

# HUMAN-CENTERED DATA UNDERSTANDING

---

Chad A. Steed | <http://csteed.com>  
Oak Ridge National Laboratory  
Computer Science and Mathematics Division

CDA Group Meeting  
July 18, 2018

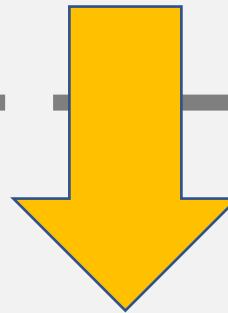
Computers are faster,  
Memory is cheap,  
Data volumes are increasing,  
And **AI** is advancing.

*Do we need a human in the loop?*

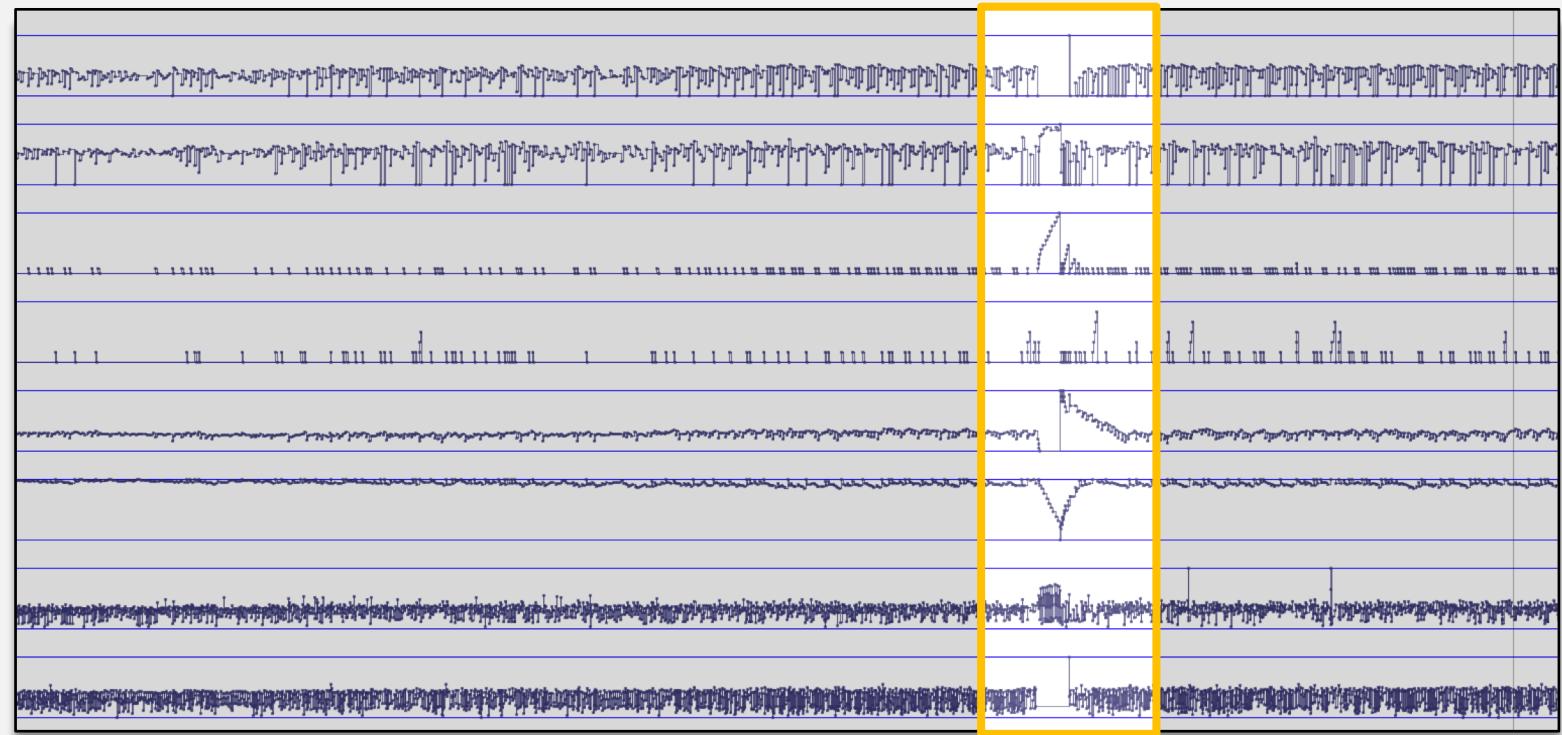


**Colonel Stanislav Petrov**  
(image courtesy *The New York Times*)

# DATA

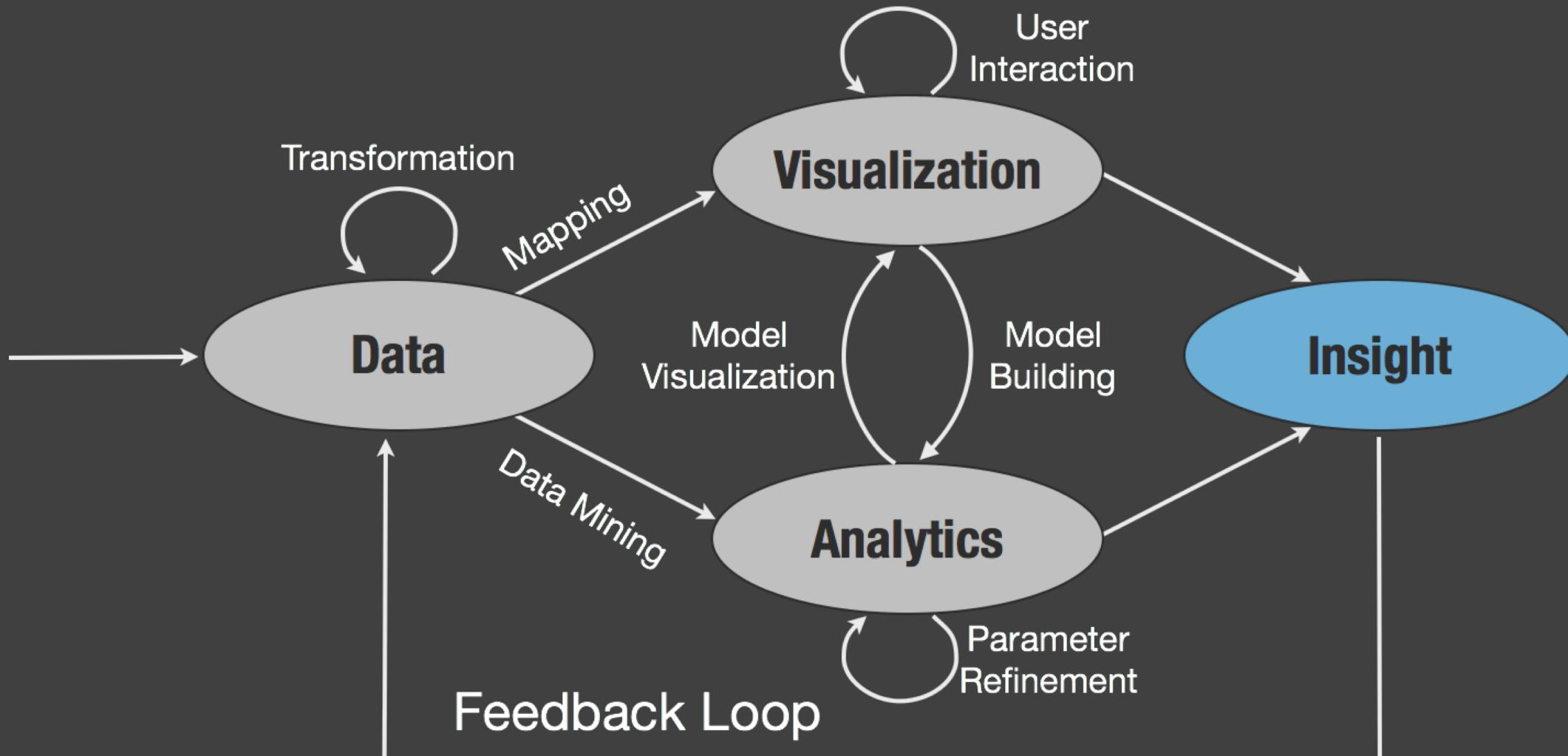


# Insight



492258	2016-06-17 21:07:03.466 OPC.Temperature.ExtraTemp1 SuperUser (OPC) 9526586 618
492259	2016-06-17 21:07:03.466 OPC.Temperature.BottomTemperatureValidation SuperUser (OPC) 9526586 618.1999999999982
492260	2016-06-17 21:07:03.486 Process.CathodeTuningControl.CathodePower  [OnChange( OPC.PowerSupply.Filament.VoltageFB)] Arcam.EBMControl.Process.CathodeTuningControl.OnCathodPowerChange() (Logic) 9526586 5.929202
492261	2016-06-17 21:07:03.516 Process.CathodeTuningControl.MeanCathodePower [OnChange(Process.CathodeTuningControl.CathodePower)] Arcam.EBMControl.Process.CathodeTuningControl.MeasureMeanPower() (Logic) 9526587 5.905572
492262	2016-06-17 21:07:03.516 OPC.PowerSupply.HighVoltage.SafetySignal [OnPositiveFlank(SafetySignalTimer.Timeout)] Arcam.EBMControl.Process.HighVoltageControl.OnTimeToSendSafetySignal() (Logic) 9526587 True

# Visual Analytics

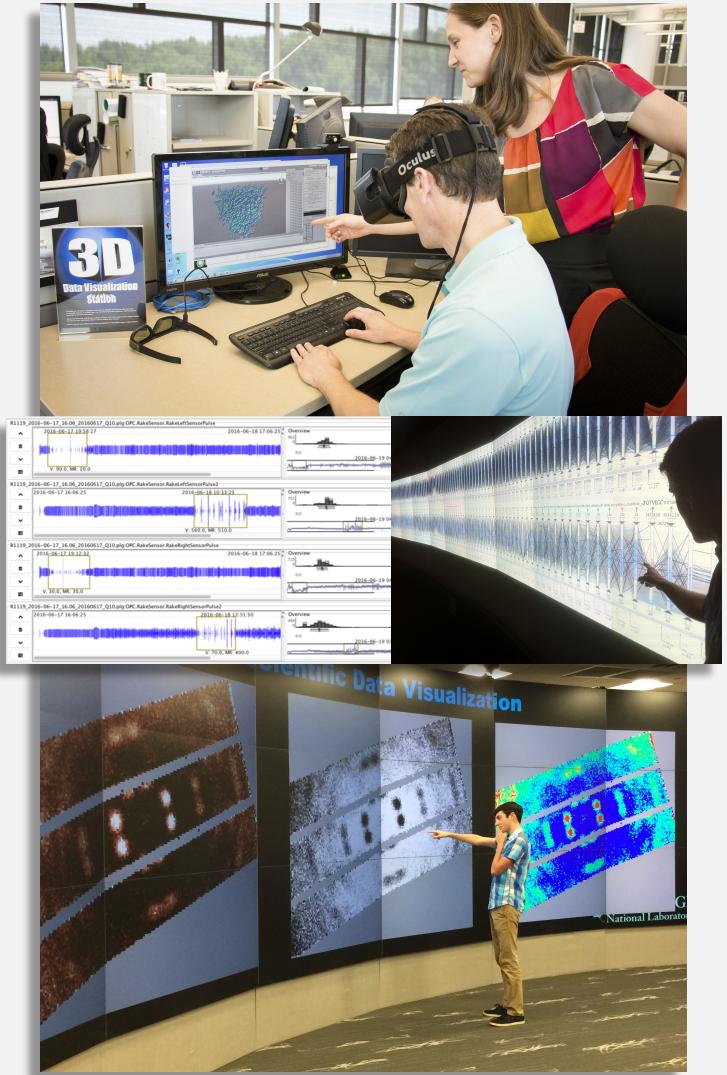


*After Keim et al., Mastering the Information Age: Solving Problems with Visual Analytics, 2011*



# My Research Focus

- **GOAL:** Improve human capabilities to discover, comprehend, and communicate insight in large and complex data through the design, implementation, and practical evaluation of visual analytics techniques
- **Specific Topics:**
  - Graphical representations of data
  - Human interactions with visualizations
  - Integrating automated machine learning



