# Evaluating Support for OpenMP Offload Features

**UDEL:** Jose Monsalve (josem@udel.edu) , Sunita Chandrasekaran (schandra@udel.edu), Kyle Friedline (utimatu@udel.edu)

**ORNL:** Swaroop Pophale (pophaless@ornl.gov), Oscar Hernandez (oscar@ornl.gov), David Berndholt (bernholdtde@ornl.gov)
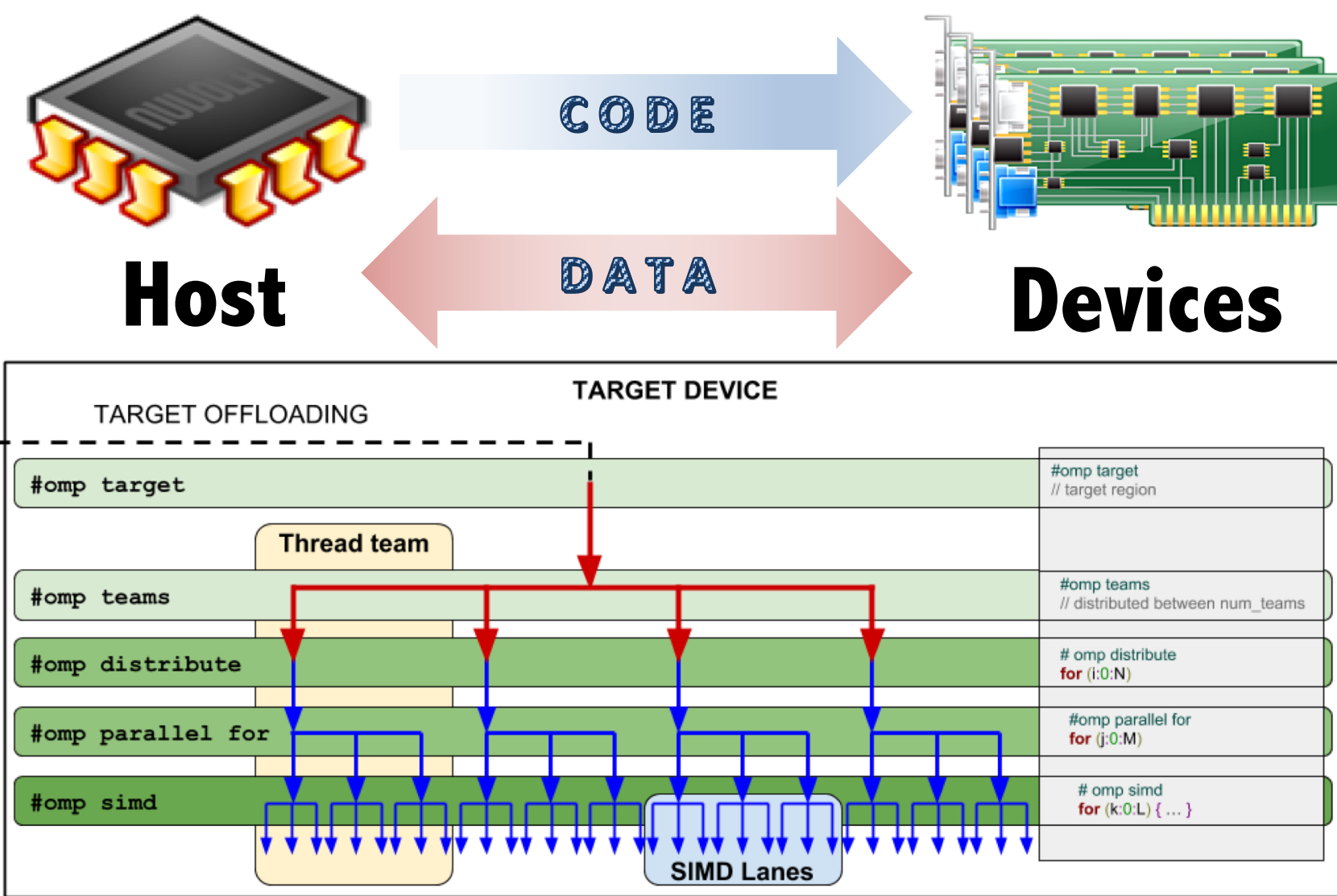
## Abstract

OpenMP has evolved to meet the rapid development in hardware platforms including heterogenous programming. DOE applications tend to push the bleeding edge of features ratified in the OpenMP specification and tend to expose the rough edges of the features' implementations. The software harness on DOE supercomputers (e.g. Titan and Summit) include Cray, Clang, Flang, XL and GCC compilers which claim partial support for the latest features in OpenMP 4.0+. This work focuses on evaluating such support across compiler implementations, focusing on OpenMP 4.5 target offload directives. Our preliminary evaluation consist of a tests suite, as well as performance comparison.

