中国软件技术大会

# ByteIR: 迈向端到端的AI编译

字节跳动 机器学习系统工程师 刘渊强

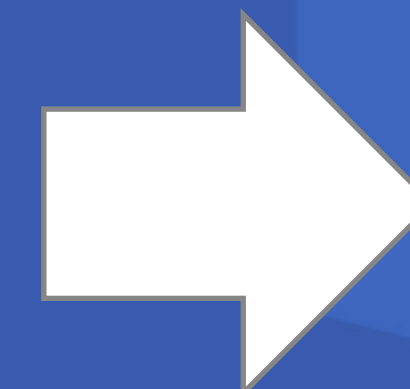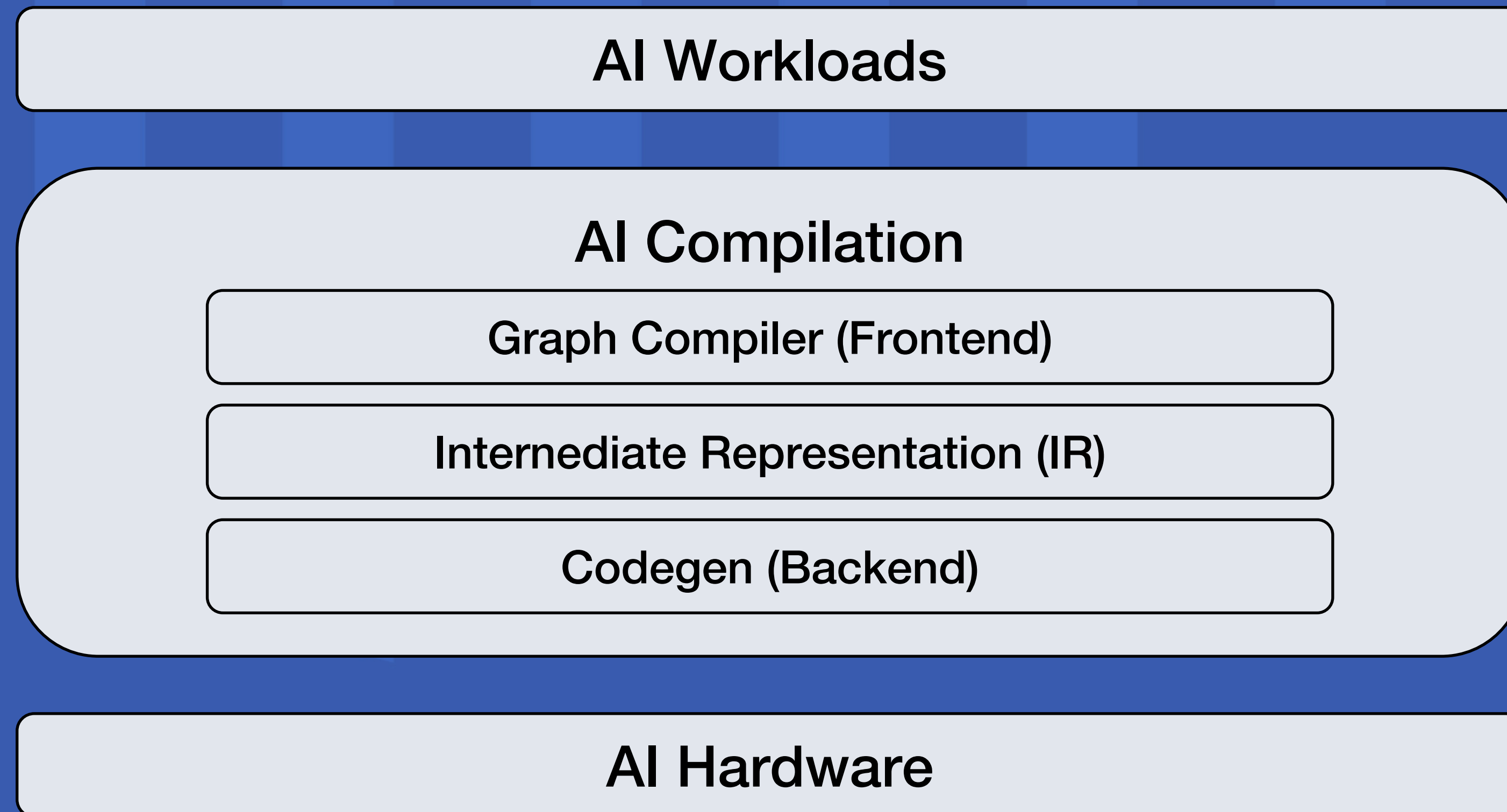2023.12

