

# CRAY VISION FOR DATA CENTRIC COMPUTING FOR EARTH SCIENCES



CRAY®



[icarpenter@cray.com](mailto:icarpenter@cray.com)



# EARTH SCIENCES: CORE DRIVERS



Performance: Drive continued improvements in fidelity of weather & climate simulations



Reliability: Maintain performance, reliability & serviceability as systems grow in size & complexity



Analysis: Derive greater value from environmental data, both observations and simulation



Vision: Manage 10-year transition/rewrite of applications to exploit 2025 machine architectures

# EARTH SCIENCES: WHY CRAY?

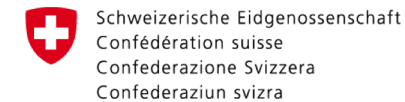
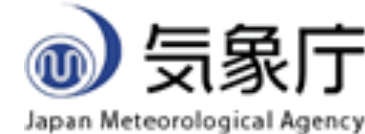


Over 80% of the world's operational weather forecast centers use Cray systems.

- Reliability
  - Operationally proven, unrivalled experience
- Performance
  - Balance performance & throughput across workflow
  - Software development environment, performance tools & application support experts
- Long-term customer partnerships
- New analysis approaches through converged systems



Ministry of Earth Sciences  
Government of India







# Supercomputers are Critical to Simulation

## > Largest ever storm prediction model

- Over 4 billion points used to simulate the landfall of Hurricane Sandy
- Urban scale grid resolution of 500m (compared to standard 3km)
- Enables the research to understanding fine grained properties of hurricanes

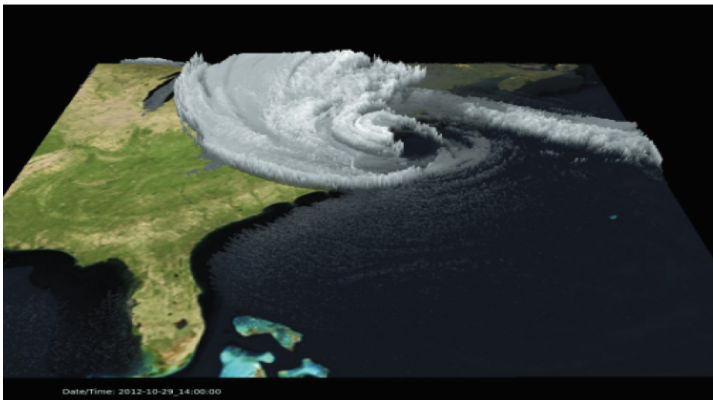


## > Studying crop devastation by whiteflies to address a major cause of hunger in East Africa

- Understanding the DNA of the species by generating phylogenetic trees
- With only 500 whiteflies in a genetic dataset, the possible relationships between these flies run into the octillions ( $10^{25}$ )

## > Key to the development of new antiretroviral drugs

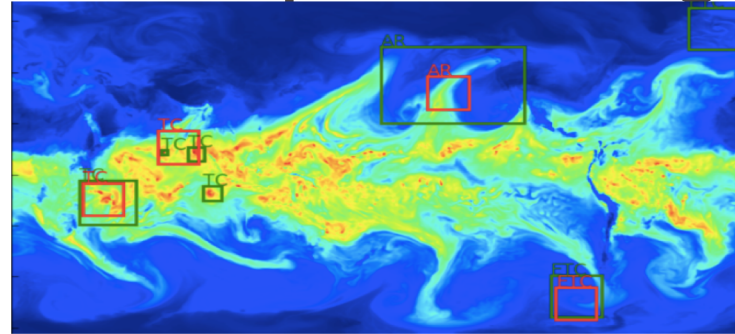
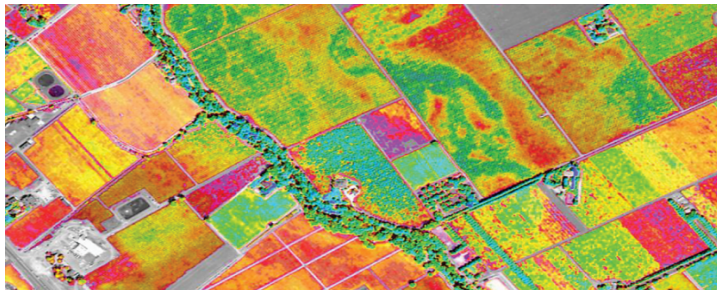
- Determined the precise chemical structure of the HIV capsid – the protein shell that protects the virus's genetic material and is a key to its virulence.
- Requires the assembly of more than 1,300 identical proteins – in atomic-level detail.



# ...and to Machine and Deep Learning

## > Crop data is key to decision makers

- Applying Deep Learning on satellite data, the two major crops can be distinguished with 95% accuracy just a few months after planting and well before harvest.
- More timely estimates could be used for a variety applications, including supply-chain logistics, commodity market future projections, and more.