# Agile, Interdisciplinary Data Science Research

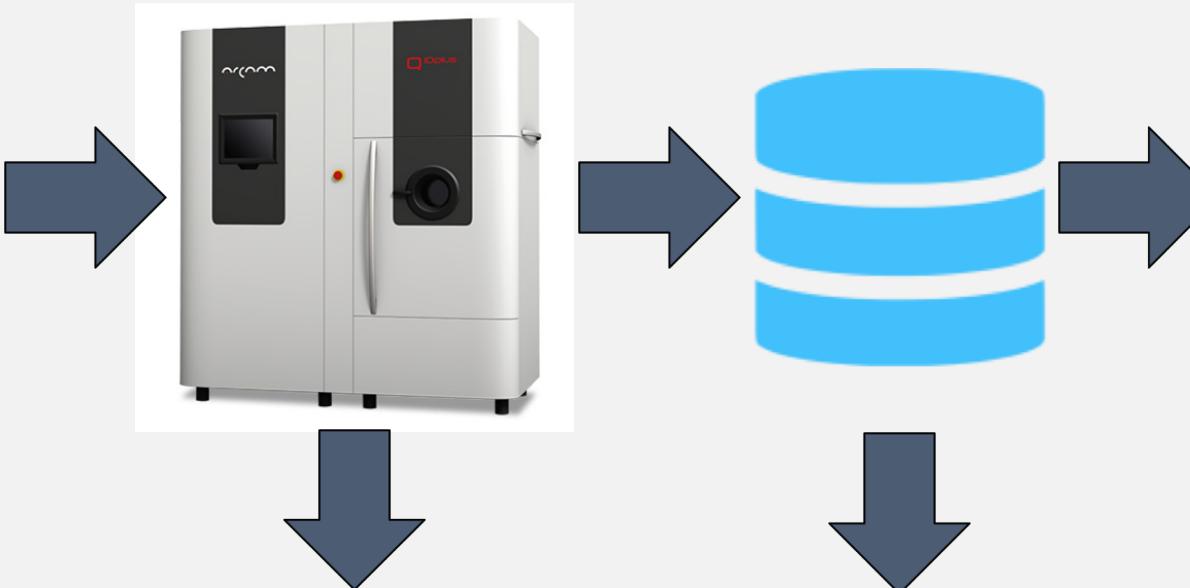
## Domain Experts

Dr. Ryan Dehoff  
Additive Manufacturing



## Scientific Domains

Large-scale 3D Printing



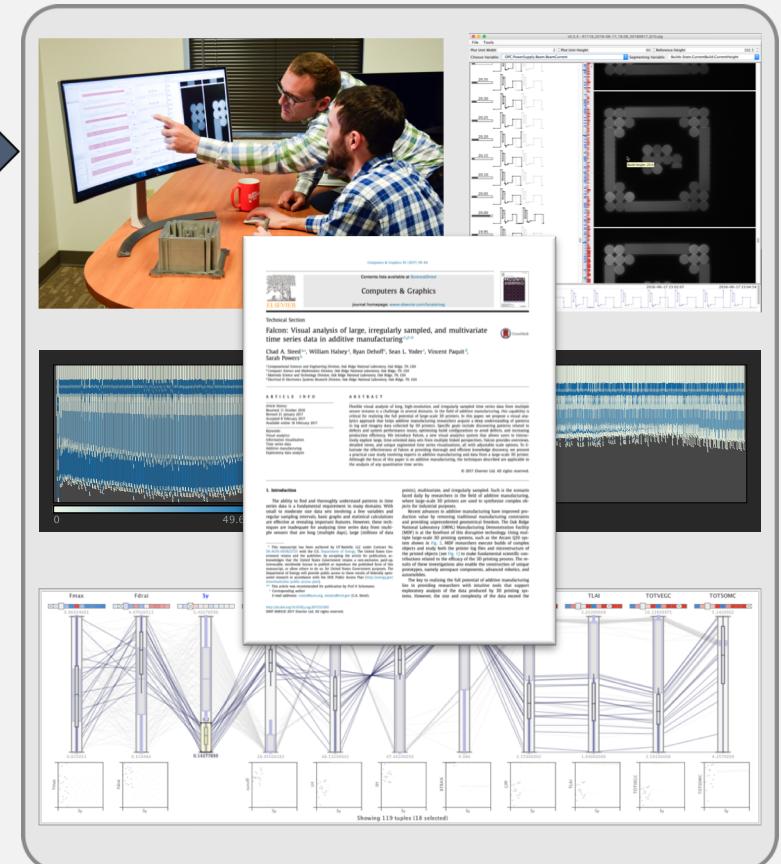
## Real-world Data

Large, Complex, Unique

492258 2016-06-17 21:07:03.466 [OPC.Temperature.ExtraTemp1|SuperUser (OPC)|9526586|618  
492259 2016-06-17 21:07:03.466 [OPC.Temperature.BottomTemperatureValidation|SuperUser (OPC)|9526586|618\_1.3999999999999992  
492260 2016-06-17 21:07:03.466 [Process.CathodeTuningControl.CathodePower|OnChange|  
OPC.PowerSupply.Filament.VoltageB|] Arcam.EBCControl.Process.  
CathodeTuningControl.OnCathodePowerChange()| (Logic)|9526586|5.92902  
492261 2016-06-17 21:07:03.516 [OPC.Temperature.ExtraTemp1|SuperUser (OPC)|9526586|618  
492262 2016-06-17 21:07:03.516 [OPC.PowerSupply.HighVoltage.SafetySignal|]  
OnChange| (Process.CathodeTuningControl.CathodePower)| Arcam.EBCControl.Process.  
CathodeTuningControl.MeasureMeanPower()| (Logic)|9526587|5.905572  
492263 2016-06-17 21:07:03.546 [OPC.Temperature.ExtraTemp1|SuperUser (OPC)|9526589|618  
492264 2016-06-17 21:07:03.546 [OPC.Temperature.ExtraTemp1|SuperUser (OPC)|9526589|598  
492265 2016-06-17 21:07:03.546 [OPC.Temperature.BottomTemperatureValidation|SuperUser (OPC)|9526589|598\_1.3999999999999991  
492266 2016-06-17 21:07:03.666 [OPC.Table.CurrentFeedback|SuperUser (OPC)|9526590|0.  
1153  
492267 2016-06-17 21:07:03.766 [OPC.PowerSupply.SmokeDetector.Counts|SuperUser (OPC)|9526594|044  
492268 2016-06-17 21:07:03.766 [OPC.Temperature.ExtraTemp1|SuperUser (OPC)|9526594|544  
492269 2016-06-17 21:07:03.766 [OPC.Temperature.BottomTemperatureValidation|SuperUser (OPC)|9526594|544\_1.4000000000000009  
492270 2016-06-17 21:07:03.786 [OPC.InternalCooling.DifferentialPressureOverFilter|  
OnChange| (OPC.InternalCooling.PressureBeforeFilter\_Unfiltered)| Arcam.  
EBCControl.Process.InternalCooling.UpdateDifferentialPressureOverFilter()| (Logic)|9526594|08296734

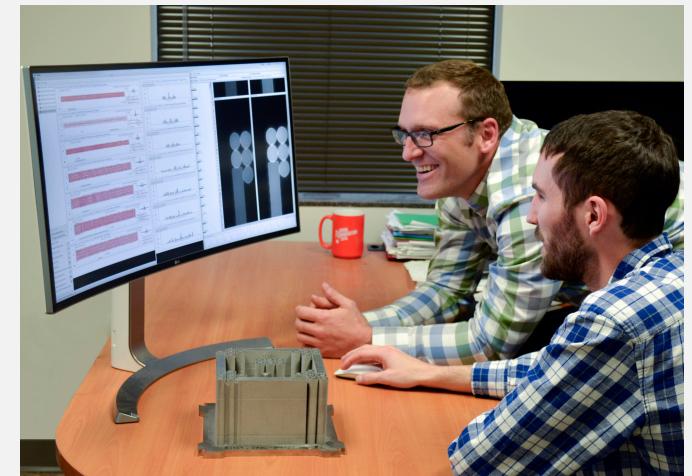
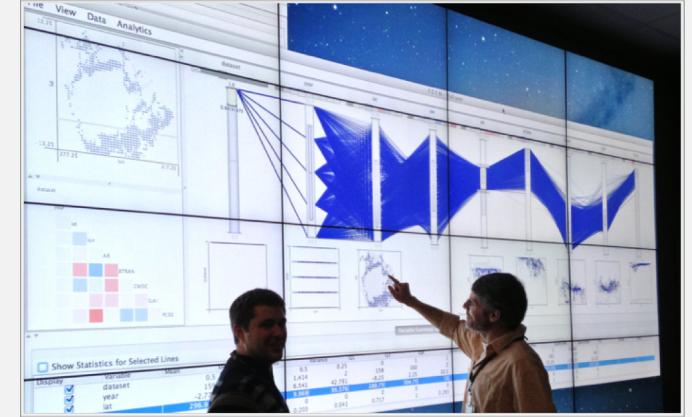
## Results / Artifacts

Tools, Techniques,  
Publications, Patents,  
Open Source Software



# Development Process

- **Domain experts are co-designers**
  - Participatory design
  - I spend time in learning domain challenges
  - They help design new features and evaluate releases
  - Domain experts keep me focus and they are invested in the process
  - We publish results together and share project artifact
- **Iterative development and evaluation cycle**
  - Design - Implement - Evaluate
  - Practical evaluations with real world scenarios
  - Proof of principle is new discoveries enabled by the techniques and transition to practice
  - Empirical evaluations are possible, but difficult at a national lab

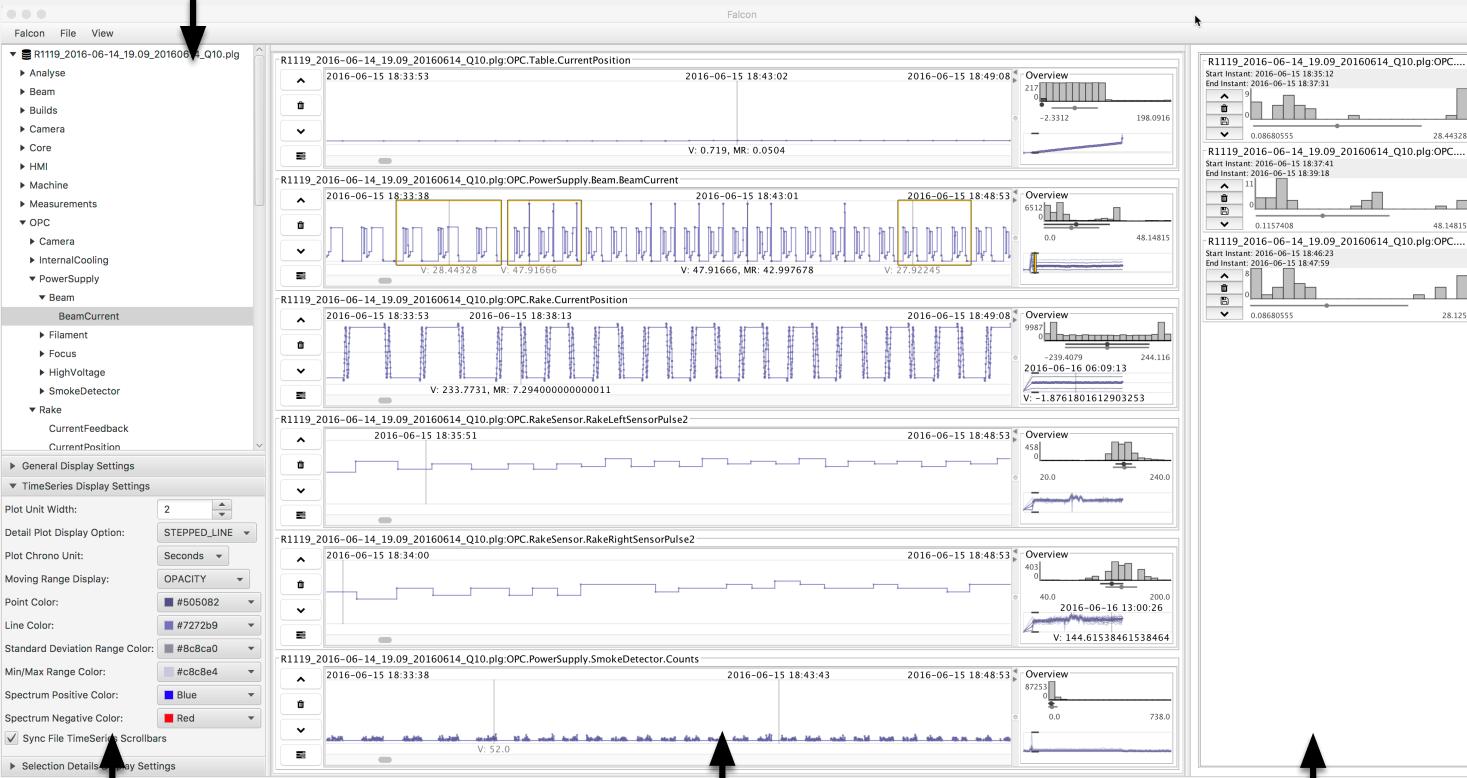


# Falcon

## Visualizing Long, Irregularly Sampled Sensor Streams

### Main Analysis Window

#### File / Variables Tree View

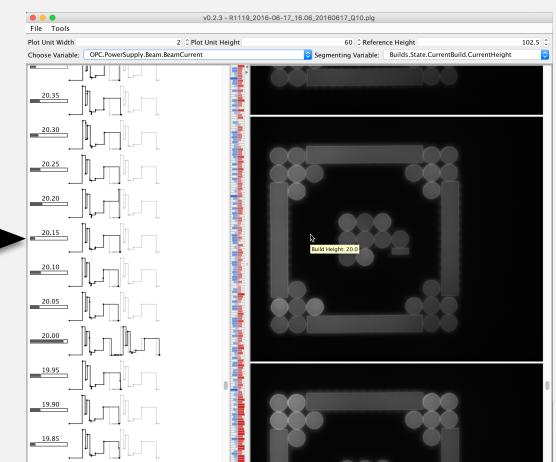
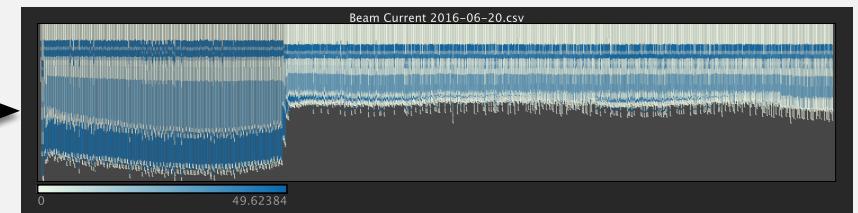


Settings Panel

Variable Visualization Panel  
(Left: detailed time series, Right: overview)

Selection Details Panel

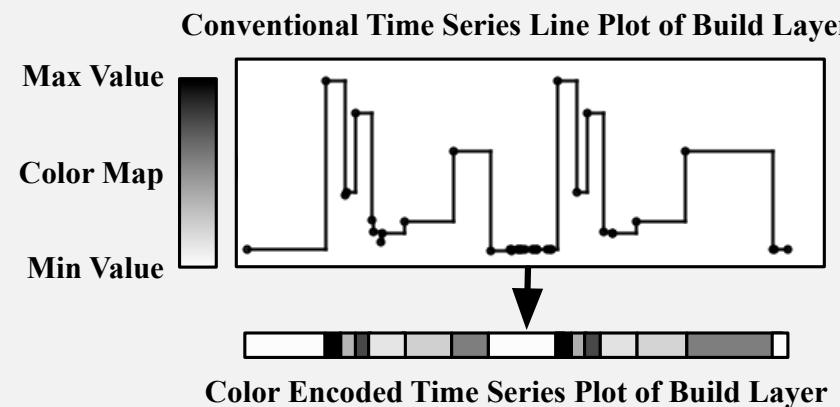
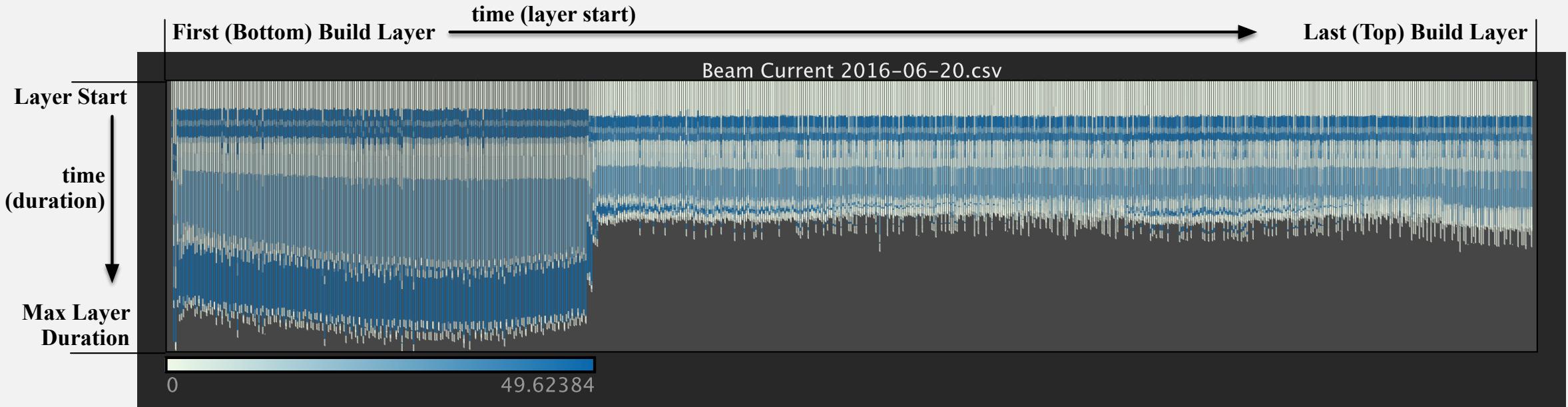
#### Waterfall Visualization



