

---

# Precision and Performance Analysis of LLVM's C Standard Math Library on GPUs

By Anton Rydahl, Joseph Huber, Ethan  
McDonough, and Johannes Doerfert  
Prepared by LLNL under Contract DE-AC52-07NA27344.

---

# About Me

- MSc in Mathematical Modelling and Computation
- Interning at Lawrence Livermore National Laboratory
- Working on libc, libc++, and OpenMP in LLVM

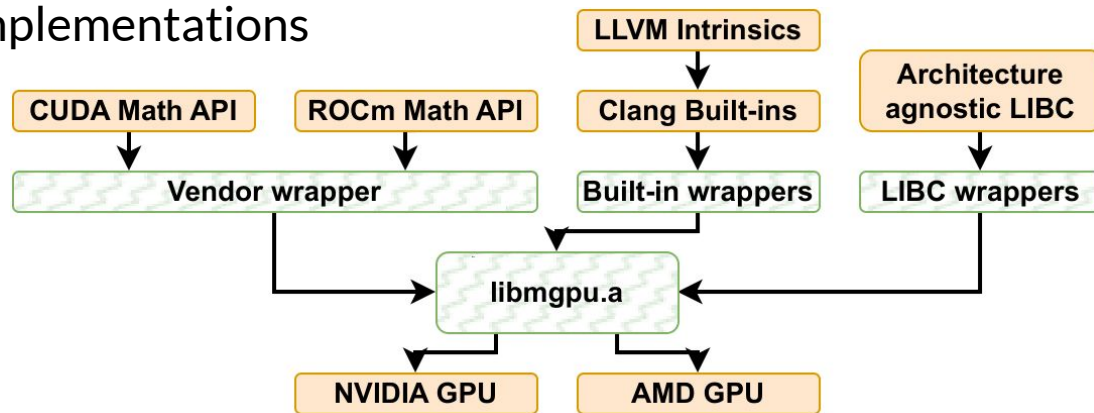


# Background

- LLVM's *libc* is being developed for GPUs
- Clang uses vendor libraries
- **Explore what LLVM infrastructure can be reused on GPUs**

# GPU Math Libraries

- NVIDIA's CUDA Math
- AMD's HIP Math
- LLVM intrinsics
- Target agnostic implementations



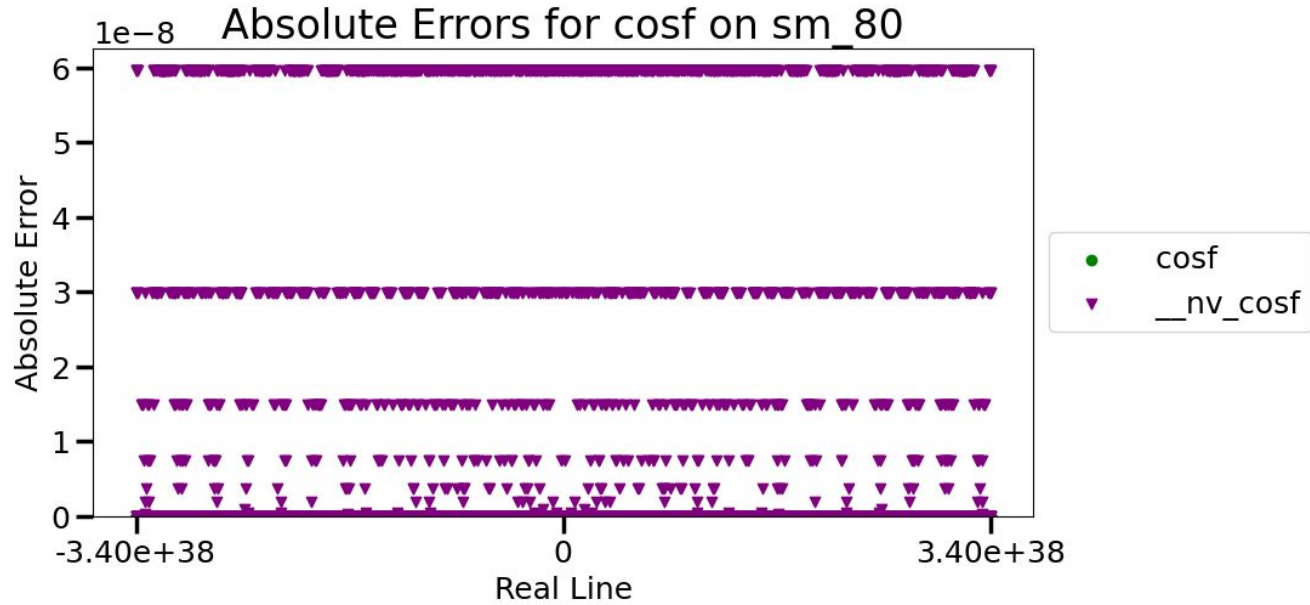
# GPU Math Libraries

```
$ llvmar x libmgpu.a remquoof.cpp.o
$ llvmojdump --offloading remquoof.cpp.o
remquoof.cpp.o:      file format elf64-x86-64
OFFLOADING IMAGE [0]:
kind      llvm ir
arch      gfx906
triple    amdgcncmd-amdhsa
producer  none
OFFLOADING IMAGE [1]:
kind      llvm ir
arch      gfx90a
triple    amdgcncmd-amdhsa
producer  none
OFFLOADING IMAGE [2]:
kind      llvm ir
arch      sm_70
triple    nvptx64-nvidia-cuda
producer  none
OFFLOADING IMAGE [3]:
kind      llvm ir
arch      sm_80
triple    nvptx64-nvidia-cuda
producer  none
```

