

# New Approach to Science Data Discovery in Petascale Systems

Tamara Sipes<sup>1,2</sup>  
<sup>1</sup>SciberQuest, Inc.  
<sup>2</sup>UCSD

**Collaborators:**

H. Karimabadi  
B. Loring  
A. Roberts  
Y. Wang  
B. Lavraud  
J. Merka  
J. Gosling  
V. Roytershteyn  
W. Daughton  
H. X. Vu  
Y. Omelchenko  
J. Raeder  
J. Dorelli  
A. Yilmaz

# Problem description

- ✘ Only about 10% of spacecraft data collected from all missions have been explored
  - Stale data
  - Low ROI
  - Biased and limited event examples make statistical studies very difficult
- ✘ New 3D kinetic codes can generate over **200 TB** of data from a **single run** with expected increase to **20 PB** by 2012!

*Sheer size and complexity mandate new approach to more fully exploit data archives.*

# Current Approach

- Dominant form of data analysis - *visual inspection of data*
- Standard data mining techniques (ANN, SVM etc.) are not directly applicable to analysis of **spacecraft data**:
  - multi-variate, multi-source, multi-scale time series data
  - multi-variate time series classification – very challenging
  - need for increased speed through parallelization or use of graphics processing units (GPUs)
  - lack of expertise to use/develop data mining techniques
- Not surprisingly, 90% of spacecraft data remains unexplored
- Our solution → Physics Mining

# Current Approach

- Scientific visualization is the dominant form of analysis of ***simulation data***, in space sciences and other fields of science
- while useful, scientific visualization lacks quantitative analysis capability
- existing visualization systems support knowledge discovery poorly, in fact, they often just *present graphics*
- typically the visualizations are not deeply integrated with methods and tools that would enable scientific insight
- Our solution → SciVis

# Definition: Physics Mining

Common Terms:

- Machine Learning
- Data Mining
- Artificial Intelligence
- Computer Vision

Physics Mining: Algorithmic approach to data analysis and derivation of analytical models from the data

# MineTool Algorithm Formula

(Karimabadi et al., 2007, 2009)

$$y_i = \mathbf{X}_i' \boldsymbol{\alpha} + \underbrace{\sum_{p=1}^P \zeta(\mathbf{X}_i) \boldsymbol{\delta}_p}_{\text{Linear Transforms}} + \underbrace{\sum_{q=1}^Q \Psi(\mathbf{X}_i, \boldsymbol{\gamma}_q) \boldsymbol{\beta}_q}_{\text{Non-linear Transformations}}$$

Reduces the problem to least squares fit rather than the difficult nonlinear optimization

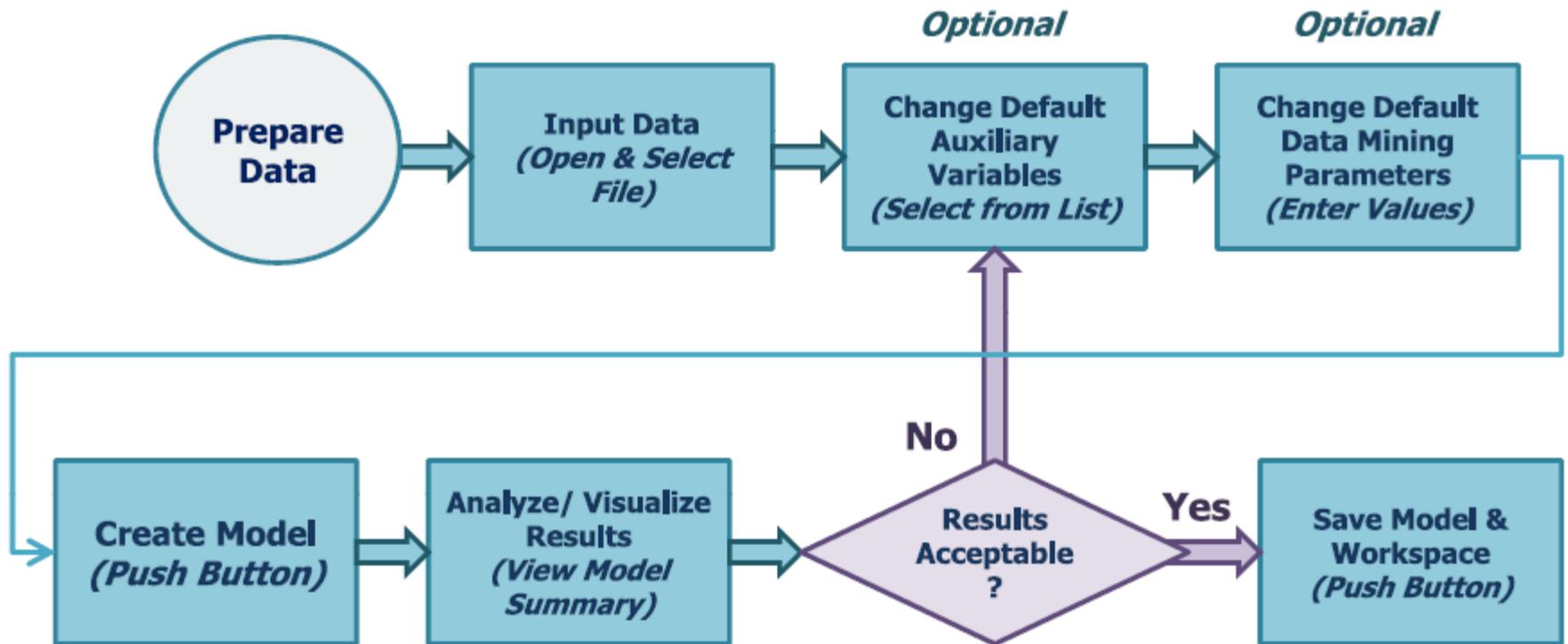


# Features of Physics Mining

Innovations in data mining include the following features:

- Applicable for both multi-variate static and time series/image data/simulation data
- Near real time processing of large data sets through use of GPUs/parallelization
- Capable of handling time series of unequal length and multi-scale events
- Unsupervised algorithms for both static and time series data
- Anomaly detection

# Process for applying MineTool



# MINETOOL USER INTERFACE DESIGN

# MineTool

SciberQuest

Open Dataset

Choose Auxiliary Variables

Set Data Mining Parameters

Visualize

Save

Select Input Dataset

Filename

Browse

Accepted formats: .txt and .csv

Dataset Variables

Number of instances: 10000

Number of variables: 4

List of Variables:

Index	✓	Variable Name	Input	Illustrative
1	<input type="checkbox"/>	ID	<input type="checkbox"/>	<input checked="" type="checkbox"/>
2	<input type="checkbox"/>	x	<input checked="" type="checkbox"/>	<input type="checkbox"/>
3	<input type="checkbox"/>	y	<input checked="" type="checkbox"/>	<input type="checkbox"/>
4	<input type="checkbox"/>	z	<input checked="" type="checkbox"/>	<input type="checkbox"/>
5	<input type="checkbox"/>	goal	<input type="checkbox"/>	<input type="checkbox"/>

Output Variable:

Select output variable

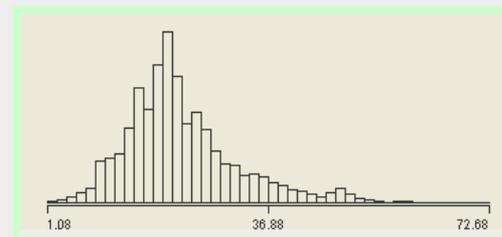


A variable is either an input, illustrative or output variable

Selected Variable Characteristics

Name: bmag  
Type: Numeric  
Missing: 0 (0%)  
Distinct: 2752  
Unique: 816 (9%)

Statistic	Value
Minimum	1.08
Maximum	72.68
Mean	22.791
StdDev	9.215



Variable Preprocessing

Preprocess Selected Variables

Normalize

Interpolate

Smooth

Remove Outliers

Nominal to Numeric

Undo

Save File

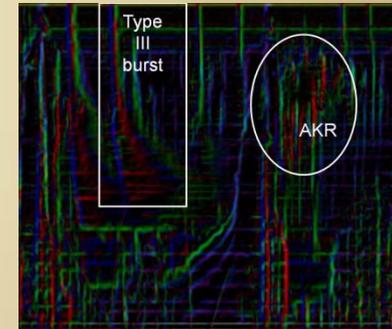
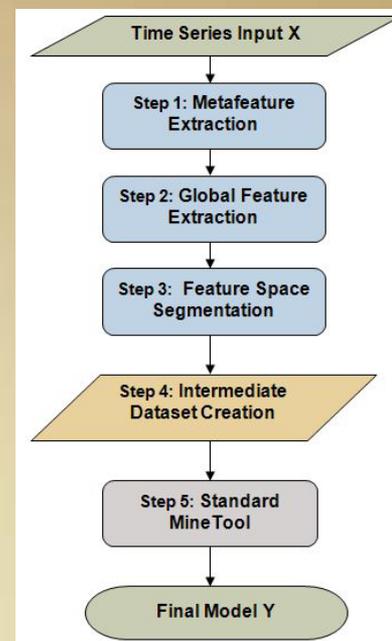
# Computer Vision

- Physics Mining can also combine advanced techniques in image analysis and computer vision
- Label images - drag and draw an ellipse bounding an event, select label
- Feature extraction from images and videos - detection of pixels and their surroundings, both spatial and spatiotemporal, that contains a variation in the observations
- Variations can be quantified by computing the changes in the locality of a pixel from its surrounding in the operator-selected event boundary

# Physics Mining Techniques...

Gee Whiz?

Or Practical Technology?



