



# GPU PERFORMANCE STUDY FOR THE WRF MODEL ON THE SUMMIT SUPERCOMPUTER

Jeff Adie, *NVIDIA*, Gökhan Sever, Rajeev Jain, *DOE Argonne NL*, and Stan Posey, *NVIDIA*

# EARTH SYSTEM NEEDS KEEP GROWING

“Future Earth system models will need over 1000 times the computational power of today” - ESPC Position Paper on HPC Needs, 2017<sup>1</sup>

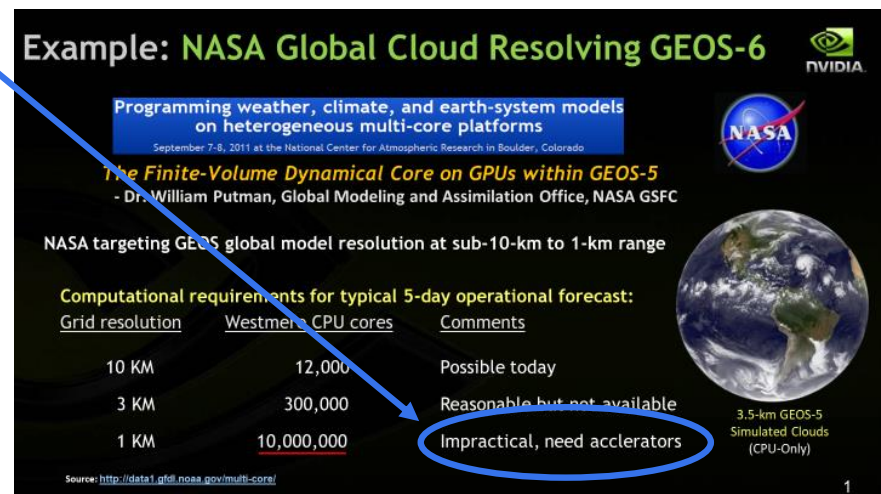
	Observations		Models	
	Volume	Type	Volume	Type
Today:	20 million = $2 \times 10^7$	98% from 60 different satellite instruments	5 million grid points 100 levels 10 prognostic variables = $5 \times 10^8$	physical parameters of atmosphere, waves, ocean
Tomorrow:	200 million = $2 \times 10^8$	98% from 80 different satellite instruments	500 million grid points 200 levels 100 prognostic variables = $1 \times 10^{11}$	physical and chemical parameters of atmosphere, waves, ocean, ice, vegetation

Factor 10 per day

Factor 2000 per time step

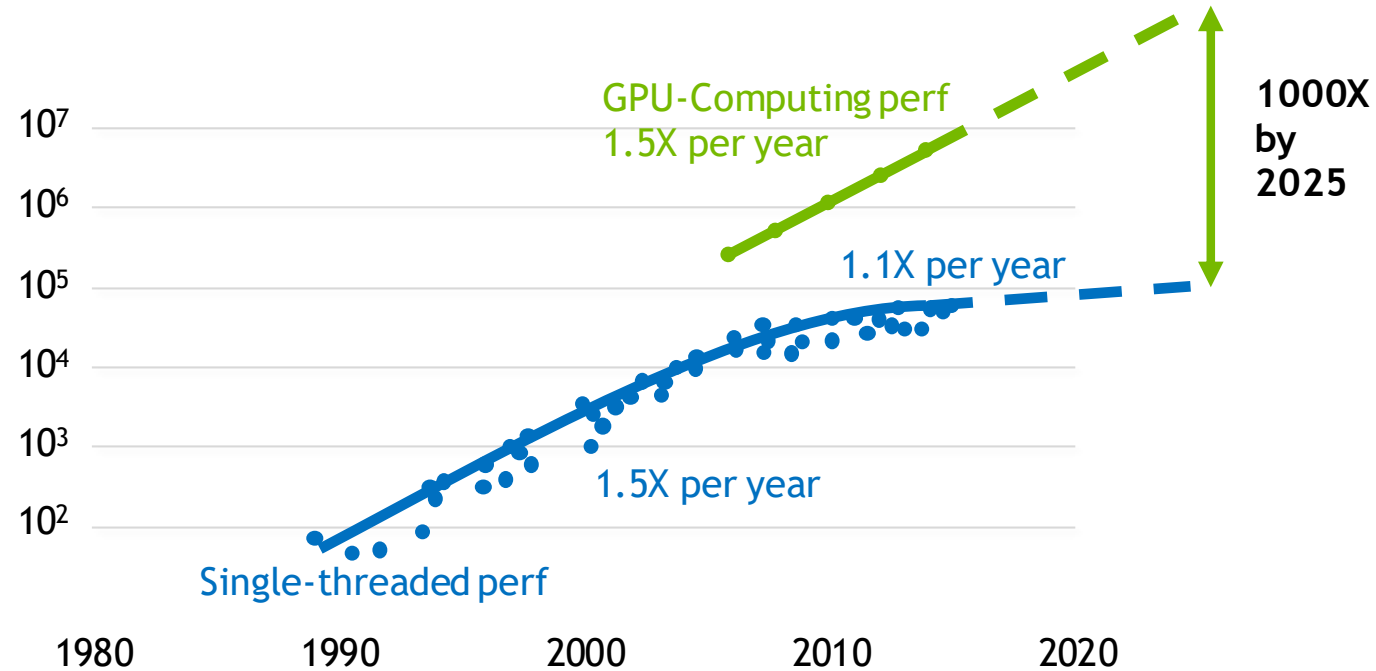
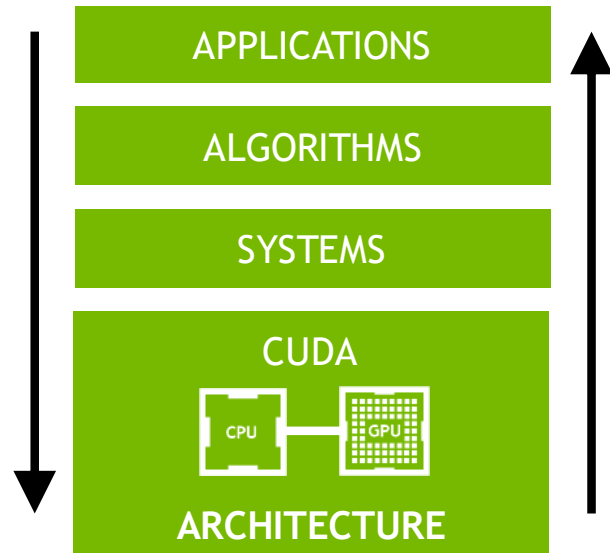
ECMWF - 2000x needed by 2032<sup>2</sup>