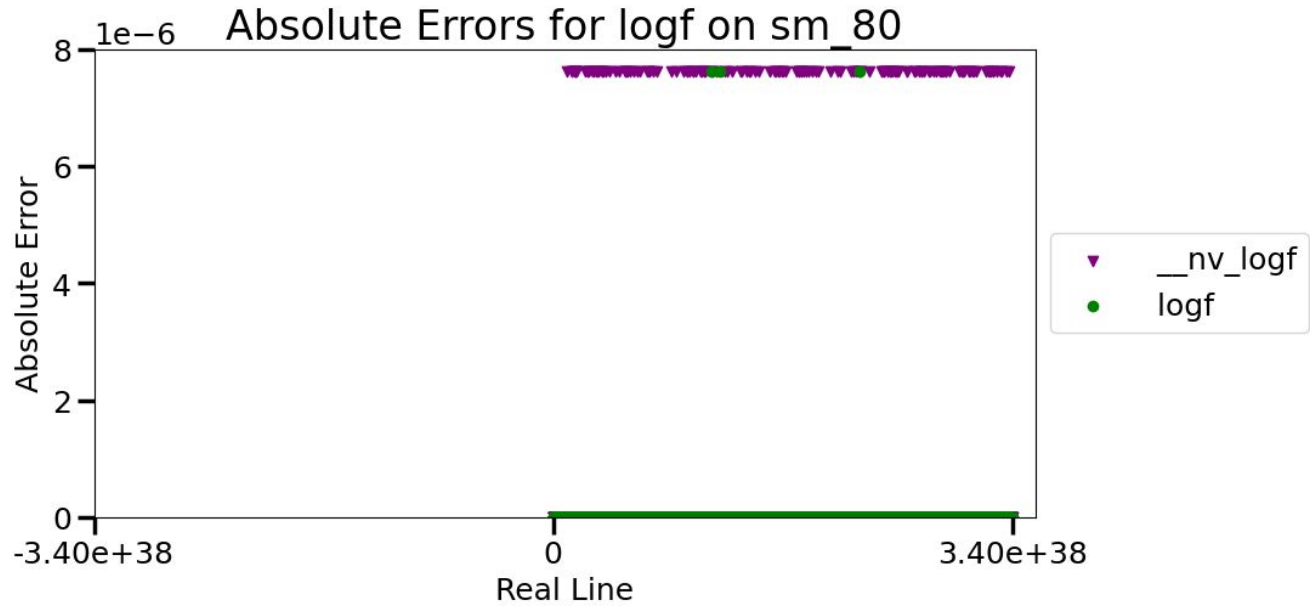
# Correctness

- Exhaustive search for univariate functions with 32 bit data types
  - Upper bound on error
- Uniformly distributed input for 64 bit data
  - Lower bound on error
- Comparing against the GNU MPFR library
  - Arbitrary precision

