# CARTS: Enabling Event-Driven Task and Data Block Compilation for Distributed HPC

**Rafael Andres Herrera Guaitero**

PhD Candidate
University of Delaware

Saturday, March 1

# Contributors

- **Joseph B. Manzano Franco**
  Pacific Northwest National Laboratory

- **Joshua D. Suetterlein**
  Pacific Northwest National Laboratory

- **Xiaoming Li**
  University of Delaware

- **Andres Marquez**
  Pacific Northwest National Laboratory

# Outline

MOTIVATION & CONTEXT

CARTS: COMPILER FOR ARTS

FUTURE WORK

# Motivation & Context - I

- **Evolving Architectures:** Modern HPC and AI/ML workloads demand architectures that efficiently leverage memory systems.

- **Hardware Heterogeneity:** The rise of multi-core CPUs, GPUs, and specialized accelerators creates complexity in managing resources.

- **Performance Pressure:** Increased demand for efficient concurrency, synchronization, and communication in large-scale systems.
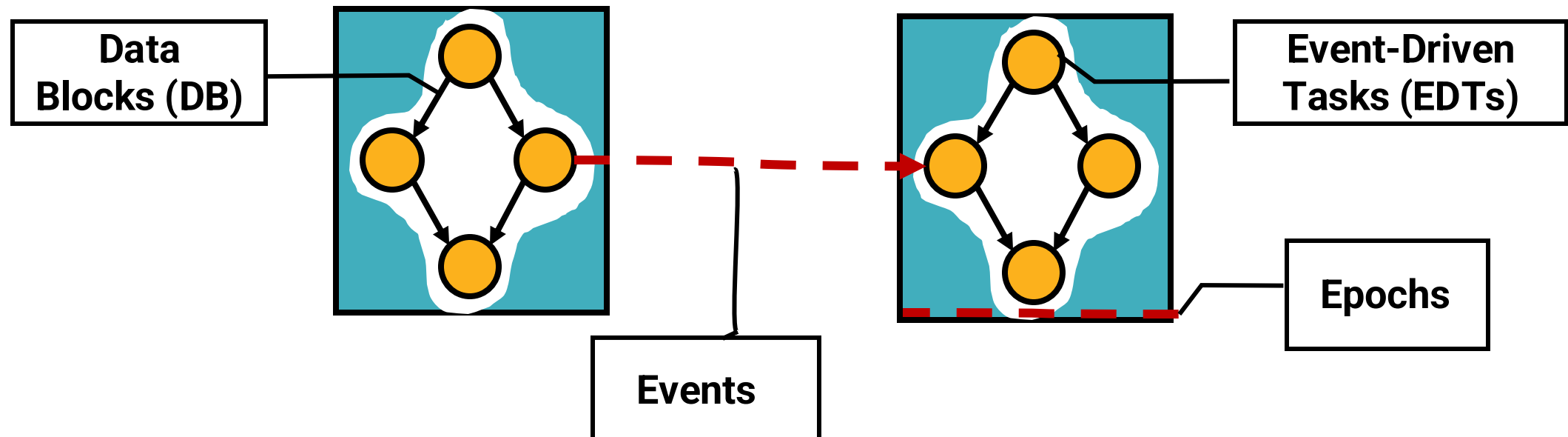
# Motivation & Context - II

- **OpenMP:** A widely adopted standard in C/C++ for specifying parallel regions, tasks, and dependencies through intuitive pragmas.

- **MLIR:** A flexible intermediate representation framework that supports multi-level optimizations and custom dialects, bridging high-level abstractions with low-level execution details.

- **ARTS (Abstract RunTime System)** is a runtime infrastructure engineered for fine-grained concurrency and efficient task scheduling in distributed systems

# Motivation & Context - III

## Asynchronous RunTime System (ARTS)

It provides users with a distributed global address space, a distributed memory model, and synchronization constructs to write efficient applications on a massively parallel system.



