# Enzyme: High-Performance Automatic Differentiation of LLVM

William S. Moses (wmoses@mit.edu), Valentin Churavy (vchuravy@mit.edu)

MIT CSAIL

## Existing Automatic Differentiation Tools

- **Differentiable DSLs** (TensorFlow, PyTorch, DiffTaichi) provide a new language where everything is differentiable. Must rewrite code in DSL.
- **Operator Overloading** (Adept, JAX) tools provide differentiable versions of existing language constructs (double => adouble, np.sum => jax.sum). May require rewriting to use non-standard utilities for support.
- **Source Rewriting** tools statically analyze code to produce a new gradient function in the source language. Requires all code available ahead of time and is difficult to use with external libraries.

## Optimization and AD

All tools for existing code operate at the source level preventing optimizations before AD without reimplementing compiler analyzes and optimization into the AD tool. While historically not considered necessary, we demonstrate in Figure 1 how crucial optimization prior to AD can be.

```
float mag(const float*); //Compute magnitude in O(N)
void norm(float* out, const float* in){
    // float res = mag(in); LICM moves mag outside loop
    for(int i = 0; i < N; i++) { out[i] = in[i] / mag(in); }
}
```

```
// LICM, then AD, O(N)              // AD then LICM, O(N^2)
float res = mag(in);               float res = mag(in);
for(int i = 0; i < N; i++) {       for(int i = 0; i < N; i++) {
  out[i] = in[i] / res;              out[i] = in[i] / res;
}                                  }
float d_res = 0;
for (int i = 0; i < N; i++) {      for (int i = 0; i < N; i++) {
  d_res += -in[i] * in[i]            float d_res = -in[i] * in[i]
         * d_out[i]/res;                     * d_out[i]/res;
  d_in[i] += d_out[i]/res;           d_in[i] += d_out[i]/res;
                                     ∇mag(in, d_in, d_res);
}                                  }
∇mag(in, d_in, d_res);             //
```

**Figure 1.** When differentiating `norm`, running LICM prior to AD is asymptotically faster than running AD followed by LICM.

## Challenges of AD on Low-Level Code

Performing AD on low-level IR presents additional challenges as source-level information is lost. For example, differentiating memcpy in Figure 2 requires the underlying type of the data to select the correct gradient.

```
void f(void* dst, void* src) { memcpy(dst, src, 8); }
```

```
// Assume double inputs            // Assume float inputs
∇f(double* dst, double* ddst,      ∇f(float* dst, float* ddst,
   double* src, double* dsrc) {       float* src, float* dsrc) {
  // Forward pass                     // Forward pass
  memcpy(dst, src, 8);               memcpy(dst, src, 8);
  // Reverse pass                     // Reverse pass
  dsrc[0] += ddst[0];                dsrc[0] += ddst[0];
  ddst[0] = 0;                       ddst[0] = 0;
                                     dsrc[1] += ddst[1];
                                     ddst[1] = 0;
}                                  }
```

**Figure 2. Top:** Call to `memcpy` for an unknown 8-byte object. **Left:** Gradient for a `memcpy` of 8 bytes of double data. **Right:** Gradient for a `memcpy` of 8 bytes of float data.

## Design

Enzyme synthesizes derivatives by:
- Running Type and Activity Analysis
- Allocating shadows of active variables
- Creating a "reverse" copy of BasicBlock's in the original code that compute the adjoints of its instructions in reverse order.

```
define @relu3(double %x)          rev_end:
entry:                              ; adjoint of return
  ; Shadows for reverse             store %d_res = 1.0
  ; alloca %d_x = 0.0               ; adjoint of %res phi node
  ; alloca %d_call = 0.0            %cmp2 = load %cmp_cache
  ; alloca %d_result = 0.0          %tmp = load %d_res
  ; Cache of %cmp                   %d_call += if %cmp2, %tmp else 0
  ; alloca %cmp_cache               store %d_res = 0.0
  %cmp = %x > 0                     br %cmp, %rev_iftrue, %rev_entry
  br %cmp, %iftrue, %end          rev_iftrue:
iftrue:                             ; adjoint of %call
  %call = @pow(%x, 3)               %df = 3 * @pow(%x, 2)
  br cond.end                       %d_x += %df * (load %d_call)
end:                                store %d_call = 0.0
  %res = ϕ[%call, if.true],         br %rev_entry
    ↪  [0, entry]                 rev_entry:
  ret %res                          %0 = load %d_x
                                    ret %0
```

**Figure 3.** Gradient synthesis of `relu(pow(x,3))`. **Left:** the original computation with comments showing the shadow allocations of active variables that would be added to the forward pass. **Right:** reverse pass generated by Enzyme. The full synthesized gradient function would combine these (with shadow allocations added), replacing the return with a branch to the reverse pass. **Bottom:** demonstration of how to call Enzyme.