Segmented Time Series Visualization

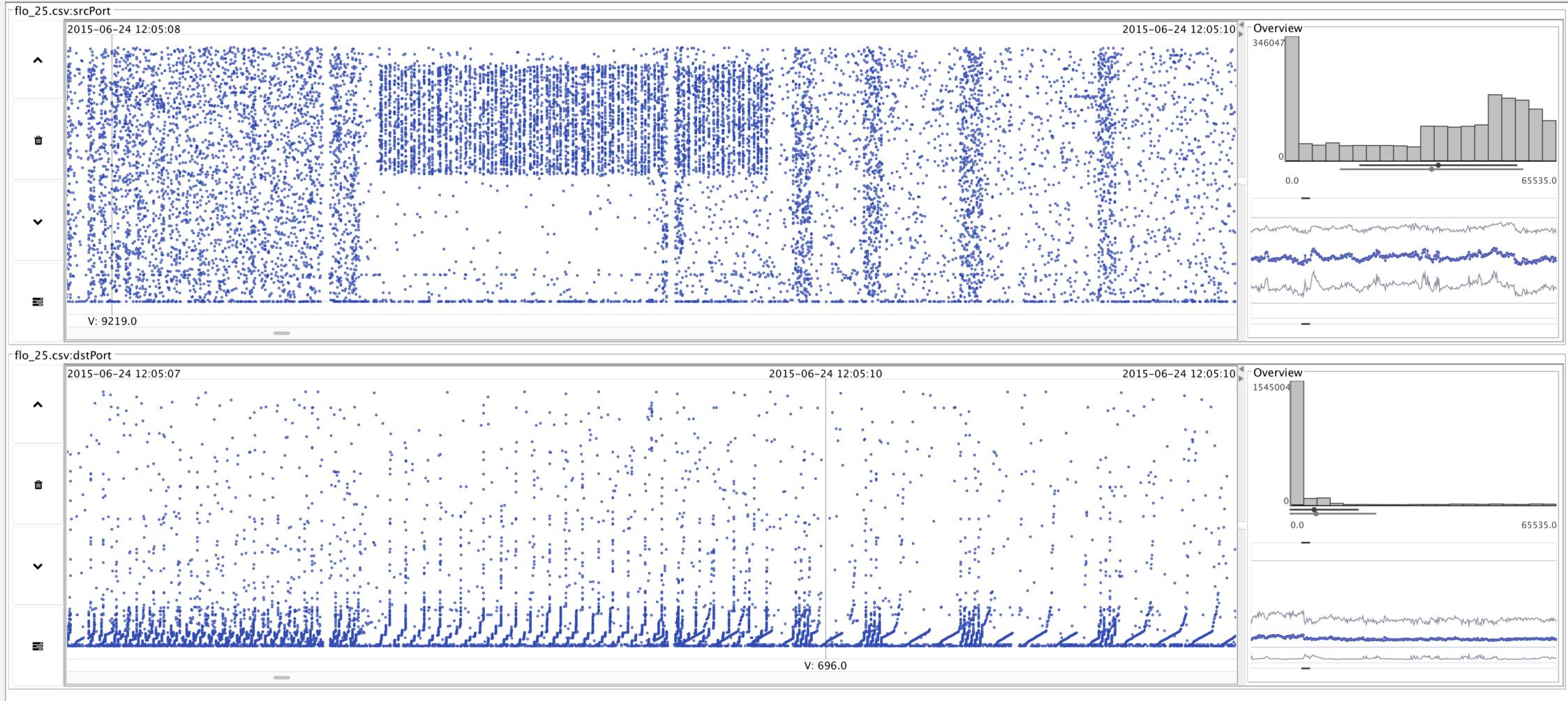
**Citation:** "Falcon: Visual Analysis of Large, Irregularly Sampled, and Multivariate Time Series Data in Additive Manufacturing". C. A. Steed, W. Halsey, R. Dehoff, S. L. Yoder, V. Paquit, and S. Powers. *Computers & Graphics*, 63:50--64, 2017.

# Falcon: Segmented Time Visualization



**Citation:** "Falcon: Visual Analysis of Large, Irregularly Sampled, and Multivariate Time Series Data in Additive Manufacturing". C. A. Steed, W. Halsey, R. Dehoff, S. L. Yoder, V. Paquit, and S. Powers. *Computers & Graphics*, 63:50--64, 2017.

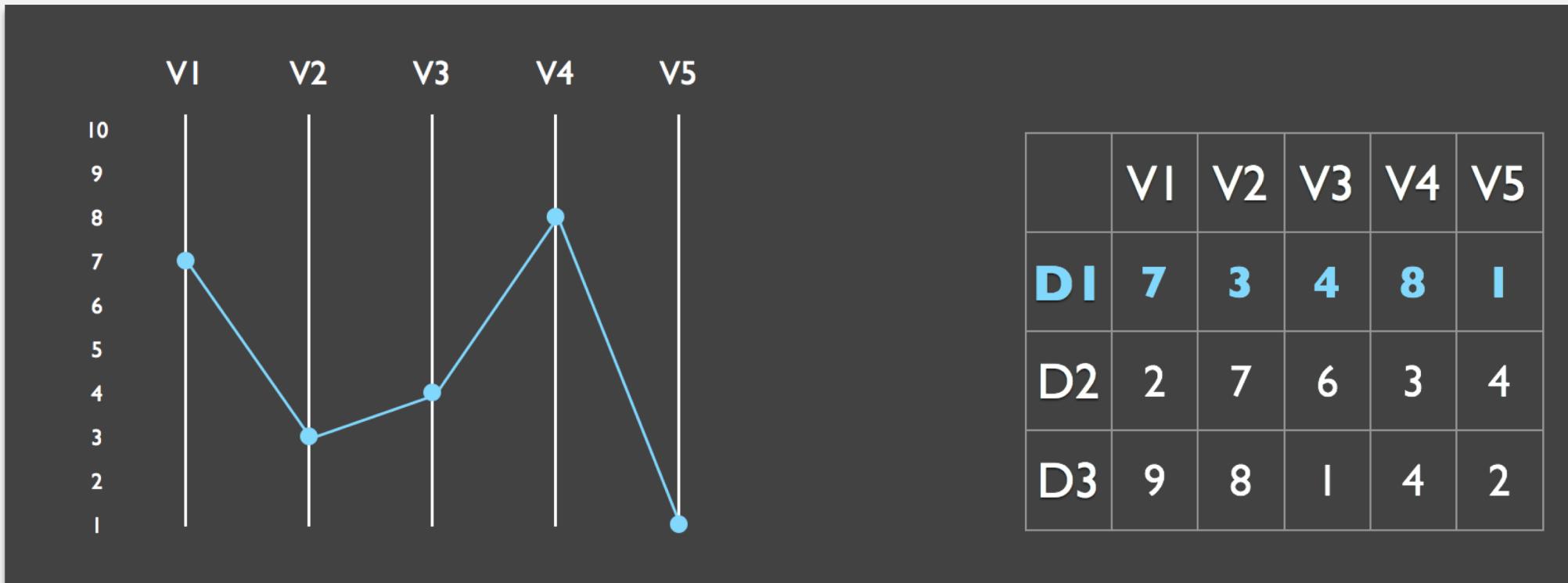
# Falcon: Network Flow Information



# Multivariate Data Visualization

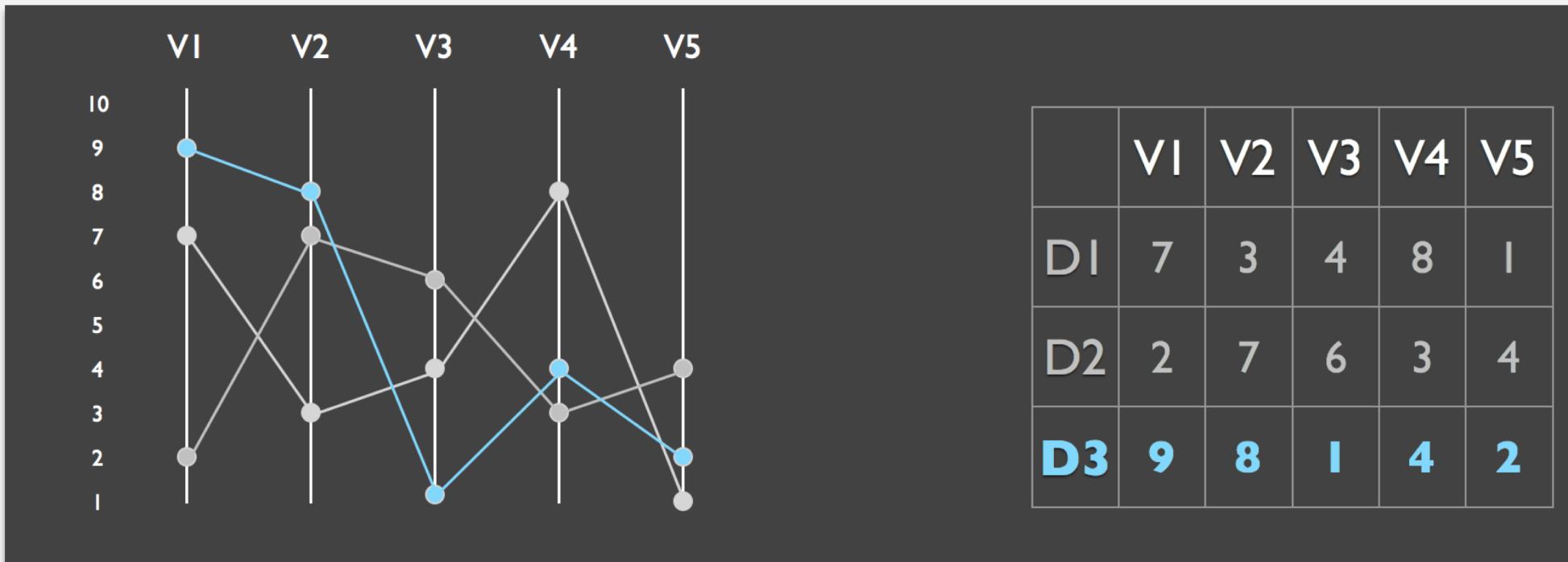
## Parallel Coordinates

Two-dimensional representation of multidimensional data sets by representing the  $N$ -dimensional data tuple  $C$  with coordinates  $(c_1, c_2, \dots, c_N)$  by points on  $N$  parallel axes which are joined with a polyline.



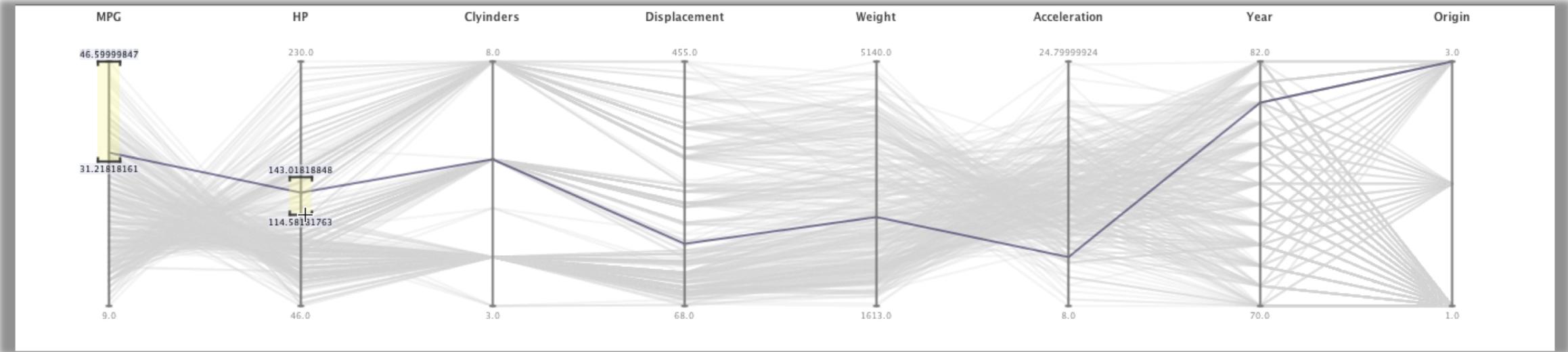
# Multivariate Data Visualization with Parallel Coordinates

- Number of variables only limited by display resolution
- Two-dimensional view of high-dimensional data
- Avoids information loss -- no data reduction



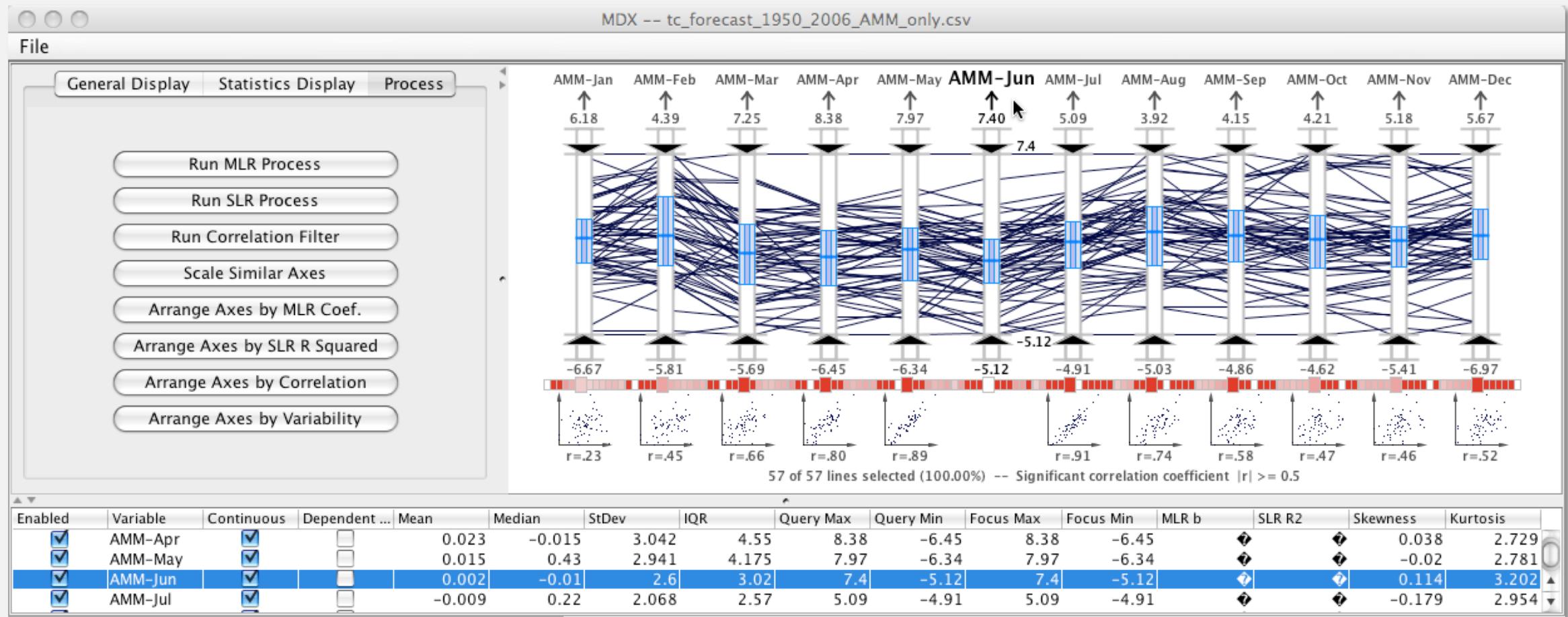
# The Power of ||-coords

- Visual queries for asking questions of the data
- Correlations (line slopes), clusters, and trends visually revealed



**Classical Parallel Coordinates Plot of the '83 ASA Cars Dataset**  
(Highlighted polyline shows Datsun 280ZX, 32.7 MPG & 132 HP)

