inputs, no lambdas, no union types, etc, etc)

- Many real-world codebases are too pythonic. Will never work with TorchScript!



Is PyTorch inherently inefficient?

All problems in computer science can be solved by another level of indirection, except performance.

All problems in computer science can be solved by another level of indirection, *including* ~~except~~ performance.



Torchy
A tracing JIT compiler for PyTorch

Most Tensors are not observed

```
w = x.mul(y)
w = w.add(y)
w.add_(x)

print(w)
```

- Function from 2 tensors to another tensor
- Intermediate values of w not observed

Tensors are only observed:

- Data access, e.g., for branching on data-dependent models
- Printing
- Some PyTorch functions query layout, size, etc for pre-dispatch optimization (a hack)

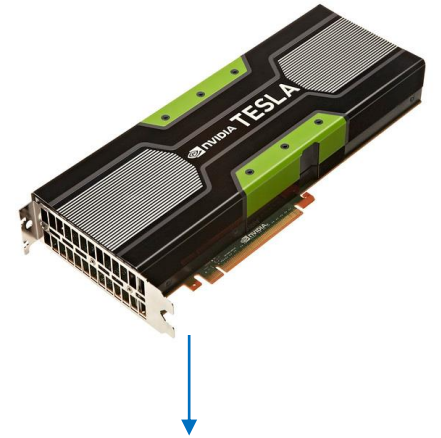
Idea: delay execution until observation

```
x = torch.tensor(((1.,2.), (3.,4.)))  
y = torch.tensor(((5.,6.), (7.,8.)))  
  
w = x.mul(y)  
w = w.add(y)  
x.add_(w)  
  
print(x)
```

Observation event!
Stop tracing and compute

Tracing JIT Compiler

```
w0 = x.mul(y)  
w1 = w0.add(y)  
x1 = x.add_(w1)
```



```
tensor([[11., 20.],  
        [31., 44.]])
```

Tracing JIT compilers

- A tremendous success for JavaScript in the past decade
- Batch operations as execution is delayed
- Detect which tensors are temporaries to help optimization

- Traces can be optimized before execution, or in background
- Traces repeat; optimization cost amortized

- Work with any codebase unmodified!

Intercepting PyTorch function calls

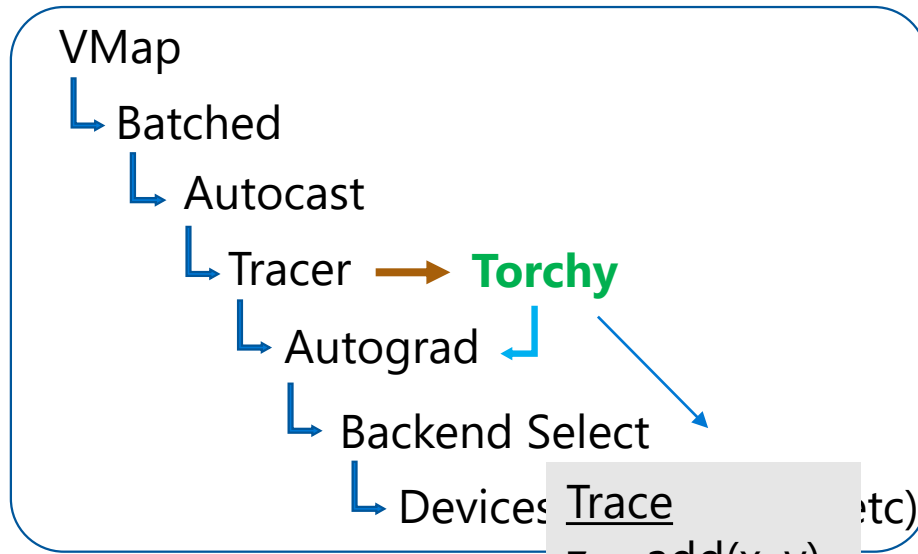
```
z = x.add(y)
```

Dispatch:
Operation = Add.Tensor
Op0 = Tensor, CPU, Float
Op1 = Tensor, CPU, Float

Global dispatcher state:
Default device = CPU
Default type = Float
Include dispatch key = Torchy

```
import torchy  
torchy.enable()
```

Waterfall dispatcher

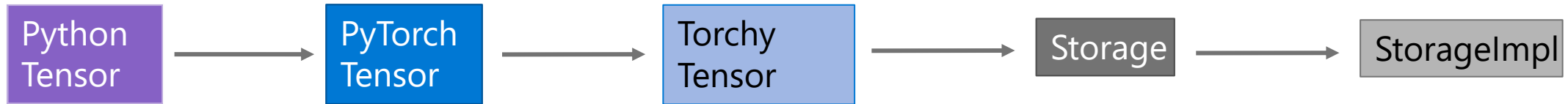


Alternative backends:

- Glow
- LLVM/MLIR
- ONNX Runtime (ORT)
- TorchScript
- XLA
- ...

