

# PyTorch Internals

Edward Z. Yang

# Who's this talk for?

You want to contribute to PyTorch,  
but the codebase seems daunting

Notable omissions: JIT, distributed

# Concepts

Tensor/Storage/Strides

Layout/Device/Dtype

Autograd

# Mechanics

Operator call stack

Tools for writing kernels

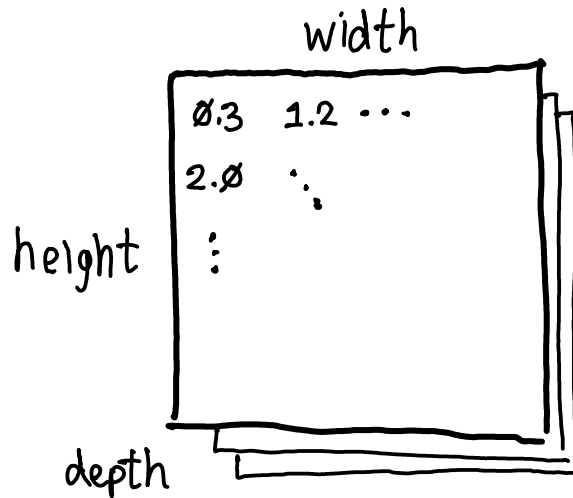
Legacy code

Efficient workflow

Concepts

Tensor

# Tensor



sizes	$(D, H, W)$	contiguous ←
strides	$(H*W, W, 1)$	
dtype	float	
device	cuda:0	
layout	strided	

# Tensor: Strided Representation

logical

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$$

`dtype=torch.int32`

maps to



physical

0x10	1
0x14	2
0x18	3
0x1c	4
0x20	

sizes

`[2, 2]`

strides

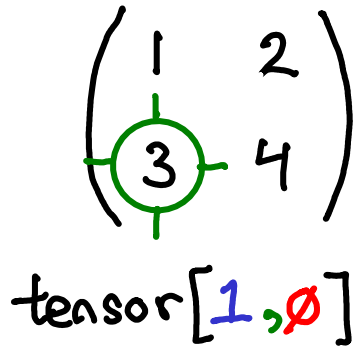
`[2, 1]`

# Tensor: Strided Representation

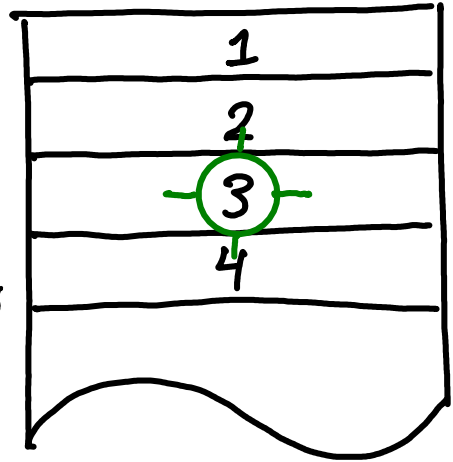
logical

$$\overset{\text{index}}{1} \times \overset{\text{stride}}{2} + \overset{\text{index}}{2} \times \overset{\text{stride}}{\emptyset} = 2$$

physical



$\emptyset \times 1 \emptyset$   
 $\emptyset \times 14$   
 $\emptyset \times 18$   
 $\emptyset \times 1C$   
 $\emptyset \times 2 \emptyset$



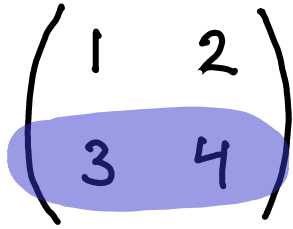
sizes  $[2, 2]$   
strides  $[2, 1]$

Upcoming: TensorAccessor



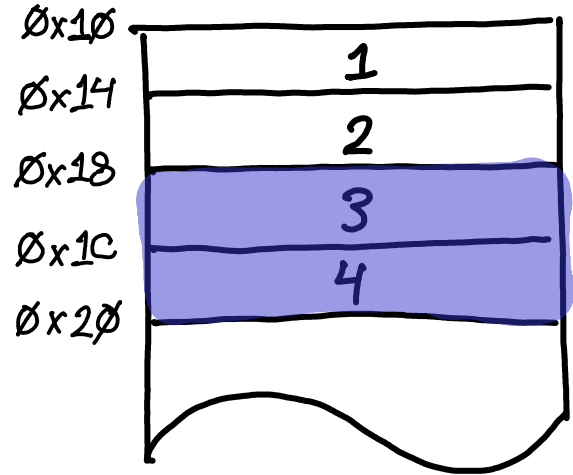
# Tensor: Strided Representation

logical



tensor[1, :]

physical



sizes  
strides  
offset

[ 2 ]  
[ 1 ]  
2

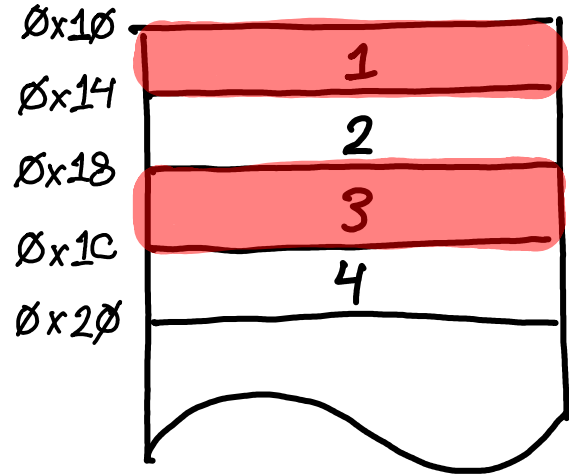
# Tensor: Strided Representation

logical

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$$

tensor[:, 0]

