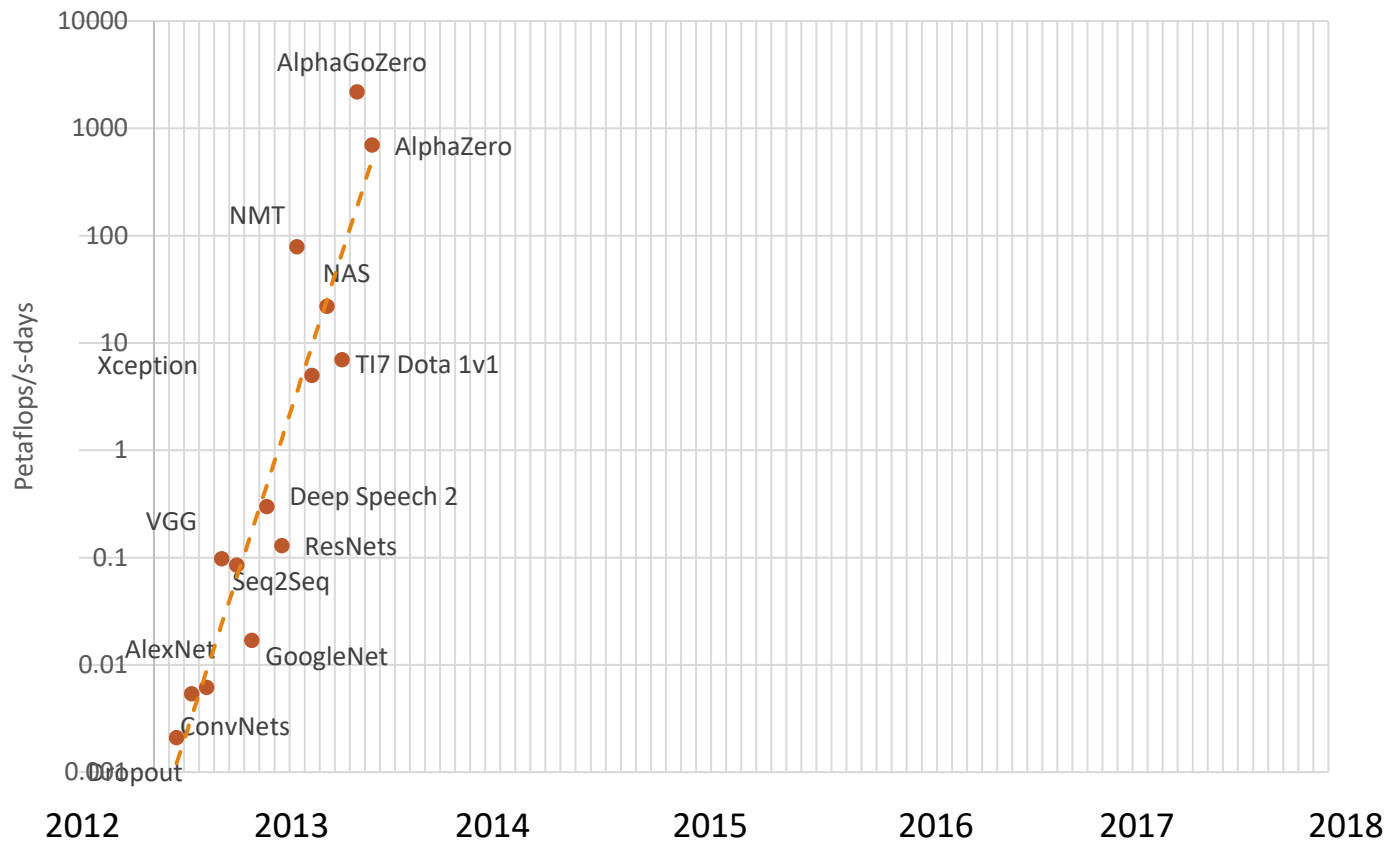# Torchy: A Tracing JIT Compiler for ⬤ PyTorch

NUNO P. LOPES

UNIVERSITY OF LISBON (IST-UL) & INESC-ID

# Computing demand for ML is exploding



FLOPS doubling
every 3.4 months!

# ML models run in frameworks

- First generation
  o Developer assembles model as a data-flow graph first
  o Hard for development & debugging
  o No support for dynamic models
  o TensorFlow 1

- Second generation / aka eager-mode or imperative
  o Instructions executed straight away
  o Easier
  o PyTorch, TensorFlow 2

# Eager-mode frameworks are amazing!

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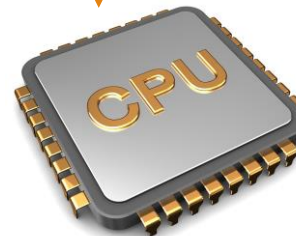
tensor([[ 5., 12.],
        [21., 32.]])

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

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

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

Is PyTorch inherently inefficient?

