



# 面向RISC-V同构融合CPU处理器的 Triton算子编译器设计与实践

进迭时空 黄竞辉

# CONTENTS

## 目录

- 01 Why Triton
- 02 RISC-V Triton-CPU软件栈
- 03 RISC-V Triton-CPU实践
- 04 RoadMap

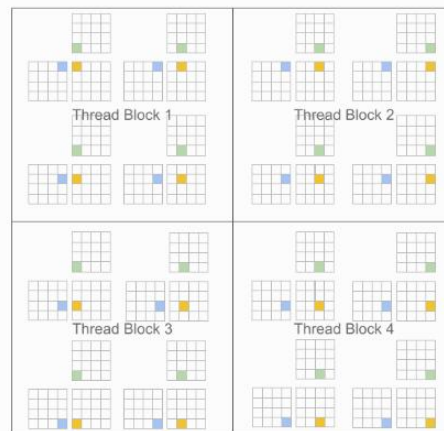
# Why Triton

- 编程模型收敛：多线程、Tile Base
- CUDA兼容性问题
- Python First or Pytorch First
- 新AI架构需要一个活跃可持续的DSL社区



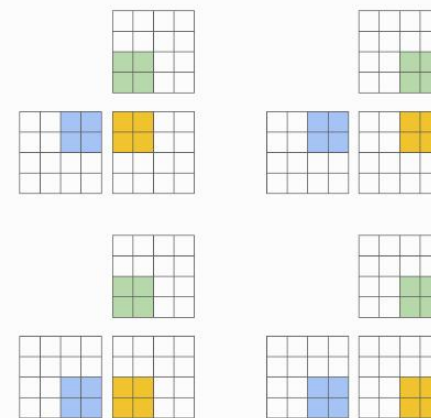
### CUDA Programming Model (Scalar Program, Blocked Threads)

```
#pragma parallel
for(int m = 0; m < M; m++){
  #pragma parallel
  for(int n = 0; n < N; n++){
    float acc = 0;
    for(int k = 0; k < K; k++){
      acc += A[m, k] * B[k, n];
    }
    C[m, n] = acc;
  }
}
```



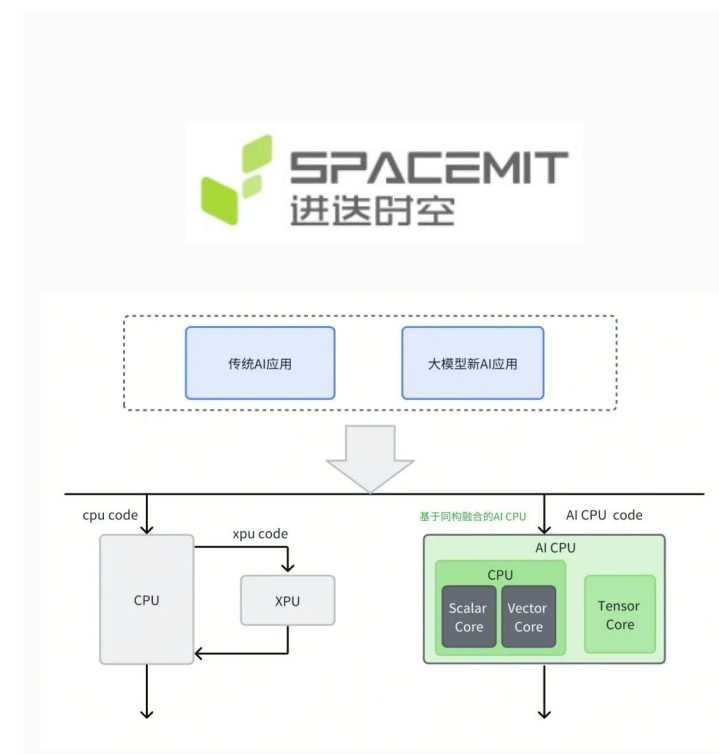
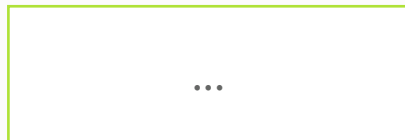
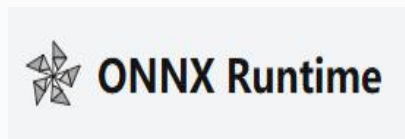
### Triton Programming Model (Blocked Program, Scalar Threads)

```
#pragma parallel
for(int m = 0; m < M; m += MB)
  #pragma parallel
  for(int n = 0; n < N; n += NB){
    float acc[MB, NB] = 0;
    for(int k = 0; k < K; k += KB)
      acc += A[m:m+MB, k:k+KB]
        ⊗ B[k:k+KB, n:n+NB];
    C[m:m+MB, n:n+NB] = acc;
  }
}
```