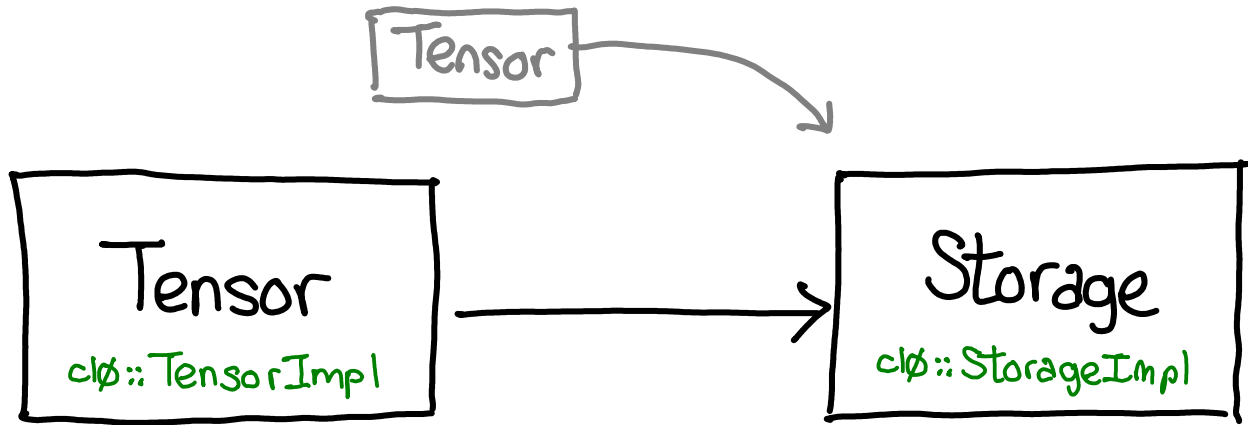
physical



sizes [2]  
strides [2]

stride  $\neq 1$  means we skip elements

# Tensor: Strided Representation



the logical view  
sizes, strides, offset

the actual physical data  
physical size, dtype

there is always a Tensor+Storage, even for "simple" cases

# Tensor operations

(slight simplification)

dtype

device type  
& layout

float

double

int

⋮

fixed set of  
supported dtype  
specializations

CPU impl

Sparse CPU impl

CUDA impl

XLA impl

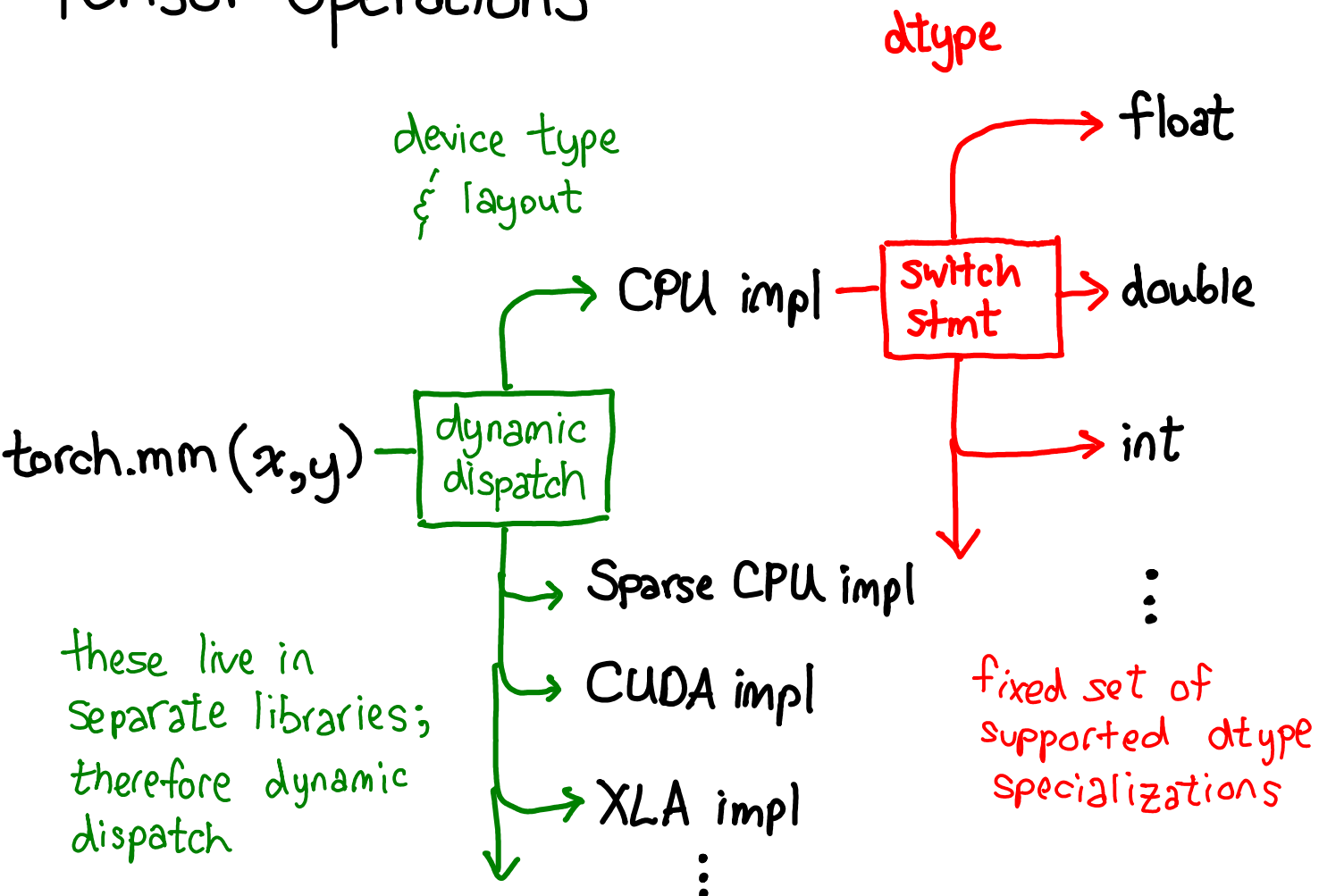
⋮

dynamic  
dispatch

Switch  
stmt

torch.mm(x,y)

these live in  
separate libraries;  
therefore dynamic  
dispatch



# Tensor extensions

sparse tensors

quantized tensors

encrypted tensors

There's more to life than strided tensors

MKLDNN tensors

TPU tensors

batch tensors

# tensor wrapper

e.g., batch tensors

(things that  
don't flow  
through autograd)