Data Blocks (DB)

Event-Driven Tasks (EDTs)

Events

Epochs

# What is the need?
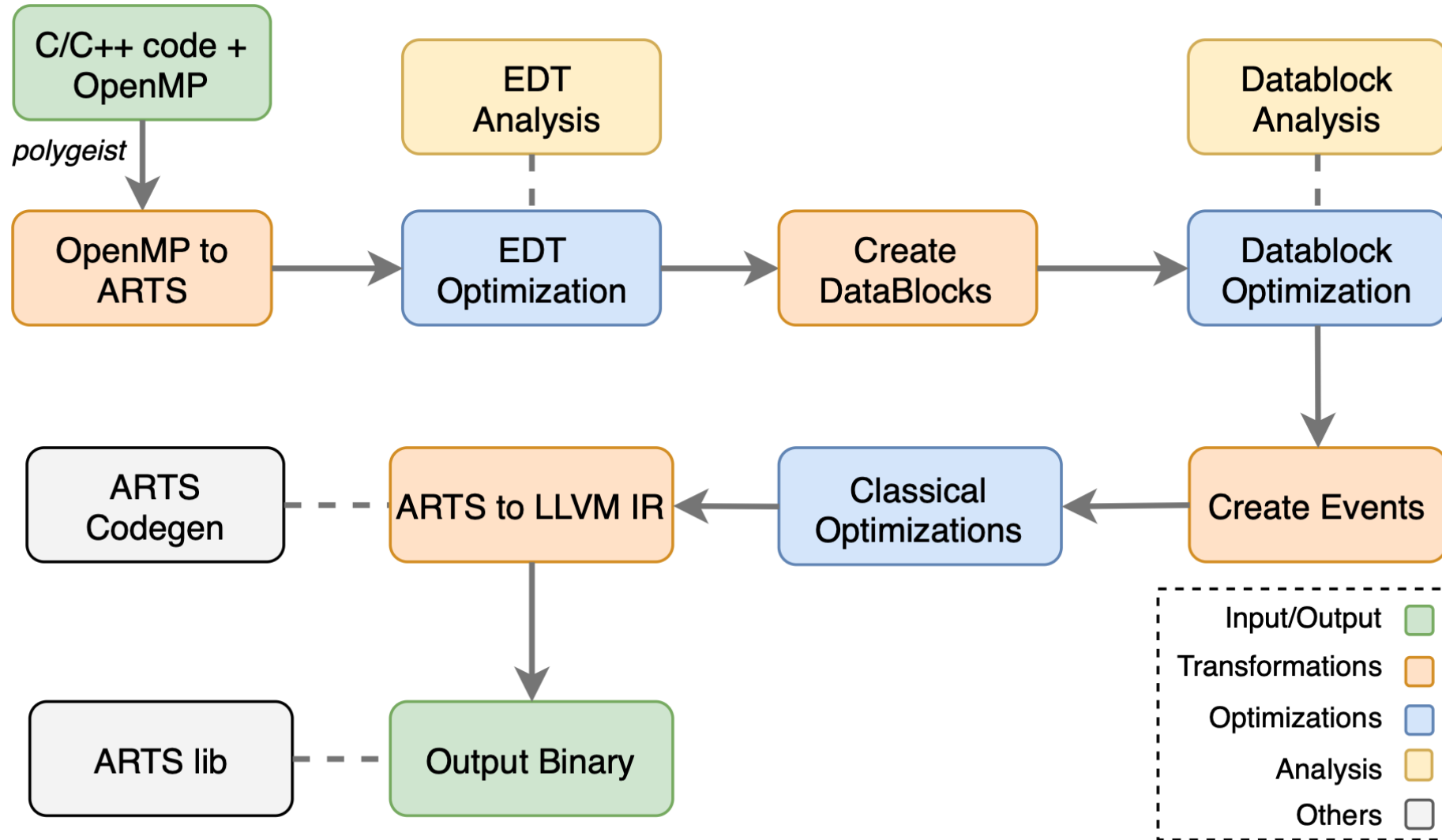
Compiler Infrastructure for ARTS

# Compiler for ARTS

A task-centric compiler pipeline that transforms OpenMP-annotated C/C++ code into an ARTS-friendly format using MLIR.