# 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)
```

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)

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

# 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)
```

Observable event!
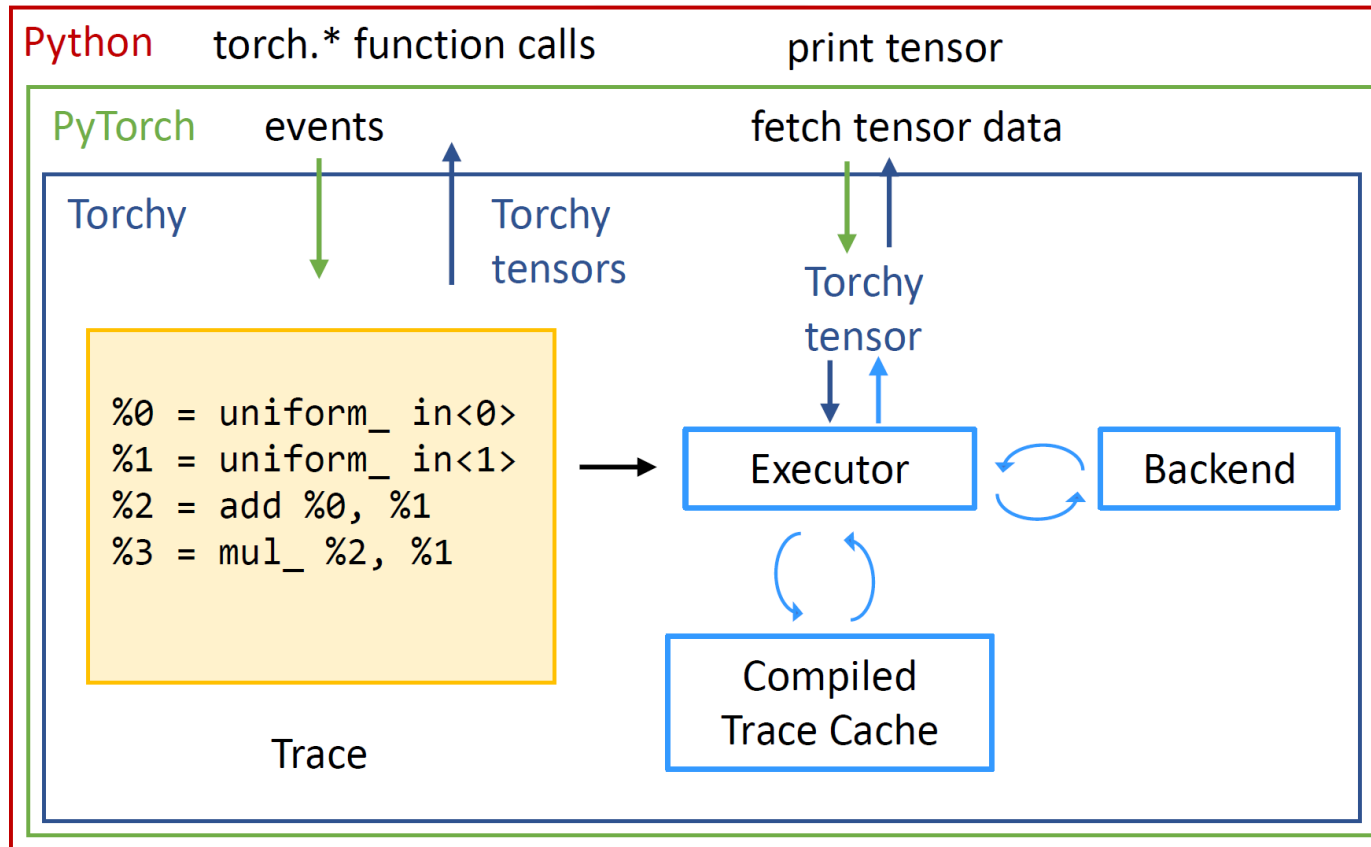Stop tracing and compute

Tracing JIT Compiler

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

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

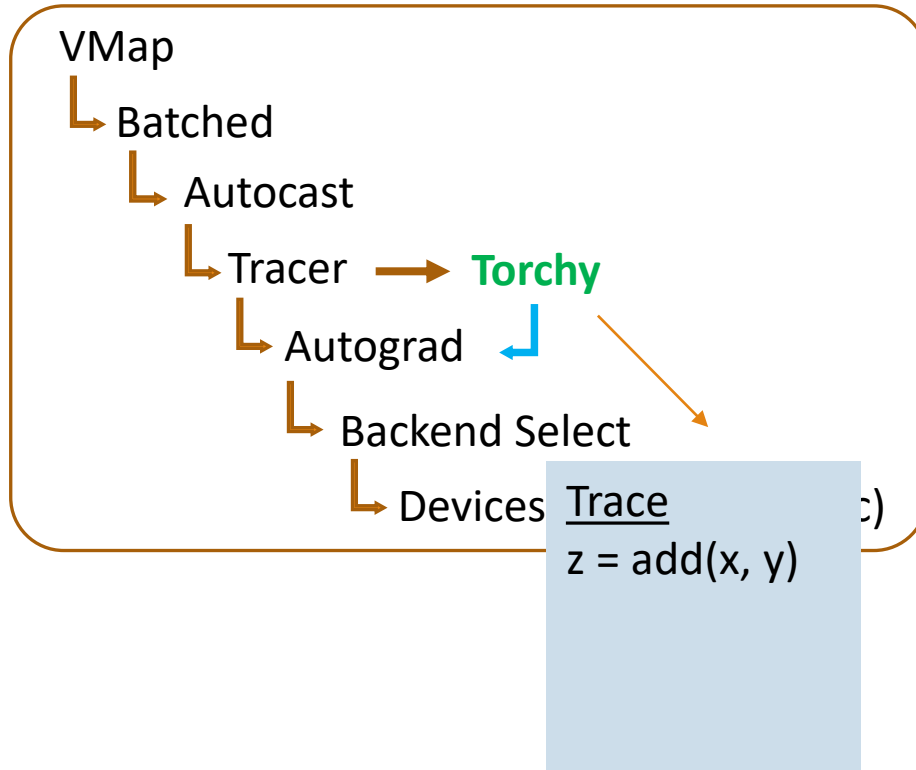# Torchy

# Intercepting PyTorch function calls

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

Waterfall dispatcher

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**

VMap
└ Batched
　└ Autocast
　　└ Tracer　→　**Torchy**
　　　└ Autograd
　　　　└ Backend Select
　　　　　└ Devices

Trace
z = add(x, y)