# Fun, But Impressive Application

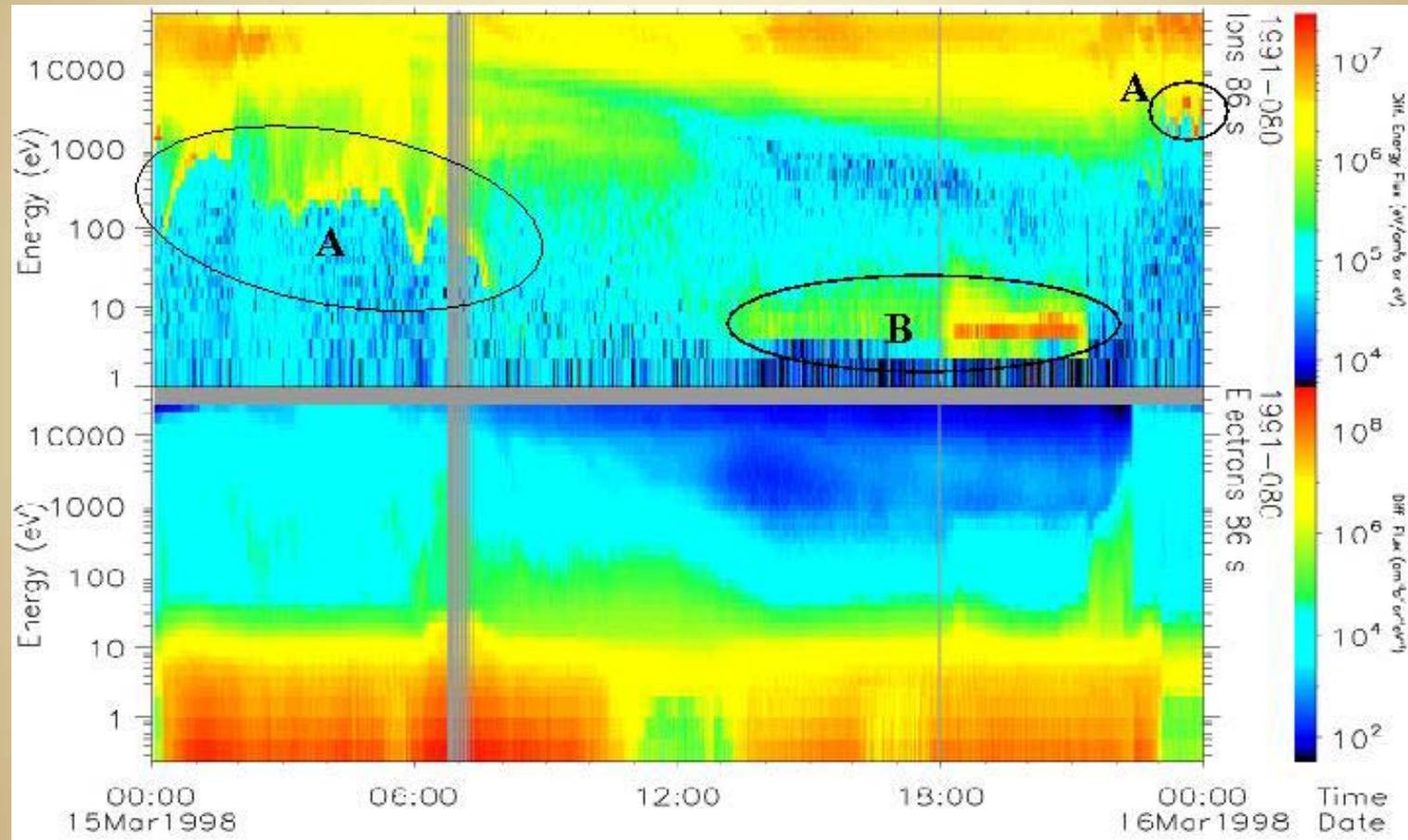
From this picture



Found “Dr. K” in this picture

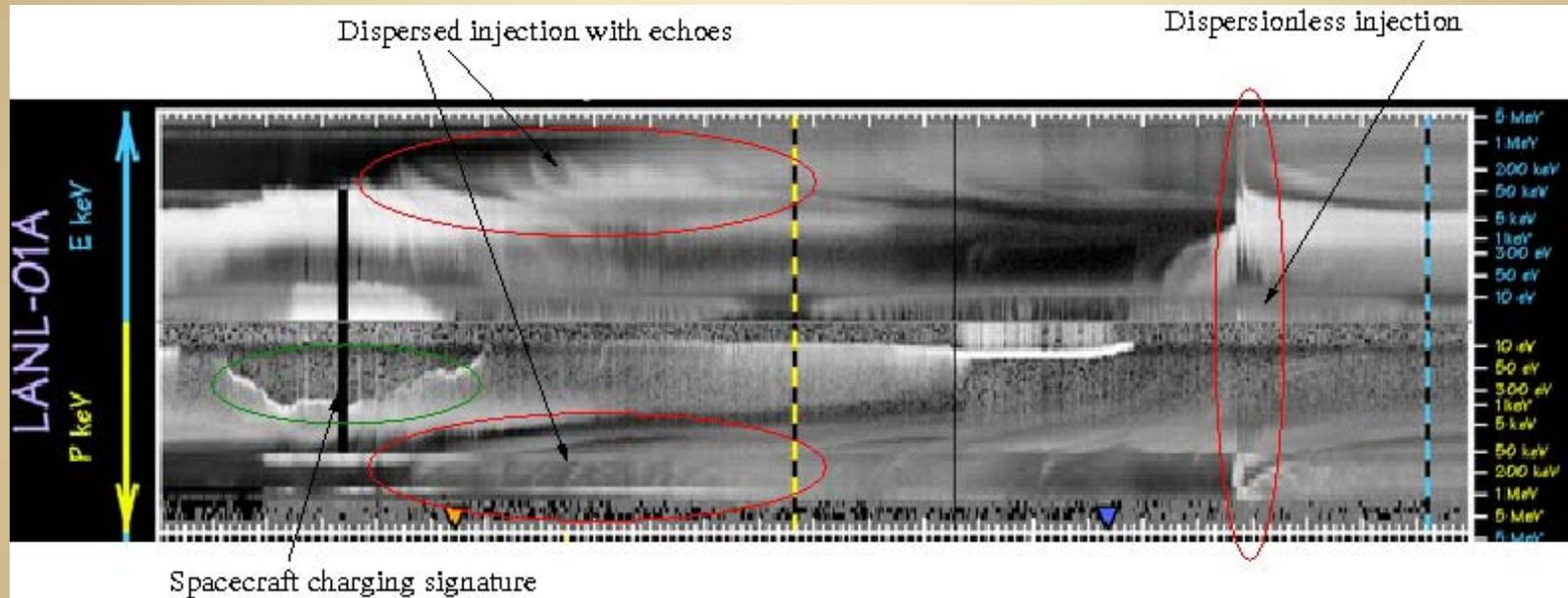


# Example: Physics Mining Applied to Spacecraft Charging



**Plasma data from MPA instrument (individual spectra stacked in time) on board spacecraft 1991-080 for 15 March, 1998 - Courtesy: R. Friedel**

# Substorm injection



Substorm Injection Signatures: electrons and ions plotted together with the zero energy in the center and energy increasing upwards for electrons, downwards for ions. Courtesy: R. Friedel

# Challenges in Global Simulations

- Multi-Scale Coupled System

- Spatial scales vary from centimeters to 200 RE
- Temporal scales vary from less than milliseconds to days
- Kinetic effects have global consequences

- Multi-physics

- electron physics: e.g., controls reconnection rate
- ion physics: e.g., dominates formation of boundaries and transport
- global features and dynamics: e.g., magnetotail/energetic particles
- coupling to the ionosphere



# Three Approaches to Global Simulations

- MHD (single fluid)
  - used extensively
  - does not resolve important ion physics
  - does not correctly capture the physics of reconnection
  - not suitable for studies of boundaries & discontinuities

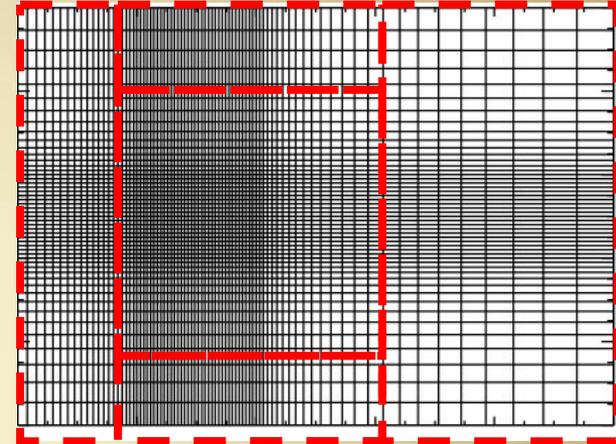
## Petascale computing has enabled:

- **Hybrid** (fluid electrons, kinetic ions)
  - fluid electrons, kinetic ions
  - it is the next stage of advance in 3D global simulations
  - resolves ion spatial scales (ion inertial length) and ion temporal scales (gyroperiod)
- **Full Particle** (kinetic electrons, kinetic ions)
  - it is the most complete description
  - only 2D is possible

# Our Approach

## Hybrid Code

- ▶ Nonuniform Mesh + Multi-time zones



- ▶ Discrete Event Technique (Omelchenko et al 2011)
  - robust (stable) and efficient (no idle computation)
  - works for arbitrary meshes
  - Predict local  $\Delta t$  for each variable  $f$  ("state") based on its estimated trajectory,  $f=f_E(t)$  and a given  $\Delta f_E$  (selected based on local CFL/reaction conditions)

## Fully Kinetic Code

- ▶ 2D global for capturing the microphysics of reconnection in a global setting

# The Hybrid Approximation

- Ions: kinetic particles
- Electrons: massless, quasi-neutral ( $en_e = q_i n_i$ ) fluid
- Electromagnetic fields:

- Faraday's law

$$(\partial/\partial t) \mathbf{B} = -c \nabla \times \mathbf{E}$$

- Ampere's law

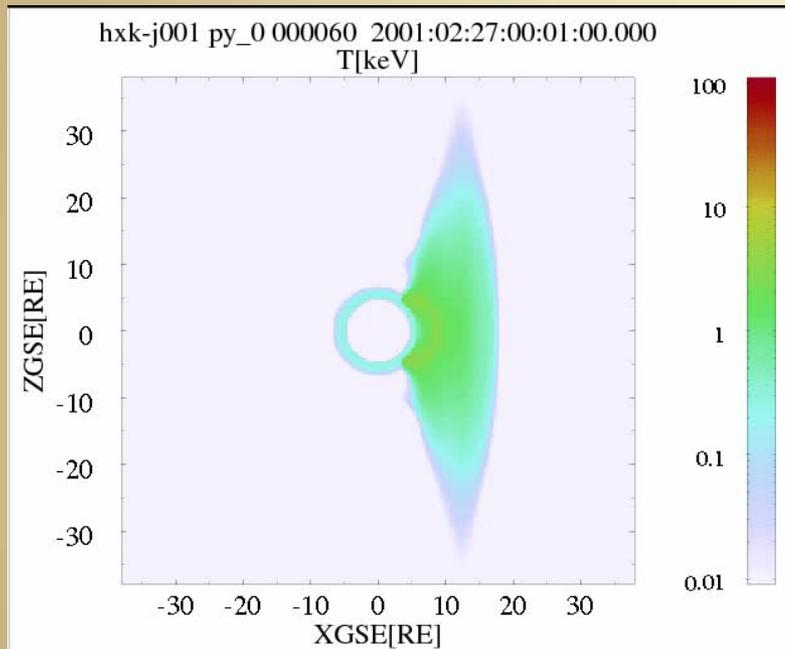
$$\nabla \times \mathbf{B} = 4 \pi \mathbf{J} / c = 4 \pi q_i n_i (\mathbf{v}_i - \mathbf{v}_e) / c$$

- Electric field from electron momentum equation

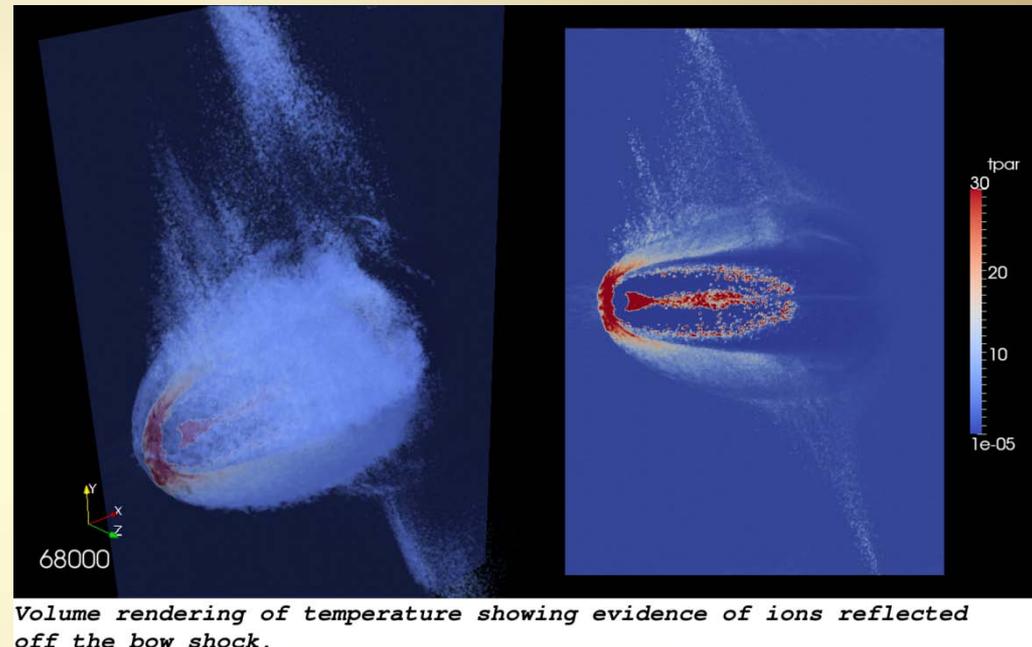
$$\mathbf{E} = -\mathbf{v}_i \times \mathbf{B} / c - \nabla p_e / (q_i n_i) - \mathbf{B} \times (\nabla \times \mathbf{B}) / (4 \pi q_i n_i) - \eta \mathbf{J}$$

