AWS
re:Invent

# The Ideal Data Science Platform



Streaming

Data Marts

Transactional Stores

External Sources

python™

R

Point-Click | Code

Data Access & Prep

Modeling, ML, DL

Batch Scoring

Real-Time Scoring

Live Applications

Edge Devices

Offer Cross-Sell Products

Predict Impending Equipment Failure

Optimize Pricing
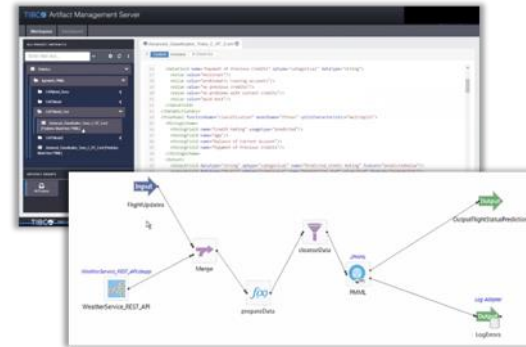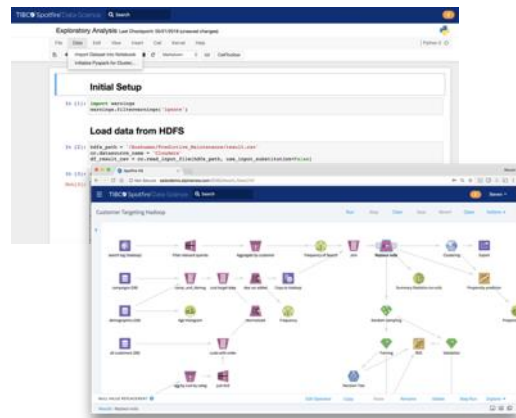
Prevent Fraud

Manage Inventory

TIBC®

# Agenda

- TIBCO Data Science and AWS Marketplace

- The TIBCO Connected Intelligence Cloud

- Anomaly Detection and Analysis

- ***Demonstration – Spatial Anomaly Analysis***

- Links and Assets

TIBC🌐®

# TIBCO Data Science

| | **Data Access/Prep** | **Modeling** | **Operations** | **Business Apps** |
|---|---|---|---|---|
| **FUNCTION** | + Distributed compute<br>+ Feature engineering<br>+ Reusable templates | + Visual composition<br>+ Multilingual notebook<br>+ Native ML & OS<br>+ Auto-ML, data prep | + Model lifecycle management<br>+ Batch automation<br>+ Real-time event processing<br>+ REST, applications, embedding | + Engineering/IoT<br>+ Customer analytics<br>+ Risk management<br>+ Supply chain |



**Medic;** e.g., researcher on epidemic monitoring

**Engineer;** e.g., aerodynamics engineer

**Marketeer;** e.g., customer engagement analyst

| | | | | |
|---|---|---|---|---|
| **USER or AUTOMATION** | *Data Scientist*<br>*Citizen Data Scientist* | *Data Scientist*<br>*Citizen Data Scientist* | *Analytics Operations*<br>*IT / Software Engineer* | *Business User*<br>*Analytics Operations*<br>*IT / Administration* |

TIBCO®

# TIBCO Data Science on AWS

*TIBCO DS on AWS Marketplace; Biggest vCPU Grid; Lightest Serverless Footprint*

# Data Science in the Cloud: Leidos Healthcare Analytics

Leidos Collaborative Advanced Analytics & Data Sharing Platform (CAADS) uses TIBCO Data Science and AWS to deliver analytics services in Healthcare



| **CDC** | **NIH** | **CMS** | **NASA** |
|---|---|---|---|
| Disease Outbreaks: Determining the cause of an HIV outbreak in the Midwest | Disease Outbreaks: Run simulations of disease propagation to guide public policy, specifically around the Zika virus | Data Governance: Analyzing and consolidating data around emerging Healthcare policies across 56 regions in the United States | Space Exploration: Analyzing human factors that affect the ability to transport astronauts on long fights (e.g., to Mars) |

TIBCO®

# TIBCO analytics transformation platform

*Powered by shared data assets*



**CONNECTED INTELLIGENCE**

ANALYTICS ACTIONS

**AUGMENT INTELLIGENCE**

- Visual Analytics
- Data Science
- Operational Analytics

**UNIFY DATA**

DATA OPERATIONS

- Spark Compute & Store
- Data Virtualization
- Data Store 1
- Data Store n
- EDW/Lake

INFORMATION MANAGEMENT

METADATA

- Metadata (Data Governance, Data Catalog)

MDM / RDM

- Master and Reference Data Management

Metadata

Master & Reference Data

Transactional Data

**INTERCONNECT EVERYTHING**

EVENTS

- Streaming / Events

INTEGRATION

- Hybrid Integration
- API Management

DATA SOURCES

- Streaming Sources
- In-memory Data Grid
- Enterprise Application 1 ... Enterprise Application n

Cloud Native    Open Platform    AI Foundation

TIBCO®

# TIBCO Data Science and AWS



**Data Sources**
- Streaming
- Data Marts
- Transactional Stores
- External Sources

**EDWs & Analytics Sandboxes**
- Amazon EMR
- Amazon Redshift
- Amazon Simple Storage Service (S3)

Real-Time (BW)

Batch (DV)

**Analytics Platform**

Federated Data Science Services

- SQL
- Spark
- Amazon SageMaker

TIBCO Data Science
- Automation
- Visual UI / Notebook
- Collaboration / Project

Manage | Execute

Visual Analytics (TIBCO Spotfire)

APIs and Microservices

Containers/Code

**Operational Systems**
- Amazon SageMaker
- Business Events & StreamBase
- Flogo
- Model Management

TIBCO Live Apps

TIBCO Cloud Integration
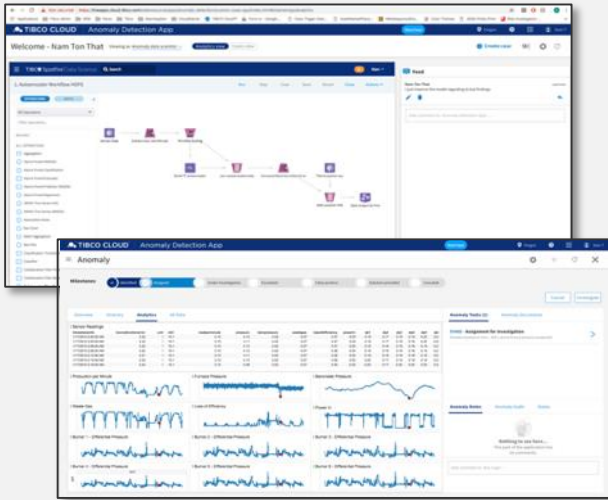
TIBCO CLOUD Live Apps