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

# Microbenchmarks



- Code with 8, 16, 32 elementwise operations

- Square matrices, n=100, 1k, 10k
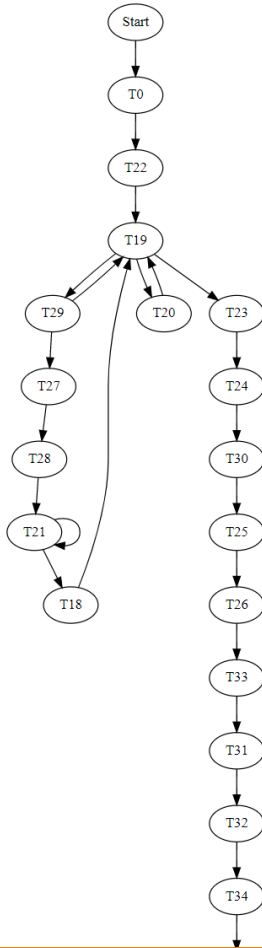
- Straight-line code & with control-flow

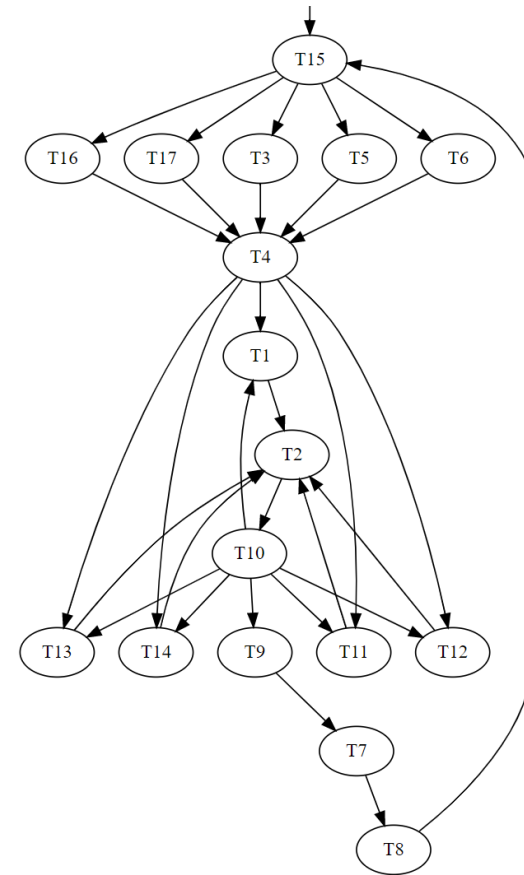# Experiments with Standard Models

- Run 1,000 inference queries over:
  - TorchVision: ResNet-18, ResNeXt, MobileNet v3 Large
  - Hungging Face: Bert Base/Large, GPT-2, RoBERTa Large
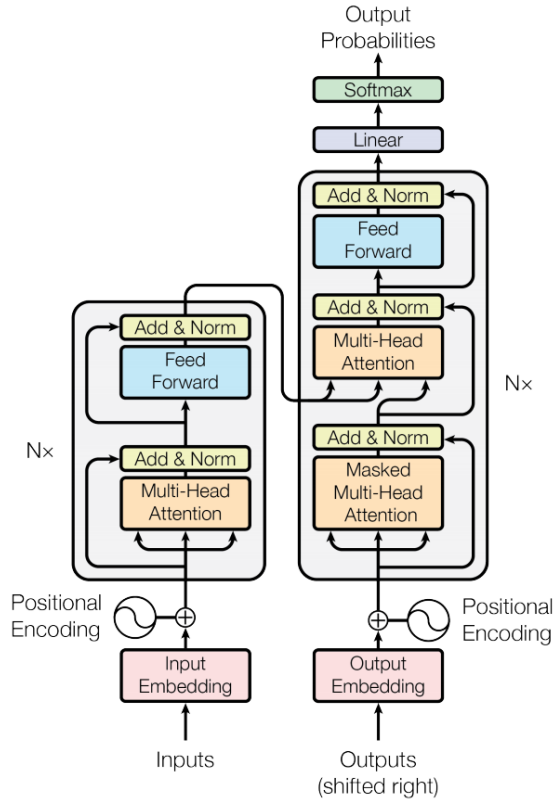
- PyTorch 1.9+