## > Quantitatively assess how extreme weather will change in the future

- A single climate simulation can produce over 100TB of data with archives reaching over 5PB.
- Deep learning techniques are ideal for pattern recognition over large data sets



## > Development of systems for connected cars and autonomous technologies

- These advances could not have been realized without the application of deep learning to object detection in image and full motion video



**SAMSUNG RESEARCH AMERICA**





# MUCH More Data to Serve, Store and Manage

CRAY

## MORE INPUT DATA

More powerful instruments



More powerful sensors



Advanced Analytics in the workflow



AI (ML/DL) in the workflow



COMPUTE

## MORE OUTPUT DATA

Bigger problem sizes



Higher fidelity models



Higher resolution models



New data-intensive algorithms



# Not Just More Data But Also Different I/O Patterns

CRAY

**Modeling &  
Simulation**

**Advanced  
Analytics**

**Artificial  
Intelligence**





# AS DATA VOLUME INCREASES



- More focus on ability to handle and use data efficiently
  - Analyze ensemble output
  - Apply AI techniques to both observations and model output to identify severe weather features
- Improve accuracy through better data assimilation methods
- Interconnect focus evolves to meet these needs
  - Less on how it enables single job MPI scalability
  - More focus on how the interconnect provides excellent data access
  - More focus on congestion management to enable MPI- and I/O-dominated jobs to get reproducible performance

# Cray Vision

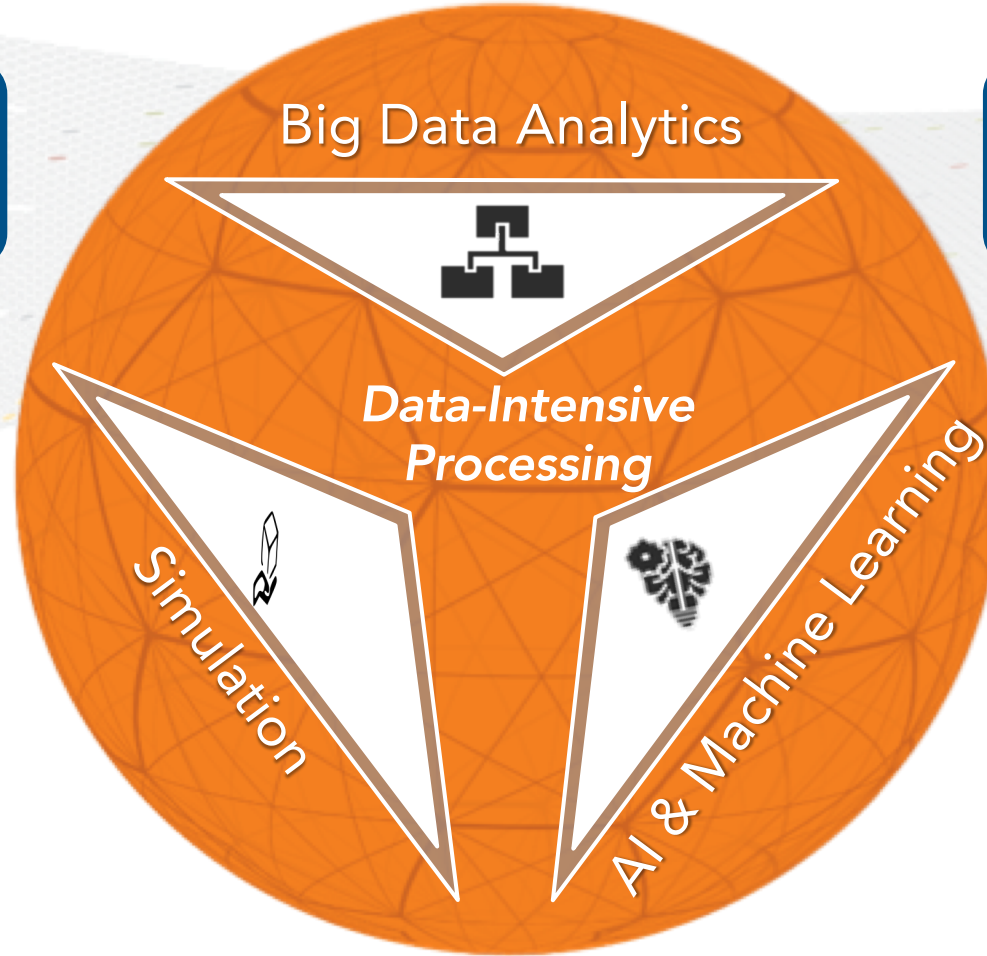
## *Convergence of Supercomputing and Artificial Intelligence*



### Supercomputing



- Large and Latency Sensitive Data Movement
- Massively Scalable Computing
- Dense Accelerated Computing
- Parallel Storage and I/O