AWS Deployments

# TIBCO Data Science Solutions on AWS

## Cloud Apps: *Visual Analytics, Data Science, Streaming, Case Management*
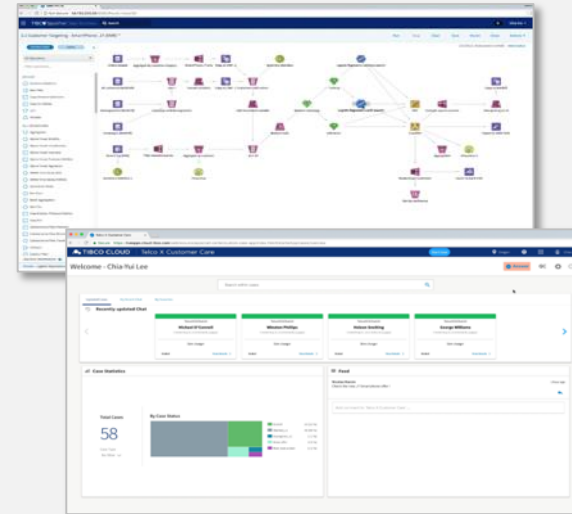
### Anomaly Detection



### Risk Management



### Customer Engagement



### Starter Set

Process Mining

IoT Analytics

Anomaly Detection

Risk Management

Customer Engagement

Blockchain – Dovetail

Partner Management

Starter Toolkit

Review Status: *TIBCO Spotfire*
*Identify issues, sweet spots*

Model: *TIBCO Data Science*
Supervised: Train
Unsupervised: Anomalies

Analyze Event Stream:
*TIBCO Flogo, Cloud Integration*
Batch and Real-Time Updates

Case Manage: *TIBCO Live Apps*
Investigate identified cases
Audit trail + recycle
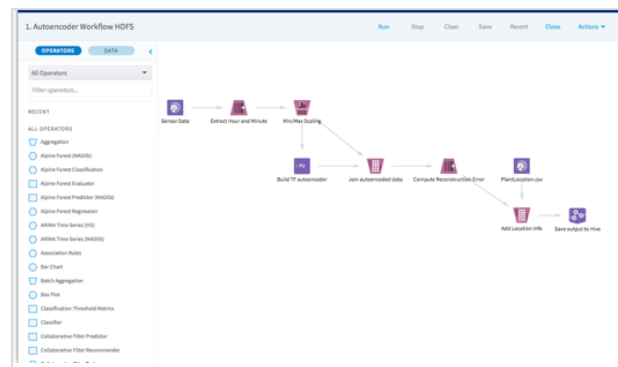
TIBCO®

# Anomaly Analysis Solution Overview

**1** Collect data from equipment, normalize, model to predict magnitude of anomaly – **TIBCO Data Science & AWS**
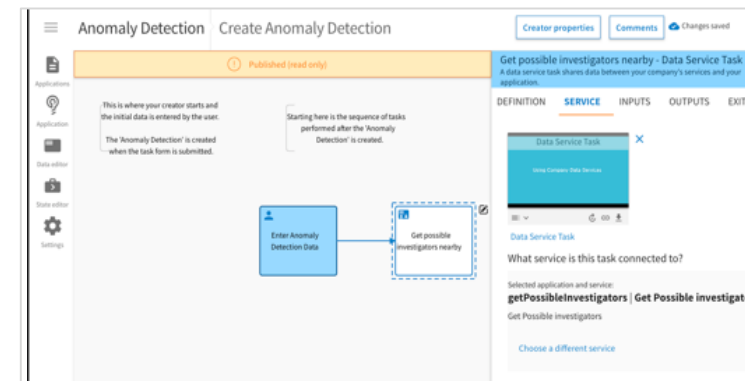


**3** Alert raised and case created – **TIBCO Cloud Live Apps**



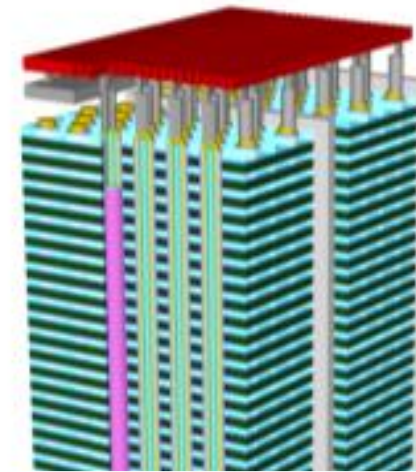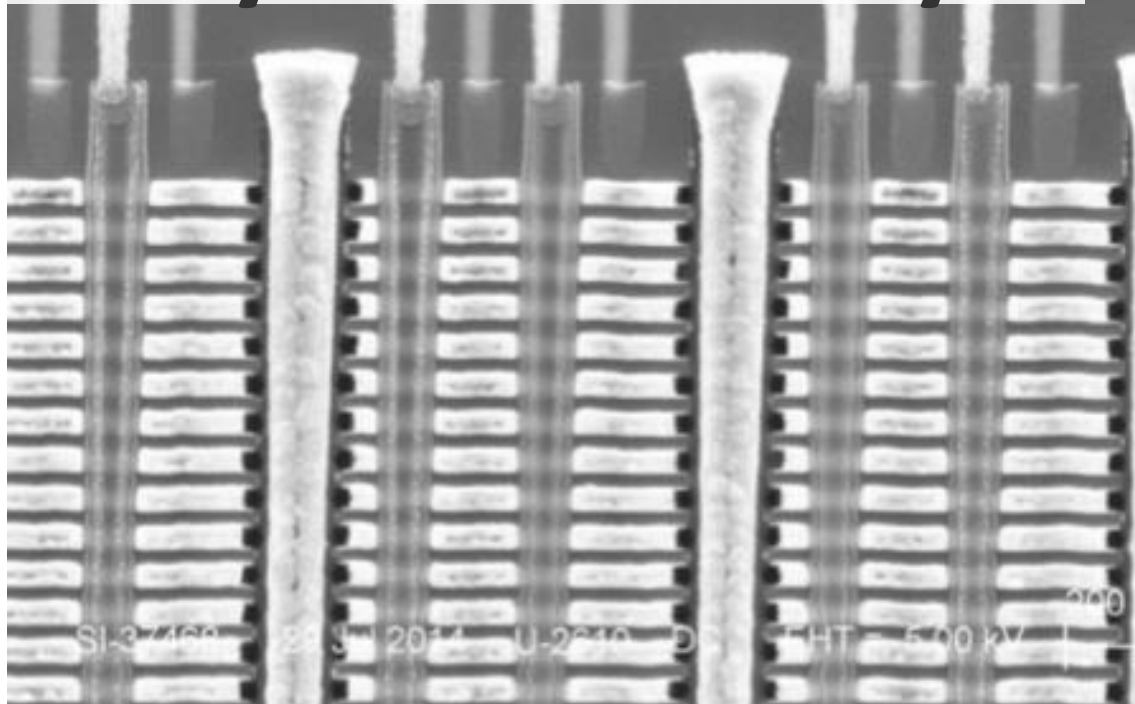**2** Model detects anomaly – **TIBCO Data Science**



