

# Dialects in LLVM IR

An Example using Cooperative Matrices

# Cooperative Matrices

- AI applications require lots of matrix-matrix multiplications
- New Vulkan shader extension introduces the cooperative matrix type
- Special data layout depending on hardware support
- Some requirements on AMD GPUs:
  - When calculating  $A \cdot B + C$ ,  $A$  needs to be transposed



# Transpose as high-level operation

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {
5   %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)
6   ret <8 x float> %a.t
7 }
```

```
558 define <8 x float> @test_transpose_simple(<8 x float> %a) {
559   %1 = call i32 @_llvm.amdgcn.mbcnt.lo(i32 -1, i32 0)
560   %2 = call i32 @_llvm.amdgcn.mbcnt.hi(i32 -1, i32 %1)
561   %3 = and i32 %2, 1
562   %4 = icmp eq i32 %3, 0
563   %5 = extractelement <8 x float> %a, i64 0
564   %6 = bitcast float %5 to i32
565   %7 = call i32 @_llvm.amdgcn.mov.dpp8.i32(i32 %6, i32 14570689)
566   %8 = bitcast i32 %7 to float
567   %9 = extractelement <8 x float> %a, i64 1
568   %10 = bitcast float %9 to i32
569   %11 = call i32 @_llvm.amdgcn.mov.dpp8.i32(i32 %10, i32 14570689)
570   %12 = bitcast i32 %11 to float
571   %13 = extractelement <8 x float> %a, i64 2
572   %14 = bitcast float %13 to i32
573   %15 = call i32 @_llvm.amdgcn.mov.dpp8.i32(i32 %14, i32 14570689)
574   %16 = bitcast i32 %15 to float
575   %17 = extractelement <8 x float> %a, i64 3
576   %18 = bitcast float %17 to i32
577   %19 = call i32 @_llvm.amdgcn.mov.dpp8.i32(i32 %18, i32 14570689)
578   %20 = bitcast i32 %19 to float
579   %21 = extractelement <8 x float> %a, i64 4
580   %22 = bitcast float %21 to i32
581   %23 = call i32 @_llvm.amdgcn.mov.dpp8.i32(i32 %22, i32 14570689)
582   %24 = bitcast i32 %23 to float
583   %25 = extractelement <8 x float> %a, i64 5
584   %26 = bitcast float %25 to i32
585   %27 = call i32 @_llvm.amdgcn.mov.dpp8.i32(i32 %26, i32 14570689)
586   %28 = bitcast i32 %27 to float
587   %29 = extractelement <8 x float> %a, i64 6
588   %30 = bitcast float %29 to i32
589   %31 = call i32 @_llvm.amdgcn.mov.dpp8.i32(i32 %30, i32 14570689)
590   %32 = bitcast i32 %31 to float
591   %33 = extractelement <8 x float> %a, i64 7
592   %34 = bitcast float %33 to i32
593   %35 = call i32 @_llvm.amdgcn.mov.dpp8.i32(i32 %34, i32 14570689)
```



## Working with new intrinsics – current approach

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {  
5   %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)  
6   ret <8 x float> %a.t  
7 }
```

## Working with new intrinsics – current approach

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {  
5     %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)  
6     ret <8 x float> %a.t  
7 }
```

```
const static char CooperativeMatrixTranspose[] = "lgc.cooperative.matrix.transpose";
```

## Working with new intrinsics – current approach

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {  
5     %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)  
6     ret <8 x float> %a.t  
7 }  
  
const static char CooperativeMatrixTranspose[] = "lgc.cooperative.matrix.transpose";  
  
std::string callName( lgcName::CooperativeMatrixTranspose);  
Value *args[] = {[0]=matrix, [1]=elementType, [2]=layout};  
  
CallInst *result = CreateNamedCall(funcName: callName, retTy: matrix->getType(),  
| | | | | | | | | | args, attribs: {[0]=Attribute::ReadOnly, [1]=Attribute::WillReturn});
```

## Working with new intrinsics – current approach

```
const static char CooperativeMatrixTranspose[] = "lgc.cooperative.matrix.transpose";  
  
std::string callName(s: lgcName::CooperativeMatrixTranspose);  
Value *args[] = {[0]=matrix, [1]=elementType, [2]=layout};  
  
CallInst *result = CreateNamedCall(funcName: callName, retTy: matrix->getType(),  
| | | | | | | | | | | | | | | | | | | | | | | args, attrs: {[0]=Attribute::ReadOnly, [1]=Attribute::WillReturn});  
  
if (auto *call: CallInst * = dyn_cast<CallInst>(Val: input)) {  
    if (auto *callee: Function * = call->getCalledFunction()) {  
        if (callee->getName().starts_with(Prefix: lgcName::CooperativeMatrixTranspose)) {  
            Value *src = call->getArgOperand(i: 0);
```

