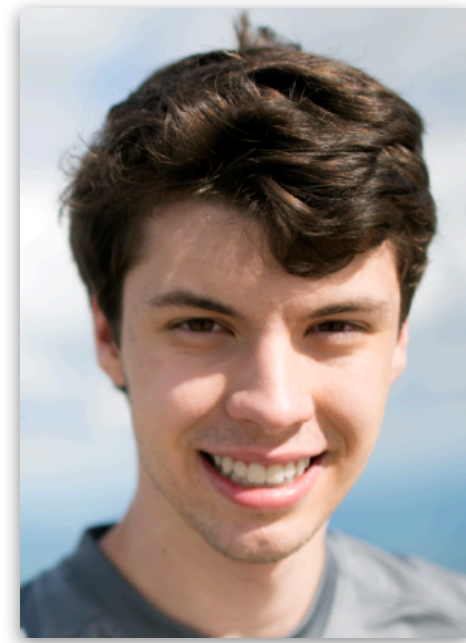


Enzyme: Fast and Effective Automatic Differentiation for Academia and Industry



William S. Moses

wmoses@mit.edu

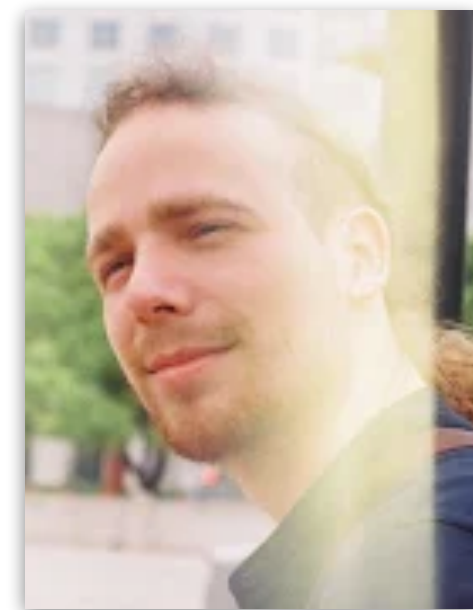
ICIAM

Aug 22, 2023



UNIVERSITY OF
ILLINOIS
URBANA-CHAMPAIGN





William S. Moses Valentin Churavy Ludger Paehler Johannes Doerfert Alex Zinenko



Jan Hückelheim

Sri Hari Krishna
Narayanan

Michel Schanen

Paul Hovland



Leila Ghaffari

Praytush Das

Tim Gymnich

Manuel Drehwald

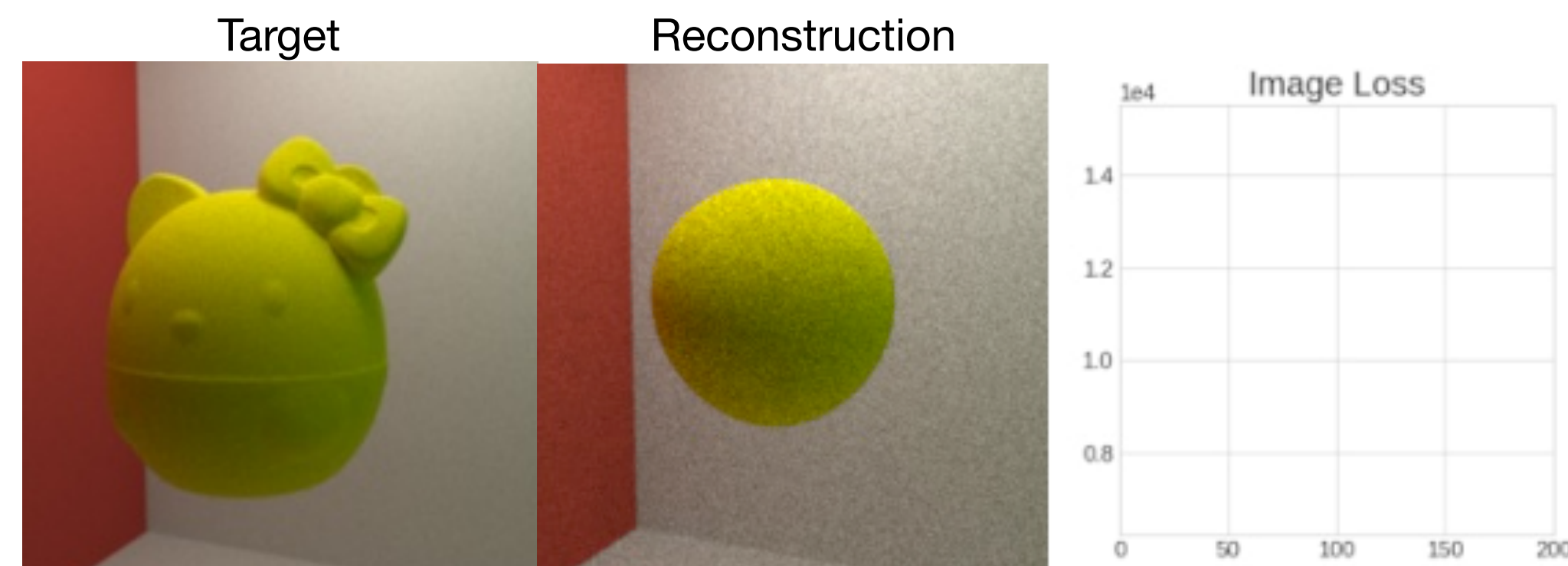
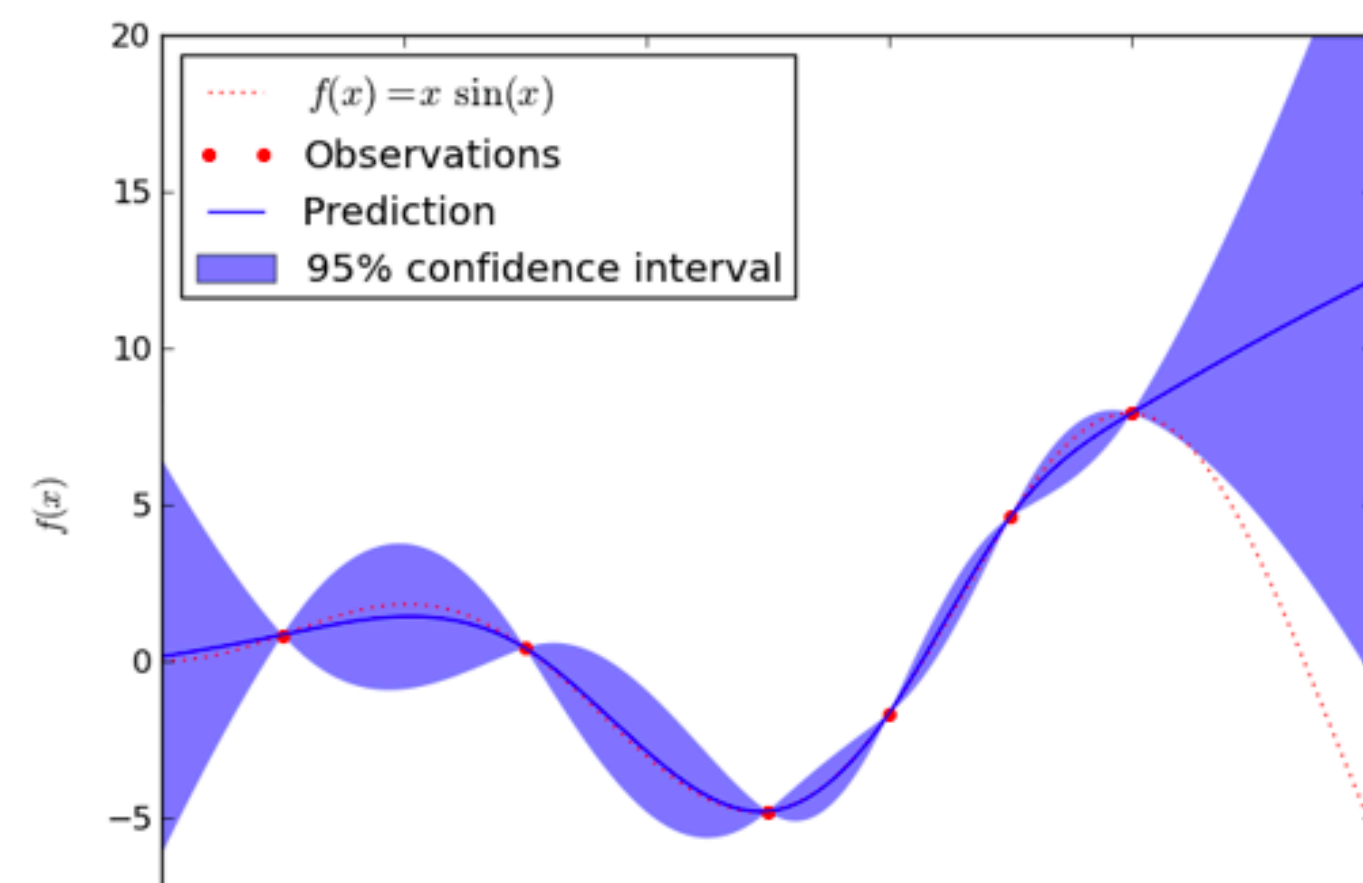
&
more

AP Calculus: Revisited

- Derivatives compute the rate of change of a function's output with respect to input(s)

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h}$$

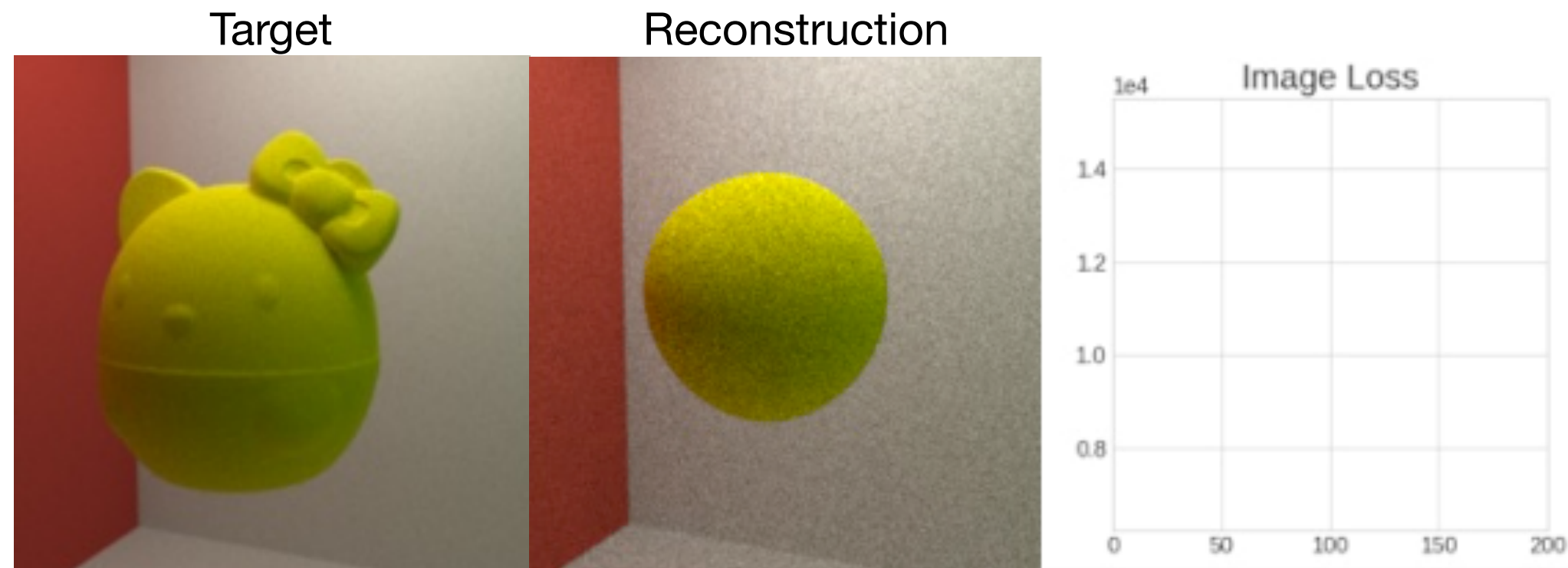
- Derivatives are used widely across science
 - Machine learning (back-propagation, Bayesian inference)
 - Scientific computing (modeling, simulation, uncertainty quantification)