device

{  
cpu  
cuda  
xla  
hip  
fpga  
⋮  
}

has an allocator

layout

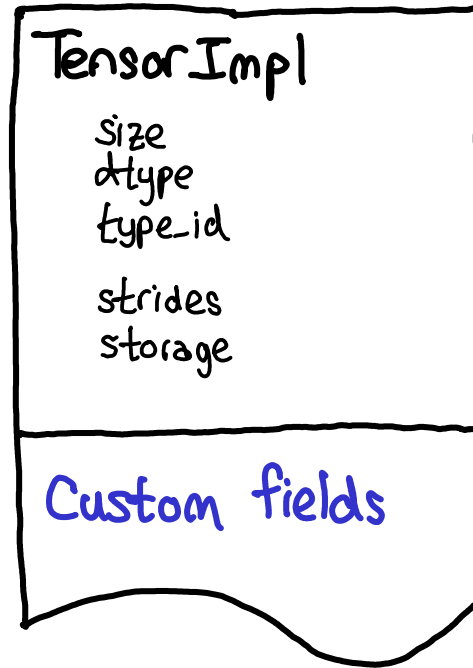
{  
strided  
sparse  
mkldnn  
⋮  
}

dtype

{  
float  
double  
int  
bool  
qint  
⋮  
}

implicitly convertible

fixed layout;  
no virtual methods!  
important for perf



} fields all tensors  
universally have

} important fields  
for strided tensors

} everything else

eg. values, indices  
for a sparse tensor

Autograd



i2h = torch.mm(W\_x, x.t())

h2h = torch.mm(W\_h, prev\_h.t())

next\_h = i2h + h2h

next\_h = next\_h.tanh()

loss = next\_h.sum()

loss.backward()

W\_h, W\_x, x, prev\_h  
requires grad

$i2h = \text{torch.mm}(W_x, x.t())$

$h2h = \text{torch.mm}(W_h, \text{prev\_h}.t())$

$\text{next\_h} = i2h + h2h$


$\text{next\_h2} = \text{next\_h}.\text{tanh}()$

$\text{loss} = \text{next\_h2}.\text{sum}()$

$W_h, W_x$ ,  $x$ ,  $\text{prev\_h}$   
requires grad

swap inputs and outputs!

$\text{grad\_loss} = \text{torch.tensor}(1, \text{dtype}=\text{loss.dtype}())$

$\text{grad\_next\_h2} += \text{grad\_loss}.\text{expand}(\text{next\_h2.size}())$   original inputs may be reused

$\text{grad\_next\_h} += \text{tanh\_backward}(\text{grad\_next\_h2}, \text{next\_h2})$

$\text{grad\_i2h}, \text{grad\_h2h} += \text{grad\_next\_h}, \text{grad\_next\_h}$

$W_h.\text{grad} += \text{mm\_mat1\_backward}(\text{grad\_h2h}, \text{prev\_h}.t(), W_h, 1)$

$W_x.\text{grad} += \text{mm\_mat1\_backward}(\text{grad\_i2h}, x.t(), W_x, 1)$

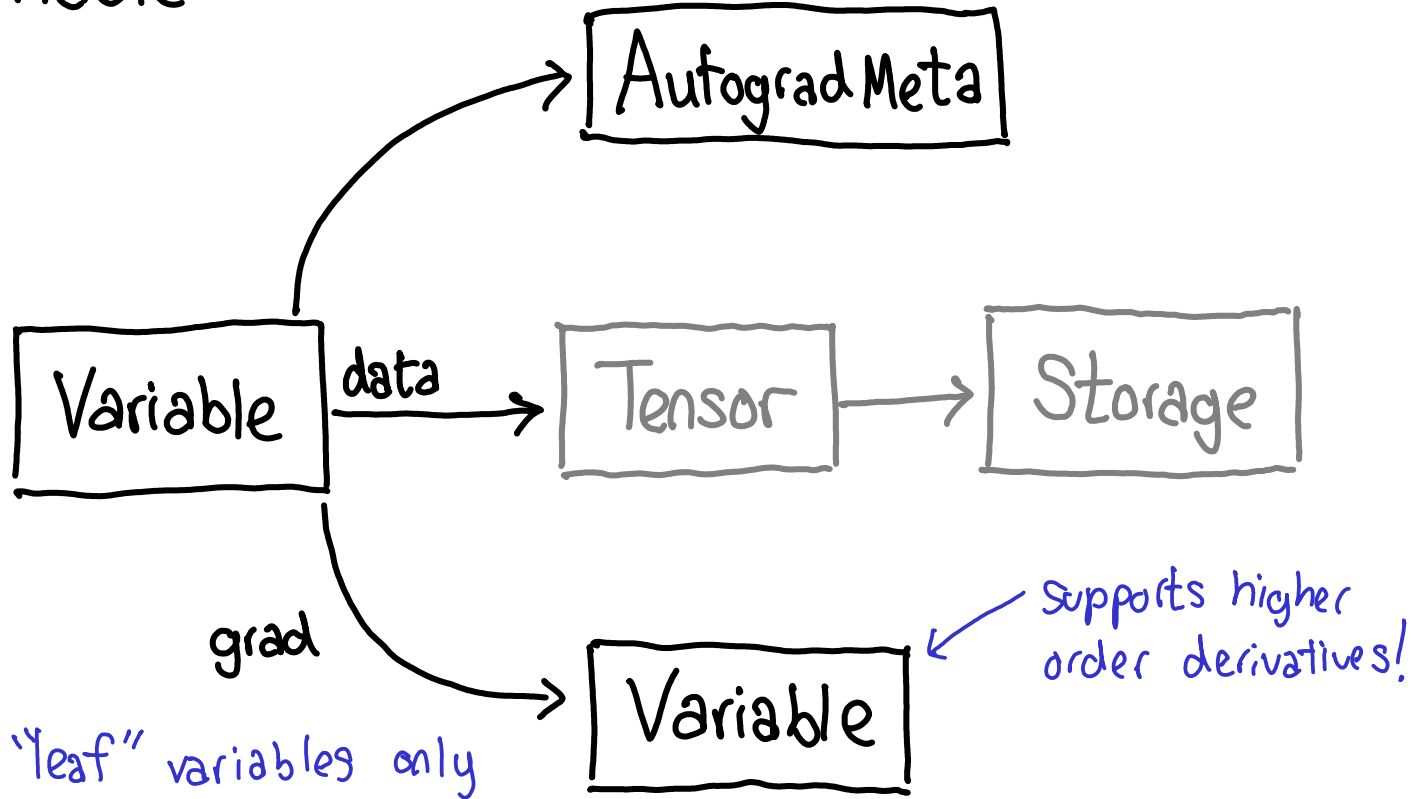
how do we get here?

