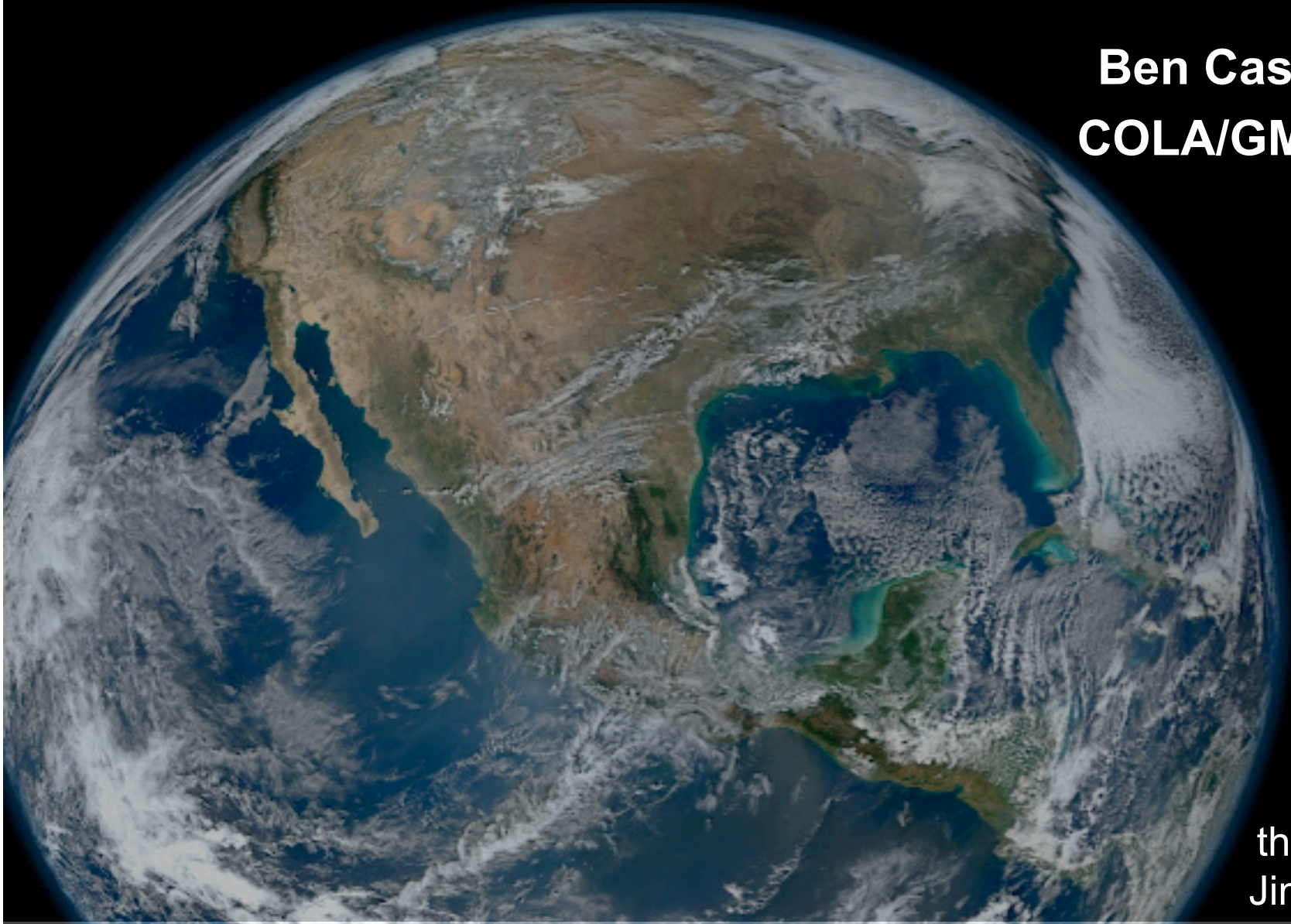


Climate Modeling and Big Data: Current Challenges and Prospects for the Future

Ben Cash
COLA/GMU

Many
thanks to
Jim Kinter!



Why Makes Climate Research a Big Data Field?

*What types of data do your consumers want?
What makes this type of data big?*

- **Society wants information about weather impacts in our current and changing climate**
 - **requires high spatial and temporal resolution**
- **Uncertainty in these impacts is addressed through multi-model ensembles and Earth system prediction**
 - **requires added data and system complexity**

Need for High Resolution

- **Improving the fidelity of climate models** has been objective of intense effort since the 1970s, but it has proven very difficult
- **Million-fold increase in computing capability** since 1980
- **Numerical weather prediction** made substantial progress by:
 - Increasing spatial resolution
 - Improving understanding of physical processes
 - Improving data assimilation methods
- **Climate models** have improved, primarily through the inclusion/refinement of more processes that are relevant to climate variability and change
- Could **enhanced spatial resolution** improve climate model fidelity (and perhaps change our understanding of climate dynamics both qualitatively and quantitatively)?

Driver: Societal Demand for Climate Information

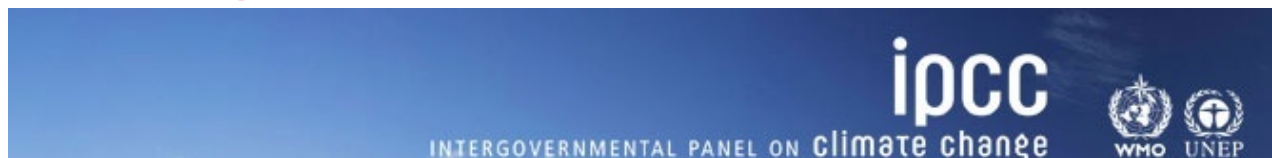
- America's Climate Choices**



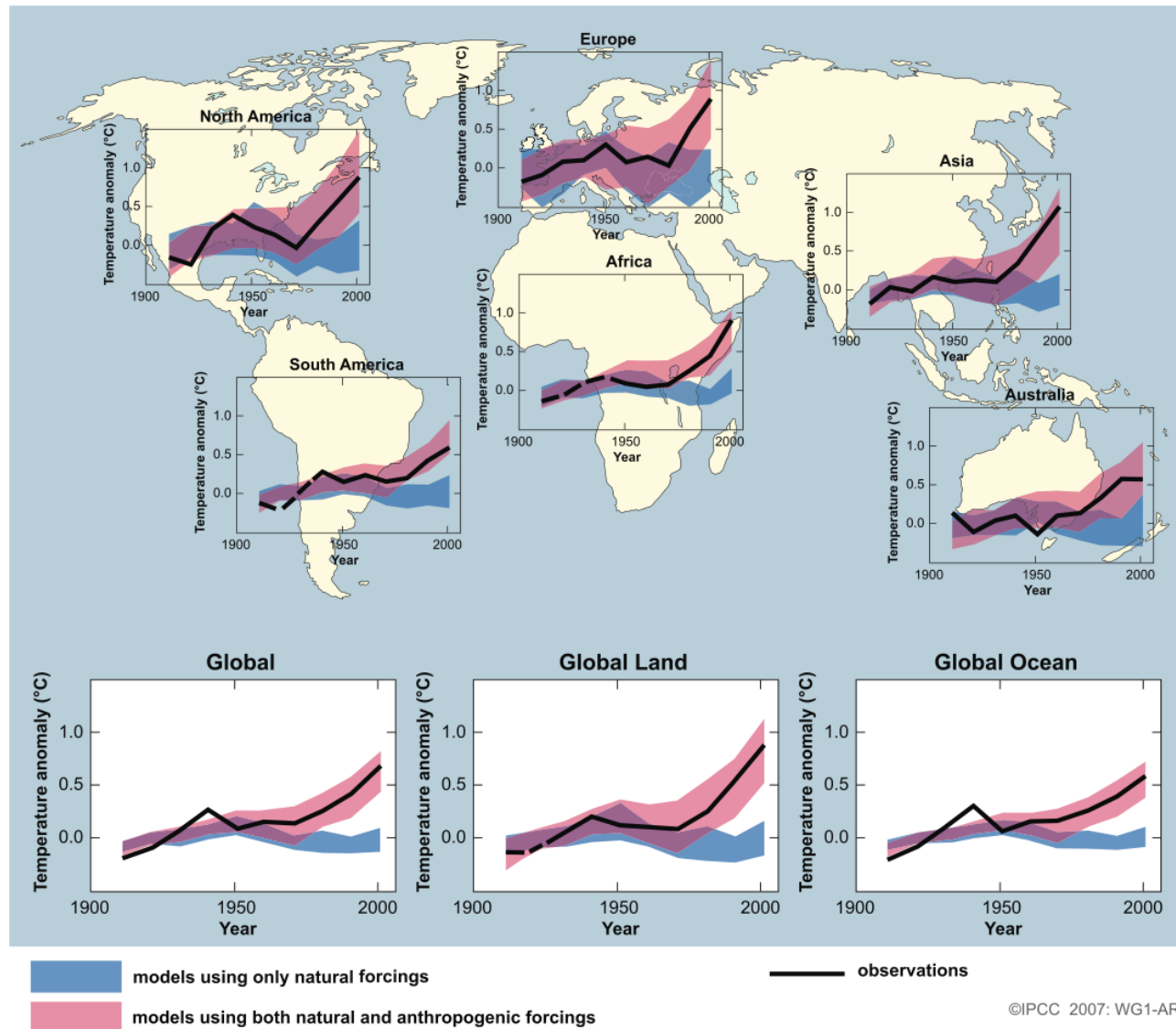
- (USGCRP) National Climate Assessment**



