

Triton

The challenges of GPU programming

The architecture of modern GPUs can be roughly divided into three major components—DRAM (Dynamic Random Access Memory), SRAM (Static Random Access Memory) and ALUs (Arithmetic Logic Units)—each of which must be considered when optimizing CUDA code:

- Memory transfers from DRAM must be coalesced into large transactions to leverage the large bus width of modern memory interfaces.
- Data must be manually stashed to SRAM prior to being re-used, and managed so as to minimize shared memory bank conflicts upon retrieval.
- Computations must be partitioned and scheduled carefully, both across and within Streaming Multiprocessors (SMs), so as to promote instruction/thread-level parallelism and leverage special-purpose ALUs (e.g., tensor cores).

Overview of Triton

Triton is a powerful open-source framework for simplifying GPU programming, especially for neural networks. By abstracting low-level CUDA intricacies, Triton enables researchers to write high-performance GPU code in a Python-like syntax without needing deep CUDA expertise.

- High Efficiency: Triton kernels achieve performance comparable to or better than hand-optimized CUDA code, e.g., cuBLAS-level FP16 matrix multiplication in under 25 lines.
- Python-like Syntax: Makes GPU programming accessible with familiar and simple constructs.

Overview of Triton

Feature	Description	CUDA	TRITON
Memory Coalescing	Ensures that memory accesses by threads are grouped into fewer, more efficient transactions to optimize global memory bandwidth. Requires explicit memory alignment and access patterns in code.	Manual	Automatic
Shared Memory Management	Refers to the explicit allocation and management of on-chip shared memory to speed up frequent memory accesses by threads in a block. In	Manual	Automatic
Scheduling (Within SMs)	Handles the execution order of threads within a Streaming Multiprocessor (SM).	Manual	Automatic
Scheduling (Across SMs)	Refers to how workloads are distributed across multiple Streaming Multiprocessors (SMs) on the GPU.	Manual	Manual

Basic Syntax of Triton

```
0 import triton
1 import triton.language as tl
2 @triton.jit
3 def add_kernel(
4     x_ptr, # *Pointer* to first input vector.
5     y_ptr, # *Pointer* to second input vector.
6     output_ptr, # *Pointer* to output vector.
7     n_elements, # Size of the vector.
8     BLOCK_SIZE: tl.constexpr, # Number of elements each program should process.
9 ):
10     # There are multiple 'programs' processing different data. We identify which program we are here:
11     pid = tl.program_id(axis=0) # We use a 1D launch grid so axis is 0.
12     # This program will process inputs that are offset from the initial data. For instance, if you had a vector of
13     # length 256 and block_size of 64, the programs would each access the elements
14     # [0:64, 64:128, 128:192, 192:256]. Note that offsets is a list of pointers:
15     block_start = pid * BLOCK_SIZE
16     offsets = block_start + tl.arange(0, BLOCK_SIZE)
17     # Create a mask to guard memory operations against out-of-bounds accesses.
18     mask = offsets < n_elements
19     # Load x and y from DRAM, masking out any extra elements in case the input is not a multiple of the block size.
20     x = tl.load(x_ptr + offsets, mask=mask)
21     y = tl.load(y_ptr + offsets, mask=mask)
22     output = x + y
23     # Write x + y back to DRAM.
24     tl.store(output_ptr + offsets, output, mask=mask)
```

Triton Semantics - Type Promotion

1. **Kind Hierarchy:** `{bool} < {integer} < {floating-point}`

Example: `(int32, bfloat16) -> bfloat16`

2. **Width:** For the same kind, the higher width is chosen.

Example: `(float32, float16) -> float32`

3. **Prefer float16:** If the same width and signedness but different types, promote to `float16`.

Example: `(float16, bfloat16) -> float16`

4. **Prefer unsigned:** For same width but different signedness, promote to unsigned.

Example: `(int32, uint32) -> uint32`

Triton Semantics - Broadcasting

Operations on tensors of different shapes expand them to compatible shapes without copying data:

- Shorter shapes are left-padded with ones.

Example: `((3, 4), (5, 3, 4)) -> ((1, 3, 4), (5, 3, 4))`

- Dimensions are compatible if equal or one is 1.

Example: `((1, 3, 4), (5, 3, 4)) -> ((5, 3, 4), (5, 3, 4))`

Triton Semantics - Rounding

- **Integer division and modulus:** Follow C semantics (rounding towards zero), unlike Python's rounding towards minus infinity.

Example: `-7 // 3 = -2` (Triton) vs. `-7 // 3 = -3` (Python).

- **Scalar operations:** Use Python semantics for integer division and modulus when all inputs are scalars.

Basic Syntax of Triton

```
0 import triton
1 import triton.language as tl
2 @triton.jit
3 def add_kernel(
4     x_ptr, # *Pointer* to first input vector.
5     y_ptr, # *Pointer* to second input vector.
6     output_ptr, # *Pointer* to output vector.
7     n_elements, # Size of the vector.
8     BLOCK_SIZE: tl.constexpr, # Number of elements each program should process.
9 ):
10     # There are multiple 'programs' processing different data. We identify which program we are here:
11     pid = tl.program_id(axis=0) # We use a 1D launch grid so axis is 0.
12     # This program will process inputs that are offset from the initial data. For instance, if you had a vector of
13     # length 256 and block_size of 64, the programs would each access the elements
14     # [0:64, 64:128, 128:192, 192:256]. Note that offsets is a list of pointers:
15     block_start = pid * BLOCK_SIZE
16     offsets = block_start + tl.arange(0, BLOCK_SIZE)
17     # Create a mask to guard memory operations against out-of-bounds accesses.
18     mask = offsets < n_elements
19     # Load x and y from DRAM, masking out any extra elements in case the input is not a multiple of the block size.
20     x = tl.load(x_ptr + offsets, mask=mask)
21     y = tl.load(y_ptr + offsets, mask=mask)
22     output = x + y
23     # Write x + y back to DRAM.
24     tl.store(output_ptr + offsets, output, mask=mask)
```