Intercepting non-dispatched events

print(x) \longrightarrow x.storage() \longrightarrow x.storage()



Is tensor materialized?

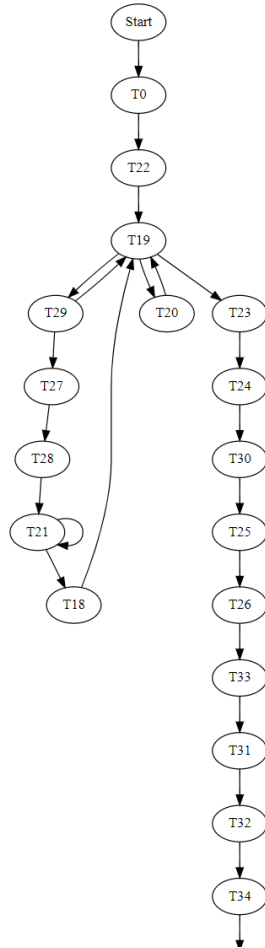
- Yes: behave like a normal tensor
- No: flush trace & act normally

Early experiments

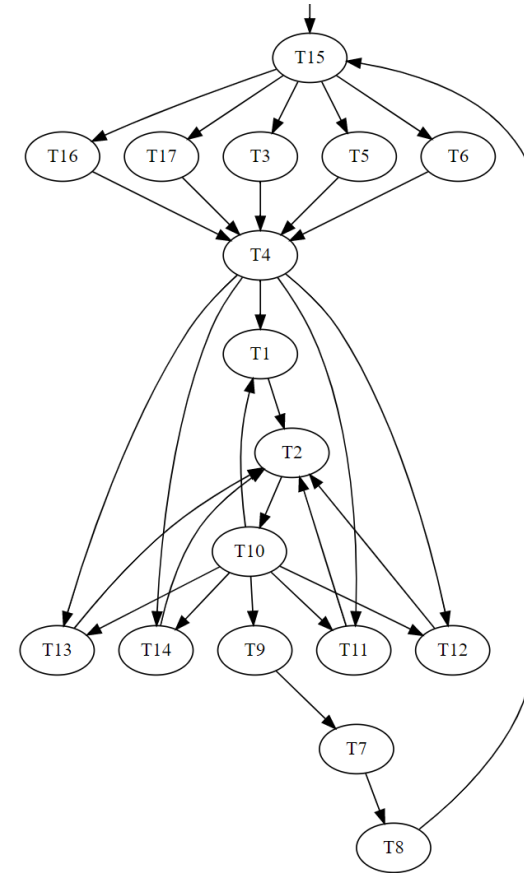
- Benchmarks:
 - TorchVision: ResNet-18, ResNeXt, MobileNet v3 large
 - Hungging Face: sentiment analysis (Bert)
- PyTorch 1.9+