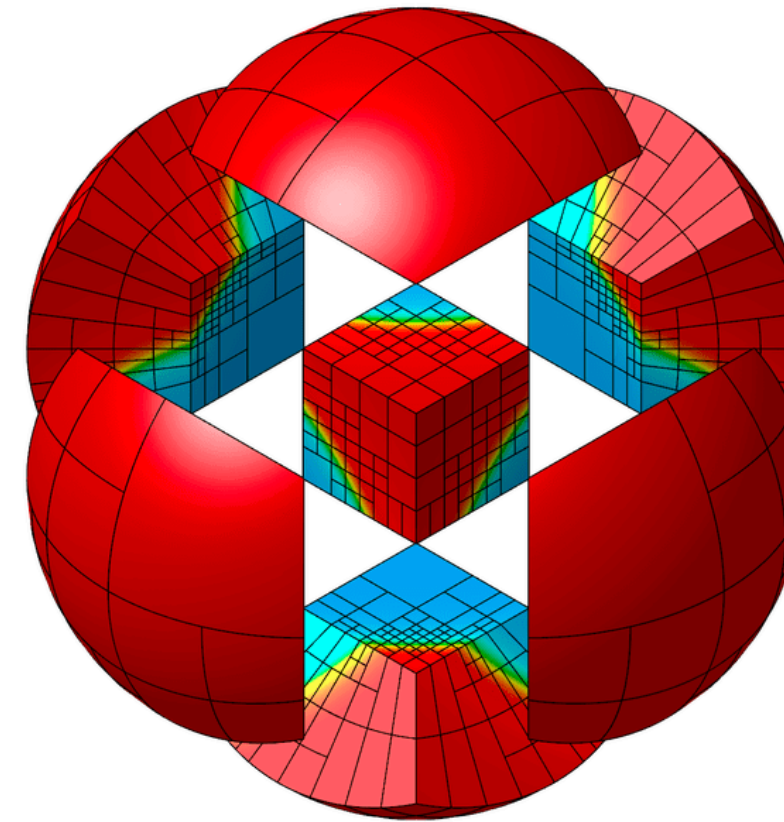
from Efficient Differentiation of Pixel Reconstruction Filters for Path-Space Differentiable Rendering, SIGGRAPH Asia 2022, Zihan Yu et al



$\partial_{\partial x}$ AD-Powered Applications



from [Efficient Differentiation of Pixel Reconstruction Filters for Path-Space Differentiable Rendering](#), SIGGRAPH Asia 2022, Zihan Yu et al



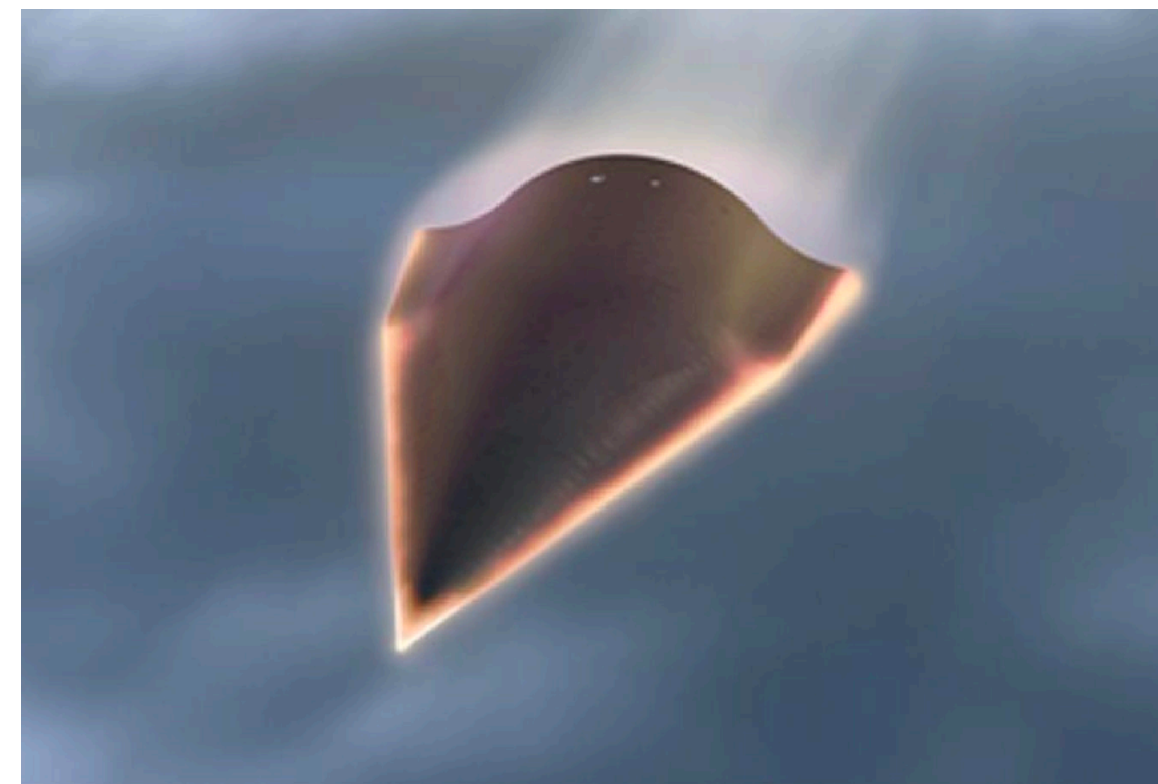
from [MFEM Team at LLNL](#)



from [Comrade: High Performance Black-Hole Imaging](#) JuliaCon 2022, Paul Tiede (Harvard)



from [CLIMA & NSF CSSI: Differentiable programming in Julia for Earth system modeling \(DJ4Earth\)](#)



from [Center for the Exascale Simulation of Materials in Extreme Environments](#)

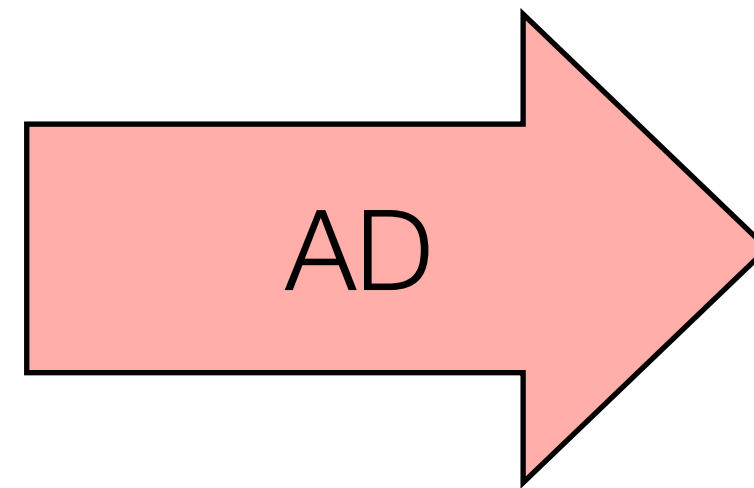


from [Differential Molecular Simulation with Molly.jl](#), EnzymeCon 2023, Joe Greener (Cambridge)

Automatic Derivative Generation

- Derivatives can be generated automatically from definitions within programs

```
double relu3(double x) {  
    if (x > 0)  
        return pow(x,3)  
    else  
        return 0;  
}
```



```
double grad_relu3(double x) {  
    if (x > 0)  
        return 3 * pow(x,2)  
    else  
        return 0;  
}
```

- Unlike numerical approaches, automatic differentiation (AD) can compute the derivative of ALL inputs (or outputs) at once, without approximation error!

```
// Numeric differentiation  
// f'(x) approx [f(x+epsilon) - f(x)] / epsilon  
double grad_input[100];  
  
for (int i=0; i<100; i++) {  
    double input2[100] = input;  
    input2[i] += 0.01;  
    grad_input[i] = (f(input2) - f(input))/0.001;  
}
```

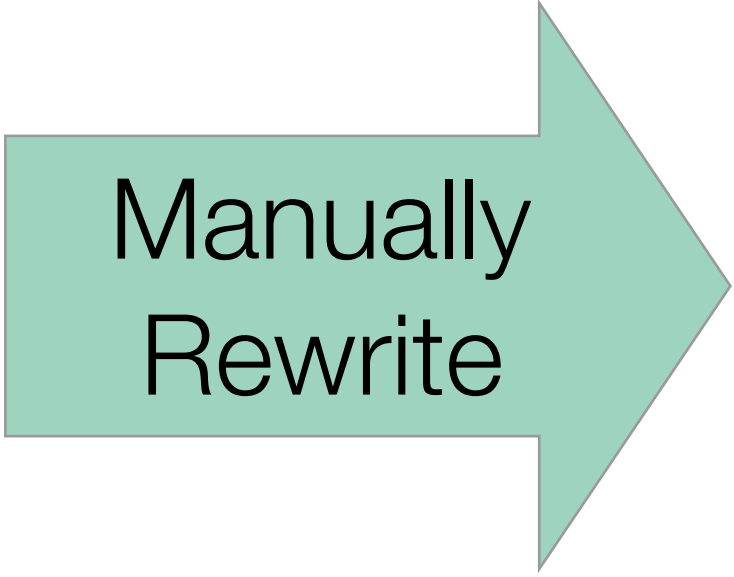
```
// Automatic differentiation  
double grad_input[100];  
  
grad_f(input, grad_input)
```

Existing AD Approaches (1/3)

- Differentiable DSL (TensorFlow, PyTorch, DiffTaichi)
 - Provide a new language designed to be differentiated
 - Requires rewriting everything in the DSL and the DSL must support all operations in original code
 - Fast if DSL matches original code well

```
double relu3(double val) {  
    if (x > 0)  
        return pow(x, 3)  
    else  
        return 0;  
}
```

Manually
Rewrite



```
import tensorflow as tf  
  
x = tf.Variable(3.14)  
  
with tf.GradientTape() as tape:  
    out = tf.cond(x > 0,  
                  lambda: tf.math.pow(x, 3),  
                  lambda: 0  
                )  
print(tape.gradient(out, x).numpy())
```

Existing AD Approaches (2/3)

- Operator overloading (Adept, JAX)
 - Differentiable versions of existing language constructs (double => adouble, np.sum => jax.sum)
 - May require writing to use non-standard utilities
 - Often dynamic: storing instructions/values to later be interpreted

```
// Rewrite to accept either
// double or adouble
template<typename T>
T relu3(T val) {
    if (x > 0)
        return pow(x, 3)
    else
        return 0;
}
```