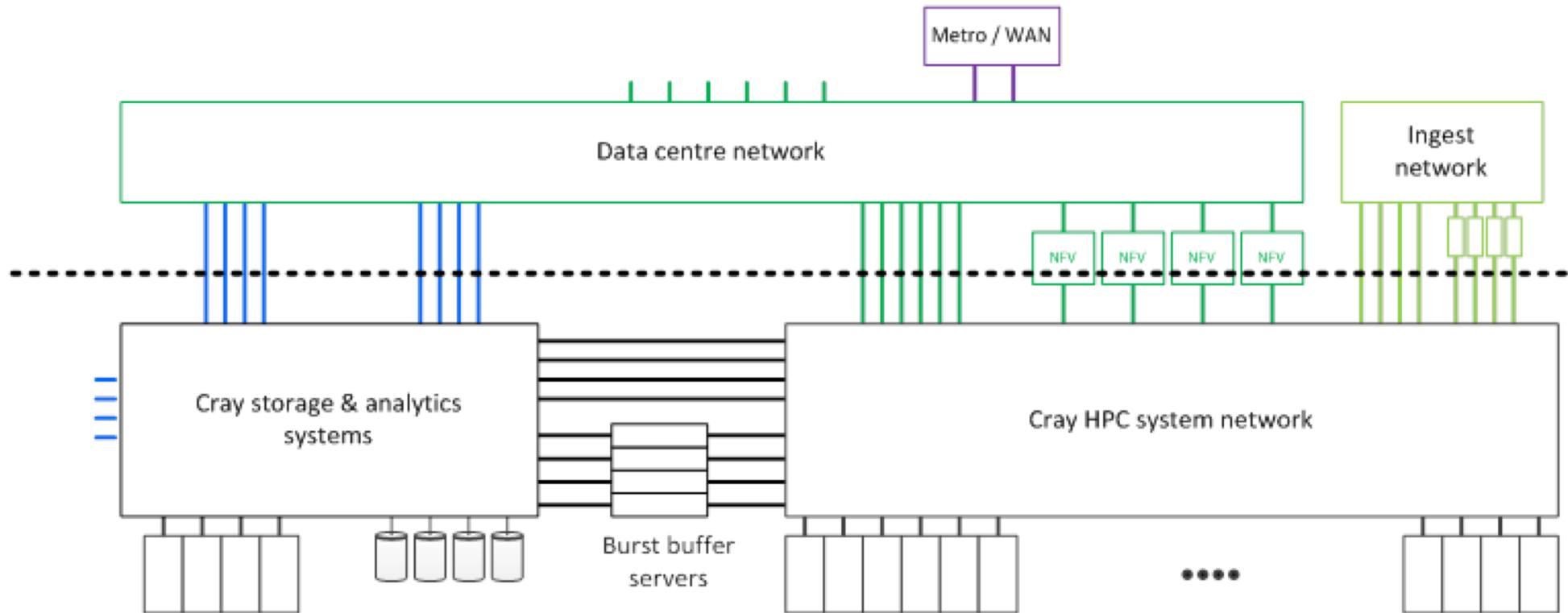
### Artificial Intelligence



- AI Workflow
- Machine Learning and Deep Learning
- Graph Analytics
- Apache Spark™ and Python-based data science



# CRAY CONVERGED ARCHITECTURE VISION



# SHASTA FOR EARTH SCIENCES

## MAJOR SLINGSHOT ENHANCEMENTS



Slingshot will be a **great interconnect** for our earth sciences customers

This is Cray's **8th** supercomputing interconnect



MPI  
PERFORMANCE



TRAFFIC  
CLASSES



CONGESTION  
MANAGEMENT



FLEXIBILITY



RELIABILITY



ETHERNET  
COMPABILITY



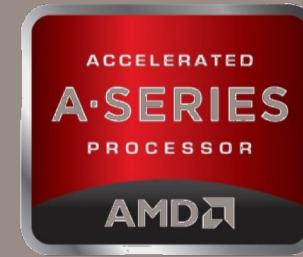
Very low average and TAIL latencies improve runtime reproducibility



Better data ingress/egress



# EXPLOSION IN PROCESSOR CHOICES



New era – end of Dennard scaling – can't get faster performance just from shrinking silicon

→ specialization for different kinds of tasks

Need for performance portability

# STORAGE TRENDS – FLASH COSTS DECLINE



## Today

- **DataWarp** (stage/de-stage via WLM, different namespace)
- **ClusterStor L300N** and **NXD** (cache - transparent to users)
- **ClusterStor L300F** (all flash, small, random I/Os, 500,000 IOPs)
- local SSDs