All grad variables are zero to start:

$\text{grad\_i2h} = \text{torch.zeros\_like}(i2h)$

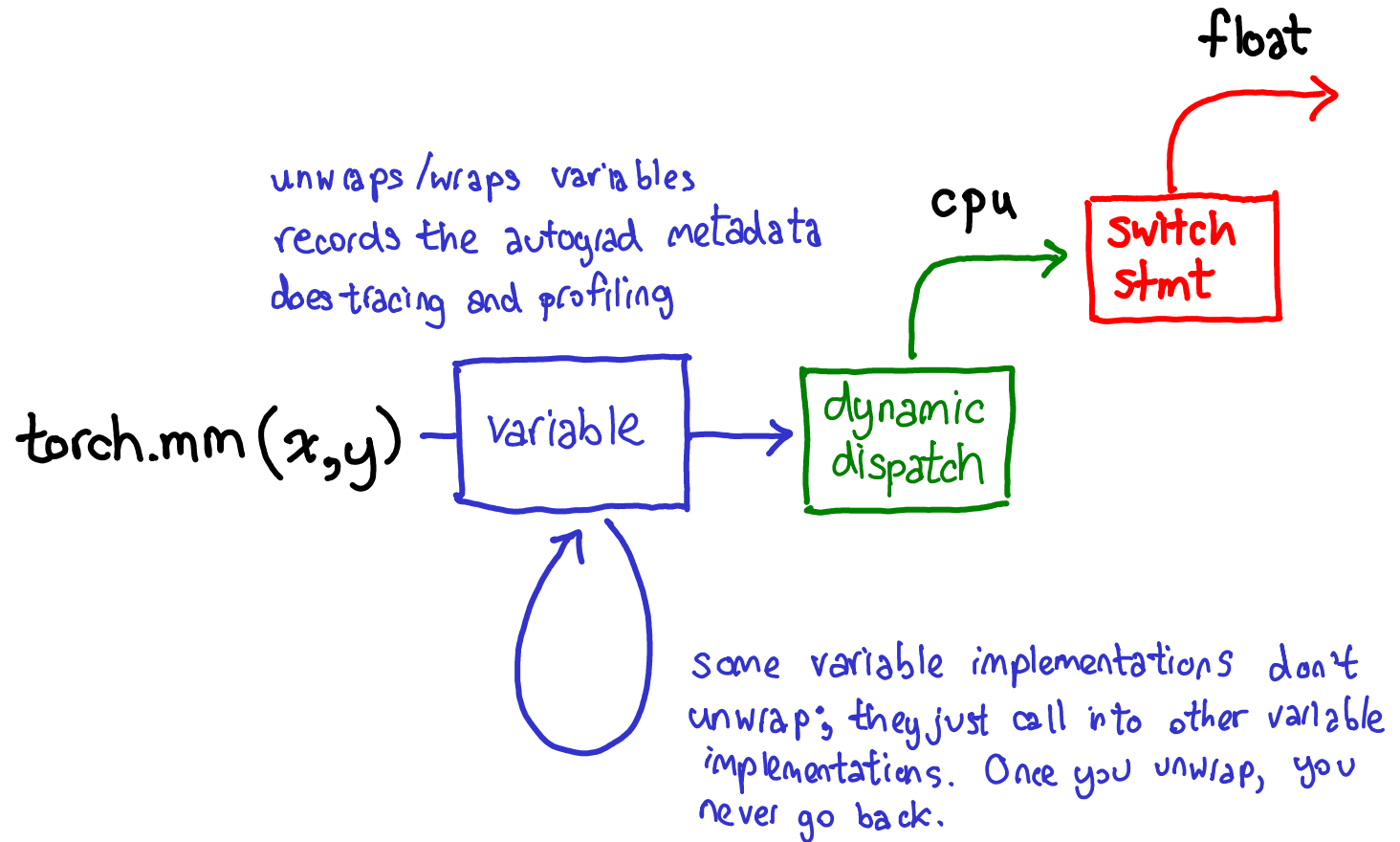
# Variable

metadata to  
generated backwards



NB: This changing soon!

# Dispatch with autograd



# AutogradMeta

Variable grad-;

shared\_ptr<Function> grad\_fn-;

weak\_ptr<Function> grad\_accumulator-;

bool requires\_grad-; // leaf only

uint32\_t output\_nr-;

class hierarchy



Function

edge\_list next\_edges-;

(e.g.,)

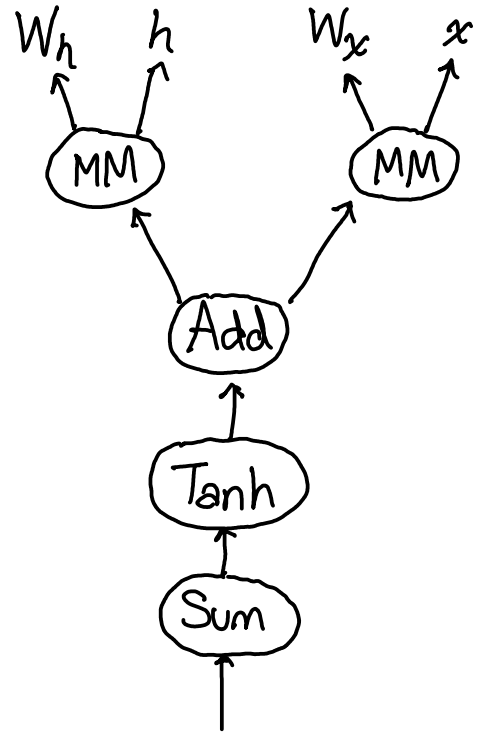
SavedVariable x;

Edge

shared\_ptr<Function> function;

uint32\_t input\_nr;

```
i2h = torch.mm(W_x, x.t())  
h2h = torch.mm(W_h, prev_h.t())  
next_h = i2h + h2h  
next_h = next_h.tanh()  
loss = next_h.sum()  
loss.backward()
```



requires\_grad=True

leaf variable →

$W_x$

grad\_acc\_

Function

Accumulate Grad

no edge for  $x$ ,  
it doesn't require  
grad!