ByteDance 字节跳动

# What is ByteIR

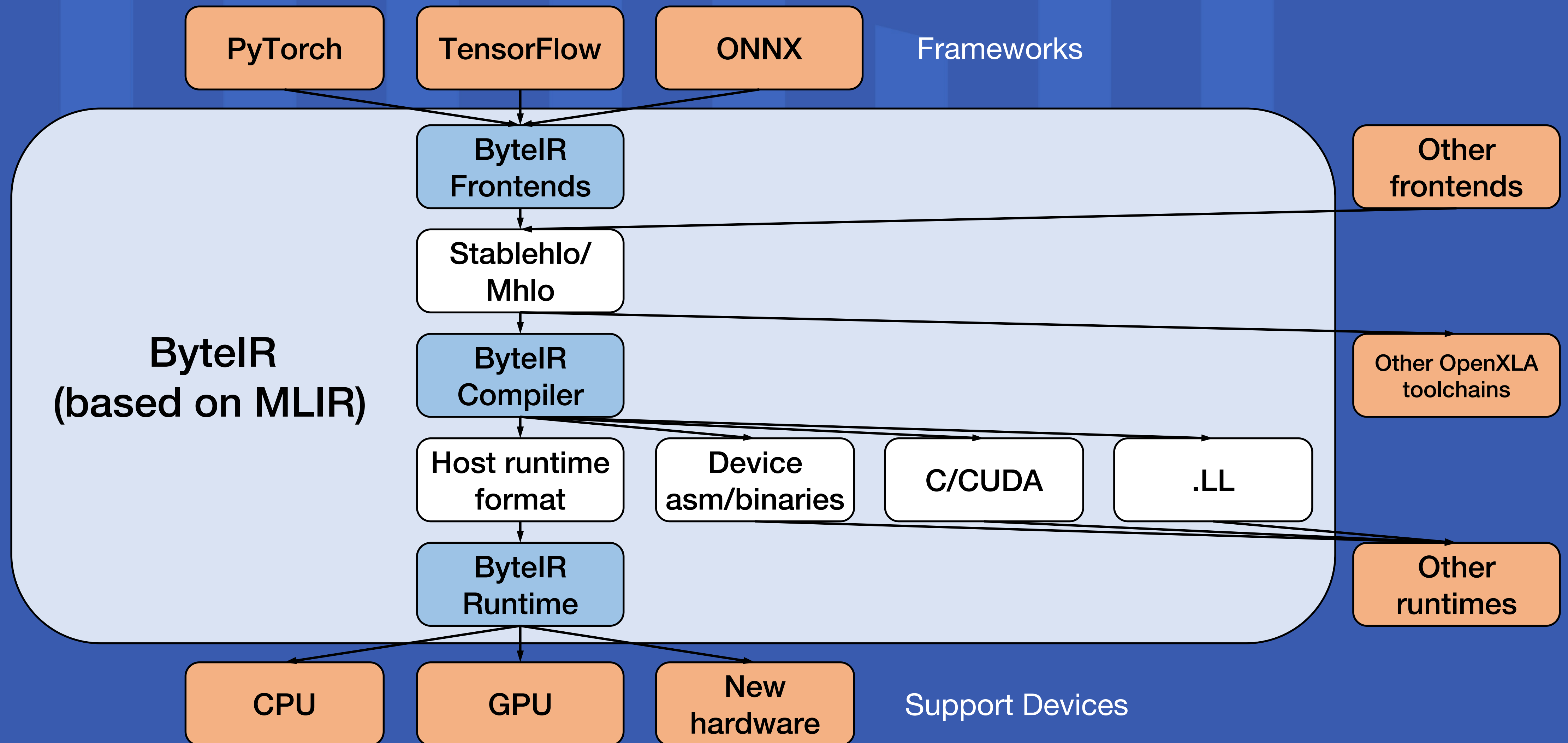**ByteIR** is NOT a specific IR (Intermediate Representation)

**ByteIR** is our solution for framework-to-hardware compilation

ByteDance 字节跳动

# AI Compilation: NN graph to HW

AI Workloads

## AI Compilation

Graph Compiler (Frontend)

Internediate Representation (IR)

Codegen (Backend)

AI Hardware

ByteDance 字节跳动

# ByteIR Architecture

PyTorch   TensorFlow   ONNX   Frameworks

Other frontends

ByteIR
(based on MLIR)

ByteIR
Frontends

Stablehlo/
Mhlo

ByteIR
Compiler

Other OpenXLA
toolchains

Host runtime
format

Device
asm/binaries

C/CUDA

.LL

ByteIR
Runtime

Other
runtimes

CPU   GPU   New
hardware   Support Devices

ByteDance 字节跳动

# Advantages of ByteIR

1. Embrace open source, upstream first

   Contribute to llvm, tensorflow, pytorch, torch-mlir, onnx-mlir, stablehlo

2. Well support for PyTorch, both inference and training

3. Friendly to new hardware (ASIC/NPU)

4. Flexible, extensible, high performance

ByteDance 字节跳动

ByteDance 字节跳动

# ByteIR Torch Frontend – Torch MLIR

**PyTorch**

LazyTensorCore ← TorchScript ← FX Graph

**Torch MLIR**

LTC MLIR Plugin

TorchScript MLIR Converter

FX Importer

Torch Dialect

Linalg, Arith, Tensor, SCF

TOSA

Stablehlo

Custom

ByteDance字节跳动

# Torch MLIR Lowering

## torch

```
class Module(torch.nn.Module):
    def __init__(self):
        super().__init__()
        self.linear = nn.Linear(10, 20)
    def forward(self, x):
        return self.linear(x)
```

torch.jit.scirpt

MLIR importer

## torch dialect

```
func.func @forward(%arg0: !torch.tensor,
%arg1: !torch.tensor) -> !torch.tensor {
    %0 = torch.aten.mm %arg0,
%arg1 : !torch.tensor, !torch.tensor -> !torch.tensor
    return %0 : !torch.tensor
}
```

Shape/dtype inference

Value semantics analysis

## torch dialect

```
func.func @forward(%arg0: !torch.vtensor<[20,
10], f32>, %arg1: !torch.vtensor<[10, 20], f32>) -
> !torch.vtensor<[20, 20], f32> {
    %0 = torch.aten.mm %arg0,
%arg1 : !torch.vtensor<[20, 10],
f32>, !torch.vtensor<[10, 20], f32> -
> !torch.vtensor<[20, 20], f32>
    return %0 : !torch.vtensor<[20, 20], f32>
}
```

