# Improving the OpenMP Offloading Driver: LTO, Libraries, and Toolchains

LLVM Performance Workshop @ CGO 2022

April 3rd 2022

Joseph Huber

ORNL is managed by UT-Battelle, LLC for the US Department of Energy

U.S. DEPARTMENT OF **ENERGY**

# Overview & Motivation

OAK RIDGE
National Laboratory

# OpenMP Offloading Overview

- Allows users to offload execution of code to another device
- Requires the compiler driver to compile & link multiple programs
- The linked image also needs to be registered

```cpp
#include <complex>

using complex = std::complex<double>;

void zaxpy(complex *X, complex *Y, complex D, int N) {
#pragma omp target teams distribute parallel for
  for (int i = 0; i < N; ++i)
    Y[i] = D * X[i] + Y[i];
}

int main() {
  const int N = 1024;
  complex X[N], Y[N], D;
#pragma omp target data map(to:X[:N]) map(tofrom:Y[:N])
  zaxpy(X, Y, D, N);
}
```
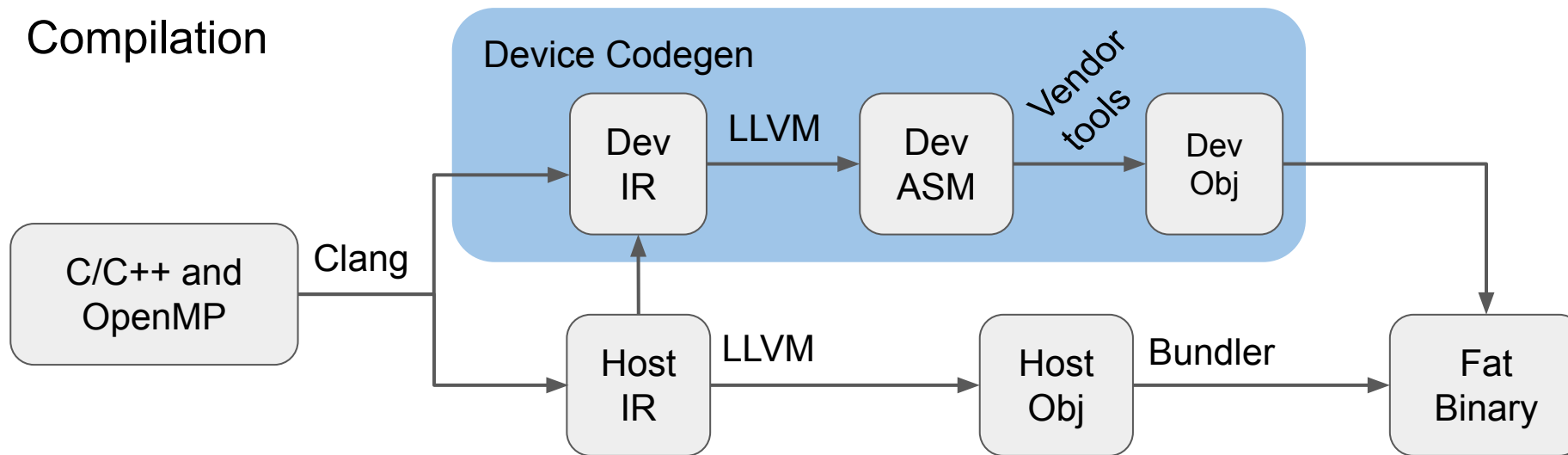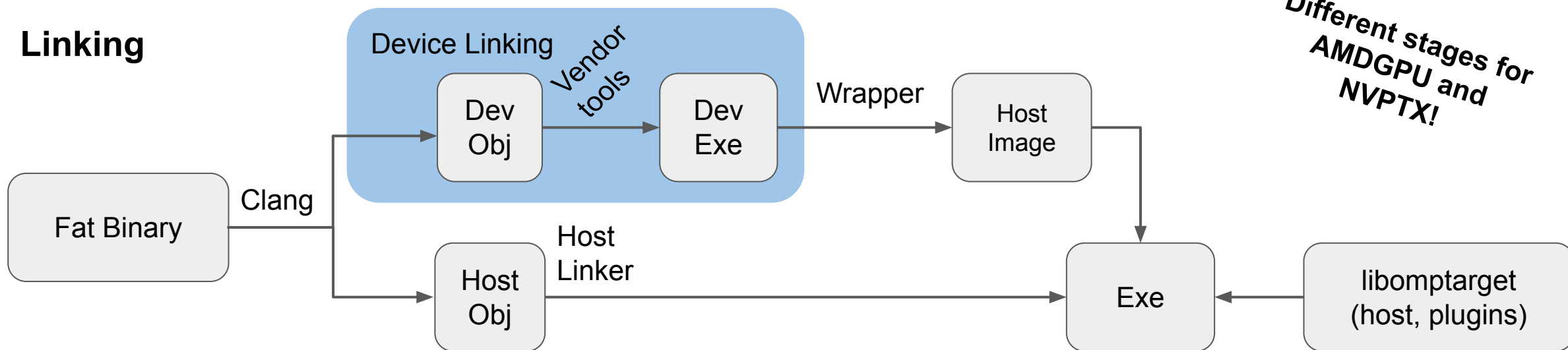
OAK RIDGE
National Laboratory