- 12 CPU cores

# Bert from the compiler's perspective
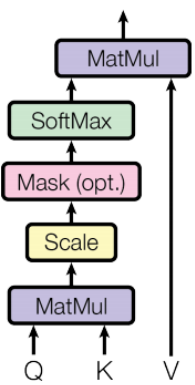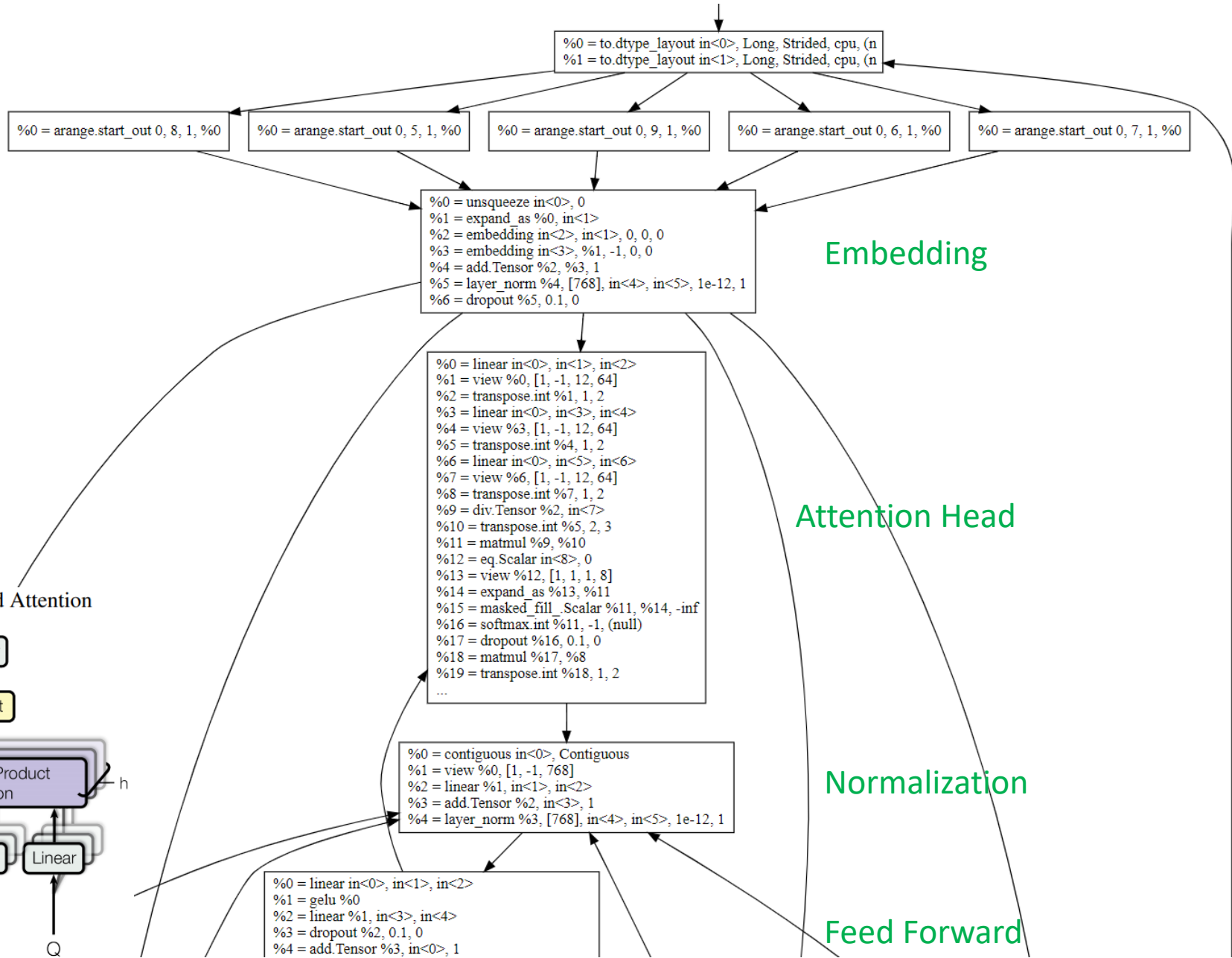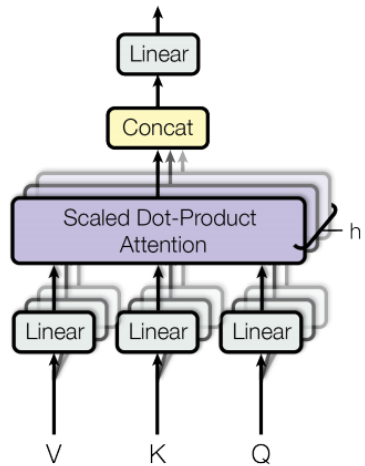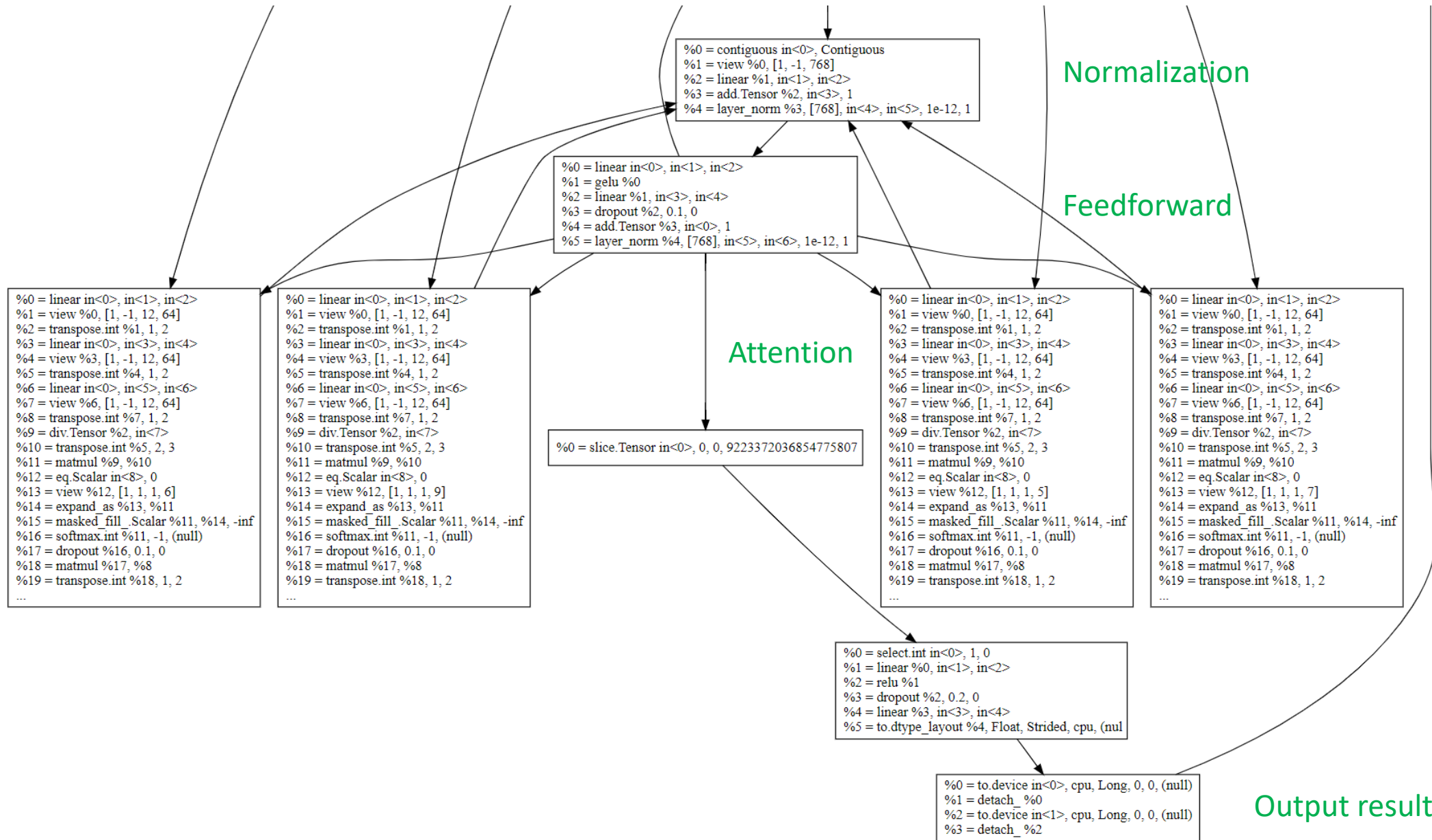