Basic Syntax of Triton - `program_id` and `num_programs`

In Triton, `pid` refers to the program ID, which represents the ID of the current program instance along a specified axis in a 3D launch grid. The `program_id` function in Triton is used to identify which program is currently executing and helps distribute work across threads or blocks in parallel computations.

axis: The dimension of the grid along which the program ID is calculated. It can be 0, 1, or 2.

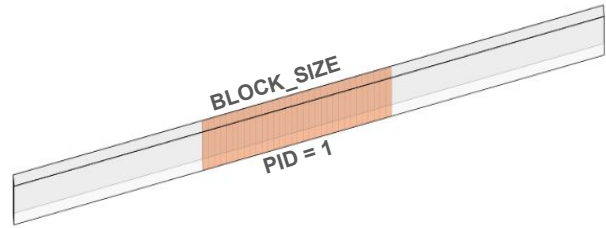
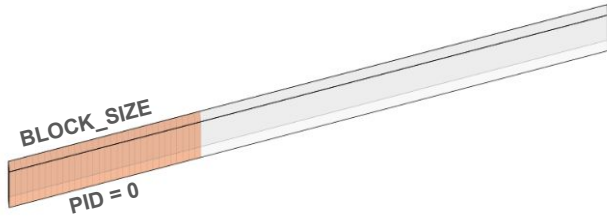
```
pid = tl.program_id(axis=0)
```

```
0 @triton.jit
1 def demo(x_ptr):
2     print(tl.num_programs(0), tl.num_programs(1), tl.num_programs(2))
3
4 x = torch.ones(2, 4, 4)
5 triton_viz.trace(demo)[(2, 2, 1)](x)
```

```
[2] [2] [1]
[2] [2] [1]
[2] [2] [1]
[2] [2] [1]
```

Basic Syntax of Triton - BLOCK_SIZE: `tl.constexpr`

```
block_start = pid * BLOCK_SIZE  
offsets = block_start + tl.arange(0, BLOCK_SIZE)  
mask = offsets < n_elements
```

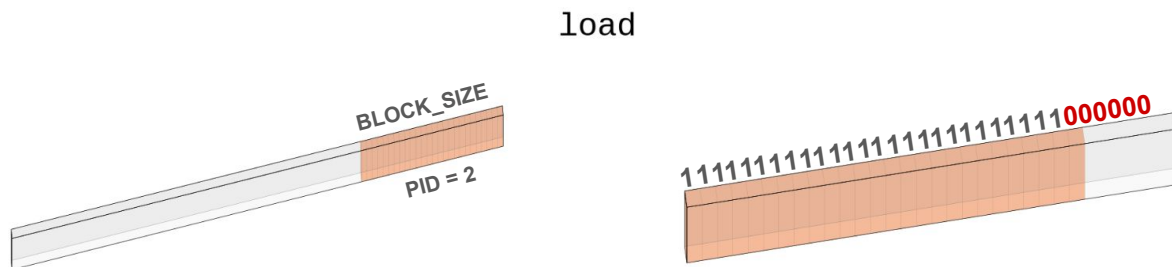


Basic Syntax of Triton - mask

The mask is used to safeguard memory operations from out-of-bounds access.

In this case:

- **mask = offsets < n_elements** determines which indices are within the array size n_elements.
- During data loading and storing (**tl.load**, **tl.store**), the mask prevents access to non-existent elements, avoiding errors.



Basic Syntax of Triton - **tl.load**, **tl.store**

```
x = tl.load(x_ptr + offsets, mask=mask)
y = tl.load(y_ptr + offsets, mask=mask)

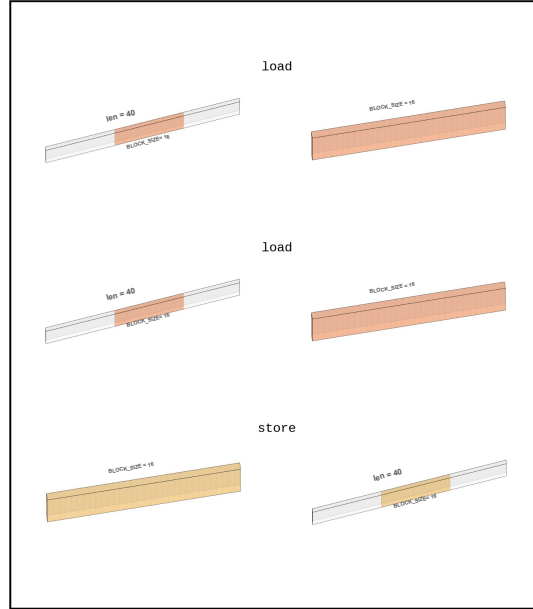
output = x + y

tl.store(output_ptr + offsets, output, mask=mask)
```