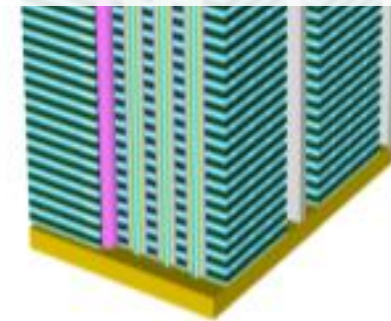**4** Case manager investigates and takes action to the equipment – **TIBCO Cloud Integration**

TIBC⊙

# Data Challenges in High-Tech Manufacturing
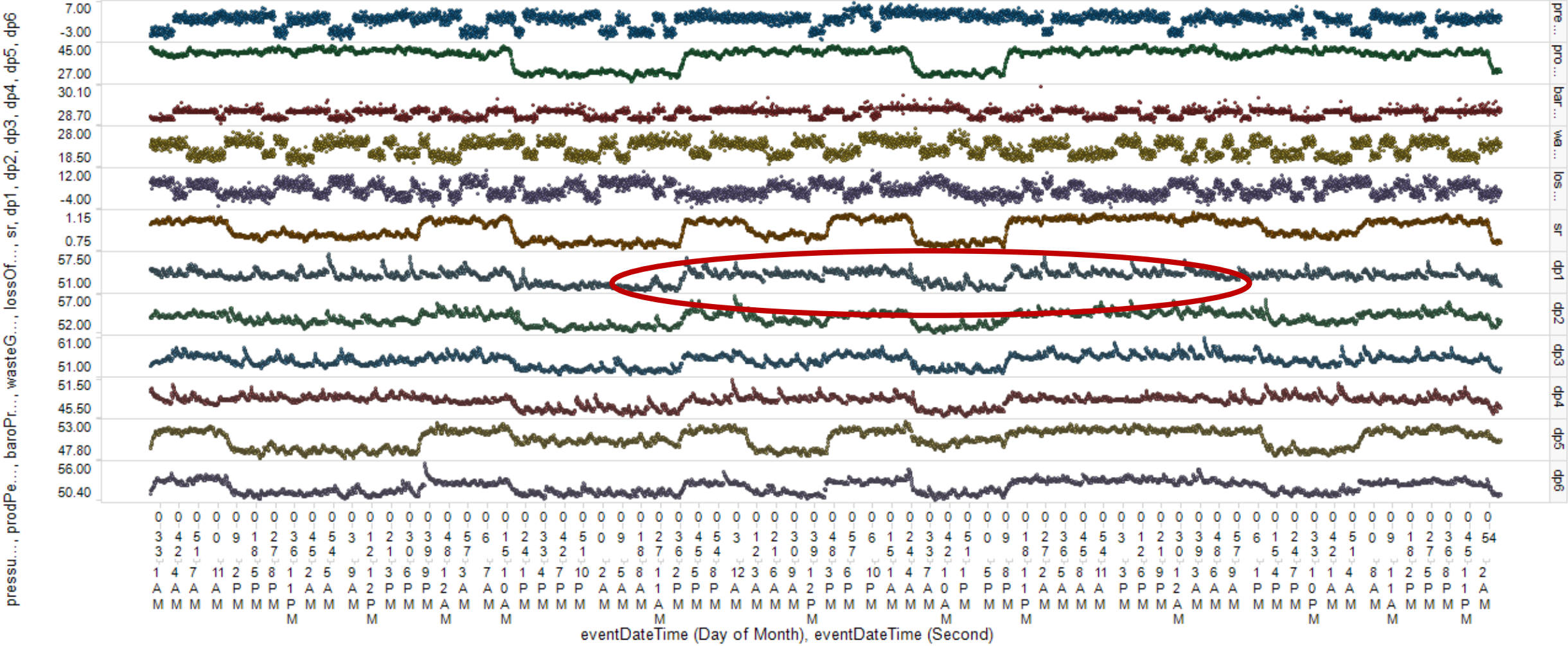
**Stack 3D Flash Memory Cell Layers Vertically**