- Intergovernmental Panel on Climate Change**

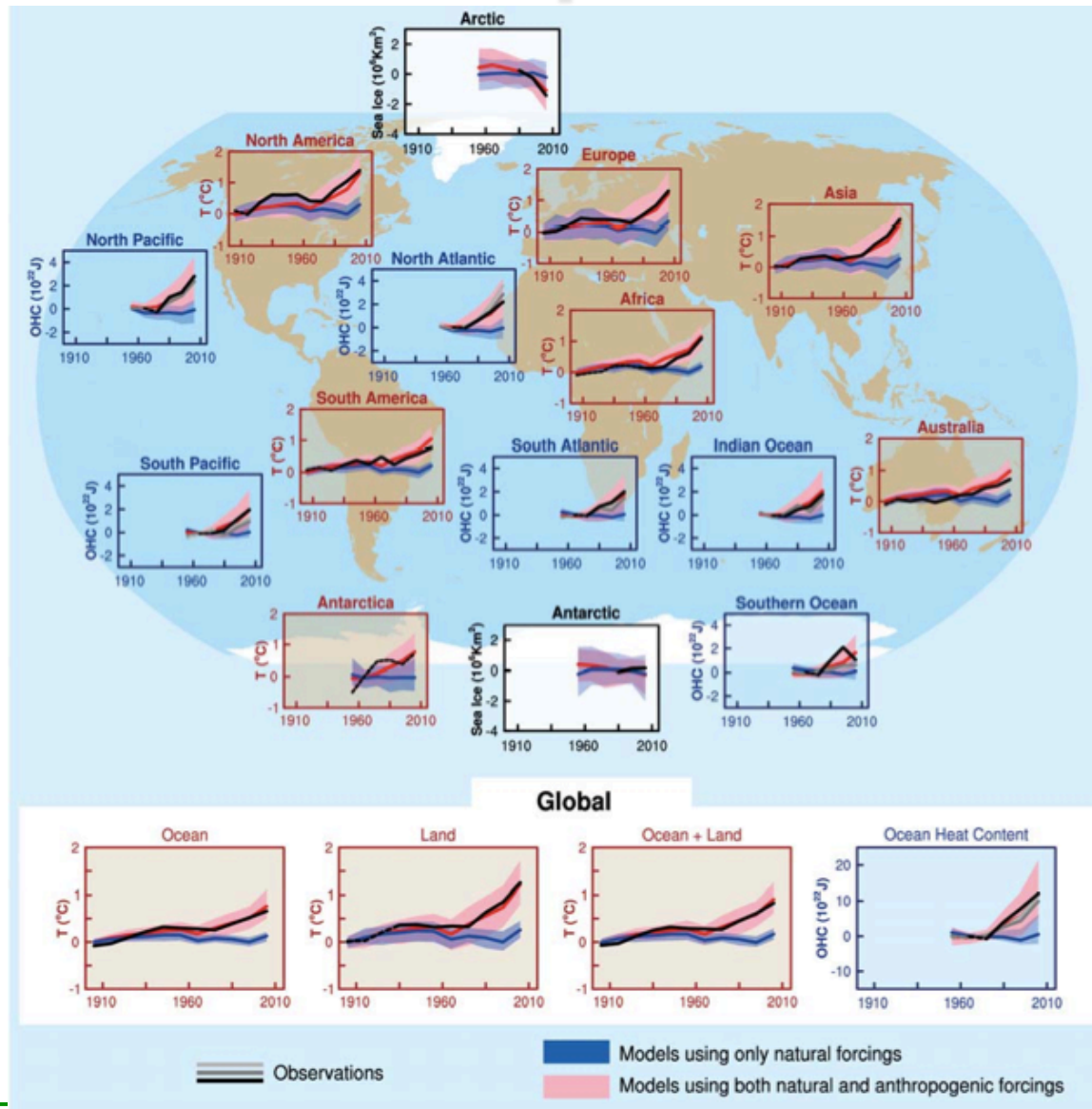


AR4 Surface Temperature Change



IPCC
AR4
WGI
2007

AR5 Surface Temperature Change

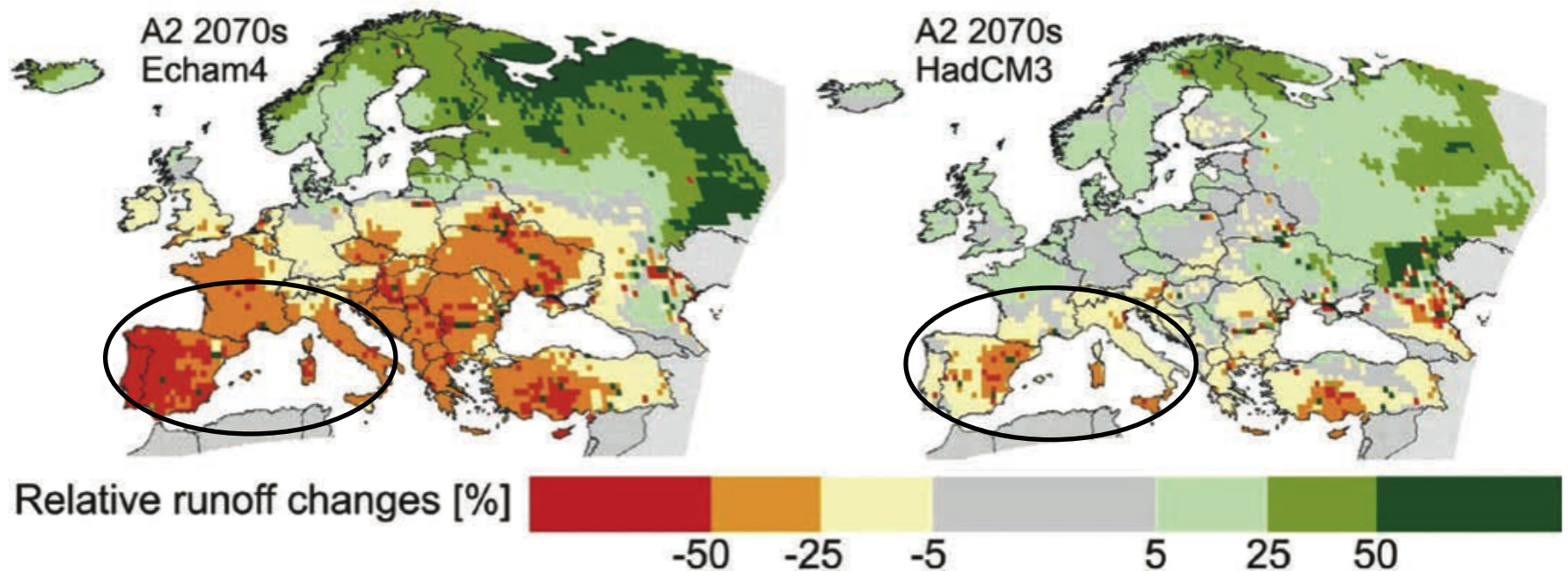


IPCC
AR5
WGI
2013

Regional Climate Change – Beyond CMIP3 Models' Ability?

ECHAM

HadCM

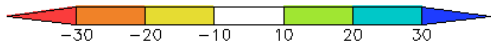
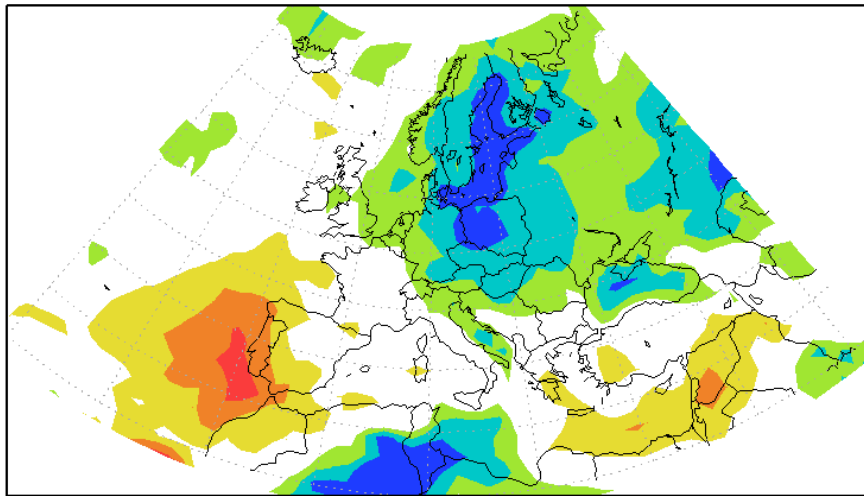


WHAT ABOUT CMIP5?