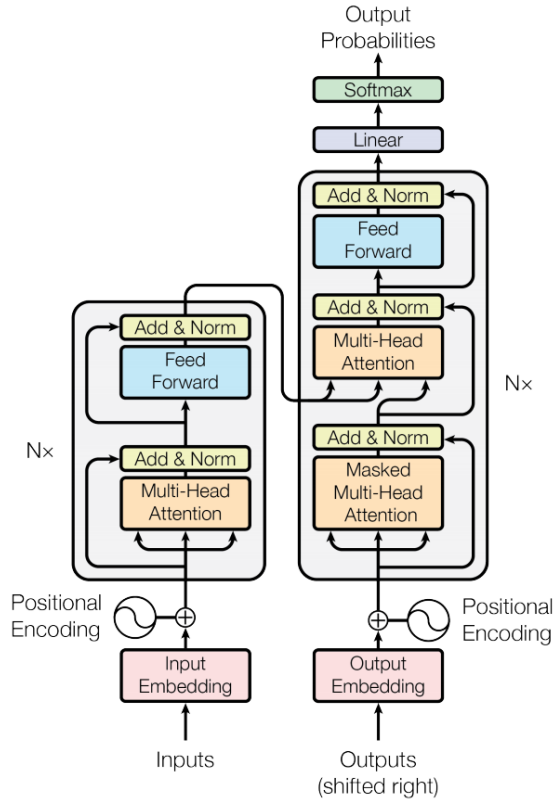
Bert from the compiler's perspective

Weight initialization

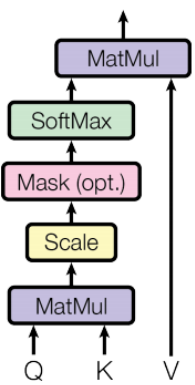


Model

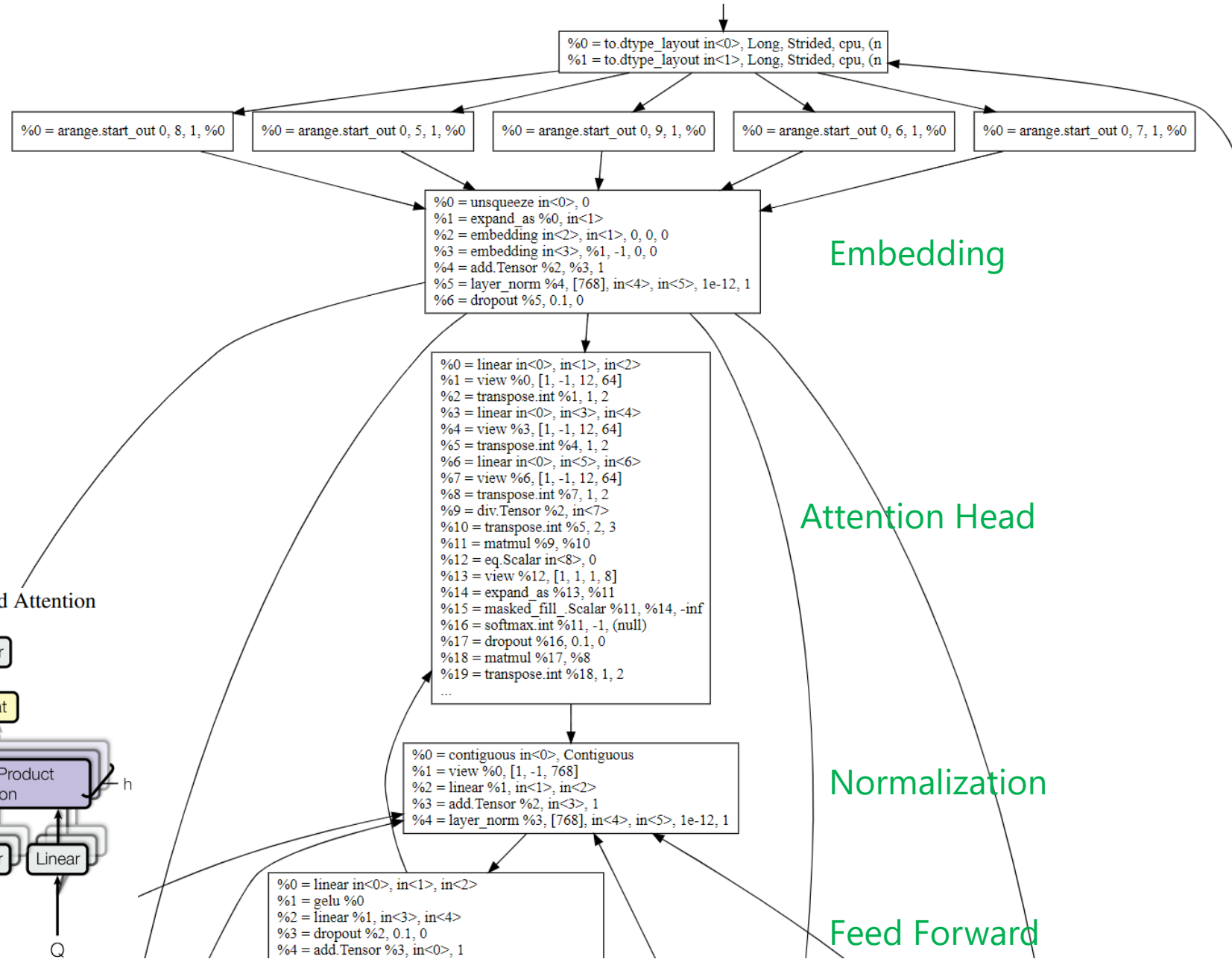
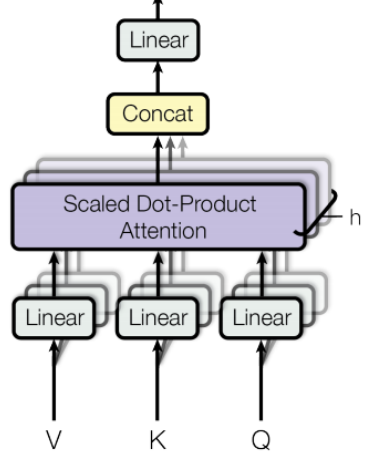


