# OvO

## OpenMP Versus Offload

Thomas Applencourt - apl@anl.gov

Argonne Leadership Computing Facility
Argonne National Laboratory
9700 S. Cass Ave
Argonne, IL 60349

# Table of contents

Argonne

# Introduction

- OpenMP versus Offload
- Not a conformance test suite!
- An OpenMP Offload Mathematical and Hierarchical Parallelism Test Suite

- 2700+ OpenMP 5.0 tests (C++ and FOR$_{mula}$TRAN$_{slation}$)[1]
- Scripts to compile, run, and check correctness
    - make
    - Bourne-Again SHell >3.2
    - Optionally $_{monthy}$Python3[2]
- *https://github.com/TApplencourt/OvO*

---

[1]More precisely C++11 and F90 standards. Some math function are defined in C++17 and C++20 standards

[2]Used to generated summaries, or re-generate tests source-code

1. OpenMP specification is extensive
2. Compilers only support a subset of it

For application developers:

- Check if a required feature is supported by a majority of compilers (else don't use it)

For compiler developers:

- Check if a required feature is supported by our compiler (else implement it)

Help:

```
1  ./ovo.sh
2  Usage:
3    ovo.sh gen
4    ovo.sh run [<test_folder>...] [--no_long] [--no_loop]
5    ovo.sh display [--detailed | --failed | --passed ] [--no_long]
   ↪ [--no_loop] [<result_folder>...]
6    ovo.sh report  [--no_long] [--no_loop]  [<result_folder>]
7    ovo.sh clean
```

# How to use OvO? "Live" Demo

```
1   $ git clone --quiet https://github.com/TApplencourt/OvO
2
3   $ CXXFLAGS='-std=c++2a' FC='gfortran' ./ovo.sh run
4   Running test_src/cpp/complex_cpp11 | Saving in ${long_path}
5   g++ -std=c++2a log_complex_float_complex_float.cpp -o
    ↪ log_complex_float_complex_float.exe
6   [...]
7   timeout -k 5s 10s ./log_complex_float_complex_float.exe
8   [...]
9
10  $ ./ovo.sh report
11  >> test_result/2020-05-12_21-02_iris15.ftm.alcf.anl.gov
12  > cpp math and complex: 259 / 259 ( 100% ) pass
13      [failures: 0 compilation, 0 offload, 0 incorrect results]
14  > cpp hierarchical parallelism: 642 / 642 ( 100% ) pass
15      [failures: 0 compilation, 0 offload, 0 incorrect results]
16  > fortran math: 79 / 79 ( 100% ) pass
17      [failures: 0 compilation, 0 offload, 0 incorrect results]
18  > fortran hierarchical parallelism: 642 / 642 ( 100% ) pass
19      [failures: 0 compilation, 0 offload, 0 incorrect results]
20  > Summary: 1622 / 1622 ( 100% ) pass
21      [failures: 0 compilation, 0 offload, 0 incorrect results]
```

Argonne

# Mathematical Functions

- Floating-Point Arithmetic is hard. I'm Not an IEEE 754 expert
- Take a best-effort approach
- Compare GPU value to assumed gold CPU[3]
- Inputs are "goods" inputs:
  - No nan, no infinities, no subnormal, …
  - Should not trigger any invalid operation, overflow, …

_____

[3]OpenMP specification doesn't define a required precision. OpenCL and CUDA® does. We choose to use a tolerance of 4 ulp, who correspond to Low Accuracy (LA) mode of Intel Math Kernel Library Vector Mathematics

Argonne

### C++

*cmath.h* (391 tests for C++11, C++17, C++20 ), *complex.h* (51 tests) and one GNU extension (*sincos*)

### FORTRAN

FORTRAN 77 Language Reference, section "Arithmetic and Mathematical Functions"[4] (67 real and 17 complex tests)

---

[4]©1999 Sun Microsystems, Inc., 901 San Antonio Road, Palo Alto, California 94303-4900 U.S.A.

|  | Hardware | C++ | FORTRAN |
|---|---|---|---|
| Cray-llvm 9.1.3 | V100 | 58% | 100% |
| Cray-Classic 9.1.3 | V100 | 49% | 100% |
| clang-AOMP 11.0.1 | Radeon VII | 70% | 44% |
| gcc 10.0.0 | P100 | 5%[5] | 59% |
| clang 11.0.0 | V100 | 70% | |
| xlC 16.1.1 | V100 | 76% | 100% |

---

[5] I guess / hope it's a problem with our installation...