Our tests not only evaluate the OpenMP implementations but also expose ambiguities in the OpenMP 4.5 specification. We see this as a synergistic effort to help identify and correct features that are required by DOE applications and prevent deployment delays later on.

## OpenMP 4.5 offloading

**OpenMP abstract machine for offloading features is host centric:** Offloading directives hint the compiler to create device executable regions of code, as well as code and data movement between host and device.

OpenMP frees the programmer from bookkeeping data allocation and movement, as well as separate compilation of code for host and device.

OpenMP 4.5 in particular provides more control to the programmer to handle data movement between host and device.

- Target region for device code generation
- Device-host data management
- Conditional execution of code in device
- Target device selection during runtime

```
#pragma omp target map(tofrom: myVar) if(myCondition) device(2)
{   myVar ++;   }
```
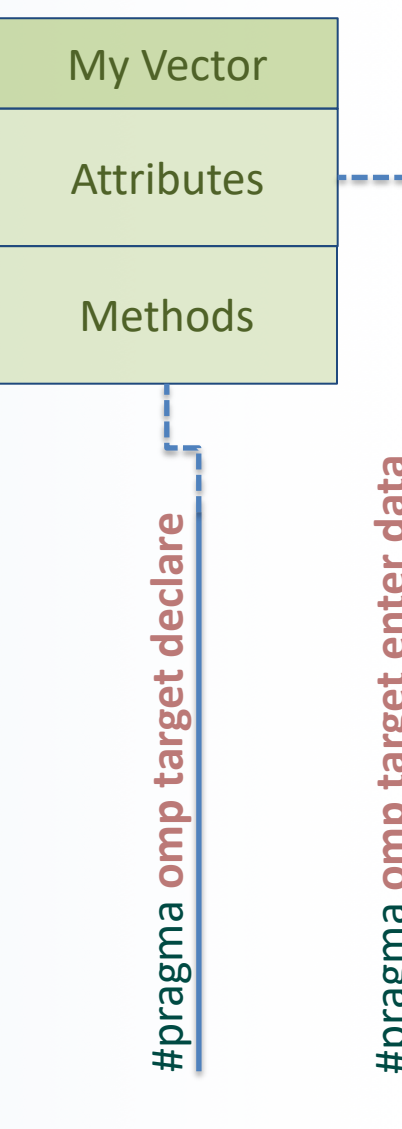


## Complex Test Cases

**Mapping Linked list to device:**
The *map_ll* function (line 6) uses *target enter data* directive to first map the head of the linked list, followed by mapping the pointer to the next node of the list and assigning it on the device. The *unmap_ll* (line 20) function explicitly copies the data using map(from:…) and *target exit data*.

**Deep Copy of Classes:**
This code came from analyzing a full scale ECP application. It uses the *declare target* directive (line 4 to 36) to ensures that procedures and global variables can be executed and data can be accessed on the device. When the C++ methods are encountered, device-specific versions of the routines are created that can be called from a target region.

Deep copy is performed through the use of target enter data (lines 43 and 44) by first mapping the class and then the individual class members. Computation is performed on the device (line 46). After computation is over, the data is copy back to the host (line 52).

My Vector — Attributes — Methods

*#pragma* **omp target declare**

*#pragma* **omp target enter data**

## Simple Test Cases

**Offloading Multiple devices:**
Each row of the matrix to each of the available devices. Use the *device* clause to select a device for data movement and computation. Target data region maps a portion of the matrix to each device (line 6). Target region does the computation (line 8).

**Mapping static attribute of a class:**
The unique value of the static VAR is default mapped inside the method of a class with a target region (lines 6-9). An OpenMP 4.5 capable compiler should capture the static variable (VAR) and map it to and from the device.

**Task dependencies:**
Task graph composed of host tasks and target tasks that have in and out dependencies between each other. Asynchronous behavior is specified using the *nowait* clause. Data map tasks are separated from computation tasks.