Weight initialization

Model

Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Add & Norm

Masked
Multi-Head
Attention

N×

N×

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

**Scaled Dot-Product Attention**

MatMul

SoftMax

Mask (opt.)

Scale

MatMul

Q    K    V

**Multi-Head Attention**

Linear

Concat

Scaled Dot-Product
Attention

h

Linear    Linear    Linear

V    K    Q

%0 = to.dtype_layout in<0>, Long, Strided, cpu, (n
%1 = to.dtype_layout in<1>, Long, Strided, cpu, (n

%0 = arange.start_out 0, 8, 1, %0

%0 = arange.start_out 0, 5, 1, %0

%0 = arange.start_out 0, 9, 1, %0

%0 = arange.start_out 0, 6, 1, %0

%0 = arange.start_out 0, 7, 1, %0

%0 = unsqueeze in<0>, 0
%1 = expand_as %0, in<1>
%2 = embedding in<2>, in<1>, 0, 0, 0
%3 = embedding in<3>, %1, -1, 0, 0
%4 = add.Tensor %2, %3, 1
%5 = layer_norm %4, [768], in<4>, in<5>, 1e-12, 1
%6 = dropout %5, 0.1, 0

**Embedding**

%0 = linear in<0>, in<1>, in<2>
%1 = view %0, [1, -1, 12, 64]
%2 = transpose.int %1, 1, 2
%3 = linear in<0>, in<3>, in<4>
%4 = view %3, [1, -1, 12, 64]
%5 = transpose.int %4, 1, 2
%6 = linear in<0>, in<5>, in<6>
%7 = view %6, [1, -1, 12, 64]
%8 = transpose.int %7, 1, 2
%9 = div.Tensor %2, in<7>
%10 = transpose.int %5, 2, 3
%11 = matmul %9, %10
%12 = eq.Scalar in<8>, 0
%13 = view %12, [1, 1, 1, 8]
%14 = expand_as %13, %11
%15 = masked_fill_.Scalar %11, %14, -inf
%16 = softmax.int %11, -1, (null)
%17 = dropout %16, 0.1, 0
%18 = matmul %17, %8
%19 = transpose.int %18, 1, 2
...

**Attention Head**

%0 = contiguous in<0>, Contiguous
%1 = view %0, [1, -1, 768]
%2 = linear %1, in<1>, in<2>
%3 = add.Tensor %2, in<3>, 1
%4 = layer_norm %3, [768], in<4>, in<5>, 1e-12, 1

**Normalization**

%0 = linear in<0>, in<1>, in<2>
%1 = gelu %0
%2 = linear %1, in<3>, in<4>
%3 = dropout %2, 0.1, 0
%4 = add.Tensor %3, in<0>, 1