# 3D MHD vs 3D Hybrid

Temperature

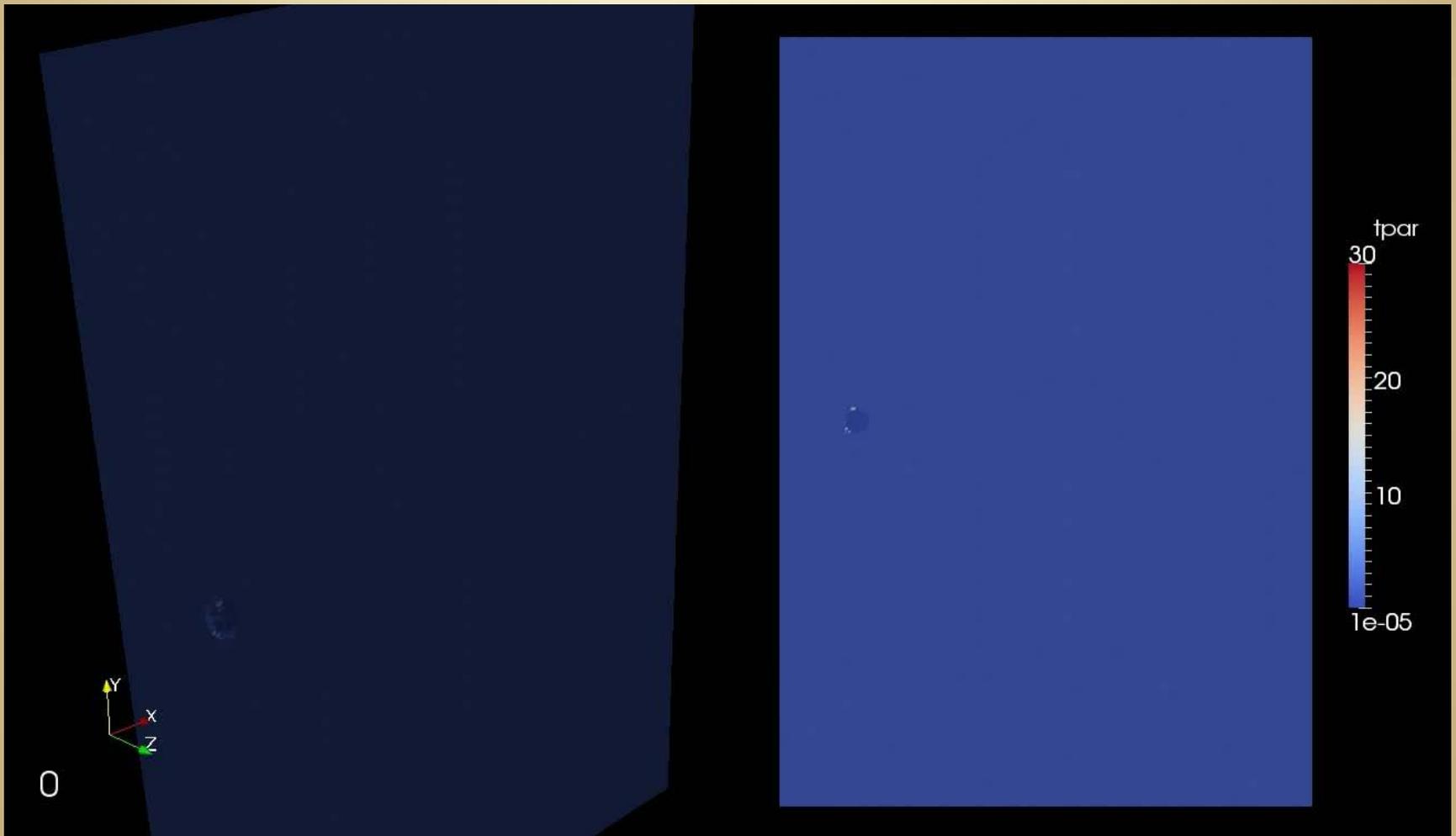


Temperature

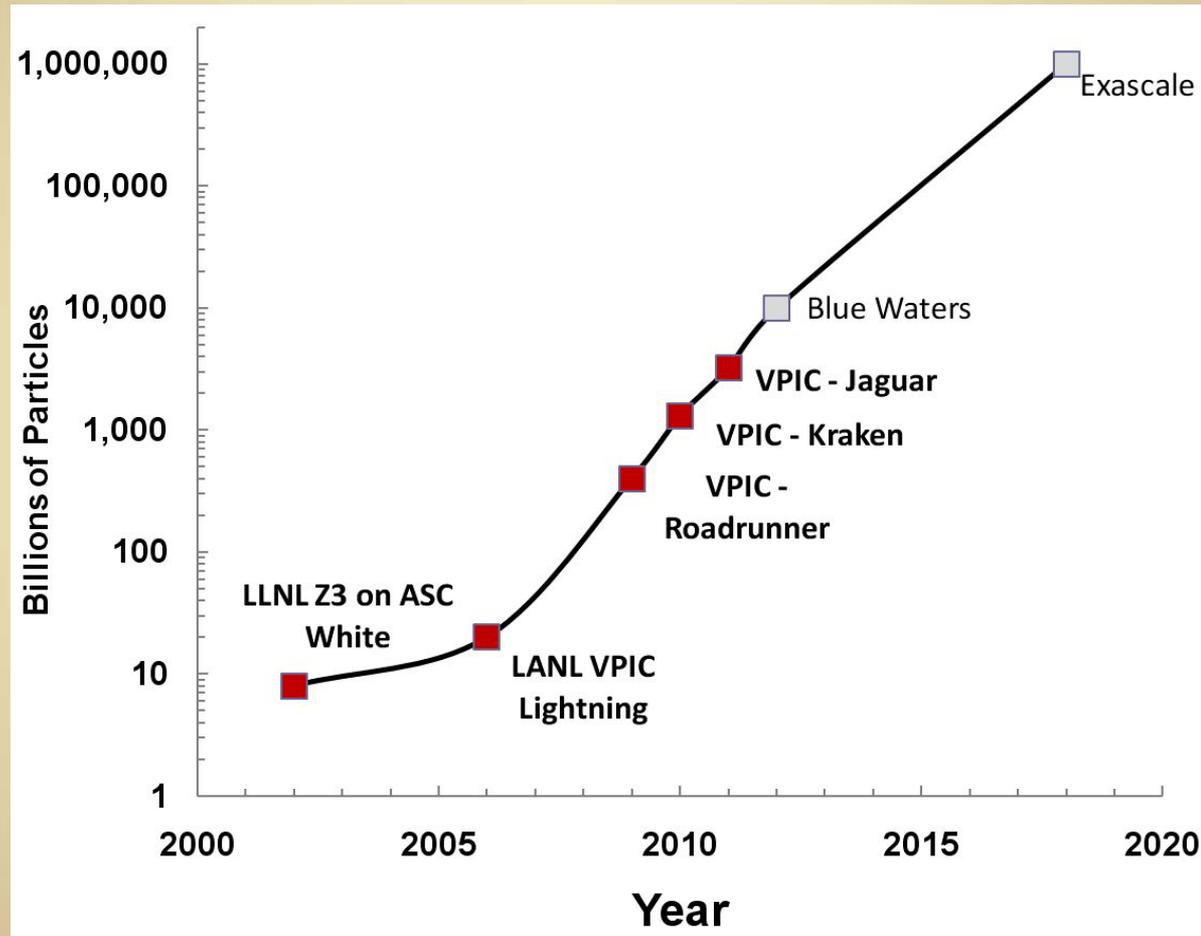


- The physics responsible for modeling backstreaming ions reflected at the bow shock are present in the hybrid simulations (right), but missing from MHD simulations (left).
- The enhanced physics of the hybrid simulation advances understanding of the real world processes.

# 3D Hybrid



# Number of Particles as a Measure of Size of Kinetic Simulations

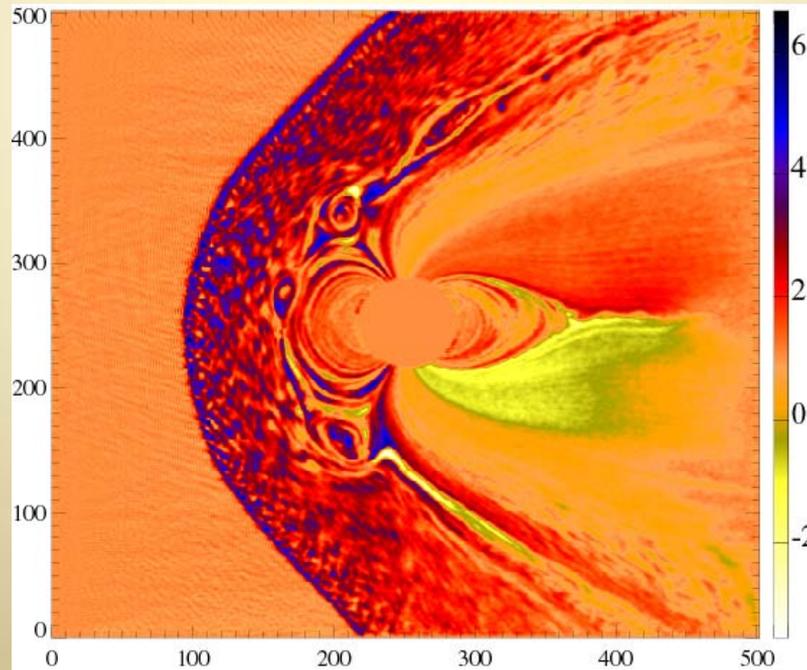
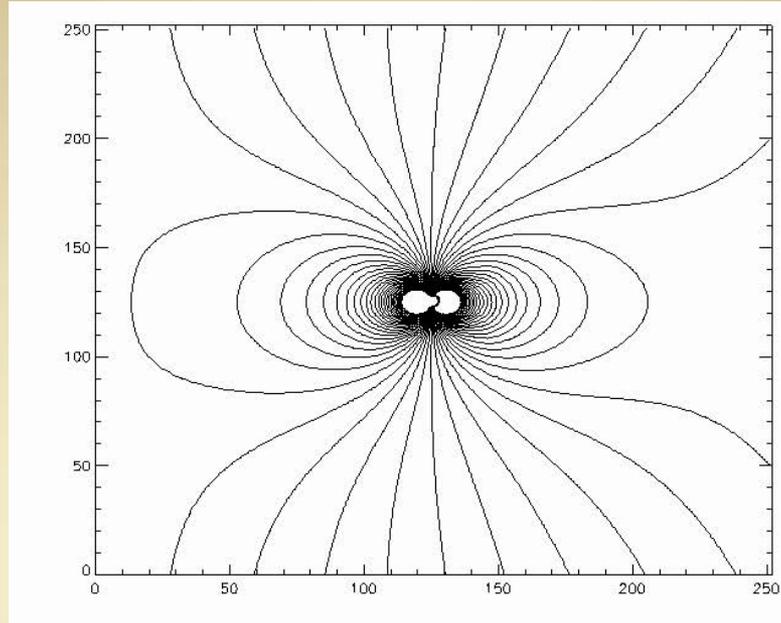


Size of largest PIC simulations in units of billion particles

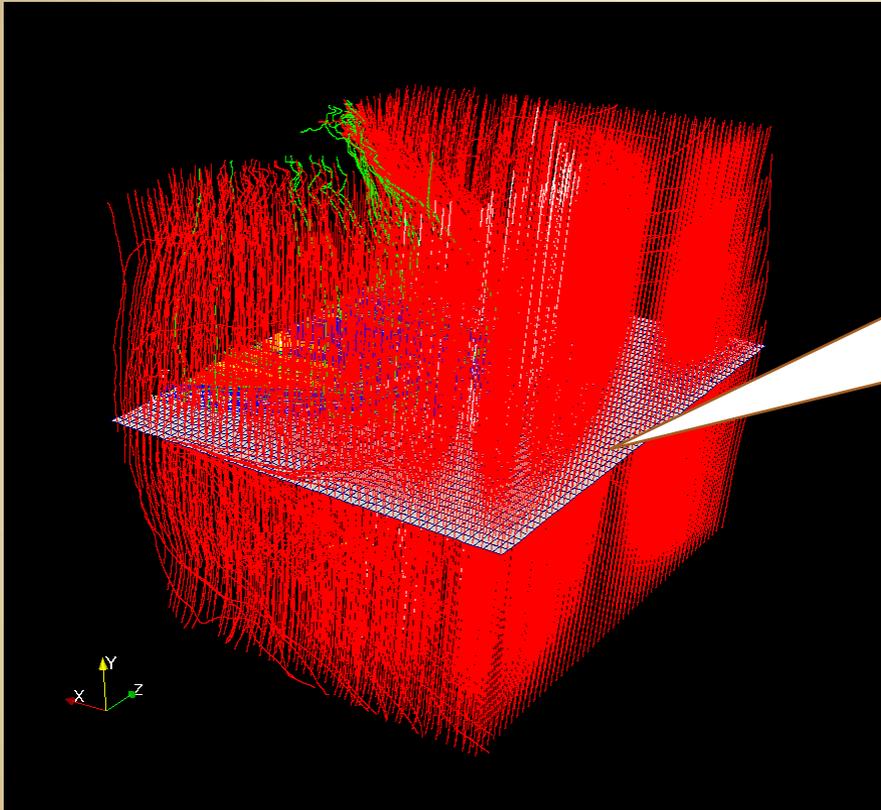
# Simulation Example: Feature Extraction and Tracking Flux Ropes in Global Hybrid Simulations

Problem: How to find flux ropes in a highly turbulent medium in a very large data set

