

Fused Multiply-Add Optimization Pass

Francisco F Cavazos Jr

Motivation

- Floating point intensive workloads often use multiply + add sequences
- FMA instruction: combines $A * B + C$ into a single instruction
- Benefits:
 - Higher Throughput
 - Reduced rounding error

What is FMA?

The operation fusedMultiplyAdd(x, y, z) computes $(x \times y) + z$ as if with unbounded range and precision, rounding only once to the destination format.

– IEEE-754-2008 (Definition 2.1.28)

Regular

```
%0 = fmul float %x, %y
%1 = fadd float %0, %z
ret float %1
```

= round(round($x * y$) + z)

FMA

```
%0 = call float @llvm.fma.f32(float %x, float %y, float %z)
ret float %0
```

= round($x * y + z$)

What is FMA?

```
1 #include <stdio.h>
2 #include <math.h>
3
4
5 float nofma(float x, float y, float z) {
6     return x * y + z;
7 }
8
9 float withfma(float x, float y, float z) {
10    return fmaf(x, y, z);
11 }
12
13 int main(void) {
14     float x = 1e8f;
15     float y = 1.0000001f;
16     float z = -1e8f;
17
18     printf("No FMA: %.10f\n", nofma(x, y, z));
19     printf("  FMA: %.10f\n", withfma(x, y, z));
20
21     return 0;
22 }
```

```
1 @.str = private unnamed_addr constant [15 x i8] c"No FMA: %.10f\0A\00", align 1
2 @.str.1 = private unnamed_addr constant [14 x i8] c"  FMA: %.10f\00", align 1
3
4 ; Function Attrs: noline nounwind ssp uwtable(sync)
5 define float @nofma(float noundef %0, float noundef %1, float noundef %2) #0 {
6     %4 = fmul float %0, %1
7     %5 = fadd float %4, %2
8     ret float %5
9 }
10
11 ; Function Attrs: noline nounwind ssp uwtable(sync)
12 define float @withfma(float noundef %0, float noundef %1, float noundef %2) #0 {
13     %4 = call float @llvm.fma.f32(float %0, float %1, float %2)
14     ret float %4
15 }
16
17 ; Function Attrs: nocallback nofree nosync nounwind speculatable willreturn memory(none)
18 declare float @llvm.fma.f32(float, float, float) #1
19
20 ; Function Attrs: noline nounwind ssp uwtable(sync)
21 define i32 @main() #0 {
22     %1 = call float @nofma(float noundef 1.000000e+08, float noundef 0x3FF000020000000, float noundef -1.000000e+08)
23     %2 = fpext float %1 to double
24     %3 = call i32 (ptr, ...) @printf(ptr noundef @.str, double noundef %2)
25     %4 = call float @withfma(float noundef 1.000000e+08, float noundef 0x3FF000020000000, float noundef -1.000000e+08)
26     %5 = fpext float %4 to double
27     %6 = call i32 (ptr, ...) @printf(ptr noundef @.str.1, double noundef %5)
28     ret i32 0
29 }
30
31 declare i32 @printf(ptr noundef, ...) #2
```

What is FMA?

We're computing:

$$1.0000001 \times (1 \times 10^8) - (1 \times 10^8) = 10$$

Results:

- No FMA = 8
- With FMA = 11.9209289551

```
Ferny tests % llc fma.ll -o fma.s; clang -o fma fma.s; ./fma
No FMA: 8.0000000000
FMA: 11.9209289551
```


What is FMA?

What we want to compute:

$$1.0000001 \times (1 \times 10^8) - (1 \times 10^8) = 10$$

What we actually compute:

$$1.000000119209289551 \times (1 \times 10^8) - (1 \times 10^8) \\ = 11.9209289551$$

Results:

- No FMA = 8
- With FMA = 11.9209289551

Project Objective

- Implement a custom LLVM pass that:
 - Detects `fmul` + `fadd` patterns
 - Replaces them with LLVM's FMA intrinsic
 - Checks hardware support
 - Ensures transformation is legal and profitable
- Evaluate across different tests

Creating the LLVM Pass

Pattern Detection

We want to detect the following pattern:

```
%mul = fmul double %x, %y
```

```
%add = fadd double %mul, %z
```

How?