## Future

- Lustre tiering and additional Lustre features

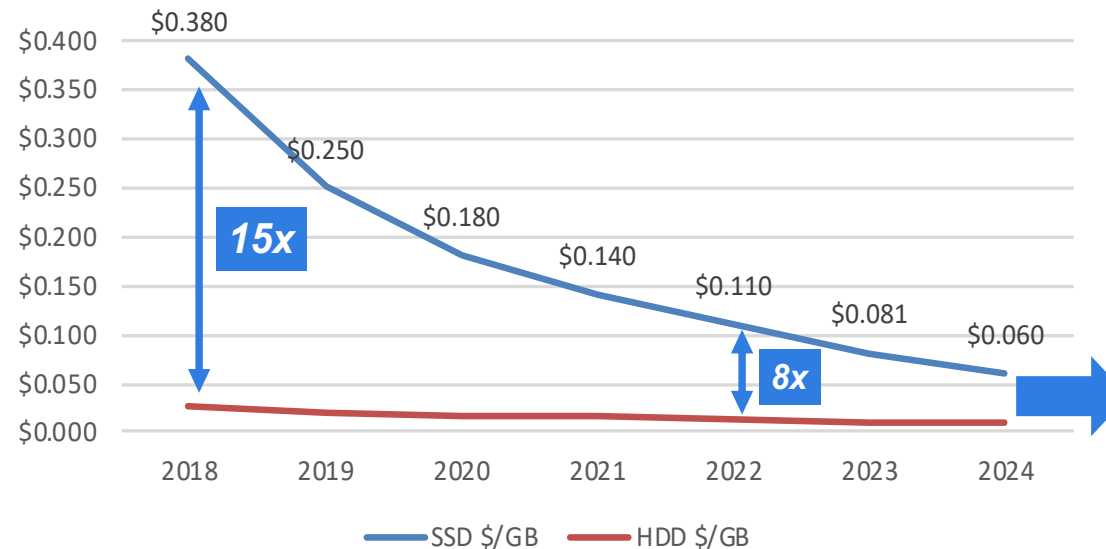
## ECONOMIES OF STORAGE MEDIA



~ 60% of the cost of any storage system is in the media

\$\$ per GB - Enterprise drives

May 2018, IDC



**SSD Media =**  
Most of performance  
+ some initial capacity

**HDD Media =**  
Some performance  
+ rest of capacity

**SOFTWARE =**  
Workflow-  
accelerating  
data placement  
on the right media  
at the right time

# EXCHANGE OF IDEAS BETWEEN SUPERCOMPUTING & BIG DATA



- Hyper-scale web companies have developed many interesting technologies:
  - Driven by data velocity, volume & variety, combined with resiliency requirements at scale
- Several are of particular interest in weather & climate science:
  - Containerization
  - Flexible & scalable data analysis platforms
  - New analysis techniques and machine learning



# CONTAINERS FOR FLEXIBILITY

The ability to create your own world supports a diverse workload on a shared supercomputer



Build with Certified Software  
Stacks



Bundle libraries and  
dependencies



Build a consistent Environment  
from desktop to supercomputer



# URIKA-XC MAKE ANALYTICS AND GRAPH “FIRST CLASS” CITIZENS ON XC SERIES SYSTEMS



## SIMULATION



**Expanding to  
Analytics and Open  
Data Science**

### **Python Open Data Science**

- **Production Supercomputing**
- Weather Forecasting
- Seismic Imaging
- Manufacturing CAE
- **Scientific Supercomputing**
- Climate Science
- Chemistry & Materials Science

### **Spark Big Data Analytics**

- Data Preparation
- Analysis
- Visualization
- Machine Learning
- Deep Learning

### **Large-Scale Graph Discovery**

- Cancer Cell Morphology
- Fraud and Insider Threat Detection

***Deployed using containers for  
portability and ease of use***

# CRAY® URIKA-CS AI SUITE



Pre-integrated and supported AI stack with popular open source AI frameworks and libraries delivered as container images for ease of development and deployment

UIs: Jupyter Notebooks, TensorBoard

Java, Scala, R, Python	MLlib, GraphX, Spark SQL, Spark Streaming	BigDL	Anaconda Python	TensorFlow™
Apache Spark™			Dask	Cray Distributed Training Framework (CrayPE ML Plugin)
Intel® MKL, Intel MKL-DNN, cuDNN, NVIDIA SDK, OpenMPI				

## Urika-CS:

- Pre-integrated software suite for AI workflow
- Distributed deep learning training for heterogeneous systems (CS500 & CS-Storm)
- Supported stack – for open source components and for distributed training framework

CRAY®  
CS-STORM

CRAY®  
CS500



# FLEXIBLE & SCALABLE DATA ANALYSIS PLATFORMS



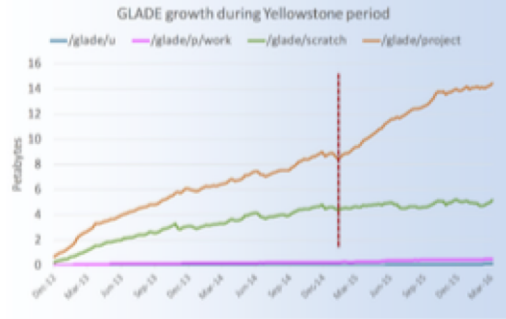
- Equipped with rich set of analysis libraries & easily extensible
- Enables higher-productivity languages (Scala, R, Python) and Jupyter interactive notebooks
- Leverage Python scalable data analysis infrastructure
  - Met Office IRIS library
  - DASK Parallel Python engine
  - Pangeo software stack
- Performance with SPARK, DASK, R, Cray Graph Engine

# NEW ANALYSIS APPROACHES: WEATHER/CLIMATE INFORMATICS

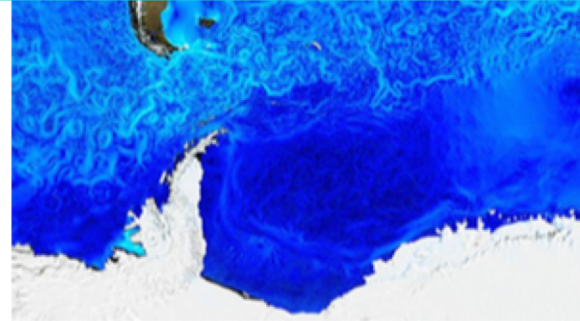


- Applying more complex analysis treatments to find new insights or increase accuracy of predictions
  - Machine/Deep Learning
  - Regression/Clustering/Optimization methods
- Challenges:
  - Meteorological Data very different from classic “Big Data” sets
    - Multivariate, with Spatial & Temporal locality
    - Can be very sparse/limited in some dimensions
    - Limited/missing observations over unpopulated areas
  - Data-driven predictions must be linked physics/first-principles based theory

# PANGEO DATA: TOWARD A BIG DATA ANALYSIS PLATFORM



**Petabyte-scale data volumes are straining CISL's infrastructure**



**Scalable analytics solutions are required to work with large datasets**

**Pangeo Goal:** create an open-source toolkit for the analysis of climate datasets, built on the **Python** language ecosystem, **Xarray** multi-dimensional array tools, and **Dask** parallel analytics system.