X

not the input of  
the forward: input  
of backward

next\_edges[0]  
input\_nr = 0

next\_edges[1]  
input\_nr = 0

i2h

grad\_fn\_

Function

MmBackward

mat1\_

mat2\_

$W_x$

$x$

output\_nr\_ = 0

output of  
forward

not necessarily  
all inputs →

next\_edges[1]  
input\_nr = 0

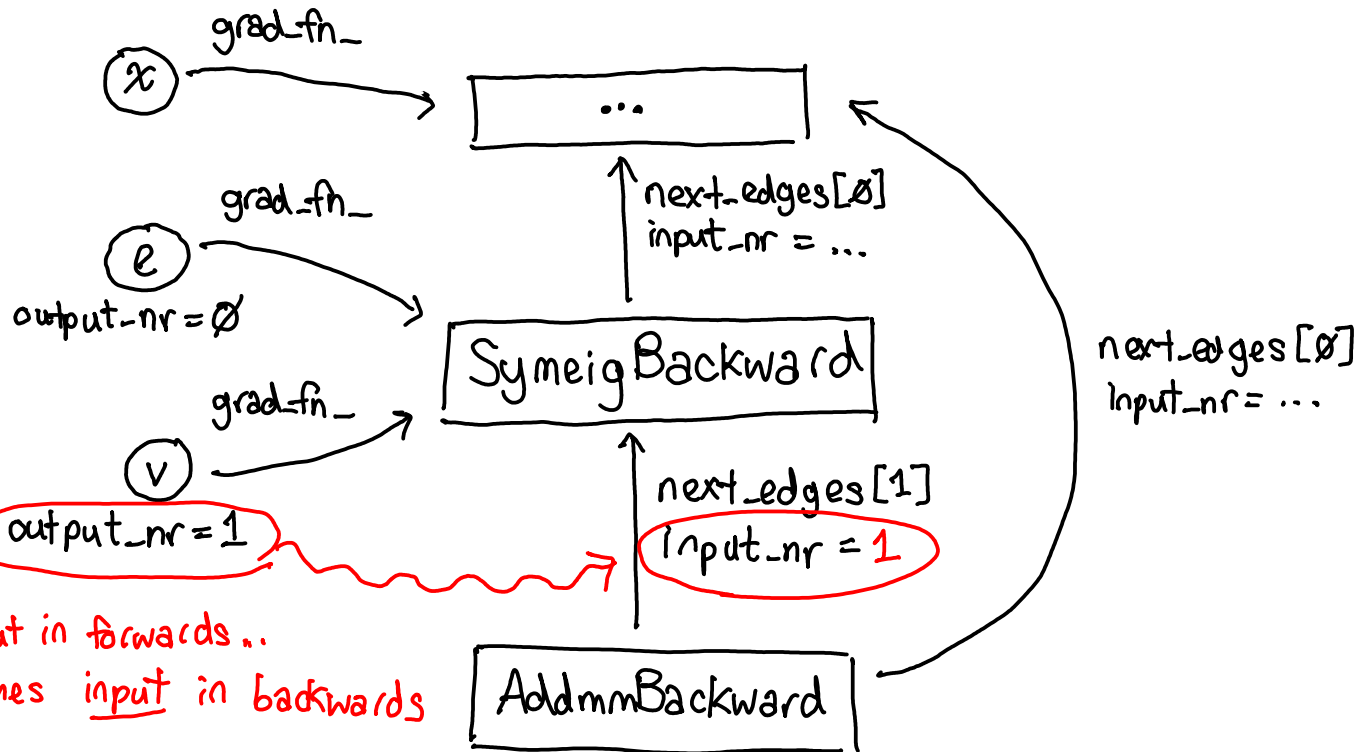
Function

AddBackward1

some fns  
have multiple  
overloads

# Multiple outputs

$e, v = \text{torch.symeig}(x, \text{eigenvectors}=\text{True})$   
 $a = \text{torch.mm}(x, v)$

