# Buddy-CAAS:
# Compiler As A Service for MLIR

**Speaker:** Hongbin Zhang | hongbin2019@iscas.ac.cn

**Authors:** Zikang Liu (ICT), Hongbin Zhang (ISCAS), Bofan Zhang (USYD)

Buddy Compiler

# Compiler As A Service for MLIR ( https://buddy.isrc.ac.cn/ )

**Buddy Compiler**

*Buddy-CAAS serves MLIR users and developers to configure the pass pipeline and demonstrate on multiple backends quickly and smoothly.*

# Compiler As A Service for MLIR ( https://buddy.isrc.ac.cn/ )

**Buddy Compiler**

*Buddy-CAAS serves MLIR users and developers to* **configure the pass pipeline** *and* **demonstrate on multiple backends** *quickly and smoothly.*

# Compiler As A Service for MLIR ( https://buddy.isrc.ac.cn/ )

Buddy Compiler

*Buddy-CAAS serves MLIR users and developers to configure the pass pipeline and demonstrate on multiple backends quickly and smoothly.*

# Compiler As A Service for MLIR ( https://buddy.isrc.ac.cn/ )

Buddy Compiler

*Buddy-CAAS serves MLIR users and developers to configure the pass pipeline and demonstrate on multiple backends quickly and smoothly.*

# Compiler As A Service for MLIR ( https://buddy.isrc.ac.cn/ )

Buddy Compiler

*Buddy-CAAS serves MLIR users and developers to* configure the pass pipeline *and* demonstrate on multiple backends *quickly and smoothly.*



Buddy Compiler

We ❤️ MLIR Playground

```
20    func.func @main(){
21        // Set up dims.
22        %cM = arith.constant 4 : index
23        %cN = arith.constant 4 : index
24        %cK = arith.constant 4 : index
25
26        // Set Init Value.
27        %cf1 = arith.constant 1.0 : f32
28
29        %A = memref.alloc(%cM, %cK) : memref<?x?xf32>
30        %B = memref.alloc(%cK, %cN) : memref<?x?xf32>
31        %C = memref.alloc(%cM, %cN) : memref<?x?xf32>
32
```

Lower Box    Translate Box    Compile Box    Link Box    Execute Box

-convert-linalg-to-loops    -lower-affine    -convert-vector-to-llvm    -convert-memref-to-llvm    -convert-arith-to-llvm    -convert-func-to-llvm    -reconcile-unrealized-casts

-convert-scf-to-cf

*Drag and Drop a Pass*

Result Output

```
1     module {
2         func.func private @printMemrefF32(memref<*xf32>)
3         func.func @matmul(%arg0: memref<?x?xf32>, %arg1: memref<?x?xf32>, %arg2: memref<?x?xf32>) {
4             %0 = llvm.mlir.constant(0 : index) : i64
5             %1 = builtin.unrealized_conversion_cast %0 : i64 to index
6             %2 = llvm.mlir.constant(1 : index) : i64
7             %3 = builtin.unrealized_conversion_cast %2 : i64 to index
8             %dim = memref.dim %arg0, %1 : memref<?x?xf32>
9             %4 = builtin.unrealized_conversion_cast %dim : index to i64
10            %dim_0 = memref.dim %arg0, %3 : memref<?x?xf32>
11            %5 = builtin.unrealized_conversion_cast %dim_0 : index to i64
12            %dim_1 = memref.dim %arg1, %3 : memref<?x?xf32>
13            %6 = builtin.unrealized_conversion_cast %dim_1 : index to i64
14            cf.br ^bb1(%1 : index)
15        ^bb1(%7: index):  // 2 preds: ^bb0, ^bb8
16            %8 = builtin.unrealized_conversion_cast %7 : index to i64
17            %9 = llvm.icmp "slt" %8, %4 : i64
18            cf.cond_br %9, ^bb2, ^bb9
19        ^bb2:  // pred: ^bb1
20            cf.br ^bb3(%1 : index)
21        ^bb3(%10: index):  // 2 preds: ^bb2, ^bb7
```