**96 Memory Cell Layers**

TIBCO®

# Hot Paths to Anomaly Detection
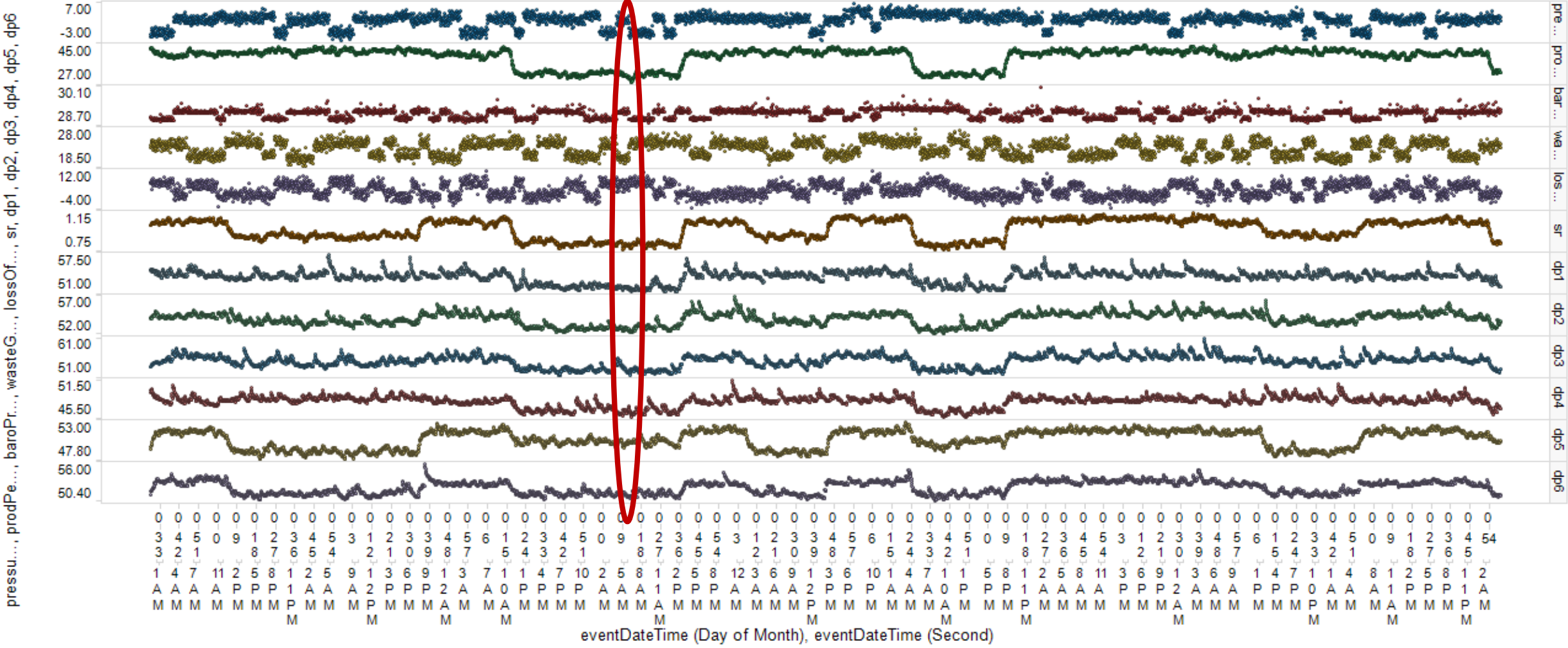
TIBC◆®

# Longitudinal Anomaly Analysis

# Cross-Sectional Anomaly Analysis



pressure, prodPerMinute, baroPressure, wasteGas, lossOfEfficiency, sr, dp1, dp2, dp3, dp4, dp5, dp6 vs. eventDateTime (Day of Month), eventDateTime (Second)
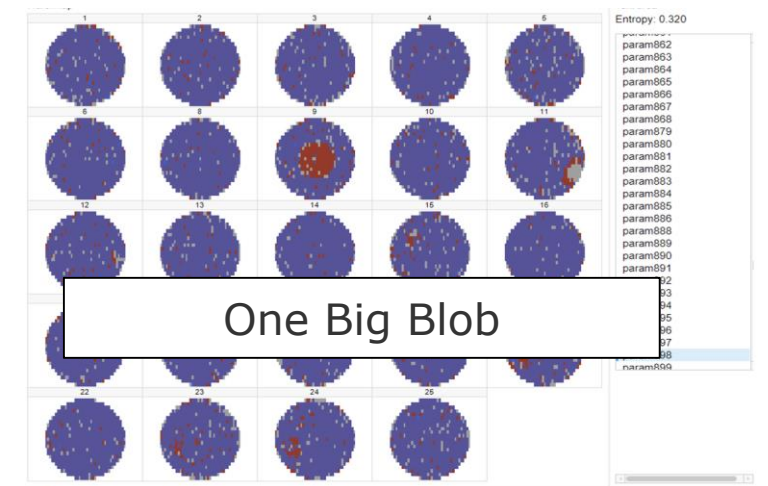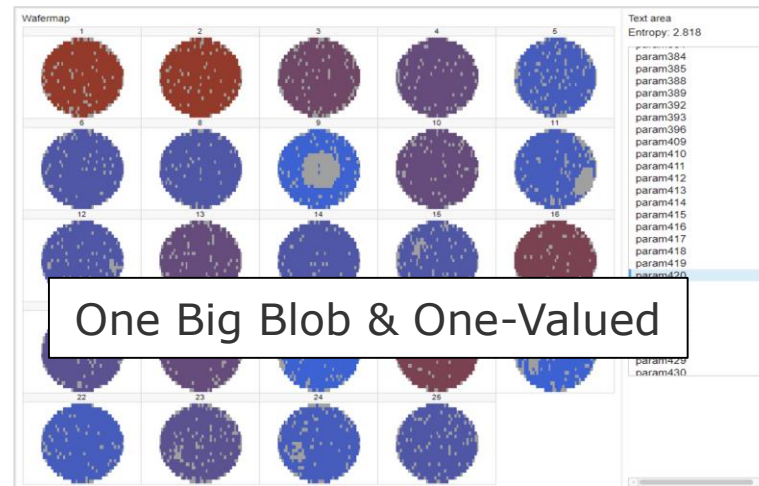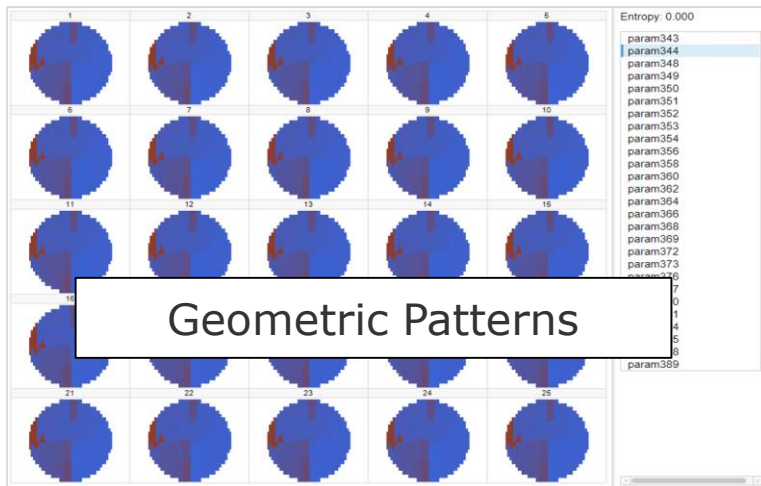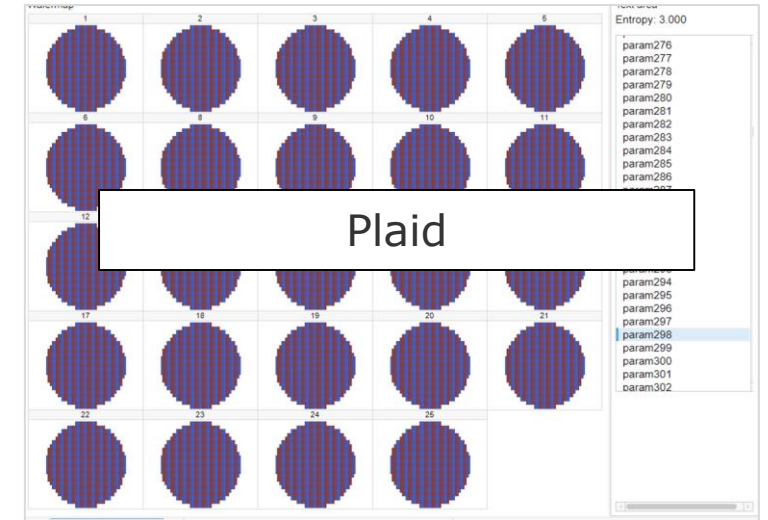
eventDateTime (Day of Month), eventDateTime (Second)

# Spatial Anomaly Analysis



One-Valued Wafer Maps

Big Red Spot

Plaid

Geometric Patterns

One Big Blob & One-Valued

One Big Blob

# Cross-Sectional Anomaly Detection

# Cross-Sectional Anomaly Detection

# Analysis: Cluster Incidents, View Signatures



© Copyright 2000-2019 TIBCO Software Inc.