**Feed Forward**
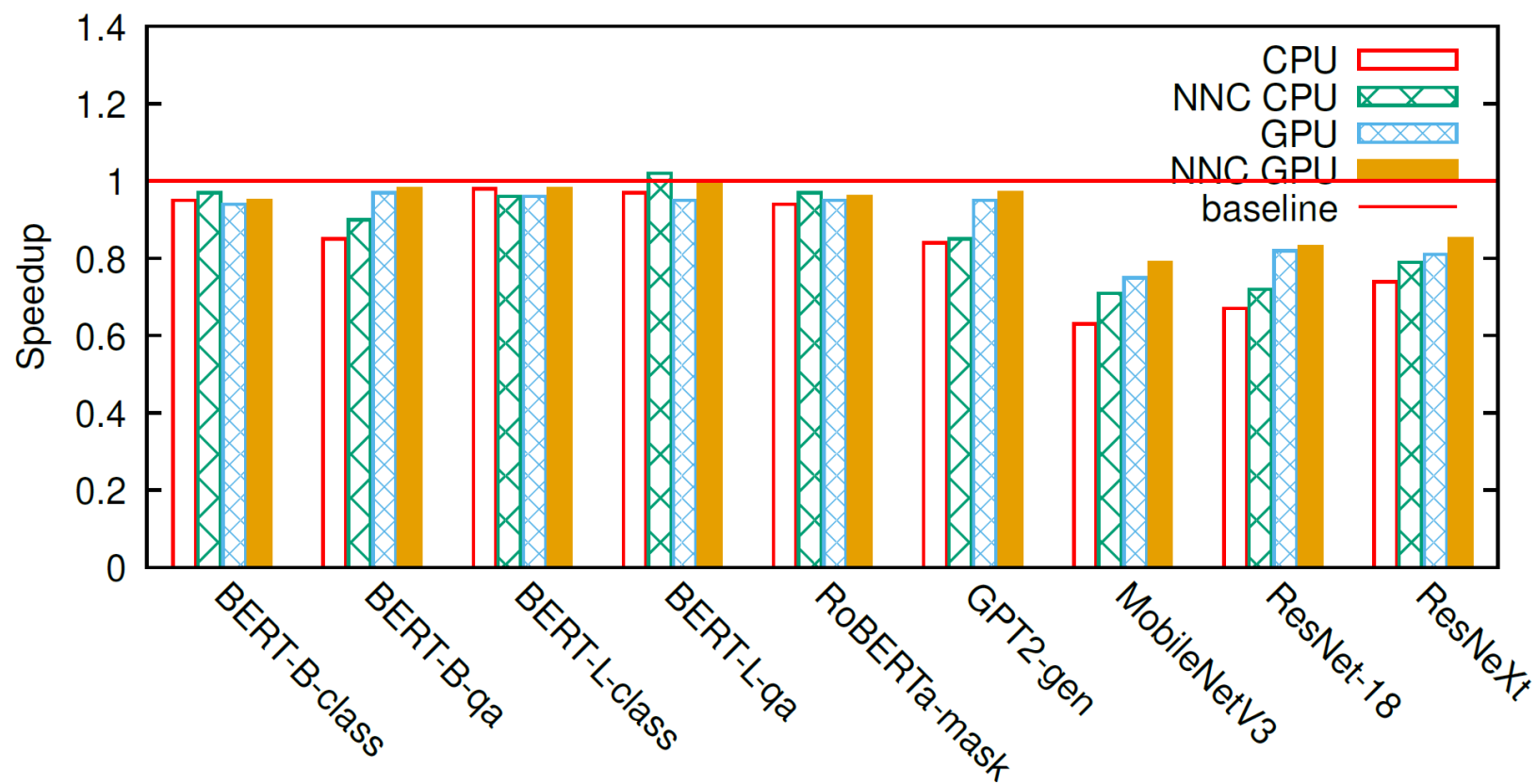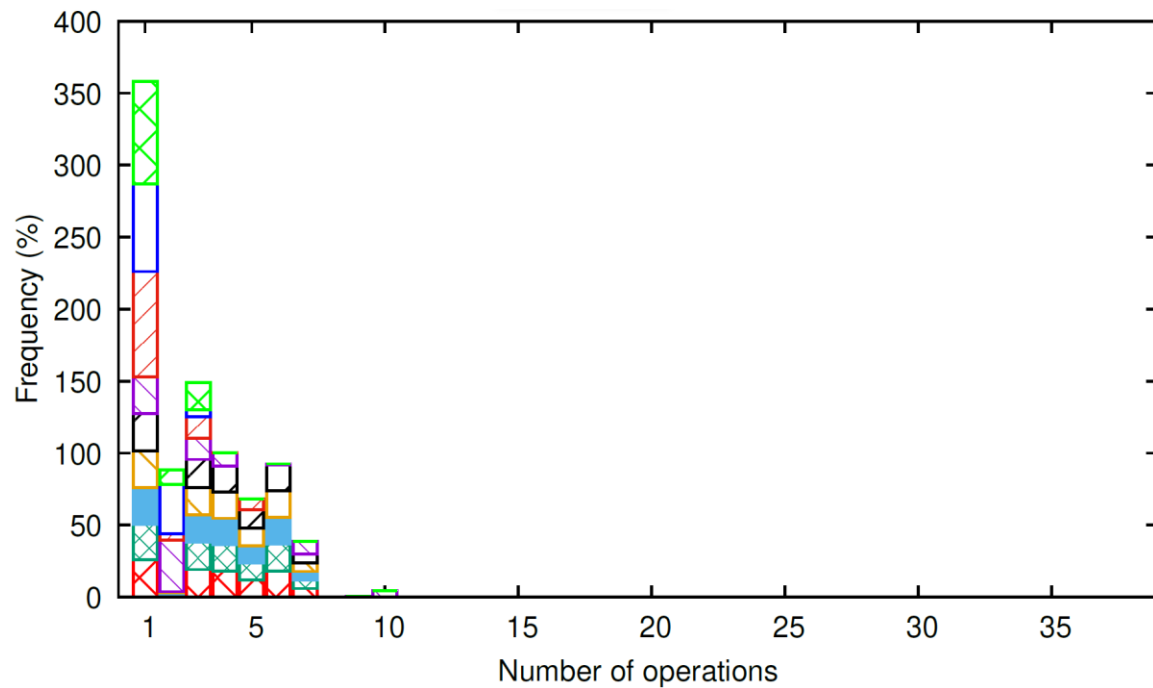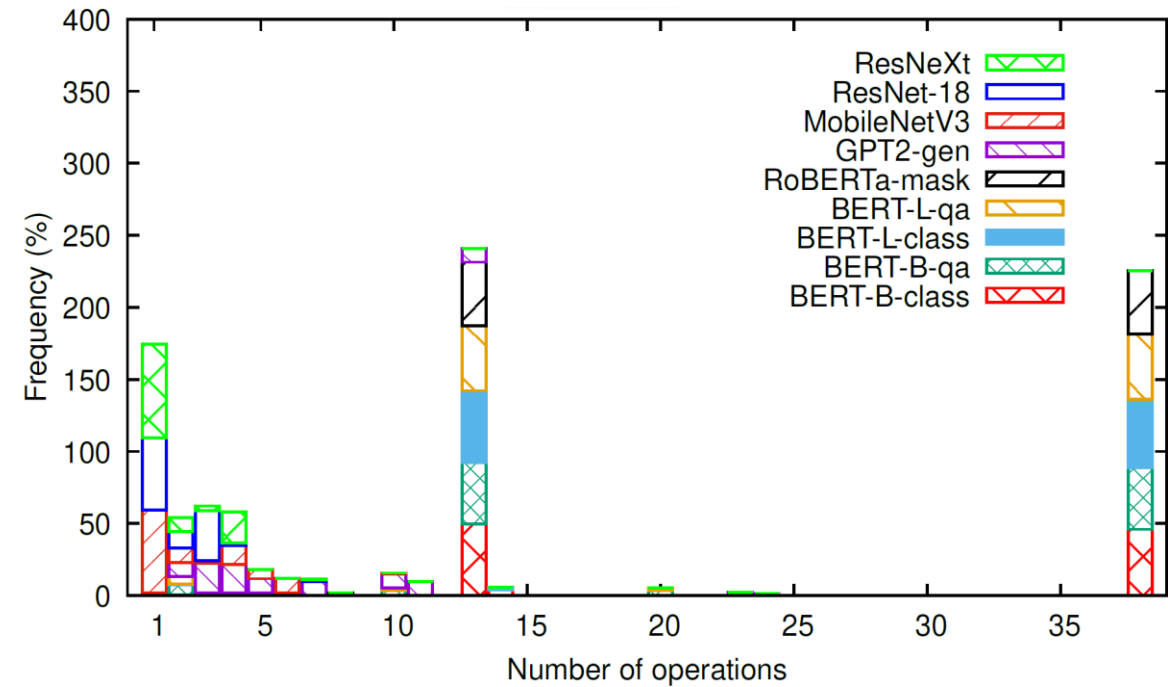
# Early Results