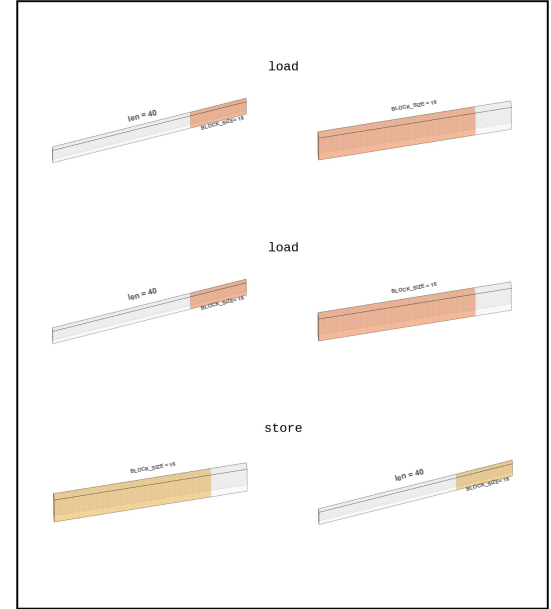
Basic Syntax of Triton - **tl.load**, **tl.store**



PID = 0



PID = 1



PID = 2

Basic Syntax of Triton

```
0 import triton
1 import triton.language as tl
2 @triton.jit
3 def add_kernel(
4     x_ptr,
5     y_ptr,
6     output_ptr,
7     n_elements,
8     BLOCK_SIZE: tl.constexpr,
9 ):
10     pid = tl.program_id(axis=0)
11
12     block_start = pid * BLOCK_SIZE
13     offsets = block_start + tl.arange(0, BLOCK_SIZE)
14     mask = offsets < n_elements
15
16     x = tl.load(x_ptr + offsets, mask=mask)
17     y = tl.load(y_ptr + offsets, mask=mask)
18
19     output = x + y
20
21     tl.store(output_ptr + offsets, output, mask=mask)
```

```
print(x_ptr.dtype)           # pointer<fp32>
print(y_ptr.dtype)           # pointer<fp32>
print(output_ptr.dtype)      # pointer<fp32>
print(type(n_elements))      # <class 'int'>
print(type(BLOCK_SIZE))      # <class 'int'>

print(pid.dtype)              # int32

print(block_start.dtype)      # int32
print(offsets.dtype)          # int32
print(mask.dtype)             # int1

print(x.dtype)                # fp32
print(output.dtype)           # fp32
print((z_ptr + offsets).dtype) # pointer<fp32>
```

Memory allocation and pointer arithmetic in C - number.

```
#include <stdio.h>
int main(void)
{
    int n = 10;
    int *ptr = &n;
    printf("address=%p \t value=%d \n", (void*)ptr, *ptr); //address=0060FEA8    value=10
    ptr++;
    printf("address=%p \t value=%d \n", (void*)ptr, *ptr); //address=0060FEAC    value=6356652
    ptr--;
    printf("address=%p \t value=%d \n", (void*)ptr, *ptr); //address=0060FEA8    value=10
    return 0;
}
```

A number in C is represented as a pointer of a specific size.
In the case of int32, it is 4 bytes.

Memory allocation and pointer arithmetic in C - array.

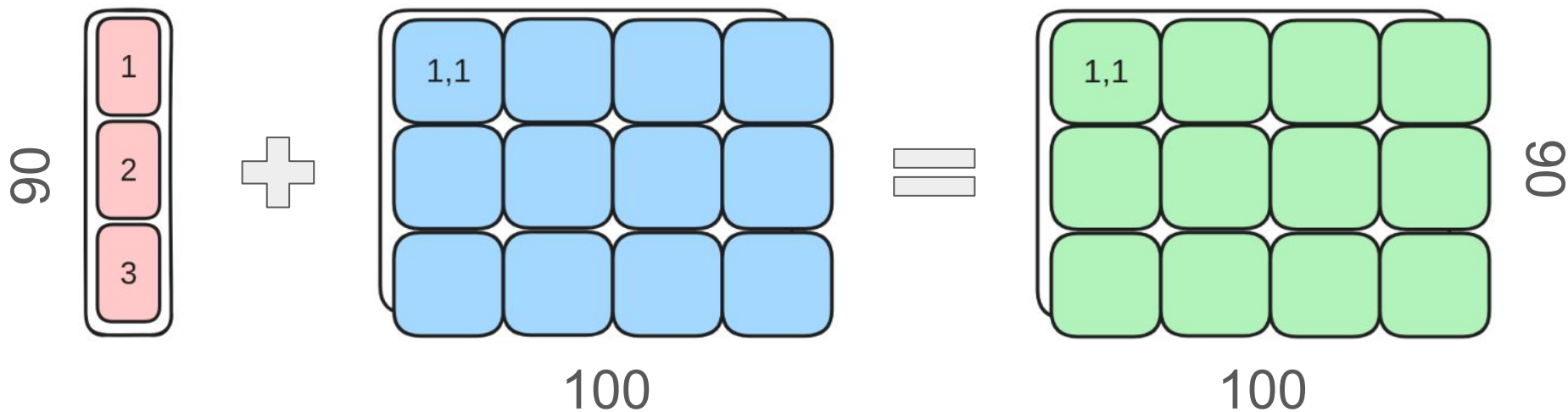
```
#include <stdio.h>
int main(void)
{
    int numbers[] = {11, 12, 13, 14};
    for(size_t i=0; i< sizeof(numbers) / sizeof(numbers[0]); ++i){
        printf("numbers[%d]: %p; value: %d\n", i, numbers+i, *(numbers+i));
    }
    return 0;
}
```

```
numbers[0]: 0x7ffffb4d10140; value: 11
numbers[1]: 0x7ffffb4d10144; value: 12
numbers[2]: 0x7ffffb4d10148; value: 13
numbers[3]: 0x7ffffb4d1014c; value: 14
```

An array in C represents a contiguous block of memory where its elements are stored. The size of the array is determined by **n*sizeof(data_type)**. Pointer arithmetic applies to this structure.