**2-D  
Tracking  
is  
Straight-  
forward**



# Information Overload : Looking for FTEs



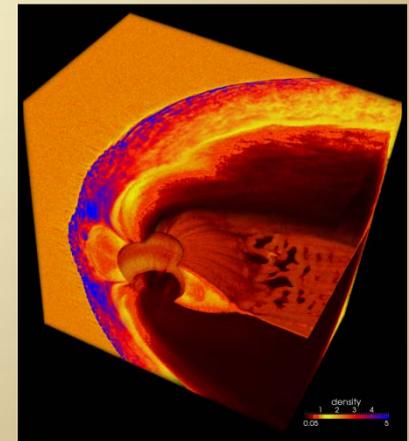
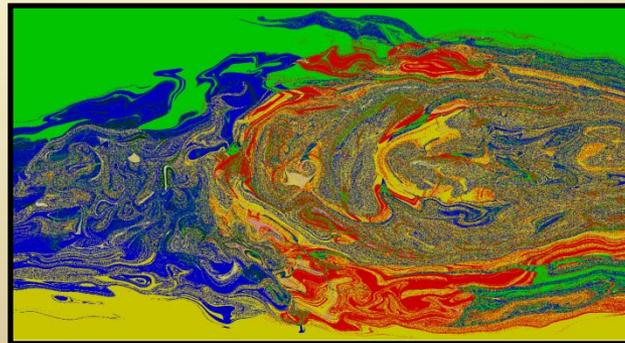
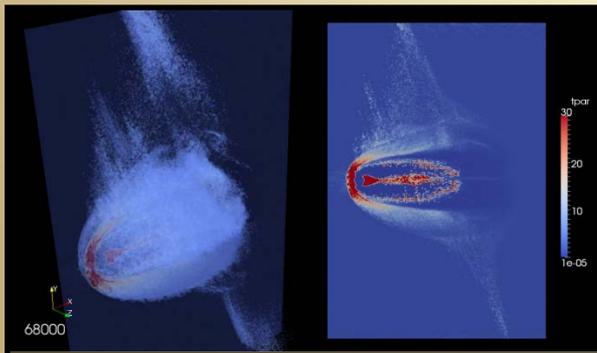
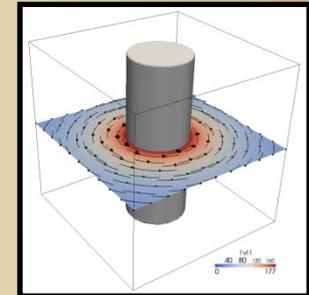
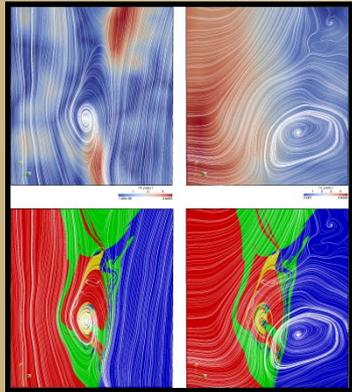
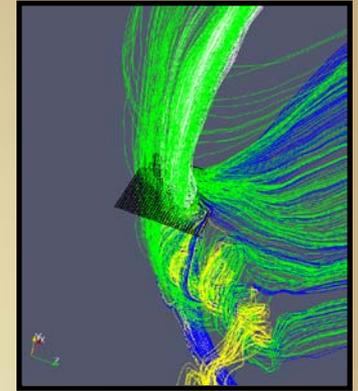
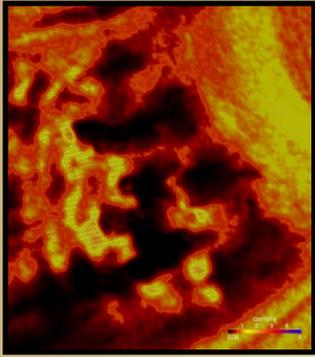
**3-D Is  
More  
Complex**

## Challenge:

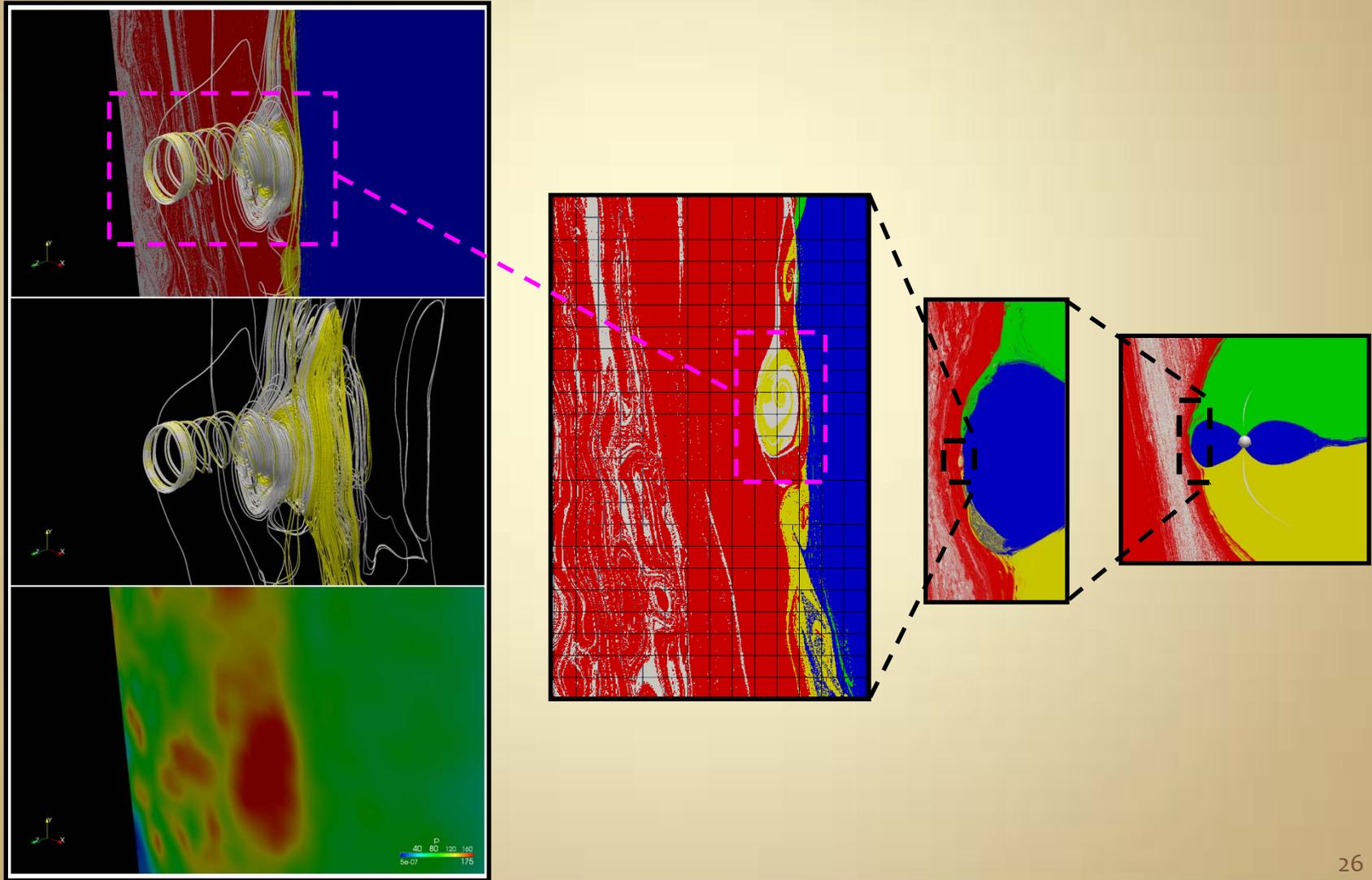
Very coarse resolution trace in a large dataset, ~100 cells per line. It is very likely that a feature will be overlooked or missed altogether.

# SciVis: Intelligent Scientific Visualization

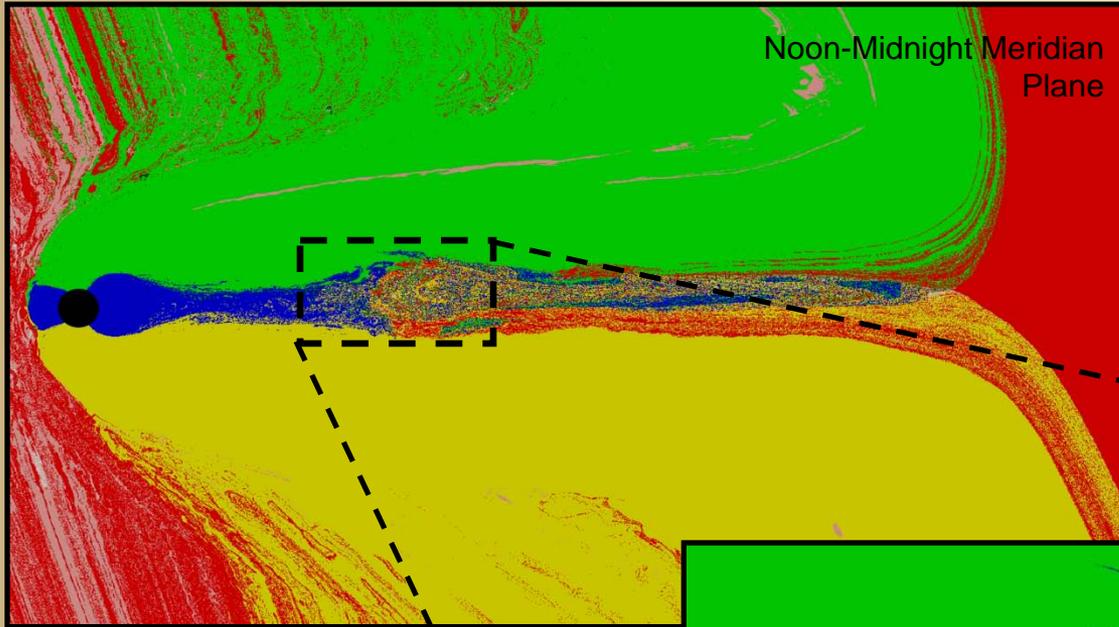
- Domain Specific Customizations for ParaView
- Incorporate Physics into Visualization Pipeline
- Autonomous Processing of Large Datasets
- Feature detection and Tracking
- Query Driven Visualization
- Incorporate Techniques from Data Mining and Computer Vision



