

# Toward Visual Analysis of Ensemble Data Sets

Or,

# You want to render what?

2009 Ultrascale Visualization Workshop

November 16, 2009

Andy Wilson Sandia National Laboratories Kristi Potter

SCI Institute, Univ. of Utah



Sandia is a multiprogram laboratory operated by Sandia Corporation, a Lockheed Martin Company, for the United States Department of Energy's National Nuclear Security Administration under contract DE-AC04-94AL85000.











# **Ensembles are...**

### • Large

- Tens of gigabytes to hundreds of terabytes

### Multivariate

- Typically 10s to 100s of state variables

# Time-varying

### Multivalued

- Think of it as PDF-valued instead of scalar-valued

### • Awkward

Raw data is frequently discarded in favor of an Excel spreadsheet



# **Ensembles help mitigate uncertainty**

#### Multiple models

- Incorporate strengths of different approaches

### • Multiple runs

- Sample an input space of uncertain parameters
- Perturb measured inputs to mitigate model/measurement error

# Multiple grids

Evaluate and demonstrate convergence

# Multiple values

- Reason about the most likely simulation outcomes





# **A Few Examples**

#### • NOAA/NCEP Short-Term Reference Ensemble

- Weather for North America
- 4 models, 21 members, 624 state variables, 30 timesteps, 36GB/run, 3 runs/day

#### Climate Simulations (Earth System Grid)

- Worldwide climate over millenia
- 30TB+ repository at LLNL, lots more elsewhere
- Varying simulation domains

#### • Parameter studies for uncertainty quantification

- Engineered systems under stress
- Weather/climate data makes a good proxy





# **Driving Questions**

• What conditions are predicted by this ensemble?

• Where and when do those conditions occur?

• What is the relative probability of some outcome?



# Major Research Issues: Data Management

• Key Insight: The user only ever needs a tiny subset of the data -- but that subset changes frequently.

#### • Examples:

- What phenomena does the ensemble predict?
  - This is usually derived or inferred from the data
- Where and when will it happen?
  - Moral equivalent of an SQL WHERE clause
- What is the relative probability of X?
  - Derive this from "where and when" by aggregating over ensemble members

This calls for data stores with database-like access and query semantics.

#include <bill\_howe\_cloud\_vis.ppt>



# **Major Research Issues: Many-Valued Data**

- Spatial PDF visualization does not (yet) appear necessary
  - Summary statistics + drill-down suffices
  - Even that much is difficult

#### The world is often not Gaussian

- Beware of mean + standard deviation!
- Watch out for multimodal distributions

#### • There's Just Too Much Data

- "Display it all and let the analyst browse" doesn't work
- Query-driven visualization becomes very important here
- We may not be able to use the supercomputers "just for vis"





# **Our Approach: Data Store**

#### Netezza NPS data warehouse appliance

- Parallel database with CPU right next to disk
- Schema exposes data values directly to the database
- All numeric queries go through SQL

#### • Research into data stores continues

- Column-store database?
- Numeric index like FastBit?
- Reordering data for multiresolution access?
- Cloud storage and processing?
- Assumption: We can move/index/reorganize the data at least once (possibly as it's being generated)





### There Is No Perfect Display

– Many displays, all linked

### An ensemble data set has several dimensions:

- Time (I)
- Space (2, 2.5 or 3)
- State Variable (1 to 1000)
- Ensemble Member (1 to 1000)

### Collapse dimensions to yield 2- or 3D display

- "Collapse" means "extract or aggregate" here





# **Summary Display**





Color = Mean

Color = Std. Dev.







# Spaghetti Plot (I)







# Spaghetti Plot (2)







# **Filmstrip Display**



- Small multiples of summary display
- Show one variable over many timesteps at low resolution
- Cameras linked with other frames
- Selection linked with other frames and other displays





#### **Quartile Chart**

![](_page_15_Figure_2.jpeg)

 Shows quick overview of distribution and clustering of data values over time

![](_page_15_Picture_4.jpeg)

![](_page_16_Picture_0.jpeg)

### **Trend Chart**

![](_page_16_Figure_2.jpeg)

![](_page_16_Picture_3.jpeg)

![](_page_17_Picture_0.jpeg)

# Discussion

#### • Exactly what should we display?

- How should we display it?

#### • What statistics are appropriate?

- Do we have enough data to support them?

- How can we indicate missing data?
- How do we store and access the data?

![](_page_17_Picture_8.jpeg)

![](_page_18_Picture_0.jpeg)

# **Future Work**

# •3D (not as difficult as it appears)

...once that's done...

- Better display metaphors
- Data fusion across different grids and time domains

![](_page_18_Picture_6.jpeg)

![](_page_19_Picture_0.jpeg)

# Acknowledgements

- Kristi Potter, Univ. of Utah
- NOAA/NCEP
- Dean Dobranich, SNL
- Valerio Pascucci, Univ. of Utah
- Chris Johnson, Univ. of Utah
- Peer-Timo Bremer, LLNL

![](_page_19_Picture_8.jpeg)

![](_page_20_Picture_0.jpeg)

# **Questions?**

![](_page_20_Picture_2.jpeg)