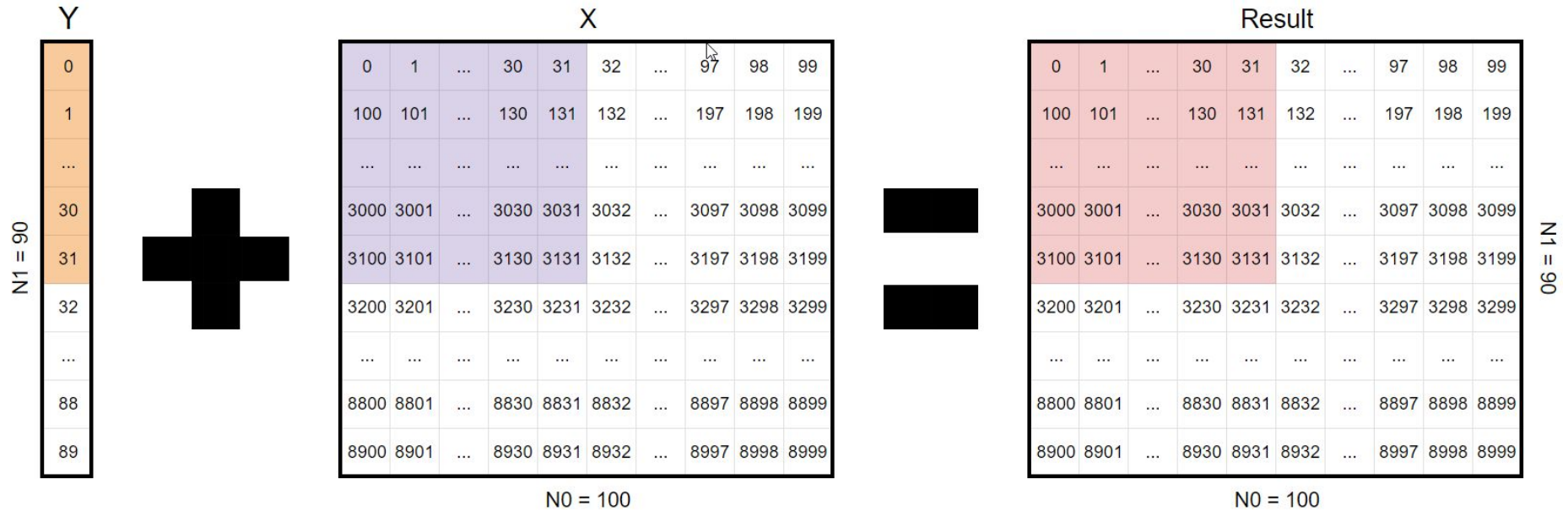
In other words, an **array in C is a pointer to the first element of the array with pointer arithmetic.**

Triton pointer arithmetic - 2D arrays



Triton pointer arithmetic - 2D arrays

$B0 = B1 = 32$
 $N0 = 100$ $N1 = 90$



Triton pointer arithmetic - 2D arrays

```
0 @triton.jit
1 def mul_relu_block_back_kernel(x_ptr, y_ptr, dz_ptr, dx_ptr, N0, N1, B0: tl.constexpr, B1: tl.2 constexpr):
2
3     pid_x = tl.program_id(0) # program ID along the X-axis (axis 0)
4     pid_y = tl.program_id(1) # program ID along the y-axis (axis 1)
5
6     offset_x = tl.arange(0, B0) + pid_x * B0 # offset_x = [0, 1, 2, ..., 31]
7     offset_y = tl.arange(0, B1) + pid_y * B1 # offset_y = [0, 1, 2, ..., 31]
8
9     # 2D array [[0, 1, 2, ..., 31], [100, 101, 102, ..., 131], ..., [3100, 3101, ..., 3131]]
10    offset = offset_x[None, :] + offset_y[:, None] * N0
11
12    mask_x = offset_x < N0 # N0 = 100; array of 32 bool values
13    mask_y = offset_y < N1 # N1 = 90; array of 32 bool values
14
15    mask = mask_x[None, :] & mask_y[:, None] # array of 32x32 bool values
16
17    y = tl.load(y_ptr + offset_y, mask_y, 0) # 1-d pointer array with length = 32
18    x = tl.load(x_ptr + offset, mask, 0) # 1-d pointer array with length = 32x32
19
20    res = x + y[:, None] # 1-d pointer array with length = 32x32
21
22    tl.store(dz_ptr + offset, res, mask)
```