Regional Climate Change: Sample from CMIP5

CNRM

CNRM Change in Precip 2071–2100 Minus 1971–2000
apr–sep

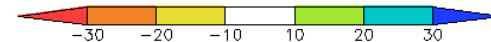
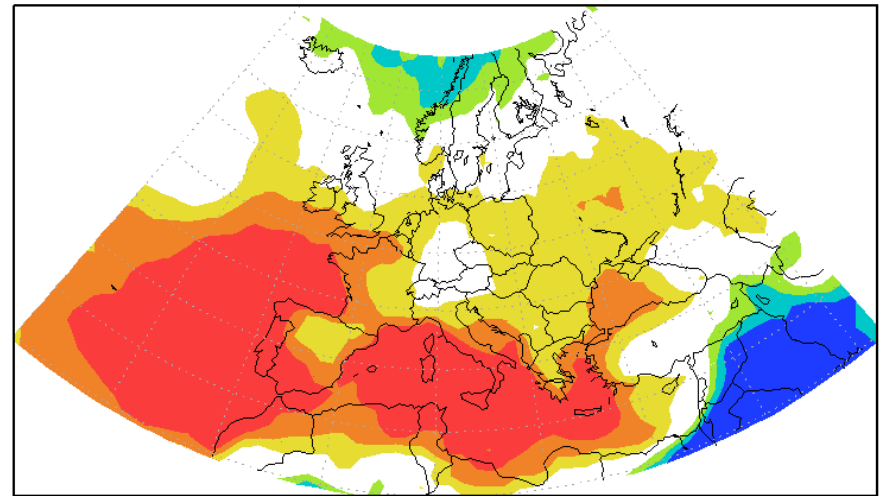


GRADS: COLA/IGES

2014-03-17-00:38 GRADS: COLA/IGES

CESM

CESM Percent Change in Precip 2071–2100 Minus 1971–2000
apr–sep



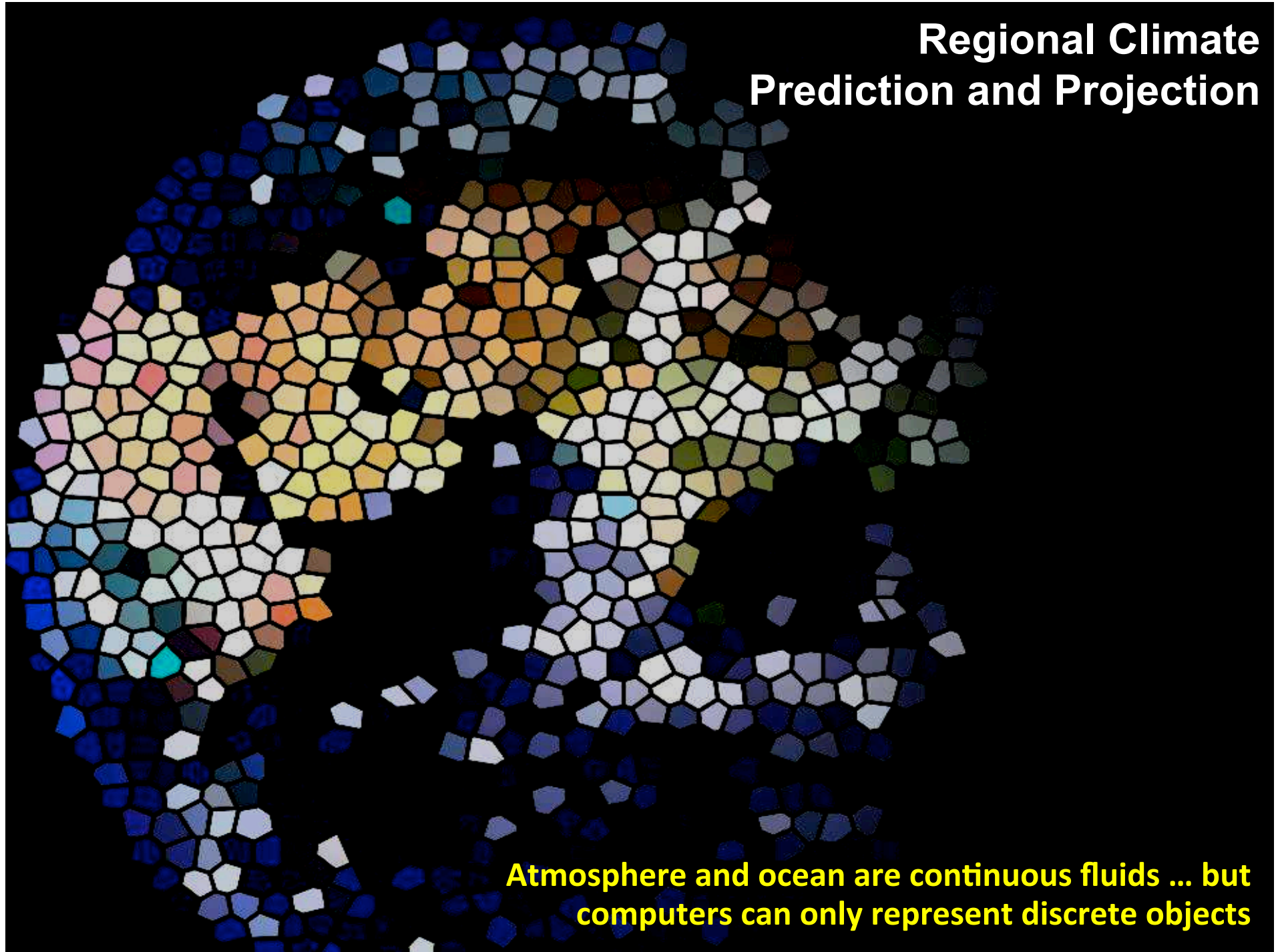
2014-03-17-00:39

A satellite image of the Earth showing the Middle East, surrounding oceans, and parts of Africa and Asia. The image is used as a background for the text.

Regional Climate Prediction and Projection

Atmosphere and ocean are continuous

Regional Climate Prediction and Projection



Atmosphere and ocean are continuous fluids ... but
computers can only represent discrete objects



Regional Climate Prediction and Projection

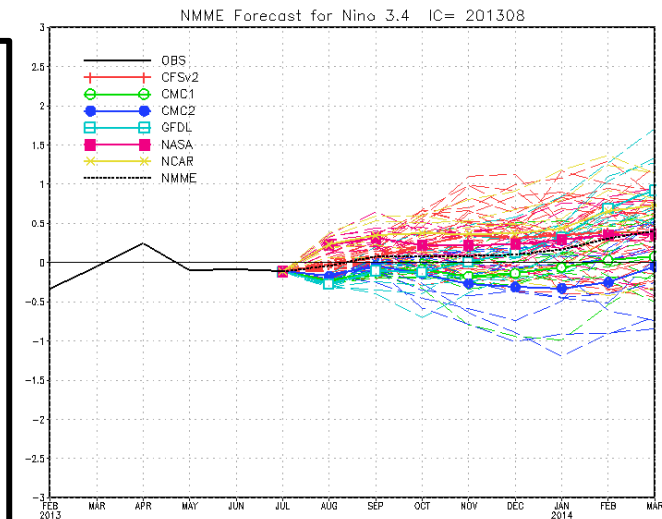
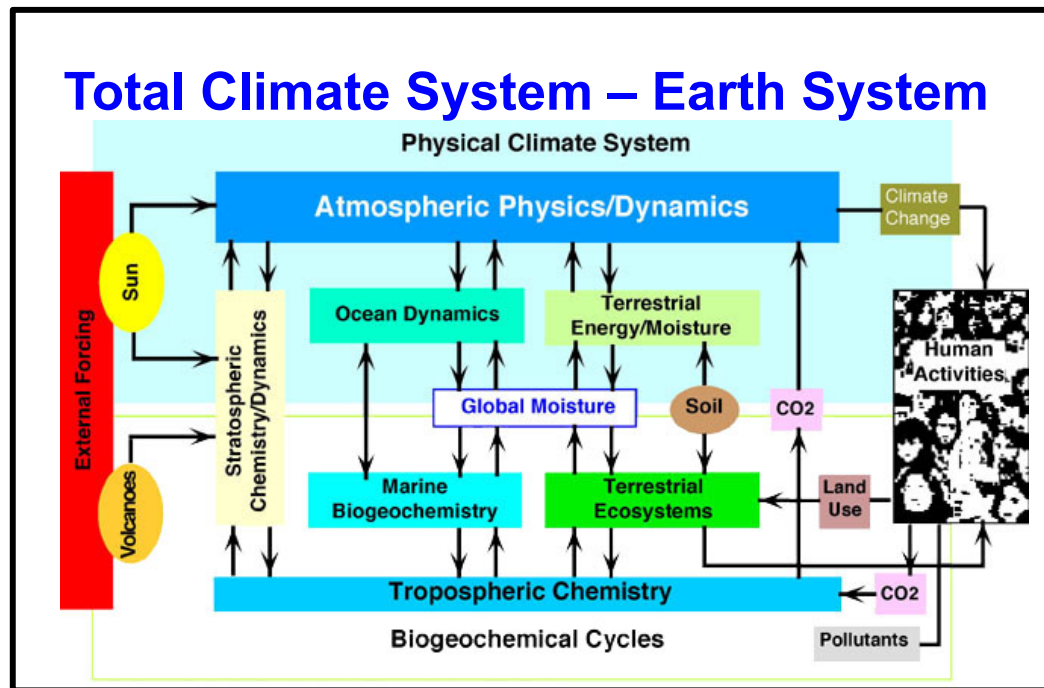
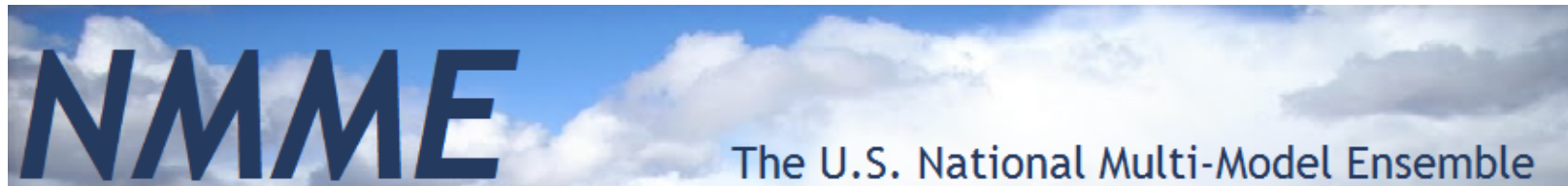
... means **increasing
models' spatial
resolution**