# Correctness

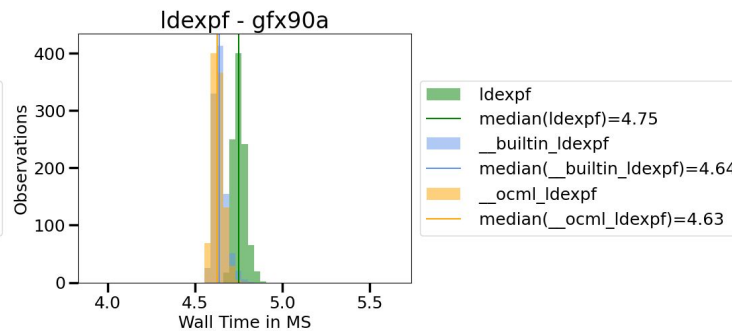
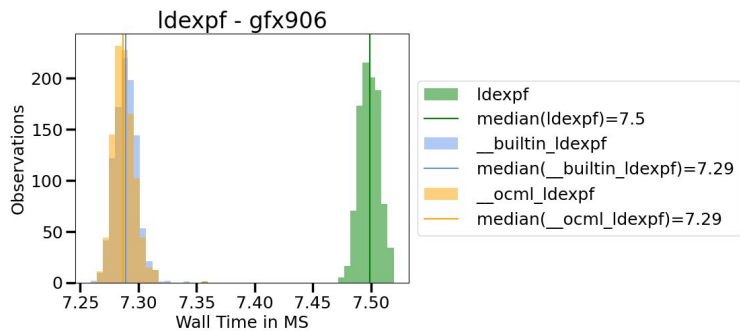


# Correctness

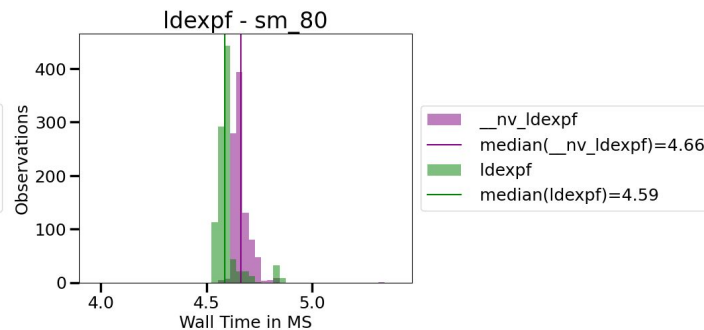
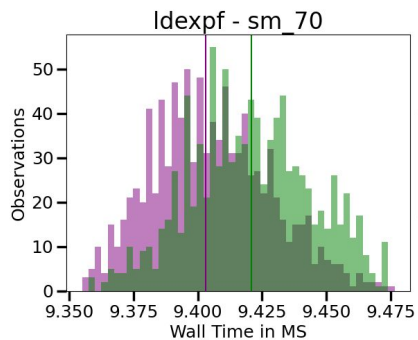




# Timings



# Timings



# Fastest Correct Set of Functions

- HIP
- CUDA
- LLVM Libc
- LLVM Builtin

Function	gfx906	gfx90a	sm_70	sm_80
log	L	L	C	C
log10	L	H	C	C
log1p	L	L	L	L
log2	L	H	C	C
logb	H	H	C	L
logbf	H	H	B	C
nearbyint	H	B	C	B
nearbyintf	H	H	B	B
nextafter	B	L	L	C
nextafterf	L	L	C	C
pow	H	H	C	C
powf	H	H	C	C
remainder	H	H	B	B
remainderf	H	H	B	B

# Fastest Correct Set of Functions

- HIP
- CUDA
- LLVM Libc
- LLVM Builtin

Function	gfx906	gfx90a	sm_70	sm_80
log	L	L	C	C
log10	L	H	C	C
log1p	L	L	L	L
log2	L	H	C	C
logb	H	H	C	L
logbf	H	H	B	C
nearbyint	H	B	C	B
nearbyintf	H	H	B	B
nextafter	B	L	L	C
nextafterf	L	L	C	C
pow	H	H	C	C
powf	H	H	C	C
remainder	H	H	B	B
remainderf	H	H	B	B

# Results of the Analysis

- Given a tolerance, we can find an optimal set of mathematical functions
  - Depends on the target architecture
  - 7 times faster than CUDA Math on sm\_80 on average
    - Influenced by outliers
    - Sensitive to inlining
  - 5 % faster than HIP Math on gfx906