**Parallelism is key: single device performance is falling behind!**



## Packages

### Pangeo Core Packages

Xarray  
Iris  
Dask  
Jupyter

### Pangeo Affiliated Packages

Guidelines for New Packages  
General Best Practices for Open Source  
Best Practices for Pangeo Projects  
Why Xarray and Dask?



- Website: <http://xarray.pydata.org/en/latest>
- GitHub: <https://github.com/pydata/xarray>

Xarray is an open source project and Python package that provides a toolkit for working with labeled multi-dimensional arrays of data. Xarray adopts the **Common Data Model** for self-describing scientific data in widespread use in the Earth sciences: **xarray.Dataset** is an in-memory representation of a **netCDF** file. Xarray provides the basic data structures used by many other Pangeo packages, as well as powerful tools for computation and visualization.

## Iris



- Website: <https://scitools.org.uk/iris/docs/latest/>
- GitHub: <https://github.com/SciTools/iris>

Iris seeks to provide a powerful, easy to use, and community-driven Python library for analysing and visualising meteorological and oceanographic data sets.

With Iris you can:

- Use a single API to work on your data, irrespective of its original format.
- Read and write (CF-)netCDF, GRIB, and PP files.
- Easily produce graphs and maps via integration with matplotlib and cartopy.

Iris is an alternative to Xarray. Iris is developed primarily by the [UK Met Office Informatics Lab](#).

## Dask



# MACHINE/DEEP LEARNING IN WEATHER/CLIMATE



- Arduous to train, but comparatively quick to run
- Data producer vs data consumer
- Complementary use cases
  - Rapid classifiers or predictors for radar/observations
  - Advanced MOS systems – D1Cast system
  - Pattern recognition in model outputs or observations
  - Emulation: Replacement of expensive parameterizations
    - train NN using expensive radiation model
    - train NN using cloud-resolving model, replace convection parameterization

# AI, MACHINE LEARNING & DEEP LEARNING



ARTIFICIAL INTELLIGENCE						
Sense		Comprehend	Predict		Act and Adapt	
ANALYTICS			MACHINE LEARNING			
Search datasets for insights			Learn patterns from the past to predict future			
Descriptive	What happened?		Unsupervised Group, cluster and organize content with domain-specific heuristic models		Supervised Train mathematical predictive models with labelled data	
Diagnostic	Why did it happen?				DEEP LEARNING	
Predictive	What will happen?		Train and use neural networks as a predictive model			
Prescriptive	How to make it happen?		Vision	Speech	Language	