... which is required for :

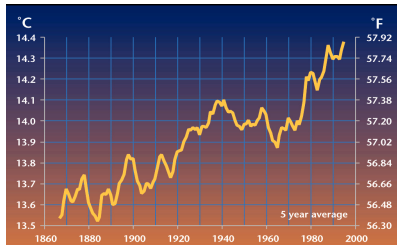
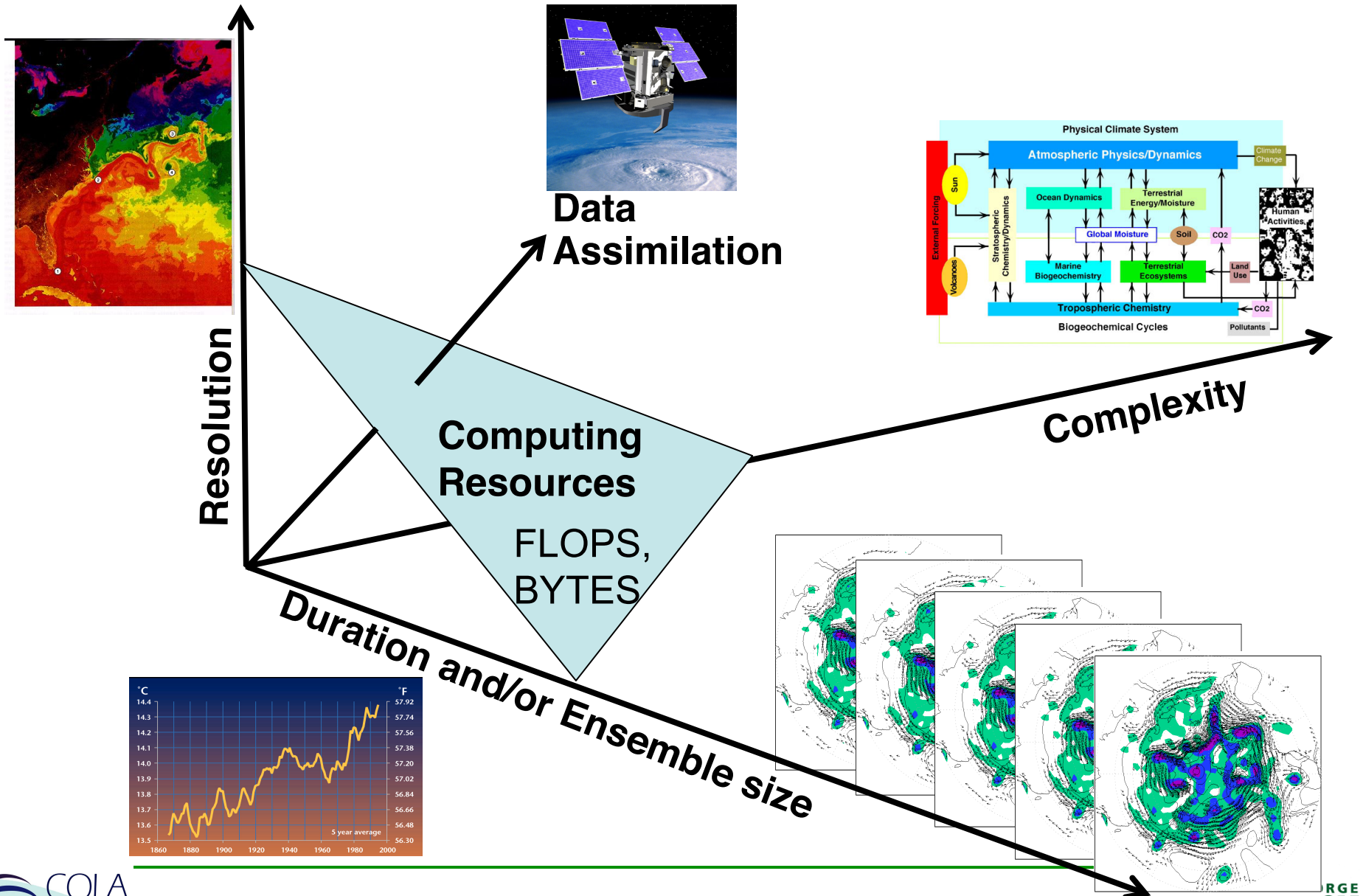
- **Accuracy of solution**
- **Representation of features**
- **Representation of processes**
 - **Interaction among scales**
 - **Meeting societal demands**

... and which **drives computational and
data resources demands**

Addressing Uncertainty: Multi-Model Ensembles & Total Climate System Prediction



Resource Demands



Adding It Up: Sample Volumes

2005-2006	<i>CMIP3 (in support of IPCC AR4)</i>	36 TB
2009-2010	<i>Project Athena</i>	1.2 PB
2010-2011	<i>CMIP5 (in support of IPCC AR5)</i>	3 PB
2012-2014	<i>Project Minerva</i>	3+ PB
2011-	<i>NMME</i>	1 PB

- COLA storage resources for 2015- 1 **PB**
- This much data breaks everything: H/W, systems management policies, networks, apps S/W, tools, and shared archive space
- **In 2012, generating 800 TB using 28 M core-hours took our group ~3 months; this would take about a week using a comparable fraction of a system with 1M cores!**

Climate Analysis

What types of methods are used to analyze your data?

Do you confront issues of privacy/security/ethics?

Do you confront issues of data standards and interoperability?

Do you confront issues of limited current capacity in the data sciences?

- Analysis methods tend to be ‘classic’ linear techniques
 - Regression, correlation, composites, principal components, etc.
 - Some published machine learning work, but relatively rare
- Tend to focus on limited number of relationships
 - Global mean temperature and greenhouse gases
 - El Nino and rainfall, etc
- Unlikely to fully reflect underlying relationships
 - Many success stories despite limitations
- “Small data” techniques
 - Generally unchanged despite increasing data volumes

Climate Analysis

- Barriers to Applying Data Science to Climate
 - Physical
 - Too much data to easily transfer from HPC facility
 - Spinning disk is not a limitless resource
 - Data needs to be in an accessible location and accessible format
 - Technical
 - Limited amount of observational data – may not be enough samples
 - Massive quantities of model data – how to handle systematic error?
 - Cultural
 - Not a standard component of climate training, lack of familiarity makes for slow or limited adoption by climate community
 - Interpretation – Without physical constraints, identified patterns can be extremely difficult to interpret. Is something a previously unknown association, or purely statistical feature?
 - Rediscovery – Extensive number crunching and analysis sometimes identifies difference between summer and winter, other well known patterns
 - Is this changing?



Climate Informatics

- 2011 First International Workshop on Climate Informatics
New York Academy of Sciences
Climate Informatics Wiki launched
- 2013 "Climate Informatics" book chapter [M et al. 2013]
- 2015 Please join us in September as Climate Informatics turns 5!

www.climateinformatics.org

Figure courtesy C. Monteleoni

Climate Analysis

- Data standards and Interoperability
 - NetCDF and GRIB are dominant data standards
 - Readable by Matlab, Fortran, etc.
 - Widely used within research community
 - Metadata much less standardized
 - CMIP5 had very specific metadata requirements, for example, which required a great deal of effort to comply with
 - More issues outside of research community
 - Unprocessed climate model output often not what is need for a given application
- Privacy and Ethics
 - Only virtual animals suffer from simulated climate change
 - Data is generally openly if not easily available

What Are the Necessary Resources?