NASA - 10M CPU cores - “impractical”<sup>3</sup>



- <https://doi.org/10.7289/V5862DH3> 3,4 – refer to AVEC Benchmark report
- Challenges of Getting ECMWF's Weather Forecast Model (IFS) to the Exascale – G. Mozdzynski, ECMWF 16th HPC Workshop
- <http://data1.gfdl.noaa.gov/multi-core/>

# RISE OF GPU COMPUTING

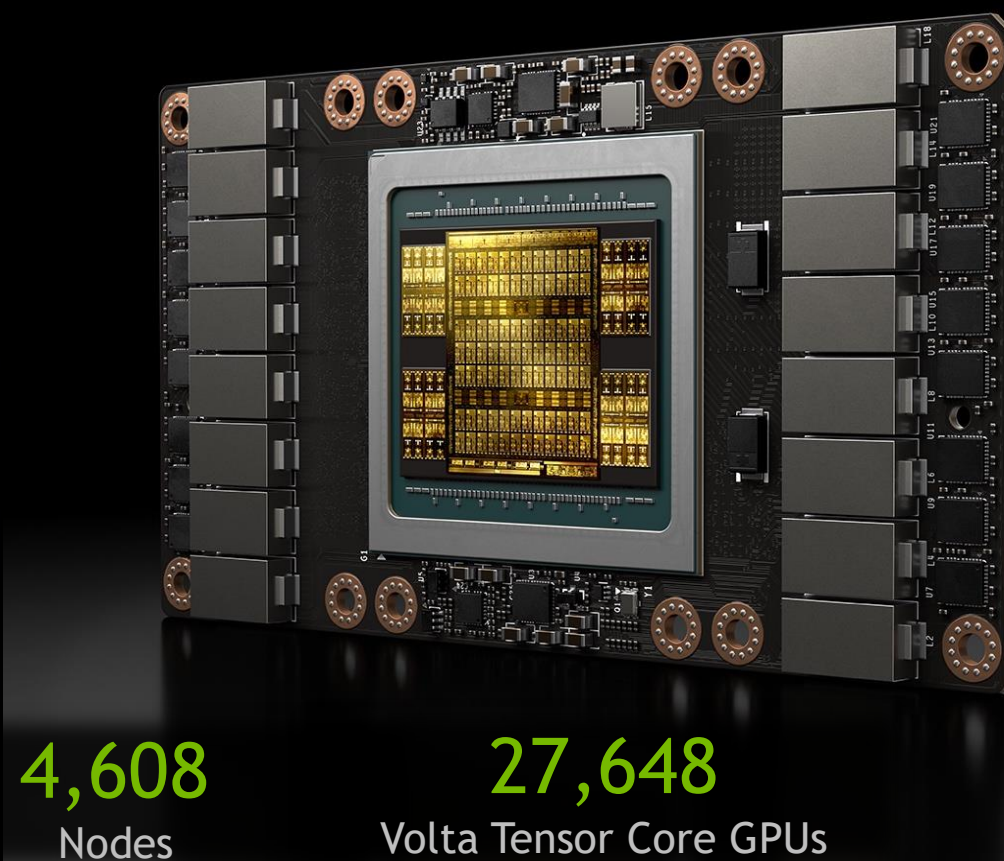


Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2015 by K. Rupp



# SUMMIT - WORLD'S FASTEST SUPERCOMPUTER

Summit Becomes First System To Scale The 100 Petaflops Milestone

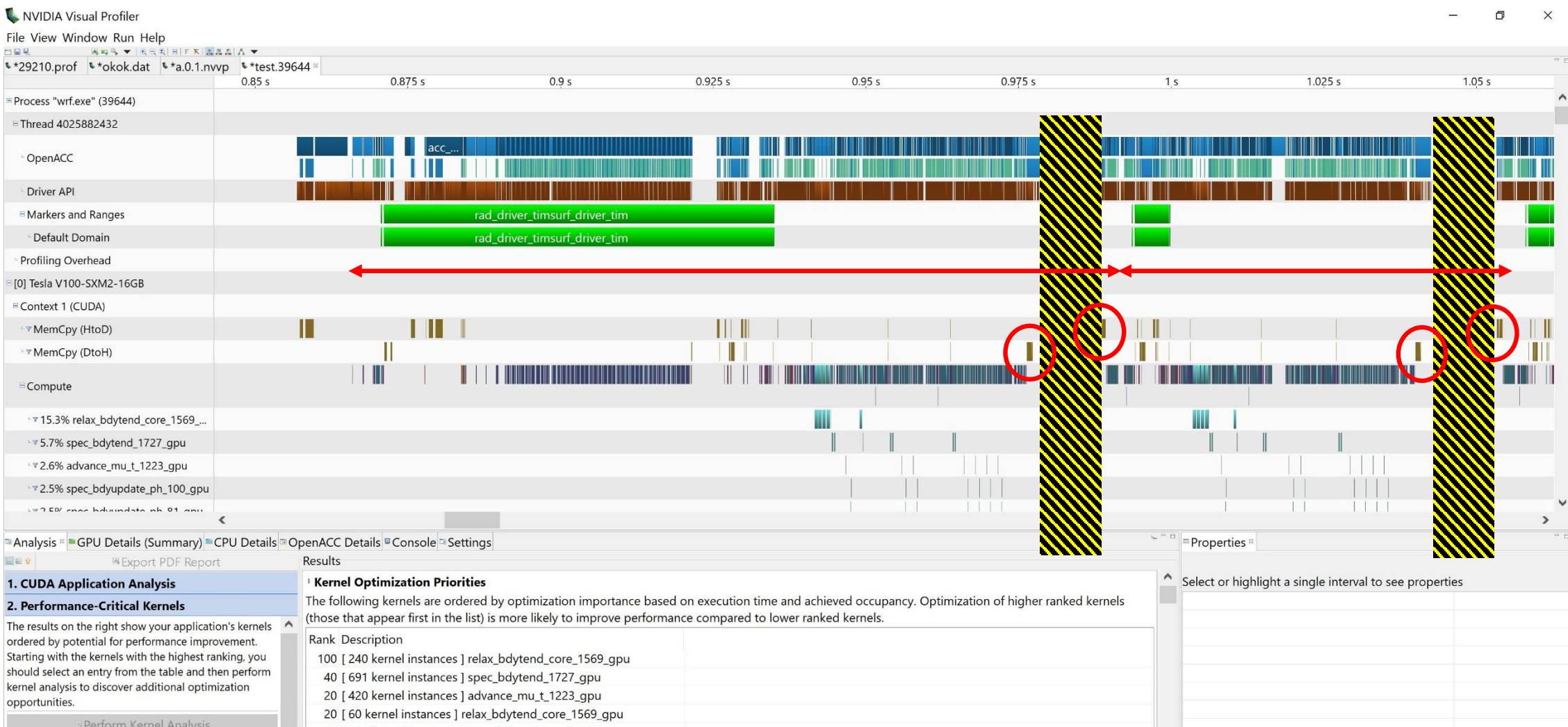


# ACCELERATED COMPUTING CHALLENGES

Not quite plug-n-play!