1. Search the LLVM IR for all `fadd`
2. For each `fadd`, check if one of the operands opcode is `fmul`
3. If true, add the `fadd` + `fmul` combo to candidate list

Pattern Detection

We want to detect the following pattern:

%mul = fmul double %x, %y

%add = fadd double %mul, %z

```
1 bool matchesFMAPattern(BinaryOperator *FAdd) {
2   Value *Op0 = FAdd->getOperand(0);
3   Value *Op1 = FAdd->getOperand(1);
4
5   // Check if either operand is an fmul instruction
6   if (auto *FMul = dyn_cast<BinaryOperator>(Op0)) {
7     if (FMul->getOpcode() == Instruction::FMul) {
8       return true;
9     }
10  }
11
12  if (auto *FMul = dyn_cast<BinaryOperator>(Op1)) {
13    if (FMul->getOpcode() == Instruction::FMul) {
14      return true;
15    }
16  }
17  return false;
18 }
```

Pattern Detection

We can use the “llvm/IR/PatternMatch.h” header functions to improve it (using match)

```
1 bool matchesFMAPattern(BinaryOperator *FAdd) {
2   Value *X, *Y, *Z;
3
4   // Pattern 1: fadd (fmul X, Y), Z
5   if (match(FAdd, m_FAdd(m_FMul(m_Value(X), m_Value(Y)), m_Value(Z)))) {
6     return true;
7   }
8
9   // Pattern 2: fadd Z, (fmul X, Y) - commutative case
10  if (match(FAdd, m_FAdd(m_Value(Z), m_FMul(m_Value(X), m_Value(Y))))) {
11    return true;
12  }
13
14  return false;
15 }
```

Legality Conditions

- Test if FMA is supported on the target machine
- Implement a profitability check (only transform if beneficial, e.g., not in precise FP mode)

Legality Conditions

- Test if FMA is supported on the target machine
 1. Get the intrinsic cost attributes for float and double FMA
 2. Get the instruction cost for the FMA's
 3. Check if the cost is valid
 4. If valid, FMA is supported. If not, FMA not supported

Legality Conditions

- Test if FMA is supported on the target machine

```
15 bool isFMASupported(TargetTransformInfo &TTI, Function &F) {  
16     LLVMContext &Ctx = F.getContext();  
17     ↓  
18     // Check for both float (f32) and double (f64) FMA support  
19     Type *DoubleTy = Type::getDoubleTy(&Ctx);  
20     Type *FloatTy = Type::getFloatTy(&Ctx);  
21     ↓  
22     // We check if the intrinsic is available and has reasonable cost  
23     IntrinsicCostAttributes ICADouble(Id: Intrinsic::fma, RTy: DoubleTy,  
24     ... Tys: {[0]=DoubleTy, [1]=DoubleTy, [2]=DoubleTy});  
25     IntrinsicCostAttributes ICAFloat(Id: Intrinsic::fma, RTy: FloatTy,  
26     ... Tys: {[0]=FloatTy, [1]=FloatTy, [2]=FloatTy});  
27     ↓  
28     InstructionCost CostDouble =  
29     ... TTI.getIntrinsicInstrCost(ICA: ICADouble, CostKind: TargetTransformInfo::TCK_CodeSize);  
30     InstructionCost CostFloat =  
31     ... TTI.getIntrinsicInstrCost(ICA: ICAFloat, CostKind: TargetTransformInfo::TCK_CodeSize);  
32     ↓  
33     // FMA is supported if the intrinsic has a valid, non-infinite cost  
34     // An invalid or very high cost indicates lack of hardware support  
35     bool DoubleSupported = CostDouble.isValid() &&  
36     ... CostDouble != InstructionCost::getInvalid() &&  
37     ... CostDouble.getValue() == 1;  
38     bool FloatSupported = CostFloat.isValid() &&  
39     ... CostFloat != InstructionCost::getInvalid() &&  
40     ... CostFloat.getValue() == 1;  
41     ↓  
42     // FMA is supported if either float or double version is supported  
43     bool Supported = DoubleSupported || FloatSupported;  
44     ↓  
45     return Supported;  
46 }
```

Legality Conditions

Implement a profitability check (only transform if beneficial, e.g., not in precise FP mode)