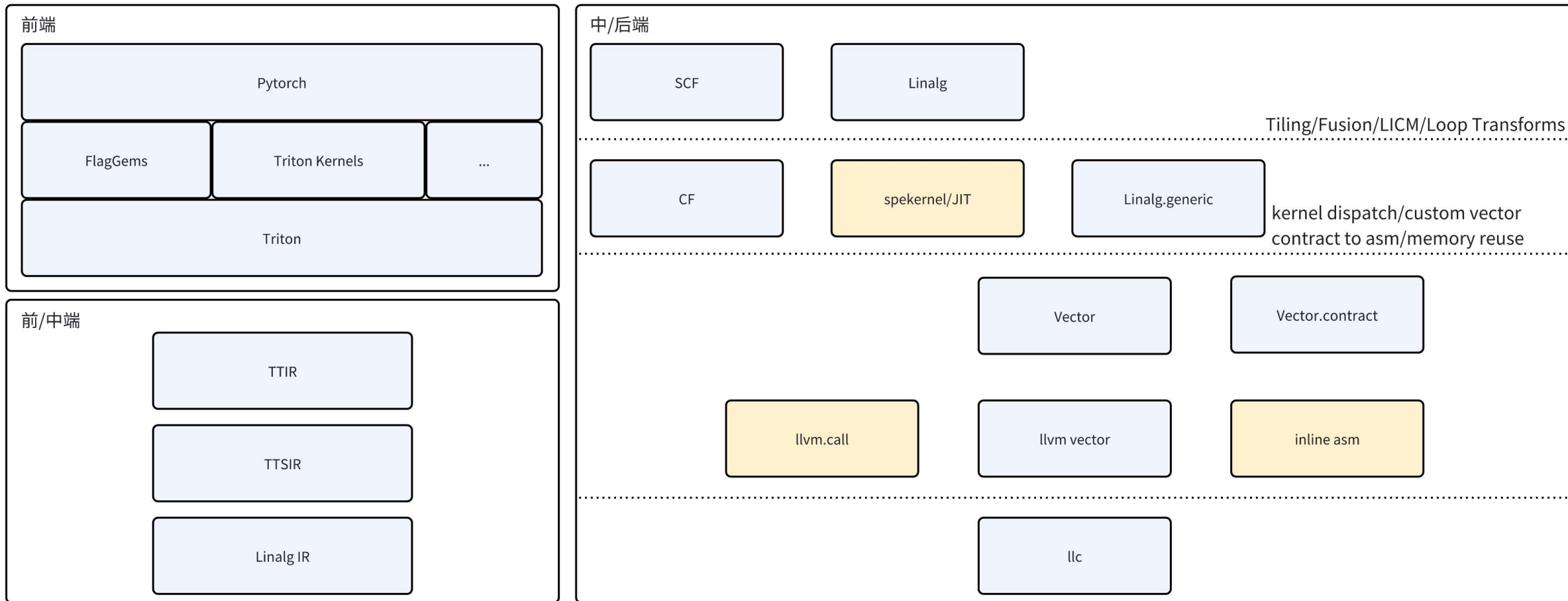


# RISCV Triton-CPU软件栈

- 统一内存，统一架构，消除Host-Device隔阂
- AICPU提供完整的语法支持
- 不止作为Pytorch加速后端，也直接提供Kernel API



# RISCV Triton-CPU软件栈



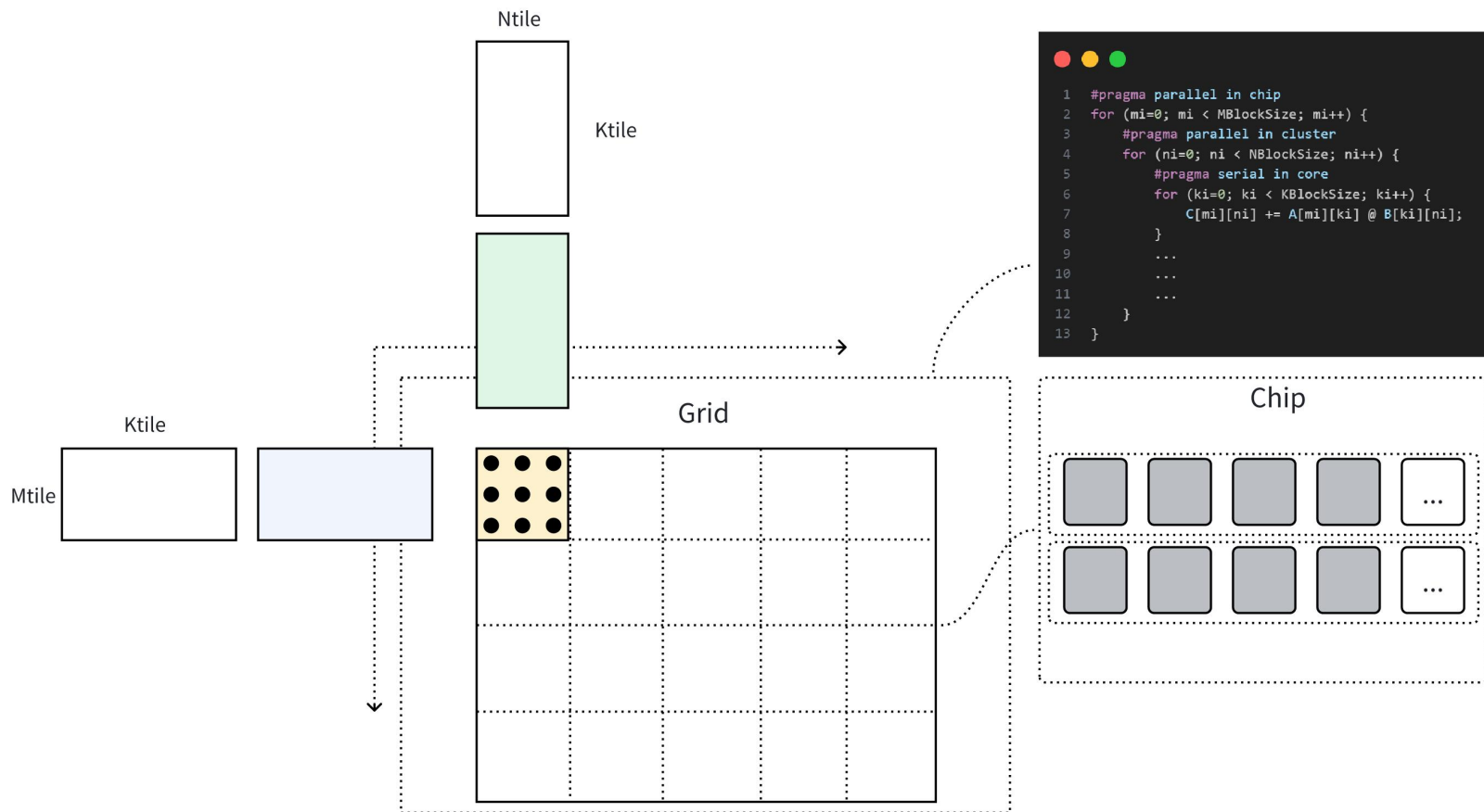
# RISCV Triton-CPU实践-网格化

- 计算网格化

- Grid ND-> Die/Cluster/Core, 提升局部性

- Cluster特化

- Kernel编程视角视为单核
- Share Memory
- 核内向量化、Tensor化



# RISCV Triton-CPU实践-块级访存

- 着重block\_ptr的使用，更贴合CPU编程思路
- 保持Load with Mask的使用方式
- 以Cluster视角的多线程kernel与单线程一致

```
1 %6 = tt.make_tensor_ptr %arg0, [%3, %4], [%5, %c1_i64], [%2, %c0_i32] {order = array<i32: 1, 0>} :  
2 %7 = arith.muli %1, %c16_i32 : i32 loc(#loc6)  
3 %8 = arith.extsi %arg4 : i32 to i64 loc(#loc7)  
4 %9 = arith.extsi %arg7 : i32 to i64 loc(#loc7)  
5 %10 = tt.make_tensor_ptr %arg1, [%4, %8], [%9, %c1_i64], [%c0_i32, %7] {order = array<i32: 1, 0>} :  
6 %11 = tt.load %6 {boundaryCheck = array<i32: 0, 1>} : !tt.ptr<tensor<16x1024xf32>> loc(#loc8)  