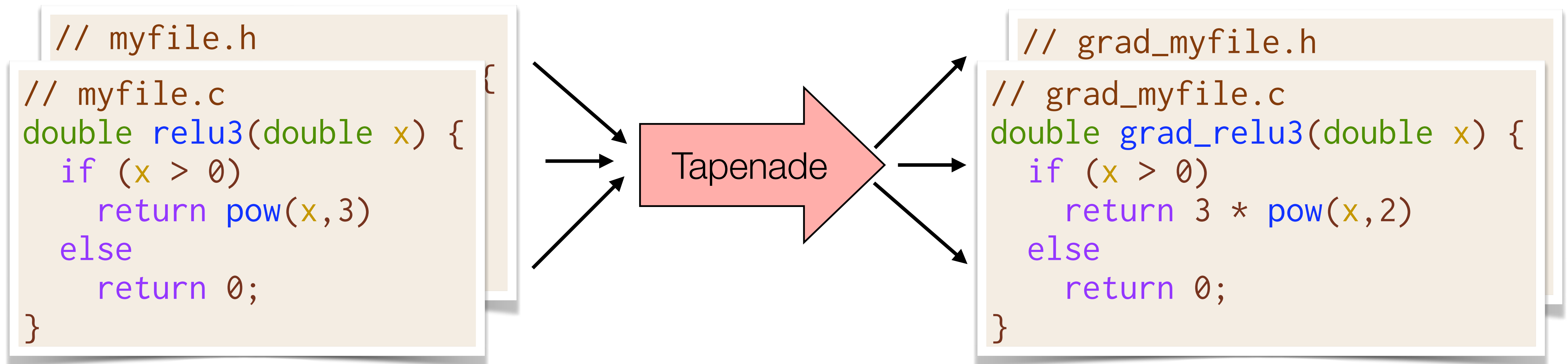
```
adept::Stack stack;
adept::adouble inp = 3.14;

// Store all instructions into stack
adept::adouble out(relu3(inp));
out.set_gradient(1.00);

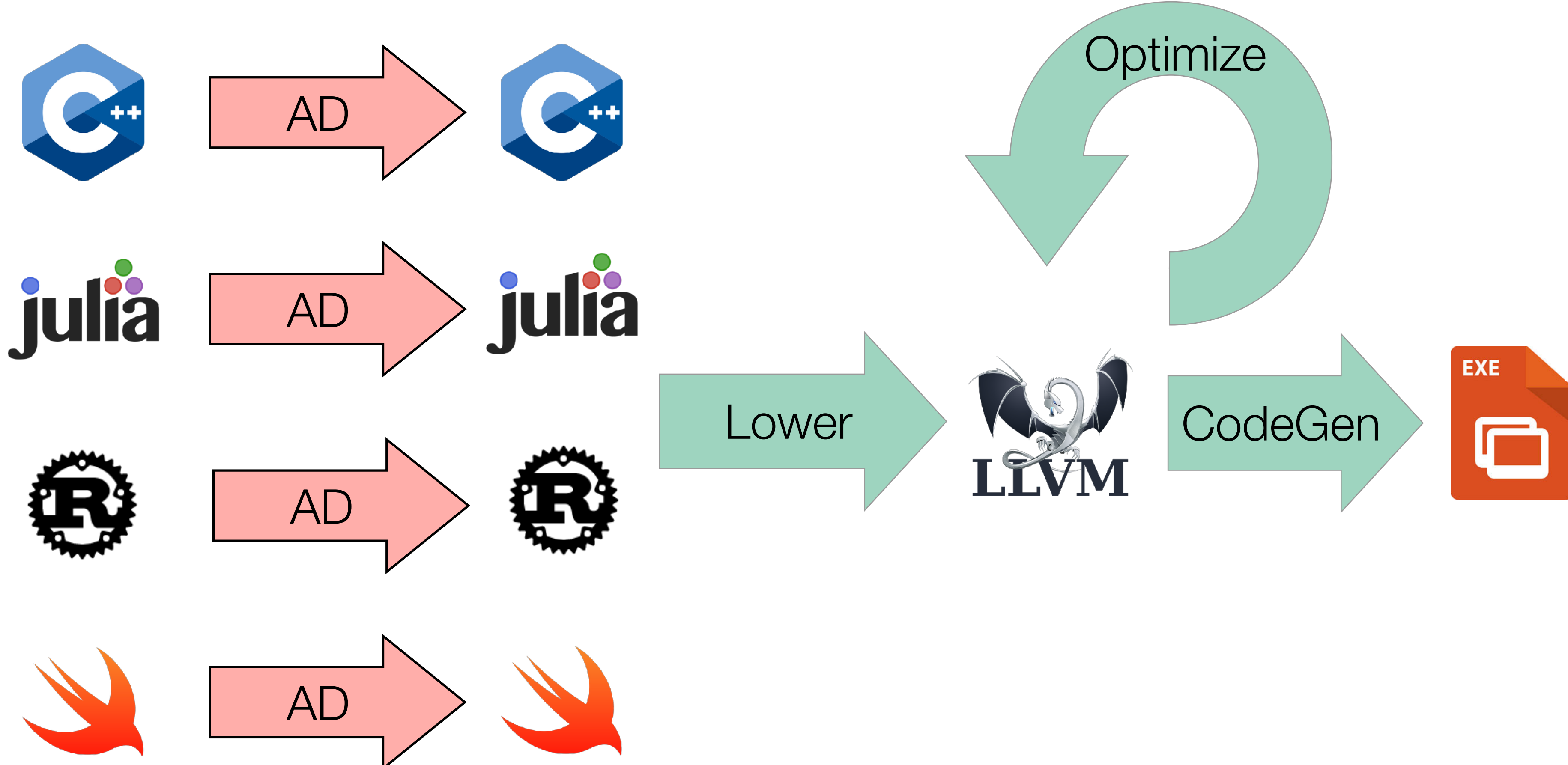
// Interpret all stack instructions
double res = inp.get_gradient(3.14);
```

Existing AD Approaches (3/3)

- Source rewriting
 - Statically analyze program to produce a new gradient function in the source language
 - Re-implement parsing and semantics of given language
 - Requires all code to be available ahead of time => hard to use with external libraries



Existing Automatic Differentiation Pipelines



Case Study: Vector Normalization

```
//Compute magnitude in O(n)
double mag(double[] x);

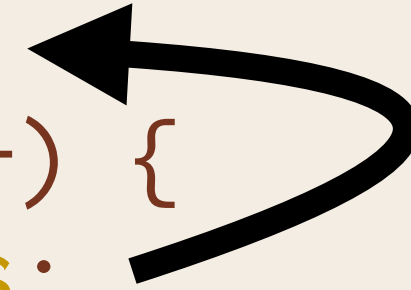
//Compute norm in O(n^2)
void norm(double[] out, double[] in) {

    for (int i=0; i<n; i++) {
        out[i] = in[i] / mag(in);
    }
}
```

Case Study: Vector Normalization

```
//Compute magnitude in O(n)
double mag(double[] x);

//Compute norm in O(n)
void norm(double[] out, double[] in) {
    double res = mag(in);
    for (int i=0; i<n; i++) {
        out[i] = in[i] / res;
    }
}
```



Optimization & Automatic Differentiation

$O(n^2)$

```
for i=0..n {  
  out[i] /= mag(in)  
}
```

Optimize

$O(n)$

```
res = mag(in)  
for i=0..n {  
  out[i] /= res  
}
```

AD

$O(n)$

```
d_res = 0.0  
for i=n..0 {  
  d_res += d_out[i]...  
}  
∇mag(d_in, d_res)
```

Optimization & Automatic Differentiation

$O(n^2)$

```
for i=0..n {  
  out[i] /= mag(in)  
}
```

Optimize

$O(n)$

```
res = mag(in)  
for i=0..n {  
  out[i] /= res  
}
```

AD

$O(n)$

```
d_res = 0.0  
for i=n..0 {  
  d_res += d_out[i]...  
}  
∇mag(d_in, d_res)
```

$O(n^2)$

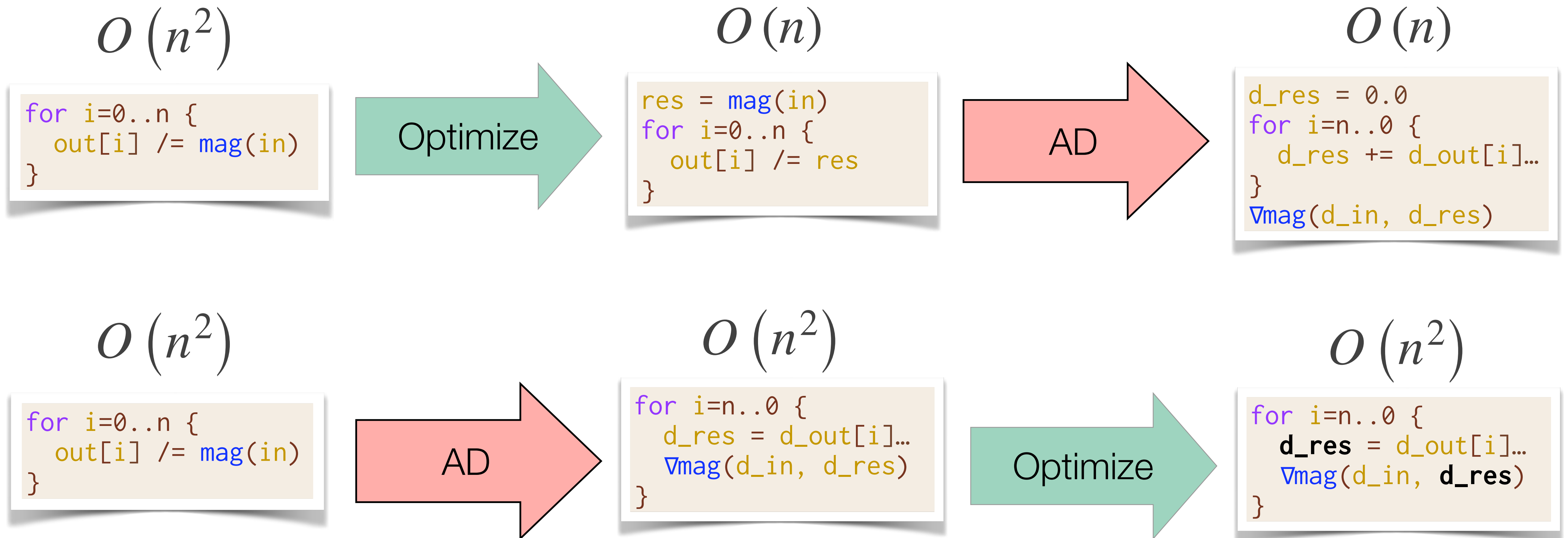
```
for i=0..n {  
  out[i] /= mag(in)  
}
```

AD

$O(n^2)$

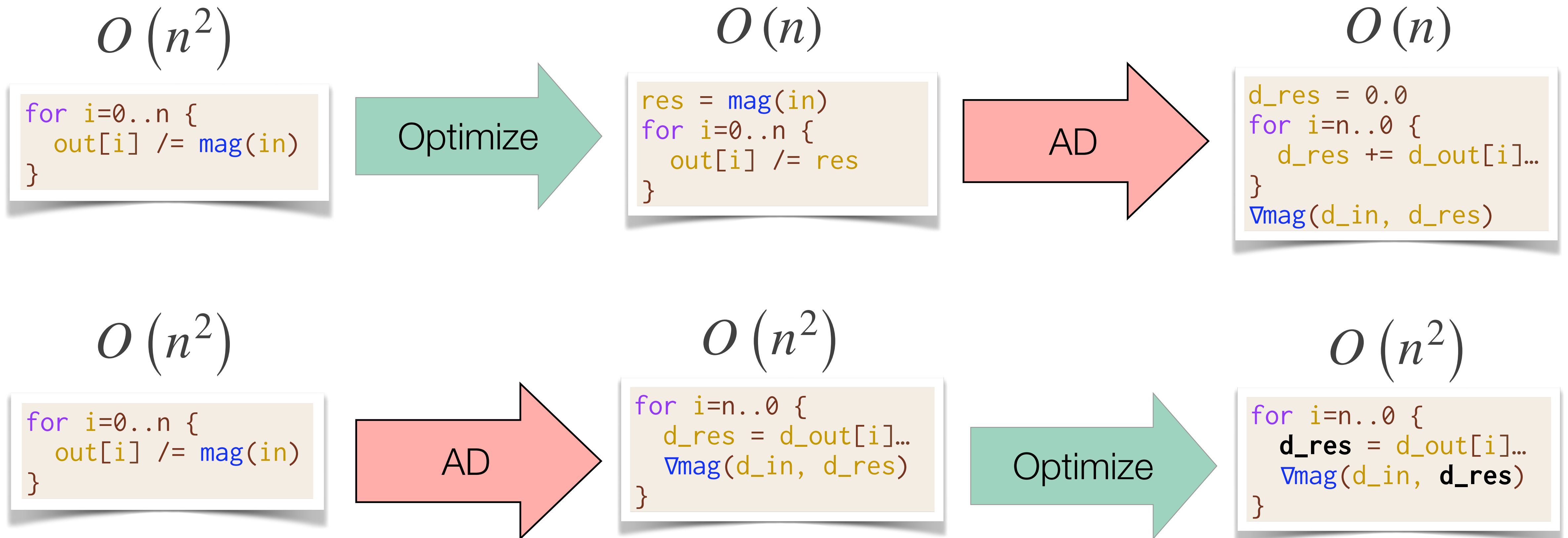
```
for i=n..0 {  
  d_res = d_out[i]...  
  ∇mag(d_in, d_res)  
}
```

