

# MLIR and PyTorch: A Compilation Pipeline Targeting Huawei's Ascend Backend

Weizhe Chen, Alexandre Singer, Fang Gao, Kai-Ting Amy Wang  
Huawei Canada Research Centre

## Summary

We present an approach to compiling **PyTorch code**, through **MLIR**, to executable **Ascend C code**. Ascend C is a programming language oriented to operator development scenarios that support **Ascend AI Processors**.

This pipeline efficiently **transforms high-level PyTorch code** to **optimized, hardware-specific Ascend C code** by bridging PyTorch, MLIR, and Ascend hardware. It allows users of Ascend Processors to benefit from MLIR optimizations and Ascend compiler engineers to contribute innovations and improvements to the MLIR community.

## Ascend C Programming Language

Ascend C is the **new programming language** launched by Huawei in 2023, targeting the Ascend architecture<sup>1</sup>.

It aims to **improve the efficiency of operator development** and to help AI engineers to develop operators and tune models at a low cost.

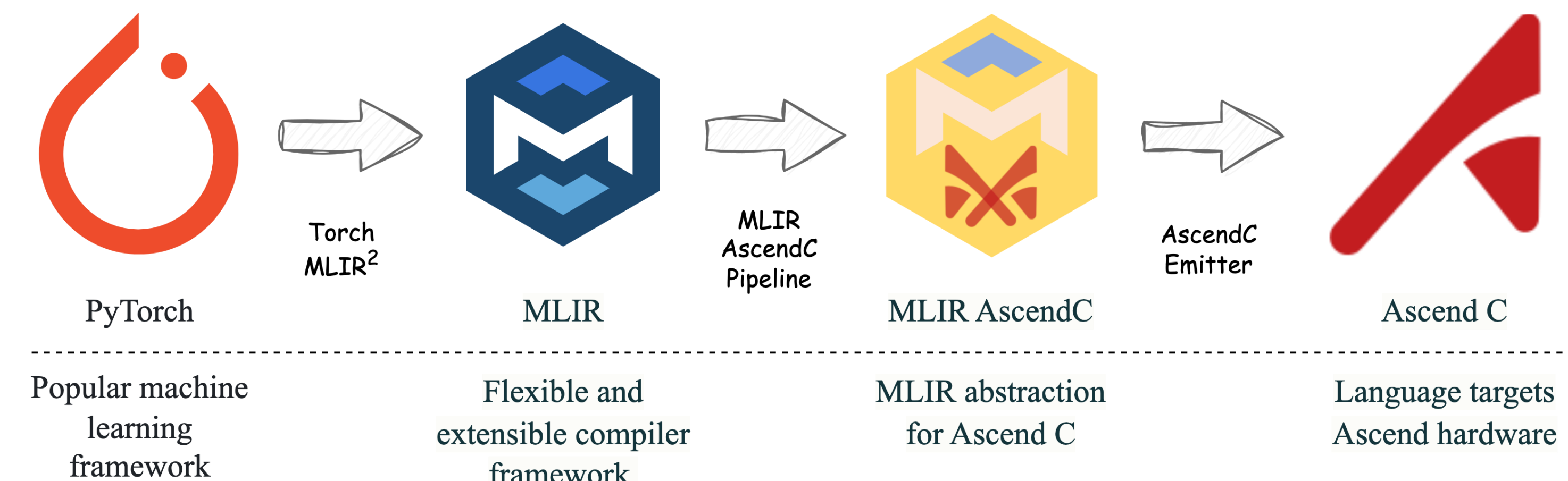
- Operator development scenarios
- Natively supports **C/C++ specifications**
- Automatic parallel scheduling
- **Multi-layer library interface**
- CPU-NPU dual debugging

```

__aicore__ inline void Compute(int32_t progress)
{
    // deque input tensors from VECIN queue
    LocalTensor<half> xLocal = inQueueX.DeQueue<half>();
    LocalTensor<half> yLocal = inQueueY.DeQueue<half>();
    LocalTensor<half> zLocal = outQueueZ.AllocTensor<half>();
    // call Add instr for computation
    Add(zLocal, xLocal, yLocal, TILE_LENGTH);
    // enqueue the output tensor to VECOUT queue
    outQueueZ.Enqueue(zLocal);
    // free input tensors for reuse
    inQueueX.FreeTensor(xLocal);
    inQueueY.FreeTensor(yLocal);
}