7 %12 = tt.load %10 {boundaryCheck = array<i32: 0, 1>} : !tt.ptr<tensor<1024x16xf32>> loc(#loc9)  
8 %13 = tt.dot %11, %12, %cst : tensor<16x1024xf32> * tensor<1024x16xf32> -> tensor<16x16xf32> loc(#loc10)  
9 %14 = arith.extsi %arg8 : i32 to i64 loc(#loc11)  
10 %15 = tt.make_tensor_ptr %arg2, [%3, %8], [%14, %c1_i64], [%2, %7] {order = array<i32: 1, 0>} : <te  
11 tt.store %15, %13 {boundaryCheck = array<i32: 0, 1>} : !tt.ptr<tensor<16x16xf32>> loc(#loc12)
```

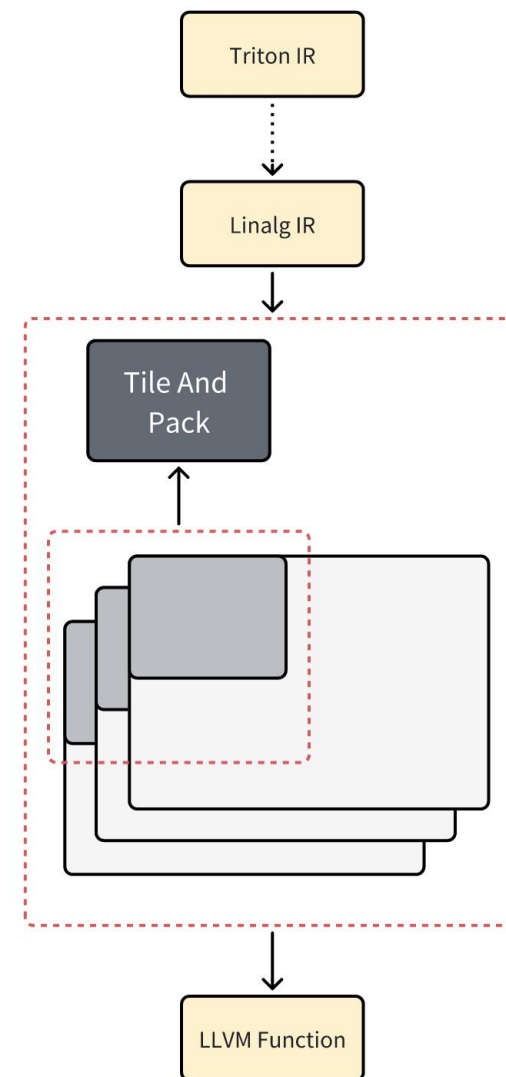
```
1 pid_m = tl.program_id(0)  
2 pid_n = tl.program_id(1)  
3  
4 a_block_ptr = tl.make_block_ptr(  
5     base=a_ptr,  
6     shape=[M, K],  
7     strides=[stride_am, stride_ak],  
8     offsets=[pid_m * BLOCK_SIZE_M, 0],  
9     block_shape=[BLOCK_SIZE_M, BLOCK_SIZE_K],  
10    order=[1, 0],  
11 )  
12  
13 b_block_ptr = tl.make_block_ptr(  
14     base=b_ptr,  
15     shape=[K, N],  
16     strides=[stride_bk, stride_bn],  
17     offsets=[0, pid_n * BLOCK_SIZE_N],  
18     block_shape=[BLOCK_SIZE_K, BLOCK_SIZE_N],  
19     order=[1, 0],  
20 )  
21  
22 accumulator = tl.zeros((BLOCK_SIZE_M, BLOCK_SIZE_N), dtype=tl.float32)  
23  
24 if EVEN_K:  
25     a = tl.load(a_block_ptr, boundary_check=(0, 1))  
26     b = tl.load(b_block_ptr, boundary_check=(0, 1))  
27     accumulator += tl.dot(a, b, allow_tf32=False)  
28 else:  
29     for k in range(0, tl.cdiv(K, BLOCK_SIZE_K)):  
30         a = tl.load(a_block_ptr, boundary_check=(0, 1))  
31         b = tl.load(b_block_ptr, boundary_check=(0, 1))  
32         accumulator += tl.dot(a, b, allow_tf32=False)  
33     a_block_ptr = tl.advance(a_block_ptr, (0, BLOCK_SIZE_K))  
34     b_block_ptr = tl.advance(b_block_ptr, (BLOCK_SIZE_K, 0))
```