Optimization & Automatic Differentiation



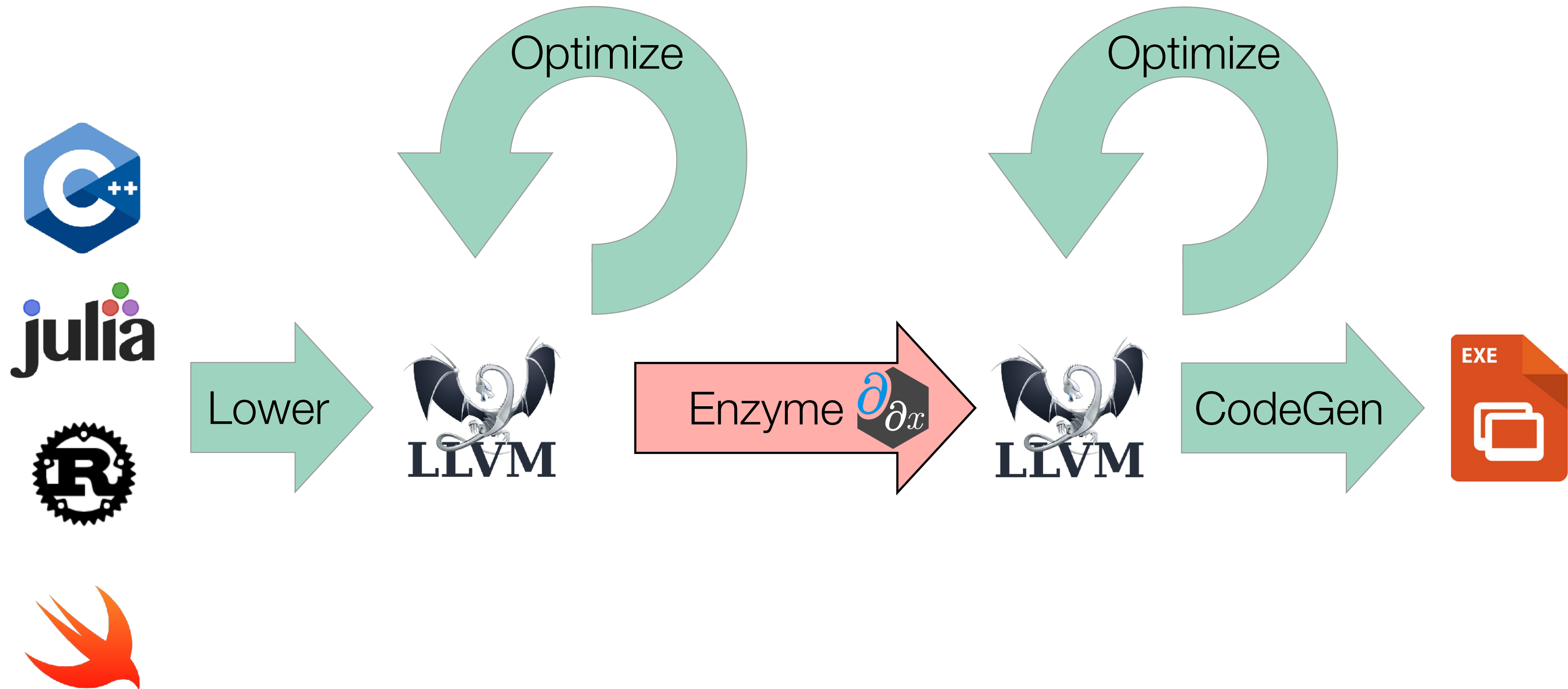
Optimization & Automatic Differentiation

Differentiating after optimization can create *asymptotically faster* gradients!



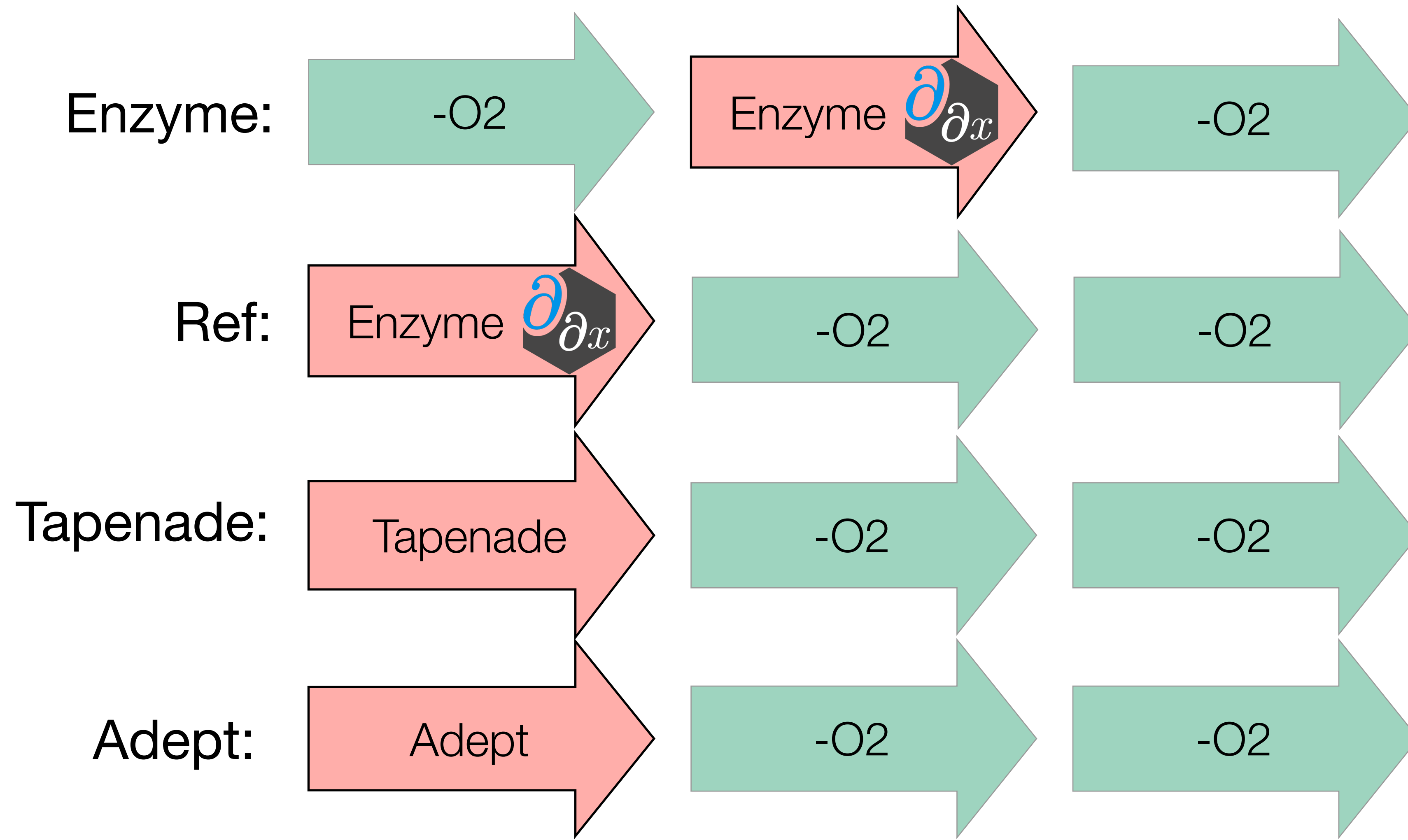
Enzyme Approach

Performing AD at low-level lets us work on ***optimized*** code!

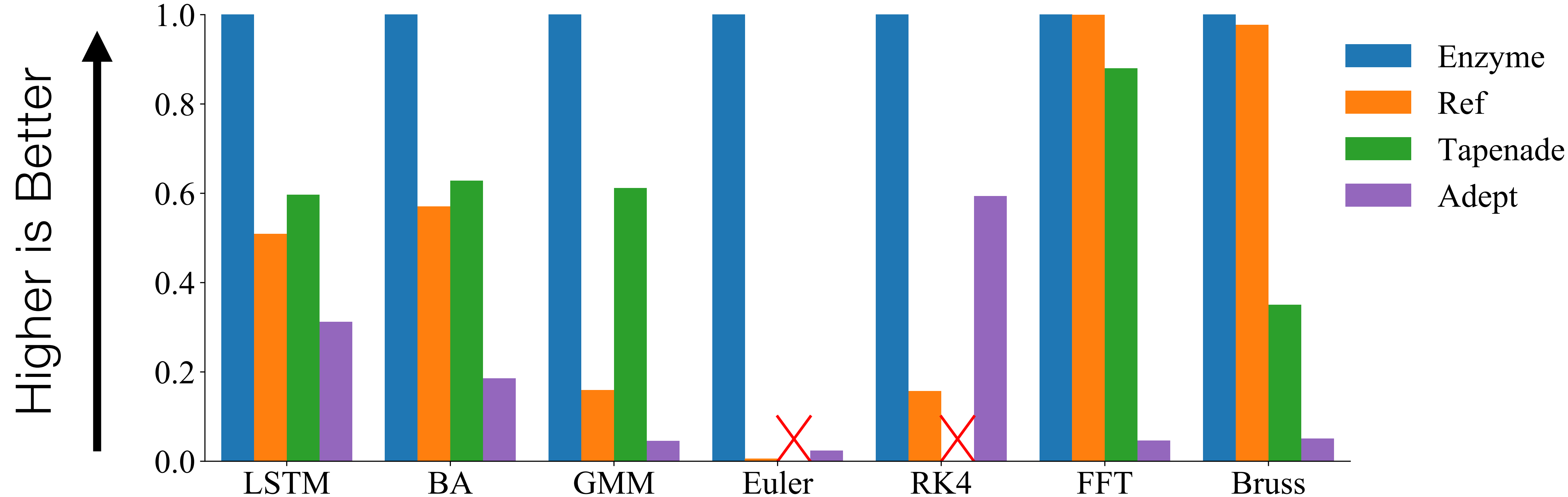


Experimental Setup

- Collection of benchmarks from Microsoft's ADBench suite and of technical interest



Speedup of Enzyme



Enzyme is **4.2x faster** than Reference!

Automatic Differentiation & GPUs

- Prior work has not explored reverse mode AD of existing GPU kernels
 1. Reversing parallel control flow can lead to incorrect results
 2. Complex performance characteristics make it difficult to synthesize efficient code
 3. Resource limitations can prevent kernels from running at all



Efficient GPU Code

- For correctness, Enzyme may need to cache values in order to compute the gradient
 - The complexity of GPU memory means large caches slow down the program by several orders of magnitude, if it even fits at all
- Like the CPU, existing optimizations reduce the overhead
- Unlike the CPU, existing optimizations aren't sufficient
- Novel GPU and AD-specific optimizations can speedup by several orders of magnitude

```
// Forward Pass
out[i] = x[i] * x[i];
x[i] = 0.0f;

// Reverse (gradient) Pass
...
grad_x[i] += 2 * x[i] * grad_out[i];
...
```

Efficient Correct GPU Code

- For correctness, Enzyme may need to cache values in order to compute the gradient
 - The complexity of GPU memory means large caches slow down the program by several orders of magnitude, if it even fits at all
- Like the CPU, existing optimizations reduce the overhead
- Unlike the CPU, existing optimizations aren't sufficient
- Novel GPU and AD-specific optimizations can speedup by several orders of magnitude

```
double* x_cache = new double[...];

// Forward Pass

out[i] = x[i] * x[i];
x_cache[i] = x[i];

x[i] = 0.0f;

// Reverse (gradient) Pass

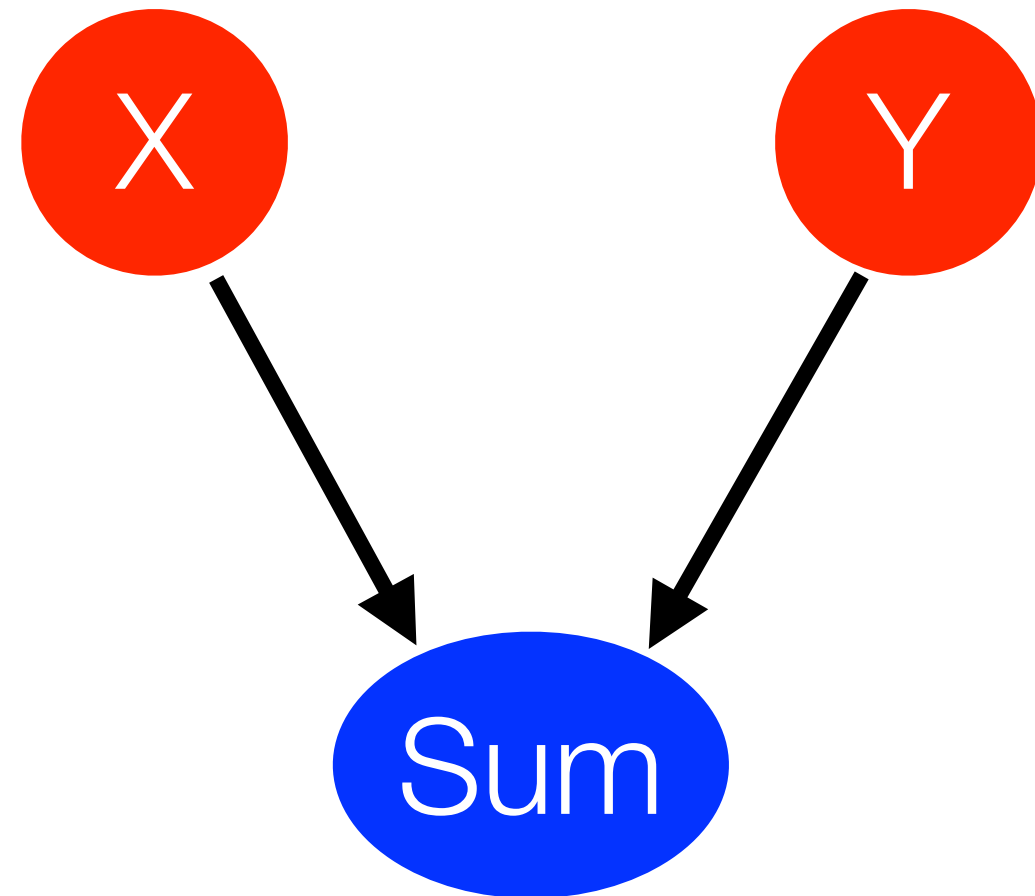
...
grad_x[i] += 2 * x_cache[i]
             * grad_out[i];
...

delete[] x_cache;
```

Cache Reduction Example

- By considering the dataflow graph we can perform a min-cut to approximate smaller cache sizes.