# Current OpenMP Offloading Driver Overview

Compilation



Linking



*Different stages for AMDGPU and NVPTX!*

OAK RIDGE
National Laboratory

# Motivation

- Why isn't the current method good enough?
- Handle device code the same as host code and support **static linking**
- Unify the required stages across **all toolchains**
- Support **Link Time Optimization** on the device
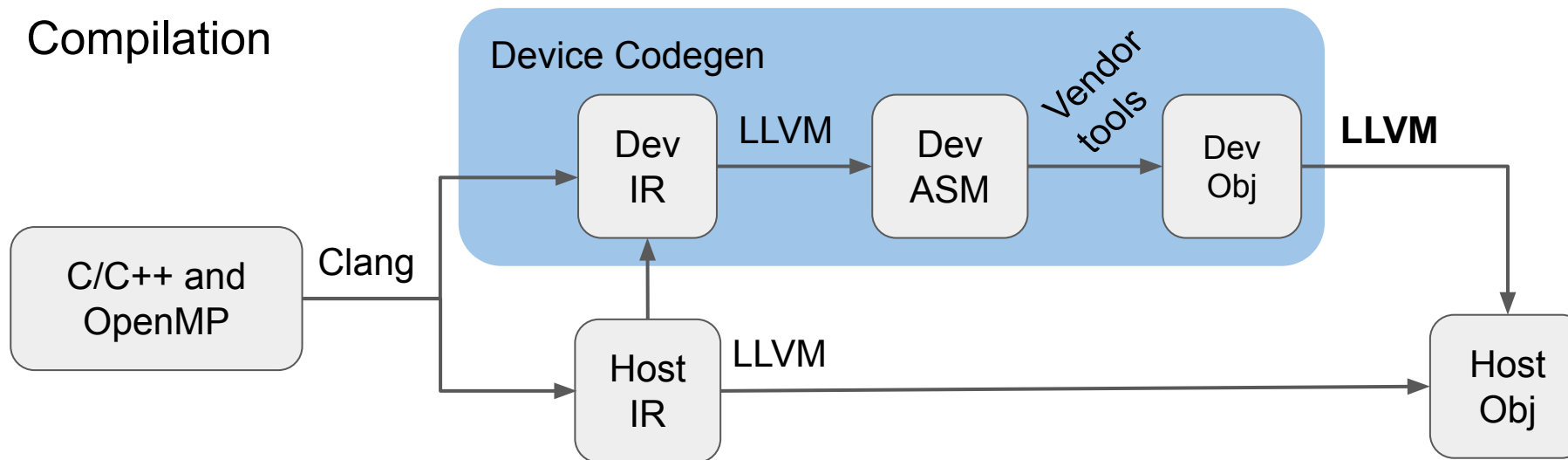- Enable offloading language **interoperability**