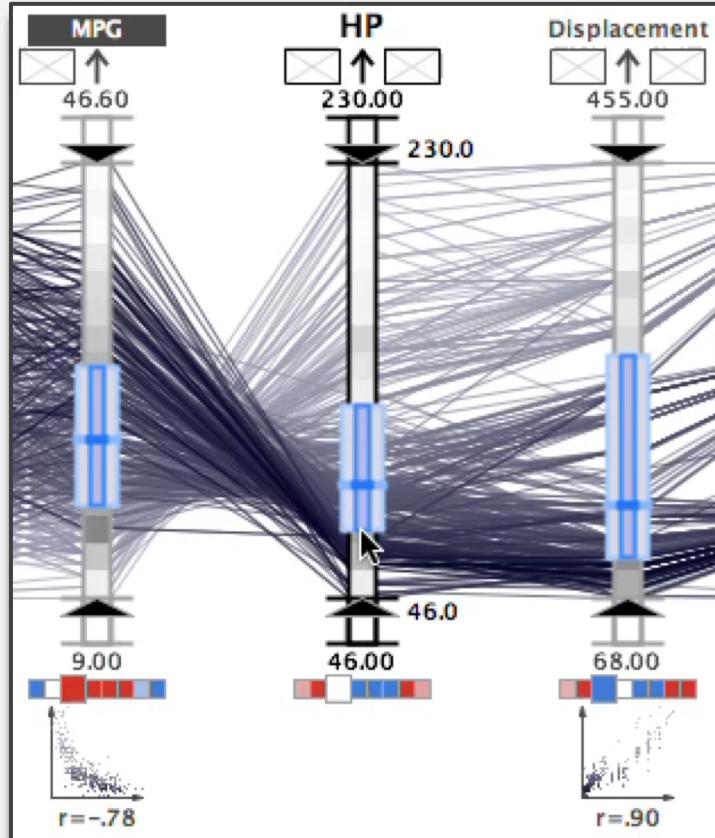
# Multidimensional Data eXplorer (MDX)



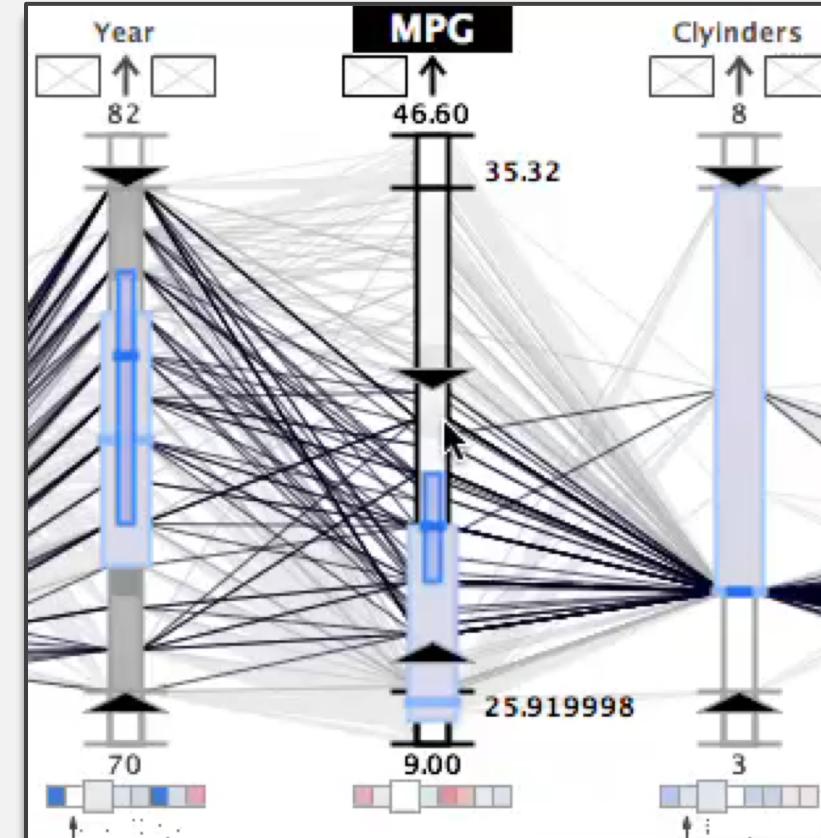
**Citation:** "Guided Analysis of Hurricane Trends using Statistical Processes Integrated with Interactive Parallel Coordinates". C. A. Steed, J. E. Swan II, T.J. Jankun-Kelly, and P. J. Fitzpatrick. Proceedings of the Symposium on Visual Analytics Science and Technology (VAST), pp. 19-26, Oct. 2009.

# New Interaction Techniques for ||-coords Dynamic Visual Queries

## Aerial Perspective



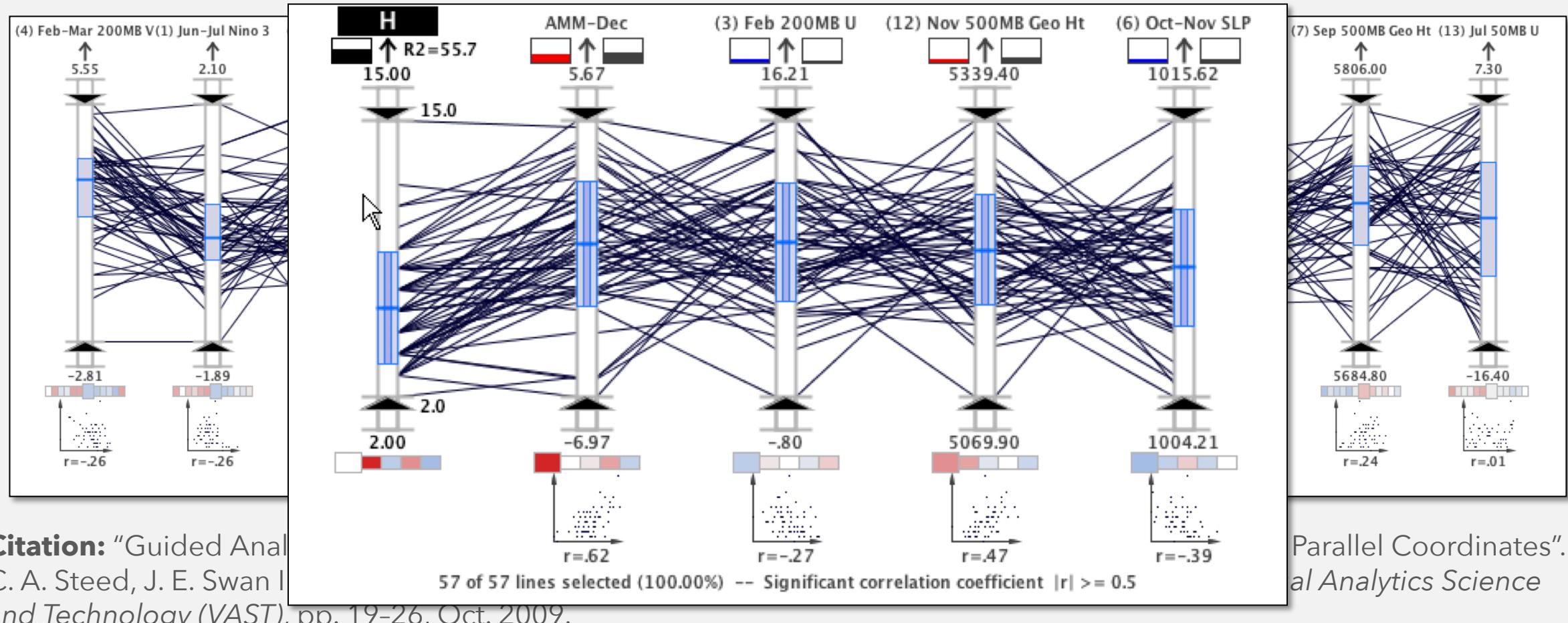
## Dynamic Axis Scaling



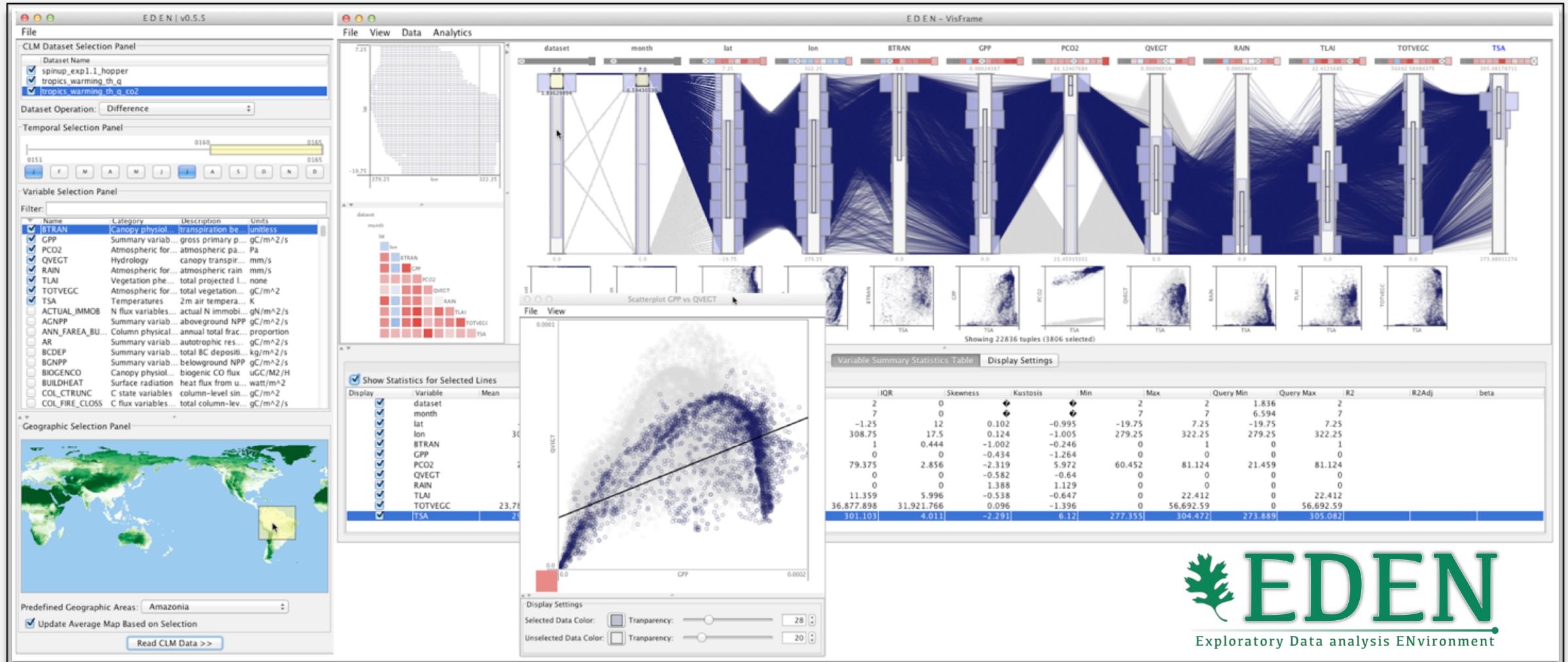
**Citation:** "Guided Analysis of Hurricane Trends using Statistical Processes Integrated with Interactive Parallel Coordinates". C. A. Steed, J. E. Swan II, T.J. Jankun-Kelly, and P. J. Fitzpatrick. Proceedings of the Symposium on Visual Analytics Science and Technology (VAST), pp. 19-26, Oct. 2009.

# MDX: Regression and Filters

$$(y - \bar{y})/\sigma_y = \sum_{i=1}^k b_i(x_i - \bar{x}_i)/\sigma_i$$

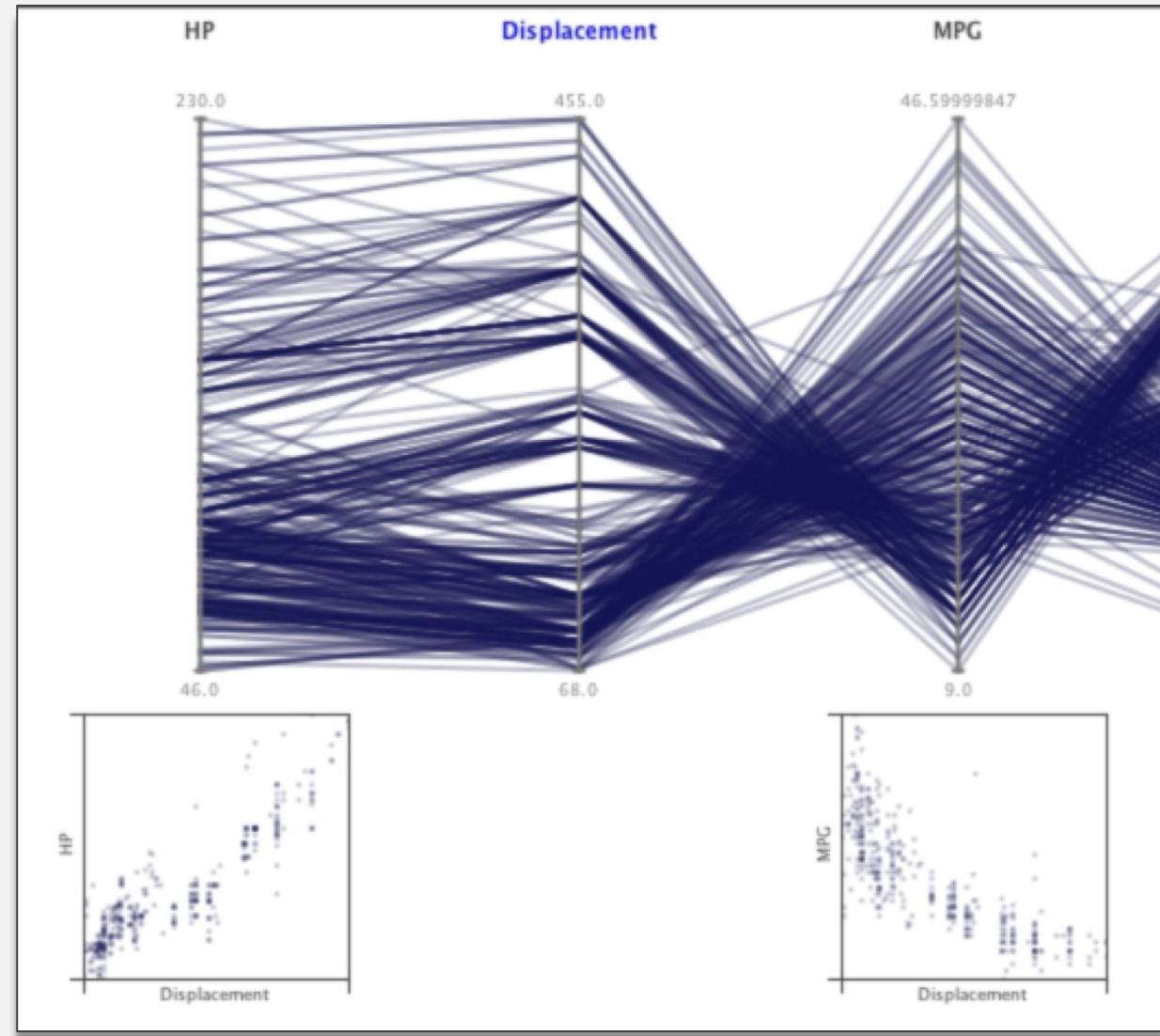
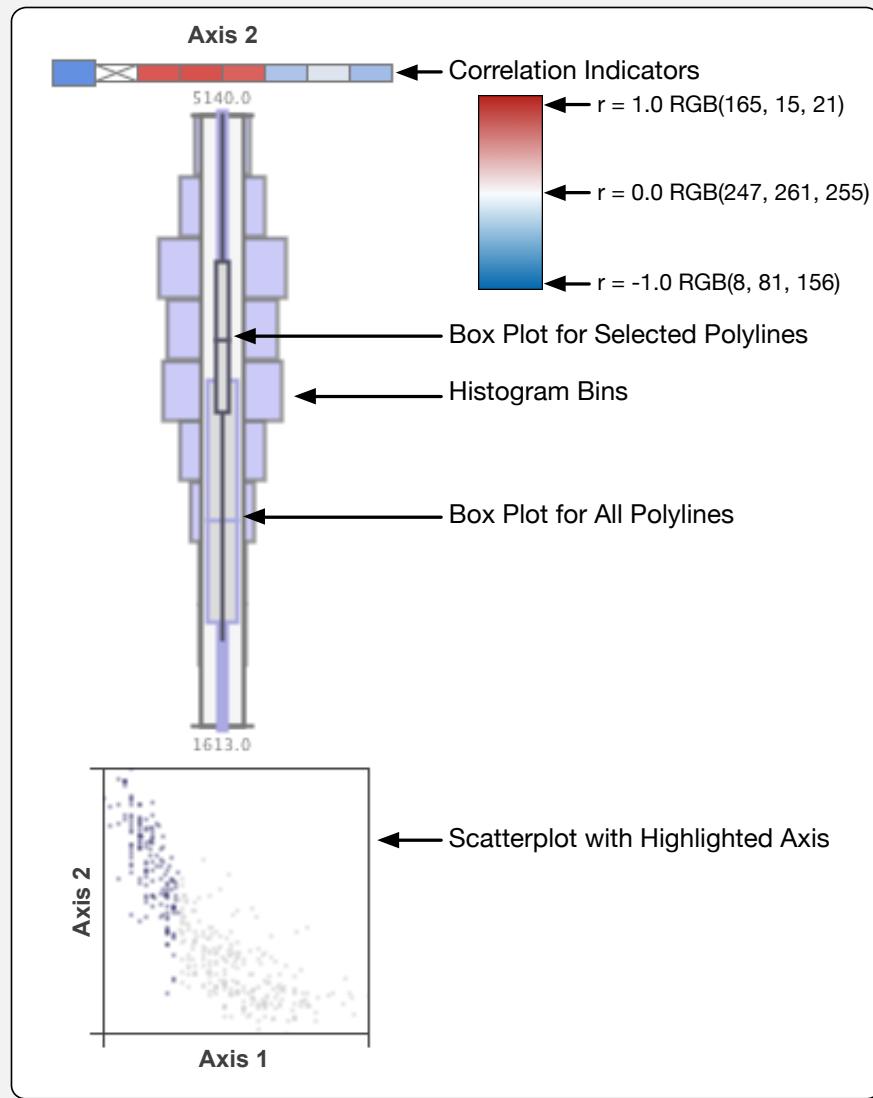