# Detection of FTE: Needle in Haystack Problem



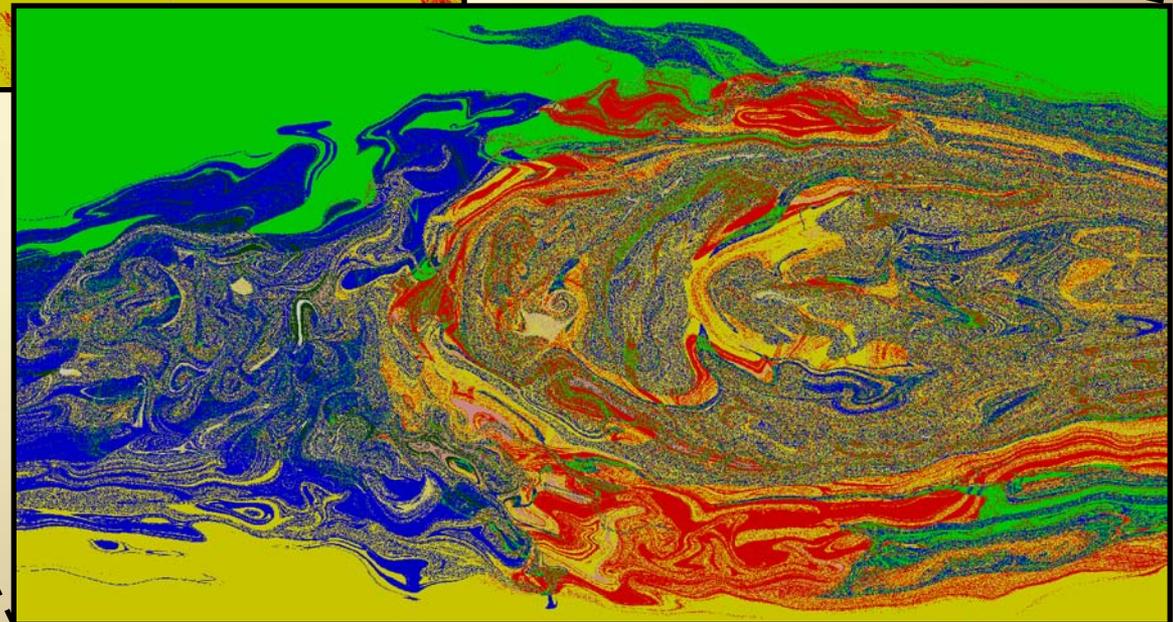
# Turbulent Magnetosphere



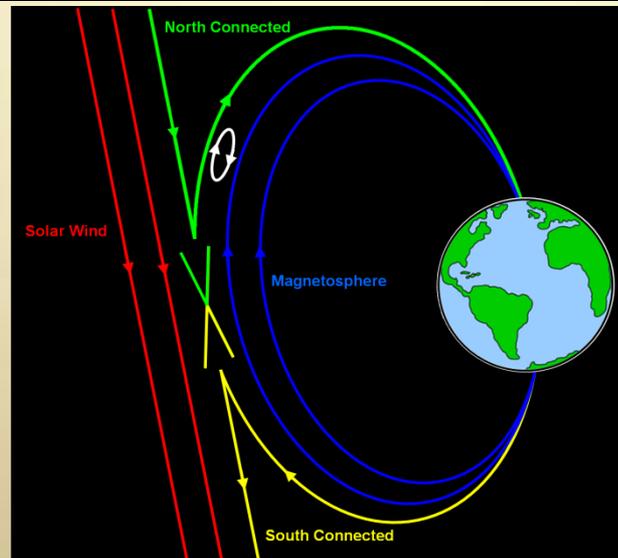
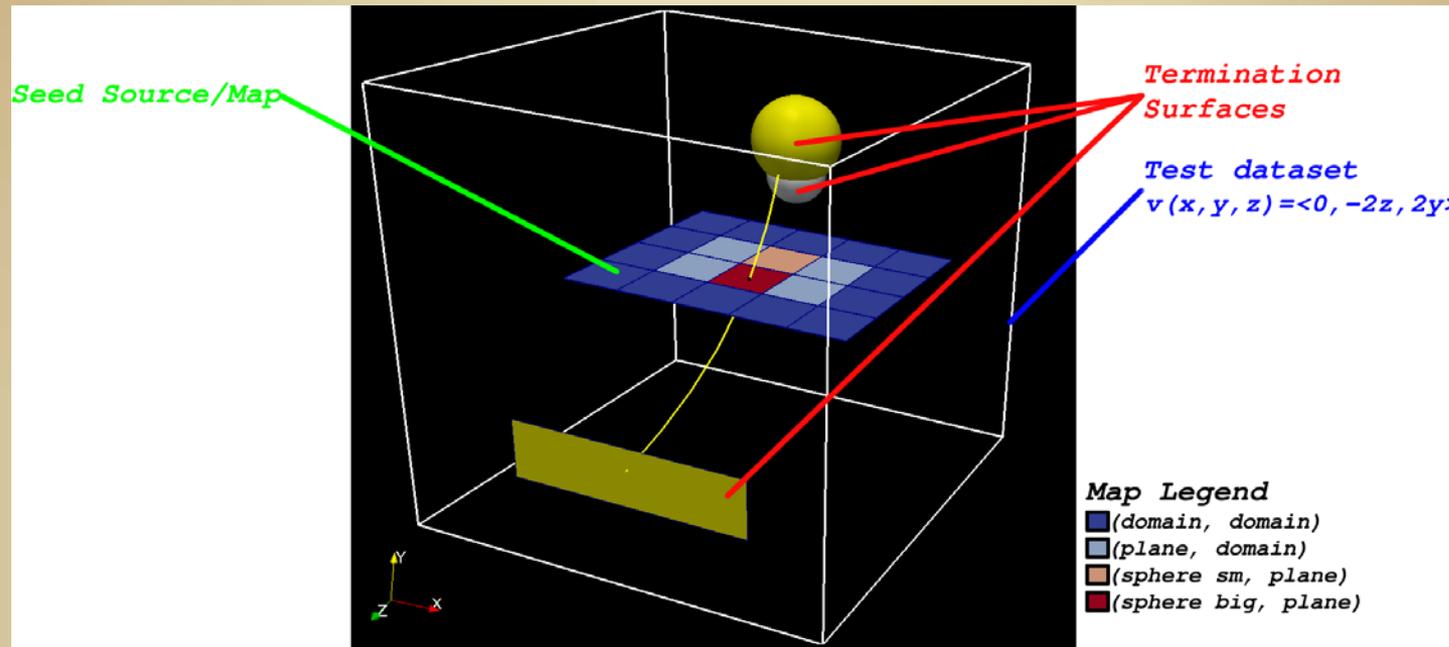
*Interactive physics-based magnetic topology visualization tool (plug in for ParaView):*

- visually spot interesting regions
- highly complex, 2D slices, tracking difficult

LEGEND		
CLASS	COLOR	DESCRIPTION
d ↔ d	Red	solar wind
o ↔ d	Light Red	
i ↔ d	Light Red	
n ↔ s	Blue	closed magnetosphere
n ↔ d	Green	northern hemisphere connected
n ↔ o	Dark Green	
n ↔ i	Dark Green	
n ↔ n	Dark Green	
s ↔ d	Yellow	southern hemisphere connected
s ↔ o	Light Yellow	
s ↔ i	Light Yellow	
s ↔ s	Orange	
o ↔ o	White	field null
i ↔ o	White	
i ↔ i	White	

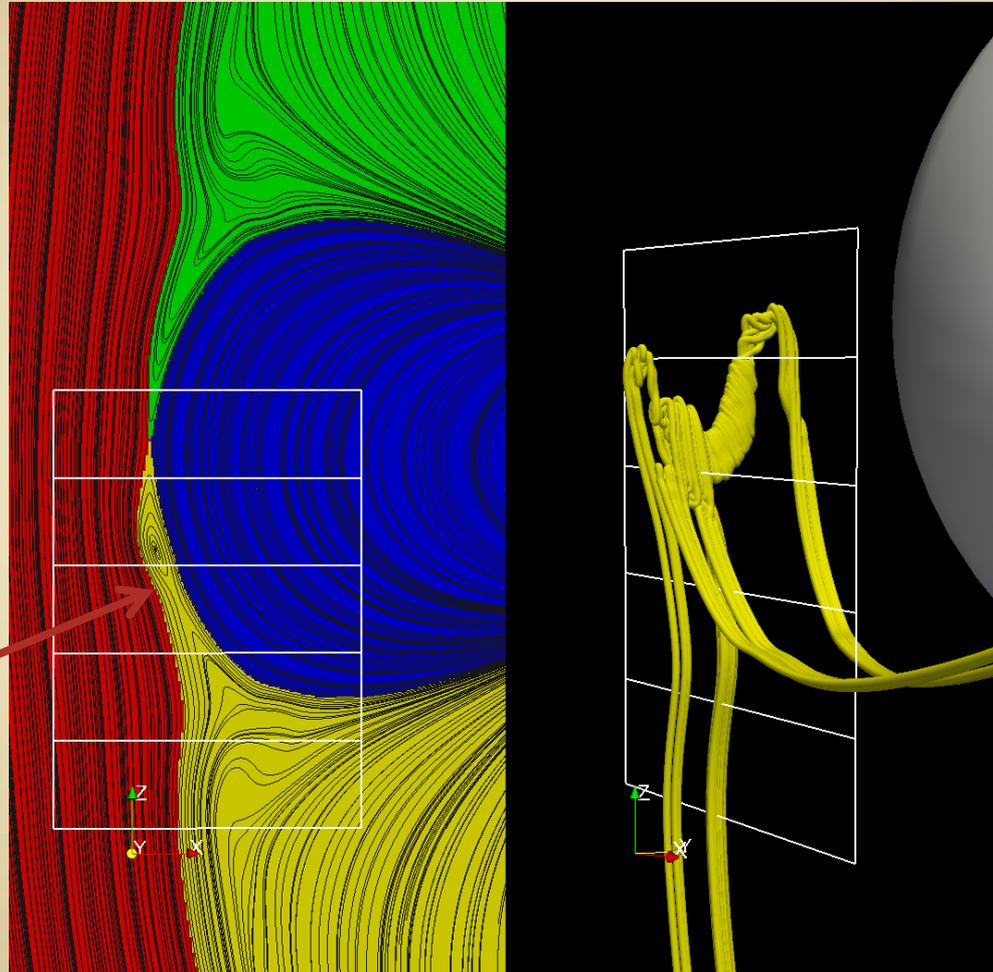


# Topomap – Key Diagnostic



- 5 topological classes of magnetic field lines:
- solar wind lines
  - magnetosphere lines
  - northern/southern ionosphere field lines
  - lines that end at a field null

# 3D Global MHD



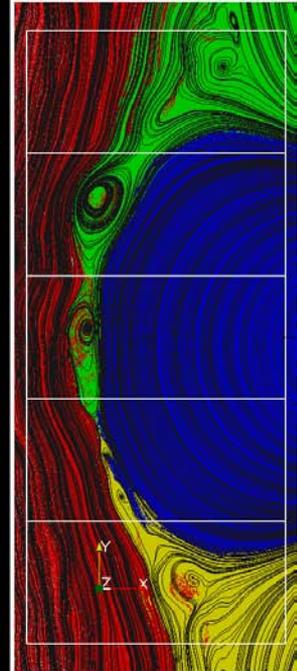
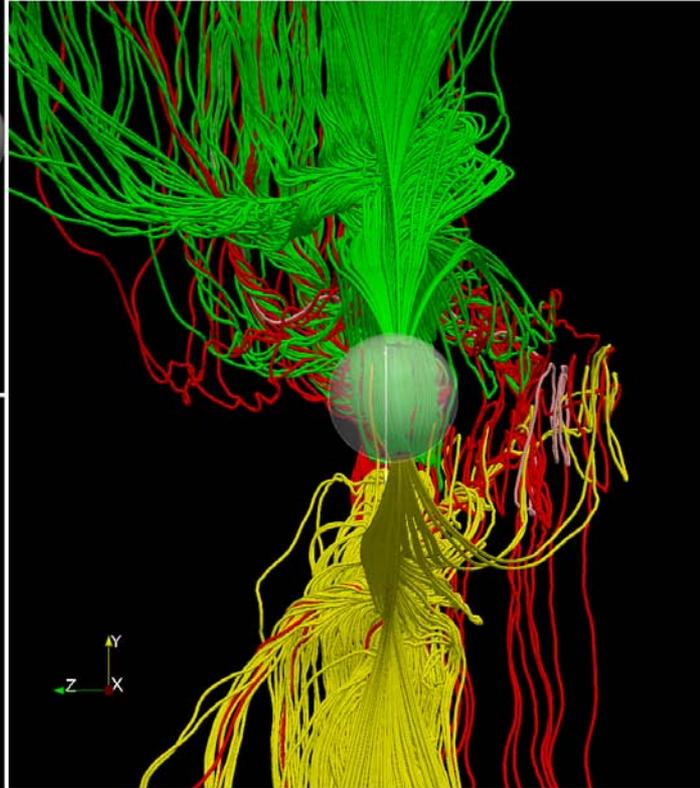
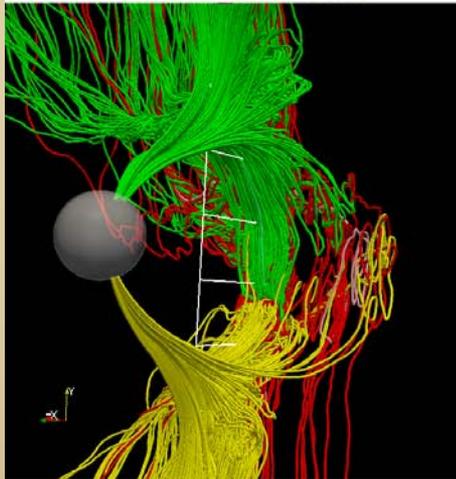
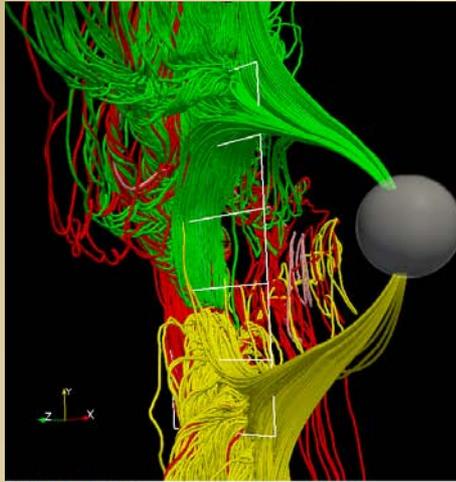
False X-line  
which may be  
misinterpreted  
as MXR

