# Use case: Optimizing the Weather Research and Forecasting Model (WRF) with OpenMP Offload and Codee

Namo Wichitrnithed
Oden Institute, UT Austin

**NERSC**
Woo-Sun Yang
Helen He
Brad Richardson

**Pacific Northwest National Laboratory**
Koichi Sakaguchi
William I. Gustafson Jr.

**Appentra Solutions S.L.**
Manuel Arenaz
Ulises Costi Blanco
Alvaro Goldar Dieste

**The Hebrew University of Jerusalem**
Jacob Shpund

# The Weather Research & Forecasting Model

- An atmospheric model written in Fortran written in the 1990's by a number of organizations such as the National Center for Atmospheric Research (NCAR) and the National Centers for Environmental Prediction (NCEP)

- Solves the 3D Euler equations using finite differences and explicit timestepping

- Used in both research and operational, real-time forecasting worldwide

- NERSC development branch: https://github.com/NERSC/WRF

# Optimization goals

- Current parallelism: domain decomposition (MPI) into patches (ims:ime, kms:kme, jms:jme) and shared memory (OpenMP) among tiles (its:ite, kts:kte, jts:jte)

- MPI + GPU approach: offloading work from each patch to a GPU

- Programming workflow
  - Profilers (gprof, perftools, Nsight)
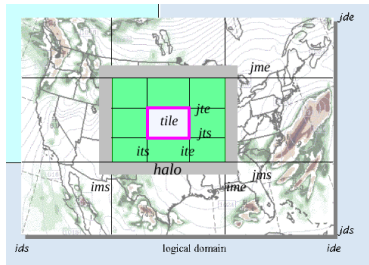  - Static code inspection (Codee)



Figure 1: WRF decomposition layer. Image from Dudhia, J. "WRF Modeling System Overview".

# Fast Spectral-Bin Microphysics (FSBM)

- Particle size spectrum divided into 33 intervals (bins)

- Computations required for each particle type and size at each grid point

- Require small timesteps (5-10 s)

- Universal; can be used for different atmospheric phenomena

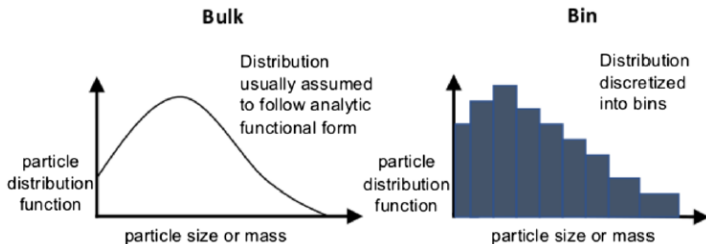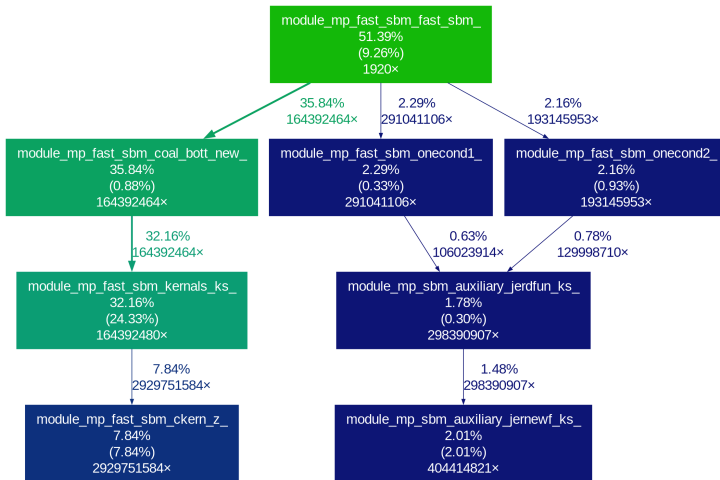- Current version in WRF: FSBM-2 (Shpund et al., 2019)



**Bulk**

Distribution usually assumed to follow analytic functional form

particle distribution function

particle size or mass

**Bin**

Distribution discretized into bins

particle distribution function

particle size or mass

Figure 2: Image from Morrison et al., 2020.

## Test case setup on Perlmutter

- Conus-12km test case
  - 425 × 300 × 50 grid
  - One-day restart
  - Time step: 5s
- Compilers: PrgEnv-nvidia/8.5.0 (NVHPC 23.9)
  - nvfortran, nvc, nvc++
  - Good GPU support for OpenMP, OpenACC, CUDA
  - GPU flags:  -mp=gpu -target-accel=nvidia80
- WRF configure option: 4 (dm+sm) PGI (pgf90/gcc)

# Finding time-consuming routines with Gprof

# Inside the FSBM routine

- Subroutine Fast_SBM() in phys/module_mp_fast_sbm.F coal_bott_new()

```fortran
do j = jts:jte
  do k = kts:kte
    do i = its:ite
      ! Collision-Coalesence process
      call COAL_BOTT_NEW(...)

      ! do stuff
    enddo
  enddo
enddo
```

# OpenMP GPU offloading

- A set of directives for C and Fortran that let the compiler generate GPU code
- Can manage parallelism and data transfer like CUDA API
- Directive-based: more portable and easier to port existing code to GPU, but less control

```
#pragma omp target teams
distribute parallel for
num_teams(3)
for (int i = 0; i < 12; ++i)
{
    C[i] = A[i] + B[i];
}
```
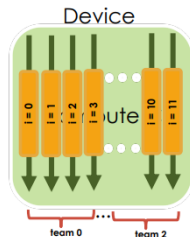