Overwritten:

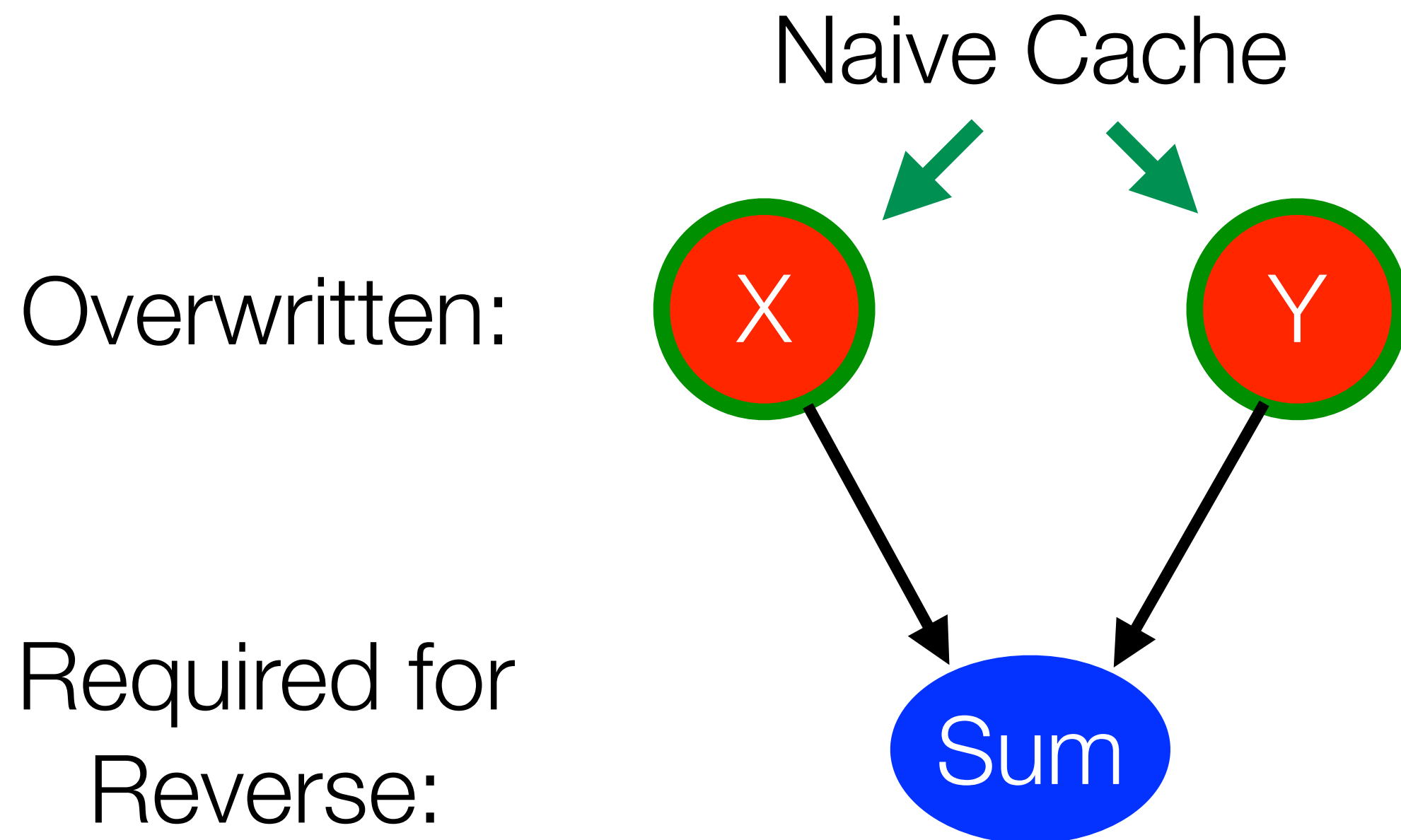


Required for
Reverse:

```
for(int i=0; i<10; i++) {  
    double sum = x[i] + y[i];  
  
    use(sum);  
}  
  
overwrite(x, y);  
grad_overwrite(x, y);  
  
for(int i=9; i>=0; i--) {  
    ...  
    grad_use(sum);  
}
```

Cache Reduction Example

- By considering the dataflow graph we can perform a min-cut to approximate smaller cache sizes.



```
double* x_cache = new double[10];
double* y_cache = new double[10];

for(int i=0; i<10; i++) {
    double sum = x[i] + y[i];
    x_cache[i] = x[i];
    y_cache[i] = y[i];
    use(sum);
}

overwrite(x, y);
grad_overwrite(x, y);

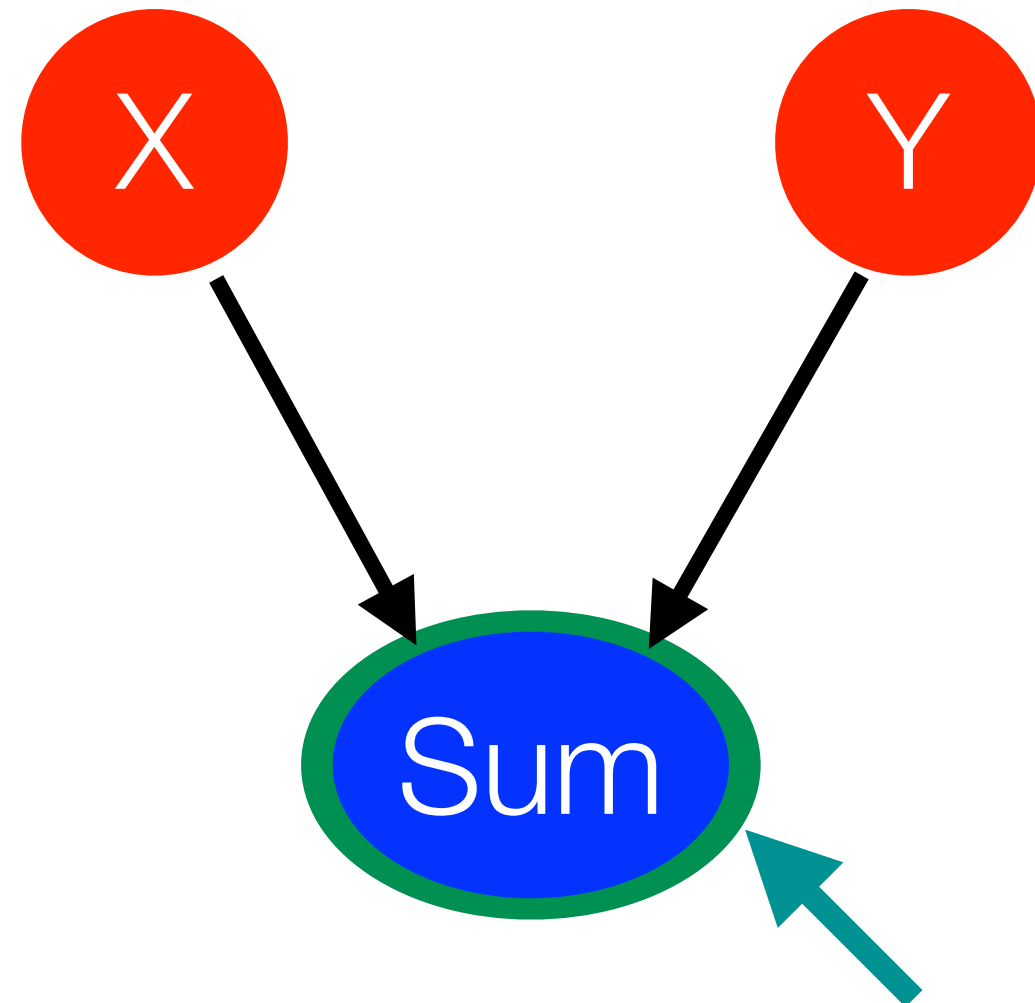
for(int i=9; i>=0; i--) {
    double sum = x_cache[i] + y_cache[i];
    grad_use(sum);
}
```

Cache Reduction Example

- By considering the dataflow graph we can perform a min-cut to approximate smaller cache sizes.

```
double* sum_cache = new double[10];  
  
for(int i=0; i<10; i++) {  
    double sum = x[i] + y[i];  
    sum_cache[i] = sum;  
  
    use(sum);  
}  
  
overwrite(x, y);  
grad_overwrite(x, y);  
  
for(int i=9; i>=0; i--) {  
    grad_use(sum_cache[i]);  
}
```

Overwritten:



Required for
Reverse:

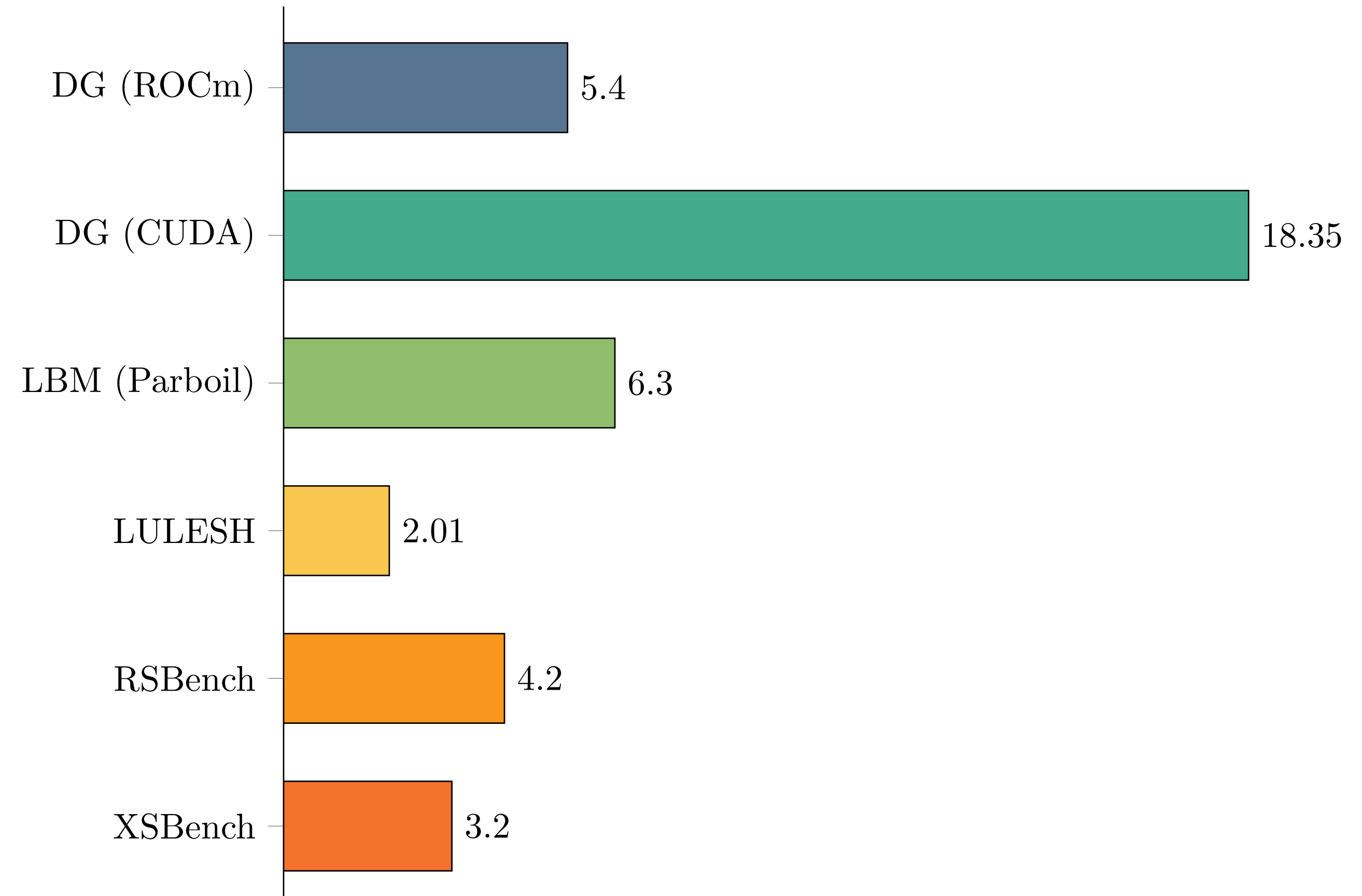
Smallest Cache

Novel AD + GPU Optimizations

- See our SC'21 paper for more (<https://c.wsmoses.com/papers/EnzymeGPU.pdf>)
Reverse-Mode Automatic Differentiation and Optimization of GPU Kernels via Enzyme. SC, 2021
- [AD] Cache LICM/CSE
- [AD] Min-Cut Cache Reduction
- [AD] Cache Forwarding
- [GPU] Merge Allocations
- [GPU] Heap-to-stack (and register)
- [GPU] Alias Analysis Properties of SyncThreads
- ...