Purely Southward

No Dipole Tilt

Very Little Mixing

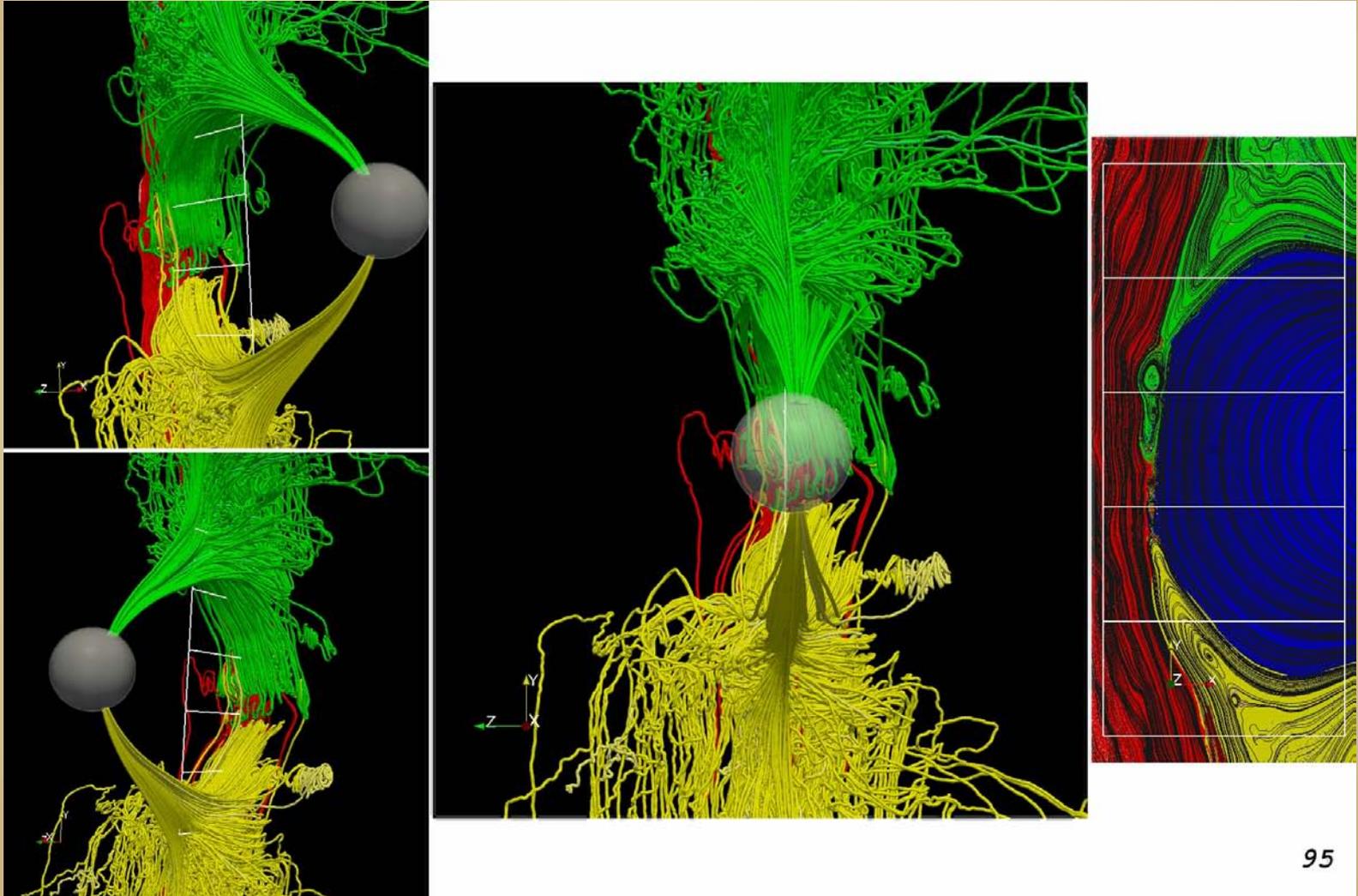
# 3D Global Hybrid



140

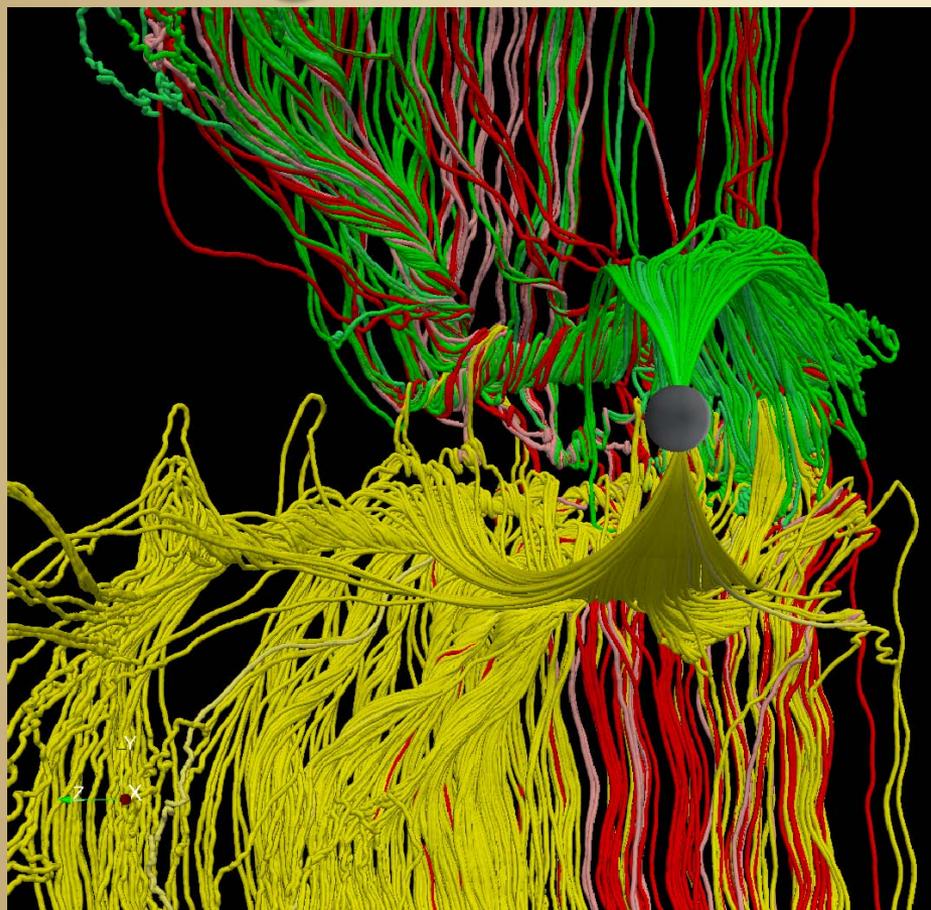
Run on 98 K cores on Kraken

# 3D Global Hybrid

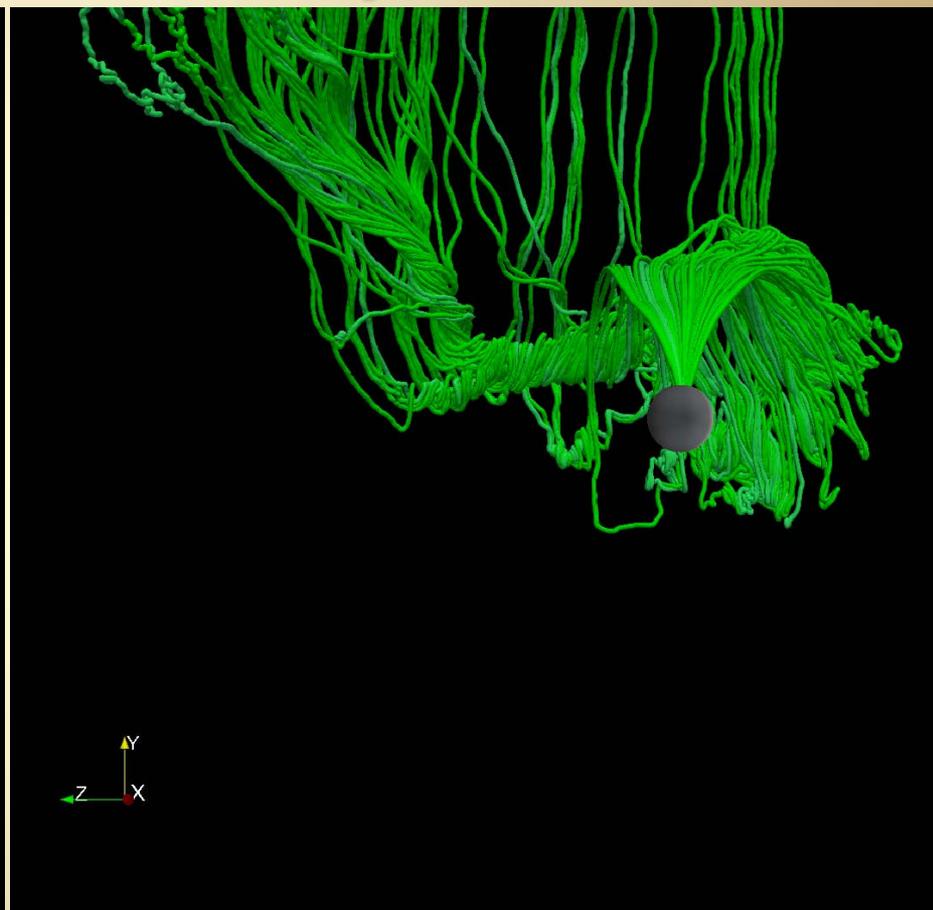


**Successful Tracking of Flux Ropes in Time**

# Mixing of Field Lines



# Shape of FTE



# Case study: Automated Detection of FTE in the Cluster Data

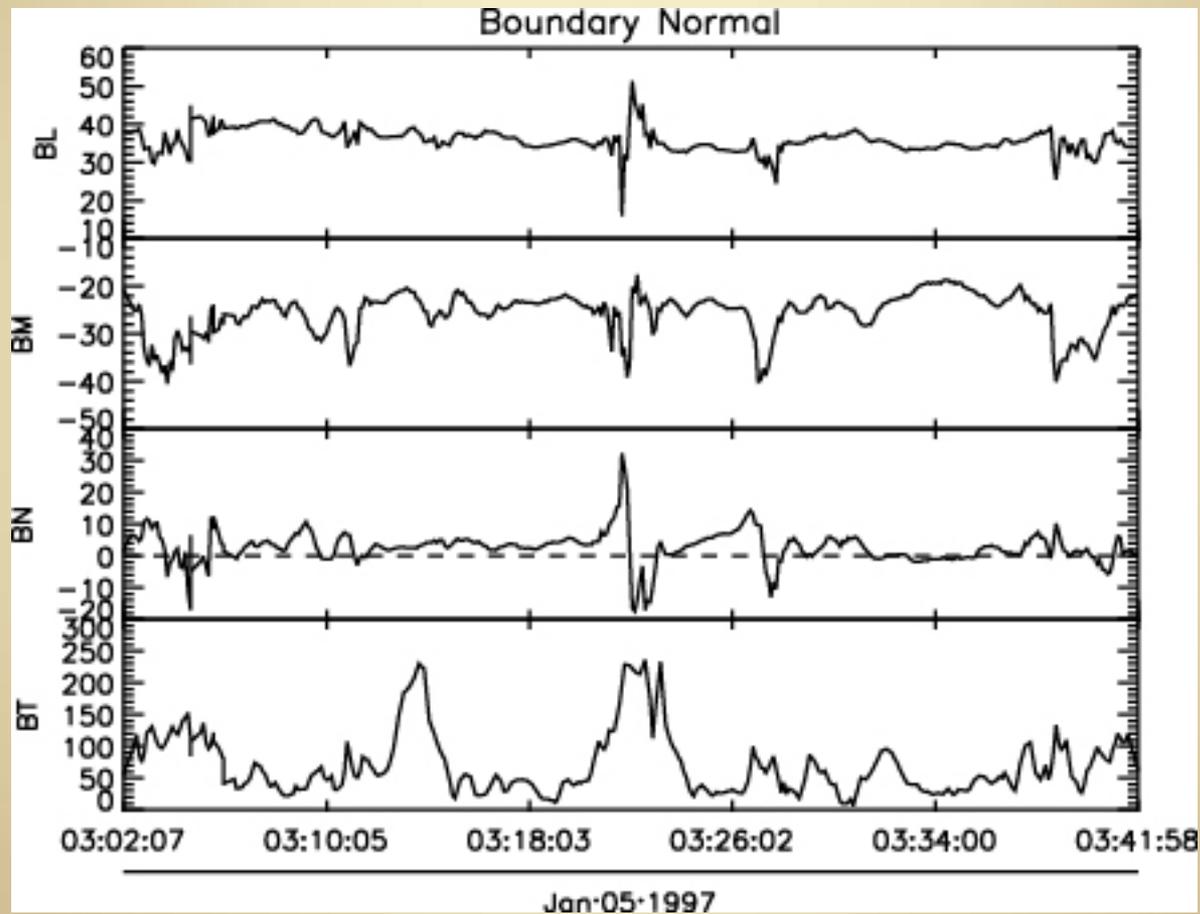
Problem: Traditional approach of discovering FTEs in cluster data is subjective, tedious, and limited.