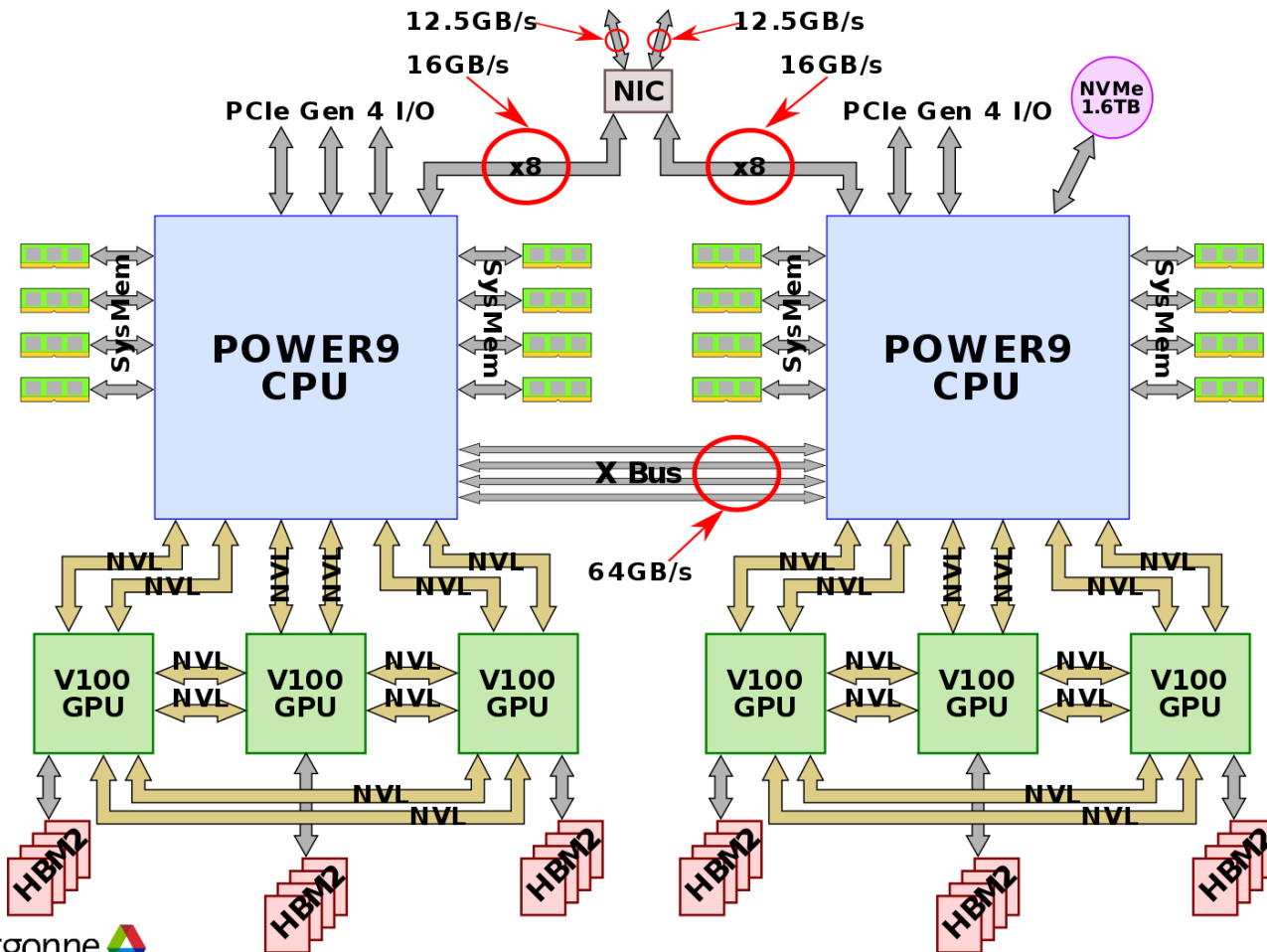
- Software needs to support accelerator (or be ported to accelerator)
  - OpenACC provides an ideal mechanism for portable, easy porting
- CPUs and accelerators have different physical memories. Data transfer has to occur, even if it is hidden from the programmer (CUDA unified memory, OpenACC managed data, CAPI, etc)
  - Minimise memory transfers (keep data on GPU for reuse wherever possible)
  - Faster memory transfers - NVLink