# Analytic workflow – methodologies + demos

**Find and cluster anomalous events [Demo #1]**

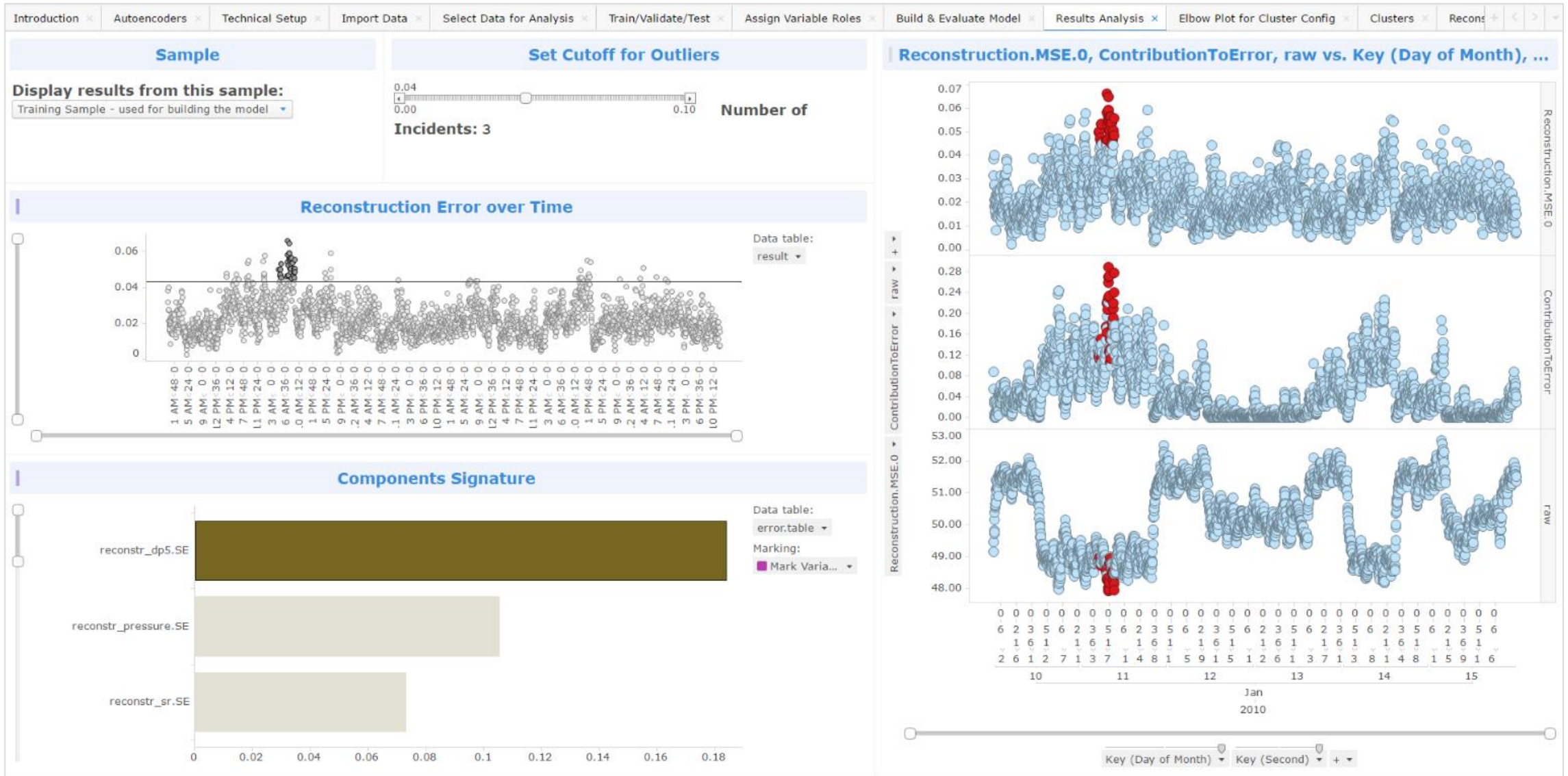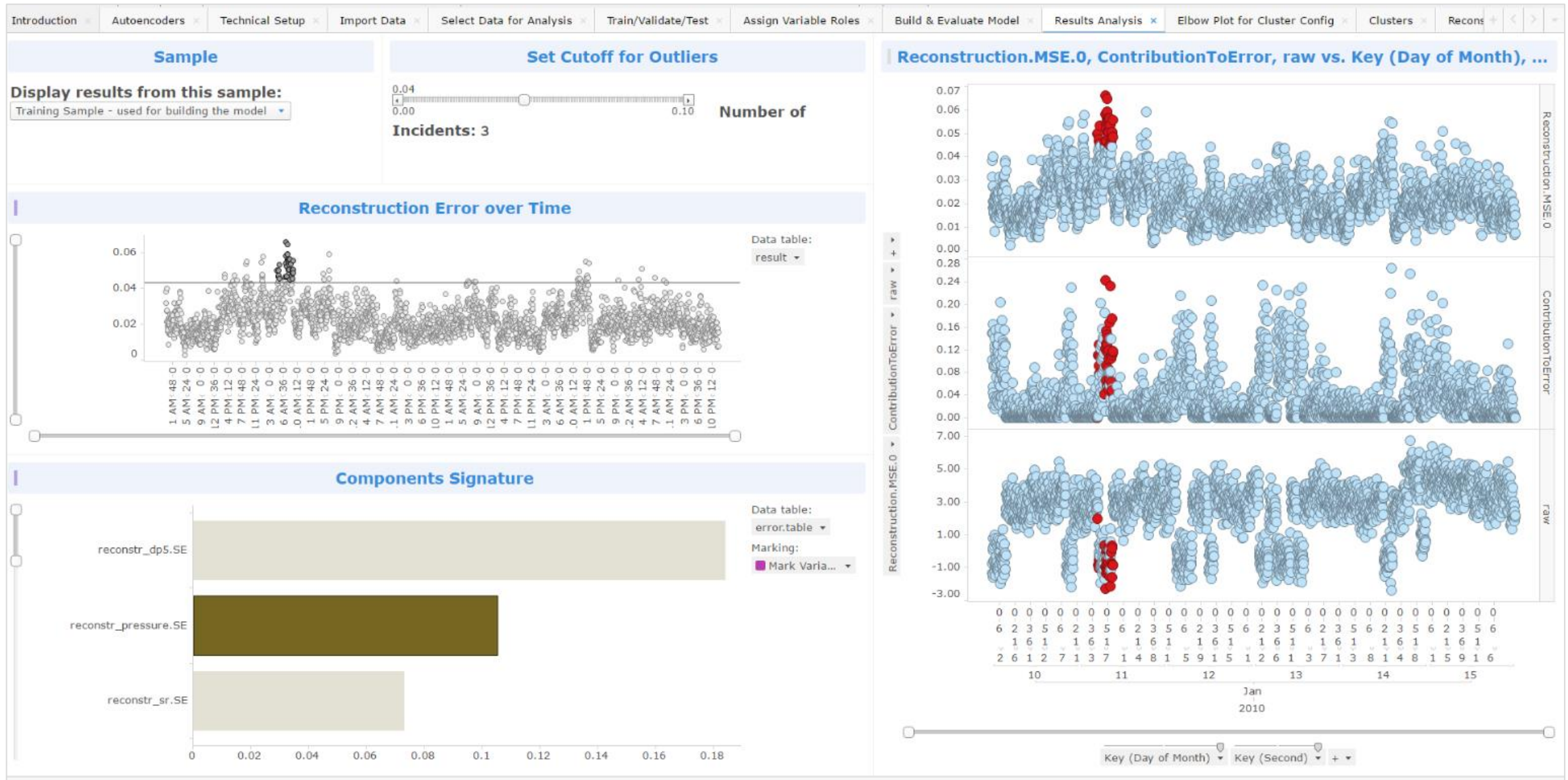- Transform wafer maps into vectorized coefficients, then cluster on quality
  - Many measured parameters, e.g., quality tests: storage fidelity, logic circuits
- Approach 1: SVD + K-means
  - Focus on failure mode parameters
- Approach 2: Bessel functions + hierarchical clustering
  - Radial basis functions
  - Rotationally invariant | Null-value tolerant | Efficient storage
  - Better than SVD + K-means for multi-parameter analysis