# RISCV Triton-CPU实践-Gemm

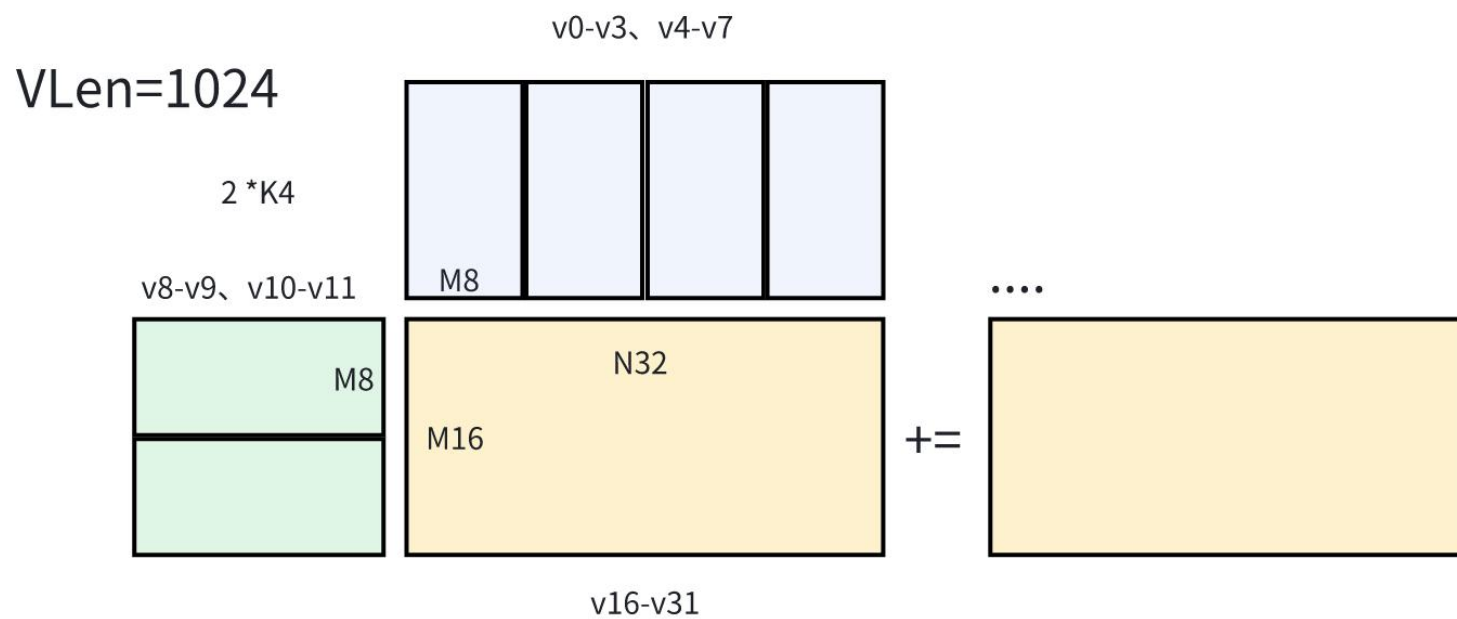
- spekernel对齐linalg.pack/mmt4d/...语义
- mmt4d->ime vmadot/vector load store
- pack->vector

```
1 %subview_14 = memref.subview %reinterpret_cast_6[0, 0] [%77, %83] [1, 1] : memref<?x?xf32, strided<[?], ?>> to memref<?x?xf32, strided<[?], ?>>
2 %subview_15 = memref.subview %reinterpret_cast_7[0, 0] [%77, %84] [1, 1] : memref<?x?xf32, strided<[?], ?>> to memref<?x?xf32, strided<[?], ?>>
3 %subview_16 = memref.subview %alloc_4[0, 0] [%77, %83] [1, 1] : memref<32x32xf32> to memref<?x?xf32, strided<[?], ?>>
4 %subview_17 = memref.subview %alloc_4[0, %83] [%77, %84] [1, 1] : memref<32x32xf32> to memref<?x?xf32, strided<[?], ?>>
5 memref.copy %subview_14, %subview_16 : memref<?x?xf32, strided<[?], ?>, offset: ?>> to memref<?x?xf32, strided<[?], ?>, offset: ?>>
6 memref.copy %subview_15, %subview_17 : memref<?x?xf32, strided<[?], ?>, offset: ?>> to memref<?x?xf32, strided<[?], ?>, offset: ?>>
7 spekernel.generic "pack" ins(%alloc_3 : memref<64x32xf32>) outs(%alloc_0 : memref<16x4x4x8xf32>)
8 spekernel.generic "pack" ins(%alloc_4 : memref<32x32xf32>) outs(%alloc_1 : memref<2x4x16x8xf32>)
9 spekernel.generic "mmt4d" ins(%alloc_0, %alloc_1 : memref<16x4x4x8xf32>, memref<2x4x16x8xf32>)
10 %85 = arith.addi %arg16, %c32 : index
11 %86 = arith.addi %arg17, %c34 : index
12 scf.yield %85, %86 : index, index
```