# ACCELERATED COMPUTING CHALLENGES





# SUMMIT NODE ARCHITECTURE



# WRF PERFORMANCE STUDY

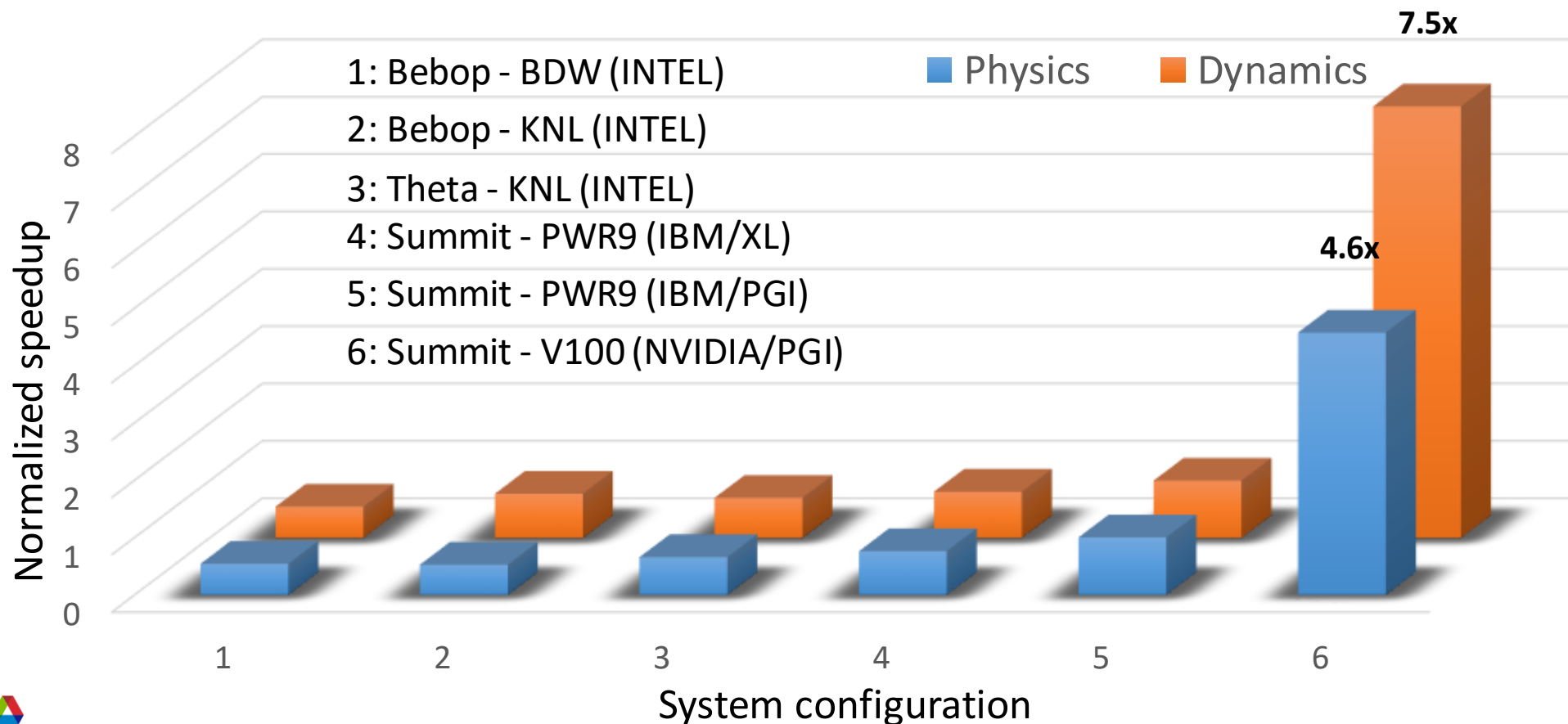
## Methodology

- Ported NVIDIA OpenACC WRF code (based on WRF V3.7.1) to POWER9
  - PGI 19.1 Compiler
  - IBM Spectrum MPI (PAMI disabled)
- Model testing using em\_les ideal case for dynamics only (384x384x42, ~6M cells)
  - Scaled vertically for GPU occupancy (384x384x252, ~37M cells)
  - Scaled horizontally for multi-node studies (3840x3840x252, ~3.7B cells)
  - Full model physics test with Thompson MP, RRTM/Dudhia, YSU PBL, Revised MM5+TDS