## Performance evaluation

### Methodology

```
OMPVV_INIT_TEST;
for  (i = 0;  i < NUM_REP;  i ++) {
    OMPVV_START_TIMER;
    #pragma omp
        OMPVV_TEST_LOAD;  // if  necessary
    OMPVV_STOP_TIMER;
    OMPVV_REGISTER_TEST;
}
OMPVV_PRINT_RESULT;
```
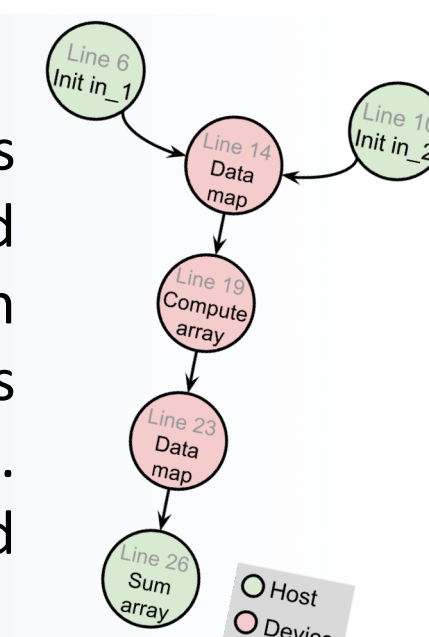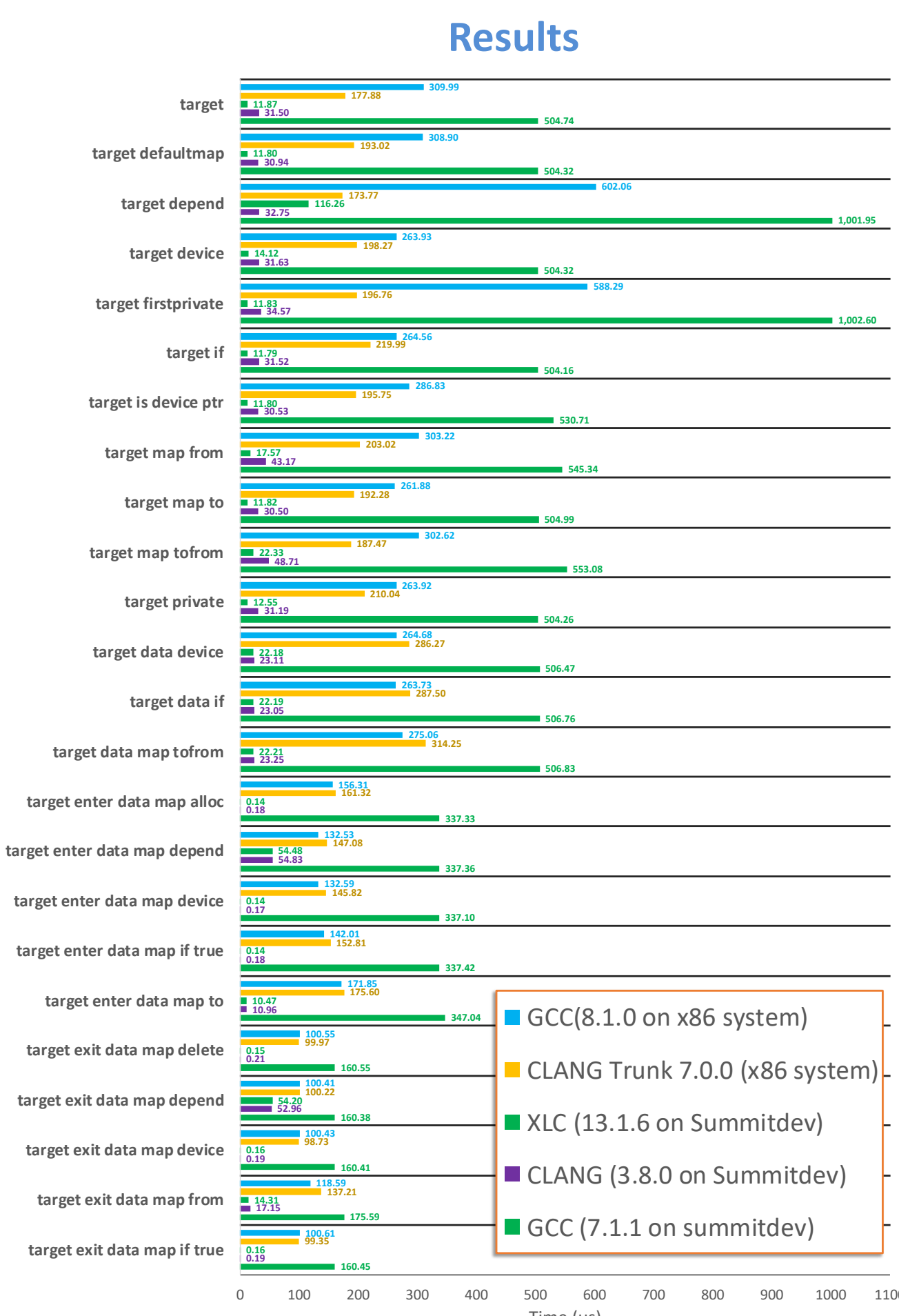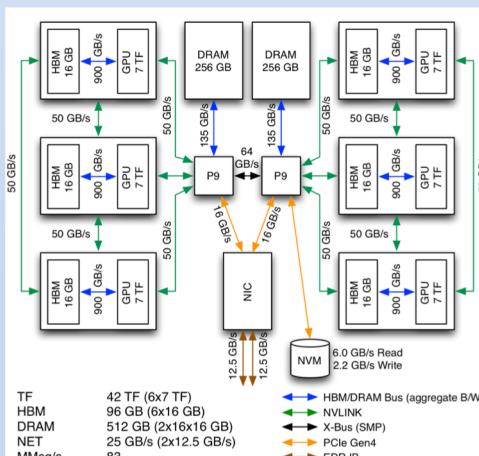
Measuring runtime overhead with multiple executions of different openMP directives and clauses. Each experiment consist of 1002 runs, removing the max and min and taking the average value.

