

# Parallel Multivariate Spatio-Temporal Clustering of Large Ecological Datasets on Hybrid Supercomputers

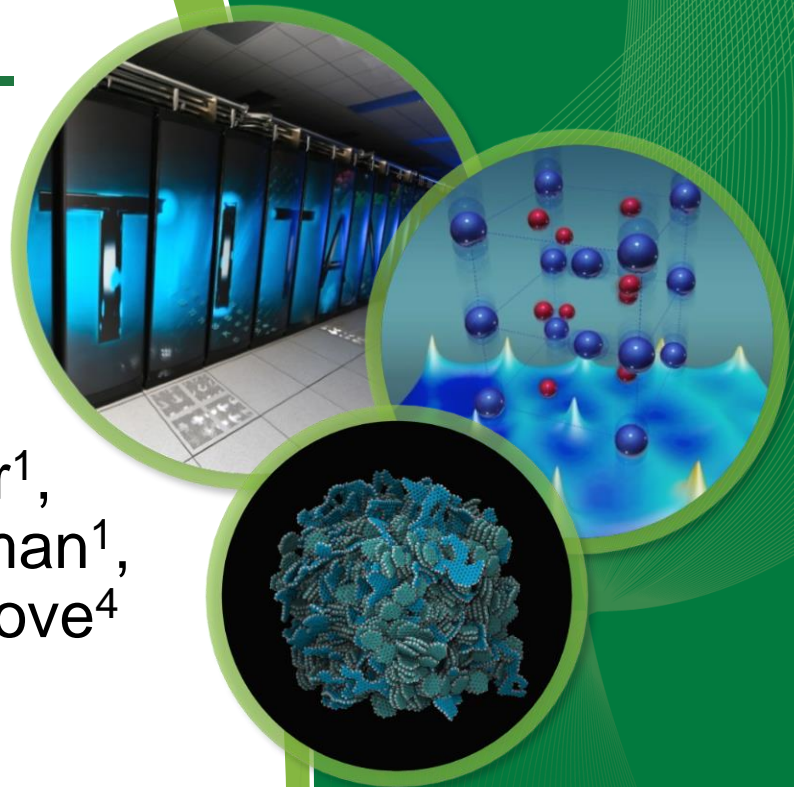
**Sarat Sreepathi**<sup>1</sup>, Jitendra Kumar<sup>1</sup>,  
Richard T. Mills<sup>2</sup>, Forrest M. Hoffman<sup>1</sup>,  
Vamsi Sripathi<sup>3</sup>, William W. Hargrove<sup>4</sup>

<sup>1</sup>Oak Ridge National Laboratory

<sup>2</sup>Argonne National Laboratory

<sup>3</sup>Intel Corporation

<sup>4</sup>USDA Forest Service



# Motivation

- Rapid proliferation of data in various domain sciences
- Earth Science
  - Advanced sensors – high fidelity data
  - Remote Sensing Platforms
    - Satellites
    - Unmanned Aircraft Systems (UAS)
    - Airborne systems
  - Observational Facilities
- Critical need for High Performance Big Data Analytics

# Applications

- Vegetation mapping and characterization
- Development of ecoregions
- Species distribution
  
- Climate zone classification
- Understand climate regime changes in future
  - Under various predicted climate change scenarios

# Parallel k-means (Baseline)

- Goal: Divide observations into k clusters
- Centralized Master-Worker paradigm
- Pick initial centroids
- Iterative method
- Workers
  - Compute distances
  - Update centroids and cluster assignments
  - Repeat till convergence is achieved
- Typical target convergence:  $< 0.5\%$  changes

# Datasets

## Great Smoky Mountains National Park (GSMNP)

- Airborne multiple return Light Detection and Ranging (LiDAR) data
  - Vertical canopy structure of the vegetation
  - 30 *m* × 30 *m* spatial resolution horizontal grid
  - 1 *m* vertical resolution to identify vegetation height from the ground surface

## Global Climate Regimes

- Bioclimatic (BioClim) data for the contemporary period
- Climate models from IPCC Third Assessment Report (CMIP3) – Parallel Climate Model (PCM) and HadCM3 model
- Two different emissions scenarios:
  - B1 (lower emissions), A1FI (high emissions)

# Global Climate Regimes: Variables

TABLE II  
VARIABLES USED FOR DELINEATION OF GLOBAL CLIMATE REGIMES.

Variable Description	Units
<b>Bioclimatic Variables</b>	
Precipitation during the hottest quarter	mm
Precipitation during the coldest quarter	mm
Precipitation during the driest quarter	mm
Precipitation during the wettest quarter	mm
Ratio of precipitation to potential evapotranspiration	–
Temperature during the coldest quarter	°C
Temperature during the hottest quarter	°C
Day/night diurnal temperature difference	°C
Sum of monthly $T_{avg}$ where $T_{avg} \geq 5^{\circ}\text{C}$	°C
Integer number of consecutive months where $T_{avg} \geq 5^{\circ}\text{C}$	–
<b>Edaphic Variables</b>	
Available water holding capacity of soil	mm
Bulk density of soil	$\text{g}/\text{cm}^3$
Carbon content of soil	$\text{g}/\text{cm}^2$
Nitrogen content of soil	$\text{g}/\text{cm}^2$
<b>Topographic Variables</b>	
Compound topographic index (relative wetness)	–
Solar interception	$(\text{kW}/\text{m}^2)$
Elevation	m