```

## Goal: PyTorch → MLIR → MLIR AscendC → Ascend



```

c = torch.add(a, b)

func.func @add_custom(%a : !ascendc.GM_ADDR, %b : !ascendc.GM_ADDR, %c : !ascendc.GM_ADDR) {
  // Compute Offset
  %total_len = arith.constant 98384 : i64
  %use_core_num = arith.constant 8 : i64
  %block_len = arith.divui %total_len, %use_core_num : i64
  %tile_num = arith.constant 16 : i64
  %buffer_num = arith.constant 2 : i64
  %tile_length_pre_buffer = arith.divui %block_len, %tile_num : i64
  %tile_length = arith.divui %tile_length_pre_buffer, %buffer_num : i64
  %block_idx = arith.constant 164 : i64
  %offset = arith.muli %block_idx, %block_len : i64

  // Struct Setup
  %gm = ascendc.create_global_tensor %a, %offset, %block_len : !ascendc.GM_ADDR -> !ascendc.Glob
  %gm = ascendc.create_global_tensor %b, %offset, %block_len : !ascendc.GM_ADDR -> !ascendc.Glob
  %gm = ascendc.create_global_tensor %c, %offset, %block_len : !ascendc.GM_ADDR -> !ascendc.Glob

  %pipe = ascendc.create_pipe : !ascendc.TPipe
  %inQueueX = ascendc.create_queue %pipe, %tile_length : !ascendc.TPipe -> !ascendc.TQueue<VECIN, 2
  %inQueueY = ascendc.create_queue %pipe, %tile_length : !ascendc.TPipe -> !ascendc.TQueue<VECIN, 2
  %outQueueZ = ascendc.create_queue %pipe, %tile_length : !ascendc.TPipe -> !ascendc.TQueue<VECOUT, 2

  // Process
  %c0 = arith.constant 0 : i64
  %c1 = arith.constant 1 : i64
  %loop_count = arith.muli %tile_num, %buffer_num : i64
  %cf = for %progress = %c0 to %loop_count step %c1 : i64 {
    %gm_offset = arith.muli %progress, %tile_length : i64