OAK RIDGE
National Laboratory

# New Driver Implementation

OAK RIDGE
National Laboratory

# New OpenMP Offloading Driver

- Embed the device objects directly in the host object
  - Data is stored in an excluded section
- Linking is done by a linker wrapper application
  - Scan each input for embedded device objects
  - Extract & link each image
  - Wrap the linked device image in a new host object file
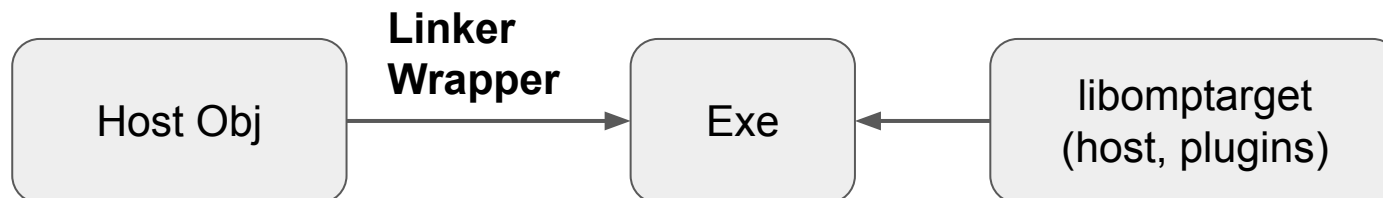  - Run the original linking job with the new wrapped image

**OAK RIDGE**
National Laboratory

# New OpenMP Offloading Driver Overview

Compilation

Device code
embedded directly
as a section in host
object

Device Codegen

| C/C++ and OpenMP | → Clang → | Dev IR | → LLVM → | Dev ASM | → Vendor tools → | Dev Obj | → **LLVM** → | Host Obj |

Host IR → LLVM → Host Obj

Linking

Device linking
complexity handled
in a single stage

| Host Obj | → **Linker Wrapper** → | Exe | ← | libomptarget (host, plugins) |

**OAK RIDGE**
National Laboratory

# Embedding OpenMP Offloading Code



**Active Toolchains**

**Section Table**

**Section Contents**

C/C++ and OpenMP → Host IR → Host Object

Device Object

Host Sections

.llvm.offloading

omp_offloading_entries

<triple and arch>

<Bitcode or Object>

<...>

Contains kernels and globals to register

```
@.llvm.embedded.object = private constant [N x i8] c"...", section ".llvm.offloading"
```