[6] *--no_loop* used

# Hierarchical Parallelism

# Goals

- One loop-nest
- All the combinations of OpenMP constructs:
  - team, distribute, parallel, for, loop, simd
- All the permutation of Type:
  - Real, Complex[7]
  - Single, Double precision
- Two big categories of tests:
  1. memory copy
  2. fold (reduction, atomic)

```
1  $find ./test_src/*/hierarchical_parallelism -type f | wc -l
2  2664
```

---

[7]OvO assume implicit mapping of `std::complex`. Not part of the C++ OpenMP standard

```
1   #pragma omp declare reduction(+: complex<float>: omp_out += omp_in)
2   std::complex<float> counter{};
3   #pragma omp target teams distribute parallel for simd
    ↪   reduction(+:counter) map(tofrom:counter)
4   for (int i = 0 ; i < N*M*L ; i++)
5       counter += std::complex<float> { 1.0f };
```

## OpenMP features

- Custom reduction[8]

- In the combined statement, the *map* clause and *reduction*
  clause share a variable (OpenMP 5.0)

## OpenMP "features"

Non-trivial type $complex<float>$ is mapped

---

[8]Doesn't required in FORTRAN. Just saying that maybe one language is superior to another...

# Quaternion of Extreme: exploded atomic

```fortran
 1   REAL counter = 0.
 2   !$OMP TARGET map(tofrom:counter)
 3   !$OMP TEAMS
 4   !$OMP LOOP
 5   DO i = 1 , L
 6       !$OMP PARALLEL
 7       !$OMP LOOP
 8       DO j = 1 , N*M
 9           !$OMP ATOMIC UPDATE
10           counter = counter +  1.
11       END DO
12       !$OMP END LOOP
13       !$OMP END PARALLEL
14   END DO
15   !$OMP END LOOP
16   !$OMP END TEAMS
17   !$OMP END TARGET
```

## OpenMP Restriction

- No *!$OMP ATOMIC* in *!$OMP SIMD*

- No *!$OMP ATOMIC* with complex

```
1    double counter{};
2    #pragma omp target map(tofrom:counter)
3    #pragma omp teams
4    {
5        const int num_teams = omp_get_num_teams();
6        double partial_counter{};
7        #pragma omp parallel reduction(+: partial_counter)
8        {
9            const int num_threads = omp_get_num_threads();
10           partial_counter += double { 1.0f/(num_teams*num_threads) } ;
11       }
12       #pragma omp atomic update
13       counter += partial_counter;
14   }
```

Trivia: What is the final value of *counter*?

```fortran
1    REAL counter = 0.
2    !$OMP PARALLEL DO REDUCTION(+: counter)
3    DO i = 1 , L
4      !$OMP TARGET TEAMS DISTRIBUTE PARALLEL DO MAP(TOFROM: counter)
5      DO j = 1 , N*M
6        !$OMP ATOMIC UPDATE
7        counter = counter +  1.
8      END DO
9      !$OMP END TARGET TEAMS DISTRIBUTE PARALLEL DO
10   END DO
11   !$OMP END PARALLEL DO
```

Argonne

# Hierarchical Parallelism Result (in no particular order)[9]

|                    | Hardware   | C++  | FORTRAN |
|--------------------|------------|------|---------|
| Cray-llvm 9.1.3    | V100       | 70%  | 99%     |
| Cray-Classic 9.1.3 | V100       | 77%  | 99%     |
| clang-AOMP 11.0.1  | Radeon VII | 50%  | 47%     |
| gcc 10.0.0         | P100       | 62%  | 81%     |
| clang 11.0.0       | V100       | 85%  |         |
| xlC 16.1.1         | V100       | 73%  | 100%    |

---

[9] *--no_loop* used

Argonne

# Conclusion

- OvO (*https://github.com/TApplencourt/OvO*)
- Nobody is perfect
- but FORTRAN is more perfect than others [10]

---

[10]Cray and xlC have near a 100% pass rate. And FORTRAN have native complex support