# Datasets

## DESCRIPTION OF DATA SETS USED IN THE CURRENT STUDY

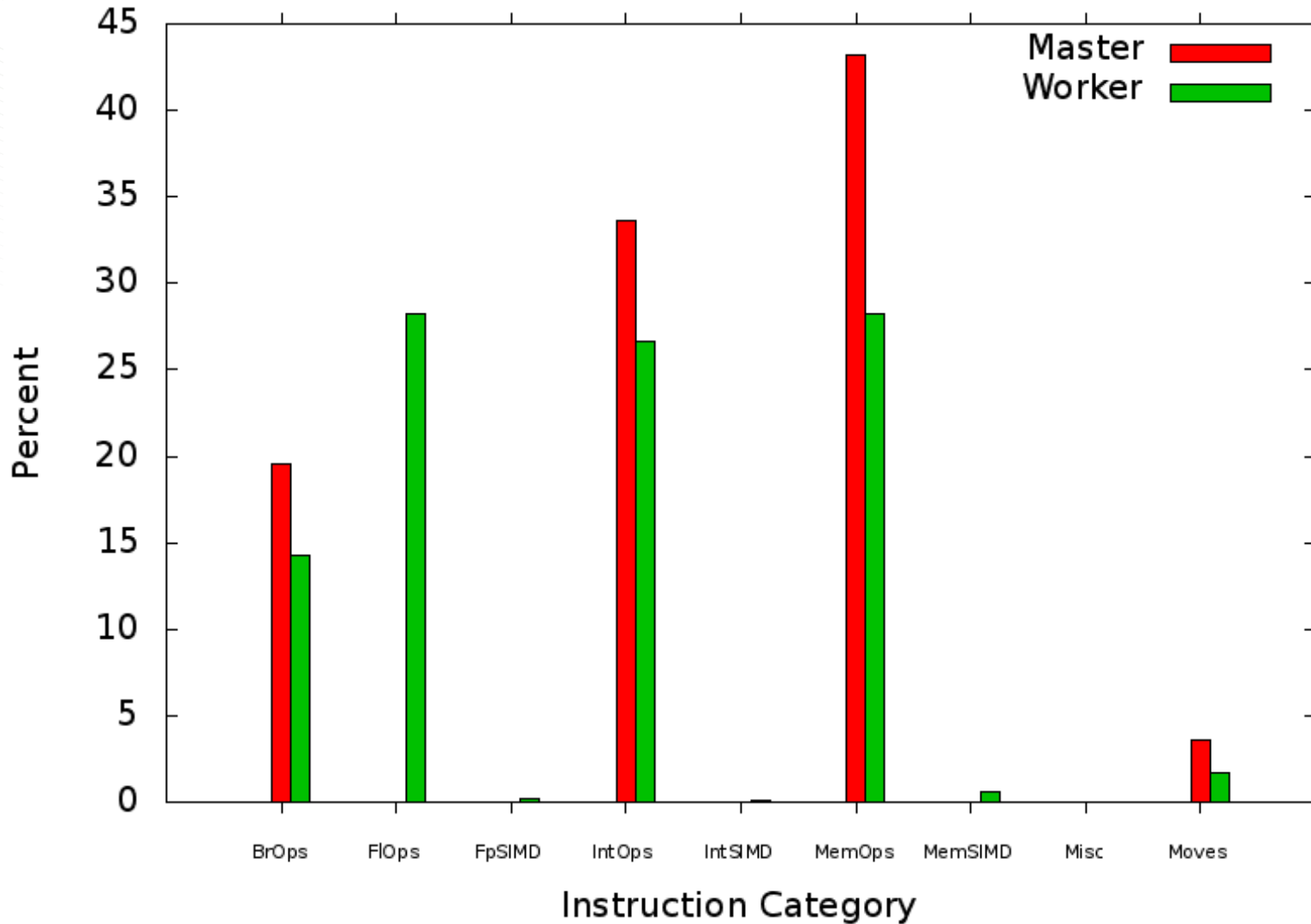
Description	Dimensions	Size
GSMNP LiDAR	$3,186,679 \times 74$	900 MB
CMIP3 Climate States	$123,471,198 \times 17$	7.9 GB

## Preprocessing

- Standardized the data set along each dimension
  - A mean of zero and standard deviation of one
- Allowing every dimension to be equally and fairly represented in the clustering algorithm

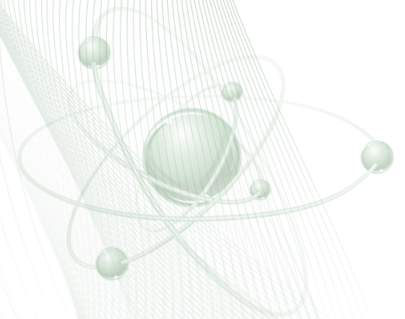
# Application Characterization: Baseline k-means

Instruction mix of baseline clustering application





# Methodology



# BLAS Formulation (Application Phase 1)

Squared Euclidean Distance:  $\mathbf{dist}_{i,j} = \|\mathbf{obs}_{i,*} - \mathbf{cent}_{i,*}\|^2$

Binomial expansion:  $\mathbf{dist}_{i,j} = \|\mathbf{obs}_{i,*}\|^2 + \|\mathbf{cent}_{i,*}\|^2 - 2 \cdot \mathbf{obs}_{i,*} \cdot \mathbf{cent}_{j,*}$

$$\mathbf{dist} = \overline{\mathbf{obs}} \cdot \mathbf{1}^T + \mathbf{1} \cdot \overline{\mathbf{cent}}^T - 2 \cdot \mathbf{obs} \cdot \mathbf{cent}^T$$

xGER

xGEMM

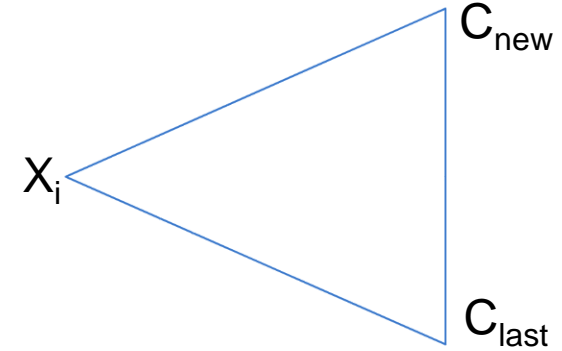
$$A := \alpha * x * y' + A$$

$$C := \alpha * op(A) * op(B) + \beta * C$$

BLAS Subroutines

# Triangular acceleration (Application Phase 2)

- Triangle inequality states :  
$$d(C_{last}, C_{new}) \leq d(X_i, C_{last}) + d(X_i, C_{new})$$
- If  $d(C_{last}, C_{new}) \geq 2d(X_i, C_{last})$ ,  
 $\Rightarrow d(X_i, C_{new}) \geq d(X_i, C_{last})$  without computing
- Distance computations can be further reduced by sorting the inter-centroid distances,  $d(C_{last}, C_{new})$
- New candidate centroids are evaluated as per sorted distance order
- Once the critical distance,  $2d(X_i, C_{last})$  is surpassed all subsequent candidate centroids can be safely discarded