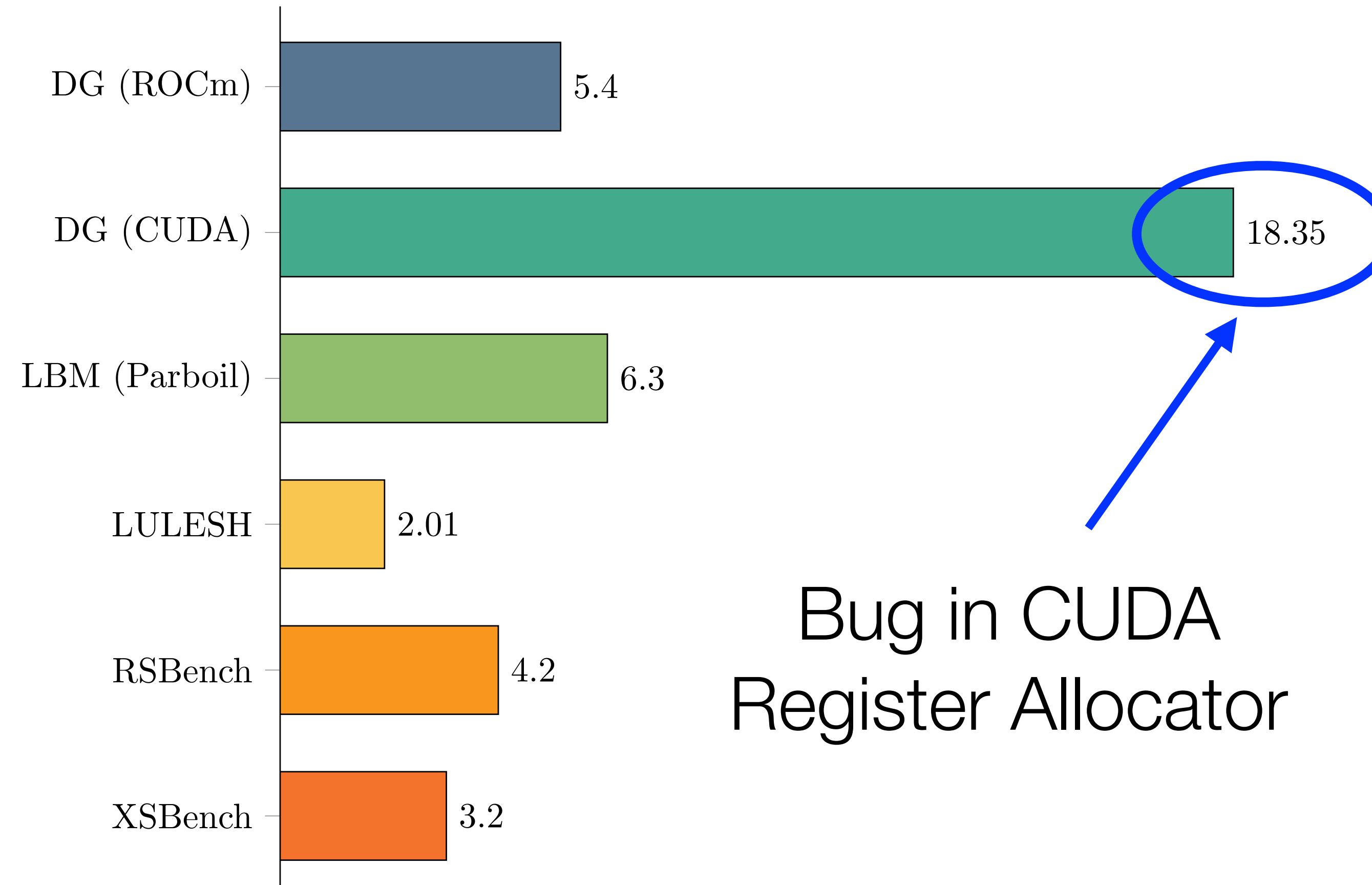
GPU Gradient Overhead

- Evaluation of both original code and gradient
 - DG: Discontinuous-Galerkin integral (Julia)
 - LBM: particle-based fluid dynamics simulation
 - LULESH: unstructured explicit shock hydrodynamics solver
 - XSBench & RSBench: Monte Carlo simulations of particle transport algorithms (memory & compute bound, respectively)



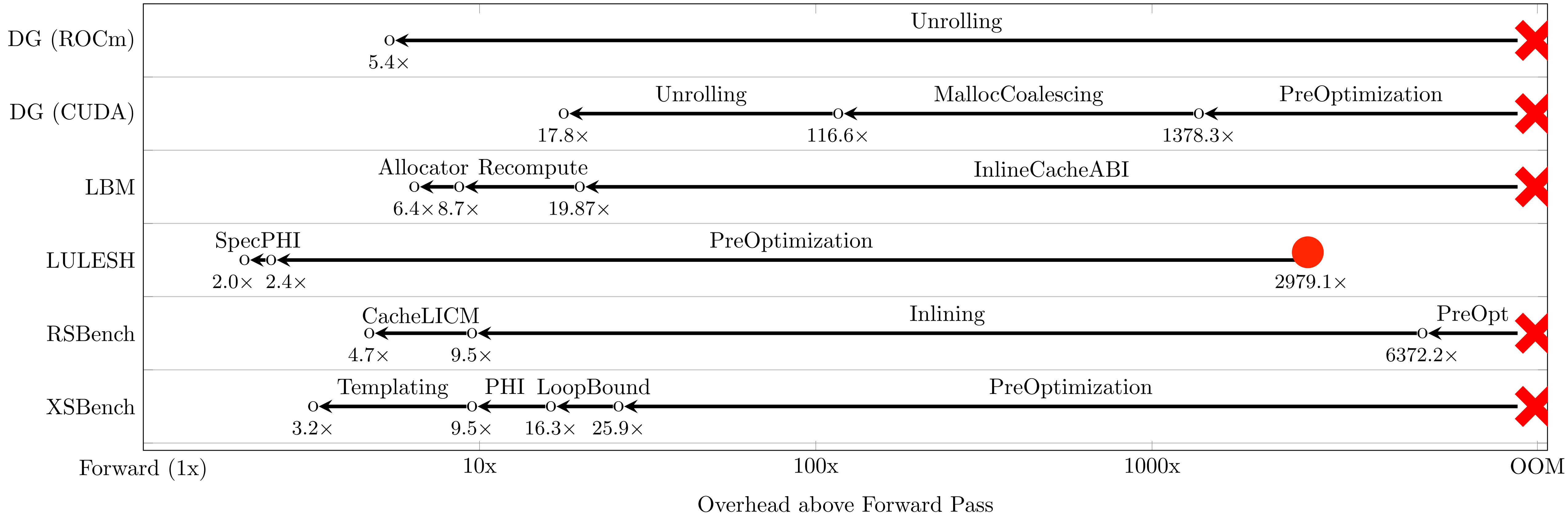
GPU Gradient Overhead

- Evaluation of both original code and gradient
 - DG: Discontinuous-Galerkin integral (Julia)
 - LBM: particle-based fluid dynamics simulation
 - LULESH: unstructured explicit shock hydrodynamics solver
 - XSBench & RSBench: Monte Carlo simulations of particle transport algorithms (memory & compute bound, respectively)