**Monitor anomalies [Demo #2]**
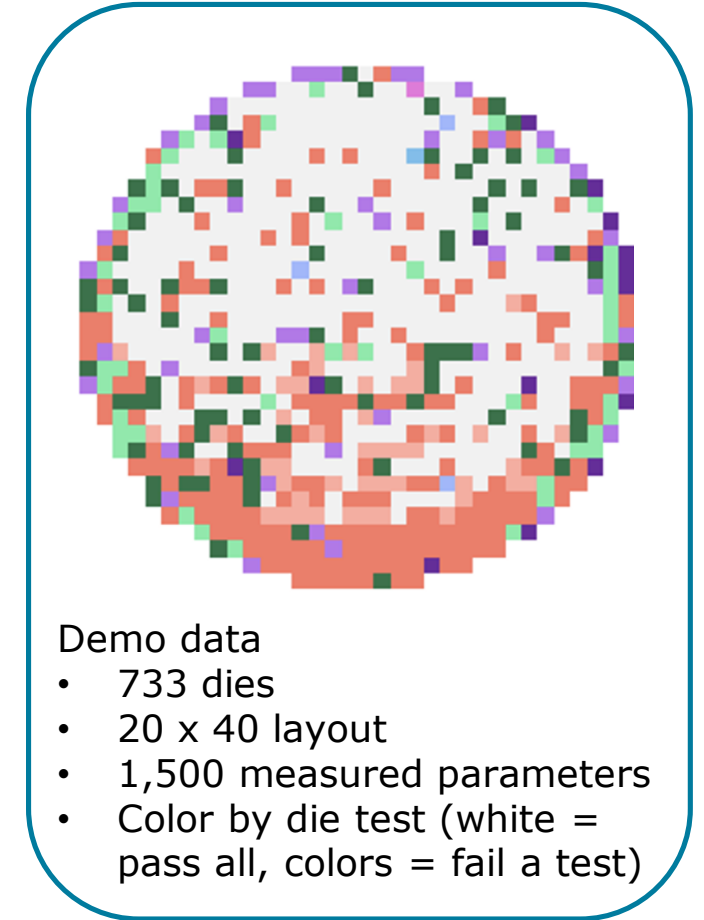
- Stream wafer data
- Vectorize and cluster

**Predict when and why anomalies occur [Demo #3]**

- Reduce dimensionality of very wide data
- Train models to determine sensor importance
- Identify responsible process parameters

*Process variable corrections/models rebasing*

- *Identify new patterns as they emerge (e.g., incident analysis)*
- *Factory monitoring staff can click to characterize the new pattern*

Demo data
- 733 dies
- 20 x 40 layout
- 1,500 measured parameters
- Color by die test (white = pass all, colors = fail a test)

TIBC⌀®

aws

TIBCO® Streaming

Capture and send data to livestream

Ingest data into Amazon Kinesis for further processing

Input

Read data from Kinesis into TIBCO Streaming

Σ

Preprocess and transform the data

Python

*(TIBCO Streaming operator for Python)*

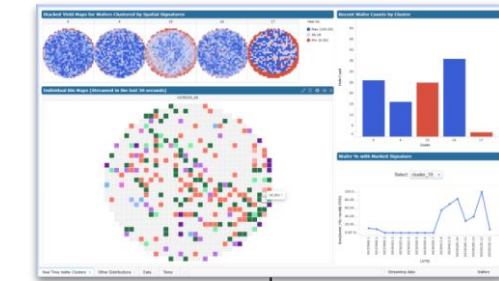Perform SVD on the bin values for wafers using Python

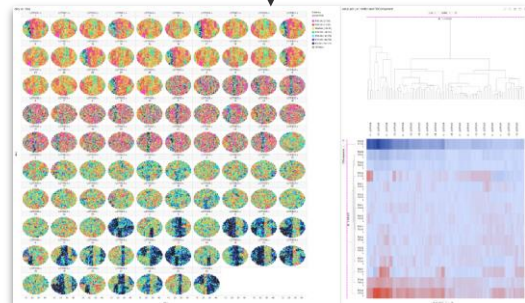TIBCO® Spotfire®

Send data for live viewing into Spotfire Data Streams

f(x)

Combine data & results

JPMML

PMML operator to cluster the data using a model trained in TIBCO Data Science
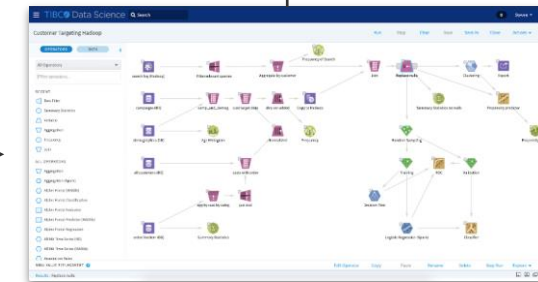
TIBCO® Data Science

Refine clusters, identify wafers of interest

Send data to Spark/Hadoop cluster for more in-depth analysis in TIBCO Data Science