## Working with new intrinsics – current approach

```
const static char CooperativeMatrixTranspose[] = "lgc.cooperative.matrix.transpose";  
const static char CooperativeMatrixTransposeWithLayout[] = "lgc.cooperative.matrix.transposeWithLayout";  
  
std::string callName(s: lgcName::CooperativeMatrixTranspose);  
Value *args[] = {[0]=matrix, [1]=elementType, [2]=layout};  
  
CallInst *result = CreateNamedCall(funcName: callName, retTy: matrix->getType(),  
| | | | | | | | | | | | | | | | | | | | | | | args, attribs: {[0]=Attribute::ReadOnly, [1]=Attribute::WillReturn});  
  
if (auto *call: CallInst * = dyn_cast<CallInst>(Val: input)) {  
    if (auto *callee: Function * = call->getCalledFunction()) {  
        if (callee->getName().starts_with(Prefix: lgcName::CooperativeMatrixTranspose)) {  
            Value *src = call->getArgOperand(i: 0);
```

# Problems

- New intrinsics are all CallInst
  - Common pattern: iterate over instructions, cast to CallInst, compare name string
  - Accessing arguments with an index
- Not developer-friendly

## Solution: llvm-dialects

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {
5   %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)
6   ret <8 x float> %a.t
7 }
```

```
def LgcCmDialect : Dialect {
  let name = "lgc.cooperative.matrix";
  let cppNamespace = "lgc::cooperativeMatrix";
}
```

## Solution: llvm-dialects

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {
5   %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)
6   ret <8 x float> %a.t
7 }
```

```
def LgcCmDialect : Dialect {
  let name = "lgc.cooperative.matrix";
  let cppNamespace = "lgc::cooperativeMatrix";
}
```

## Solution: llvm-dialects

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {
5   %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)
6   ret <8 x float> %a.t
7 }
```

```
def LgcCmDialect : Dialect {
  let name = "lgc.cooperative.matrix";
  let cppNamespace = "lgc::cooperativeMatrix";
}
```

## Solution: llvm-dialects

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {
5   %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)
6   ret <8 x float> %a.t
7 }
```

```
def LgcCmDialect : Dialect {
  let name = "lgc.cooperative.matrix";
  let cppNamespace = "lgc::cooperativeMatrix";
}
```

## Solution: llvm-dialects

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {
5   %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)
6   ret <8 x float> %a.t
7 }
```

```
def CoopMatrixTransposeOp : Op<LgcCmDialect, "transpose", [Memory<[(read)]>, WillReturn[]]> {
  let arguments = (ins (FixedVectorType $element_type, $num_elements):$matrix, I32:$elem_type, I32:$layout);
  let results = (outs (eq $matrix):$result);

  let summary = "transpose a matrix";
  let description = [
    Given a cooperative matrix, the element type and the matrix layout, returns
    the transposed matrix.
  ];
}
```

## Solution: llvm-dialects

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {  
5   %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)  
6   ret <8 x float> %a.t  
7 }
```

```
def CoopMatrixTransposeOp : Op<LgcCmDialect, "transpose", [Memory<[(read)], WillReturn]> {
```

## Solution: llvm-dialects

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {  
5   %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)  
6   ret <8 x float> %a.t  
7 }
```

```
def CoopMatrixTransposeOp : Op<LgcCmDialect, "transpose", [Memory<[(read)], WillReturn]> {
```

```
CallInst *result = CreateNamedCall(funcName: callName, retTy: matrix->getType(),  
| | | | | | | | | | | | | | | | | | | | | | args, attrs: {[0]=Attribute::ReadOnly, [1]=Attribute::WillReturn});
```

## Solution: llvm-dialects

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {
5   %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)
6   ret <8 x float> %a.t
7 }
```

```
def CoopMatrixTransposeOp : Op<LgcCmDialect, "transpose", [Memory<[(read)]>, WillReturn] > {
  let arguments = (ins (FixedVectorType $element_type, $num_elements):$matrix, I32:$elem_type, I32:$layout);
  let results = (outs (eq $matrix):$result);
```

## Solution: llvm-dialects

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {
5   %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)
6   ret <8 x float> %a.t
7 }
```

```
def CoopMatrixTransposeOp : Op<LgcCmDialect, "transpose", [Memory<[(read)], WillReturn] > {
  let arguments = (ins [FixedVectorType $element_type, $num_elements]:$matrix, I32:$elem_type, I32:$layout);
  let results = (outs (eq $matrix):$result);
```

## Solution: llvm-dialects

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {
5   %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)
6   ret <8 x float> %a.t
7 }
```

```
def CoopMatrixTransposeOp : Op<LgcCmDialect, "transpose", [Memory<[(read)], WillReturn] > {
  let arguments = (ins (FixedVectorType $element_type, $num_elements):$matrix, I32:$elem_type, I32:$layout);
  let results = (outs (eq $matrix):$result);
```

## Solution: llvm-dialects

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {
5   %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)
6   ret <8 x float> %a.t
7 }
```