Calling a Triton Kernel from Python

```
0 import triton
1 import triton.language as tl
2 @triton.jit
3 def add_kernel(
4     x_ptr,
5     y_ptr,
6     output_ptr,
7     n_elements,
8     BLOCK_SIZE: tl.constexpr,
9 ):
10     pid = tl.program_id(axis=0)
11
12     block_start = pid * BLOCK_SIZE
13     offsets = block_start + tl.arange(0, BLOCK_SIZE)
14     mask = offsets < n_elements
15
16     x = tl.load(x_ptr + offsets, mask=mask)
17     y = tl.load(y_ptr + offsets, mask=mask)
18
19     output = x + y
20
21     tl.store(output_ptr + offsets, output, mask=mask)
```

```
23 def add(x: torch.Tensor, y: torch.Tensor):
24     assert x.is_cuda and y.is_cuda
25     output = torch.empty_like(x)
26     n_elements = x.numel()
27     grid = lambda meta: (triton.cdiv(n_elements, meta['BLOCK_SIZE']),)
28     add_kernel[grid](x, y, output, n_elements, BLOCK_SIZE=1024)
29     return output
```

```
31 if __name__ == "__main__":
32     x = torch.tensor([1, 2, 3, 4], dtype=torch.float32, device='cuda')
33     y = torch.tensor([10, 20, 30, 40], dtype=torch.float32, device='cuda')
34     output = add(x, y)
35     print(output) # tensor([11., 22., 33., 44.], device='cuda:0')
```

- ! 1) How many values can a pid take?
- 2) What does the mask array look like?

Calling a Triton Kernel from Python

```
23 def add(x: torch.Tensor, y: torch.Tensor):
24     # We need to preallocate the output.
25     output = torch.empty_like(x)
26     assert x.device == DEVICE and y.device == DEVICE and output.device == DEVICE
27     n_elements = output.numel()
28     # The SPMD launch grid denotes the number of kernel instances that run in parallel.
29     # It is analogous to CUDA launch grids. It can be either Tuple[int], or Callable(metaparameters) -> Tuple[int].
30     # In this case, we use a 1D grid where the size is the number of blocks:
31     grid = lambda meta: (triton.cdiv(n_elements, meta['BLOCK_SIZE']), )
32     # - Each torch.tensor object is implicitly converted into a pointer to its first element.
33     # - `triton.jit`ed functions can be indexed with a launch grid to obtain a callable GPU kernel.
34     # - Don't forget to pass meta-parameters as keywords arguments.
35     add_kernel[grid](x, y, output, n_elements, BLOCK_SIZE=1024)
36     # We return a handle to z but, since `torch.cuda.synchronize()` hasn't been called, the kernel is still
37     # running asynchronously at this point.
38     return output
```

Basic Syntax of Triton

- **Jit** Decorator for JIT-compiling a function using the Triton compiler.
- **autotune** Decorator for auto-tuning a triton.jit'd function.
- **heuristics** Decorator for specifying how the values of certain meta-parameters may be computed.
- **Config** An object that represents a possible kernel configuration for the auto-tuner to try.

```
0  @triton.autotune(  
1      configs=[  
2          triton.Config({"BLOCK_SIZE": 1024}, num_warps=4),  
3          triton.Config({"BLOCK_SIZE": 2048}, num_stages=1),  
4      ],  
5      key=["n_elements"],  
6  )  
7  @triton.jit  
8  def _quantize_global(  
9      x_ptr,  
10     absmax_inv_ptr,  
11     output_ptr,  
12     n_elements,  
13     BLOCK_SIZE: tl.constexpr,  
14 ):  
15     pid = tl.program_id(axis=0)  
16     block_start = pid * BLOCK_SIZE  
17     offsets = block_start + tl.arange(0, BLOCK_SIZE)  
18     mask = offsets < n_elements  
19     x = tl.load(x_ptr + offsets, mask=mask)  
20     absmax_inv = tl.load(absmax_inv_ptr)  
21     output = tl.libdevice.llrint(127.0 * (x * absmax_inv))  
22     tl.store(output_ptr + offsets, output, mask=mask)  
23  
24 def quantize_global(x: torch.Tensor):  
25     absmax = x.abs().max().unsqueeze(0)  
26     absmax_inv = 1.0 / absmax  
27     output = torch.empty(*x.shape, device="cuda", dtype=torch.int8)  
28     assert x.is_cuda and output.is_cuda  
29     n_elements = output.numel()  
30     grid = lambda meta: (triton.cdiv(n_elements, meta["BLOCK_SIZE"]),)  
31     _quantize_global[grid](x, absmax_inv, output, n_elements)  
32     return output, absmax
```

Basic Syntax of Triton - Config

- **kwargs** – a dictionary of meta-parameters to pass to the kernel as keyword arguments.
- **num_warps** – the number of warps to use for the kernel when compiled for GPUs. For example, if num_warps=8, then each kernel instance will be automatically parallelized to cooperatively execute using $8 * 32 = 256$ threads.
- **num_stages** – the number of stages that the compiler should use when software-pipelining loops.

```
configs=[  
    triton.Config({"BLOCK_SIZE": 1024}, num_warps=4),  
    triton.Config({"BLOCK_SIZE": 2048}, num_stages=1),  
]
```

Basic Syntax of Triton - heuristics

Decorator for specifying how the values of certain meta-parameters may be computed. This is useful for cases where auto-tuning is prohibitively expensive, or just not applicable.

```
@triton.heuristics(values={'BLOCK_SIZE': lambda args: triton.next_power_of_2(args['x_size'])})
```

values (**dict**[**str**, **Callable**[[**list**[**Any**]], **Any**]]) – a dictionary of meta-parameter names and functions that compute the value of the meta-parameter. each such function takes a list of positional arguments as input.

Basic Syntax of Triton - autotune

Decorator for auto-tuning a triton.jit function.