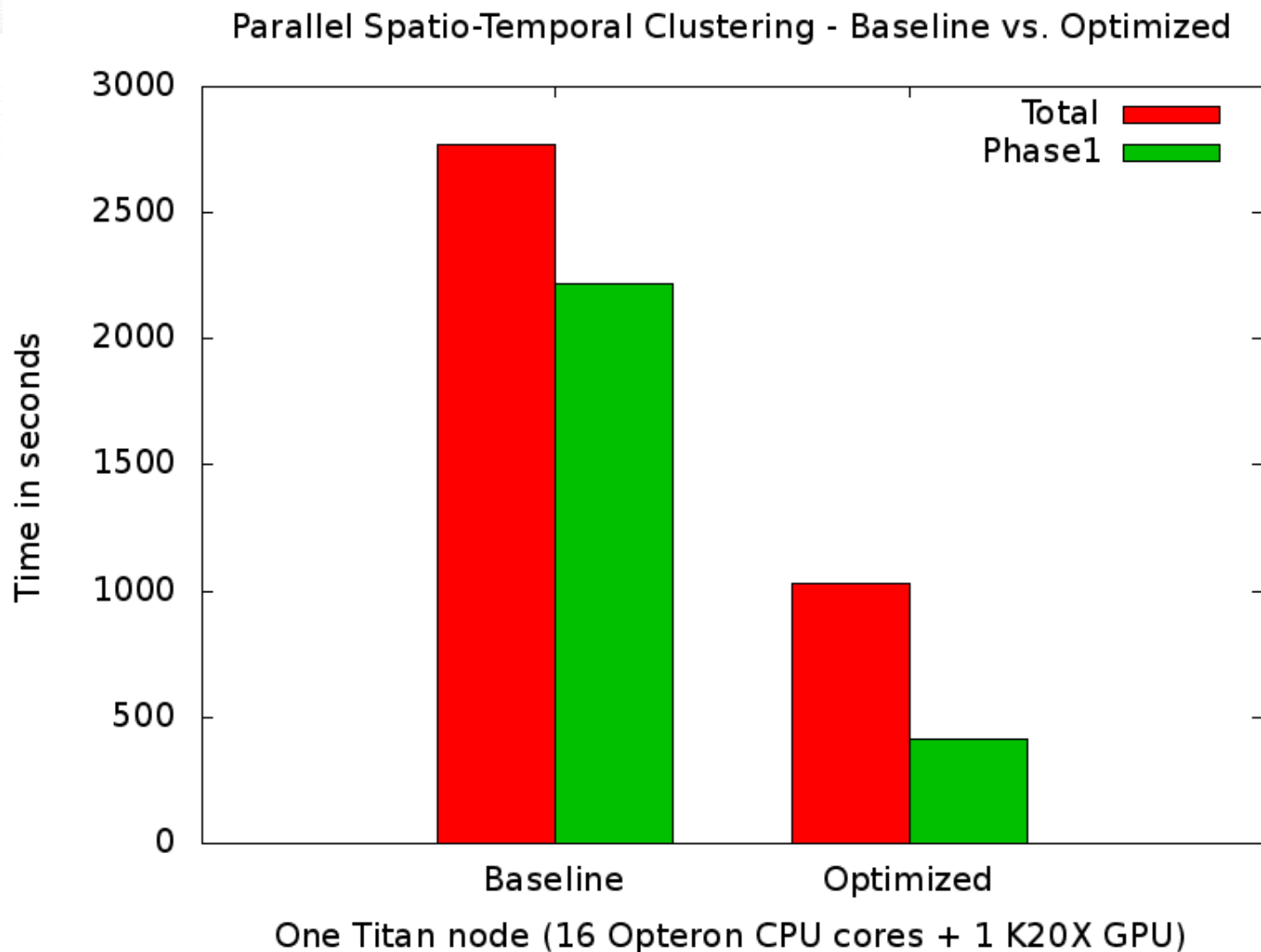


# Computational Environment: Titan

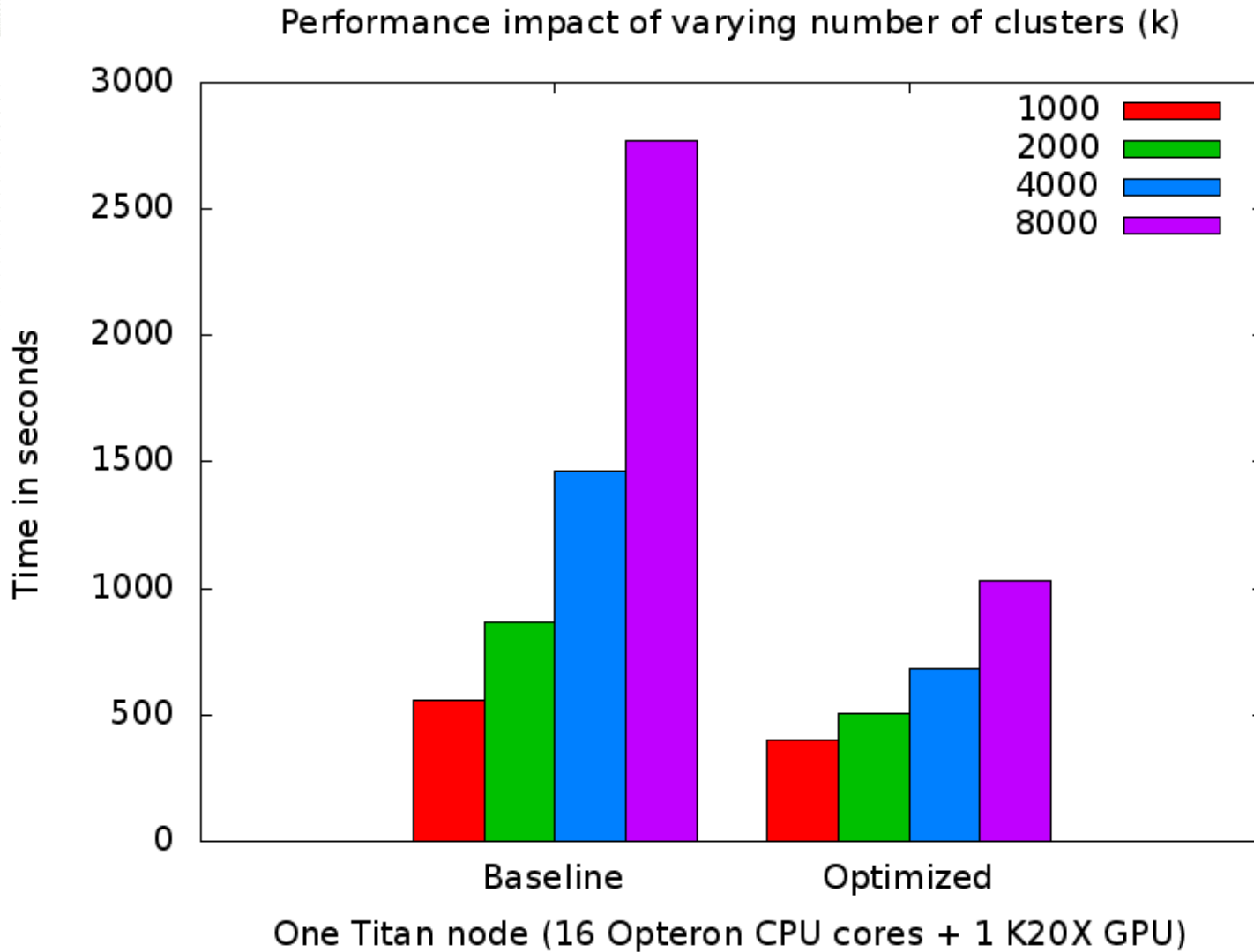
- Cray XK7 system
- Each node
  - 16-core AMD Opteron CPUs
  - NVIDIA Kepler K20X GPUs
  - 32 GB memory
- Total of 18,688 nodes
  - 299,008 CPU cores and 18,688 GPUs.
- Software
  - CPU (MKL + OpenMP)
  - GPU (cuBLAS + OpenACC)
  - MPI for communication