*What infrastructure, funding, and policies are needed to generate this data?
Does your project involve partnerships or other types of sustained
organizational relationships?*

- **Athena Project** (2009-2012) – long simulations with atmosphere-only models having various levels of spatial resolution, up to and including cloud-system resolving grids (Dedicated XT4 at NICS)
- **Minerva Project** (2012-2014) – seasonal predictions with coupled models having large-scale vs. mesoscale resolutions in the atmosphere and land surface (Dedicated Advanced Scientific Discovery on NCAR Yellowstone)

Project Athena

- 2008 World Modeling Summit: **dedicate petascale supercomputers to climate modeling**
- U.S. National Science Foundation **offered to dedicate the Athena supercomputer for 6 months** in 2009-2010 as a pilot study
- An **international collaboration** (*Project Athena*) was formed by groups in the U.S., Japan and the U.K. to use Athena to take up the challenge
- COLA, ECMWF, JAMSTEC, NICS, Cray





Project Athena Resources

- The Cray XT4 – **Athena** – the first NICS machine in 2008
 - 4512 nodes: AMD 2.3 GHz quad-core CPUs + 4 GB RAM
 - #30 on June 2009 Top 500 list
 - **18,048 cores** + 17.6 TB aggregate memory
 - **165 TFLOPS peak** performance
 - Dedicated to this project during October 2009 – March 2010 → **72 million core-hours!**
- Other resources made available to project:
 - **85 TB Lustre file system**
 - **258 TB auxilliary Lustre file system** (called *Nakji*)
 - *Verne*: **16-core** 128-GB system (data analysis) during production phase (2009-2010)
 - *Nautilus*: SGI UV with **1024 Nehalem EX cores**, 8 GPUs, 4 TB memory, 960 TB GPFS disk (data analysis) in 2010-11



Project Athena Resources

- The Cray XT4 – **Athena** – the first NICS machine in 2008
 - 4512 nodes: AMD 2.3 GHz quad-core CPUs + 4 GB RAM
 - **#30 on June 2009 Top 500 list**
 - **18,048 cores + 17.6 TB aggregate memory**
 - **165 TFLOPS peak performance**
 - Dedicated to this project during October 2009 – March 2010 → **72 million core-hours!**
- Other resources made available to project:
 - **85 TB Lustre file system**
 - **258 TB auxilliary Lustre file system (called *Nakji*)**
 - *Verne*: **16-core** 128-GB system (data analysis) during production phase (2009-2010)
 - *Nautilus*: SGI UV with **1024 Nehalem EX cores**, 8 GPUs, 4 TB memory, 960 TB GPFS disk (data analysis) in 2010-11



Project Minerva

- Opportunity to continue successful *Athena* collaboration between COLA and ECMWF, and to address limitations in the *Athena* experiments
- Explore the impact of **increased atmospheric resolution** on model fidelity and prediction skill in a ***coupled, seamless framework*** by using a state-of-the-art coupled operational long-range prediction system to systematically evaluate the prediction skill and reliability of a robust set of hindcast **ensembles** at low, medium and high atmospheric resolutions
- **NCAR Advanced Scientific Discovery Program** to inaugurate *Yellowstone* (72 K-core IBM iDataPlex)



**Many thanks to
NCAR for
resources and
sustained
support!**



Project Minerva Resources

- **NCAR Yellowstone**

- In 2012, NCAR-Wyoming Supercomputing Center (NWSC) debuted *Yellowstone*, the successor to *Bluefire*
- IBM iDataplex, 72,280 cores, **1.5 petaflops peak** performance
- #17 on June 2013 Top500 list
- **10.7 PB disk capability**
- High capacity HPSS data archive
- Dedicated large memory and floating point accelerator clusters (*Geyser* and *Caldera*)
- **10x increase in FLOPS, 100x increase in storage over Athena**

- **Accelerated Scientific Discovery (ASD) program**

- NCAR accepted a small number proposals for early access to Yellowstone, as it has done in the past with new hardware installs
- **3 months of near-dedicated access before being opened to general user community**
- Allocated 21 M core-hours on Yellowstone
- **Used ~28 M core-hours** (Our jobs squeaked in under core size that “broke” the system)
- Allocated 250 TB... then 400 TB.... then 500 TB



Project Minerva Resources

- **NCAR Yellowstone**

- In 2012, NCAR-Wyoming Supercomputing Center (NWSC) debuted *Yellowstone*, the successor to *Bluefire*
- **IBM iDataplex, 72,280 cores, 1.5 petaflops peak performance**
- **#17 on June 2013 Top500 list**
- **10.7 PB disk capability**
- High capacity HPSS data archive
- Dedicated large memory and floating point accelerator clusters (*Geyser* and *Caldera*)
- **10x increase in FLOPS, 100x increase in storage over Athena**

- **Accelerated Scientific Discovery (ASD) program**

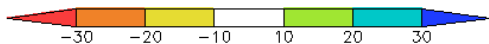
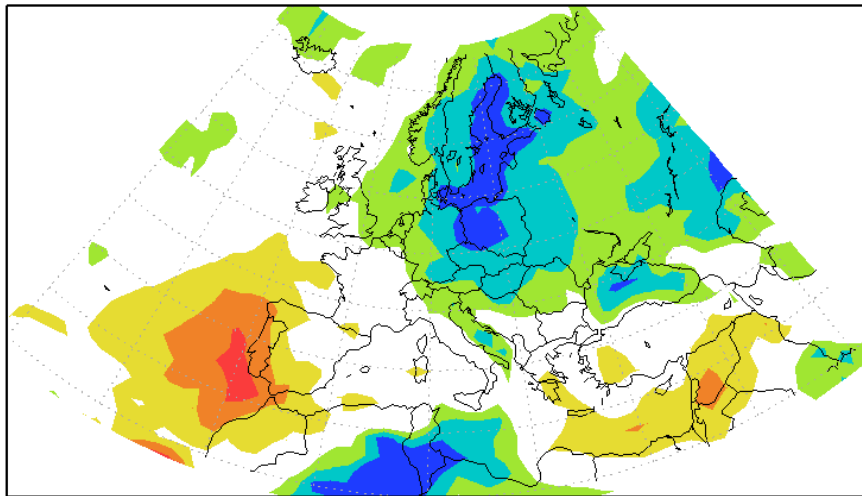
- NCAR accepted a small number proposals for early access to Yellowstone, as it has done in the past with new hardware installs
- **3 months of near-dedicated access before being opened to general user community**
- Allocated 21 M core-hours on Yellowstone
- **Used ~28 M core-hours** (Our jobs squeaked in under core size that “broke” the system)
- Allocated 250 TB... then 400 TB.... then 500 TB

What Have We Learned?

*How has your big data work changed your field?
What advice do you have for others running big data projects?*

CNRM

CNRM Change in Precip 2071–2100 Minus 1971–2000
apr–sep

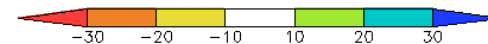
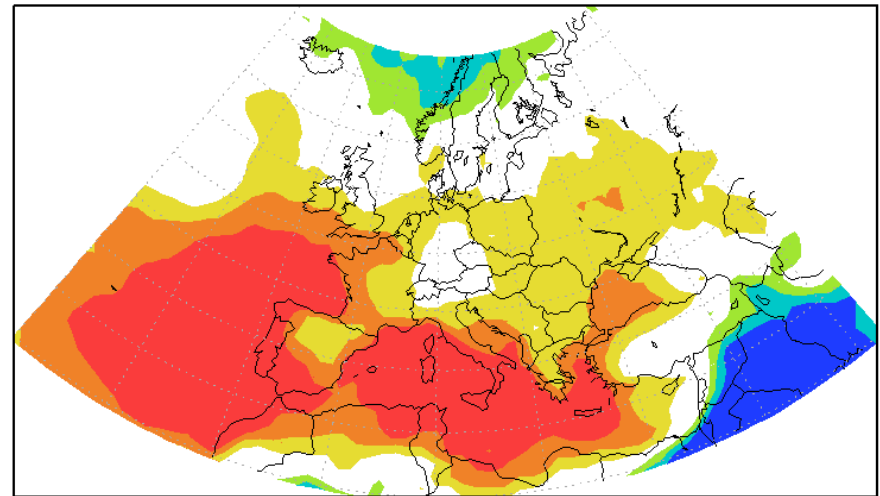


GRADS: COLA/IGES

2014-03-17-00:38 GRADS: COLA/IGES

CESM

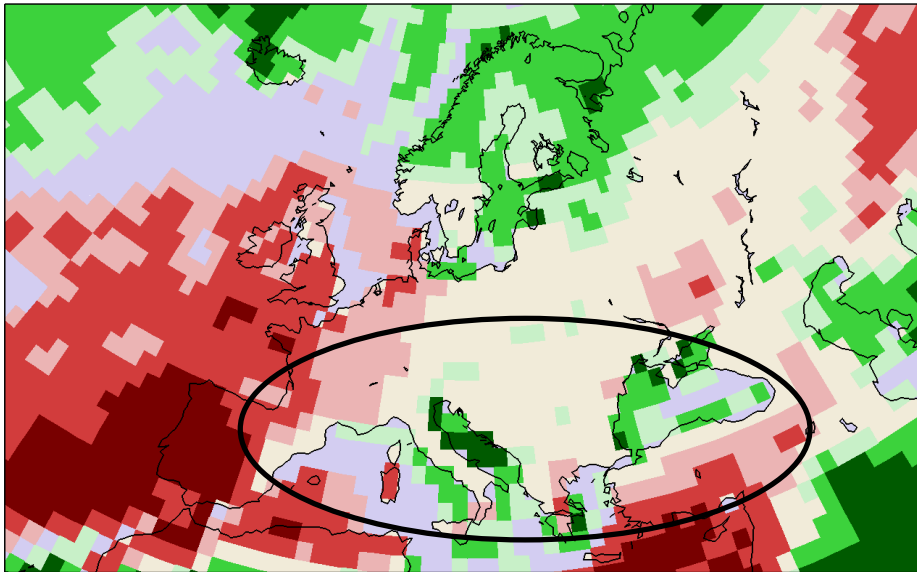
CESM Percent Change in Precip 2071–2100 Minus 1971–2000
apr–sep