OAK RIDGE
National Laboratory

# OpenMP Offloading Linker Wrapper

Device Linking

Dev BC → LTO → Dev Obj

Dev Obj → Vendor Linker → Dev Exe

Extract

Links every extracted object with a compatible triple & architecture

Wrapper → Host Image

Linker Input

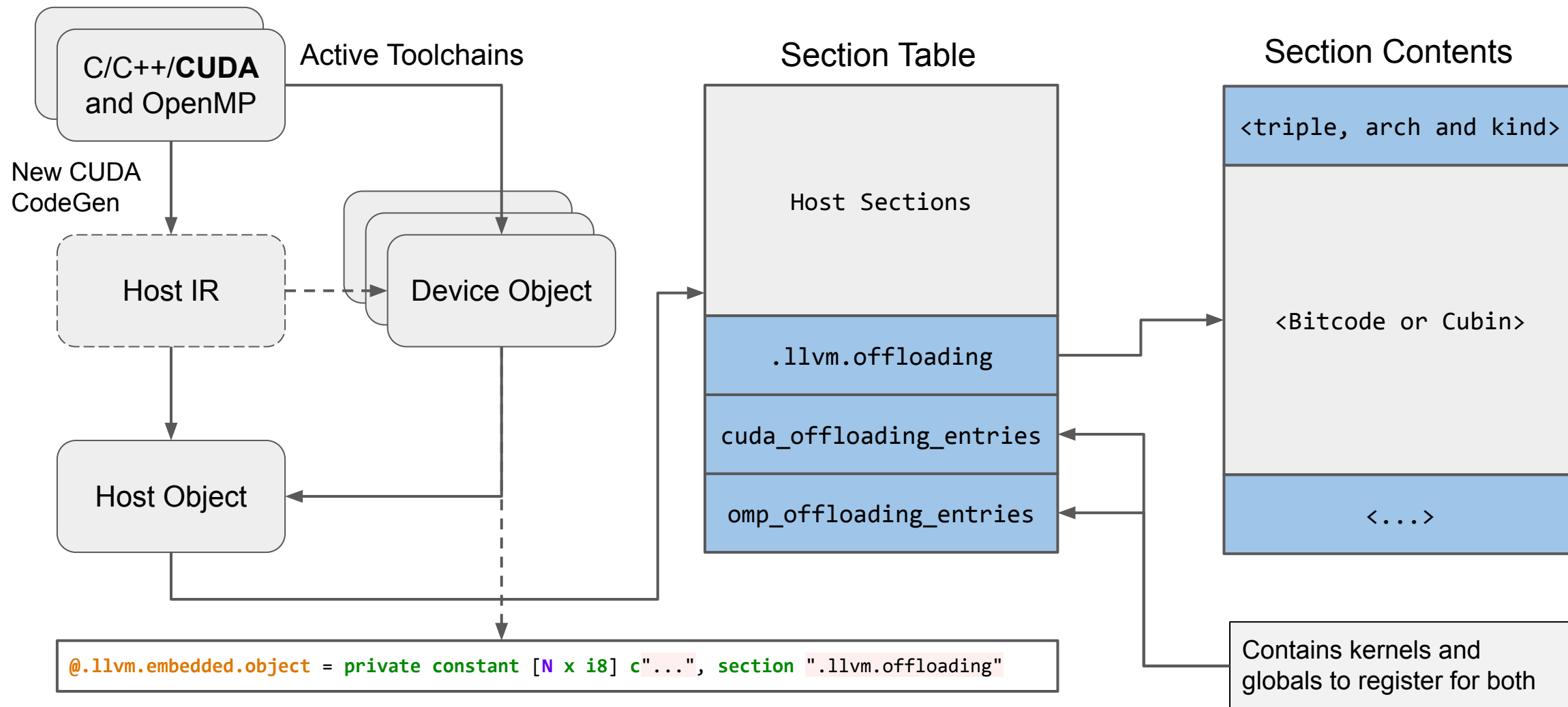Wrapped Linker job

Exe

# Benefits

- Offloading binaries behave like host binaries
  - Static libraries and relocatable linking works as expected
- Fewer stages required to create an offloading program
- Much simpler driver code
- Fully functional LTO on the device **-foffload-lto**
  - Greatly improves performance on some applications
  - The OpenMPOpt pass greatly benefits from whole program visibility
    - (See Optimizing OpenMP GPU Execution in LLVM @ LLVM Dev2021)
- Will be the default method for OpenMP Offloading very soon!