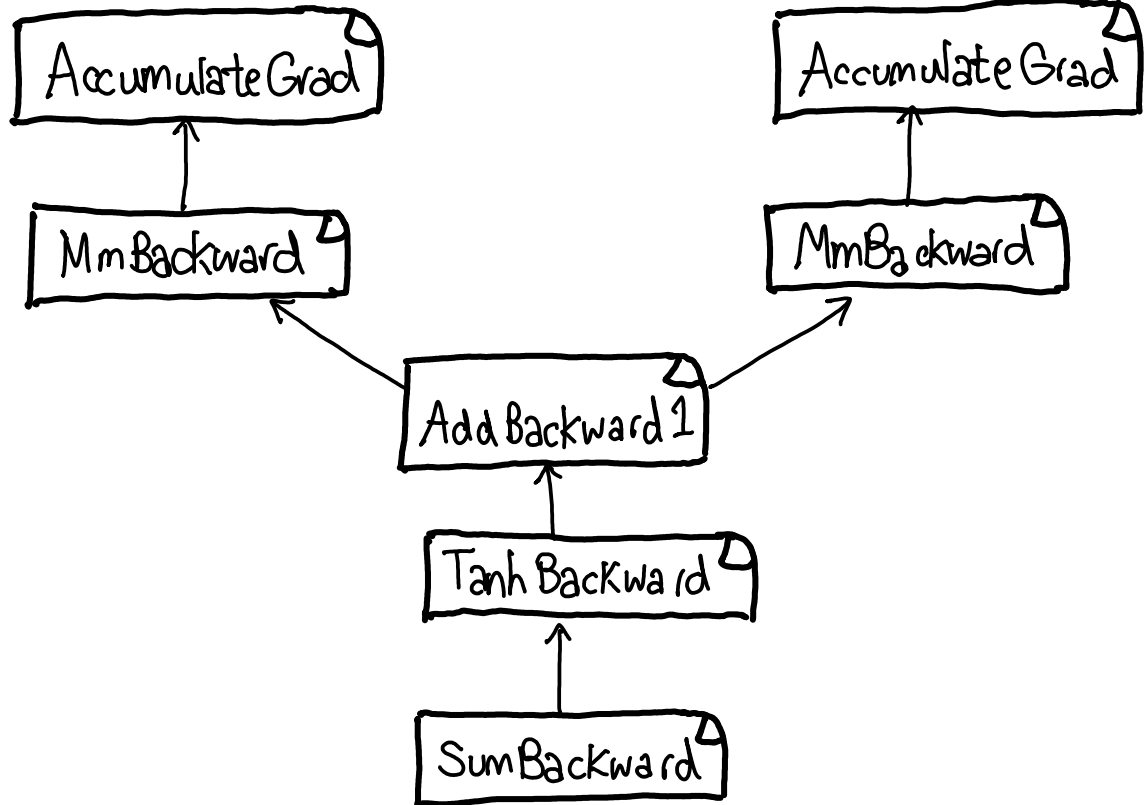




# Autograd engine

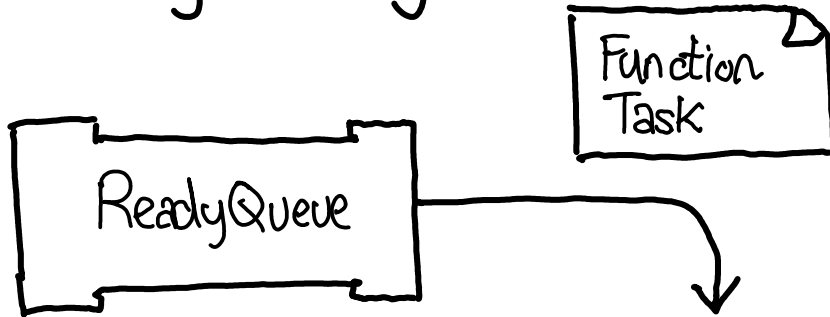
"just a parallel graph executor"

FunctionTasks:



# Autograd engine

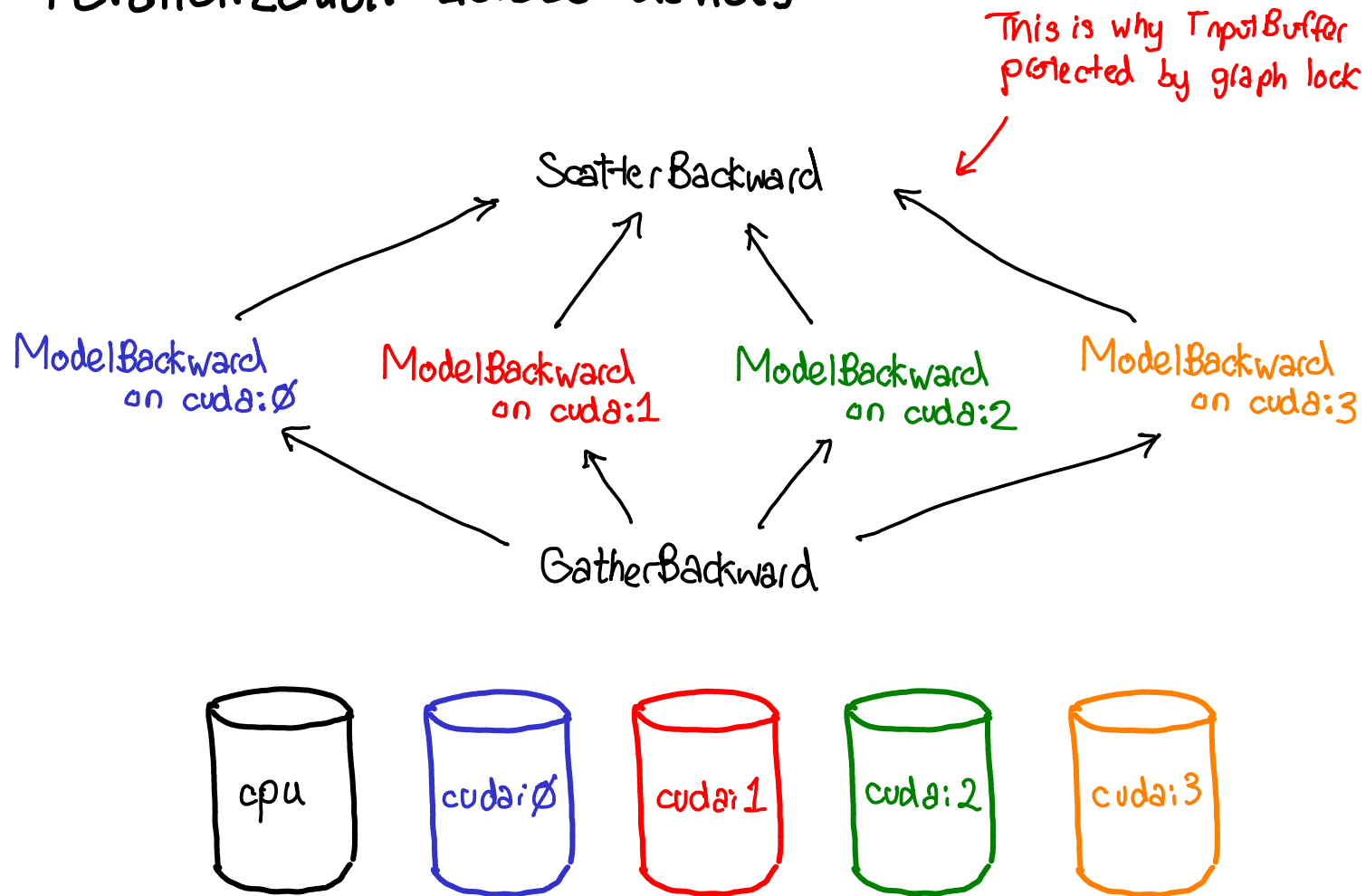
There's a preprocessing scan of the graph to determine how many dependencies each FunctionTask has



1. Evaluate function on InputBuffer
  - a. Call pre-hooks on inputs
  - b. Call the function for real
  - c. Call post-hooks on outputs
2. Release saved variables
3. **With graph lock**, accumulate grad into InputBuffers of next functions, and add newly ready FunctionTasks to ReadyQueue



# Parallelization across devices



Mechanics

Finding your way around

# Source structure

tests  
&  
code  
generation

test/  
tools/

torch/

Python; the frontend!

torch/csrc/

Python bindings, this functionality

api/

C++ API

autograd/

jit/

TorchScript

aten/src/ATen/

Tensor operator implementations

native/ Modern

../TH, THC, ...