Retrain models for clustering and root-cause analysis

TIBCO®

# Anomaly Analysis

Find and explore anomalous events

Monitor anomalies

Predict when and why they occur

TIBC☉®

# Anomaly Detection and Analysis
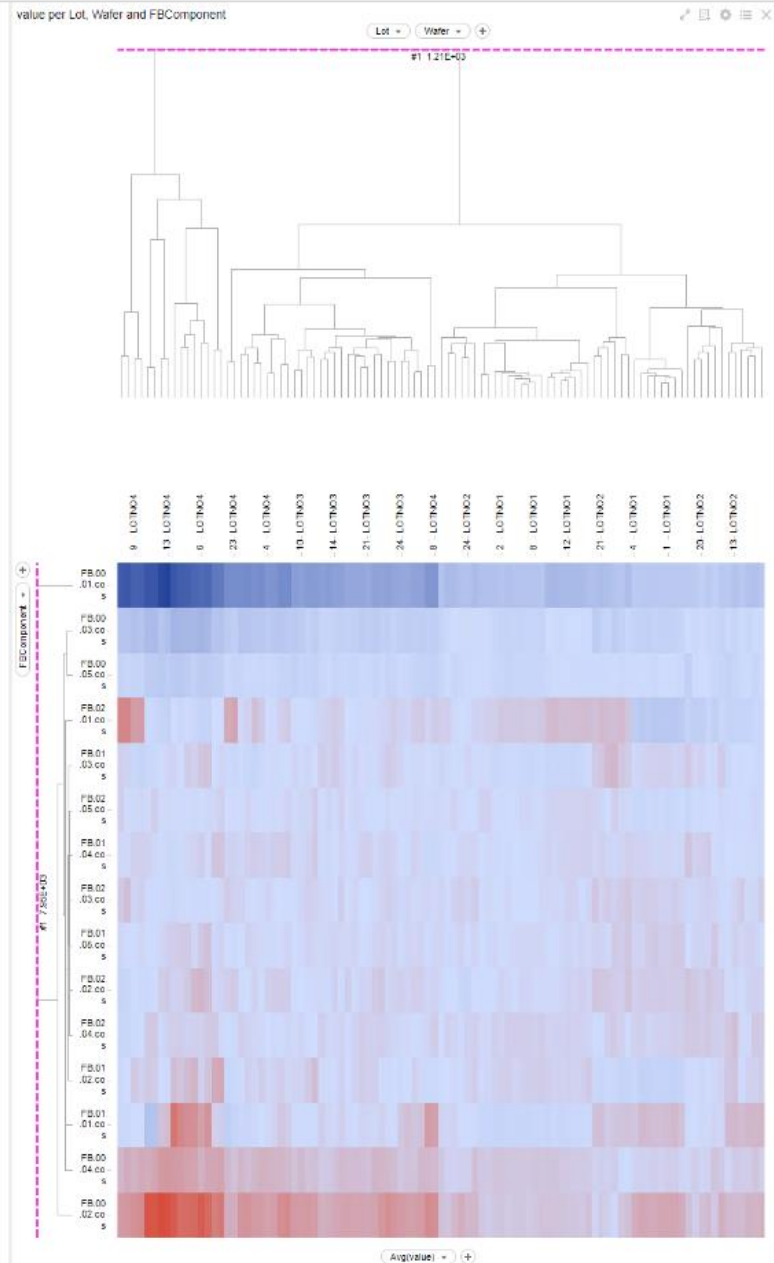
Find and explore anomalous events

Monitor anomalies

Predict when and why they occur

TIBC○®

# Anomaly Detection and Analysis

Find and explore anomalous events

Monitor anomalies

Predict when and why they occur

TIBC⦿®

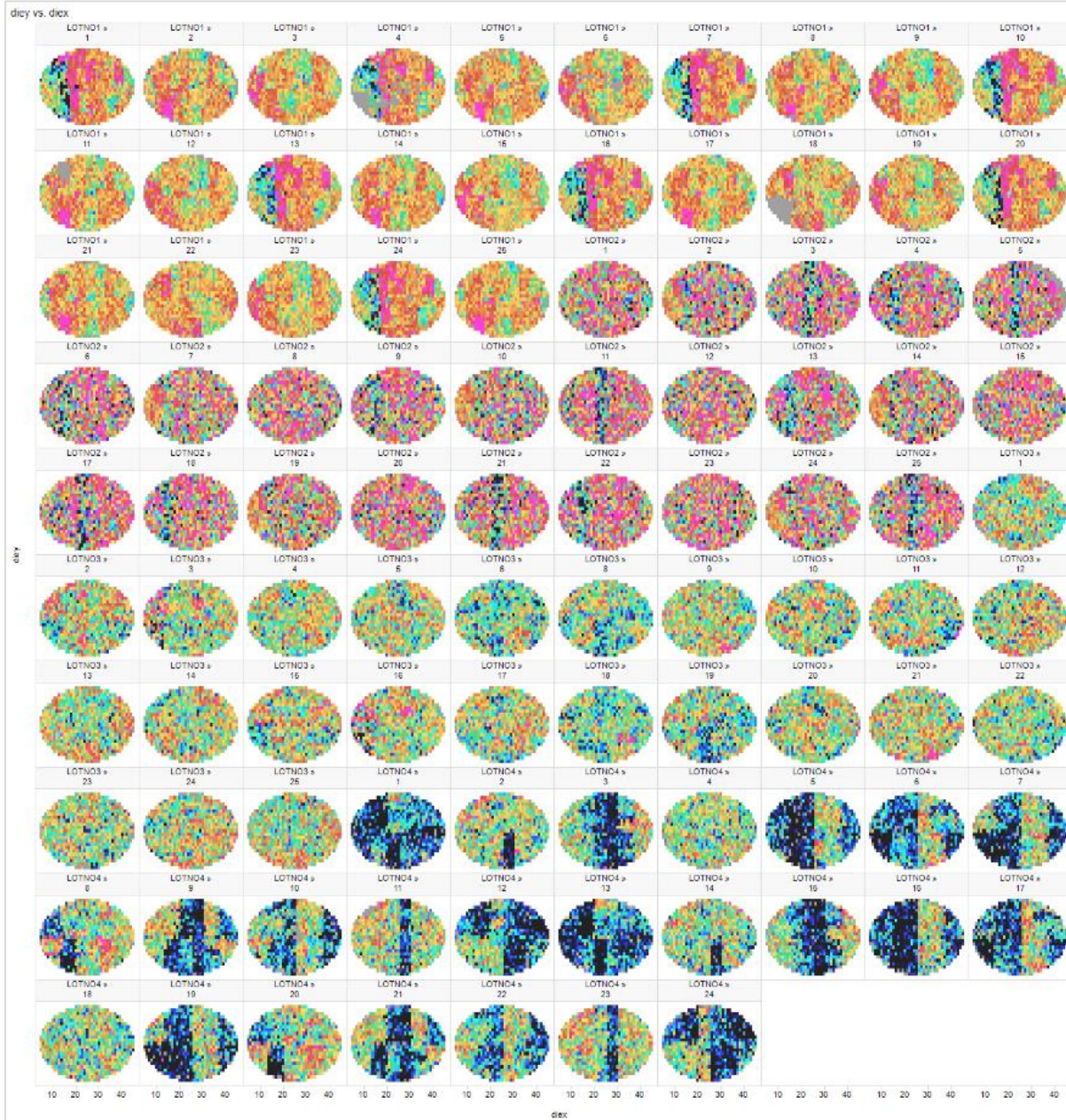# Anomaly Detection and Analysis

Find and explore anomalous events

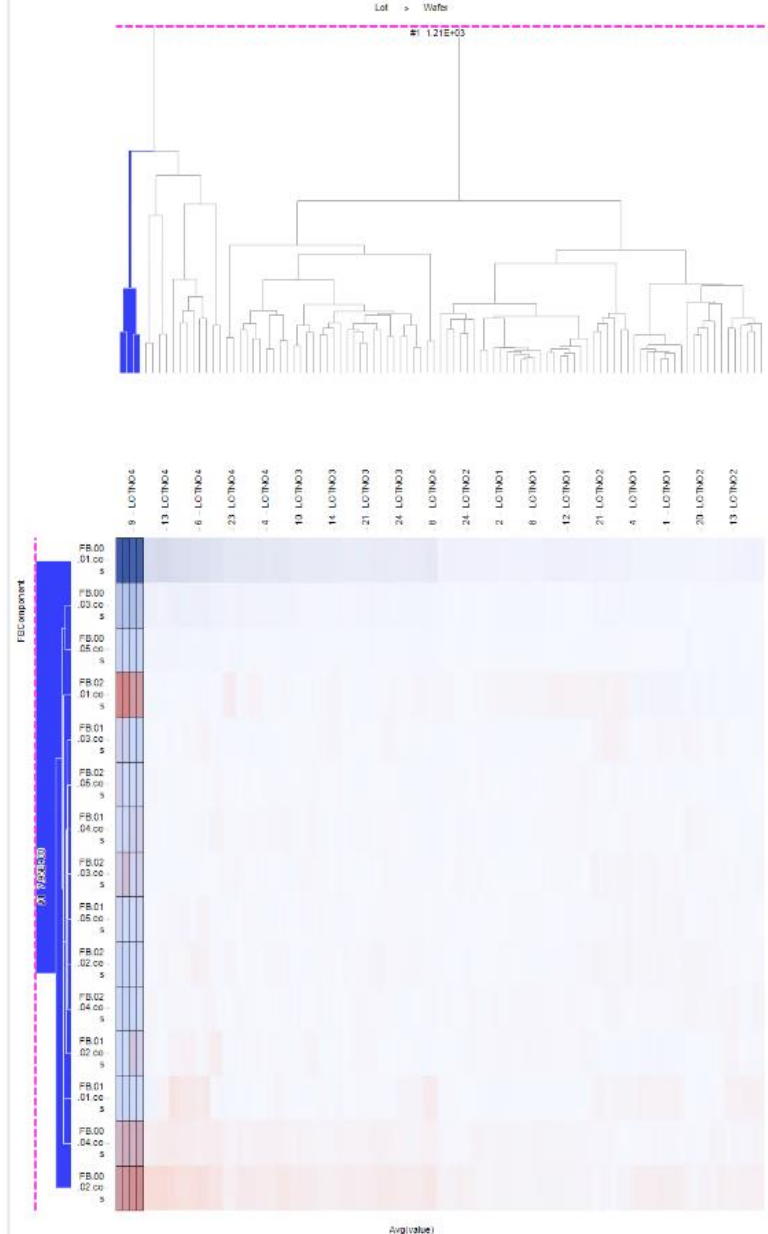Monitor anomalies

Predict when and why they occur

TIBC◯®

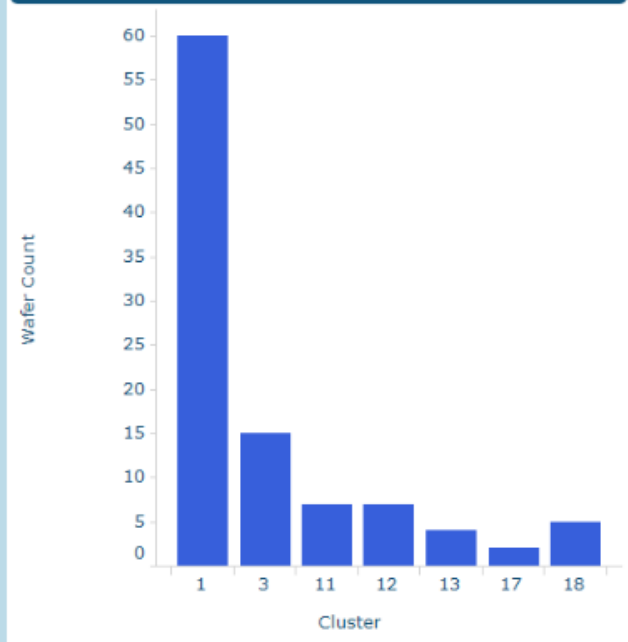# Digital Twin for Semiconductor Yield

Digital Twins for Semiconductor Manufacturing Yield: Wide-and-Big Data Analysis
Build Models to Relate Product Yield Failure Modes ( $Y_i$ ) with Process Parameters ( $P_i$ )



*~ 1K Steps, > 1M Process Parameters*

Start → Process Step 1 ($P_1$) · · · Process Step i ($P_i$) → Measurement Step 1 ($P_{i+1}$) → Process Step i+2 ($P_{i+2}$) · · · Process Step n ($P_n$)

Equipment Sensor Data — Physical Measurement Data — Equipment Names

**Product Test: Pass / Fail, Failure Modes**

*Run continuously to capture new relationships as they emerge*