# RISCV Triton-CPU实践-Gemm

- Gemm Analysis
  - 寄存器用满, M/N First->K
  - Core上缓存用满, K First->M/N



# RISCV Triton-CPU实践-Gemm Lower

- Gemm Codegen

```
1  %3 = scf.for %arg3 = %c0 to %c512 step %c128 iter_args(%arg4 = %arg2) -> (tensor<256x512xf32>) {
2  %extracted_slice = tensor.extract_slice %arg1[0, %arg3] [256, 128] [1, 1] : tensor<256x512xf32> to tensor<256x128xf32>
3  %extracted_slice_0 = tensor.extract_slice %arg4[0, %arg3] [256, 128] [1, 1] : tensor<256x512xf32> to tensor<256x128xf32>
4  %4 = spekernel.generic "pack" ins(%extracted_slice : tensor<256x128xf32>) outs(%1 : tensor<4x8x32x32xf32>) (%c256, %c128, %c4, %c16)
5  %5 = scf.for %arg5 = %c0 to %c256 step %c16 iter_args(%arg6 = %extracted_slice_0) -> (tensor<256x128xf32>) {
6  %extracted_slice_1 = tensor.extract_slice %arg0[%arg5, 0] [16, 256] [1, 1] : tensor<256x256xf32> to tensor<16x256xf32>
7  %extracted_slice_2 = tensor.extract_slice %arg6[%arg5, 0] [16, 128] [1, 1] : tensor<256x128xf32> to tensor<16x128xf32>
8  %6 = spekernel.generic "pack" ins(%extracted_slice_1 : tensor<16x256xf32>) outs(%0 : tensor<1x8x16x32xf32>) (%c16, %c256, %c16, %c4)
9  %7 = spekernel.generic "pack" ins(%extracted_slice_2 : tensor<16x128xf32>) outs(%2 : tensor<1x4x16x32xf32>) (%c16, %c128, %c16, %c4)
10 %8 = spekernel.generic "mmt4d" ins(%6, %4 : tensor<1x8x16x32xf32>, tensor<4x8x32x32xf32>) outs(%7 : tensor<1x4x16x32xf32>) (%c16, %c4, %c16, %c4)
11 %9 = spekernel.generic "unpack" ins(%8 : tensor<1x4x16x32xf32>) outs(%extracted_slice_2 : tensor<16x128xf32>) (%c1, %c4, %c16, %c4)
12 %inserted_slice_3 = tensor.insert_slice %9 into %arg6[%arg5, 0] [16, 128] [1, 1] : tensor<16x128xf32> into tensor<256x128xf32>
13   scf.yield %inserted_slice_3 : tensor<256x128xf32>
14 }
15 %inserted_slice = tensor.insert_slice %5 into %arg4[0, %arg3] [256, 128] [1, 1] : tensor<256x128xf32> into tensor<256x512xf32>
16   scf.yield %inserted_slice : tensor<256x512xf32>
17 }
18 return %3 : tensor<256x512xf32>
19 }
```