Karimabadi, H., T. Sipes, Y. Wang, B. Lavraud, A. Roberts, JGR  
2009

# Comparison of traditional and automated approaches

<b>Traditional Approach Using Visual Data Analysis</b>	<b>Query Driven Analysis Using MineTool</b>
Transform to normal boundary coordinates using a model of MP	Automate the manual processes to place focus on answering questions:
Choose a specific criterion for FTE  D. Sibeck: clear bipolar signatures and well defined $B_M$ component  Y. Wang: clear bipolar signature in $B_N$ and $ B $ enhancements	<ul style="list-style-type: none"> <li>• Can FTEs be found in GSM coordinates?</li> <li>• Can FTEs be found based on plasma data alone?</li> <li>• What combination of variables are good indicators of FTEs?</li> </ul>
Visually inspect the data  ... ..  <b><i>Huge investment with limited return.</i></b>	<ul style="list-style-type: none"> <li>• Are there FTEs with no enhancement in <math> B </math>?</li> <li>• Does <math>B_M</math> have a well defined signature in all FTEs?</li> <li>• How complete is our taxonomy of FTEs?</li> </ul> <b><i>Huge return with limited investment.</i></b>

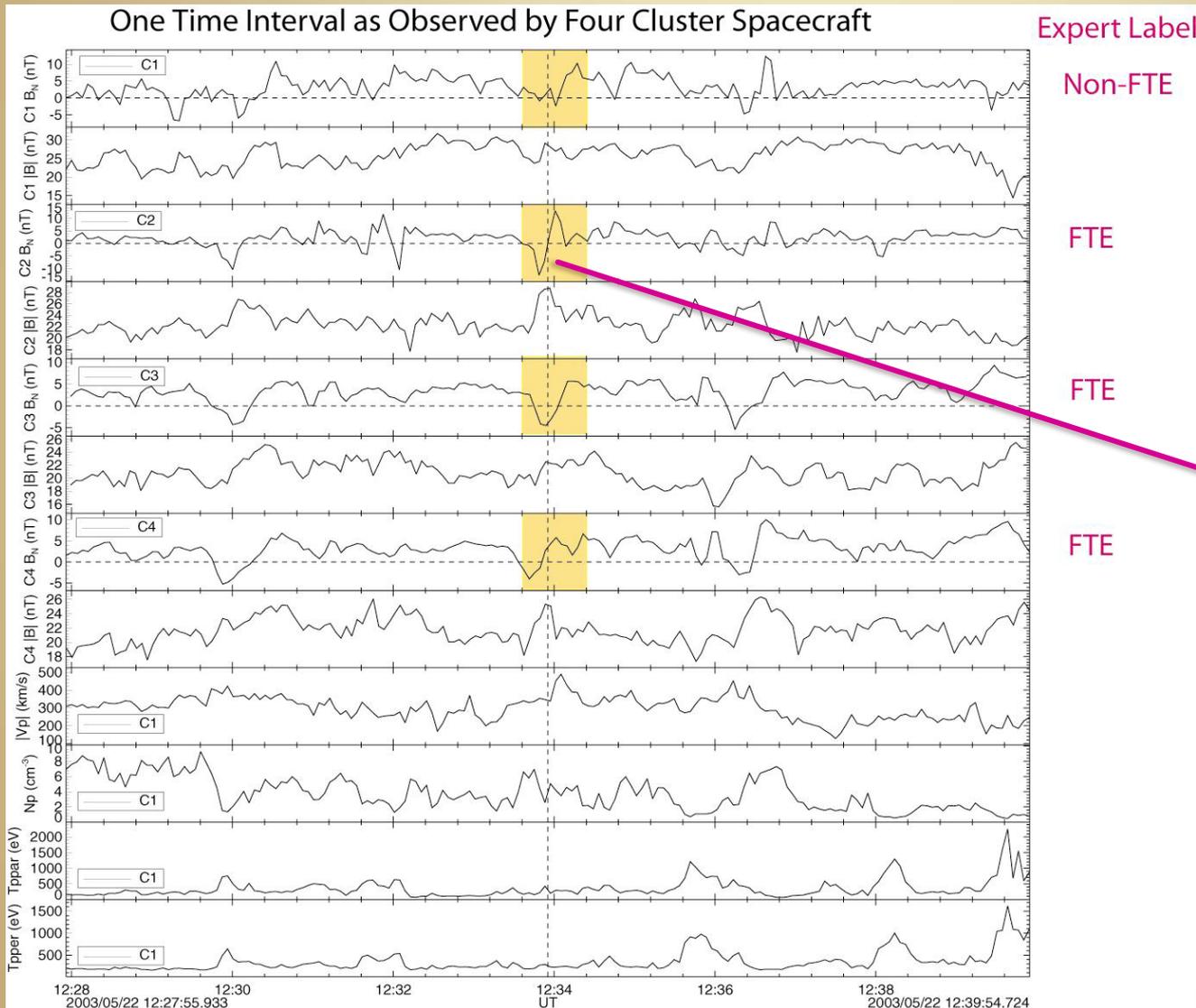
# Example of an FTE



# FTE Cluster Data Characteristics

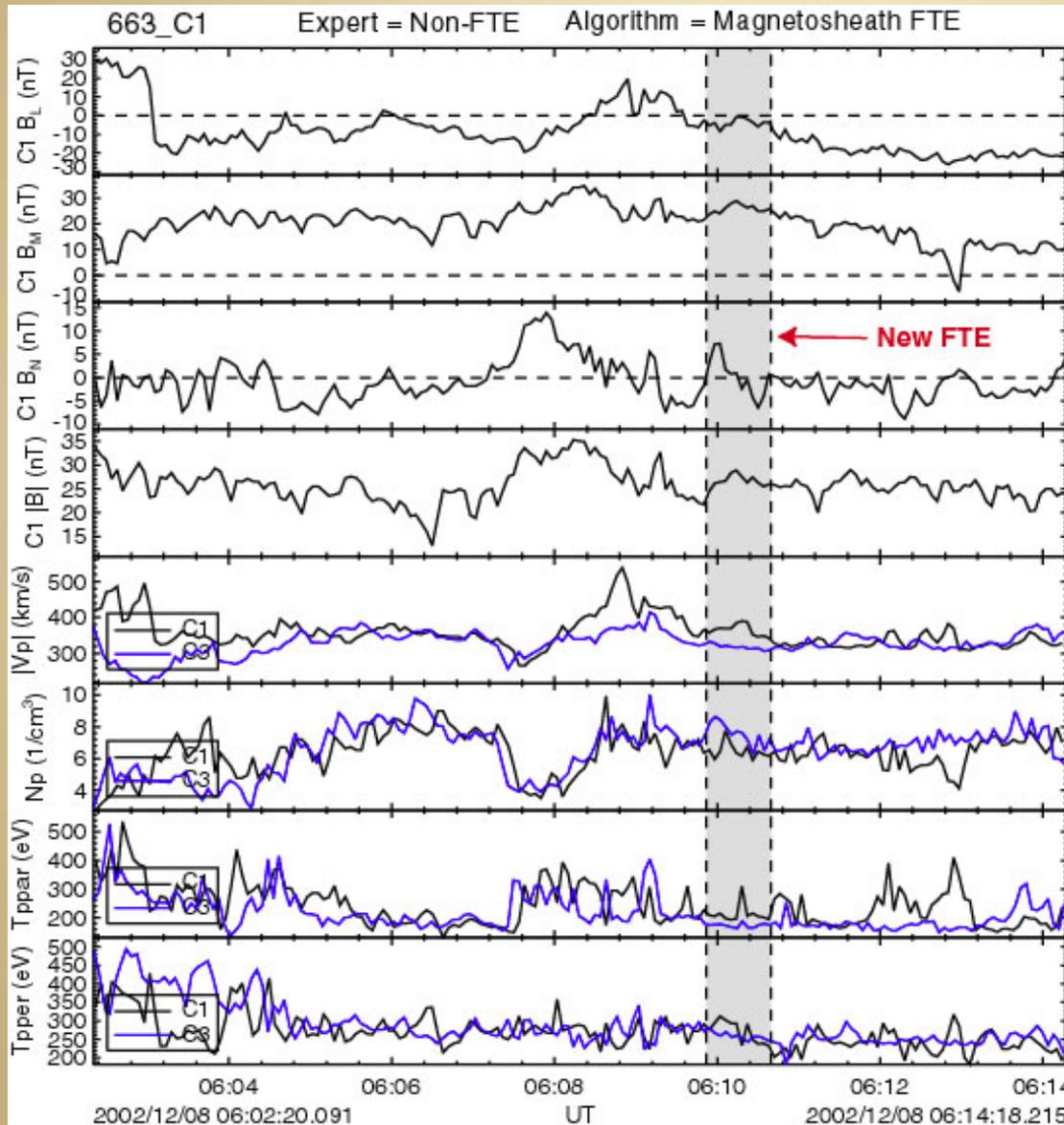
- Dataset containing magnetospheric and IMF condition measurements:
  - ▣ Bx, By, Bz, BL, BM, BN, Bmag, Np, Vx, Vy, Vz, Vmag, Tpar, Tper
- 615 events (~30% nonFTEs)
- Labeled data:
  - ▣ magnetosheath FTEs, magnetospheric FTEs and nonFTEs
- Each time series was 12 minutes in duration

# Efficiency Gains with Physics Mining



Physics mining only asks that you indicate the existence of an event - no need to provide start and end times, as in traditional approach.

# Successful Discovery of FTE in Cluster Data



Physics  
mining  
successfully  
identified new  
FTE that was  
missed by  
visual  
inspection

# Physics Mining – Confidence Study

Three-Tier Classification		Classification of FTE vs Non-FTE				Classification of Msheath vs Msphere FTE			
MODEL	Correctly classified (P <sub>3</sub> )	Correctly classified (FTE vs non-FTE) (P <sub>1</sub> )	# of false non-FTEs		# of false FTEs		Correctly classified (msth vs msph FTE) (P <sub>2</sub> )	# of false msth FTEs	# of false msph FTEs
			msth	msph	msth	msph			
T <sub>⊥</sub> /T <sub>∥</sub>	59%	83%	7	0	23	4	67%	50	0
N <sub>p</sub>	77%	92%	2	0	13	1	79%	30	1
B <sub>tot</sub>	91%	95%	0	0	8	2	94%	0	9
B <sub>z</sub>	92%	93%	9	0	14	0	98%	0	3
T <sub>p</sub>	92%	97%	0	0	6	0	93%	0	10
B <sub>L</sub>	92%	96%	0	0	7	1	95%	3	5
T <sub>∥</sub>	93%	97%	0	0	5	1	94%	0	9
T <sub>⊥</sub>	93%	99%	0	0	2	0	91%	11	2
V <sub>x</sub>	93%	98%	0	0	4	0	93%	7	4
B <sub>x</sub>	93%	94%	0	0	11	1	99%	0	2
V <sub>z</sub>	94%	97%	0	0	6	0	95%	0	7
B <sub>y</sub>	94%	98%	0	0	4	1	95%	0	8
B <sub>N</sub>	95%	99%	0	0	2	0	94%	3	6
V <sub>y</sub>	96%	98%	0	0	3	2	98%	0	3
B <sub>M</sub>	96%	99%	0	0	2	1	97%	0	5

# Physics Mining – Confidence Study

Three-Tier Classification		Classification of FTE vs Non-FTE				Classification of Msheath vs Msphere FTE			
MODEL	Correctly classified	Correctly classified(FTE vs non-FTE)	# of false non-FTEs		# of false FTEs		Correctly classified (msth vs msph FTE)	# of false msth FTEs	# of false msph FTEs
			msth	msph	msth	msph			
$B_x, B_{tot}$	93%	95%	0	0	10	0	97%	0	4
$B_x, B_y, B_z, B_{tot}$	95%	95%	0	0	10	0	100%	0	0
$B_x, B_y, B_z, B_{tot}, T_p, N_p$	96%	96%	0	0	8	0	99%	0	1
$B_N, B_{tot}, N_p$	96%	97%	0	0	6	0	99%	0	2
$B_N, B_{tot}$	96%	98%	0	0	5	0	98%	0	3
$B_N, B_{tot}, N_p$	96%	97%	0	0	6	0	99%	0	2
$B_x, B_y, B_z, B_{tot}, V_x, V_y, V_z, T_p, N_p$	96%	96%	0	0	8	0	100%	0	0
$T_{\perp}, T_{\parallel}$	97%	99%	1	0	1	0	97%	0	5