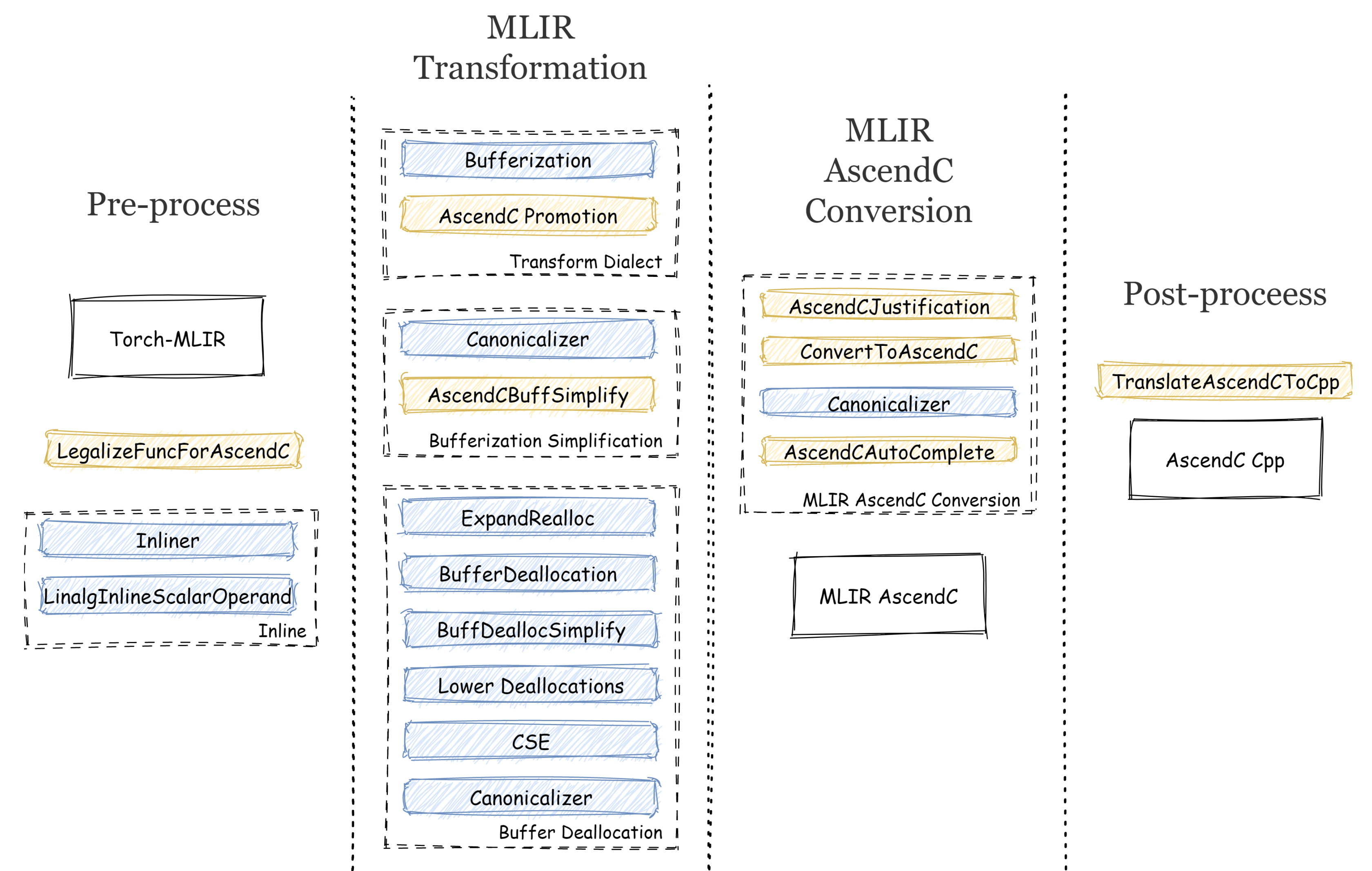
    // Copy in
    %xLocal = ascendc.alloc_tensor(%inQueueX) : !ascendc.TQueue<VECIN, 2> -> !ascendc.LocalTensor
    %yLocal = ascendc.alloc_tensor(%inQueueY) : !ascendc.TQueue<VECIN, 2> -> !ascendc.LocalTensor
    ascendc.data_copy %xLocal, %gm_offset, %tile_length : !ascendc.GlobalTensor<16384xf16> ->
    ascendc.data_copy %yLocal, %gm_offset, %tile_length : !ascendc.GlobalTensor<16384xf16> ->
    ascendc.enqueue %xLocal, %inQueueX : !ascendc.LocalTensor<16384xf16> -> !ascendc.TQueue<VECIN, 2>
    ascendc.enqueue %yLocal, %inQueueY : !ascendc.LocalTensor<16384xf16> -> !ascendc.TQueue<VECIN, 2>
    // Compute
    %xLocal_deque = ascendc.deque(%inQueueX) : !ascendc.TQueue<VECIN, 2> -> !ascendc.LocalTensor
    %yLocal_deque = ascendc.deque(%inQueueY) : !ascendc.TQueue<VECIN, 2> -> !ascendc.LocalTensor
    %xLocal = ascendc.alloc_tensor(%outQueueZ) : !ascendc.TQueue<VECOUT, 2> -> !ascendc.LocalTensor
    ascendc.add %xLocal, %yLocal_deque, %xLocal, %tile_length : (!ascendc.LocalTensor
    ascendc.enqueue %xLocal, %outQueueZ : !ascendc.LocalTensor<16384xf16> -> !ascendc.TQueue<VECOUT, 2>
    ascendc.free_tensor %xLocal_deque : !ascendc.LocalTensor<16384xf16>
    ascendc.free_tensor %yLocal_deque : !ascendc.LocalTensor<16384xf16>
    // Copy out
    %xLocal_deque = ascendc.deque(%outQueueZ) : !ascendc.TQueue<VECOUT, 2> -> !ascendc.LocalTensor
    ascendc.data_copy %xLocal_deque, %gm_offset, %tile_length : !ascendc.LocalTensor
    ascendc.free_tensor %xLocal_deque : !ascendc.LocalTensor<16384xf16>
    %cf_yield
  }
  return
}