- **Type Analysis**: A new interprocedural analysis that derives the underlying types of data by manipulating "TypeTree's" of LLVM Values. TypeTrees are initialized with Constant, TBAA info. Each instruction is given a type propagation rule that until fixpoint. We provide compile-time error if a necessary type cannot be deduced statically.
- **Activity Analysis** determines what instructions could impact derivative computation to avoid computing unnecessary adjoints. Build off Type & Alias Analysis to get better results. E.g. all read-only function that returns an integer are inactive since they cannot propagate adjoints through the return or to any memory location.
- **Shadow Memory** is used to store the derivatives of values. Shadow versions of data structures created inside the differentiated function are created automatically. Data structures passed as arguments to the differentiated function must also have shadow arguments passed.
- **Cache**: Some adjoint instructions require values from the forward pass (e.g. $\nabla(x * y) \longrightarrow x * dy + y * dx$). Memory is automatically allocated for all such values. Enzyme optimizes the cache by recomputing instead of caching and avoids caching unnecessary or equivalent values, whenever possible.
- **Extensibility:** Custom gradients are supported by attaching metadata specifying the corresponding gradient function. Multisource AD is supporting by leveraging LLVM's LTO support and libraries with embedded bitcode. This ensures that any potentially active call has a definition available to differentiate.

## Evaluation

Performing AD after optimization yields a 4.5× speedup over AD before optimization. This accounts for much, but not all, of Enzyme's improvement over prior art (different cache and activity analysis implementations).
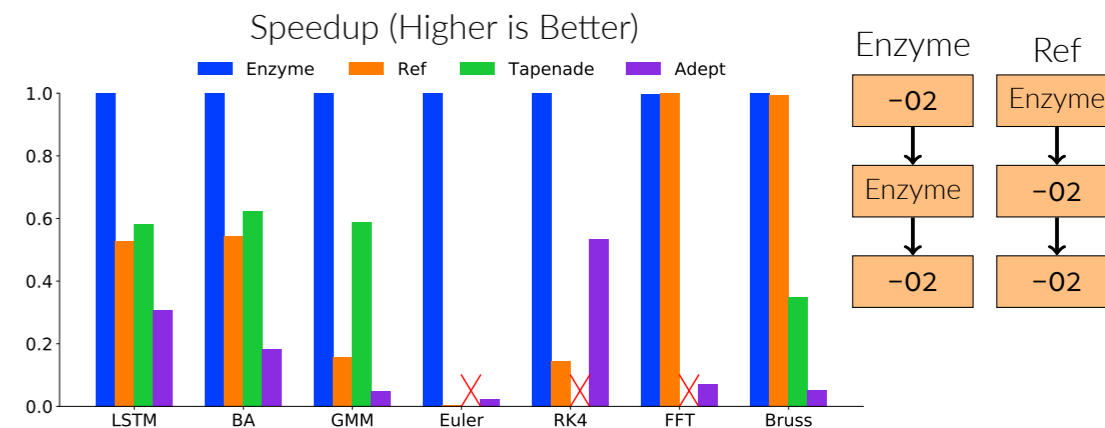


**Figure 4.** Relative speedup of AD systems on ADBench+ benchmarks, higher is better. A red X denotes programs that an AD system does not produce a correct gradient. A value of 1.0 denotes the fastest system, whereas 0.5 denotes taking twice as long.

## Usage

A user can use gradient functions by calling `__enzyme_autodiff` with the function to be differentiated as the first argument. When the Enzyme optimization pass is run, it will replace any calls to `__enzyme_autodiff` with a call a newly-generated gradient function.

```
%grad = call double @__enzyme_autodiff(@relu3, double %x)
```
```
               opt -load=/path/to/LLVMEnzyme.so -enzyme
```
```
%grad = call double @grad_relu3(double %x)
```

**Figure 5.** Convention for invoking Enzyme.

Enzyme is built as an LLVM compiler plugin for versions 7 and later ease incorporation into an existing tools. We have demonstrated taking derivatives of C/C++ via Clang, PyTorch, and Tensorflow. We've also demonstrated dynamic language support by using Enzyme to differentiate Julia.

For more information about installing and using Enzyme, please visit `https://enzyme.mit.edu` and come to our student research competition presentation!

## Acknowledgements