**OAK RIDGE**
National Laboratory

# Future Work & Interoperability
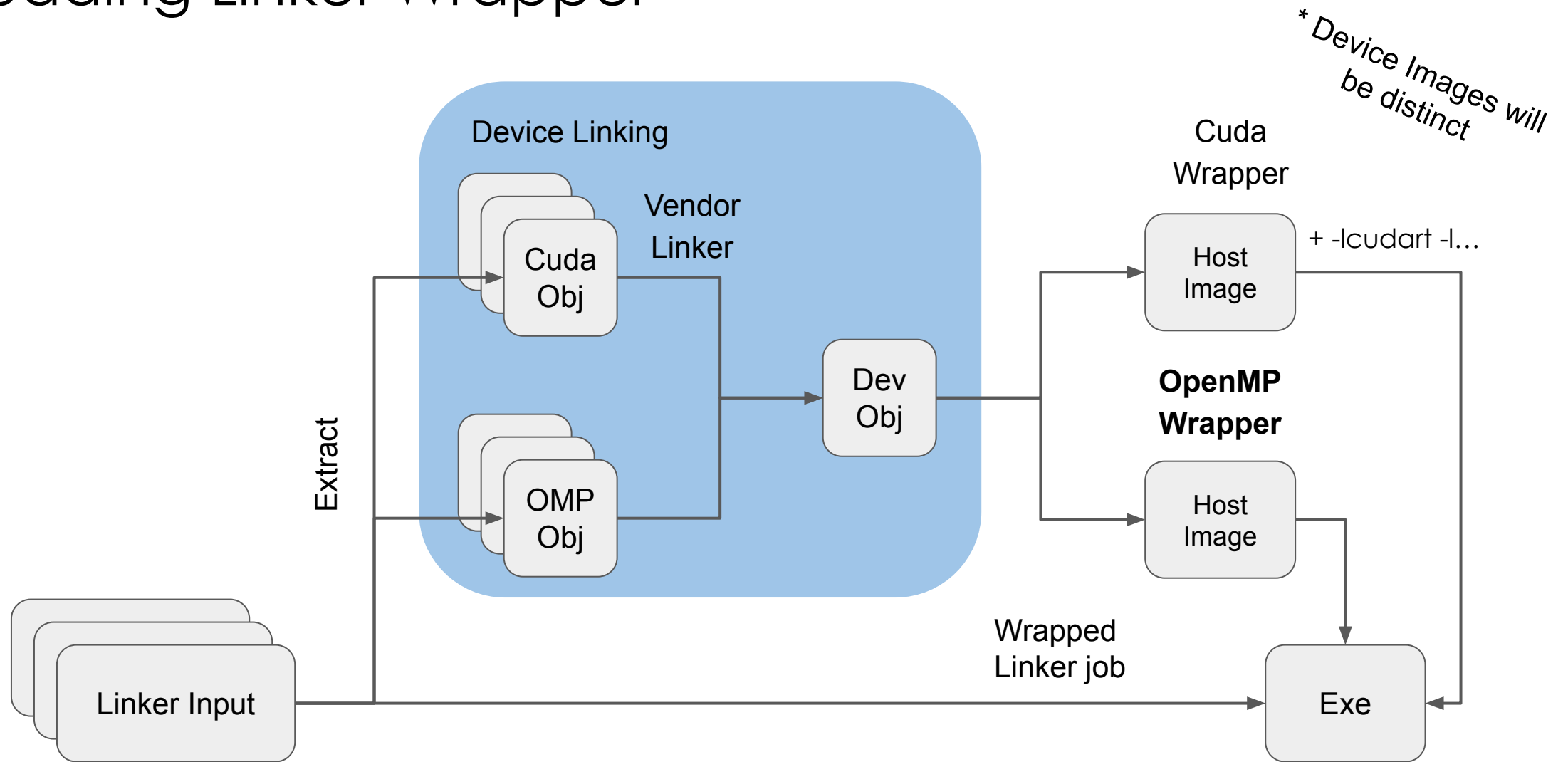
Oak Ridge
National Laboratory

# Extending the New Driver

- The new driver can be adapted for CUDA / HIP as well
  - Change code-generation to support the offloading sections
  - Implement a wrapper for CUDA / HIP code
- Allows for redistributable device code (RDC) support in Clang
- The linker wrapper will link all compatible object files
- Allows for OpenMP to call CUDA code and vice-versa
  - Needs additional code for full interoperability

**OAK RIDGE**
National Laboratory

# Embedding CUDA & OpenMP Offloading Code



```
@.llvm.embedded.object = private constant [N x i8] c"...", section ".llvm.offloading"
```

OAK RIDGE
National Laboratory

# Offloading Linker Wrapper

Device Linking

Vendor Linker

Cuda Obj

OMP Obj

Dev Obj

Extract

Linker Input

Cuda Wrapper

Host Image

**OpenMP Wrapper**

Host Image

\* Device Images will be distinct

+ -lcudart -l...