```
def CoopMatrixTransposeOp : Op<LgcCmDialect, "transpose", Memory<[(read)], WillReturn[]> {
  let arguments = (ins (FixedVectorType $element_type, $num_elements):$matrix, I32:$elem_type, I32:$layout);
  let results = (outs (eq $matrix):$result);
```

## Solution: llvm-dialects

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {
5   %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)
6   ret <8 x float> %a.t
7 }
```

```
def CoopMatrixTransposeOp : Op<LgcCmDialect, "transpose", [Memory<[(read)], WillReturn[]]> {
  let arguments = (ins (FixedVectorType $element_type, $num_elements):$matrix, I32:$elem_type, I32:$layout);
  let results = (outs (eq $matrix):$result);
```

## Solution: llvm-dialects

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {
5   %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)
6   ret <8 x float> %a.t
7 }
```

```
def CoopMatrixTransposeOp : Op<LgcCmDialect, "transpose", [Memory<[(read)], WillReturn] > {
  let arguments = (ins (FixedVectorType $element_type, $num_elements):$matrix, I32:$elem_type, I32:$layout);
  let results = (outs (eq $matrix) $result);
```

Many other constraints possible, see the  
[Constraints.md](#) in the [llvm-dialects repo](#)

## Solution: llvm-dialects

```
4 define <8 x float> @test_transpose_simple(<8 x float> %a) {
5   %a.t = call <8 x float> @lgc.cooperative.matrix.transpose.v8f32(<8 x float> %a, i32 1, i32 0)
6   ret <8 x float> %a.t
7 }
```

```
def CoopMatrixTransposeOp : Op<LgcCmDialect, "transpose", [Memory<[(read)], WillReturn[]> {
  let arguments = (ins (FixedVectorType $element_type, $num_elements):$matrix, I32:$elem_type, I32:$layout);
  let results = (outs (eq $matrix):$result);

  let summary = "transpose a matrix";
  let description = [
    Given a cooperative matrix, the element type and the matrix layout, returns
    the transposed matrix.
  ];
}
```

## Solution: llvm-dialects

```
def CoopMatrixTransposeOp : Op<LgcCmDialect, "transpose", [Memory<[(read)], WillReturn] > {  
    let arguments = (ins (FixedVectorType $element_type, $num_elements):$matrix, I32:$elem_type, I32:$layout);  
    let results = (outs (eq $matrix):$result);
```

## Solution: llvm-dialects

```
def CoopMatrixTransposeOp : Op<LgcCmDialect, "transpose", [Memory<[(read)], WillReturn] > {
    let arguments = (ins (FixedVectorType $element_type, $num_elements):$matrix, I32:$elem_type, I32:$layout);
    let results = (outs (eq $matrix):$result);

    CoopMatrixTransposeOp *result = create<CoopMatrixTransposeOp>(
        matrix, elementType, layout);

    std::string callName(s: lgcName::CooperativeMatrixTranspose);
    Value *args[] = {[0]=matrix, [1]=elementType, [2]=layout};

    CallInst *result = CreateNamedCall(funcName: callName, retTy: matrix->getType(),
        args, attrs: {[0]=Attribute::ReadOnly, [1]=Attribute::WillReturn});
}
```

## Solution: llvm-dialects

```
def CoopMatrixTransposeOp : Op<LgcCmDialect, "transpose", [Memory<[(read)], WillReturn] > {
    let arguments = (ins (FixedVectorType $element_type, $num_elements):$matrix, I32:$elem_type, I32:$layout);
    let results = (outs (eq $matrix):$result);

if (isa<cooperativeMatrix::CoopMatrixTransposeOp>(Val: input)) {
    - - - - -
}

if (auto *call: CallInst * = dyn_cast<CallInst>(Val: input)) {
    if (auto *callee: Function * = call->getCalledFunction()) {
        if (callee->getName().starts_with(Prefix: lgcName::CooperativeMatrixTranspose)) {
            Value *src = call->getArgOperand(i: 0);
        }
    }
}
```

## Solution: llvm-dialects

```
CoopMatrixTransposeOp *result = create<CoopMatrixTransposeOp>(
    matrix, elementType, layout);

static const auto visitor = You, 12 hours ago • Uncommitted changes
    llvm_dialects::VisitorBuilder<LowerCooperativeMatrix>()
    .add(fn: &LowerCooperativeMatrix::visitCooperativeMatrixTranspose)
    .build();

visitor.visit(&payload: *this, &: module);

void LowerCooperativeMatrix::visitCooperativeMatrixTranspose(
    cooperativeMatrix::CoopMatrixTransposeOp &coopTranspose) {
```

## Solution: llvm-dialects

```
def CoopMatrixTransposeOp : LgcCmOp<"transpose", [Memory<[(read)]>, WillReturn]> {
    let arguments = (ins (FixedVectorType $element_type, $num_elements):$matrix, I32:$elem_type, I32:$layout);
    let results = (outs (eq $matrix):$result);
```

## Solution: llvm-dialects

```
| J> {  
I32:$elem_type, I32:$layout);
```

```
transpose.getElemType()  
transpose.getLayout()
```

```
transpose->getArgOperand(i: 1)  
transpose->getArgOperand(i: 2)
```

## Additional features

- More complex constraints
- OpSet/OpMap working on classes of Operations
- ContextExtension

## Future Work

- Custom Type Definitions
- Human readable structs in metadata

Thank you for listening!