1. **configs** (list[triton.Config]):

- A list of `triton.Config` objects defining configurations.

2. **key** (list[str]):

- Names of argument(s) that trigger re-evaluation of configurations.

3. **prune_configs_by**:

- A dictionary of pruning functions:
 - **'perf_model'**: Predicts running time of different configurations.
 - **'top_k'**: Selects top-k configurations based on performance.
 - **'early_config_prune' (optional)**: Prunes configurations based on initial criteria.

4. **reset_to_zero** (list[str]):

- A list of argument names to reset to zero before evaluating configurations.

5. **restore_value** (list[str]):

- A list of argument names to restore to their original values after configurations are evaluated.

6. **pre_hook** (lambda args, reset_only):

- A function called before kernel execution.
 - **args**: Dictionary of all arguments passed to the kernel.
 - **reset_only**: Boolean indicating whether only reset logic should be applied.

7. **post_hook** (lambda args, exception):

- A function called after kernel execution.
 - **args**: Dictionary of all arguments passed.
 - **exception**: Exception raised by the kernel if any error occurs.

8. **do_bench** (lambda fn, quantiles):

- A benchmark function to measure kernel execution time at various quantiles.

If the environment variable **TRITON_PRINT_AUTOTUNING=1**, Triton will print a message to stdout after autotuning each kernel, including the time spent autotuning and the best configuration.

Calling a Triton Kernel from Python

```
0 import triton
1 import triton.language as tl
2 @triton.jit
3 def add_kernel(
4     x_ptr,
5     y_ptr,
6     output_ptr,
7     n_elements,
8     BLOCK_SIZE: tl.constexpr,
9 ):
10     pid = tl.program_id(axis=0)
11
12     block_start = pid * BLOCK_SIZE
13     offsets = block_start + tl.arange(0, BLOCK_SIZE)
14     mask = offsets < n_elements
15
16     x = tl.load(x_ptr + offsets, mask=mask)
17     y = tl.load(y_ptr + offsets, mask=mask)
18
19     output = x + y
20
21     tl.store(output_ptr + offsets, output, mask=mask)
```

```
23 def add(x: torch.Tensor, y: torch.Tensor):
24     assert x.is_cuda and y.is_cuda
25     output = torch.empty_like(x)
26     n_elements = x.numel()
27     grid = lambda meta: (triton.cdiv(n_elements, meta['BLOCK_SIZE']),)
28     add_kernel[grid](x, y, output, n_elements, BLOCK_SIZE=1024)
29     return output
```

```
31 if __name__ == "__main__":
32     x = torch.tensor([1, 2, 3, 4], dtype=torch.float32, device='cuda')
33     y = torch.tensor([10, 20, 30, 40], dtype=torch.float32, device='cuda')
34     output = add(x, y)
35     print(output) # tensor([11., 22., 33., 44.], device='cuda:0')
```

Basic Syntax of Triton - autotune

```
0 import torch
1 import triton
2 import triton.language as tl
3
4 @triton.autotune(
5     configs=[
6         triton.Config({"BLOCK_SIZE": 128}, num_stages=3, num_warps=8),
7         triton.Config({"BLOCK_SIZE": 256}, num_stages=3, num_warps=8),
8         triton.Config({"BLOCK_SIZE": 512}, num_stages=4, num_warps=4),
9         triton.Config({"BLOCK_SIZE": 64}, num_stages=4, num_warps=4),
10        triton.Config({"BLOCK_SIZE": 32}, num_stages=4, num_warps=4),
11        triton.Config({"BLOCK_SIZE": 32}, num_stages=2, num_warps=2),
12        triton.Config({"BLOCK_SIZE": 16}, num_stages=4, num_warps=4),
13    ],
14    key=[], # No arguments are used to select configurations
15 )
16 @triton.jit
17 def add_kernel(
18     x_ptr,
19     y_ptr,
20     output_ptr,
21     n_elements,
22     BLOCK_SIZE: tl.constexpr,
23 ):
24     pid = tl.program_id(axis=0)
25     block_start = pid * BLOCK_SIZE
26     offsets = block_start + tl.arange(0, BLOCK_SIZE)
27     mask = offsets < n_elements
28     x = tl.load(x_ptr + offsets, mask=mask)
29     y = tl.load(y_ptr + offsets, mask=mask)
30     output = x + y
31     tl.store(output_ptr + offsets, output, mask=mask)
32
```

```
33 def add(x: torch.Tensor, y: torch.Tensor):
34     assert x.is_cuda and y.is_cuda
35     output = torch.empty_like(x)
36     n_elements = x.numel()
37     grid = lambda meta: (triton.cdiv(n_elements, meta['BLOCK_SIZE']),)
38     add_kernel[grid](x_ptr=x, y_ptr=y, output_ptr=output, n_elements=n_elements)
39     return output
40
41 if __name__ == "__main__":
42     x = torch.tensor([1, 2, 3, 4], dtype=torch.float32, device='cuda')
43     y = torch.tensor([10, 20, 30, 40], dtype=torch.float32, device='cuda')
44     output = add(x, y)
45     print(output) # tensor([11., 22., 33., 44.], device='cuda:0')
```

NB: When using autotune in Triton, you do not need to explicitly specify the following parameters, as they are automatically determined: BLOCK_SIZE, num_stages, num_warps.

Thus, these parameters are handled internally by the configuration when using autotune.