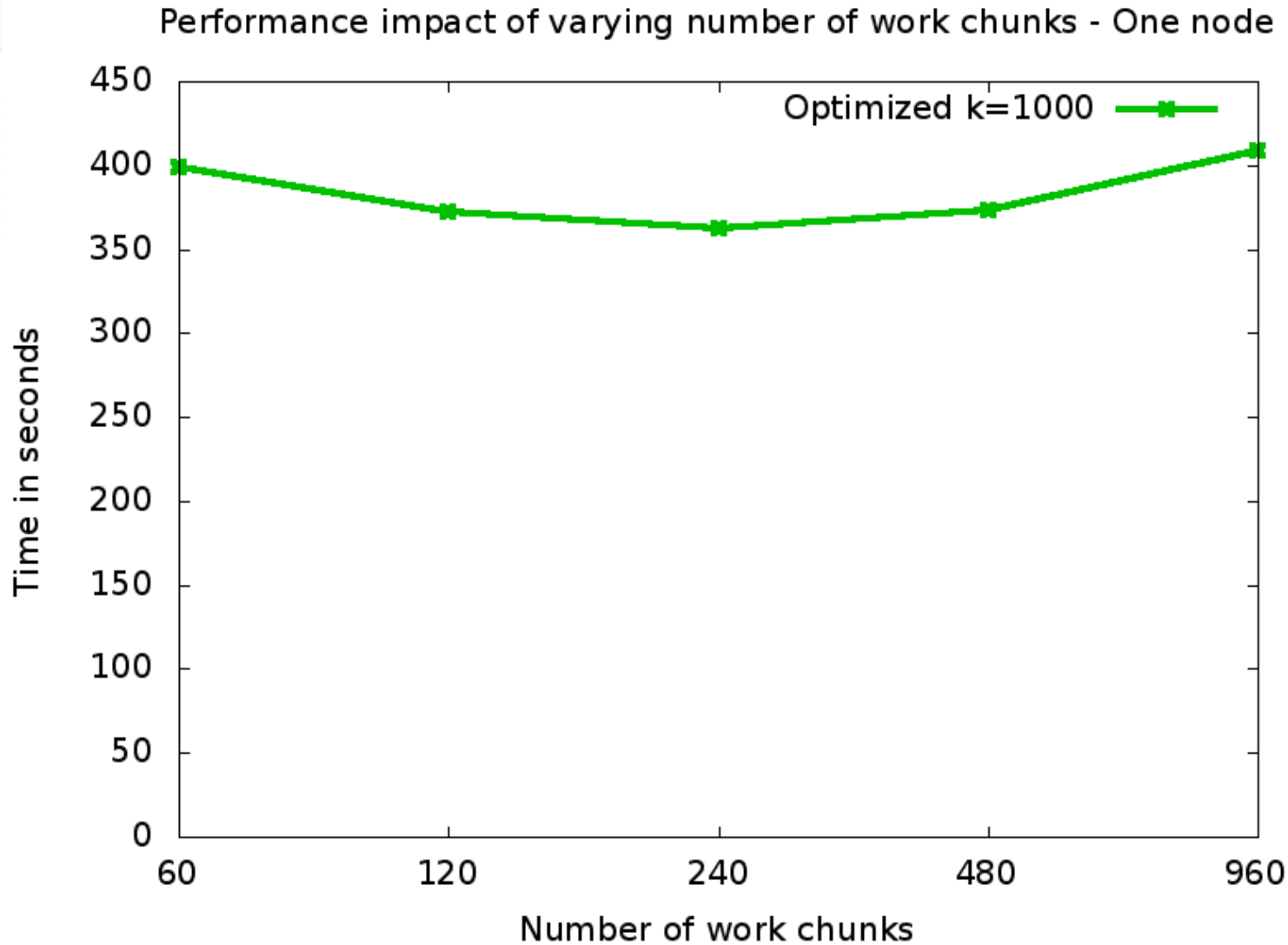
# Performance Comparison



# Performance: Varying Number of Clusters (k)

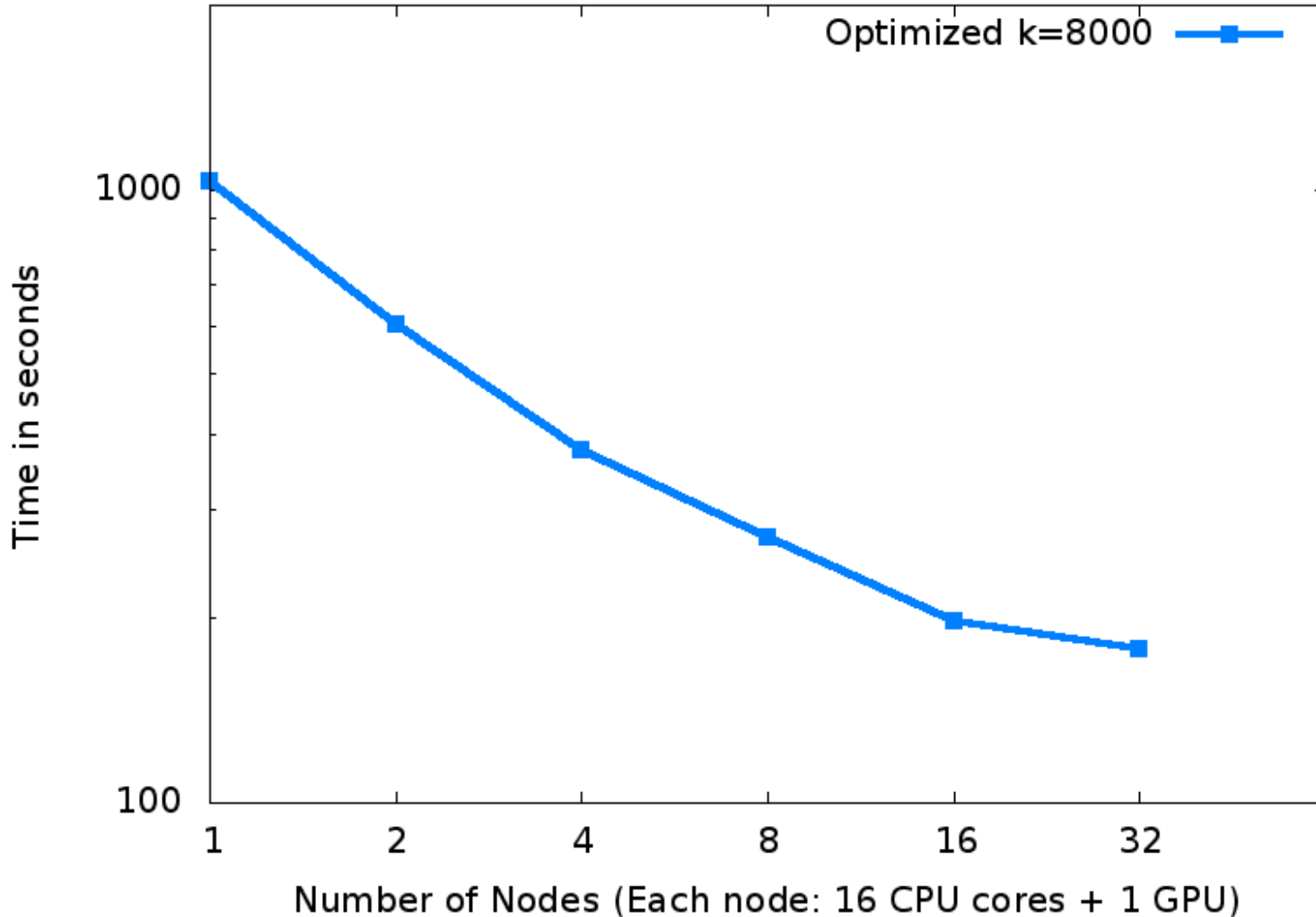


# Performance: Impact of no. of work chunks



# Performance: Strong Scaling

Parallel Spatio-Temporal Clustering - Strong Scaling on Titan

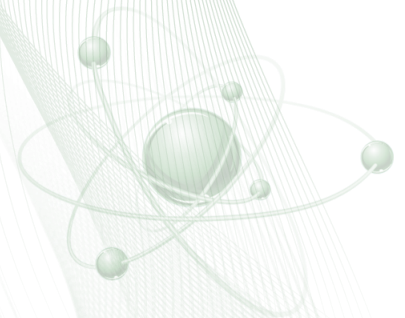




# Limitations and Future Work

- Centralized master: inherent scaling limits
  - Decentralized version in development
- Larger datasets: Exceeds available on-node memory
  - Cluster assignment table and intermediate data structures
  - Short term: Decentralized version should partially address
  - Long term: Looking into NVRAM
- Application phases
  - Heuristic for switching
  - Combination
- Ported to KNL (paper in preparation)