Programmer → C/C++ + OpenMP → CARTS → LLVM IR + ARTS API Calls

We extended *Polygeist* by providing support to more OpenMP constructs (e.g. tasks) and other C/C++ related operations

1 - C/C++ + OpenMP support

C/C++ code + OpenMP

*polygeist*

**OpenMP to ARTS**

EDT Analysis

Transform standard MLIR + OpenMP standard dialects to the ARTS Dialect.

Datablock Analysis

Datablock Optimization

ARTS Codegen

ARTS to LLVM IR

Classical Optimizations

Create Events

ARTS lib

Output Binary

Input/Output
Transformations
Optimizations
Analysis
Others

2 – Mapping to ARTS

**EDT Analysis**

**EDT Optimization**

C/C++ code + OpenMP

*polygeist*

OpenMP to ARTS

Create DataBlocks

Datablock Analysis

Datablock Optimization

ARTS Codeg...

Classical Optimizations

Create Events

- Remove unused EDTs
- Convert a 'parallel' EDT with only a 'single' EDT in the region into a 'sync' EDT

ARTS lib

Output Binary

Input/Output
Transformations
Optimizations
Analysis
Others

3 – EDTs Analysis and Optimization

```
C/C++ code +
OpenMP
```
*polygeist*

```
OpenMP to
ARTS
```

```
EDT
Analysis
```

```
EDT
Optimization
```

**Create
DataBlocks**

```
Datablock
Analysis
```

```
Datablock
Optimization
```

```
ARTS
Codegen
```

ARTS to LL...

reate Events

```
ARTS lib
```

```
Output Binary
```

Analyze data access patterns inside EDTs and create Datablocks.

Input/Output
Transformations
Optimizations
Analysis
Others

4 – Data Blocks Identification

- **Identify:** Parent EDT, EDT Users
- **Aliasing Relationships:** Set of DBs that are equal or alias a given DB.
- **Child DBs:** Identify direct child DBs derived from the current DB.
- **Dependent DBs:** List all DBs that depend on the current DB.
- **Usage Count:** Track how frequently the DB is used throughout the program.
- **Attributes:** Add appropriate attributes to the EDT

Datablock Analysis

Datablock Optimization

Create Events

Input/Output
Transformations
Optimizations
Analysis
Others

5 – Data Blocks Analysis

- **Adaptive DB Configuration:** Adjust types and sizes based on data access patterns.
- **Pruning Redundant DBs:** Remove "out-only" DBs to streamline resource usage.
- **Isolated Producer Identification:** Detect producer DBs without consumers, except when embedded within an epoch.