Decomposition

## stablehlo dialect

```
func.func @forward(%arg0: tensor<20x10xf32>,
%arg1: tensor<10x20xf32>) ->
tensor<20x20xf32> {
    %0 = "stablehlo.dot"(%arg0, %arg1) :
tensor<20x10xf32>, tensor<10x20xf32> ->
tensor<20x20xf32>
    return %0 : tensor<20x20xf32>
}
```

Lowering

ByteDance 字节跳动

# ByteIR Torch Frontend



Able to provide coarse-grained operators

# Corner Cases

```
def forward(self, x):
    x += 1
    return x
```

Success, but no effect on input

```
def forward(self, x, y):
    x += y
    z = x.view(-1, 4)
    x += 1
    return x, z
```

Failed, really overwrite same memory

```
def forward(self, x):
    y = torch.as_strided(x, size, stride)
    return y
```

Failed, no abstraction of pointer/storage

- Don't support dynamic if/for/while

……

ByteDance字节跳动

# ByteIR Compiler Overview

```
          ┌──────────────┐      ┌──────────────────┐
          │  Opt. passes │─────▶│   Stablehlo/Mhlo │
          └──────────────┘      └──────────────────┘
                 ╱                    │            ╲
                ╱                     ▼             ╲
┌──────────────┐          ┌──────────────┐      ┌──────────────┐
│   Lowering   │          │   Lowering   │      │   Lowering   │
└──────────────┘          └──────────────┘      └──────────────┘
       │                         │                     │
       ▼                         ▼                     ▼
┌──────────────┐      ┌──────────────┐  ┌─────────┐  ┌──────────────┐
│     CAT*     │      │  Opt. passes │─▶│ Linalg  │  │     ByRE     │
└──────────────┘      └──────────────┘  │(tensor) │  └──────────────┘
       │                                └─────────┘         │
       ▼                                     │              ▼
┌──────────────┐                     ┌──────────────┐ ┌──────────────┐
│  AITemplate  │                     │ Bufferization│ │Serialization │
│ backend for  │                     └──────────────┘ └──────────────┘