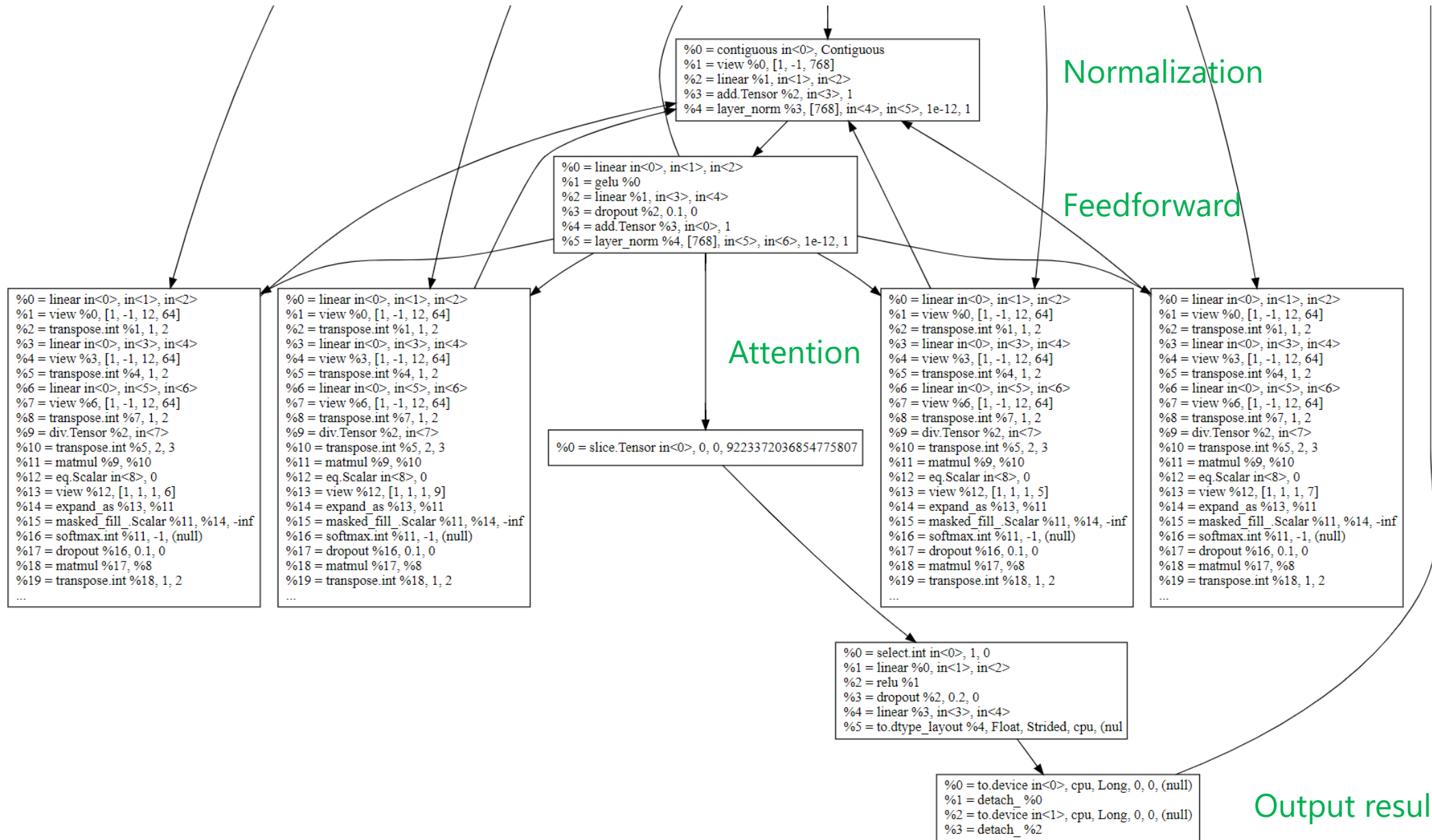


Scaled Dot-Product Attention



Multi-Head Attention





Summary: Torchy

- Acceleration for dynamic PyTorch programs
- Converts programs into small-ish straight-line programs (traces)
- Optimizes and runs each trace with the best backend
- Zero code changes! Just 'pip install torchy'