C/C++ code + OpenMP

*polyg*

EDT

Datablock Analysis

Datablock Optimization

ARTS Codegen

ARTS to LLVM IR

Classical Optimizations

Create Events

ARTS lib

Output Binary

Input/Output
Transformations
Optimizations
Analysis
Others

6 – Data Blocks Optimization

C/C++ code +
OpenMP

*polygeist*

OpenMP to
ARTS

EDT
Analysis

EDT
Optimization

Create
DataBlocks

Datablock
Analysis

Datablock
Optimization

Create Events

Allocate, create and load events in the EDT using the DB analysis information

ARTS lib

Output Binary

Input/Output
Transformations
Optimizations
Analysis
Others

7 – Events Creation

- Common Subexpression Elimination (CSE)
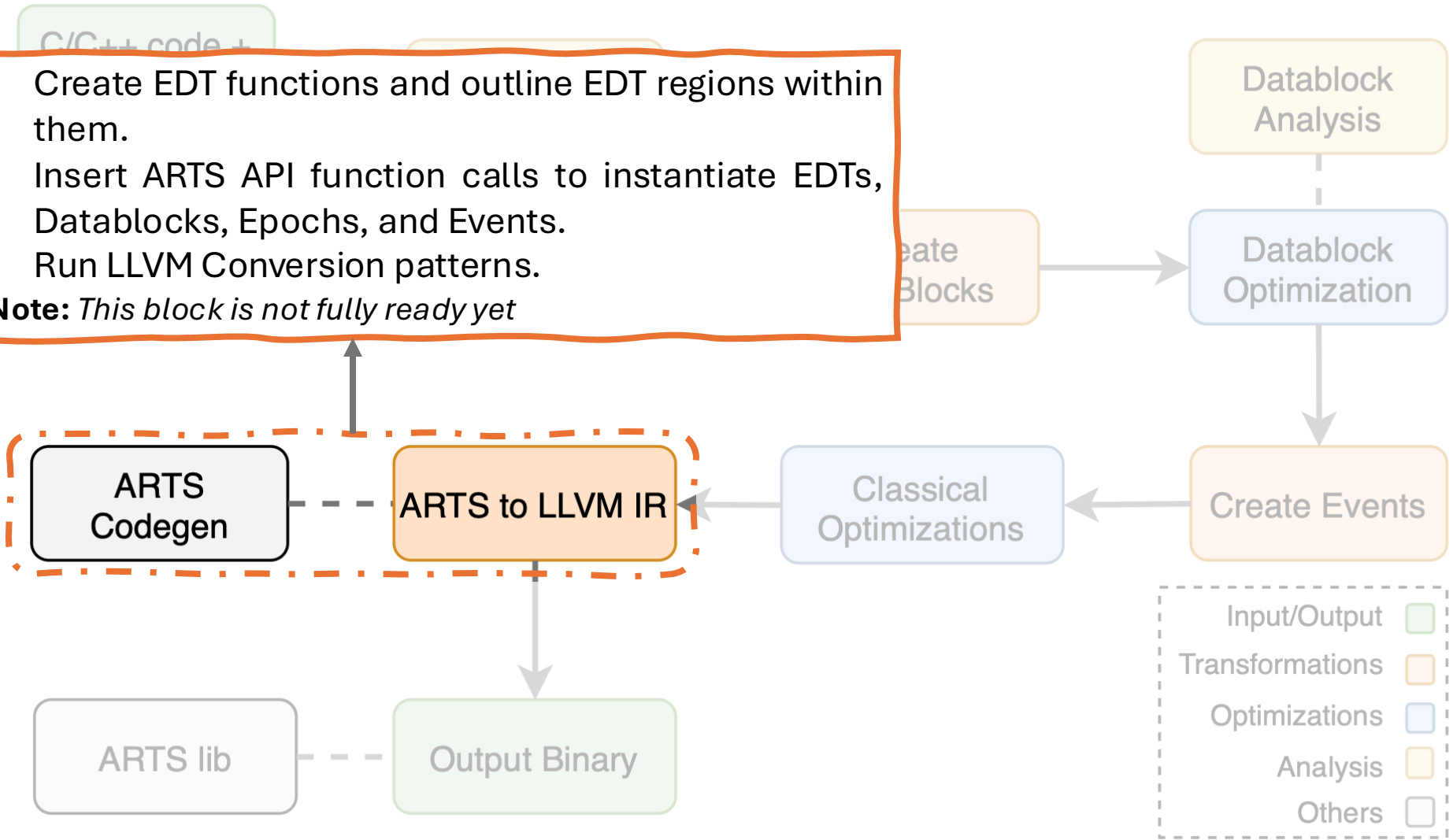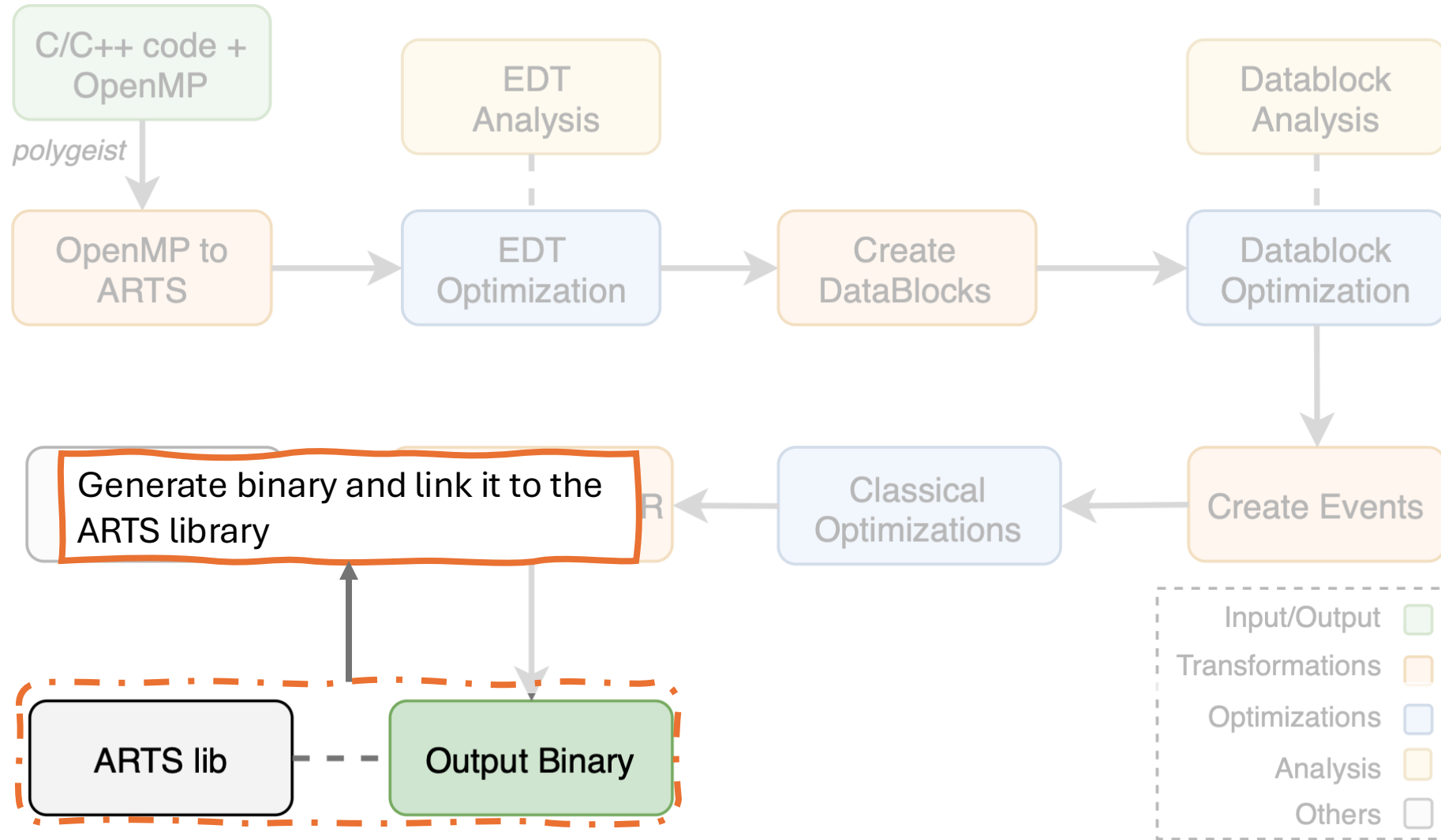- Dead Code Elimination (DCE)
- Canonicalization