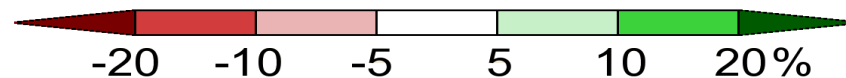
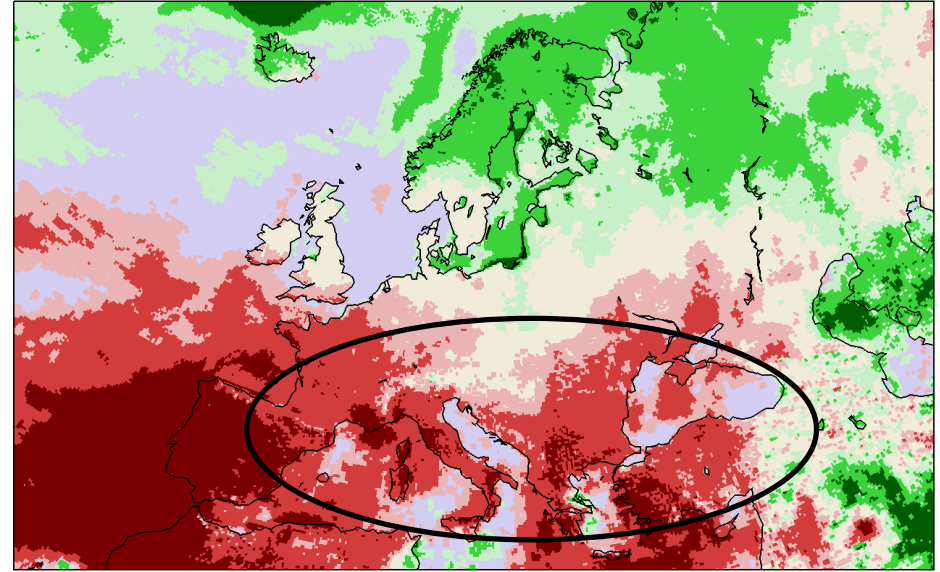
2014-03-17-00:39

Europe Growing Season (Apr-Oct) Precipitation Change: 20th C to 21st C

T159 (125-km)



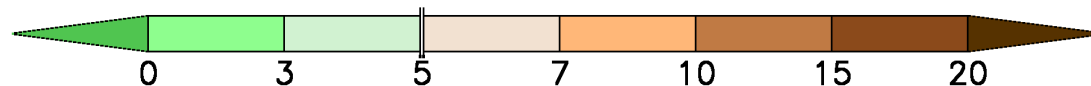
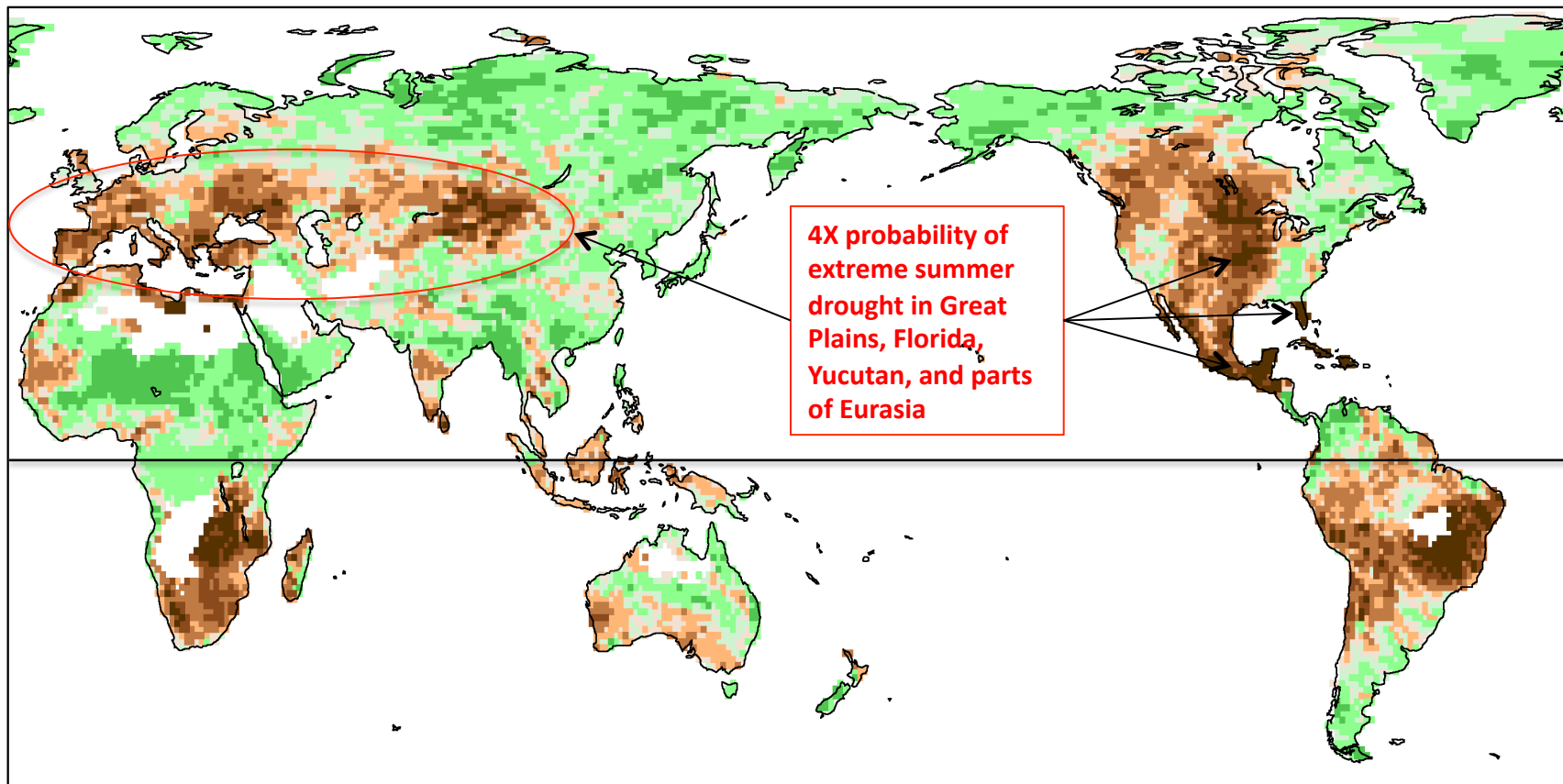
T1279 (16-km)



“Time-slice” runs of the ECMWF IFS global atmospheric model with observed SST for the 20th century and CMIP3 projections of SST for the 21st century at two different model resolutions

The continental-scale pattern of precipitation change in April – October (growing season) associated with global warming is similar, but the regional details are quite different, particularly in southern Europe.