Figure 3: Image from
https://www.olcf.ornl.gov/
wp-content/uploads/2021/08/
ITOpenMP_Day1.pdf

# Inside the FSBM routine

- Subroutine Fast_SBM() in phys/module_mp_fast_sbm.F

- Parallelization granularity: number of grid points, assuming no race conditions inside coal_bott_new()

```
1    do j = jts:jte
2      do k = kts:kte
3        do i = its:ite
4          ! Collision-Coalesence process
5          call COAL_BOTT_NEW(...)
6
7          ! do stuff
8        enddo
9      enddo
10   enddo
11
12
```

## Inside the Kernals_KS subroutine

- Global collision tables, e.g. `cwlg`, `cwls` are being modified at each grid point `(i,k,j)`

```
1    do n = 1,33
2      do m = 1,33
3        ckern_1 = ...
4        ckern_2 = ...
5        ! water - graupel
6        cwlg(m,n) = (ckern_2 + (ckern1-ckern_2* ..)) * ...
7
8        ckern_1 = ...
9        ckern_2 = ...
10       ! water - snow
11       cwls(m,n) = (ckern_2 + (ckern1-ckern_2* ..)) * ...
12
13       ! 18 more arrays
14
15      enddo
16    enddo
17
```

# Setting up Codee for WRF

```
1    # Capture compilation flags in JSON file
2    bear -- ./compile -j 8 wrf
3
4    # Initial screening report
5    codee screening --config compile_commands.json
6
7    # Checks report
8    codee checks --config compile_commands.json
9
10   # Example: in-place OpenMP offload insertion
11   codee rewrite --offload omp --in-place \
12     module_mp_fast_sbm.f90:6293:4 \
13     --config compile_commands.json
14
```

# Codee analysis of Kernals_KS

- Codee implies there are no loop-carried dependencies, so the individual entries can be computed independently

```
1     ! Codee: Loop modified
2     !$omp target teams distribute parallel do &
3     !$omp private(n) map(from: cwlg, cwls, ...) ...
4     do n = 1,33
5       ! Codee: Loop modified
6       !$omp simd
7       do m = 1,33
8         ckern_1 = ...
9         ckern_2 = ...
10        cwlg(m,n) = (ckern_2 + (ckern1-ckern_2* ..)) * ...
11
12        ckern_1 = ...
13        ckern_2 = ...
14        cwls(m,n) = (ckern_2 + (ckern1-ckern_2* ..)) * ...
15
16        ! 18 more arrays
17      enddo
18    enddo
19
```

# Removing the global arrays

- Replace looking up m, n entry with computing as needed

```
1     pure real function get_cwlg(..., m, n)
2     pure real function get_cwls(..., m, n)
3
```

- No more shared arrays between grid points
- Speedup: around 1.4x
  - A lot of collision types are not used in FSBM
  - Not every entry m,n are used

## Offloading the main loop

- Each grid point can now be assigned to a thread

- Further memory optimization allows a full collapse(3)

```fortran
1    !$omp target teams distribute parallel do collapse(3)
2    do j = jts:jte
3      do k = kts:kte
4        do i = its:ite
5          ! Collision - Coalesence
6          call COAL_BOTT_NEW(...)
7
8          ! do stuff
9        enddo
10      enddo
11    enddo
12
```
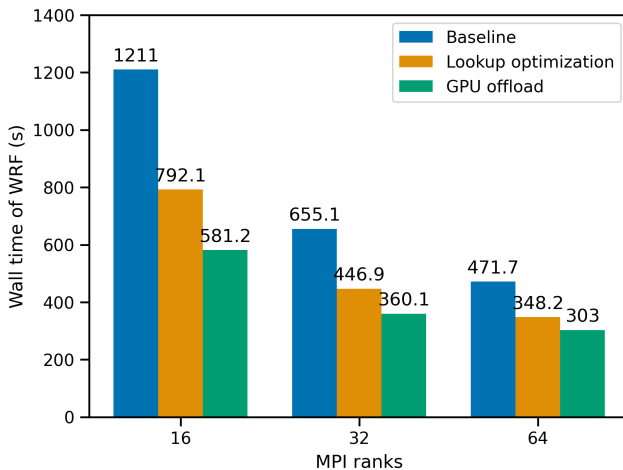
U.S. DEPARTMENT OF ENERGY | Office of Science

# Speedup results

- 10-minute runs
- 1 OpenMP thread per MPI rank
- 1 GPU per MPI rank

| Routine | Total speedup |
|:---:|:---:|
| coal_bott_new loop | 66.6x |
| fast_sbm | 2.99x |
| Overall | 2.20x |

# Strong scaling

- No. of GPUs is fixed to 16, and no. of MPI ranks is varied from 16 to 64

# Summary

- Accelerated a big part of the FSBM routine to GPUs through loop restructuring and OpenMP device offload
  - Codee's dependency analysis functionality exposed independence between computations among different grid points
- Achived an overall speedup of 2.2x for the 1 GPU per rank case for the CONUS-12km case
- A combination of runtime profiling and static code analysis is a very helpful aid in optimization efforts, especially for those not fully familiar with the context of the code

QUESTIONS