# Exploratory Data analysis ENvironment

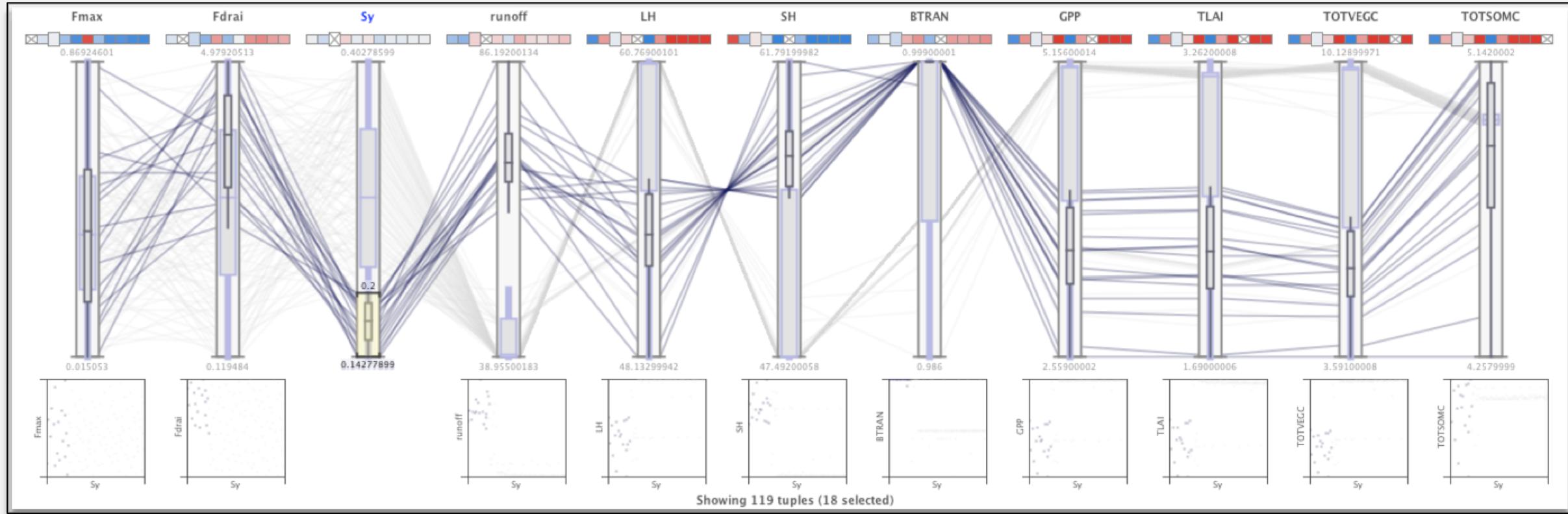


**Citation:** "Big Data Visual Analytics for Exploratory Earth System Simulation Analysis". C.A. Steed, D. M. Ricciuto, G. Shipman, B. Smith, P. E. Thornton, D. Wang, and D. N. Williams. *Computers & Geosciences*, 61:71-82, 2013.

# Guided Analysis Using Information “Scent”

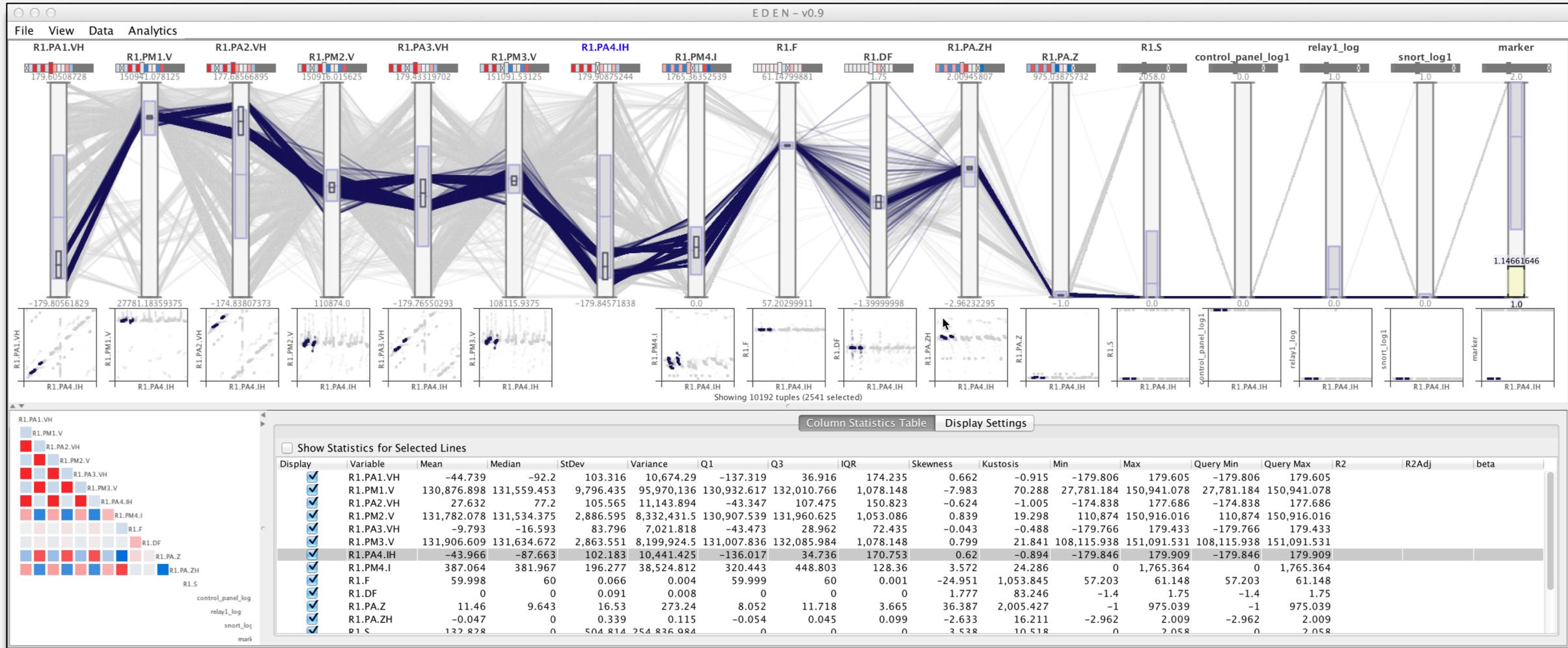


# EDEN: Climate Ensemble Analysis

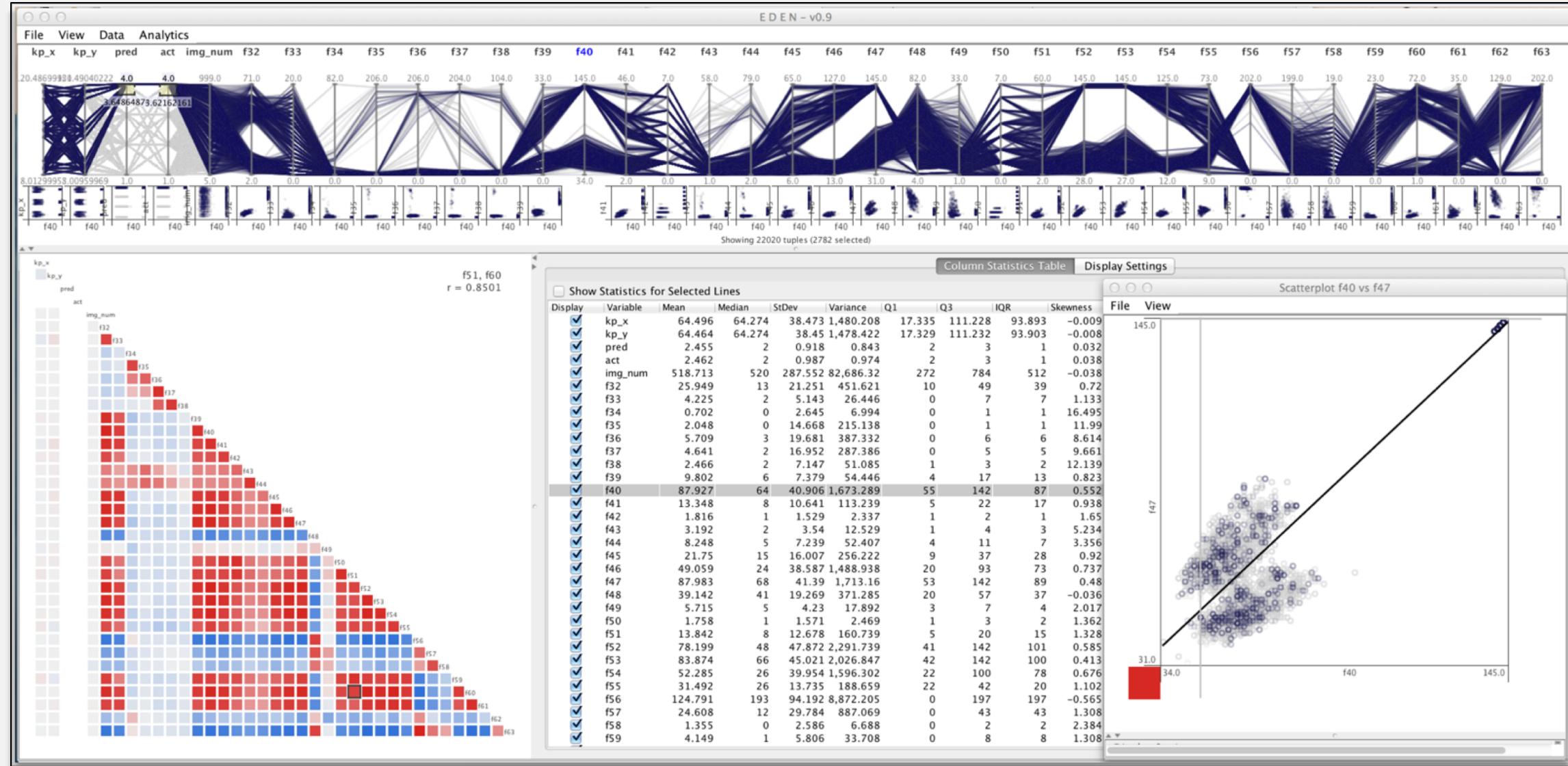


**Citation:** "Practical Application of Parallel Coordinates for Climate Model Analysis". Chad A. Steed, Galen Shipman, Peter Thornton, Daniel Ricciuto, David Erickson, and Marcia Branstetter. Proceedings of the *International Conference on Computational Science*, pp. 877-886, June 2012.