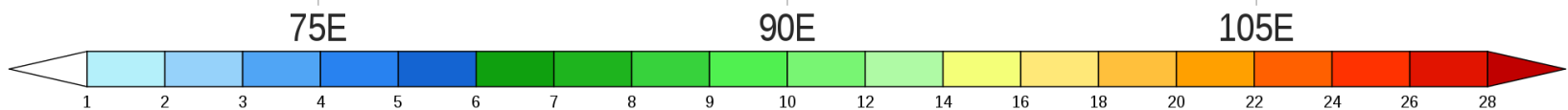
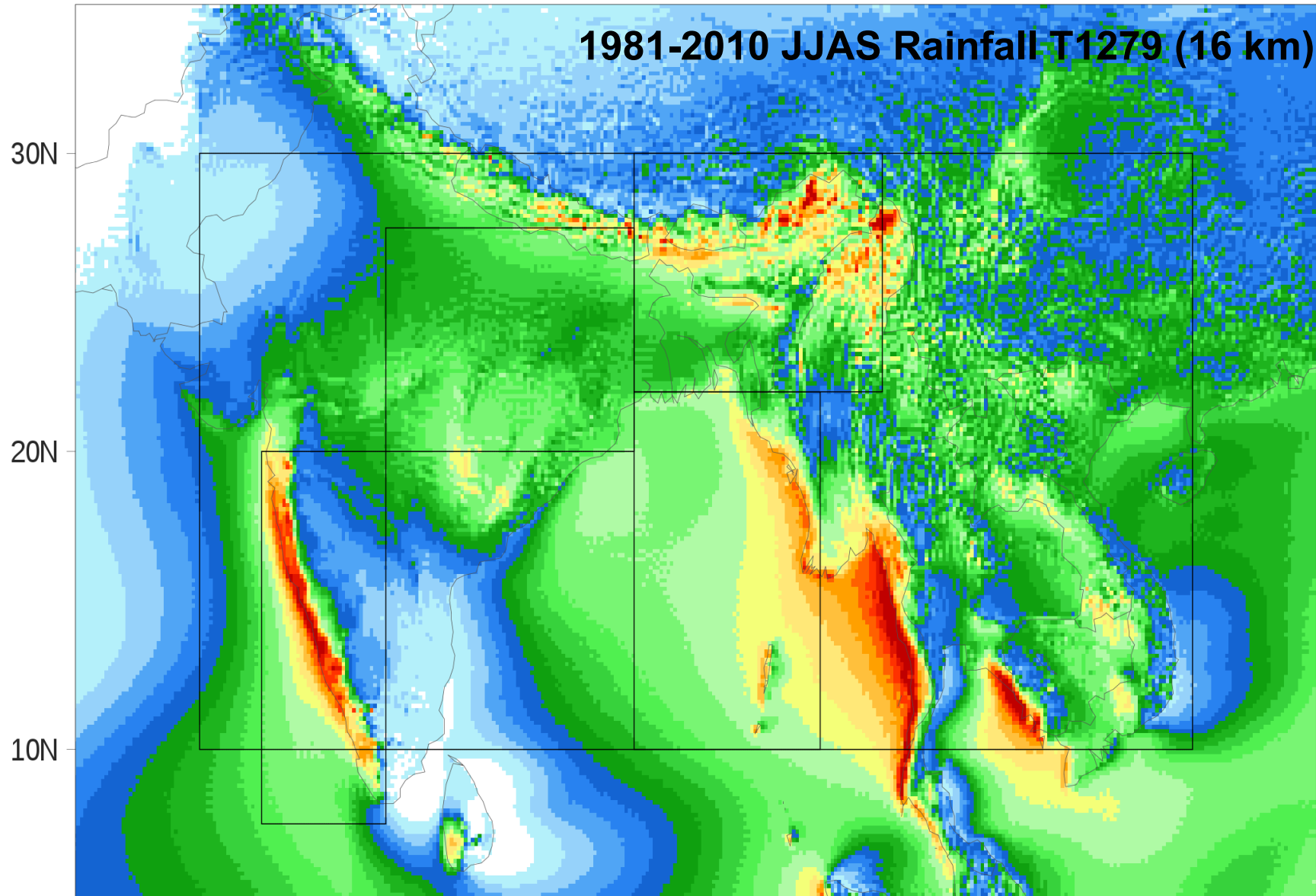
Future Change in Extreme Summer Drought Late 20th C to Late 21st C

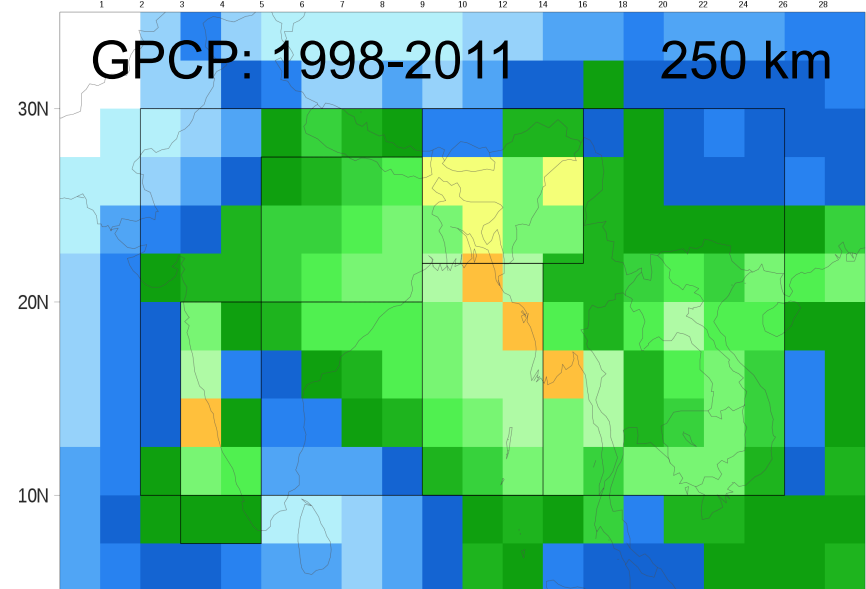
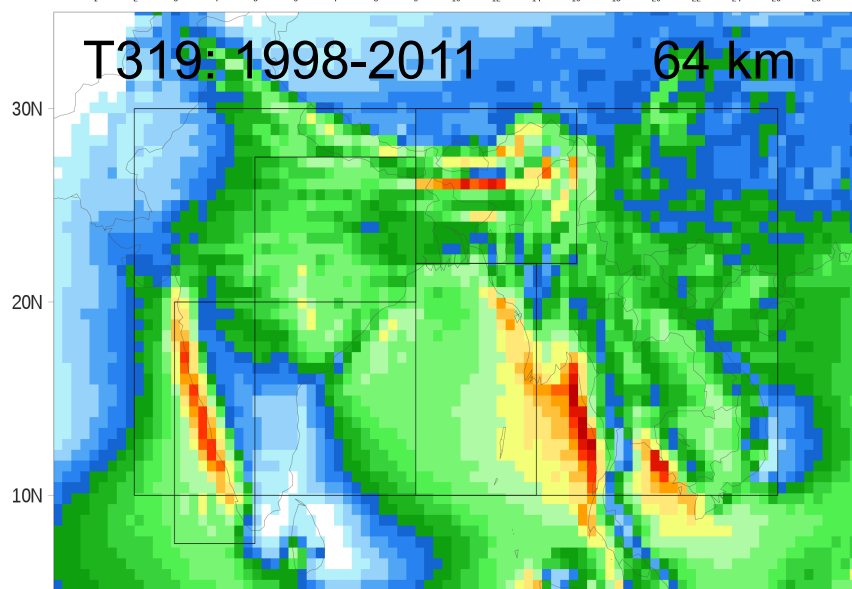
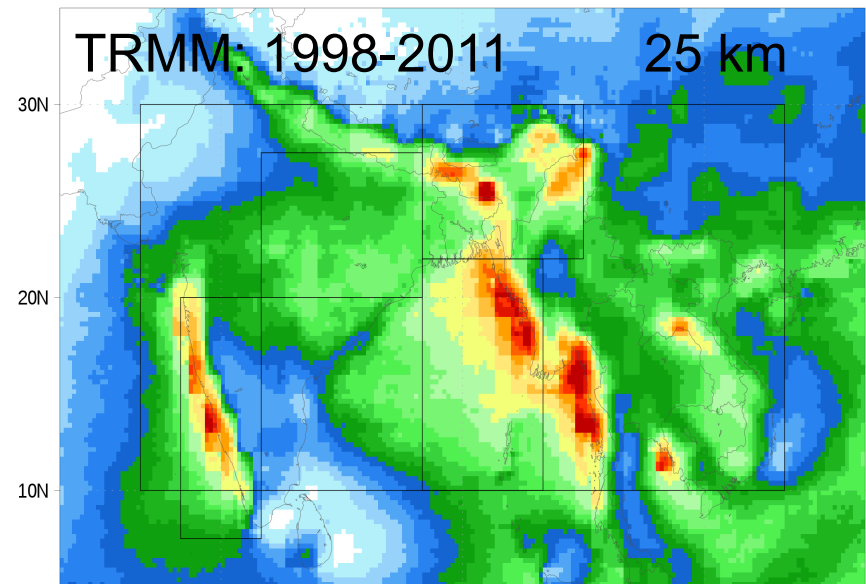
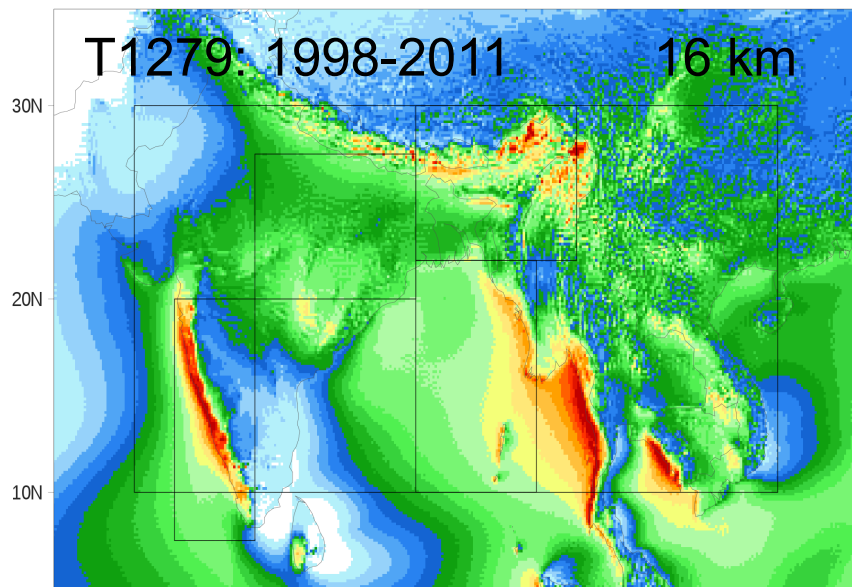


10th Percentile Drought: Number of years out of 47 in a simulation of future climate (2071-2117) for which the June-August mean rainfall was less than the 5th driest year of 47 in a simulation of current climate (1961-2007).