# Summary: Torchy

- Acceleration for dynamic PyTorch programs through JIT compilation

- Converts programs into small-ish straight-line programs (traces)

- Optimizes and runs each trace with the best backend

- Zero code changes! Just run 'pip install'
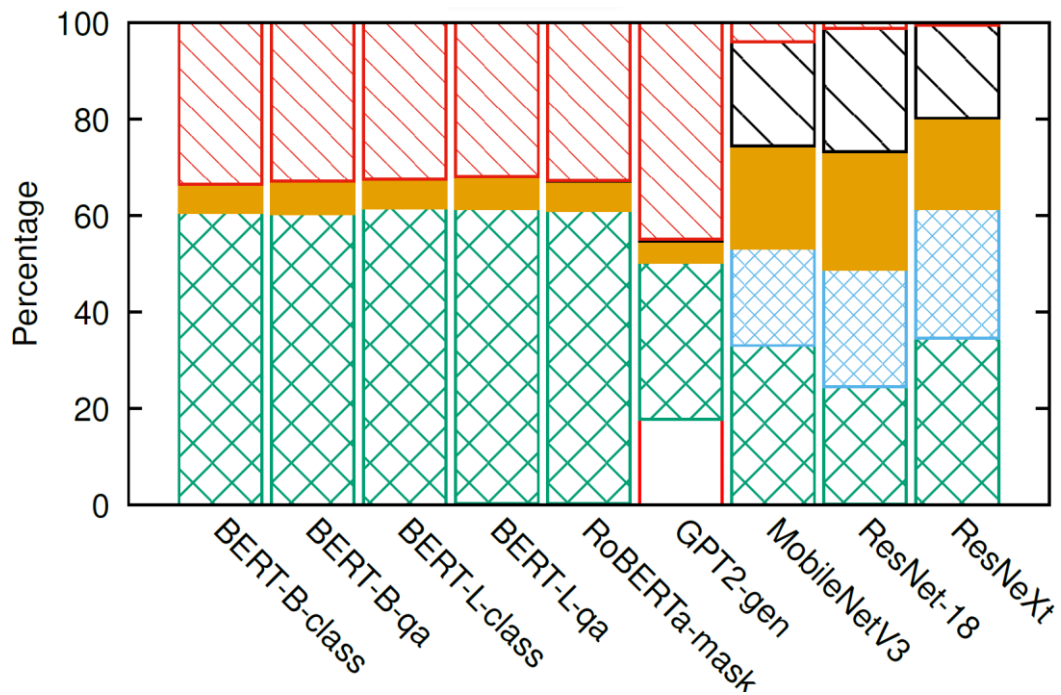
# Trace sizes