8 – Other Optimizations

- Create EDT functions and outline EDT regions within them.
- Insert ARTS API function calls to instantiate EDTs, Datablocks, Epochs, and Events.
- Run LLVM Conversion patterns.

**Note:** *This block is not fully ready yet*

Datablock Analysis

Datablock Optimization

Create Blocks

ARTS Codegen

ARTS to LLVM IR

Classical Optimizations

Create Events

ARTS lib

Output Binary

Input/Output
Transformations
Optimizations
Analysis
Others

C/C++ code + OpenMP

*polygeist*

OpenMP to ARTS

EDT Analysis

EDT Optimization

Create DataBlocks

Datablock Analysis

Datablock Optimization

Create Events

Classical Optimizations

Generate binary and link it to the ARTS library

ARTS lib

Output Binary

Input/Output

Transformations

Optimizations

Analysis

Others

# Future Work

Test the infrastructure with different benchmarks.

Provide support to more OpenMP Constructs (e.g. for, barriers, locks...)

Advanced Transformation Passes.

Memory-Centric Optimizations based on a Memory cost model

Feedback-Directed Compilation.

Domain-Specific Extensions

# **Acknowledgments**

# Thank you!
*Any questions?*

**Rafael Andres Herrera Guaitero**

rafaelhg@udel.edu

Phd Candidate at University of Delaware