Triton autotuning for function add_kernel finished after 2.38s; best config selected: BLOCK_SIZE: 32, num_warps: 4, num_ctas: 1, num_stages: 4, maxnreg: None; tensor([11., 22., 33., 44.], device='cuda:0')

Basic Syntax of Triton - inline_asm_elementwise

Generate vectorized assembly code for tensors. The input tensors are assumed to have the same shape after broadcasting.

The output type can be a tuple of types, with each element being a separate tensor. The assembly processes blocks of data, handling multiple elements at once. An empty return value is not allowed – at least one tensor must be returned for compatibility.

```
0 import torch
1 import triton
2 import triton.language as tl
3 @triton.jit
4 def kernel(A, B, C, D, BLOCK: tl.constexpr):
5     a = tl.load(A + tl.arange(0, BLOCK))
6     b = tl.load(B + tl.arange(0, BLOCK))
7
8     (c, d) = tl.inline_asm_elementwise(
9         asm="""
10         {
11             .reg .b8 tmp<4>;
12             mov.b32 {tmp0, tmp1, tmp2, tmp3}, $0;
13             cvt.u32.u8 $0, tmp0;
14             cvt.u32.u8 $1, tmp1;
15             cvt.u32.u8 $2, tmp2;
16             cvt.u32.u8 $3, tmp3;
17         }
18         cvt.rn.f32.s32 $4, $0;
19         cvt.rn.f32.s32 $5, $1;
20         cvt.rn.f32.s32 $6, $2;
21         cvt.rn.f32.s32 $7, $3;
22         max.f32 $4, $4, $9;
23         max.f32 $5, $5, $10;
24         max.f32 $6, $6, $11;
25         max.f32 $7, $7, $12;
26         """,
27         constraints=(
28             "r,r,r,r,r,r,r,r,r,r,r,r",
29             "r,r,r,r,r,r"),
30         args=[a, b],
31         dtype=(tl.int32, tl.float32),
32         is_pure=True,
33         pack=4,
34     )
35     tl.store(C + tl.arange(0, BLOCK), c)
36     tl.store(D + tl.arange(0, BLOCK), d)
37
```

```
38 def run_kernel(A, B):
39     assert A.is_cuda and B.is_cuda
40
41     N = A.numel()
42     BLOCK = 128
43     C = torch.empty_like(A, dtype=torch.int32, device='cuda')
44     D = torch.empty_like(A, dtype=torch.float32, device='cuda')
45
46     grid = lambda meta: (triton.cdiv(N, meta['BLOCK']),)
47     kernel[grid](A, B, C, D, BLOCK=BLOCK)
48     return C, D
49
50 if __name__ == "__main__":
51     A = torch.tensor([40, 50, 3, 4], dtype=torch.uint8, device='cuda')
52     B = torch.tensor([10.0, 20.0, 30.0, 40.0], dtype=torch.float32, device='cuda')
53     C, D = run_kernel(A, B)
54     print("A (input):", A)
55     print("B (input):", B)
56     print("C (int32 part):", C)
57     print("D (float32 part):", D)
```

```
A (input): tensor([40, 50, 3, 4], device='cuda:0', dtype=torch.uint8)
B (input): tensor([10., 20., 30., 40.], device='cuda:0')
C (int32 part): tensor([40, 50, 3, 4], device='cuda:0', dtype=torch.int32)
D (float32 part): tensor([40., 50., 30., 40.], device='cuda:0')
```

Basic Syntax of Triton - static_print

```
0 @triton.jit
1 def add_kernel(
2     x_ptr,
3     y_ptr,
4     output_ptr,
5     n_elements,
6     BLOCK_SIZE: tl.constexpr,
7 ):
8     pid = tl.program_id(axis=0)
9     block_start = pid * BLOCK_SIZE
10    offsets = block_start + tl.arange(0, BLOCK_SIZE)
11    mask = offsets < n_elements
12    x = tl.load(x_ptr + offsets, mask=mask)
13    y = tl.load(y_ptr + offsets, mask=mask)
14    output = x + y
15    tl.store(output_ptr + offsets, output, mask=mask)
16
17    tl.static_print(f"BLOCK_SIZE:\t\t", BLOCK_SIZE)
18    tl.static_print("pid:\t\t\t", pid)
19    tl.static_print("block_start:\t\t", block_start)
20    tl.static_print("tl.arange(0, BLOCK_SIZE):\t", tl.arange(0, BLOCK_SIZE))
21    tl.static_print("mask:\t\t\t", mask)
22    tl.static_print("x:\t\t\t", x)
23    tl.static_print("x_ptr:\t\t\t", x_ptr)
24    tl.static_print("output:\t\t\t", output)
25    tl.static_print("output_ptr+offsets\t\t", output_ptr + offsets)
```

BLOCK_SIZE:	32
pid:	int32[]
block_start:	int32[]
tl.arange(0, BLOCK_SIZE):	int32[constexpr[32]]
mask:	int1[constexpr[32]]
x:	fp32[constexpr[32]]
x_ptr:	pointer<fp32>[]
output:	fp32[constexpr[32]]
output_ptr+offsets	pointer<fp32>[constexpr[32]]