Wrapped Linker job

Exe

OAK RIDGE
National Laboratory

# Generic Offloading Libraries

- Create a static library with code for every offloading target
  - Allow more compatible architectures to be linked
- No longer need to specially compile & link device bitcode libraries

Section Contents

| |
|---|
| `<nvptx64 sm_80>` |
| `<nvptx64 sm_70>` |
| `<nvptx64 any>` |
| `<amdgcn gfx908>` |
| `<...>` |

**OAK RIDGE**
National Laboratory

# Linker Wrapper in the Linker

- Currently we rely on Clang to call the linker wrapper with the appropriate arguments
- Prevents offloading code from being truly agnostic
- Embed the linker wrapper functionality inside a linker plugin or LLD

**OAK RIDGE**
National Laboratory

# Application Experiences

OAK RIDGE
National Laboratory

# MiniMDock

- Protein-ligand docking mini-application
- A call to an external function prevented a crucial optimization
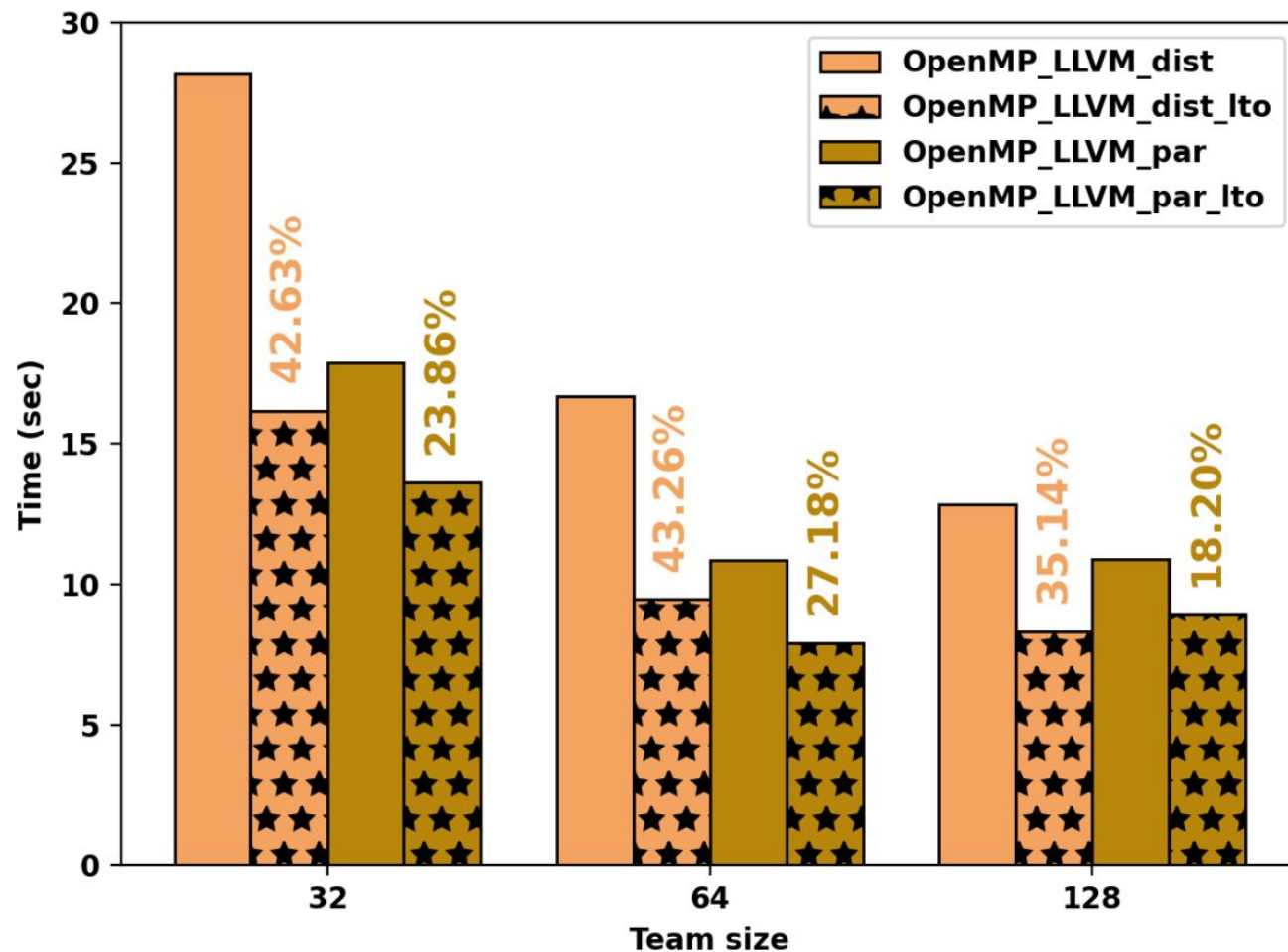- Device-side LTO allows us to see the whole program