# Applications

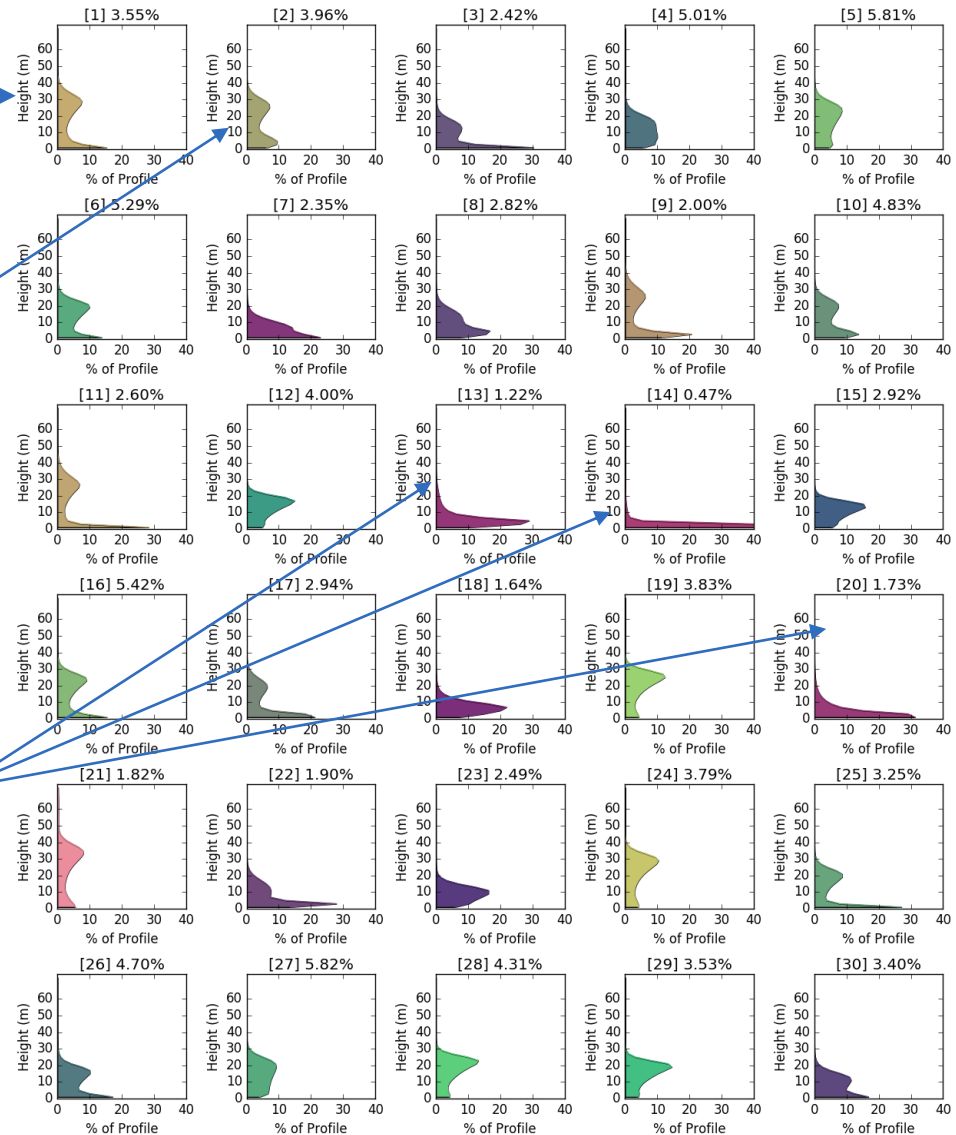


# GSMNP: 30 representative vertical structures (cluster centroids) identified

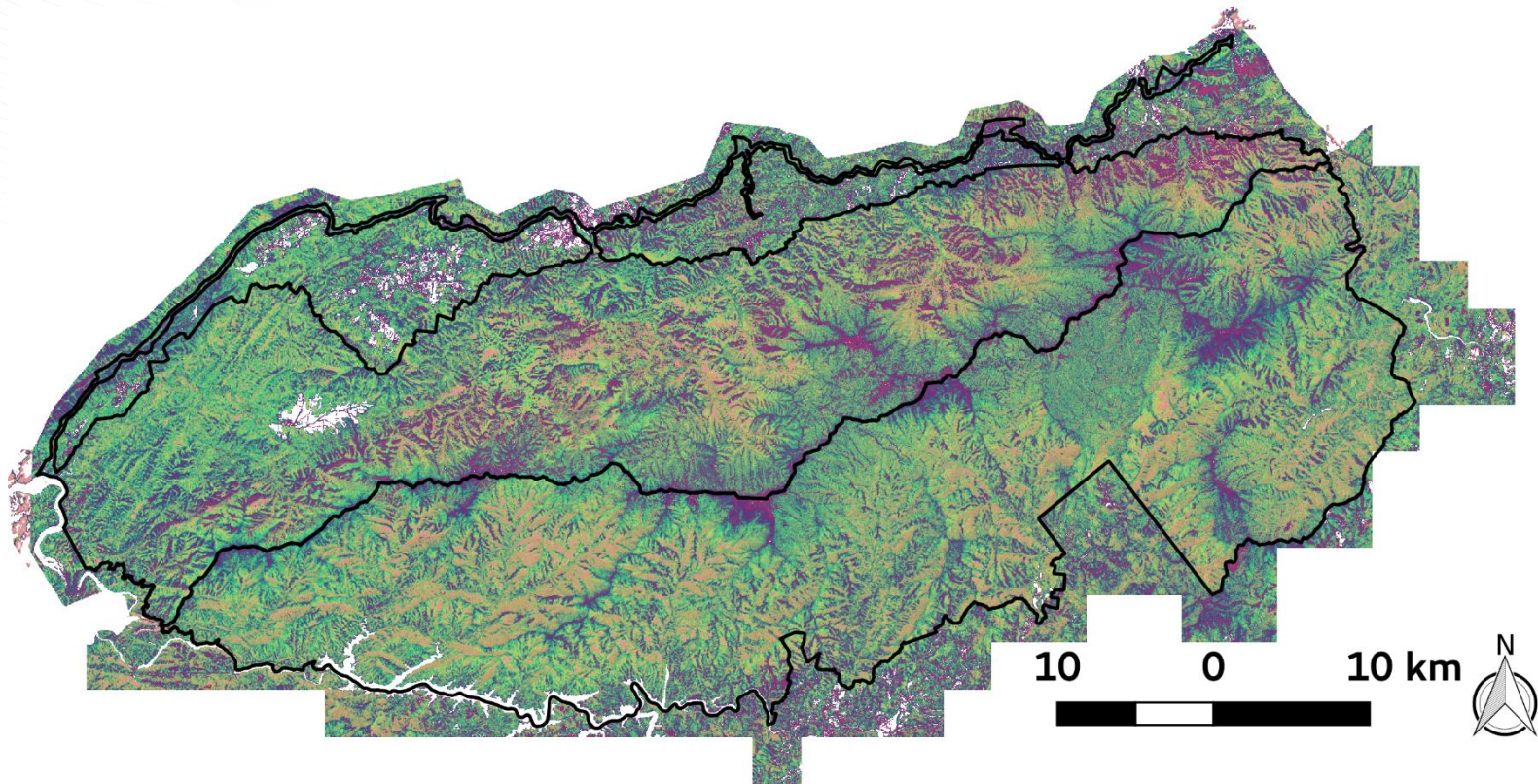
tall forests with low understory vegetation