# Physics Mining – Confidence Study

## Three-Tier Classification

## Classification of FTE vs Non-FTE

## Classification of Msheath vs Msphere FTE

MODEL	Correctly classified (P <sub>3</sub> )	Correctly classified (FTE vs non-FTE) (P <sub>1</sub> )	# of false non-FTEs		# of false FTEs		Correctly classified (msth vs msph FTE) (P <sub>2</sub> )	# of false msth FTEs	# of false msph FTEs
			msth	msph	msth	msph			
B <sub>L</sub> , B <sub>M</sub> , B <sub>N</sub> , B <sub>tot</sub> , T <sub>p</sub> , N <sub>p</sub>	97%	99%	0	0	3	0	97%	0	4
B <sub>L</sub> , B <sub>M</sub> , B <sub>N</sub> , B <sub>tot</sub>	98%	98%	0	0	4	0	99%	0	1
B <sub>x</sub> , B <sub>y</sub> , B <sub>z</sub> , B <sub>tot</sub> , V <sub>x</sub> , V <sub>y</sub> , V <sub>z</sub> , N <sub>p</sub>	98%	99%	0	0	3	0	99%	0	2
B <sub>L</sub> , B <sub>M</sub> , B <sub>N</sub> , B <sub>tot</sub> , N <sub>p</sub> , V <sub>x</sub> , V <sub>y</sub> , V <sub>z</sub> , T <sub>p</sub>	98%	98%	0	0	4	0	100%	0	0
B <sub>L</sub> , B <sub>M</sub> , B <sub>N</sub> , B <sub>tot</sub> , N <sub>p</sub>	99%	99%	0	0	3	0	100%	0	0
B <sub>x</sub> , B <sub>tot</sub>	93%	95%	0	0	10	0	97%	0	4
B <sub>x</sub> , B <sub>y</sub> , B <sub>z</sub> , B <sub>tot</sub>	95%	95%	0	0	10	0	100%	0	0

# Questions answered using query driven analysis approach

Traditional Approach Using Visual Data Analysis	Query Driven Analysis Using MineTool
Transform to normal boundary coordinates using a model of MP	<ul style="list-style-type: none"> <li>• Can FTEs be found in GSM coordinates? <b>YES</b></li> </ul>
Choose a specific criterion for FTE  D. Sibeck: clear bipolar signatures and well defined $B_M$ component  Y. Wang: clear bipolar signature in $B_N$ and $ B $ enhancements	<ul style="list-style-type: none"> <li>• Can FTEs be found based on plasma data alone? <b>YES</b></li> <li>• What combination of variables are good indicators of FTEs? <b>Surprising results.</b></li> <li>• Are there FTEs with no enhancement in <math> B </math>? <b>YES</b></li> </ul>
Visually inspect the data  <p style="text-align: center;"><b>3 years later...</b></p>	<ul style="list-style-type: none"> <li>• Does <math>B_M</math> have a well defined signature in all FTEs? <b>NO. Signatures can be visually indiscernible to human eye.</b></li> <li>• How complete is our taxonomy of FTEs? <b>Not diverse enough for statistical study.</b></li> </ul> <p style="text-align: center;"><b>1 hour later, science advances</b></p>

# What Problems Can It Solve?

- Feature extraction and tracking
  - Substorm injection detection, FTEs, MP crossings
- Classification
  - Flux ropes in the SW
- Dependency analysis
  - What triggers FTEs?
- Modeling
  - Developing 3D model of magnetopause from s/c data
- Anomaly detection
  - Spacecraft charging

# Advantage of Intelligent Archiving

## ✘ Current Archive Search

- + Search general description of data set for type of data.
- + Requires analyzing data to discover events of interest.

## ✘ Intelligent Archive Search

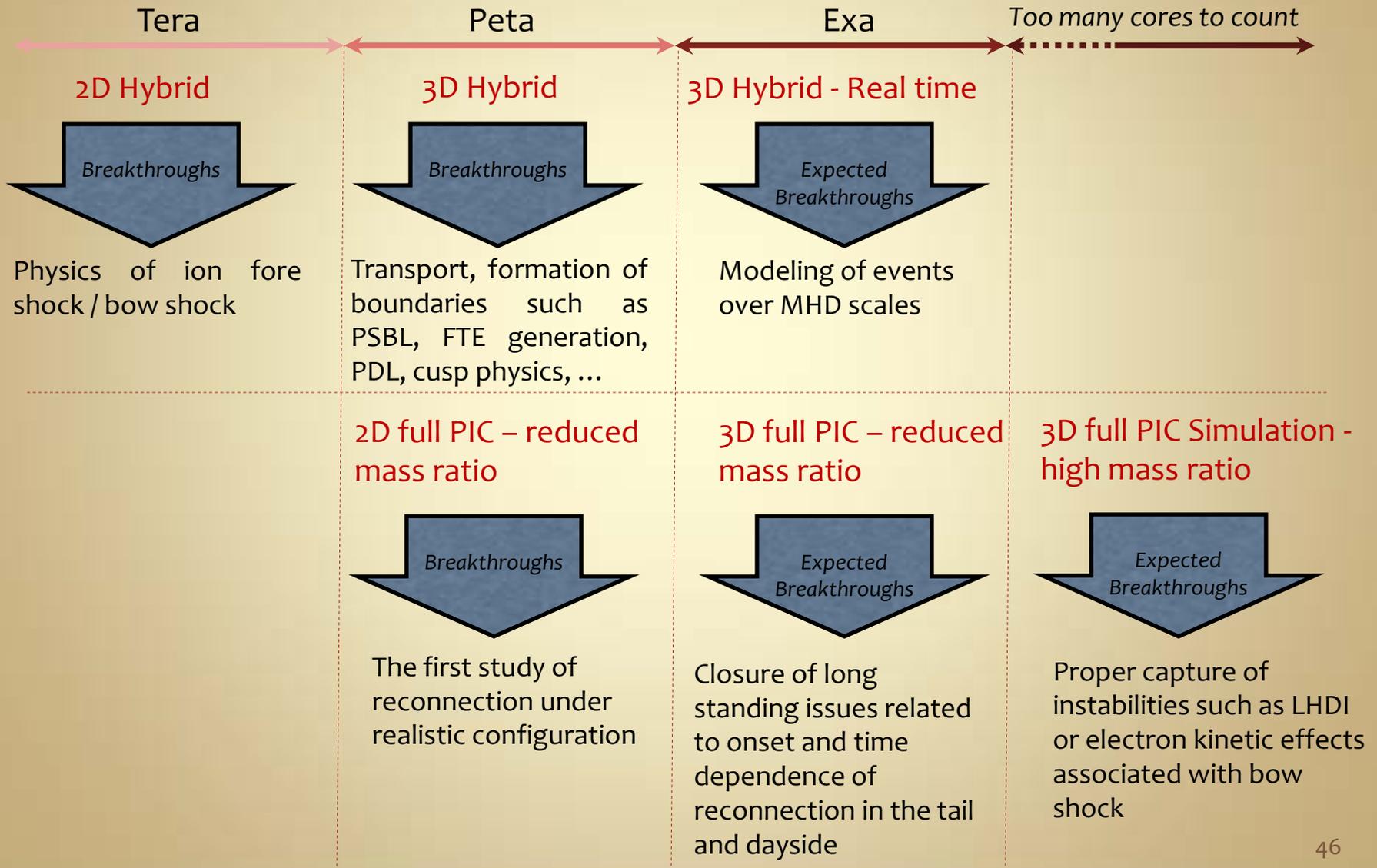
- + Automatically generate metadata about events in the data set using data mining techniques.
- + Search for data **content** - what the data set contains.
- + Events of interest are already identified and accessible to the user through a query.

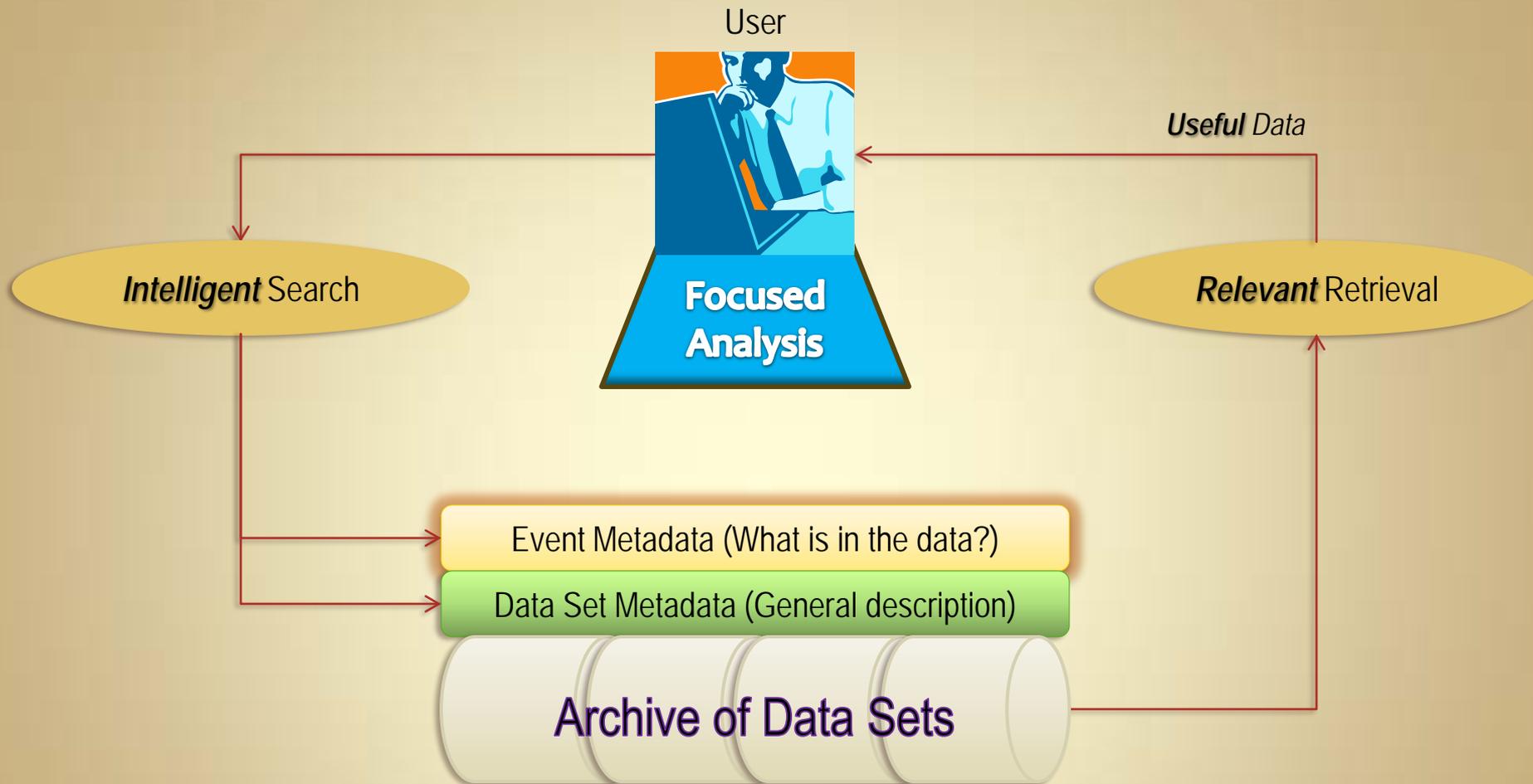
***Knowledge is extracted directly from the data about events, promoting intelligent search on archives.***

# Summary

- Knowledge discovery from increasingly complex and large data sets is a major bottleneck to progress in space sciences today.
- Physics mining techniques are powerful tools that optimize knowledge extraction from large data sets.
- MineTool has been successfully applied to several space physics problems.
- Intelligent archiving results in significant increase in ROI by extracting knowledge about the data content.

# Expected Scientific Breakthroughs Using Global Simulations on Exascale





**The creation of event metadata associated with the data sets is essential to extend reach and richness of data, and to drive collaboration by making events visible to the community of users. Enhancing the metadata allows Intelligent Search, significantly increasing value of the archive.**