│  Nvidia GPU  │                            │
└──────────────┘                            ▼
                      ┌──────────────┐ ┌──────────────┐
                      │  Opt. passes │─▶│Linalg(memref)│
                      └──────────────┘ └──────────────┘
                                              │
                                              ▼
                                      ┌──────────────┐
                                      │   Lowering   │
                                      └──────────────┘
                                              │
                                              ▼
                      ┌──────────────┐ ┌──────────────┐
                      │  Opt. passes │─▶│scf/affine/vec/│
                      └──────────────┘ └──────────────┘
                             ╱              │          ╲
                            ╱               ▼           ╲
               ┌──────────────┐  ┌──────────────┐  ┌──────────────┐
               │ GPU Lowering │  │   Lowering   │  │   C Emitter  │
               └──────────────┘  └──────────────┘  └──────────────┘
                      │                  │                 │
                      ▼                  ▼                 ▼
               ┌──────────────┐  ┌──────────────┐  ┌──────────────┐
               │     GPU      │  │     LLVM     │  │      C       │
               └──────────────┘  └──────────────┘  └──────────────┘
                    ╱    ╲
                   ╱      ╲
      ┌──────────────┐  ┌──────────────┐
      │   Lowering   │  │ CUDA Emitter │
      └──────────────┘  └──────────────┘
             │                 │
             ▼                 ▼
      ┌──────────────┐  ┌──────────────┐
      │  NVVM/LLVM   │  │     CUDA     │
      └──────────────┘  └──────────────┘
```
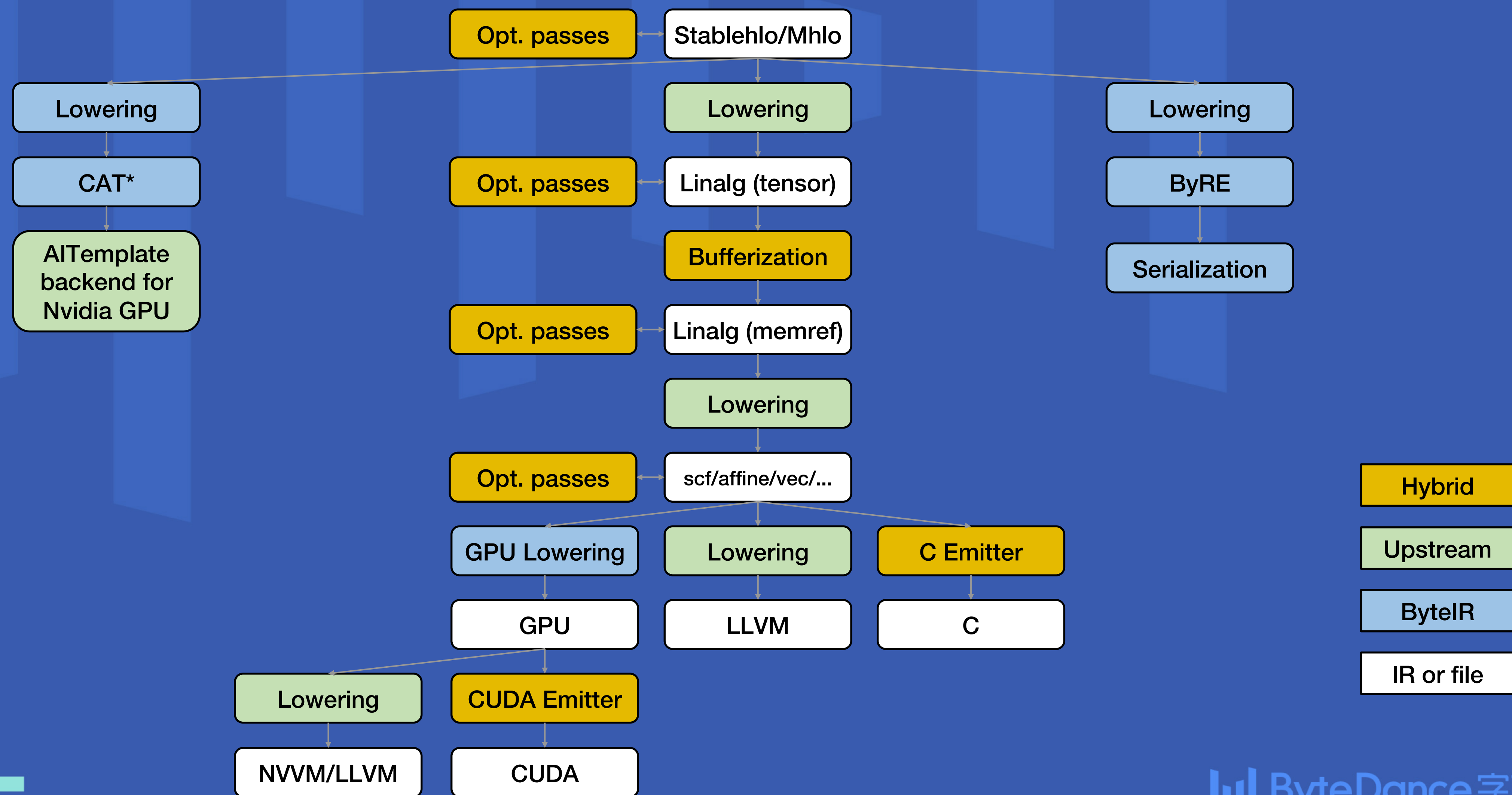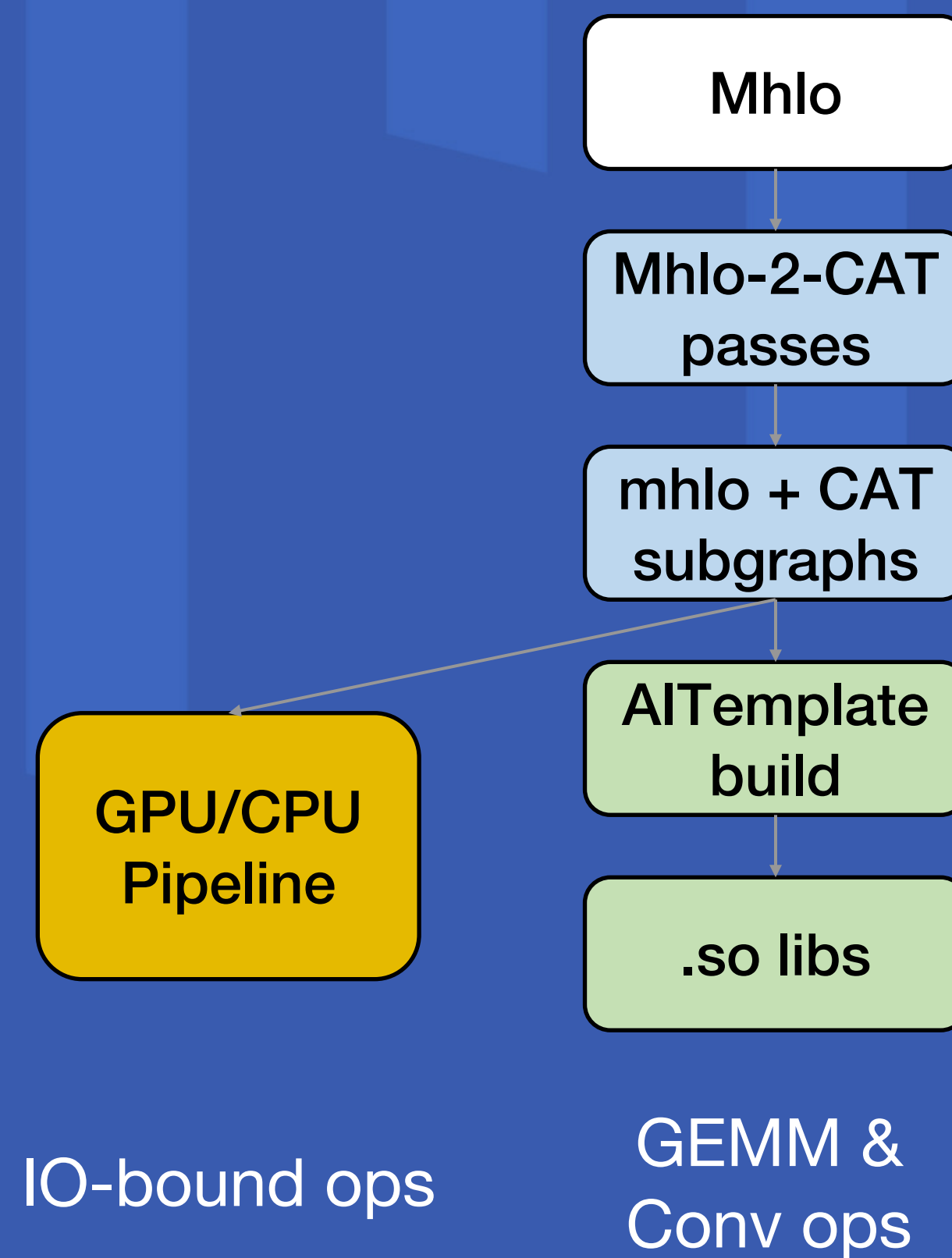
Legend:
- Hybrid
- Upstream
- ByteIR
- IR or file