forests with slightly lower mean height with dense understory vegetation

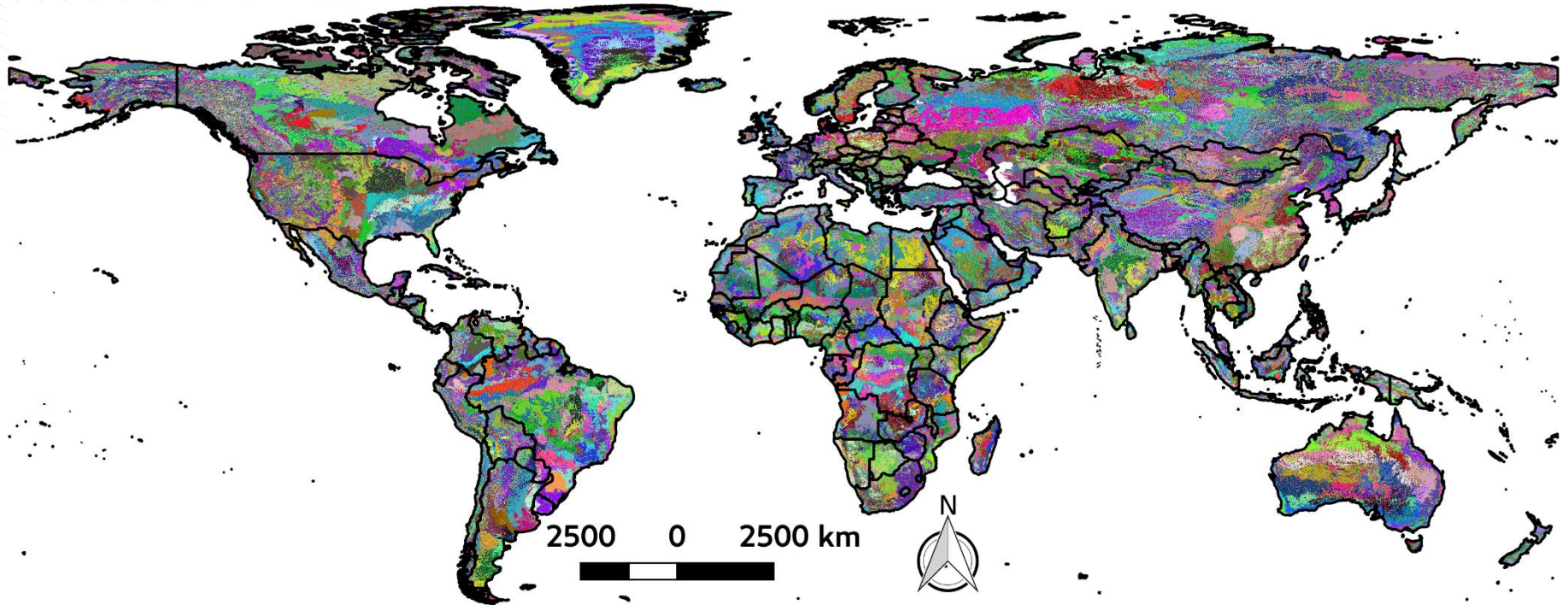
low height grasslands and heath balds that are small in area but distinct landscape type