1. Check if fast-math flags are enabled
2. Check count of uses for FMUL
3. Check if in precise FP mode
4. Default case: don't optimize

Legality Conditions

- Implement a profitability check
 1. Check if fast-math flags are enabled
 2. Check count of uses for FMUL
 3. Check if in precise FP mode
 4. Default case: don't optimize

```
48 bool isProfitable(BinaryOperator *FAdd, BinaryOperator *FMul) {  
49     // Check 1: Fast math flags  
50     // FMA changes the rounding behavior, so we should only apply it if:  
51     // -- Fast math flags allow it (contract or fast)  
52     // -- Or if both instructions already allow reassociation  
53     FastMathFlags FAddFlags = FAdd->getFastMathFlags();  
54     FastMathFlags FMulFlags = FMul->getFastMathFlags();  
55     ↓  
56     // Precise FP mode check  
57     // If we're in precise FP mode (no fast math flags at all), don't transform  
58     if (!FAddFlags.any() && !FMulFlags.any()) {  
59         return false;  
60     }  
61     ↓  
62     // If either operation allows contraction, FMA is profitable  
63     if (FAddFlags.allowContract() || FMulFlags.allowContract()) {  
64         return true;  
65     }  
66     ↓  
67     // If both operations are marked as fast math, FMA is profitable  
68     if (FAddFlags.isFast() && FMulFlags.isFast()) {  
69         return true;  
70     }  
71     ↓  
72     // Check 2: Single use of multiply result  
73     // If the multiply result is only used once (by this add), combining is  
74     // beneficial. This reduces register pressure and eliminates an intermediate  
75     // value  
76     if (FMul->hasOneUse()) {  
77         return true;  
78     }  
79     ↓  
80     // Default: if we have some fast math flags but not contract/fast, be  
81     // conservative and don't transform  
82     return false;  
83 }
```

Transformation Implementation

Transform all candidate ``fadd`` + ``fmul`` into FMA instructions

1. Extract operators and operands from pattern
2. Get appropriate FMA intrinsic type based on pattern
3. Create FMA intrinsic call using extracted operands
4. Preserve flags from original instruction in FMA call
5. Replace the ``fadd`` with the new FMA call
6. Remove dead ``fmul`` code

Transformation Implementation

Transform all candidate `fadd` + `fmul` into FMA instructions

1. Extract operators and operands from pattern

```
86  · Value *X = nullptr, *Y = nullptr, *Z = nullptr;↓
87  · BinaryOperator *FMul = nullptr;↓
88  ↓
89  · // Extract the operands from the pattern↓
90  · // Pattern 1: fadd (fmul X, Y), Z↓
91  · if (match(V: FAdd, P: m_FAdd(L: m_FMul(L: m_Value(&V: X), R: m_Value(&V: Y)), R: m_Value(&V: Z)))) {↓
92  ·   FMul = dyn_cast<BinaryOperator>(Val: FAdd->getOperand(i_nocapture: 0));↓
93  · }↓
94  · // Pattern 2: fadd Z, (fmul X, Y) -- commutative case↓
95  · else if (match(V: FAdd, P: m_FAdd(L: m_Value(&V: Z), R: m_FMul(L: m_Value(&V: X), R: m_Value(&V: Y)))) {↓
96  ·   FMul = dyn_cast<BinaryOperator>(Val: FAdd->getOperand(i_nocapture: 1));↓
97  · }↓
98  ↓
99  · // Failed to extract operands↓
100 · if (!FMul || !X || !Y || !Z) {↓
101 ·   return;↓
102 · }
```

Transformation Implementation

Transform all candidate `fadd` + `fmul` into FMA instructions

1. Extract operators and operands from pattern
2. Get appropriate FMA intrinsic type based on pattern

```
104  // Get the type of the operation (float or double) ↓  
105  Type *Ty = FAdd->getType(); ↓  
106  ↓
```