Without shape inference

With shape inference (partial)
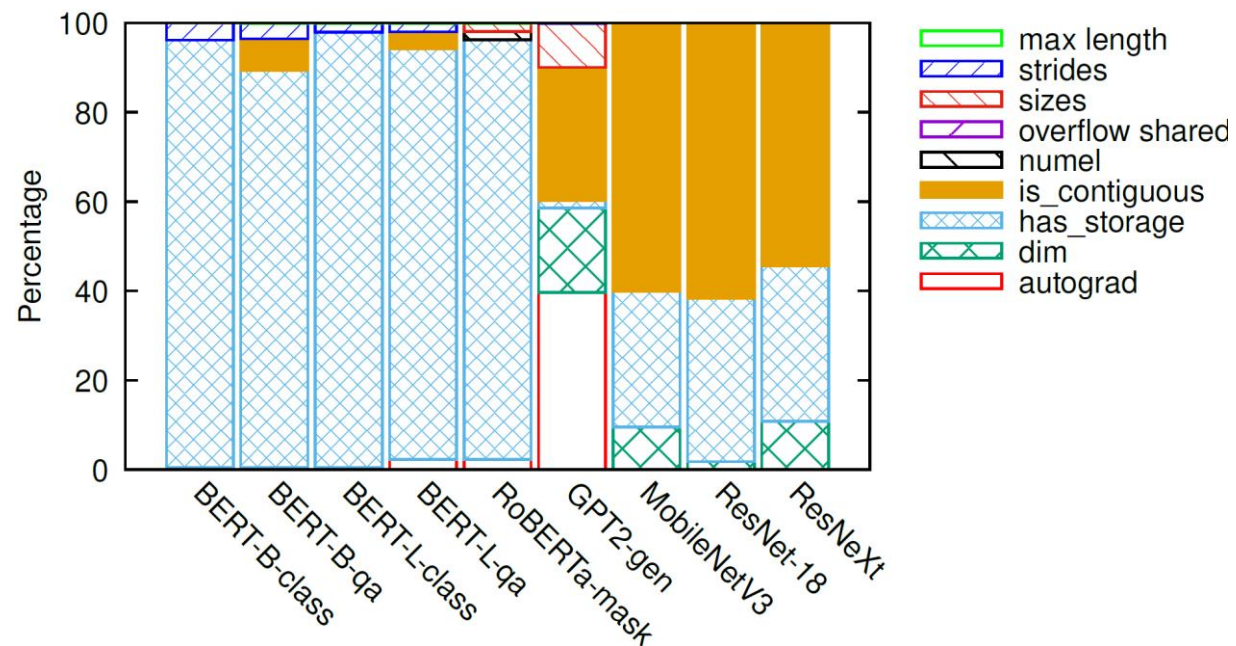
# Tracing JIT compilers

- A tremendous success for JavaScript in the past decade

- Peek into the future 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!

# Flush Reasons



Without shape inference
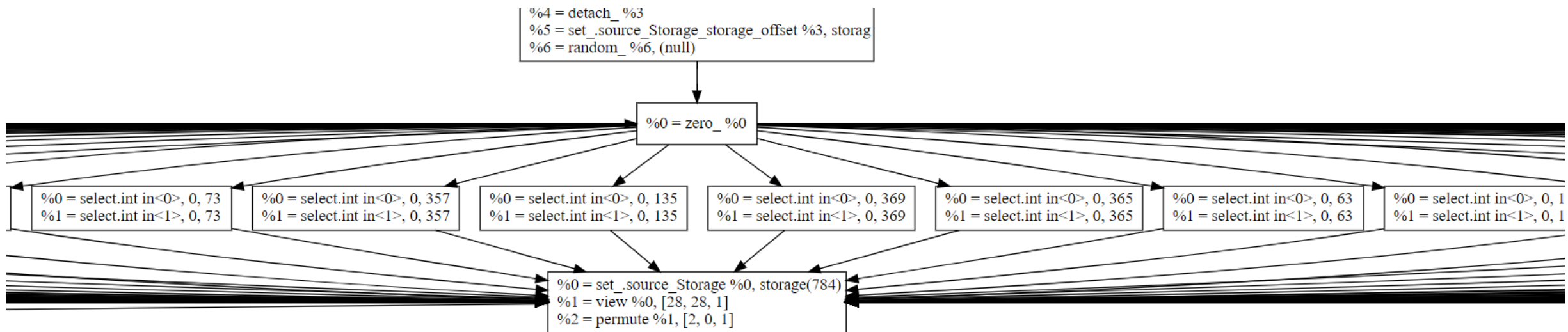
With shape inference (partial)

# How many traces?

| Model | Unique Traces wo/ inference |
|---|---|
| BERT Base | 87 |
| BERT Large | 87 |
| RoBERTa Large | 87 |
| GPT-2 | 246 |
| ResNet-18 | 1048 |
| ResNeXt | 1065 |
| MobileNet v3 large | 1124 |

# ResNet trace explosion

# Life of a PyTorch function call

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

Waterfall dispatcher

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
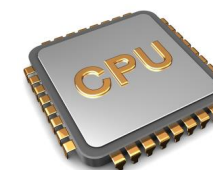
VMap
↳ Batched
  ↳ Autocast
    ↳ Tracer
      ↳ Autograd
        ↳ Backend Select
          ↳ Devices (CPU, CUDA, etc)

# What's a Tensor?

Multiple tensor implementations:
- Dense
- Sparse
- Batched
- …

| Python Tensor | → | PyTorch Tensor | Ref counted → | TensorImpl | → | Storage | Ref counted → | StorageImpl |

Other Languages ⇢ PyTorch Tensor

Interpretation of the data (aka view)

The data

pybind / C Python API

PyTorch C++'14 code

# 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

# 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) – by design!

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

# 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

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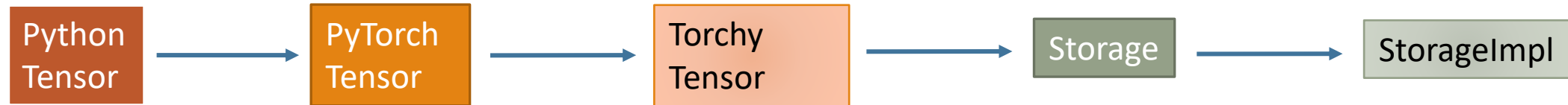
- There are 4 possible different traces depending on the input!
- But TorchScript Tracing only supports single-trace functions.

# Intercepting non-dispatched events

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

| Python Tensor | $\rightarrow$ | PyTorch Tensor | $\rightarrow$ | Torchy Tensor | $\rightarrow$ | Storage | $\rightarrow$ | StorageImpl |

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

# Eager-frameworks "hacks"

```
x = torch.tensor(((1.,2.), (3.,4.)))

z = x.transpose(0, 1)

z[0,0] = 42

print(z)
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!