# GSMNP: Spatial distribution of the 30 vegetation clusters across the national park

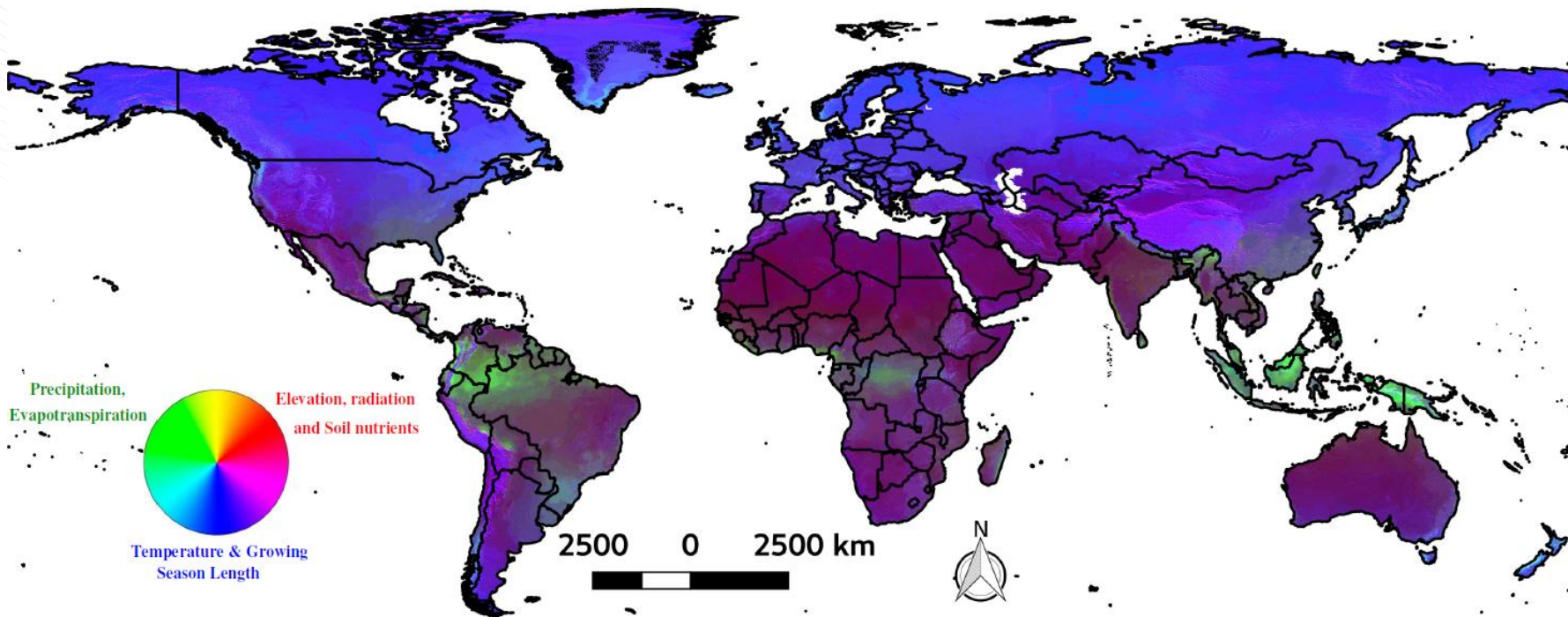


# Global Climate Regimes: 1000 clusters Contemporary using Random color scheme

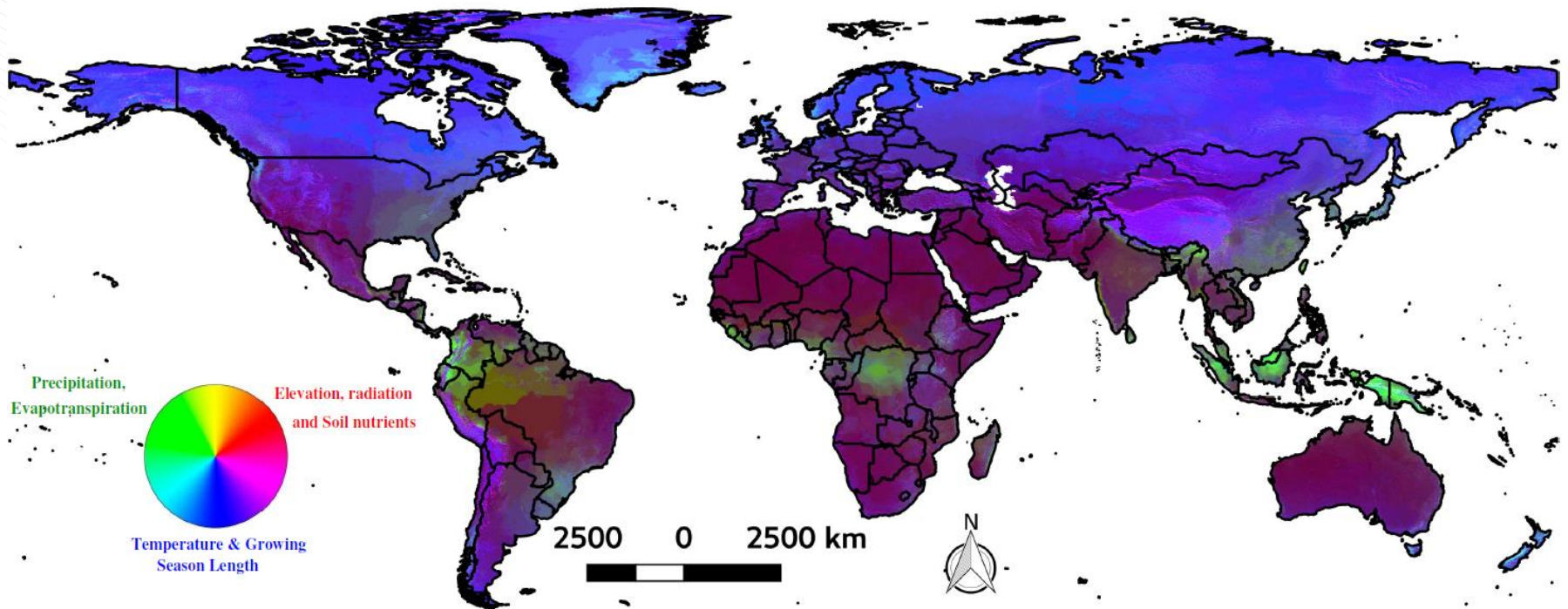


# Global Climate Regimes: 1000 clusters

## Contemporary using Similarity color scheme



# Global Climate Regimes: 1000 clusters 2100 using Similarity color scheme



# Conclusions

- Parallel k-means clustering implementation for hybrid supercomputers
- BLAS formulation to accelerate Euclidean distance calculations
- Demonstrated up to 2.7x speedup over baseline CPU version in specific problem configurations on Titan
- Demonstrated capability to process large datasets
- Two Earth science applications
  - Great Smoky Mountains National Park: identification of vegetation structure
  - Global Climate Regimes: understanding global patterns of climate, vegetation and terrestrial ecology



# Acknowledgments



U.S. Department of Agriculture, U.S. Forest Service,  
Eastern Forest Environmental Threat Assessment Center.



## SciDAC

Scientific Discovery through Advanced Computing

Partial support from SUPER, a DOE Scientific Discovery through Advanced Computing (SciDAC) project

Computer time was provided by the Innovative and Novel Computational Impact on Theory and Experiment (INCITE) program.

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.