Basic Syntax of Triton - device_print

```
0 @triton.jit
1 def add_kernel(
2     x_ptr,
3     y_ptr,
4     output_ptr,
5     n_elements,
6     BLOCK_SIZE: tl.constexpr,
7 ):
8     pid = tl.program_id(axis=0)
9     block_start = pid * BLOCK_SIZE
10    offsets = block_start + tl.arange(0, BLOCK_SIZE)
11    mask = offsets < n_elements
12    x = tl.load(x_ptr + offsets, mask=mask)
13    y = tl.load(y_ptr + offsets, mask=mask)
14    output = x + y
15    tl.store(output_ptr + offsets, output, mask=mask)
16
17    tl.device_print(f"BLOCK_SIZE:\t\t\t", BLOCK_SIZE)
18    tl.device_print("pid:\t\t\t\t", pid)
19    tl.device_print("block_start:\t\t\t", block_start)
20    tl.device_print("tl.arange(0, BLOCK_SIZE):\t", tl.arange(0, BLOCK_SIZE))
21    tl.device_print("mask:\t\t\t\t", mask)
22    tl.device_print("x:\t\t\t\t\t", x)
23    tl.device_print("x_ptr:\t\t\t\t\t", x_ptr)
24    tl.device_print("output:\t\t\t\t\t", output)
25    tl.device_print("output_ptr+offsets\t\t", output_ptr + offsets)
```

```
pid (0, 0, 0) idx () BLOCK_SIZE:           : 32
pid (0, 0, 0) idx () pid:                   : 0
pid (0, 0, 0) idx () block_start:           : 0

pid (0, 0, 0) idx ( 0) tl.arange(0, BLOCK_SIZE): : 0
pid (0, 0, 0) idx ( 1) tl.arange(0, BLOCK_SIZE): : 1
...
pid (0, 0, 0) idx ( 0) tl.arange(0, BLOCK_SIZE): : 31

pid (0, 0, 0) idx ( 0) mask:                 : 4294967295
pid (0, 0, 0) idx ( 1) mask:                 : 4294967295
pid (0, 0, 0) idx ( 2) mask:                 : 4294967295
pid (0, 0, 0) idx ( 3) mask:                 : 4294967295
pid (0, 0, 0) idx ( 4) mask:                 : 0
pid (0, 0, 0) idx ( 5) mask:                 : 0

pid (0, 0, 0) idx ( 0) x:                    : 1.000000
pid (0, 0, 0) idx ( 1) x:                    : 2.000000
pid (0, 0, 0) idx ( 2) x:                    : 3.000000
pid (0, 0, 0) idx ( 3) x:                    : 4.000000
pid (0, 0, 0) idx ( 4) x:                    : 0.000000
...
pid (0, 0, 0) idx (31) x:                    : 0.000000

pid (0, 0, 0) idx () x_ptr:                  : 0x7f9122600000

pid (0, 0, 0) idx ( 0) output:                : 11.000000
pid (0, 0, 0) idx ( 1) output:                : 22.000000
pid (0, 0, 0) idx ( 2) output:                : 33.000000
pid (0, 0, 0) idx ( 3) output:                : 44.000000
pid (0, 0, 0) idx ( 4) output:                : 0.000000
...
pid (0, 0, 0) idx (31) output:                : 0.000000

pid (0, 0, 0) idx ( 0) output_ptr+offsets    : 0x7f9122600400
pid (0, 0, 0) idx ( 1) output_ptr+offsets    : 0x7f9122600404
pid (0, 0, 0) idx ( 2) output_ptr+offsets    : 0x7f9122600408
pid (0, 0, 0) idx ( 3) output_ptr+offsets    : 0x7f912260040c
...
pid (0, 0, 0) idx (31) output_ptr+offsets    : 0x7f912260047c
```

Basic Syntax of Triton - device_print

On CUDA, printf's are streamed through a buffer of limited size (on one host, we measured the default as 6912 KiB, but this may not be consistent across GPUs and CUDA versions). If you notice some printf's are being dropped, you can increase the buffer size by calling

```
triton.runtime.driver.active.utils.set_printf_fifo_size(size_bytes)
```

CUDA may raise an error if you try to change this value after running a kernel that uses printf's. The value set here may only affect the current device (so if you have multiple GPUs, you'd need to call it multiple times).

References

1. (Documentation) <https://triton-lang.org/main/index.html>
2. (Documentation Tutorials) <https://triton-lang.org/main/getting-started/tutorials/index.html>
3. (Git repo) <https://github.com/triton-lang/triton>
4. (Official publication) <https://www.eecs.harvard.edu/~htk/publication/2019-mapl-tillet-kung-cox.pdf>
5. (Triton puzzles) <https://github.com/triton-lang/triton>
6. (Solved Triton puzzles) <https://github.com/alexzhang13/Triton-Puzzles-Solutions>
7. (Tutorial Part1) <https://readmedium.com/understanding-the-triton-tutorials-part-1-6191b59ba4c>
8. (Tutorial Part2) <https://readmedium.com/understanding-triton-tutorials-part-2-f6839ce50ae7>
9. (Triton autocompiler) <https://dev-discuss.pytorch.org/t/torchinductor-a-pytorch-native-compiler-with-define-by-run-ir-and-symbolic-shapes/747>
10. (Pointer arithmetic in C) <https://metanit.com/c/tutorial/>
11. (Use example) <https://github.com/bitsandbytes-foundation/bitsandbytes/tree/main/bitsandbytes/triton>
12. (Practical tasks) https://docs.google.com/spreadsheets/d/1bD2Nr1UEzj2hsJVqkB0uNLkP6dkepcqYqN_7HK2F5HQ/edit?usp=sharing
13. (OpenAI overview): <https://openai.com/index/triton/>