# Integrating AITemplate

We introduce **CAT** (Composable Algebra Template) dialect

```
┌─────────────┐
│    Mhlo     │
└─────────────┘
       │
┌─────────────┐
│  Mhlo-2-CAT │
│    passes   │
└─────────────┘
       │
┌─────────────┐
│  mhlo + CAT │
│  subgraphs  │
└─────────────┘
    ╱      │
┌─────────┐ ┌─────────────┐
│ GPU/CPU │ │  AITemplate │     AITemplate: Meta's inference framework converting
│ Pipeline│ │    build    │     PyTorch NN to CUDA code (based on CUTLASS)
└─────────┘ └─────────────┘
            ┌─────────────┐
            │   .so libs  │
            └─────────────┘

 IO-bound ops    GEMM &
                 Conv ops
```

ByteDance 字节跳动

# Mhlo-2-CAT Passes

Mhlo-2-CAT:
- Convert Mhlo ops to CAT ops (one CAT op corresponds to one AIT op)
- Eliminate redundant transpose/permuate ops



Mhlo dialect

CAT dialect

# Linalg Tiling and Fusion

Mhlo op:
```
func.func @fuse_element(%arg0 ..., %arg1 ...)
    %0 = mhlo.some_elemwise_binary_1(%arg0, %arg1)
    %1 = mhlo.some_elemwise_binary_2(%0, %arg1)
    return %1
```

⬇ Linalg transformation

Linalg op:
```
func.func @fuse_element(%arg0 ..., %arg1 ...)
    %1 = linalg.generic {indexing_maps = ...,
            iterator_types = ...} ins(%arg0, %arg1) outs... {
^bb0(%in0, %in1, %out)
        %2 = arith.some_elemwise_binary_1(%in0, %in1)
        %3 = arith.some_elemwise_binary_2(%2, %in1)
        linalg.yield %3 ...
    }
    return %1
```

# ByetIR's Linalg Extension

More ops:
- Alias, Diag, Scan, Scatter, Softmax, TopK
- support transformations of extended ops

Enhanced fusion transformations:
- producer-consumer & input-sharing fusion
- tiling along reduction axis correction
- intermediates as outputs within fusion
- intermediate tensor dim simplification
- map ops to generic ops conversion
- ...

Other introduced transformations:
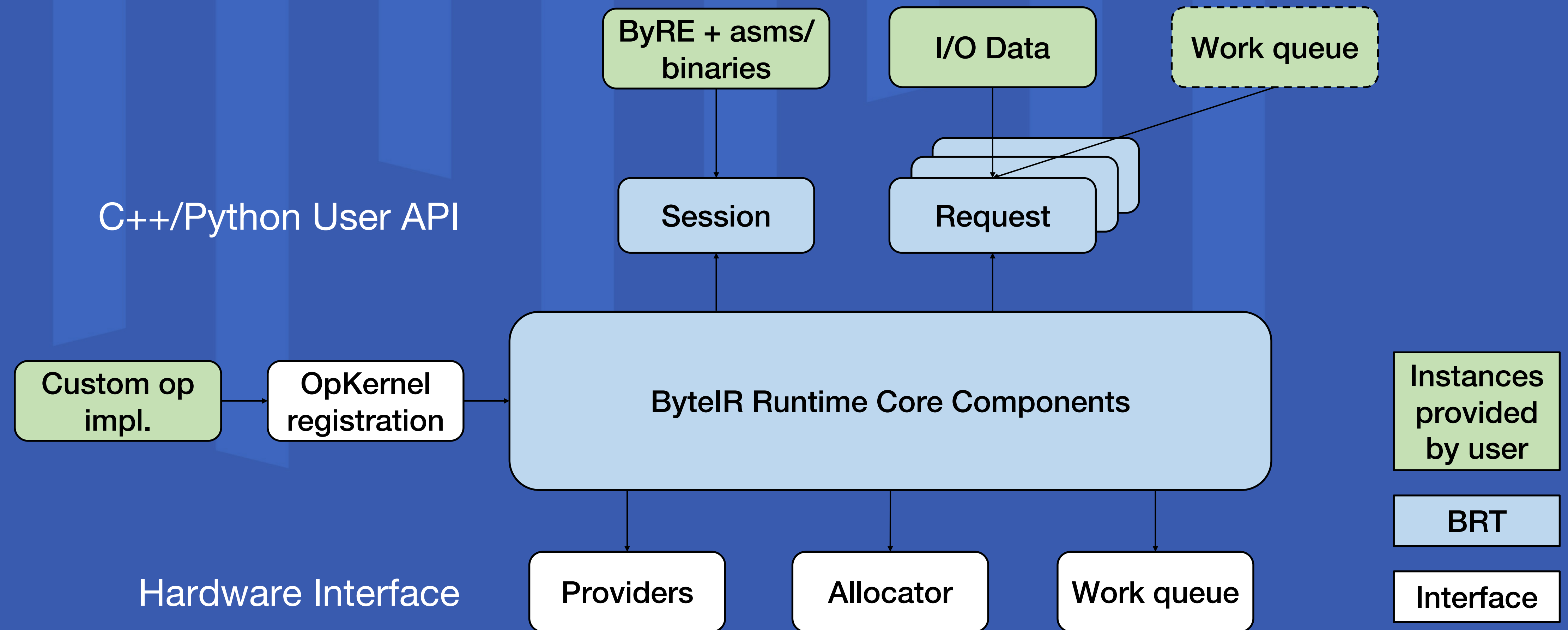- Collapse dims transformation
- Fuse operands transformation
- ...