### Testing systems

Two different systems and 5 different compilers and versions were tested. Summitdev features IBM S822LC nodes with POWER8 processors. Each node has a total of 160 hardware threads, 256 Gb of DRAM and as target devices, 4 NVIDIA Tesla P100 GPUs. The second system is an in-house cluster where each node features two Intel(R) Xeon(R) CPU E5-2670 with 32 Hardware threads, 64 Gb of DRAM and one NVIDIA K20.

### Compilers

**Summitdev:**
- GCC 7.1.1
- Clang 3.8.0
- XLC 13.1.6

**In house cluster:**
- GCC 8.1.0
- Clang 7.0.0 (Trunk version)

### Results



- GCC(8.1.0 on x86 system)
- CLANG Trunk 7.0.0 (x86 system)
- XLC (13.1.6 on Summitdev)
- CLANG (3.8.0 on Summitdev)
- GCC (7.1.1 on summitdev)

## Specification coverage evaluation

We are currently developing a test suite to asses the level of coverage of the OpenMP 4.5 specifications by the different compiler implementations. We have put together a methodology that guarantees full coverage of the specification as well as correct test implementation. We currently have released over 64 tests and we are currently in the process of releasing 33 more tests that are under review.

### Systems

| System | Model | Processors | Cores/node | Threads/node | Memory | Accelerator | Complers |
|---|---|---|---|---|---|---|---|
| Titan | Cray XK7 | AMD Opteron 6274 | 16 | 16 | 32 GB | 1 NVIDIA K20X | CCE 8.7.2 |
| Summitdev | IBM S822LC | 2x Power8 | 20 | 160 | 256 GB | 4 NVIDIA P100 | GCC 7.1.1 Clang 3.8.0 XLC 13.1.6 |
| Summit | IBM AC922 | 2x Power9 | 42 | 168 | 512 GB | 6 NVIDIA V100 | Clang 3.8.0 XLC 13.1.7 |

### Results summary

| OpenMP 4.5 Feature | GCC 7.1.1 | gfortran 7.1.1 | Summitdev Clang 3.8.0 | XLC 13.1.6 | XLF 15.1.7 | Summit Clang 3.8.0 | XLC 13.1.7 | Titan CCE 8.7.2 |
|---|---|---|---|---|---|---|---|---|
| target | 14/14 | 13/13 | 14/14 | 13/14 | 12/13 | 12/14 | 11/14 | 13/14 |
| target data | 5/6 | 4/4 | 6/6 | 6/6 | 2/4 | 6/6 | 6/6 | 3/6 |
| target enter/exit data | 6/7 | - | 6/7 | 6/7 | - | 6/7 | 6/7 | 5/7 |
| target enter data | 6/7 | - | 6/7 | 6/7 | - | 6/7 | 6/7 | 5/7 |
| target update | 5/5 | - | 5/5 | 4/5 | - | 5/5 | 4/5 | 4/5 |
| target teams distribute | 10/11 | - | 8/11 | 10/11 | - | - | - | 9/11 |
| target teams distribute parallel for | 13/14 | - | 11/14 | 11/14 | - | - | - | 10/14 |