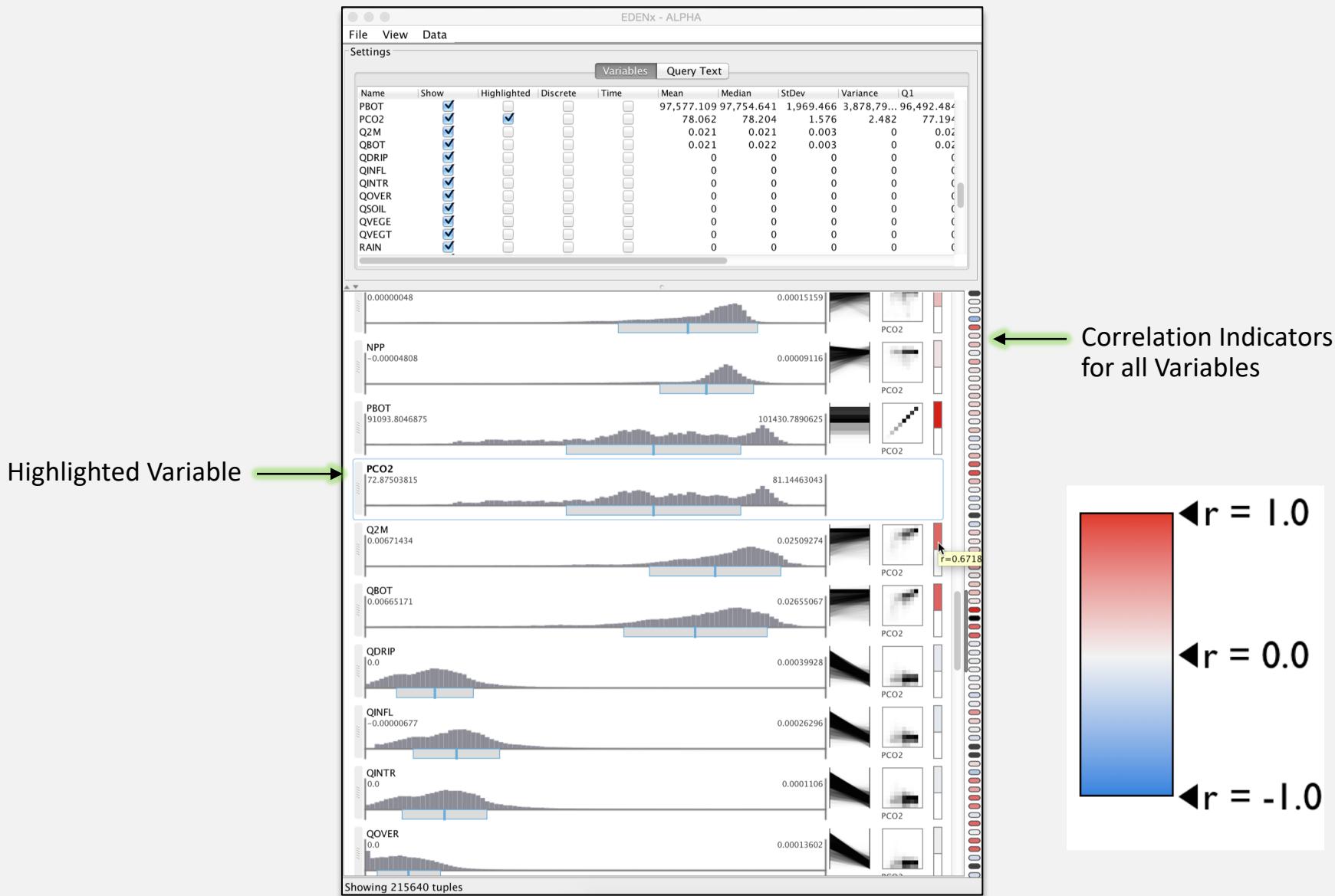
# EDEN: Cyber Physical Data Analysis



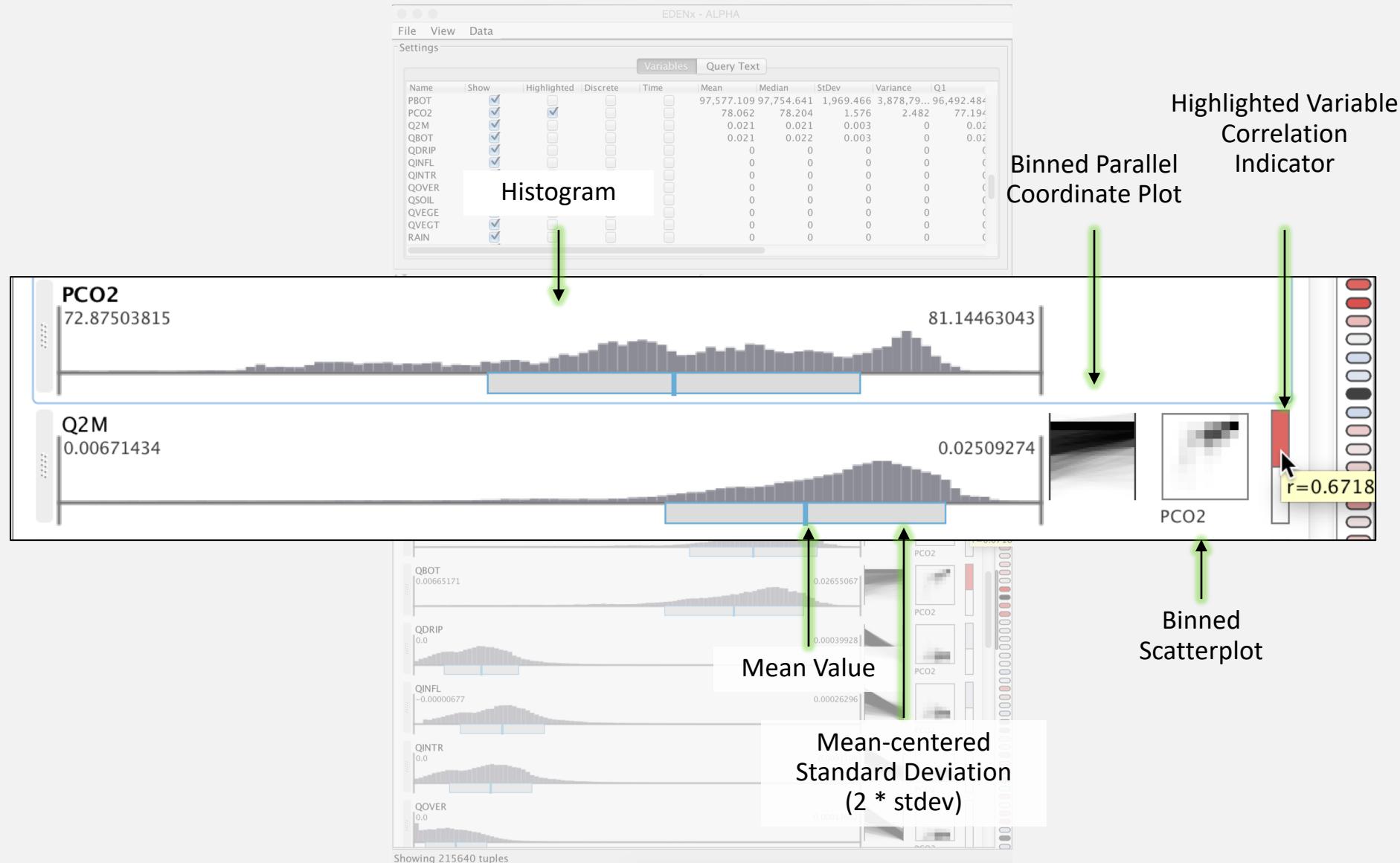
# EDEN: Machine Learning Feature Space



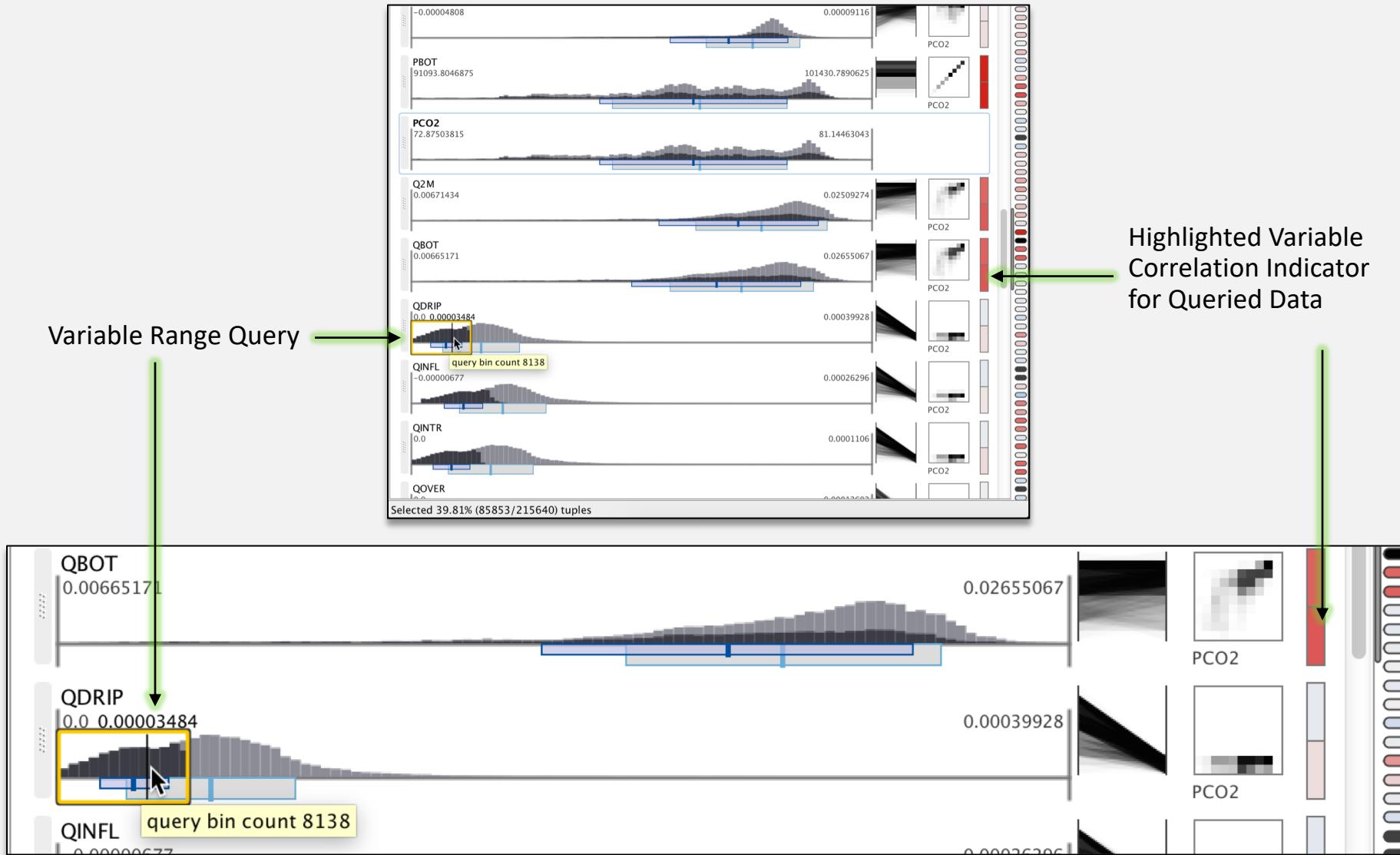
# EDENx: Handling Larger Data Sets



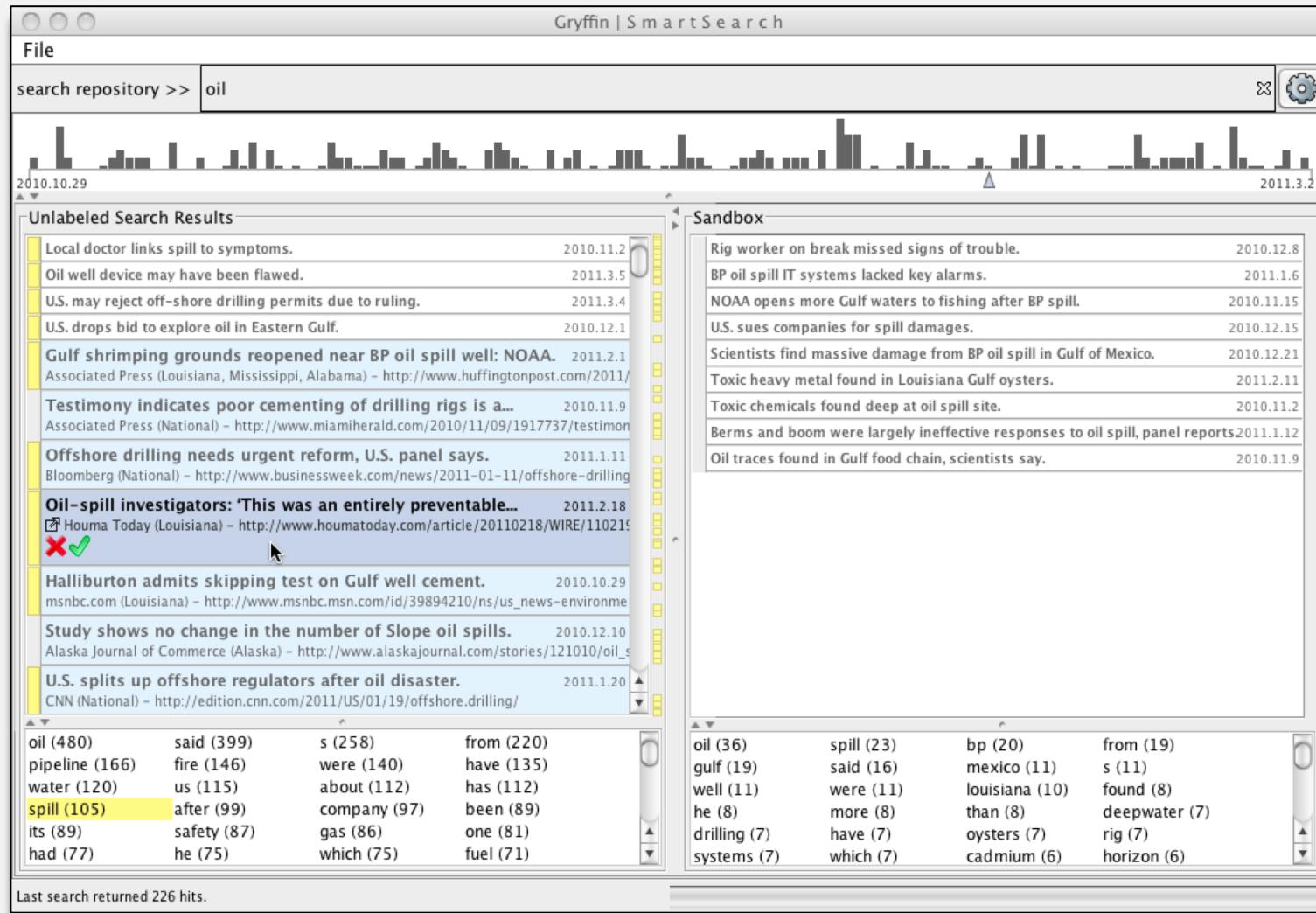
# EDENx: Handling Larger Data Sets



# EDENx: Handling Larger Data Sets



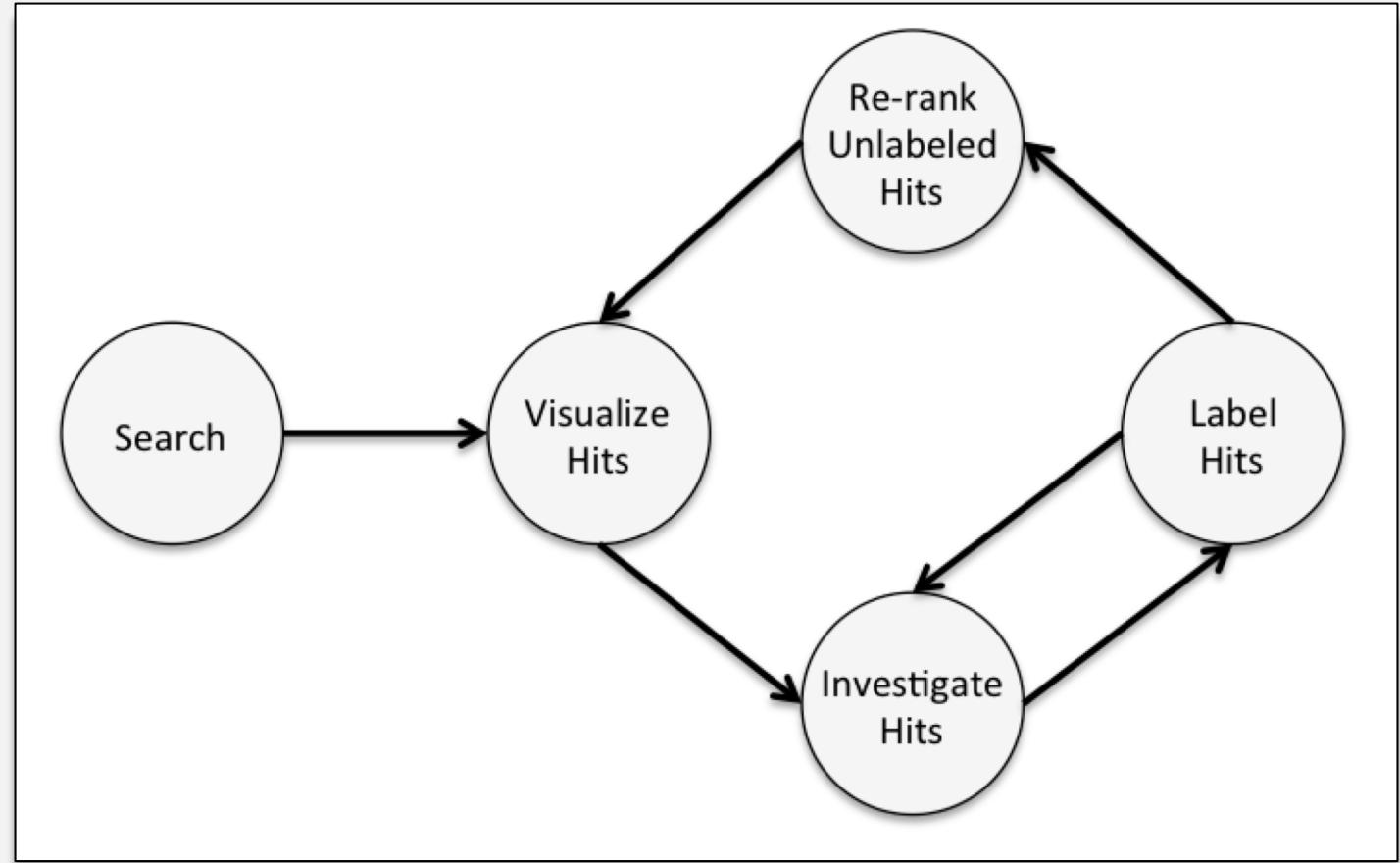
# Gryffin: Guided Text Search



**Citation:** "Guided Text Analysis Using Adaptive Visual Analytics". C. A. Steed, C. Symons, F. DeNap, and T. E. Potok. Proceedings of the Visualization and Data Analysis Conference, Jan. 2012.