Figure Generated by Mathialakan Thavappiragasam

# Thermo4PFM

- Library to evaluate alloy compositions in Phase-Field models
- Application built with static libraries
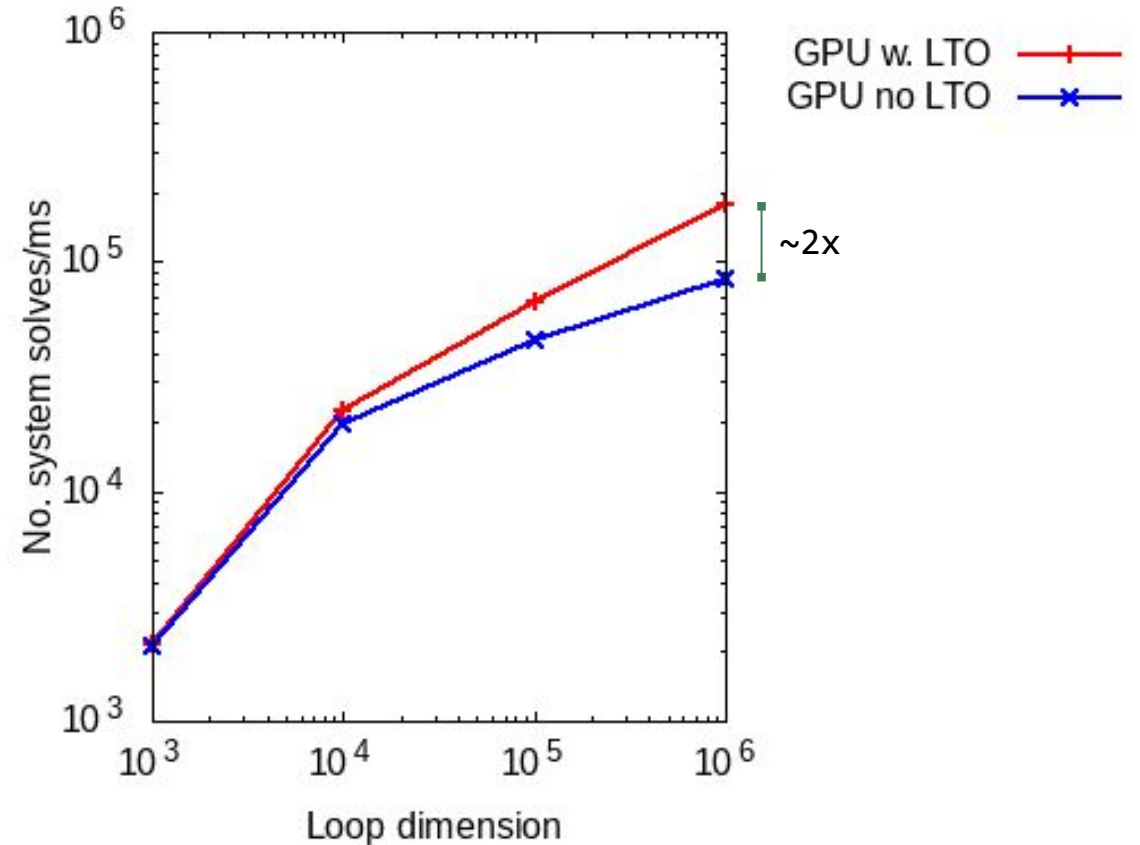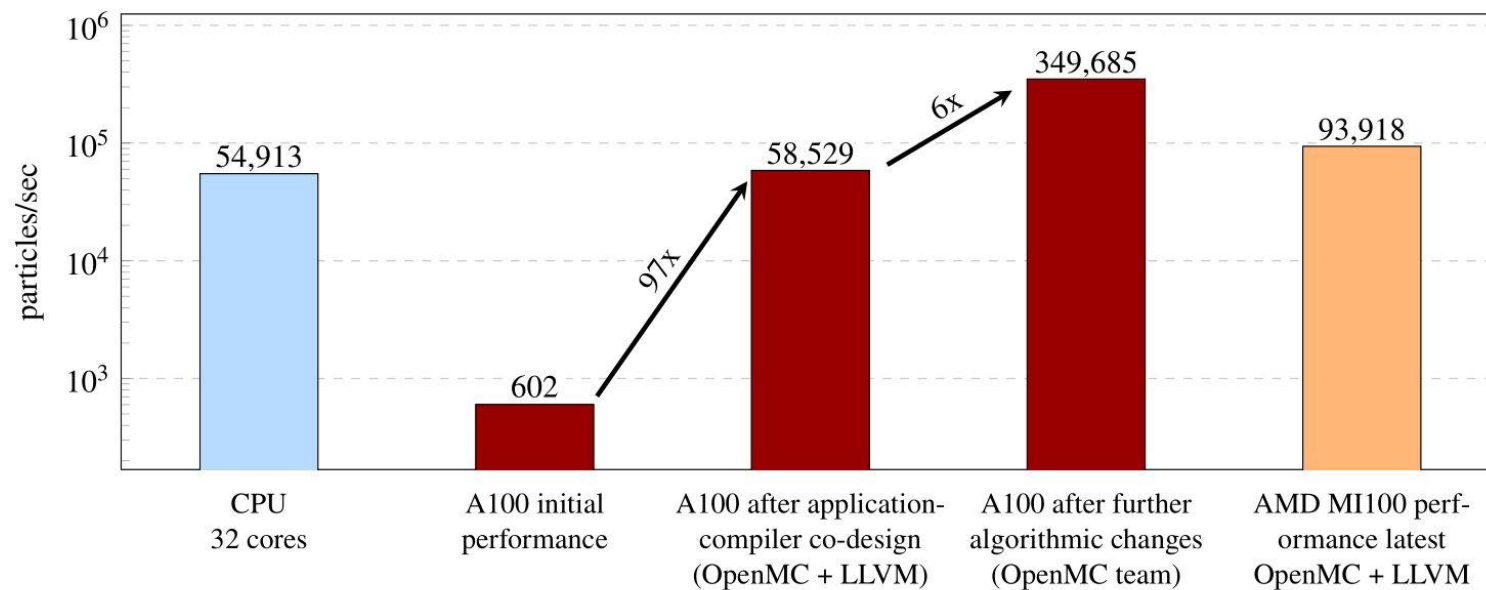- Large amount of files with device code