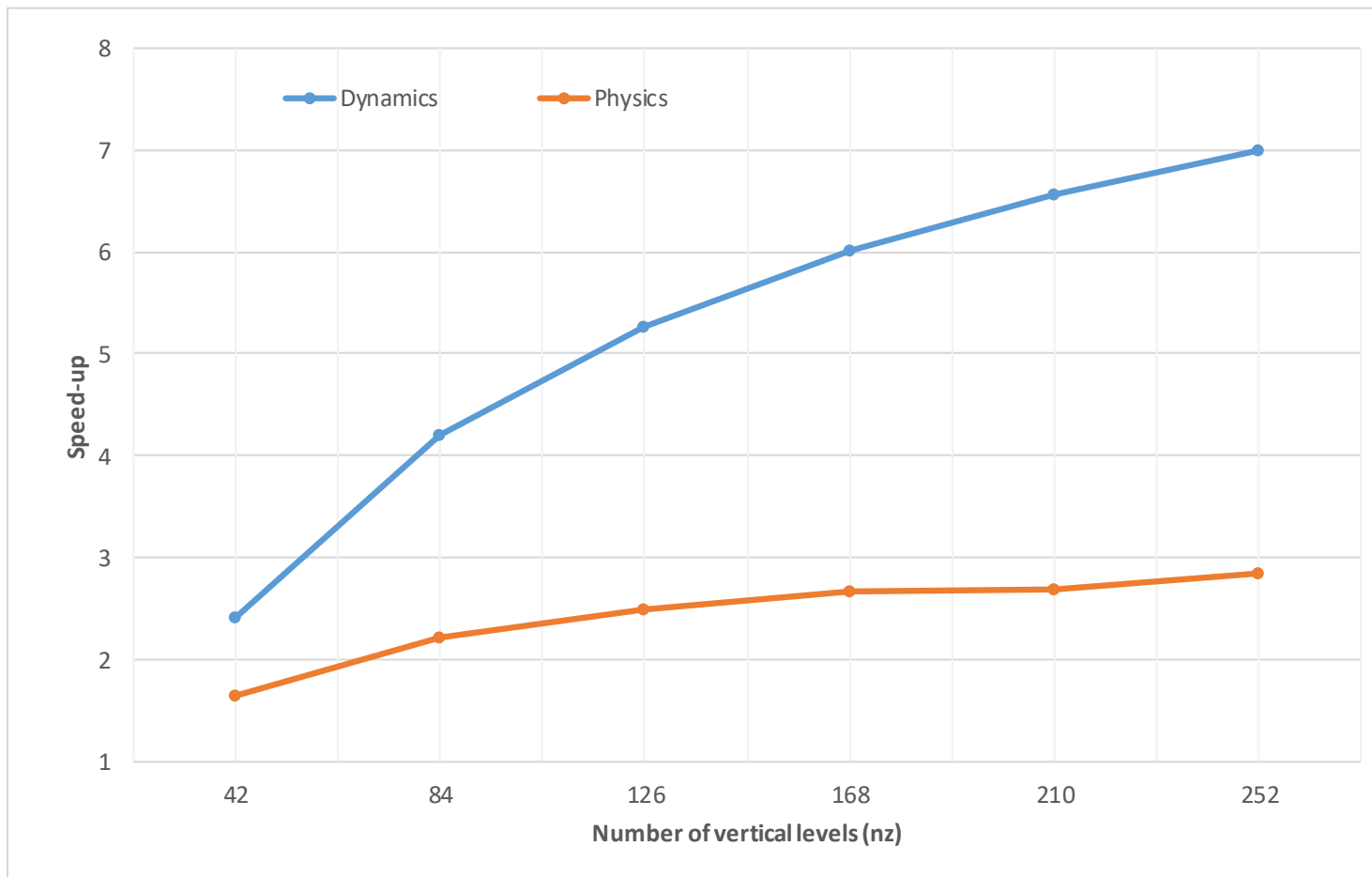
# WRF SINGLE NODE PERFORMANCE

Benchmarks of single-node setups on DOE clusters



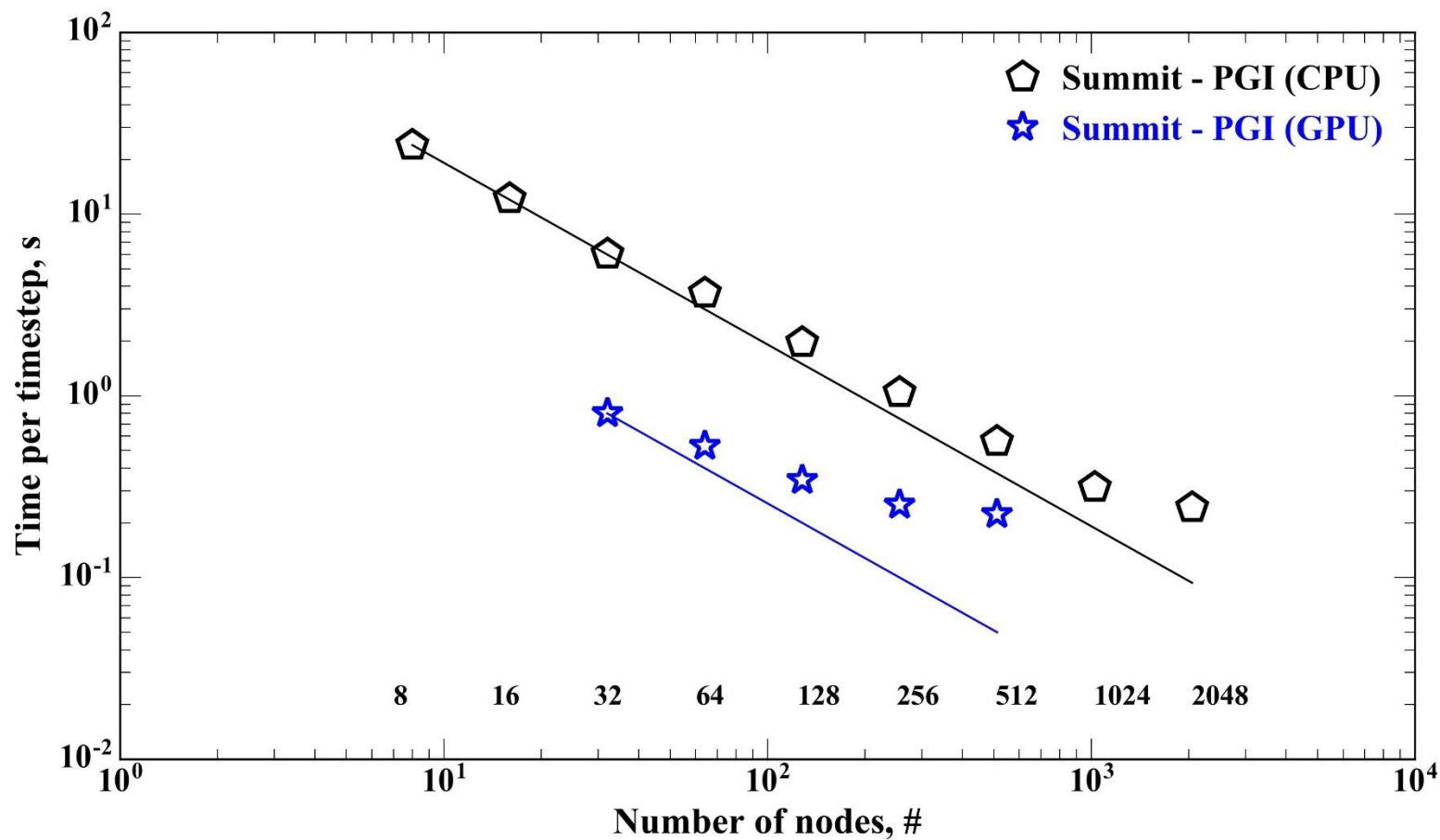
# WRF SINGLE NODE PERFORMANCE

Speedup sensitivity to amount of work



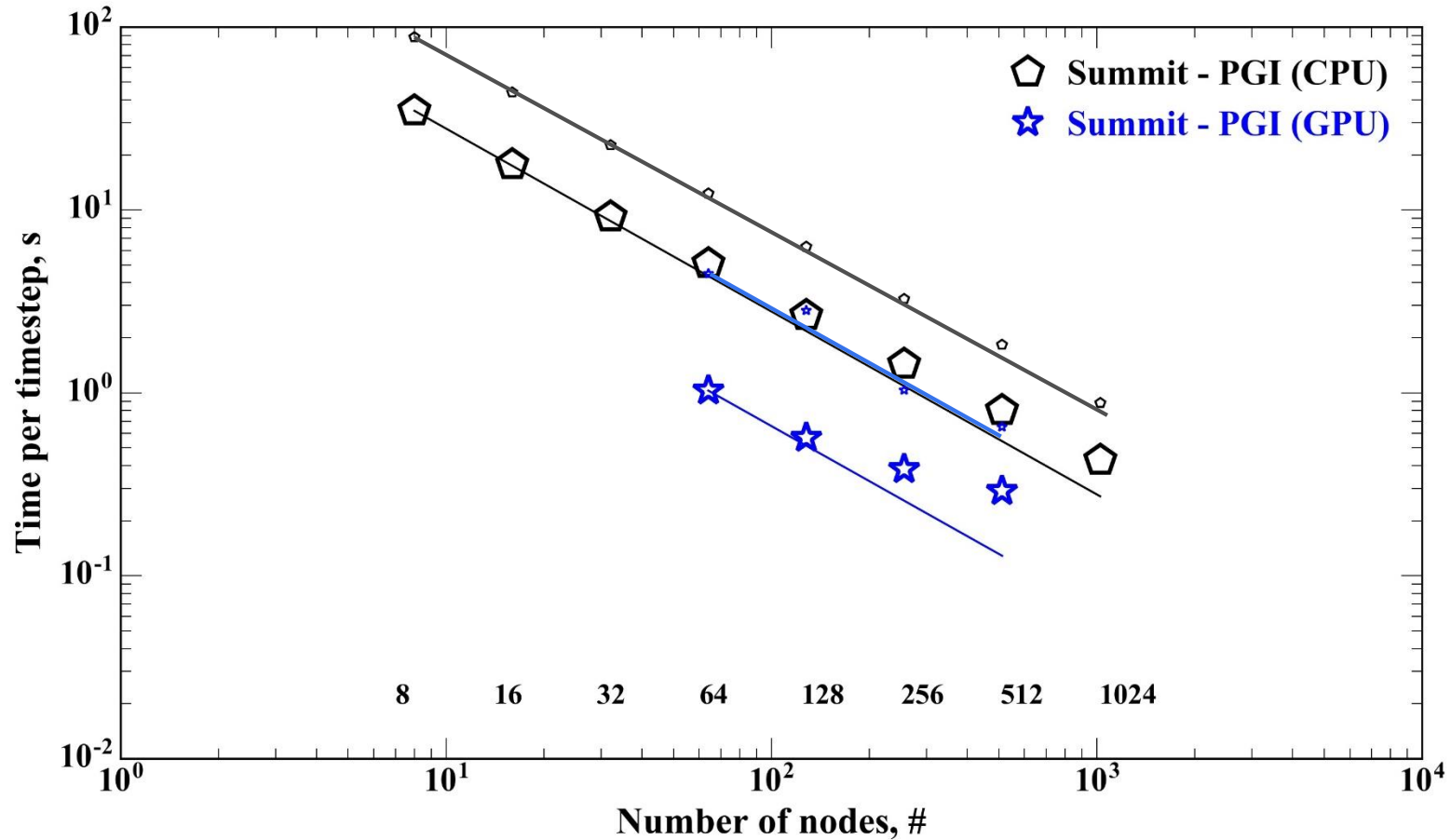
# WRF SCALING

Dynamics only, 4B cells



# WRF SCALING

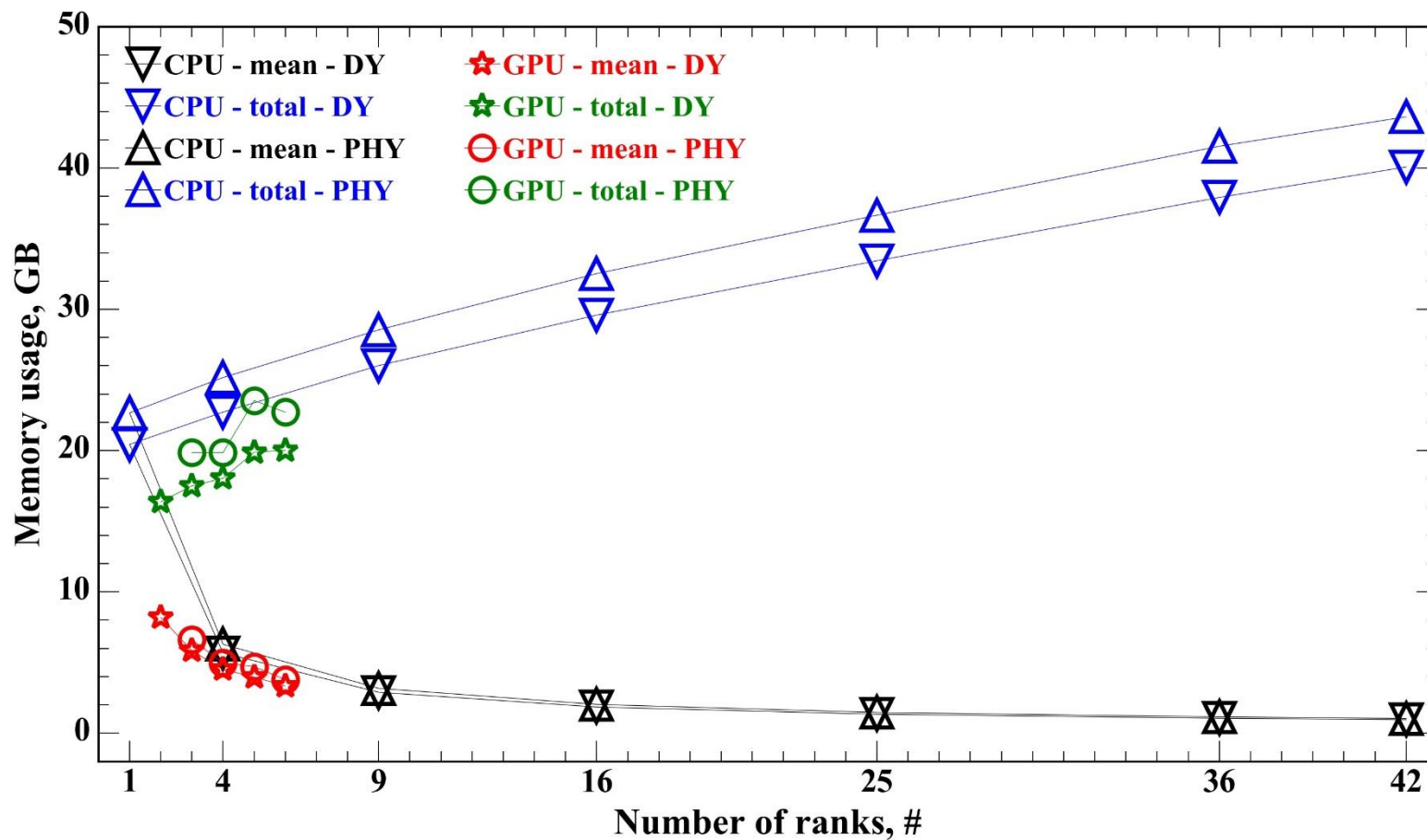
Full Model, 4B cells