# Gryffin: Intelligent User Interface

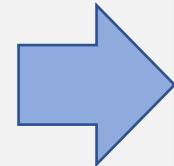
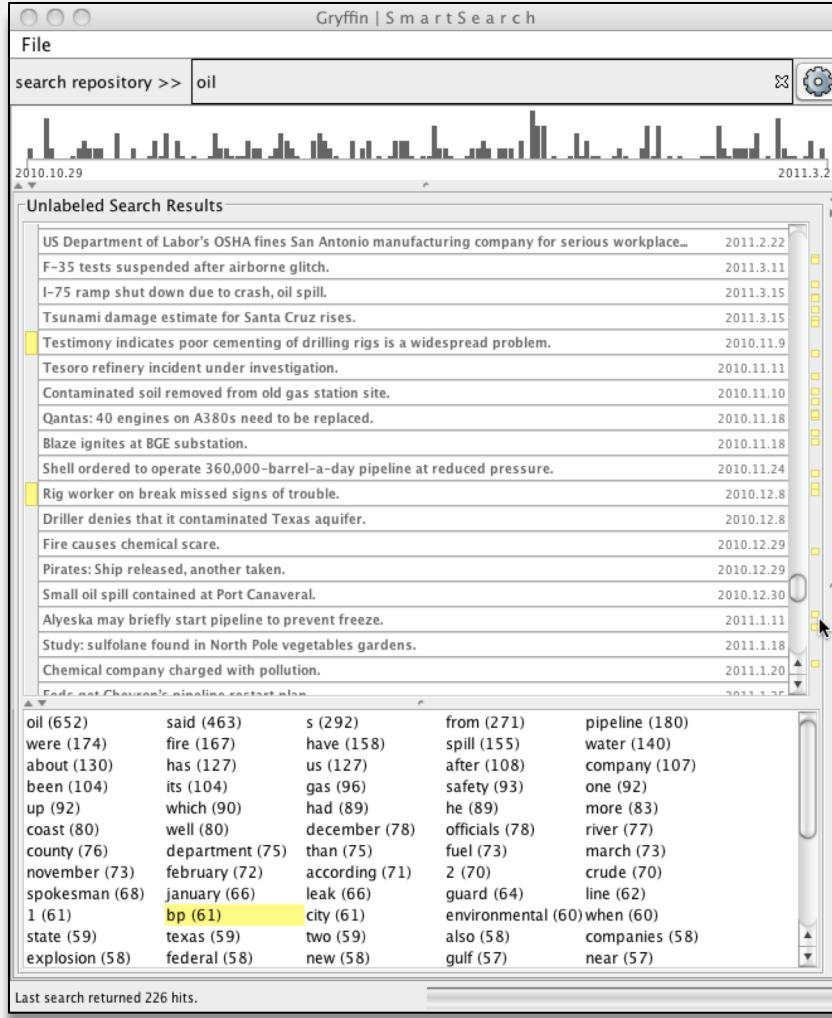
- **Workflow:** User focuses on the top of the list and labels documents as relevant / irrelevant
- A graph-based machine learning algorithm, developed by my colleague Chris Symons, uses labeled documents to re-rank list using document similarity



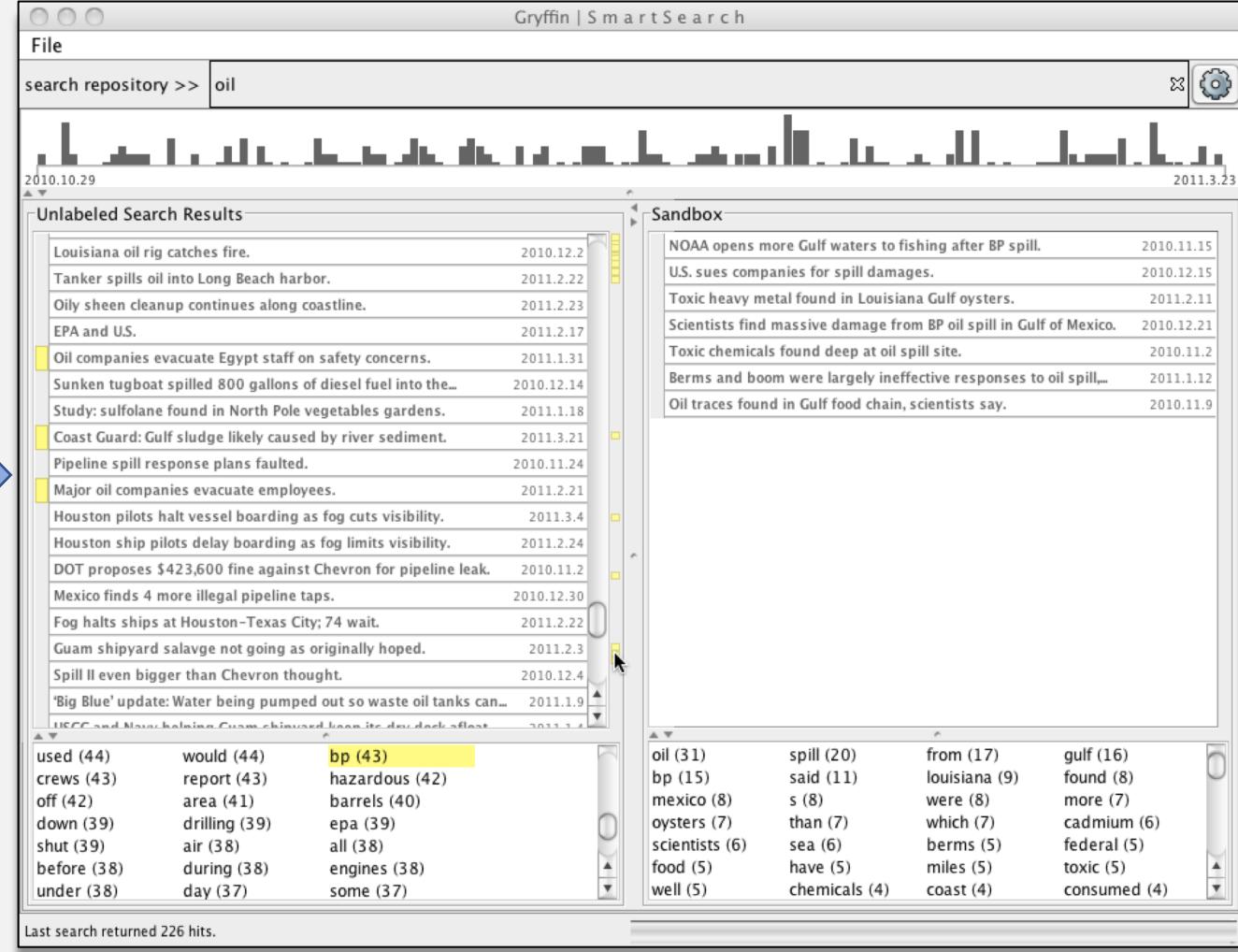
**Citation:** "Guided Text Analysis Using Adaptive Visual Analytics". C. A. Steed, C. Symons, F. DeNap, and T. E. Potok. Proceedings of the *Visualization and Data Analysis Conference*, Jan. 2012.

# Gryffin: Deepwater Horizon Case Study

## After Initial Search

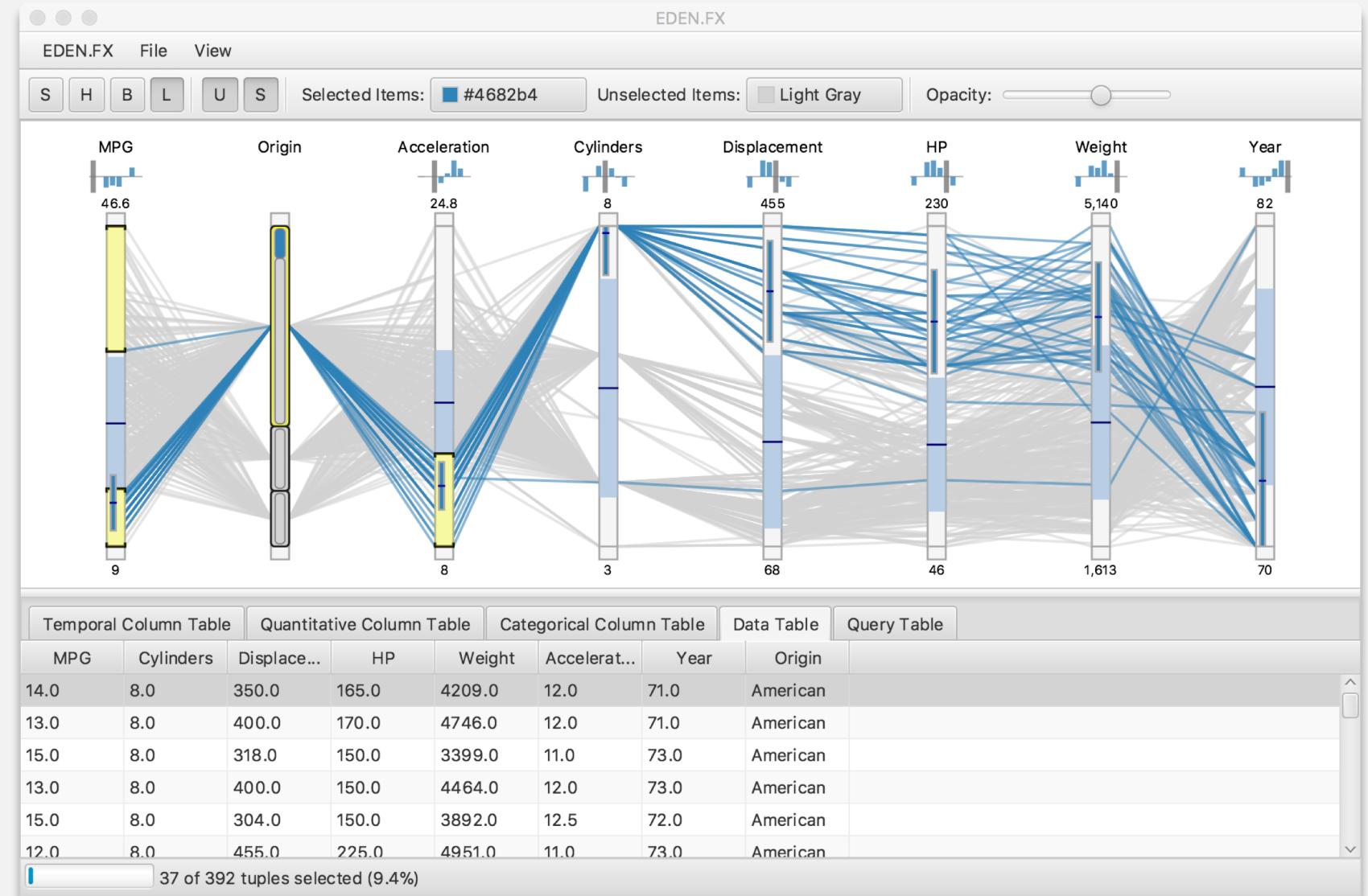


## Label – Learn – Re-Rank



# EDEN.fx: Scalable Multivariate Visual Analytics

- Combining EDEN, Falcon, Gryffin, and other systems
- Categorical, temporal axes
- Focus on scalability through level of detail rendering



# EDEN.fx: Climate Simulation Analysis

## Scientific Achievement

EDEN enables exploratory data analysis for new DOE E3SM climate simulation and observational data using techniques that combine interactive data visualization and statistical analytics.

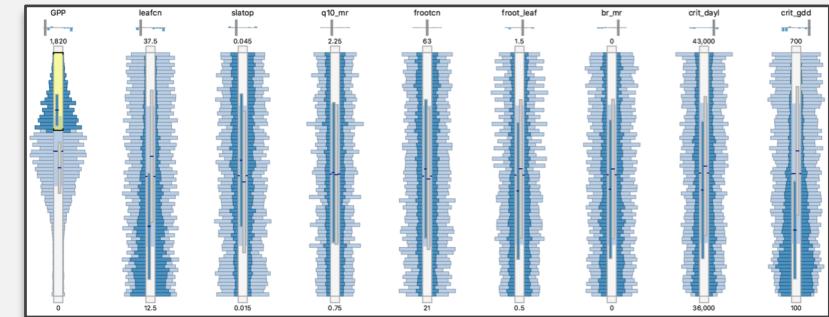
## Significance and Impact

We are collaborating closely with Dan Ricciuto (ORNL) on his SciDAC BER application, "An Integrated System for Optimization of Sensor Networks to Improve Climate Model Predictions".

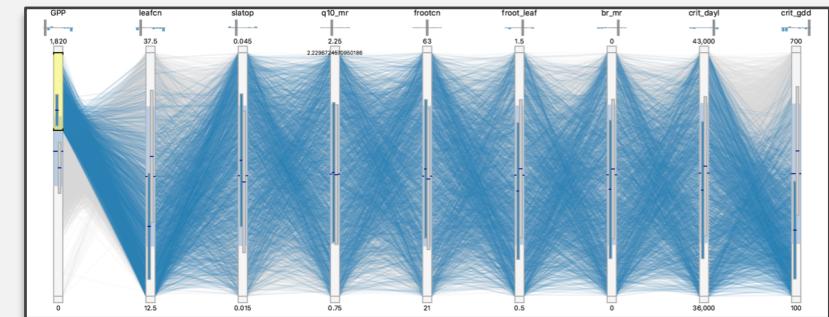
EDEN gives climate scientists the ability to consider more variables from large scale, land model parameter sensitivity analyses and ultimately improve DOE model accuracy.

## Research Details

- The screenshots at left show one insight found during parameter sensitivity analysis for realistic values of GPP, a model output that measures photosynthesis in plants.
- The plots helped scientists see that high values of GPP are associated with low leaf carbon to nitrogen ratio values (leafcn) and low critical growing degree days (crit\_gdd).
- Based on this insight, climate scientists will generate new ensembles covering smaller ranges of the leafcn and crit\_gdd parameter space for more accurate surrogate models.