```

## MLIR-AscendC Dialect

- Simplified and **one-to-one abstraction of the Ascend C language**.
- Fully express Ascend C operators using MLIR.
  - **Type:** AscendC\_TPipe, AscendC\_TQueue, AscendC\_GMADDR, AscendC\_TensorType
  - **Op:** constructors, memory management, data process, vector computation, system variable access
  - **Attribute:** TPosition Enum Attribute

## MLIR-AscendC Pipeline



### 1. Pre-Process

The first stage will prepare the MLIR code for further transformation and conversion.

It will create a memref operand as a replacement for the original tensor operand and combine two through inlining to **prepare for bufferization**. It will then modify function arguments to customized types to **prepare for Ascend C conversion**.

```

func.func @add(%arg0: !ascendc.GM_ADDR, %arg1: !ascendc.GM_ADDR, %arg2: !ascendc.GM_ADDR) {
  %xloc = memref.alloc() : memref<16384xf16, #ascendc.TPosition<GM>>
  %yloc = memref.alloc() : memref<16384xf16, #ascendc.TPosition<GM>>
  %zloc = memref.alloc() : memref<16384xf16, #ascendc.TPosition<GM>>
  %x = bufferization.to_tensor %xloc restrict writable : memref<16384xf16, #ascendc.TPosition<GM>>
  %y = bufferization.to_tensor %yloc restrict writable : memref<16384xf16, #ascendc.TPosition<GM>>
  %z = tensor.empty() : tensor<16384xf16>
  %x3 = linalg.generic (indexing_maps = [#map, #map, #map], iterator_types = ["parallel"]) ins(%x, %y) : tensor<16384xf16>, tensor<16384xf16> -> outs(%z) : tensor<16384xf16> {
    %bb0(%in: f16, %in2: f16, %out: f16):
      %x4 = arith.addf %in, %in2 : f16
      linalg.yield %x4 : f16
  } -> tensor<16384xf16>
  %xloc1 = memref.alloc() : memref<16384xf16, #ascendc.TPosition<GM>>
  %yloc1 = memref.alloc() : memref<16384xf16, #ascendc.TPosition<GM>>
  %zloc1 = memref.alloc() : memref<16384xf16, #ascendc.TPosition<GM>>
  %x1 = bufferization.materialize_in_destination %x3 in writable %xloc1 : (tensor<16384xf16>, memref<16384xf16, #ascendc.TPosition<GM>>) -> ()
  return
}

```

### 3. MLIR AscendC Conversion

The third stage converts standard MLIR to MLIR AscendC, which adds target-specific information.

Specifically, it **justifies the usage of Ascend C** by inserting TPipe and adjusting TPosition. Then, it will **recursively convert into MLIR AscendC** and add the synchronization methods used by Ascend C.

```

func.func @add(%arg0: !ascendc.GM_ADDR, %arg1: !ascendc.GM_ADDR, %arg2: !ascendc.GM_ADDR) {
  %c16384_i32 = arith.constant 16384 : i32
  %x0 = ascendc.create_global_tensor %arg0, %c16384_i32 : !ascendc.GM_ADDR -> !ascendc.GlobalTensor<16384xf16>
  %x3 = ascendc.create_pipe : !ascendc.TPipe
  %x8 = ascendc.create_queue %x3, %c16384_i32 : !ascendc.TPipe -> !ascendc.TQueue<VECIN, 1>
  %y0 = ascendc.create_global_tensor %arg1, %c16384_i32 : !ascendc.GM_ADDR -> !ascendc.GlobalTensor<16384xf16>
  %y3 = ascendc.create_pipe : !ascendc.TPipe
  %y8 = ascendc.create_queue %y3, %c16384_i32 : !ascendc.TPipe -> !ascendc.TQueue<VECIN, 1>
  %z0 = ascendc.create_global_tensor %arg2, %c16384_i32 : !ascendc.GM_ADDR -> !ascendc.GlobalTensor<16384xf16>
  ascendc.enqueue %x0, %x8 : !ascendc.LocalTensor<16384xf16> -> !ascendc.TQueue<VECIN, 1>
  ascendc.enqueue %y0, %y8 : !ascendc.LocalTensor<16384xf16> -> !ascendc.TQueue<VECIN, 1>
  ascendc.data_copy %x8, %y8, %c16384_i32 : !ascendc.LocalTensor<16384xf16> -> !ascendc.GlobalTensor<16384xf16>
  ascendc.enqueue %x8, %y8 : !ascendc.LocalTensor<16384xf16> -> !ascendc.TQueue<VECIN, 1>
  ascendc.free_tensor %x0, %y0 : !ascendc.LocalTensor<16384xf16>
  ascendc.free_tensor %x8, %y8 : !ascendc.LocalTensor<16384xf16>
  ascendc.free_tensor %z0 : !ascendc.LocalTensor<16384xf16>
  return
}

```

### 2. MLIR Transformation

The second stage performs MLIR transformation including **bufferization, promotion and deallocation** as Ascend C has memref-like operands with customized address space.