Benefits:
- Extreme IO-bound op fusion
- Lower overhead for fused ops
  (Exploit GPU DRAM bandwidth)

ByteDance 字节跳动

# ByteIR Runtime (BRT) Overview

# BRT Interface for Hardware

Provider | A collection of op implementation
- e.g., mm, d2h/h2d memcpy

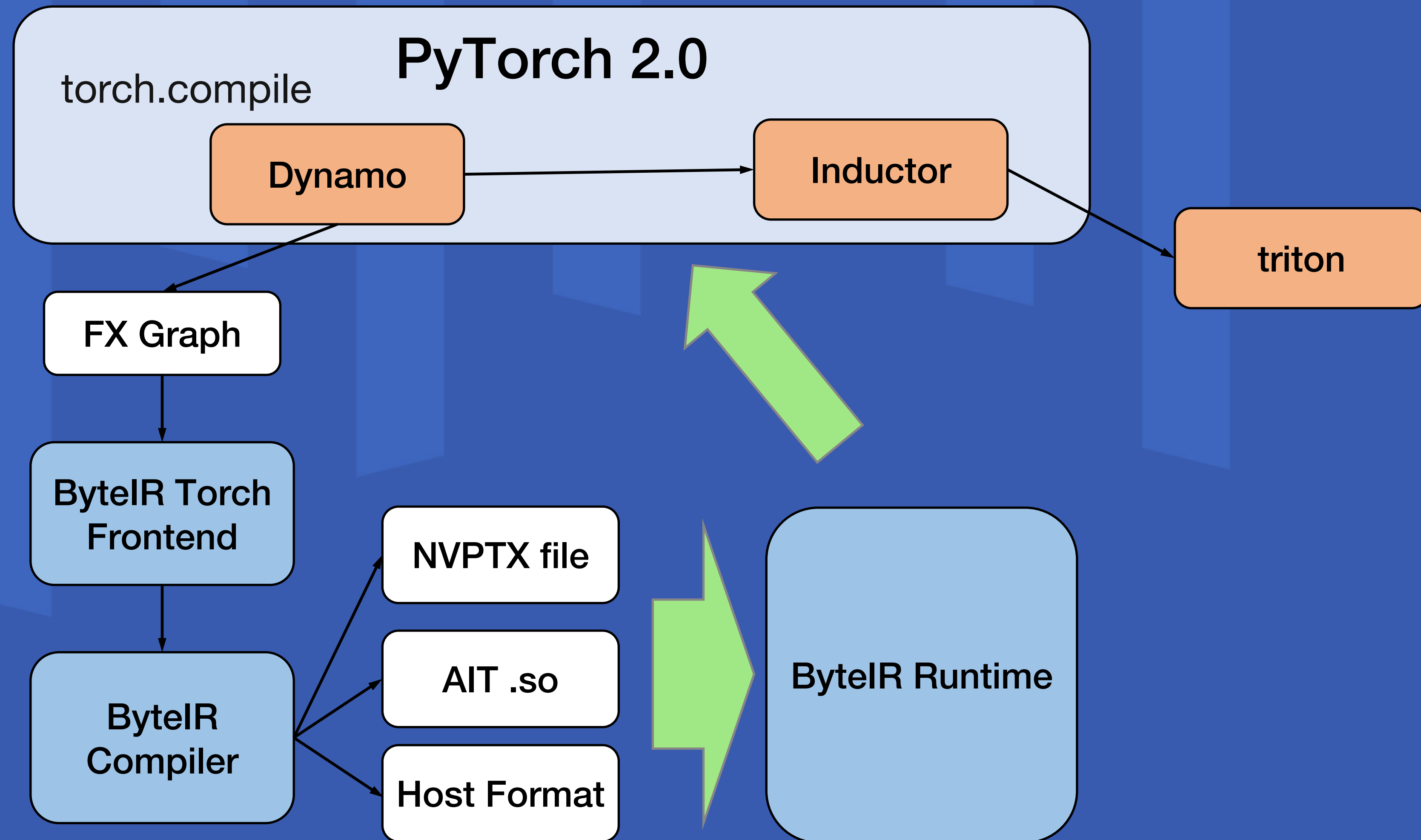Work Queue | Abstraction(like CUDAStream) for execution order

Allocator | Memory Allocate/Free

ByteDance 字节跳动

# ByteIR LLM Training Compilation Pipeline

# Optimization in FW/BW Partition

**Forward**

**Backward**

**Forward**

**Backward**

X

**Recompute**

Transpose
perm: (0, 1, 3, 2)

Reshape
shape: (0*1, 2, 3)

BMM

Transpose
perm: (0, 2, 1)

BMM

X

Transpose
perm: (0, 1, 3, 2)

Reshape
shape: (0*1, 2, 3)

BMM

**Fusion**

Transpose
perm: (0, 1, 3, 2)

Reshape
shape: (0*1, 2, 3)

Transpose
perm: (0, 2, 1)

BMM

One case of our partition strategy

ByteDance字节跳动

# ByteIR PyTorch Compile Example

```python
from byteir import byteir_compile_fx

model = make_model(model_name)

# 1, compile with byteir
optimized_model = torch.compile(model, backend=byteir_compile_fx)

# 2, execution as usual
data = make_data(optimized_model, model_name, device)
model.zero_grad(set_to_none=True)
with torch.cuda.amp.autocast(enabled=True, dtype=torch.float16):
    # forward compile
    loss = compute_loss(optimized_model, data)
    # backward compile
    loss.backward()
```

Demo Link: https://github.com/bytedance/byteir/tree/main/frontends/torch-frontend/examples/demo