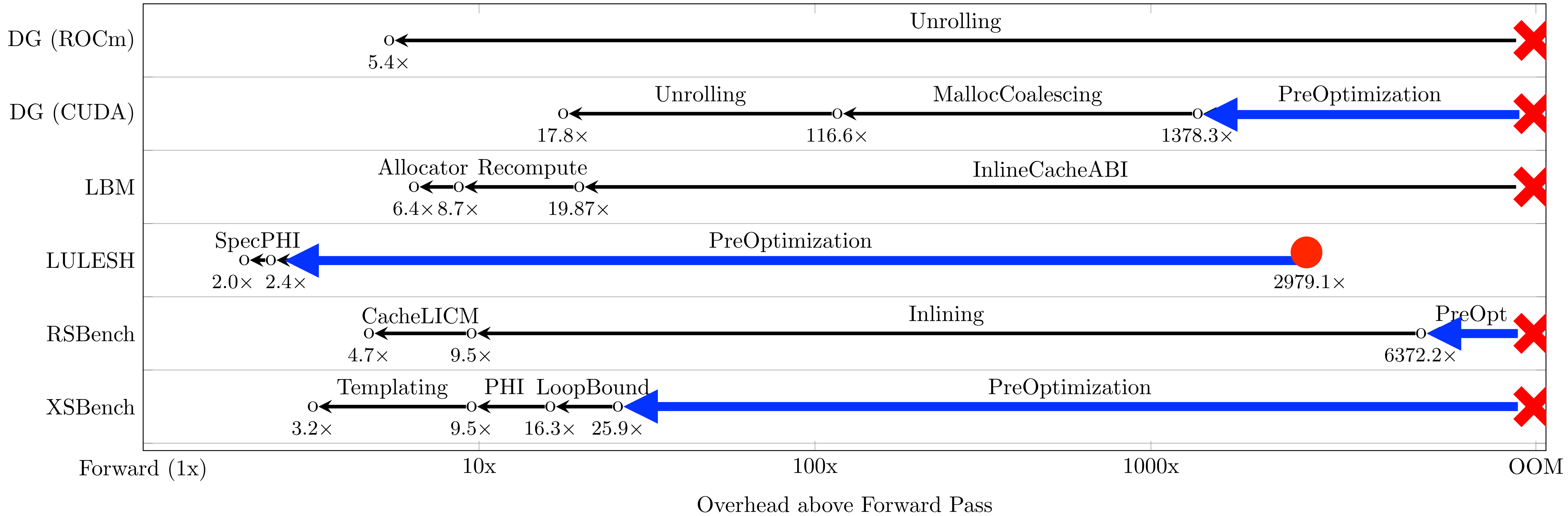


Bug in CUDA
Register Allocator

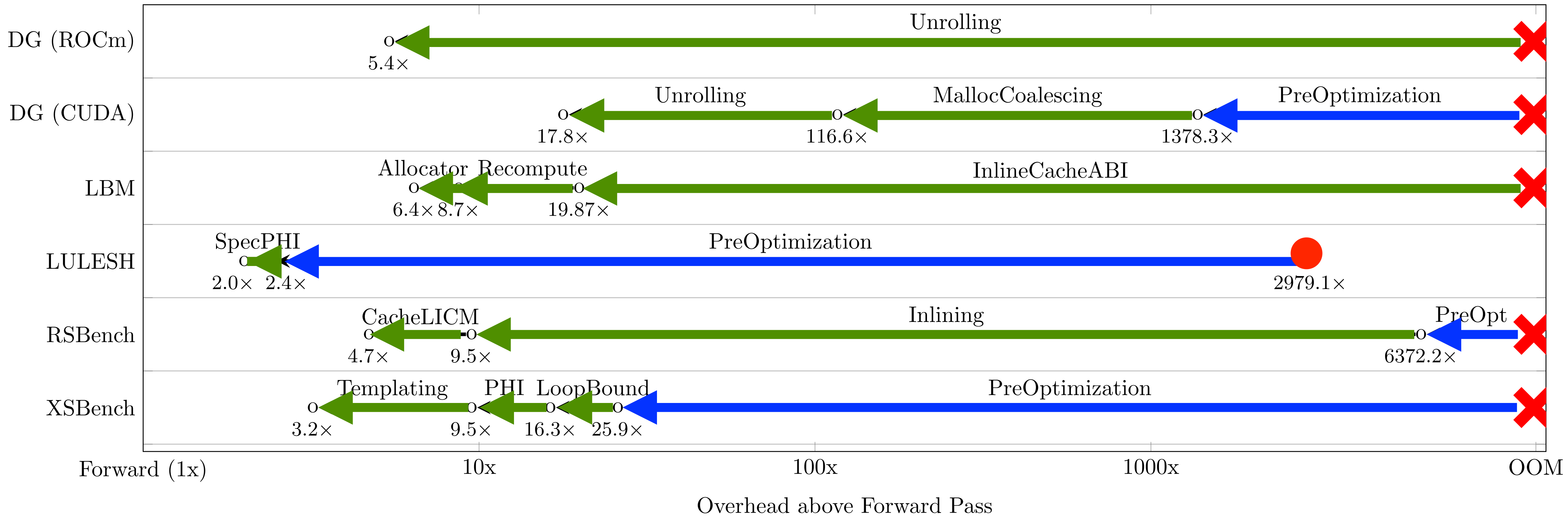
Ablation Analysis of Optimizations



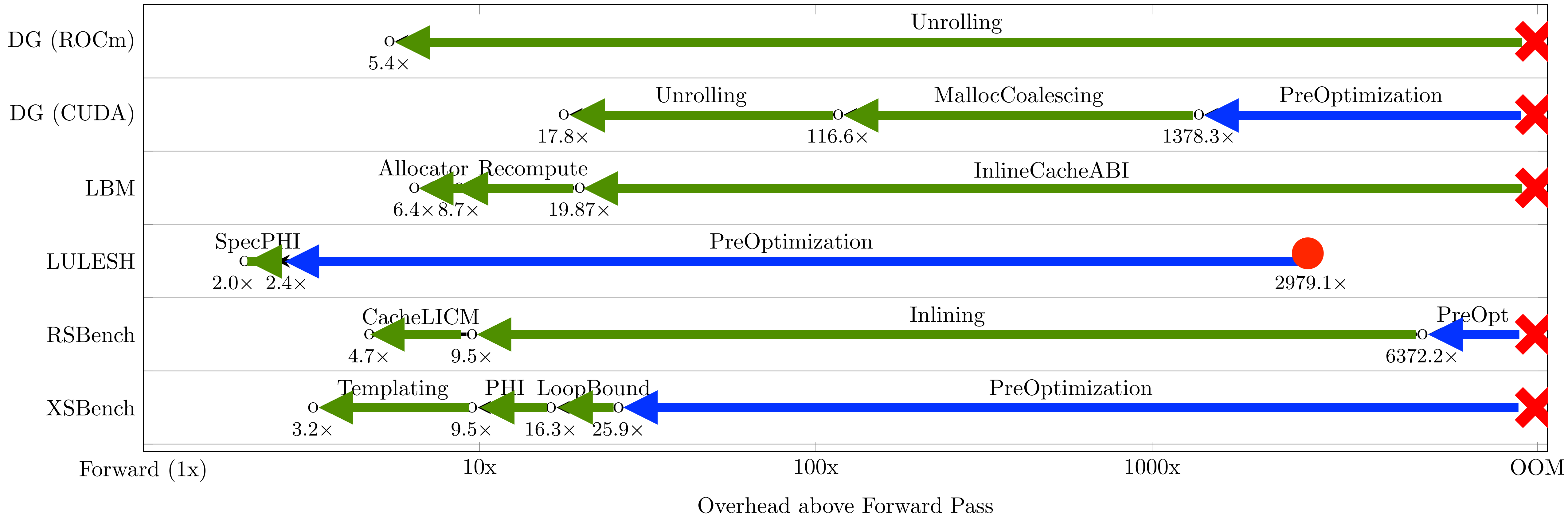
Ablation Analysis of Optimizations



Ablation Analysis of Optimizations



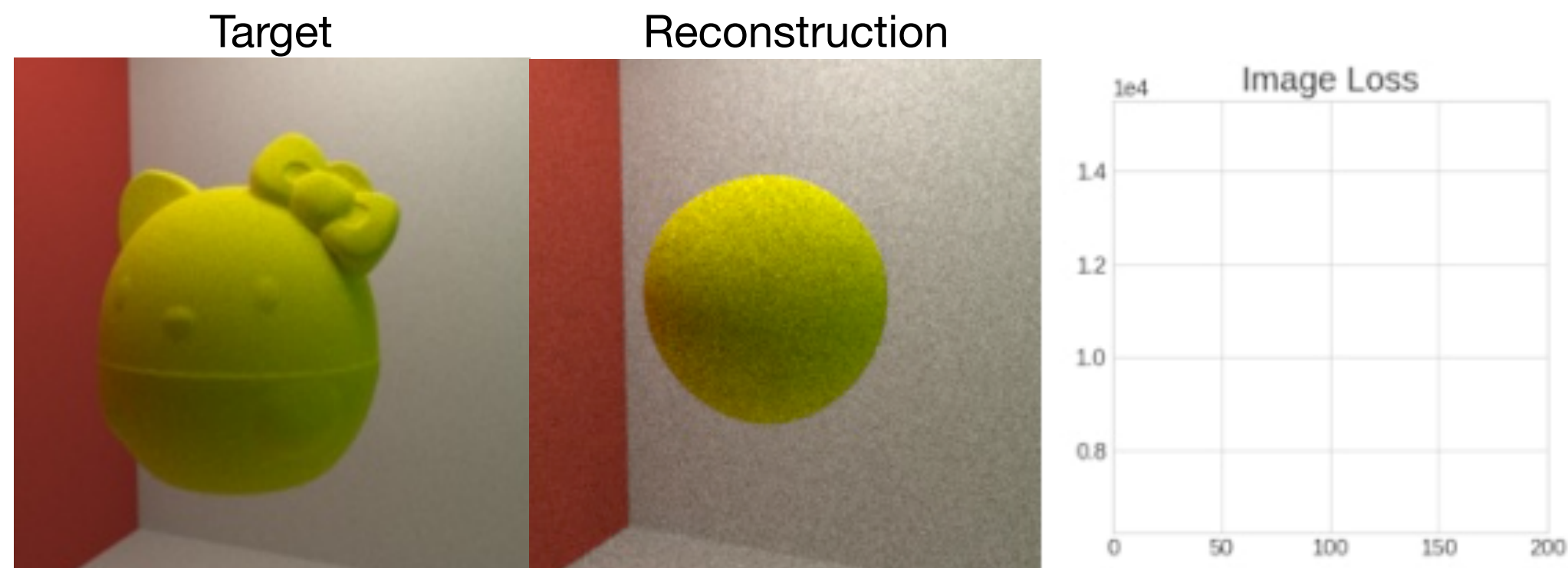
Ablation Analysis of Optimizations



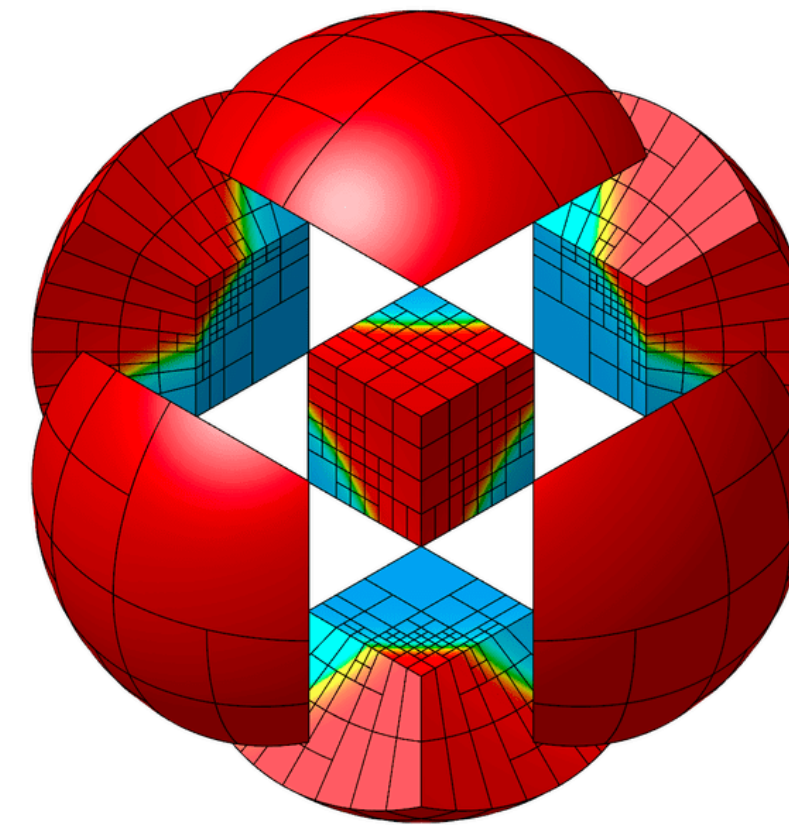
GPU AD is Intractable Without Optimization!



Enzyme-Powered Applications



from [Efficient Differentiation of Pixel Reconstruction Filters for Path-Space Differentiable Rendering](#), SIGGRAPH Asia 2022, Zihan Yu et al



from [MFEM Team at LLNL](#)



>100x speedup!

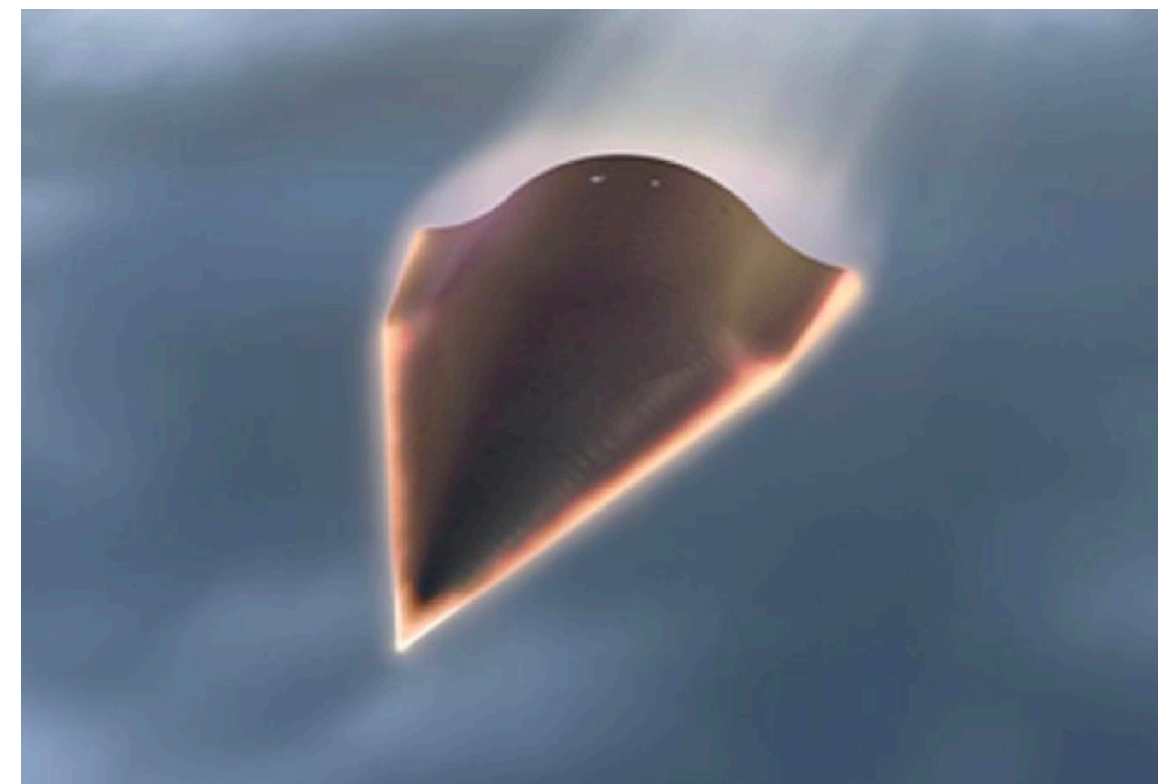
Prior:
5 days (cluster)

Enzyme-Based:
1 hour (laptop)