Legacy

c10/

Core abstractions (e.g. TensorImpl)

# Anatomy of an operator call

How to figure this out?  
Build PyTorch with `DEBUG=1`,  
set a breakpoint on  
`at::native::add`, and look  
at the backtrace!

`torch.add(x, y)`

- Python argument parsing

`THPVariable_add`  
`python_variable_methods.cpp`

- VariableType

`VariableType::add`

- Type

`TypeDefault::add (or, CPUFloatType)`

- Native function

`at::native::add`  
`BinaryOps.cpp`

- TH function

`THTensor_add`  
`generic/THTensorMath.cpp`, e.g.

```

static PyObject * THPVariable_add(PyObject* self_, PyObject* args, PyObject* kwargs)
{
    HANDLE_TH_ERRORS
    static PythonArgParser parser({
        "add(Tensor input, Scalar alpha, Tensor other, *, Tensor out=None)|deprecated",
        "add(Tensor input, Tensor other, *, Scalar alpha=1, Tensor out=None)",
    }, /*traceable=*/true);

    ParsedArgs<4> parsed_args;
    auto r = parser.parse(args, kwargs, parsed_args);

    if (r.idx == 0) {
        if (r.isNone(3)) {
            return wrap(dispatch_add(r.tensor(0), r.scalar(1), r.tensor(2)));
        } else {
            return wrap(dispatch_add(r.tensor(0), r.scalar(1), r.tensor(2), r.tensor(3)));
        }
    } else if (r.idx == 1) {
        if (r.isNone(3)) {
            return wrap(dispatch_add(r.tensor(0), r.tensor(1), r.scalar(2)));
        } else {
            return wrap(dispatch_add(r.tensor(0), r.tensor(1), r.scalar(2), r.tensor(3)));
        }
    }
    Py_RETURN_NONE;
    END_HANDLE_TH_ERRORS
}

```

auto-generated! just skim

argument parser

release the GIL

AutoNoGIL no\_gil;  
return self.add(other, alpha);

rewrap into PyObject

actual binding  
on torch.\_C.VariableFunctions

```

static PyMethodDef torch_functions[] = {
    ...
    {"add", (PyCFunction)THPVariable_add, METH_VARARGS | METH_KEYWORDS | METH_STATIC, N
    ...

```

File: torch/csrc/autograd/generated/python\_torch\_functions\_dispatch.h



```
inline Tensor Tensor::add(const Tensor & other, Scalar alpha) const {
    return type().add(*this, other, alpha);
}
```

File: aten/src/ATen/core/TensorMethods.h

CPUFloatType

virtual method

might be overridden by subclass

Switch (e.g. CUDA) device

```
Tensor TypeDefault::add(const Tensor & self, const Tensor & other, Scalar alpha) const {
    const OptionalDeviceGuard device_guard(device_of(self));
    return at::native::add(/* native_actuals */ self, other, alpha);
}
```

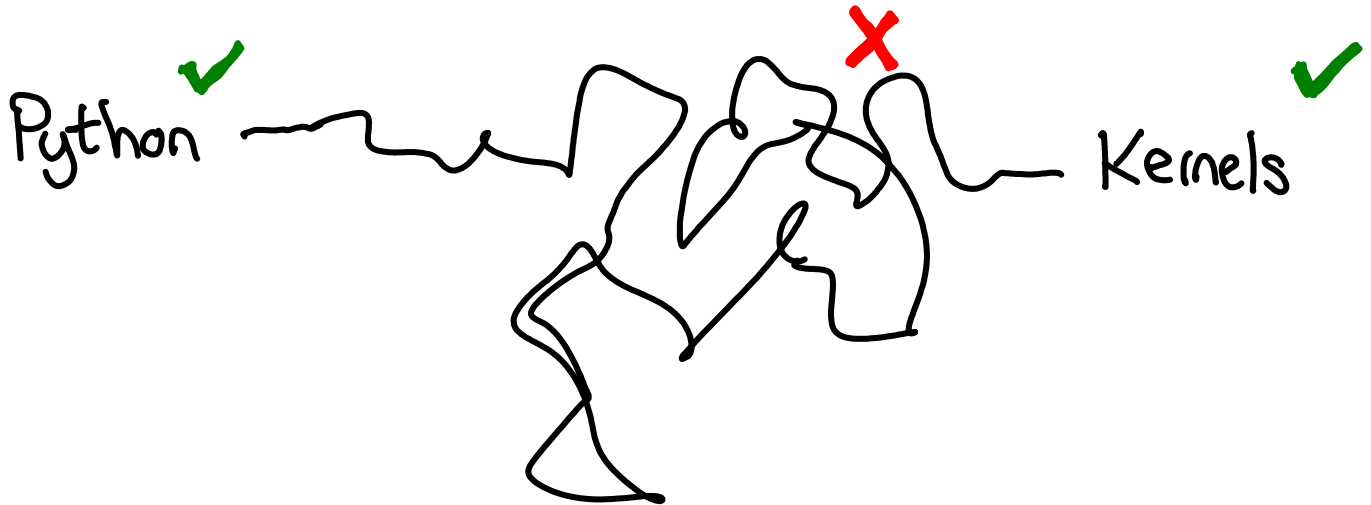
File: build/aten/src/ATen/TypeDefault.cpp

```
Tensor add(const Tensor& self, const Tensor& other, Scalar alpha) {
    Tensor result;
    if (other.is_sparse()) {
        result = at::empty({0}, self.options());
        return native::add_out(result, self, other, alpha);
    }
    auto iter = TensorIterator::binary_op(result, self, other);
    add_stub(iter->device_type(), *iter, alpha);
    return iter->output();
}
```

the kernel proper

File: aten/src/ATen/native/BinaryOps.cpp

I don't recommend spending too much time mucking around backtraces



Writing kernels

# Anatomy of a kernel

(+ Infra wiring up!)

```
Tensor my_op_out_cpu(Tensor& result, const Tensor& self, const Tensor& other) {
```

```
    TORCH_CHECK(result.is_cpu() && self.is_cpu() && other.is_cpu());
```

```
    TORCH_CHECK(self.dim() == 1);
```

```
    TORCH_CHECK(self.sizes() == other.sizes());
```