# DETECTING EXTREME WEATHER

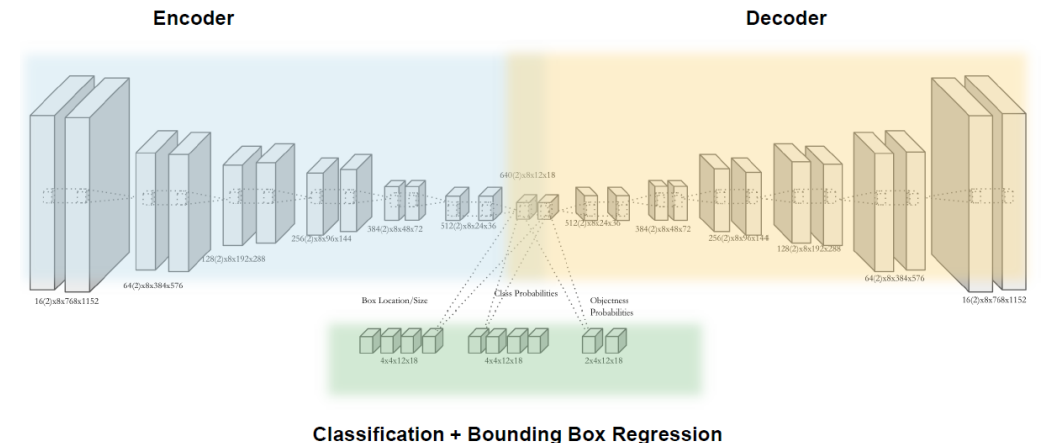
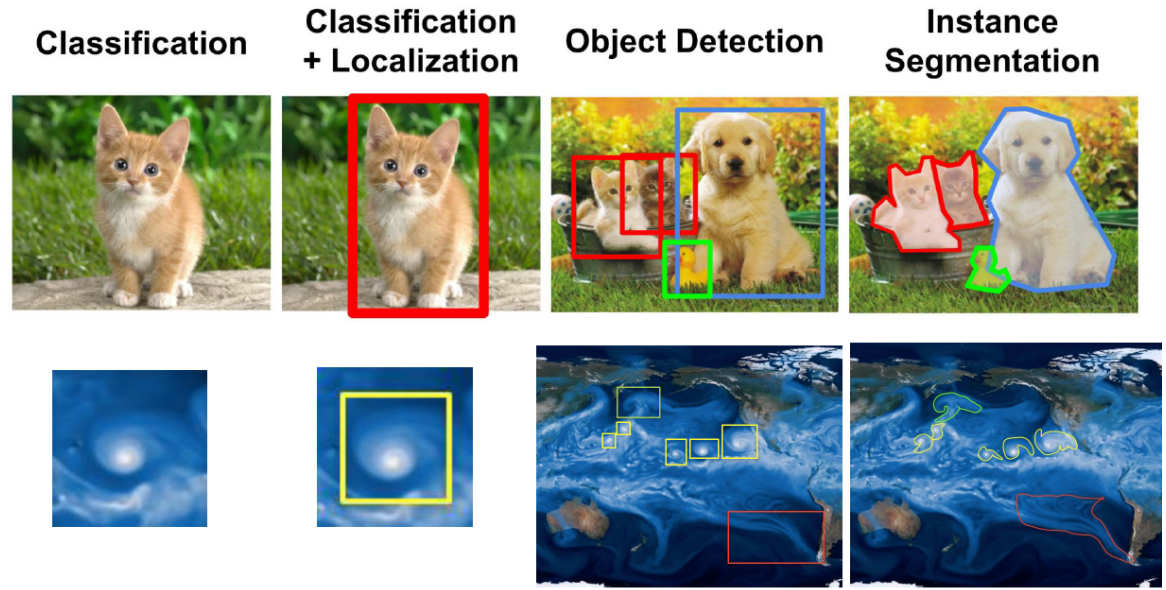


## Challenge:

- Climate simulations run at 10000x faster than real-time and high resolution required to reproduce extreme weather events – generate 100'sTB of data

## Supercomputer + AI Solution:

- Semi-supervised convolutional architectures can identify extreme weather events such as Tropical Cyclones, Atmospheric Rivers, Weather Fronts with 90% accuracy
- 15-PetaFLOP Deep Learning system used to scale training of a single model to ~9600 Xeon-Phi nodes; obtaining peak performance of 11.73-15.07 PFLOP/s



Full Paper - <https://arxiv.org/pdf/1708.05256.pdf>

# ML APPLIED TO WEATHER MODEL OUTPUT



Can neural networks to enhance the prediction of damaging hail? If so, why?

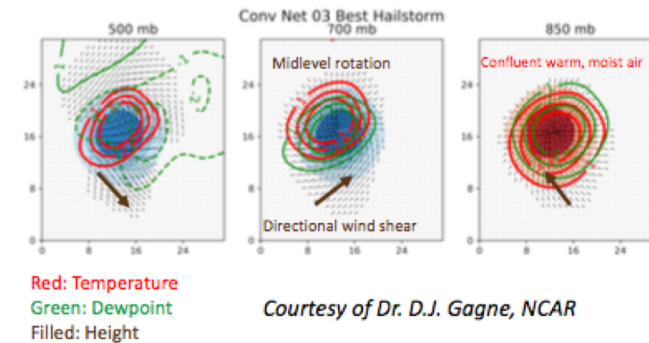
NN identifies physically relevant features: temperature, dewpoint, elevation, confluent warm moist air, mid-level rotation.

*D.J. Gagne, NCAR*

Rich Loft, 2018  
ECMWF Workshop on  
HPC in Meteorology

## Interpretability: Do Neural Nets dream of electric hail storms?

Neural network identify physically relevant features for hailstorm prediction from core weather fields. Running the network in reverse reveals these features.



Courtesy of Dr. D.J. Gagne, NCAR



NCAR  
UCAR

21st Century Earth System Modeling *air • planet • people*



# IMPROVING SATELLITE DATA UTILIZATION



## Through Deep Learning at NOAA

- Satellites provide more data than can be assimilated, ~3% of available data is used today
- Use DL object detection to identify areas of atmospheric instability from satellite observation data, focus extraction of observations on these regions of interest

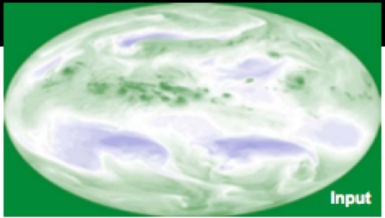
Run on Theia – Cray CS-Storm system

100 nodes, each with 8 NVIDIA Tesla P100 GPUs


### Using Satellite Data for Training

- Water Vapor Channel from GOES 10, 11,12,13,14, and 15
- Storm centers from IBTracks Dataset
- Data normalized to range from -1 to +1
- Trained 2010-2013 Validated 2014, Test 2015
- Images resized and cropped to 1024x512
- Image segmentation 25x25 pixel box segmentation centered on storm
- Only use storms classified as Tropical Storm or greater on Saffir Simpson Scale
  - 34 knots and above