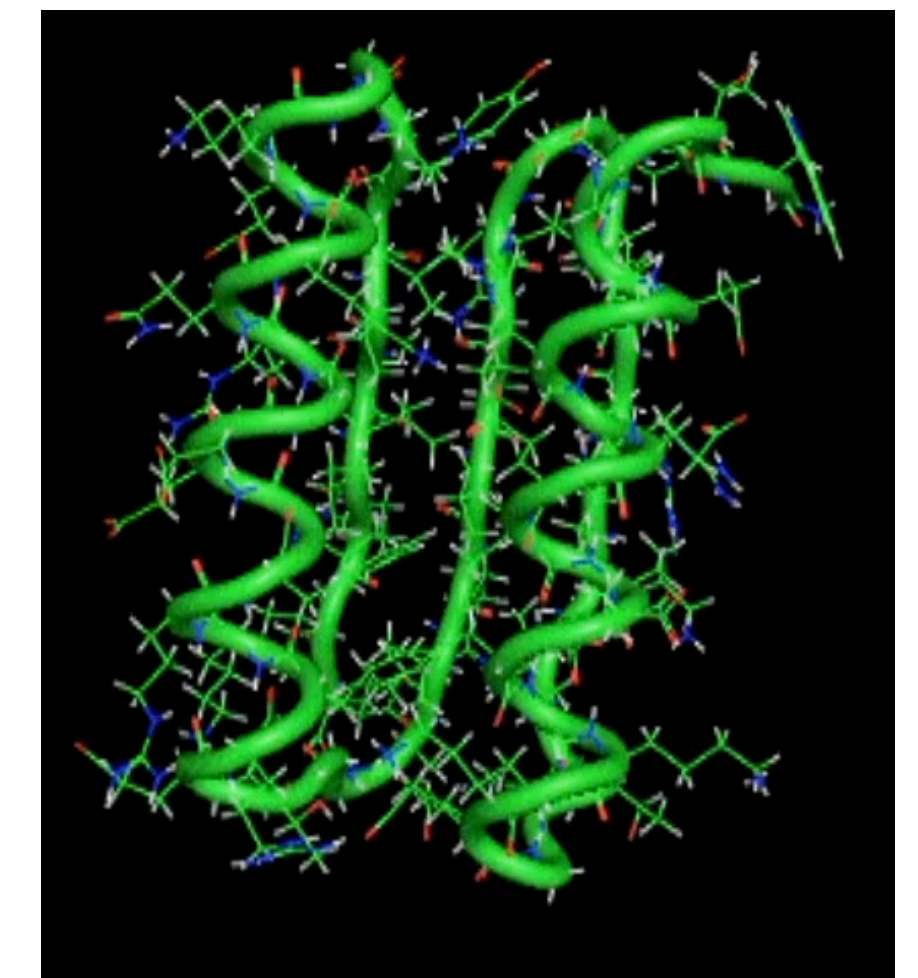
from [Comrade: High Performance Black-Hole Imaging](#) JuliaCon 2022, Paul Tiede (Harvard)



from [CLIMA & NSF CSSI: Differentiable programming in Julia for Earth system modeling \(DJ4Earth\)](#)



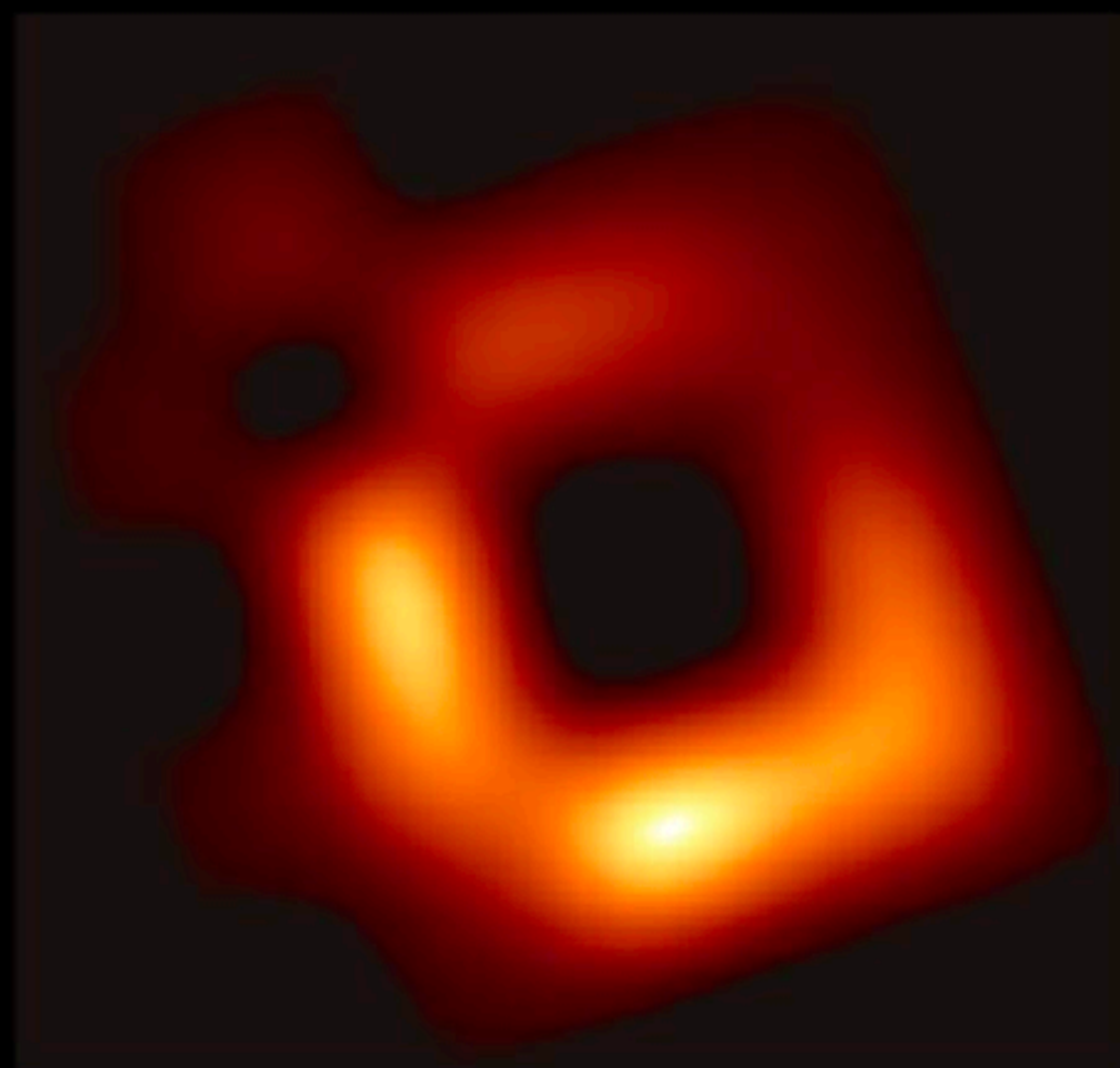
from [Center for the Exascale Simulation of Materials in Extreme Environments](#)



from [Differential Molecular Simulation with Molly.jl](#), EnzymeCon 2023, Joe Greener (Cambridge)

Accelerated Black Hole Imaging with Julia & Enzyme

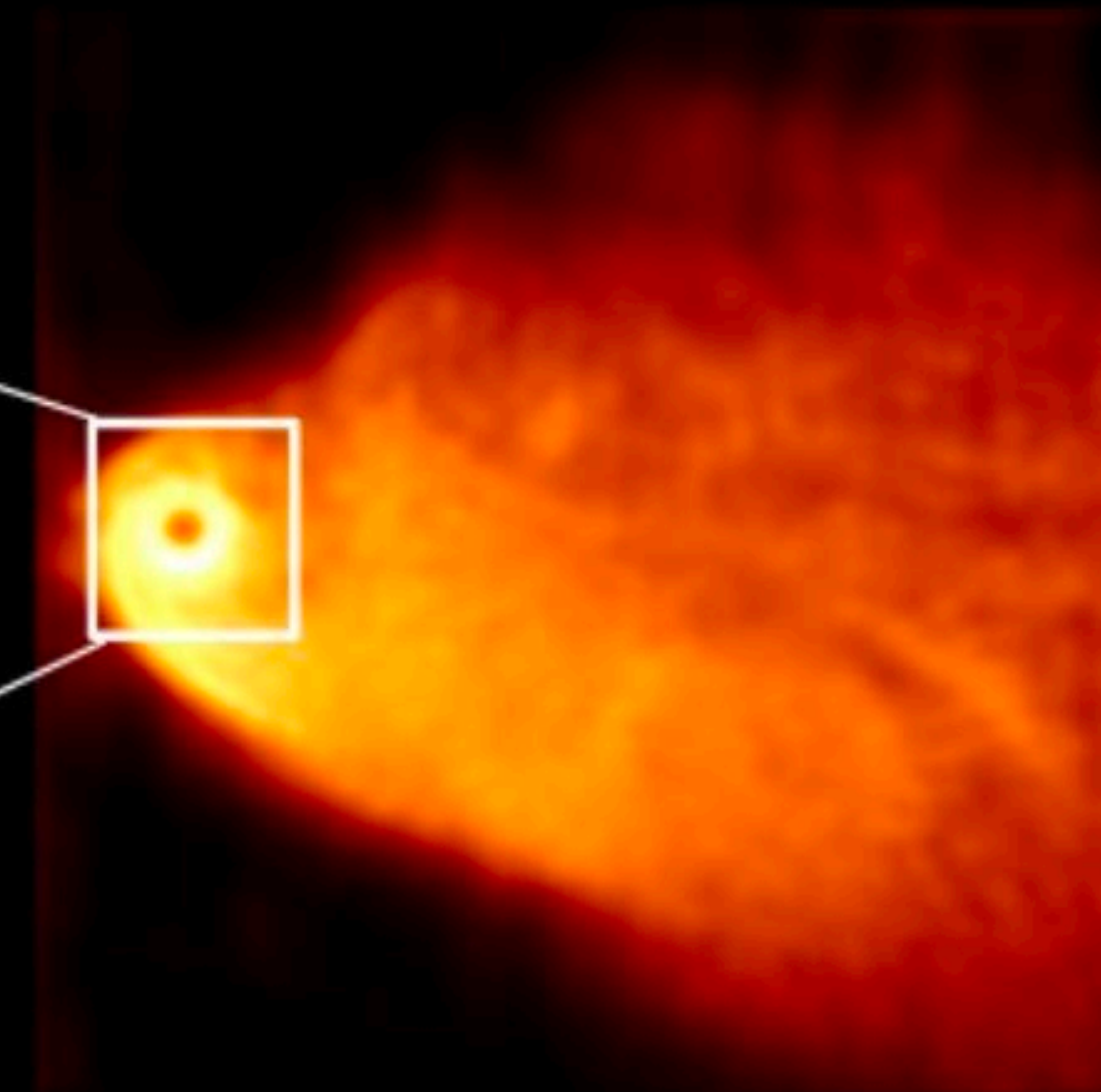
EHT Tools M87 2017
Image Analysis: ~ **1 week** (cluster)



Julia+Enzyme M87 2017
Image Analysis: **1 hour** (1 thread)



Julia+Enzyme next-generation images
Image Analysis: **1-2 days (8 threads)**
(100x increase in computational complexity)



Simulation

Comrade.jl: Julia Bayesian Black Hole Imaging



- Tool for performing reverse-mode (and forward mode) AD of statically analyzable LLVM IR
- Differentiates code in a variety of parallel frameworks (OpenMP, MPI, Julia Tasks, GPU), and languages (C, C++, Fortran, Julia, Rust, Swift, etc)
- Parallel and AD-specific optimizations crucial for performance
- Keep similar scalability as non-differentiated code
- Open source (enzyme.mit.edu & join our mailing list)!
- Ongoing work to support Mixed Mode, Batching, Checkpointing, and more



A Growing Enzyme Community (EnzymeCon 2023)

- 40 attendees spanning developers, users, and everywhere in between.
- 17 great talks from AD internals, to algorithms, to climate science, to physics, and beyond (<https://enzyme.mit.edu/conference>).
- Talks live streamed to YouTube (to be split individually soon):
 - [Day 1 Link](#)
 - [Day 2 Link](#)



Acknowledgements

- Thanks to James Bradbury, Alex Chernyakhovsky, Lilly Chin, Hal Finkel, Marco Foco, Laurent Hascoet, Mike Innes, Tim Kaler, Charles Leiserson, Yingbo Ma, Chris Rackauckas, TB Schardl, Lizhou Sha, Yo Shavit, Dhash Shrivathsa, Nalini Singh, Vassil Vassilev, and Alex Zinenko
- William S. Moses was supported in part by a DOE Computational Sciences Graduate Fellowship DESC0019323. Valentin Churavy was supported in part by the Defense Advanced Research Projects Agency (DARPA) under Agreement No. HR0011-20-9-0016, and in part by NSF Grant OAC-1835443. Ludger Paehler was supported in part by the German Research Council (DFG) under grant agreement No. 326472365.
- This research was supported in part by LANL grant 531711; in part by the Applied Mathematics activity within the U.S. Department of Energy, Office of Science, Advanced Scientific Computing Research Program, under contract number DE-AC02-06CH11357; in part by the Exascale Computing Project (17-SC-20-SC). Research was sponsored by the United States Air Force Research Laboratory and was accomplished under Cooperative Agreement Number FA8750-19-2-1000.
- The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the United States Air Force or the U.S. Government.