Transformation Implementation

Transform all candidate `fadd` + `fmul` into FMA instructions

1. Extract operators and operands from pattern
2. Get appropriate FMA intrinsic type based on pattern
3. Create FMA intrinsic call using extracted operands

```
110  · // Get the FMA intrinsic for the appropriate type ↓
111  · Module *M = FAdd->getModule(); ↓
112  · Function *FMAIntrinsic = ↓
113  ·   ··· Intrinsic::getOrInsertDeclaration(M, id: Intrinsic::fma, Tys: Ty); ↓
114  ↓
115  · // Create the FMA intrinsic call: fma(X, Y, Z) computes X * Y + Z ↓
116  · Value *FMACall = Builder.CreateCall(Callee: FMAIntrinsic, Args: {[0]=X, [1]=Y, [2]=Z}); ↓
```

Transformation Implementation

Transform all candidate `fadd` + `fmul` into FMA instructions

1. Extract operators and operands from pattern
2. Get appropriate FMA intrinsic type based on pattern
3. Create FMA intrinsic call using extracted operands
4. Preserve flags from original instruction in FMA call

```
118  ··// Preserve fast-math flags from the original instructions↓
119  ··if (auto *FMAInst: Instruction * = dyn_cast<Instruction>(Val: FMACall)) {↓
120  ···FastMathFlags FMF = FAdd->getFastMathFlags();↓
121  ···FMF |= FMul->getFastMathFlags(); // Combine flags from both operations↓
122  ···FMAInst->setFastMathFlags(FMF);↓
123  ··}↓
```

Transformation Implementation

Transform all candidate `fadd` + `fmul` into FMA instructions

1. Extract operators and operands from pattern
2. Get appropriate FMA intrinsic type based on pattern
3. Create FMA intrinsic call using extracted operands
4. Preserve flags from original instruction in FMA call
5. Replace the `fadd` with the new FMA call

```
125  ..// Replace all uses of the FAdd with the FMA call↓  
126  .. FAdd->replaceAllUsesWith(V: FMACall);↓
```

Transformation Implementation

Transform all candidate `fadd` + `fmul` into FMA instructions

1. Extract operators and operands from pattern
2. Get appropriate FMA intrinsic type based on pattern
3. Create FMA intrinsic call using extracted operands
4. Preserve flags from original instruction in FMA call
5. Replace the `fadd` with the new FMA call
6. Remove dead code

```
128     · · · · · // If the FMul has no other uses, we can remove it ↓  
129     · · if (FMul->hasOneUse()) ↓  
130         · · · FMul->eraseFromParent(); ↓  
131     ↓  
132     · · · · · // Remove the original FAdd instruction ↓  
133     · · FAdd->eraseFromParent(); ↓  
134     } ↓
```

Testing

- Compiled with `-O0`

```
1 double test_basic_fma(double x, double y, double z) {  
2   double mul_result = x * y;  
3   return mul_result + z;  
4 }
```

```
1 ; Function Attrs: noinline nounwind ssp uwtable(sync)  
2 define double @test_basic_fma(double noundef %0, double noundef %1, double noundef %2)  
3   %4 = fmul double %0, %1  
4   %5 = fadd double %4, %2  
5   ret double %5  
6 }
```

```
1 ; Function Attrs: noinline nounwind ssp uwtable(sync)  
2 define double @test_basic_fma(double noundef %0, double noundef %1, double noundef %2)  
3   %4 = call double @llvm.fma.f64(double %0, double %1, double %2)  
4   ret double %4  
5 }
```

clang -emit-llvm
opt -passes=mem2reg

opt -passes=fma

Testing

- Compiled with `-O0`

```
1 float test_basic_fma_float(float x, float y, float z) {  
2   float mul_result = x * y;  
3   return mul_result + z;  
4 }
```

```
1 ; Function Attrs: noline nounwind ssp uwtable(sync)  
2 define float @test_basic_fma_float(float noundef %0, float noundef %1, float noundef %2)  
3   %4 = fmul float %0, %1  
4   %5 = fadd float %4, %2  
5   ret float %5  
6 }
```

```
1 ; Function Attrs: noline nounwind ssp uwtable(sync)  
2 define float @test_basic_fma_float(float noundef %0, float noundef %1, float noundef %2)  
3   %4 = call float @llvm.fma.f32(float %0, float %1, float %2)  
4   ret float %4  
5 }
```

clang -emit-llvm
opt -passes=mem2reg

opt -passes=fma

Next Steps

- Properly benchmark to see if there's performance improvements
- Test on different architectures (only have tested on my mac)
- Could add support for different FMA patterns ($-x * y + z$ OR $x * y - z$)