# RISCV Triton-CPU实践-Gemm

- Gemm Codegen

```
1 // load b
2 %0 = vector.load [...] : memref<?x?x32x4xf32>, vector<4x32xf32>
3 // load a
4 %1 = vector.load [...] : memref<?x?x16x4xf32>, vector<2x32xf32>
5 // vfmadot -> 2x8x4 @ 4x8x4 => 2x4x8x8
6 %2 = vector.contract {...} %1, %0 : vector<2x32xf32>, vector<4x32xf32> into vector<16x32xf32>
7
```

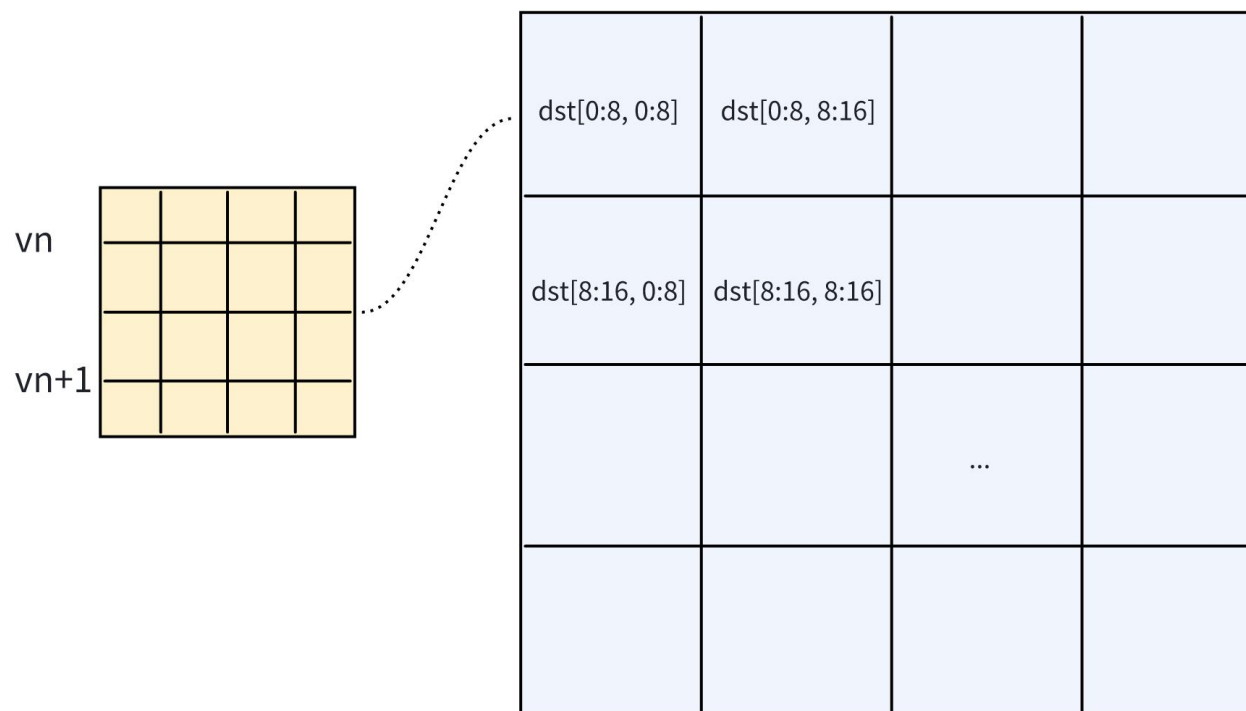
```
1 // pack A -> [m_block_size, k_block_size, mb, kb]
2 // pack B -> [n_block_size, k_block_size, nb, kb]
3 // packA @ packC -> [m_block_size, k_block_size, mb, nb]
4 for (mi : m_block_size) {
5     for (ni : n_block_size) {
6         // cf loop
7         for (ki : k_block_size) {
8             // vector.contract
9             C[mi, ni, :, :] += A[mi, ki, :, :] @ B[ni, ki, :, :]
10        }
11    }
12 }
```

```
1 vle32.v ...
2 vle32.v ...
3 ...
4 vmadot $(acc:0), $(lhs:0), $(rhs:0)
5 vmadot $(acc:2), $(lhs:0), $(rhs:1)
6 ...
7 vmadot $(acc:8), $(lhs:1), $(rhs:0)
8 vmadot $(acc:10), $(lhs:1), $(rhs:1)
9 ...
```