**Figure Generated by Jean-Luc Fattebert**

**OAK RIDGE**
National Laboratory

# OpenMC

- Monte-Carlo particle transport application
- Needed a CMake Unity build to get reasonable performance
- Device-side LTO gives the same performance and compiles several times faster



The data for the figure was generated and generously provided by John Tramm.

OAK RIDGE
National Laboratory

# MiniQMC

- Quantum Monte-Carlo mini-application
- Made heavy use of static libraries
  - Can now compile without a CMake workaround
- No performance difference with and without LTO

`-D USE_OBJECT_TARGET=ON` is used to workaround static linking issue.

**OAK RIDGE**
National Laboratory

# Conclusion & Closing Thoughts

OAK RIDGE
National Laboratory

# Conclusion & Closing thoughts

- The new driver greatly improves the usability of OpenMP Offloading in LLVM
  - Allows interoperability
  - Multiple devices embedded in the same binary
- Device-side LTO gives real-world applications significant performance increases
- A unified offloading Toolchain is possible

**OAK RIDGE**
National Laboratory

# Questions?

**OAK RIDGE**
National Laboratory

# Current Drawbacks

- Currently still relies on Clang to call the linker wrapper
- Cannot embed device code without host LLVM IR
  - Need a new phase to perform an objcopy
- Cannot properly handle incremental compilation
  - e.g. clang foo.c -S

OAK RIDGE
National Laboratory