# Compiler As A Service for MLIR ( https://buddy.isrc.ac.cn/ )

**Buddy Compiler**

*Buddy-CAAS serves MLIR users and developers to configure the pass pipeline and demonstrate on multiple backends quickly and smoothly.*



**Buddy Compiler**

We ♥ MLIR Playground

```
14    linalg.matmul
15      ins(%a, %b: memref<?x?xf32>, memref<?x?xf32>)
16      outs(%c:memref<?x?xf32>)
17    return
18  }
19
20  func.func @main(){
21    // Set up dims.
22    %cM = arith.constant 4 : index
23    %cN = arith.constant 4 : index
24    %cK = arith.constant 4 : index
25
26    // Set Init Value.
27    %cf1 = arith.constant 1.0 : f32
28
29    %A = memref.alloc(%cM, %cK) : memref<?x?xf32>
30    %B = memref.alloc(%cK, %cN) : memref<?x?xf32>
31    %C = memref.alloc(%cM, %cN) : memref<?x?xf32>
32
33    linalg.fill
```

**Choose the Execute Box**    Lower Box    Translate Box    Compile Box    Link Box    **Execute Box**

-cpu    rv64,x-v=true,vlen=128                                           ✏ ⊙ ✅ *Run!!!*

Log    **Check Your Result** 👀                                                    ︿

```
Unranked Memref base@ = 0x24f50 rank = 2 offset = 0 sizes = [4, 4] strides = [4, 1] data =
[[5, 5, 5, 5],
[5, 5, 5, 5],
[5, 5, 5, 5],
[5, 5, 5, 5]]
```

# Compiler As A Service for MLIR ( https://buddy.isrc.ac.cn/ )

**Buddy Compiler**

*Buddy-CAAS serves MLIR users and developers to configure the pass pipeline and demonstrate on multiple backends quickly and smoothly.*

Buddy-CAAS serves MLIR users and developers to **configure the pass pipeline** and **demonstrate on multiple backends** quickly and smoothly.



Buddy Compiler

We ❤ MLIR Playground

```
15        ins(%a, %b: memref<?x?xf32>, memref<?x?xf32>)
16        outs(%c:memref<?x?xf32>)
17      return
18    }
19
20    func.func @main(){
21      // Set up dims.
22      %cM = arith.constant 4 : index
23      %cN = arith.constant 4 : index
24      %cK = arith.constant 4 : index
25
26      // Set Init Value.
27      %cf1 = arith.constant 1.0 : f32
28
29      %A = memref.alloc(%cM, %cK) : memref<?x?xf32>
30      %B = memref.alloc(%cK, %cN) : memref<?x?xf32>
31      %C = memref.alloc(%cM, %cN) : memref<?x?xf32>
32
33      linalg.fill
```
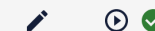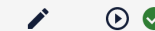
*Choose the Compile Box*     Lower Box     Translate Box     Compile Box     Link Box     Execute Box

-mtriple     riscv64     -target-abi     lp64d     -mattr=+m,+d,+v     -riscv-v-vector-bits-min=128

Result Output     *Check the Generated Assembly Code*

```
35        bge t3, a4, .LBB0_4
36    .LBB0_7:                        #   Parent Loop BB0_2 Depth=1
37                                    #     Parent Loop BB0_5 Depth=2
38                                    # =>    This Inner Loop Header: Depth=3
39        mul t4, t1, a5
40        add t4, t3, t4
41        slli    t4, t4, 2
42        add t4, a1, t4
43        flw ft0, 0(t4)
44        mul t4, t3, a0
45        add t4, t4, t2
46        slli    t4, t4, 2
47        add t4, a6, t4
48        flw ft1, 0(t4)
49        mul t4, t1, a7
50        add t4, t4, t2
```

**Buddy Compiler**

# Thanks

## https://buddy.isrc.ac.cn/