# RISCV Triton-CPU实践-Gemm

- Gemm Codegen



# RISCV Triton-CPU实践-Conv

- Conv Codegen

```
1 // kernel_shape=[3,3]
2 // %acc: vector<16x32xf32>
3 // load b
4 %0 = vector.load [...] : memref<?x?x32x4xf32>, vector<9x32xf32>
5 // load a h0
6 %1 = vector.load [...] : memref<?x?x16x4xf32>, vector<3x32xf32>
7 // vfmadotn -> 3x8x4 @ 9x8x4 => 2x4x8x8
8 %2 = vector.contract {...} %1, %0 %acc : vector<3x32xf32>, vector<9x32xf32> into vector<16x32xf32>
9 // load a h1
10 %3 = vector.load [...] : memref<?x?x16x4xf32>, vector<3x32xf32>
11 // vfmadotn -> 3x8x4 @ 9x8x4 => 2x4x8x8
12 %4 = vector.contract {...} %3, %0 %2 : vector<3x32xf32>, vector<9x32xf32> into vector<16x32xf32>
13 // load a h2
14 %5 = vector.load [...] : memref<?x?x16x4xf32>, vector<3x32xf32>
15 // vfmadotn -> 3x8x4 @ 9x8x4 => 2x4x8x8
16 %6 = vector.contract {...} %5, %0 %4 : vector<3x32xf32>, vector<9x32xf32> into vector<16x32xf32>
17
```

```
1 vle32.v ...
2 vle32.v ...
3 ...
4 vmadot $(acc:0), $(lhs:0), $(rhs:0)
5 vmadot1 $(acc:0), $(lhs:0), $(rhs:0)
6 vmadot2 $(acc:0), $(lhs:0), $(rhs:0)
7 ...
8 vmadot $(acc:2), $(lhs:0), $(rhs:1)
9 vmadot1 $(acc:2), $(lhs:0), $(rhs:1)
10 vmadot2 $(acc:2), $(lhs:0), $(rhs:1)
11 ...
12 vse32.v
13 vse32.v
```



# RoadMap

## 开源

- Triton前中端以及后端Compiler逐步开源

## 更彻底的Codegen

- 矩阵计算序列与Vector融合调优
- AME
- triton inline asm

## 分布式支持

- 关注到triton-distribute, 分布式通信原语及硬件支持





# Thanks

**以RISC-V架构数智未来**  
RISC-V ARCHITECTURE FOR INTELLIGENT FUTURE

进迭时空（杭州）科技有限公司  
SPACEMIT (Hangzhou) Technology Co., Ltd

电话：0571-89000775

网址：[www.spacemit.com](http://www.spacemit.com)

地址：浙江省杭州市余杭区五常街道未来天空中心b座701