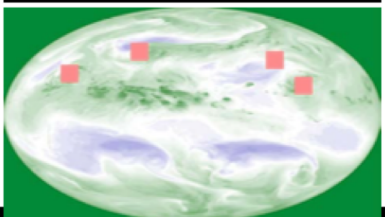
~ 4500 Labeled Data



Input



Labeled Data



NOAA - Earth System Research Laboratory

*Jebb Stewart, 2018 ECMWF workshop on HPC in Meteorology*

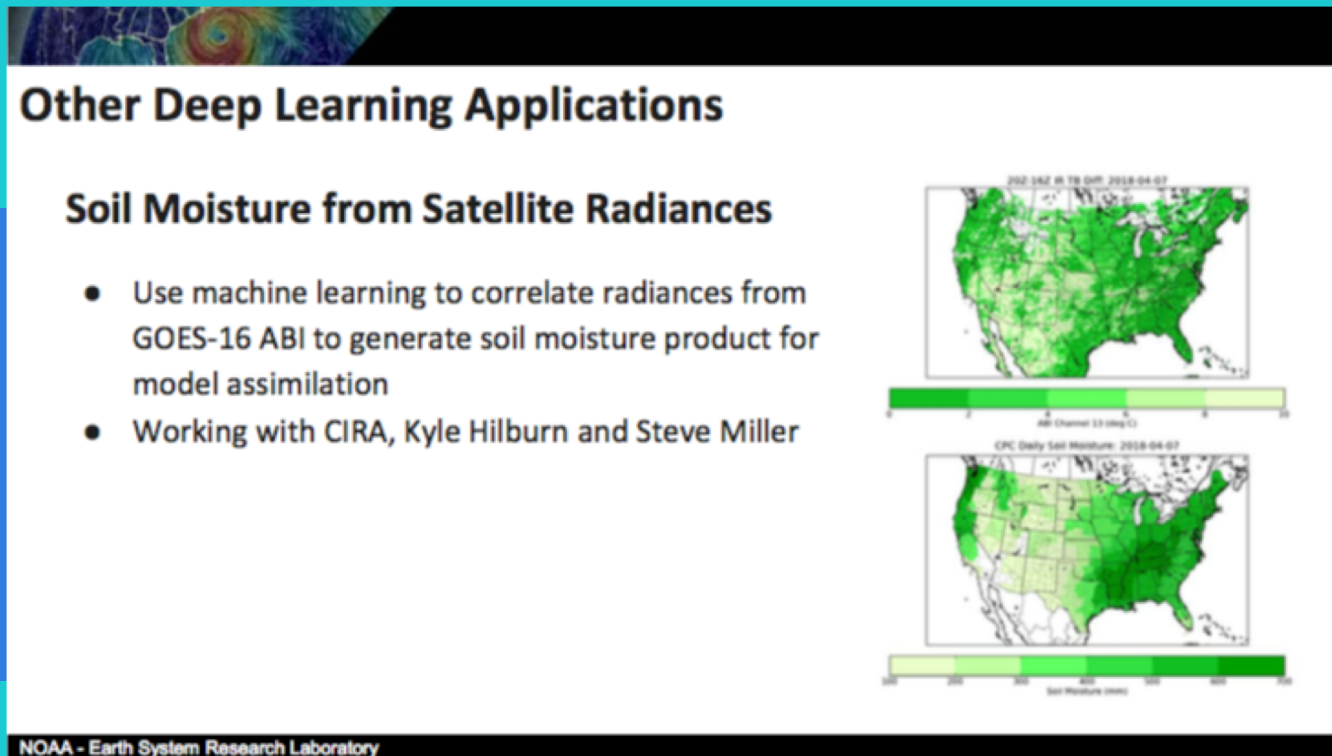
# SOIL MOISTURE

## Through Deep Learning at NOAA



Run on Theia – Cray CS-Storm system

100 nodes, each with 8 NVIDIA Tesla P100 GPUs



*Jebb Stewart, 2018 ECMWF workshop on HPC in Meteorology*

# MET OFFICE INSTALLS Urika-XC



“As in many industries, we are challenged with increasing data volumes and are turning to large-scale analytics, machine learning and deep learning applications to drive new insights and innovation,” said Charles Ewen, director of technology and CIO at the Met Office. “The Met Office already has one of the world’s largest Cray XC supercomputing systems. Now with our implementation of Cray’s Urika-XC software, we are applying AI and analytics to deliver ever-more accurate and detailed weather forecasts and climate change analyses, while also developing new commercial products.”

## UK MET OFFICE CHOOSES CRAY AI SOLUTION TO UNLOCK BUSINESS VALUE FROM WEATHER DATA

### Weather Center to Use Cray Urika-XC AI and Analytics Software to Develop Tailored Forecasts and Specialized Commercial Weather Products

SEATTLE, Sept. 26, 2018 (GLOBE NEWSWIRE) — Global supercomputer leader Cray Inc. (Nasdaq:CRAY) today announced that the [Met Office](#), the United Kingdom's National Weather Service, has expanded its Cray® XC40™ supercomputer with artificial intelligence (AI) and analytics capabilities. The Met Office added [Cray's Urika®-XC](#) AI and analytics software suite to its supercomputer to unlock the highest levels of business value from the massive volumes of weather data it processes daily.