← Error checking

```
    result.resize(self.sizes());
```

← Output allocation

```
    AT_DISPATCH_FORALL_TYPES(
```

```
        self.scalar_type(), "my_op_cpu", [&] {
```

```
            my_op_cpu_kernel<scalar_t>(result, self, other);
```

```
        }
```

```
    );
```

```
}
```

← Dtype dispatch

```
template <typename scalar_t>
```

```
void my_op_cpu_kernel(Tensor& result, const Tensor& self, const Tensor& other) {
```

```
    auto result_accessor = result.accessor<scalar_t, 1>();
```

```
    auto self_accessor = self.accessor<scalar_t, 1>();
```

```
    auto other_accessor = other.accessor<scalar_t, 1>();
```

```
    parallel_for(0, self.size(0), 0, [&](int64_t start, int64_t end) {
```

```
        ... self_accessor[i] ...
```

```
    });
```

```
}
```

Parallelization

← Data access (vectorized?!)

# Scaffolding

Every operator has a schema

↙ Signature

```
- func: add(Tensor self, Tensor other, *, Scalar alpha=1) -> Tensor
```

```
variants: function, method
```

```
dispatch:
```

```
  CPU: add
```

```
  CUDA: add
```

```
  SparseCPU: add
```

```
  SparseCUDA: add
```

```
  MklDnnCPU: mklDnn_add
```

↖ control Python/C++ API

} dispatch

[aten/src/ATen/native/native\\_functions.yaml](#)

This schema controls our code generation

# Scaffolding

Also a derivative:

```
- name: add(Tensor self, Tensor other, *, Scalar alpha)  
  self: grad  
  other: maybe_multiply(grad, alpha)
```

signature

} one grad per Tensor argument

the gradient flowing into add

grad\_self = grad\_add

# Error checking

## Low level

Hand-write your error message.

```
#include <cl0/util/Exception.h>
```

```
TORCH_CHECK(self.dim() == 1, "Expected self to be 1-D tensor, but "  
"was ", self.dim(), "-D tensor")
```

Condition

Intersperse values directly; they are formatted with operator <<

## Higher level

```
#include <ATen/TensorUtils.h>
```

```
TensorArg result_arg{result, "result", 0}  
          self_arg{self, "self", 1},  
          other_arg{other, "other", 2};
```

Metadata about each argument


```
CheckedFrom c = "my_op_cpu";  
checkDim(c, self_arg, 1);
```

Pre-canned error message

Not complete!

# Output allocation

no-op if already the right size



```
Tensor& abs_out(Tensor& result, const Tensor& self) {  
    result.resize_(self.sizes());  
    // ... the real implementation ...  
}
```

preallocated  
output

```
Tensor abs(const Tensor& self) {  
    Tensor result = at::empty({0}, self.options());  
    abs_out(result, self);  
    return result;  
}
```

↑ assume \_out  
will resize

generate  
for me

```
Tensor& abs_(Tensor& self) {  
    return abs_out(self, self);  
}
```

inplace

↑ assumes \_out kernel handles  
aliasing inputs



# Dtype dispatch

↙ not actually all types

```
AT_DISPATCH_ALL_TYPES(  
    self.scalar_type(), "my_op_cpu", [&] {  
        my_op_cpu_kernel<scalar_t>(result, self, other);  
    }  
);
```

↖ specialize lambda for each scalar\_type

```
AT_DISPATCH_ALL_TYPES(TYPE, NAME, ...)  
AT_DISPATCH_FLOATING_TYPES(TYPE, NAME, ...)  
AT_DISPATCH_INTEGRAL_TYPES(TYPE, NAME, ...)  
AT_DISPATCH_ALL_TYPES_AND(SCALARTYPE, TYPE, NAME, ...)
```

# Data access

## Point access

```
auto x_accessor = x.accessor<float, 3>();  
float val = x_accessor[0][0][0];
```

CUDA: `packed_accessor` (gotcha: 32-bit!)

## Elementwise

```
#include <ATen/native/TensorIterator.h>  
auto iter = TensorIterator::Builder()  
    .add_output(output)  
    .add_input(input)  
    .build();  
binary_kernel(iter, [](float a, float b) {  
    return a + b;  
});
```

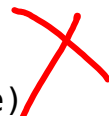
CUDA: `gpu_binary_kernel`

## Vectorized

```
#include <ATen/native/cpu/Loops.h>  
binary_kernel_vec(iter,  
    [=](scalar_t a, scalar_t b) -> scalar_t {  
        return a + alpha * b; },  
    [=](Vec256<scalar_t> a, Vec256<scalar_t> b) {  
        return vec256::fmadd(b, alpha_vec, a);  
    });
```

# Parallelization (CPU only)

```
#pragma omp parallel for private(i) if (input_size > grain_size)
for (i = 0; i < input_size; i++) {
    ...
}
```



```
at::parallel_for(0, input_size, grain_size, [&](int64_t begin, int64_t end) {
    for (auto i = start; i < end; i++) {
        // ...
    }
});
```



critical? implement the locks yourself!

# Legacy code

TH, THC, THNN, THCUNN

- C code

THTensor\_wrap()

- Manual recounting

- Preprocessor shenanigans

Write new code in C++, consider porting old code when you can!

Workflow efficiency

- Don't edit headers
- Don't test by CI
- Do setup ccache
- Get a beefy workstation for builds

# How to get involved!

<https://github.com/pytorch/pytorch/issues>

high priority

module

small

Port some code from TH to ATen!

Help us improve our documentation!

Help minimize reproducers on bugs!

Help answer questions on the forums

Help us discuss design of new features!

Scratch your own itches!