# WRF MEMORY USAGE

## Single Node



# WRF PERFORMANCE STUDY

## Preliminary Conclusions

- Single node shows 7.5x speedup on Dycore, 5x on full model with Physics
- Multi node scaling tested to 512 nodes (3,072 GPUs)
  - Scaling above 64 nodes (384 GPUs) limited by model size (GPU Occupancy)
- Currently not using CAPI between CPUs and GPUs
  - GPU memory limits model size
  - Cannot use PAMI for MPI acceleration
  - Requires code changes to the RSL\_LITE comms

# WRF PERFORMANCE STUDY

## Future Work

- Perform strong and weak scaling tests
- Get CAPI working
  - Remove GPU Memory size limitation
  - Allow for PAMI support for faster communications
- Port UCM (Urban Canopy Model) to GPU
- Look at supporting later WRF versions (V4)

# WRF PERFORMANCE STUDY

## Acknowledgements

This research used resources of the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725.

The computational resources on Summit for the account # CSC307 are provided via the UrbanECP project

Thanks to Bob Walkup (IBM) and Todd Hutchinson (TWC/IBM) for help with the WRF configuration for POWER.

We appreciate the help from the OLCF support team for various technical problems.

Akira Kyle's script (wrf\_stats.py) is utilized to extract average benchmark timings.