South Asian Monsoon in Minerva

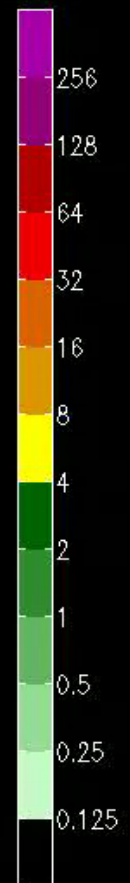
1981-2010 JJAS Rainfall T1279 (16 km)



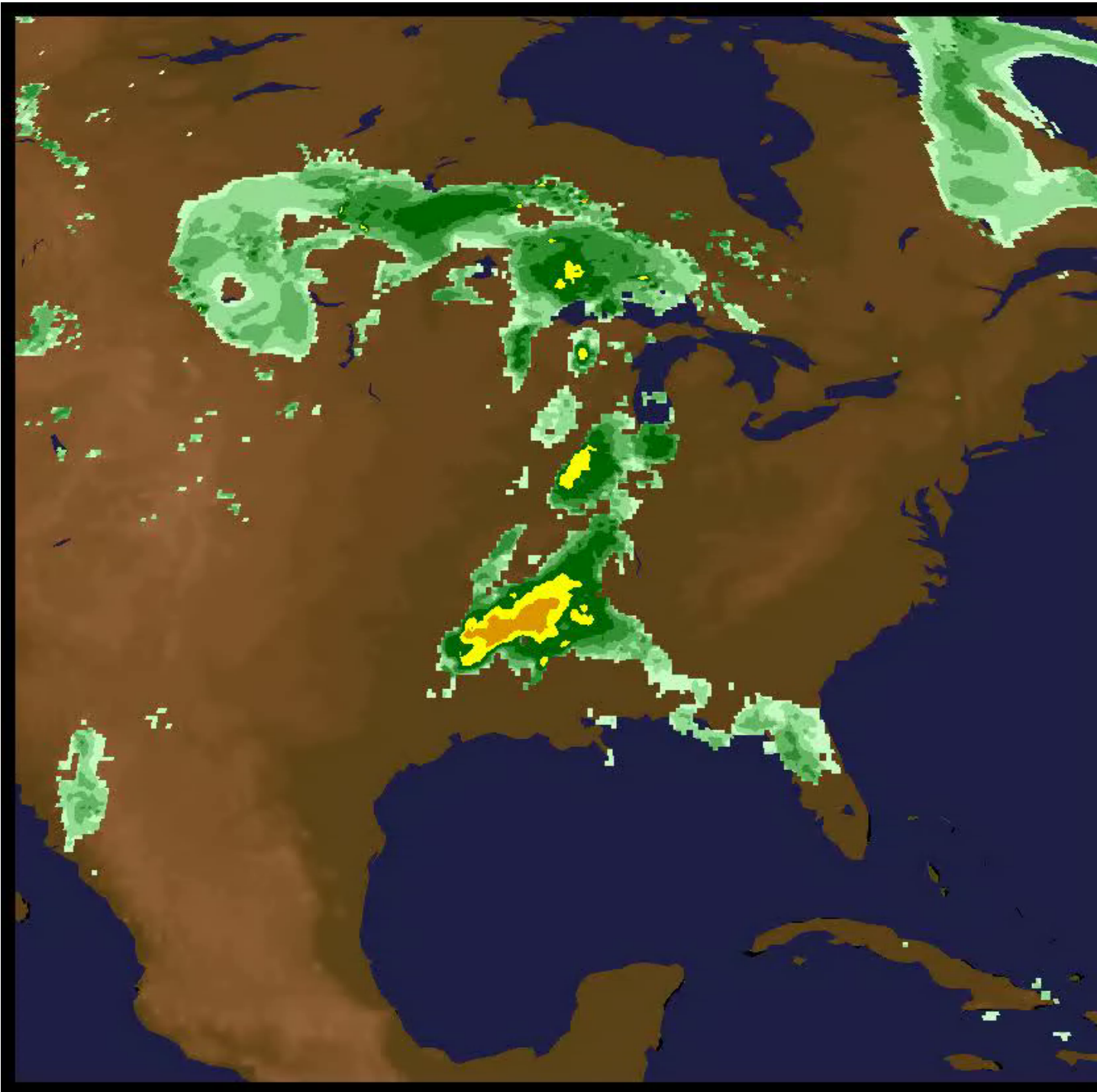


02Z MAY 01, 2010

Precipitation
mm/hr



**Precipitation:
Summer 2010
(IFS 16km)**



What Have We Learned?

- Many features of atmospheric circulation improve substantially with global models having improved spatial resolution up to at least a 16-km grid spacing (mesoscale-permitting). **Athena**
- Some (mostly tropical) features of atmospheric circulation are insensitive to spatial resolution, with current parameterizations. **Athena, Minerva**
- Realism of atmospheric mesoscale features in the extratropics substantially and significantly improves coupled seasonal forecast skill, both deterministic and probabilistic. **Minerva**
- Resolving (or at least permitting) mesoscale ocean eddies significantly improves many features of coupled simulations. (**PetaApps**, not shown)
- Significantly better simulations of ISI variability can be achieved with improved representation of convection. (**SP-CCSM**, not shown)
- **Validating high-resolution, high-complexity data pushes and in some cases exceeds observational capabilities (All)**
- **Spatial resolution alone is not a panacea**, and it is expensive
 - ~100X the HPC cost of conventional resolution models.
 - Time-to-solution constraints drive demand for greater parallelism.
 - Output volume and intensity ~30-40X compared to conventional resolution models.

What Have We Learned?

- **Dedicated usage** of a relatively big supercomputer **greatly enhances productivity**
 - **Experience with Athena and ASD period demonstrates tremendous progress can be made with dedicated access**
- Dedicated computing campaigns provide demonstrably **more efficient utilization**
 - **Noticeable decrease in efficiency once scheduling multiple jobs of multiple sizes was turned over to a scheduler**
- In-depth exploration
 - Data saved at much higher frequency
 - Multiple ensemble members, increased vertical levels, etc.
- Dedicated simulation projects like *Athena* and *Minerva* **generate enormous amounts of data** to be archived, analyzed and managed.
 - **Data management is a tremendous challenge.**
 - Other than machine instability, data management and post-processing were solely responsible for halts in production.

Resource Summary

- High Performance Computing
 - Climate applications can occupy arbitrarily large numbers of cores for the foreseeable future
 - Personally occupied 60,000+ cores on Yellowstone during ASD
 - Highest resolution in *Athena* and *Minerva* is still far from where we want to be
 - Convection still parameterized – 1km and below is target
 - From 16km to 1k is a 4000x increase in required computing power
 - Exascale machines will be necessary
 - Codes will need to scale to 10^5 - 10^6 cores
 - Codes will need to be fault tolerant as odds of core failure go to 1
 - Significant computer science challenges ahead

Resource Summary

- High Performance Storage
 - HPC for *Athena* and *Minerva* were sufficient to generate massive, cutting edge data sets
 - Initial storage capability for *Athena* was entirely inadequate
 - *Minerva* (Yellowstone) had vastly increased storage capacity
 - Design of Yellowstone partially informed by *Athena* experience
 - remained the main limiter during the computational phase
 - Archival systems are critical, but use should be avoided
 - Data on tape is not available for analysis
 - Moving data in and out of archive has been largest single bottleneck
 - Centers must be designed with disk and analysis capacity commensurate to their computational level
 - Data storage at this level is beyond individual research groups
 - Must be treated as something to be provided along with FLOPS

Resource Summary

- High Performance People
 - Big data projects require large investments of human hours
 - Only large institutions might have all necessary resources in house
 - Model development
 - Run-time management
 - Data management
 - Data analysis
 - Significant commitment of personnel
- Always more data than eyeballs
 - Sheer volume of data to create, manage, and analyze is strong incentive for collaborations

Looking to the Future: Coping with the “Exaflood”

- Current challenge is managing petabytes
- Higher resolution and complexity will push this to exabytes
- We will need systematic, repeatable data solutions
 - Climate scientists currently handle Big Data with largely ad hoc solutions
- Some methods that can help:
 - Data compression: 2-3X without loss of “science” content
 - Cannot allow compression/decompression steps to overwhelm time to solution
 - Remote server-side analysis: Analyze the data where they reside
 - Avoid moving multi-TB over WAN
 - This has been tried with limited success – scalability, data security, multi-site storage, familiar diagnostics are harder to obtain at very high resolution
 - Workflow management
 - In-line post-processing
 - Automatic generation of metrics and diagnostics
 - Automatic generation of visualizations

Looking to the Future: Coping with the “Exaflood”

- Have we wrung all the “science” out of the data sets, given that we can only keep a small percentage of the total data volume on spinning disk? **How can we tell?**
- Must move from ad hoc solutions to systematic, repeatable solutions
- Transform Noah’s Ark → a Shipping Industry

“We need exaflood insurance.”

- Jennifer Adams (COLA)