**EDEN Parallel Histogram Summary View**

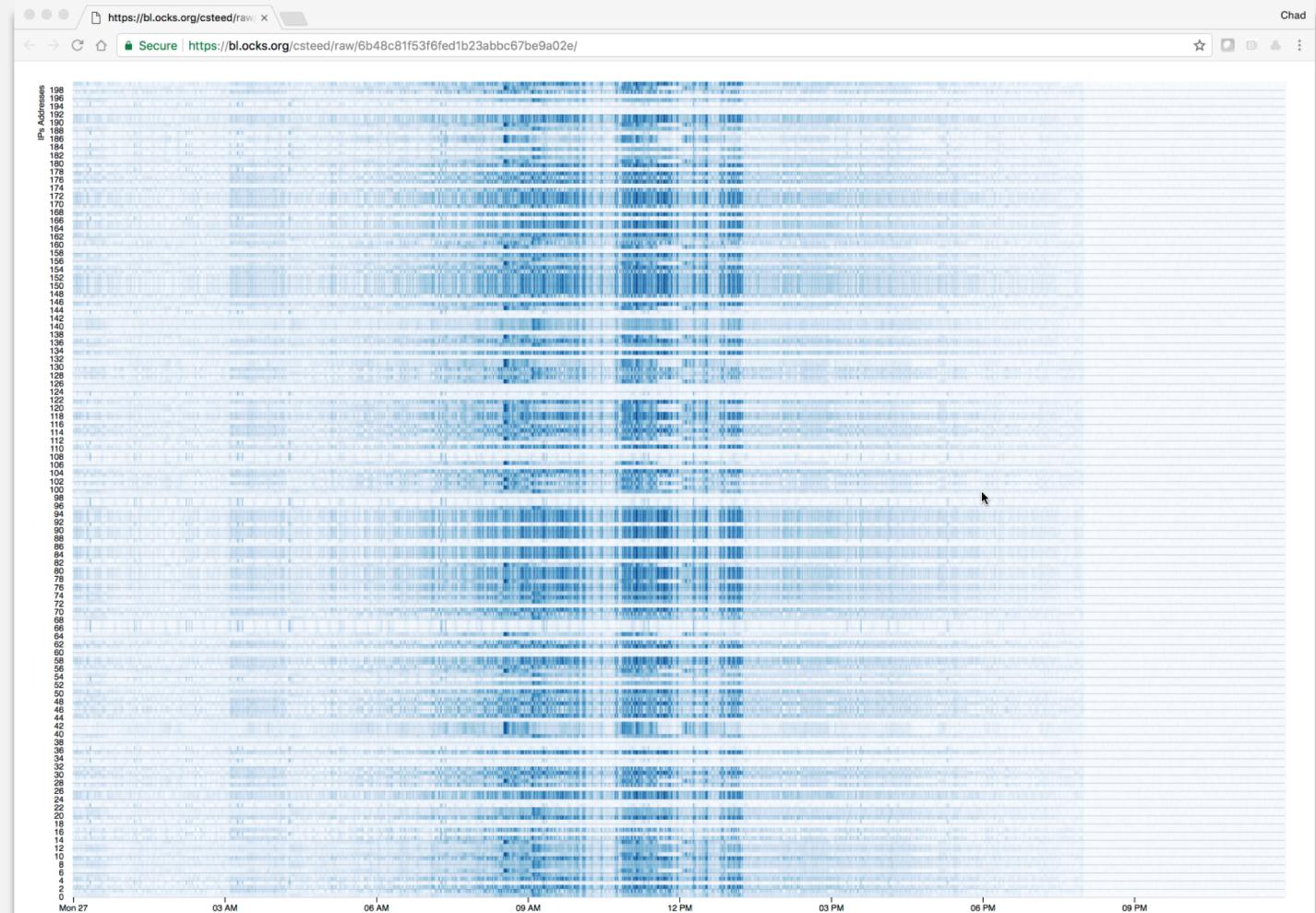


**EDEN Parallel Histogram Summary View**

The left axis, GPP, is a model output and the other 8 axes are parameters used to run model ensembles (each line is an ensemble run). In both views, the upper range of GPP values are selected revealing an association with low values of crit\_gdd and leafcn. New model runs will be executed using more constrained values for these parameters.

# Web-based Visualization: Overview IP/Stream Visualization

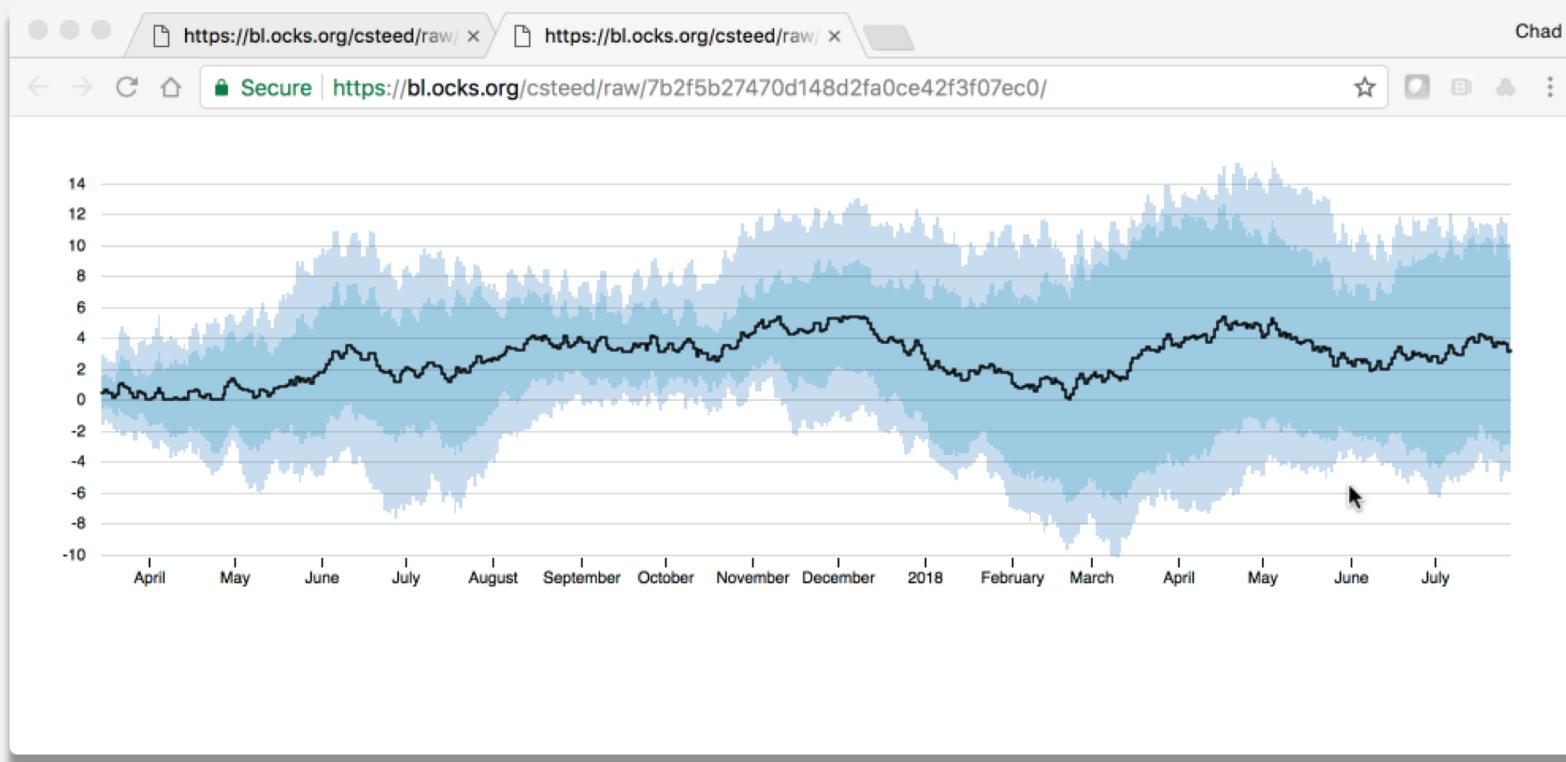
- **Goal:** Visualize aggregated time series data for hundreds of IP addresses
- Color scale encodes temporal bin statistics
- Support for zooming in to access more details (smaller time bins)



<https://bl.ocks.org/csteed>

# Web-based Visualizations: Temporal Summary Visualization

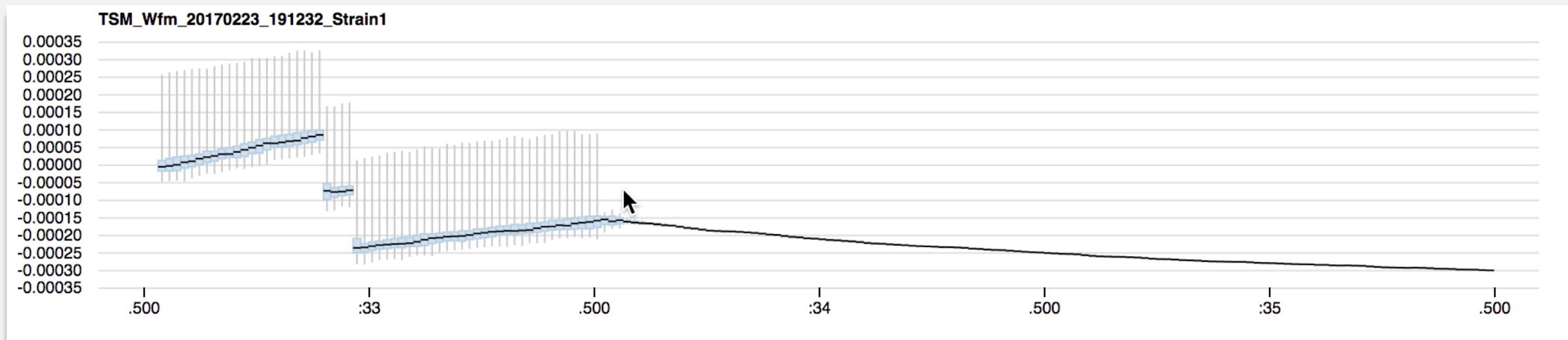
- **Goal:** Show time series typical values and variation in a single visualization
- Middle line encodes mean/median value
- Range polygons encode dispersion statistics
  - Standard deviation
  - Min/Max range
  - Interquartile range



<https://bl.ocks.org/csteed>

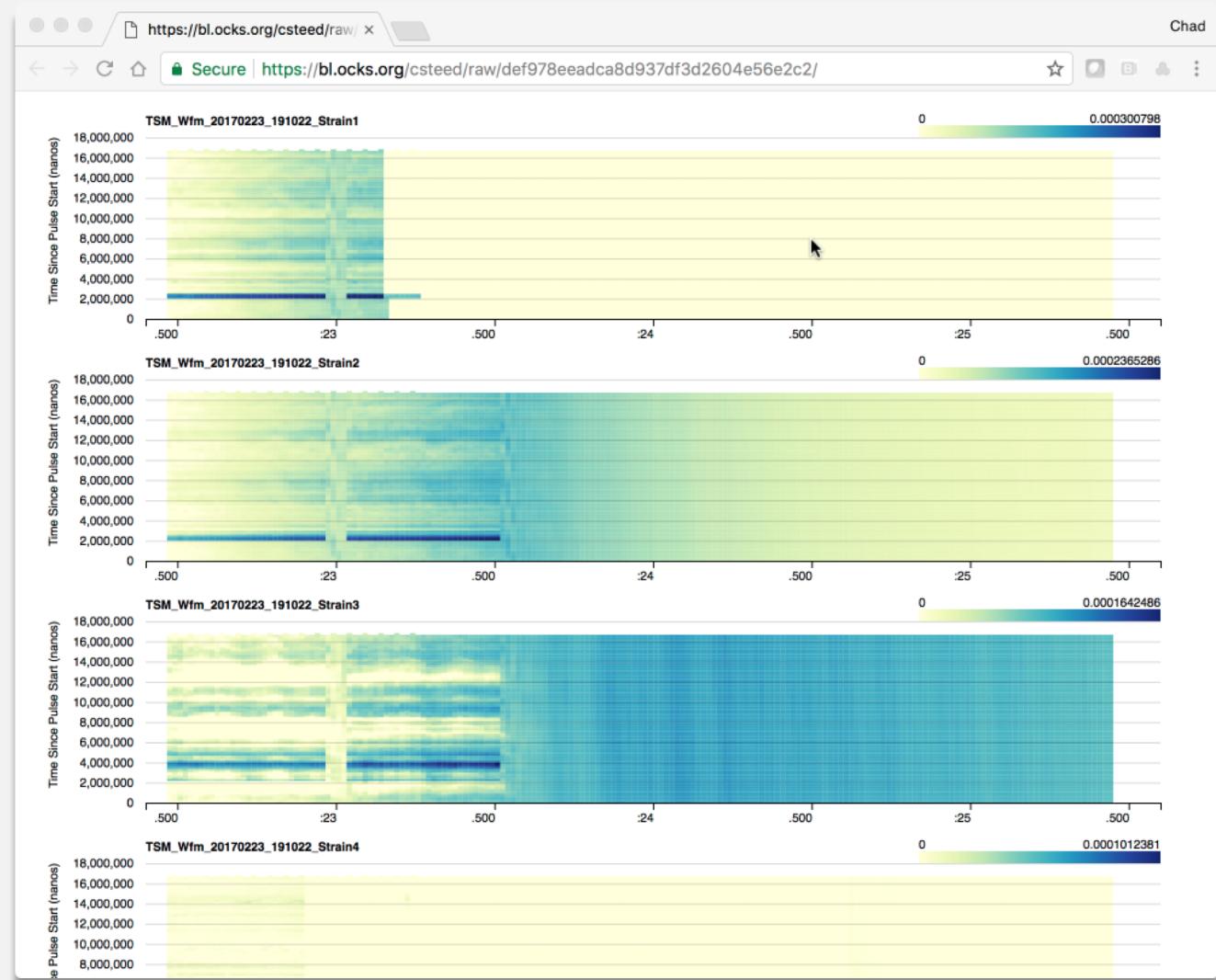
# Web-based Visualizations: Segmented Time Series Data

- **Goal:** Provide overview visualization of large scale sensor stream data as segments
- Data from strain sensors on ORNL SNS facility measuring neutron pulse activity
- Each vertical line represent range for pulse/segment
- Box represents standard deviation range or IQR for pulse/segment
- Black horizontal line represents mean or median for pulse/segment
- Interactively zoom in/out to see more/less details



# Web-based Visualizations: Temporal Summary Visualization

- **Goal:** Show detail time series visualizations for SNS pulse/segment data
- Mesoscale visualization between overview and raw data view
- Large scale data set at millisecond resolution binned and represented as color filled boxes



<https://blocks.org/csteed>

# Future Research Topics

- **Anticipatory Visual Analysis Models**
  - Human-in-the-loop
  - ML-based recommenders to guide analysis and reduce search space
- **Explainable Artificial Intelligence**
  - Data visualization to understand and fine-tune ML algorithms
- **New Data Visualization Techniques**
  - Focus on new ways to graphically represent abstract information (e.g., text, network flow, graphs, system performance)
- **Applications to Pressing National Challenges**
  - Cyber physical data analysis (e.g., streams, temporal + multivariate)
  - Analysis of network flow traffic (e.g., overview + detail)
  - Climate simulation and materials analysis