The Met Office is using Cray's Urika-XC suite to explore the use of new methods, such as machine learning, in extracting insights from observational and model data to better develop and customize commercial products, such as tailored forecasts. The Urika-XC suite was designed to run on Cray XC systems to eliminate the need for organizations to install new purpose-built analytics hardware and enable customers to run simulation and big data workloads on the same system.

“As in many industries, we are challenged with increasing data volumes and are turning to large-scale analytics, machine learning and deep learning applications to drive new insights and innovation,” said Charles Ewen, director of technology and CIO at the Met Office. “The Met Office already has one of the world's largest Cray XC supercomputing systems. Now with our implementation of Cray's Urika-XC software, we are applying AI and analytics to deliver ever-more accurate and detailed weather forecasts and climate change analyses, while also developing new commercial products.”

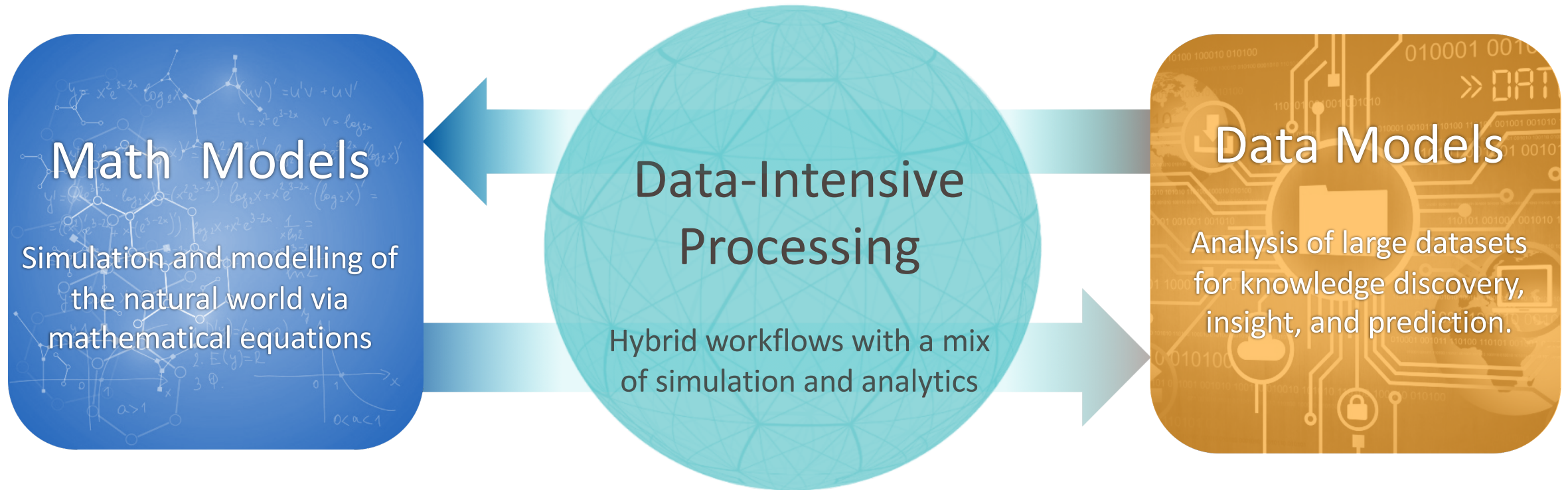
With the power of a Cray supercomputer, the Met Office is able to take in [215 billion weather observations](#) from all over the world every day and uses an advanced atmospheric model to create tailored forecasts and briefings that are delivered to customers including governments, environmental agencies, the military and the general public as well as businesses and other organizations.

“Cray and the Met Office share a long, productive and successful relationship, and we're pleased that one of the world's most prestigious weather agencies is taking the next step using Cray AI and analytics solutions to augment their capabilities to drive new business opportunities,” said Per Nyberg, vice president of market development, artificial intelligence and cloud at Cray. “The Met Office's decision demonstrates its confidence in Cray's innovation and creativity in helping tackle some of the planet's biggest weather and climate challenges. At Cray, we believe big data analytics, modeling and simulation are converging into new workflows leading to powerful insights for customers.”

# MIXED SIMULATION/ANALYTICS WORKLOADS

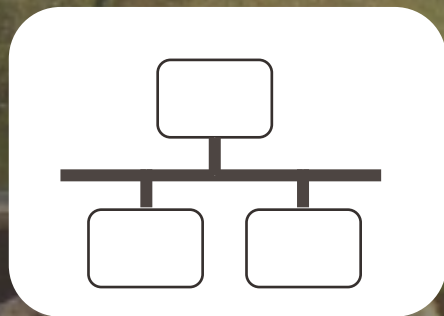


## Computational Modeling





# Run Any Workflows



Converged  
Network



Containers  
&  
Virtualization



Extensibility





# THANK YOU

QUESTIONS?



[icarpenter@cray.com](mailto:icarpenter@cray.com)



[cray.com](http://cray.com)



[@cray\\_inc](https://twitter.com/cray_inc)



[linkedin.com/company/cray-inc-/](https://www.linkedin.com/company/cray-inc/)