ByteDance 字节跳动

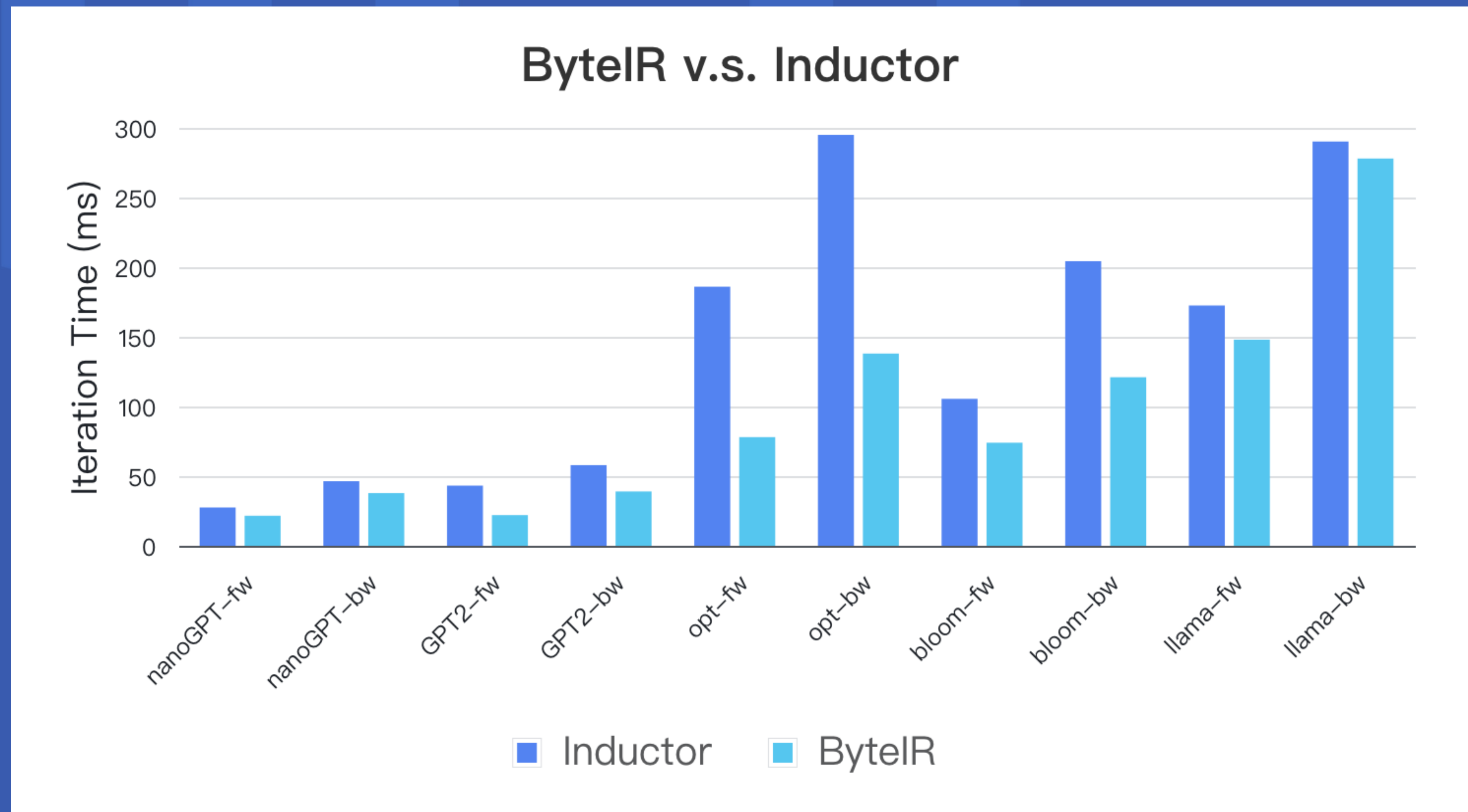# ByteIR Runtime Example

```python
import brt

session = brt.Session()
session.load(byre_model_path)
req = session.new_request_context(torch.cuda.current_stream())
inputs, outputs = [], []

# init input/output data
for offset in session.get_input_arg_offsets():
    inputs.append(torch.randn(session.get_static_shape(offset),
            dtype=dtype, device="cuda"))
    req.bind_arg(offset, inputs[-1].data_ptr())
...

req.finish_io_binding()
req.run()
req.sync()
```

Demo Link: https://github.com/bytedance/byteir/tree/main/frontends/torch-frontend/examples/demo

# Performance

- Flash Attention 2

- Elementwise Fusion

- AITemplate

- Reduce Codegen



ByteIR v.s. Inductor

ByteDance 字节跳动

# Conclusion & Future Work

We introduce ByteIR: a framework-to-hardware compiler solution

- Friendly to PyTorch and GPU/ASIC

- PyTorch 2.0 training/inference demo on LLM

Future Work:

- Distributed support

- TensorCore MMA Codegen

Website: https://byteir.ai
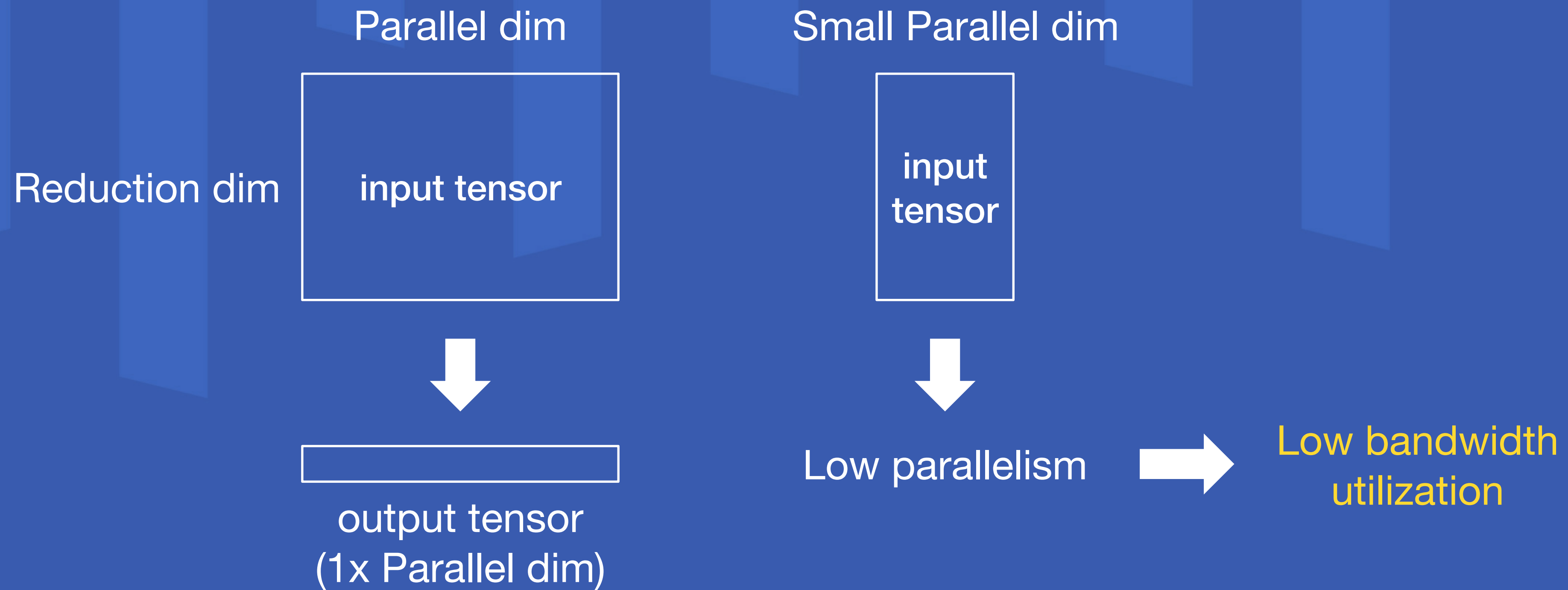
Github: https://github.com/bytedance/byteir

ByteDance 字节跳动

# THANKS

## https://byteir.ai

ByteIR Github: https://github.com/bytedance/byteir

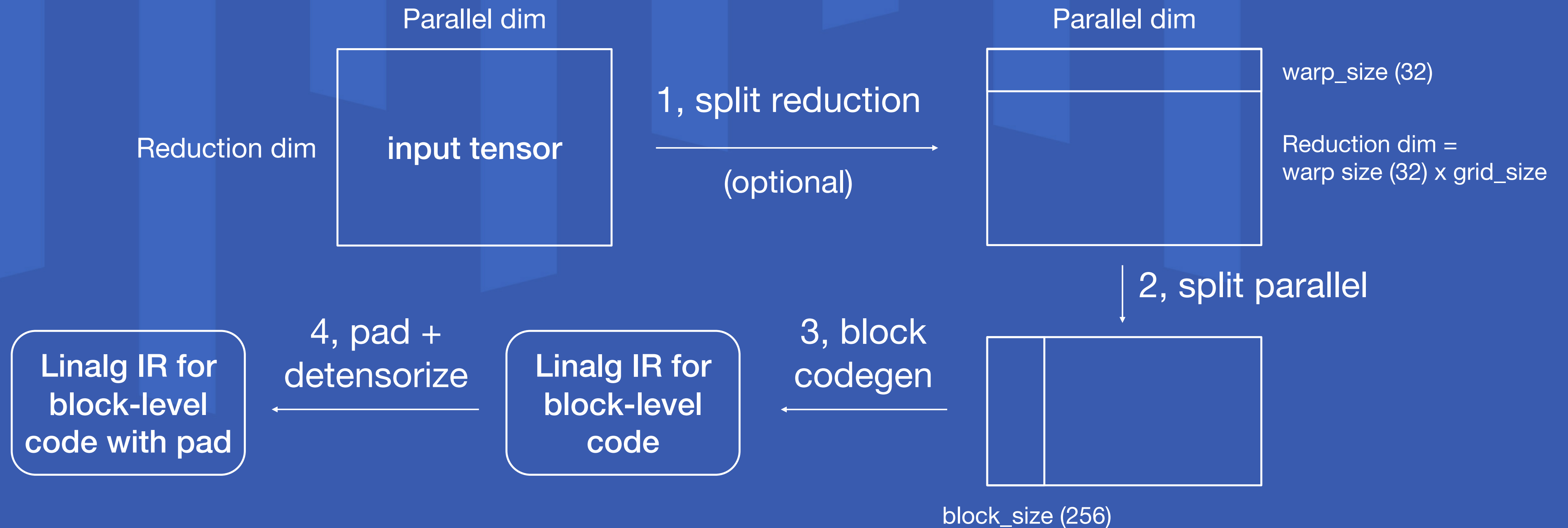Personal Github: https://github.com/qingyunqu

# Appendix 1: Reduce Op Optimization (Fusion)

**Optimization 1**: Fusing reduce op with producer ops

# Appendix 2: Reduce Op Optimization (Tiling)

**Optimization 2**: Parallelizing reduce dimension

Parallel dim

Reduction dim | input tensor

1, split reduction

(optional)

Parallel dim

warp_size (32)

Reduction dim =
warp size (32) x grid_size

2, split parallel

Linalg IR for block-level code with pad ← 4, pad + detensorize ← Linalg IR for block-level code ← 3, block codegen ←

block_size (256)

*We use utilize our LinalgExt transformations to achieve best tiling efficiency*

ByteDance 字节跳动