We perform one-shot bufferize, customized promotion, and utilize the upstream pipeline for buffer deallocation to prepare for MLIR AscendC conversion.

```

func.func @add(%arg0: !ascendc.GM_ADDR, %arg1: !ascendc.GM_ADDR, %arg2: !ascendc.GM_ADDR) {
  %xloc = memref.alloc() : memref<16384xf16, #ascendc.TPosition<VECIN>>
  %yloc = memref.alloc() : memref<16384xf16, #ascendc.TPosition<VECIN>>
  %zloc = memref.alloc() : memref<16384xf16, #ascendc.TPosition<VECIN>>
  %x = memref.copy %xloc, %arg0 : memref<16384xf16, #ascendc.TPosition<GM>> -> memref<16384xf16, #ascendc.TPosition<VECIN>>
  %y = memref.copy %yloc, %arg1 : memref<16384xf16, #ascendc.TPosition<GM>> -> memref<16384xf16, #ascendc.TPosition<VECIN>>
  %z = memref.copy %zloc, %arg2 : memref<16384xf16, #ascendc.TPosition<GM>> -> memref<16384xf16, #ascendc.TPosition<VECIN>>
  %xloc1 = memref.alloc() : memref<16384xf16, #ascendc.TPosition<VECIN>>
  %yloc1 = memref.alloc() : memref<16384xf16, #ascendc.TPosition<VECIN>>
  %zloc1 = memref.alloc() : memref<16384xf16, #ascendc.TPosition<VECIN>>
  %x1 = memref.dealloc %xloc, %xloc1 : memref<16384xf16, #ascendc.TPosition<VECIN>>
  %y1 = memref.dealloc %yloc, %yloc1 : memref<16384xf16, #ascendc.TPosition<VECIN>>
  %z1 = memref.dealloc %zloc, %zloc1 : memref<16384xf16, #ascendc.TPosition<VECIN>>
  %x1 = memref.dealloc %xloc1, %x1 : memref<16384xf16, #ascendc.TPosition<VECIN>>
  %y1 = memref.dealloc %yloc1, %y1 : memref<16384xf16, #ascendc.TPosition<VECIN>>
  %z1 = memref.dealloc %zloc1, %z1 : memref<16384xf16, #ascendc.TPosition<VECIN>>
  %x1 = memref.dealloc %x1, %x1 : memref<16384xf16, #ascendc.TPosition<VECIN>>
  %y1 = memref.dealloc %y1, %y1 : memref<16384xf16, #ascendc.TPosition<VECIN>>
  %z1 = memref.dealloc %z1, %z1 : memref<16384xf16, #ascendc.TPosition<VECIN>>
  return
}

```

### 4. Post-Process

The fourth stage will emit the Ascend C language, which is the backend for the pipeline.

It will extend the current EmitC<sup>3</sup> abilities to include Ascend C specifications and produce executable code that targets the Ascend architecture.

```

#include "kernel_operator.h"
using namespace AscendC;

extern "C" __global__ __aicore__ void add(GM_ADDR v1, GM_ADDR v2, GM_ADDR v3) {
  int32_t i4 = 16384;
  GlobalTensor<half> gmTensor5;
  gmTensor5.SetGlobalBuffer((__gm_half*)v1, i4);
  GlobalTensor<half> gmTensor6;
  gmTensor6.SetGlobalBuffer((__gm_half*)v2, i4);
  GlobalTensor<half> gmTensor7;
  gmTensor7.SetGlobalBuffer((__gm_half*)v3, i4);
  TPipe pipe;
  TQueuePosition<VECOUT, 1> outQueue8;
  pipe.InitBuffer(outQueue8, i4 * sizeof(half));
  LocalTensor<half> localTensor9 = outQueue8.AllocTensor<half>();
  TQueuePosition<VECIN, 1> inQueue9;
  pipe.InitBuffer(inQueue9, i4 * sizeof(half));
  LocalTensor<half> localTensor10 = inQueue9.AllocTensor<half>();
  TQueuePosition<VECIN, 1> inQueue12;
  pipe.InitBuffer(inQueue12, i4 * sizeof(half));
  LocalTensor<half> localTensor11 = inQueue12.AllocTensor<half>();
  DataCopy(localTensor10, gmTensor5, i4);
  DataCopy(localTensor11, gmTensor6, i4);
  inQueue12.Enqueue(localTensor11);
  inQueue10.Enqueue(localTensor11);
  LocalTensor<half> localTensor14 = inQueue10.DeQueue<half>();
  LocalTensor<half> localTensor15 = inQueue12.DeQueue<half>();
  Add(localTensor9, localTensor14, localTensor15);
  outQueue8.Enqueue(localTensor9);
  inQueue10.FreeTensor(localTensor10);
  inQueue12.FreeTensor(localTensor12);
  DataCopy(gmTensor7, localTensor9, i4);
  outQueue8.FreeTensor(localTensor9);
  return;
}

```

## References

<sup>1</sup> Ascend. Ascend C [Computer software]. <https://www.hiascend.com/en>  
<sup>2</sup> LLVM. Torch-MLIR [Computer software]. <https://github.com/llvm/torch-mlir>  
<sup>3</sup> LLVM. MLIR EmitC [Computer software]. <https://github.com/llvm/mlir-include/mlir/Dialect/EmitC>