Process Optimization
Process Monitoring
Process Control

Digital Twin Yield Models

$$Y_i = f(P_1, P_2, P_3, \ldots P_n)$$

Yield, Yield Components or Clusters

$$Y = f(Y_1, Y_2, Y_3, \ldots Y_n)$$

TIBCO®

# The Extreme Challenge of Big & Wide Data

- **Not just big data** – many rows: lots, wafers, die, units

- **Also wide data** – many columns: > 1M process parameters
  - Sensor traces
    - Time series for every sensor on each machine in each run
  - Physical measurements
    - Film thickness, critical dimensions, layer-to-layer overlay, defect classes & counts
  - Equipment and process attributes
    - Machine and component IDs, process recipe info
  - Supplies
    - Chemical batch IDs, QA sample data

"Today [semiconductor] fabs collect more than 5 billion sensor data points each day. The challenge is to turn massive amounts of data into valuable information."

*—Ann Kellehere, VP of the Technology and Manufacturing Group, Intel*

TIBCO

# Solution Architecture

**Data Prep, Feature Engineering & Selection**

**Further Feature Selection & Model Building**

**Visualization of Results**



- In-database parallelized computing
- Leverages Hadoop, Apache Spark

- In-memory dedicated fast server

- Interactive in-memory visualization environment

TIBCO®

# Performance Benchmarks & Conclusions

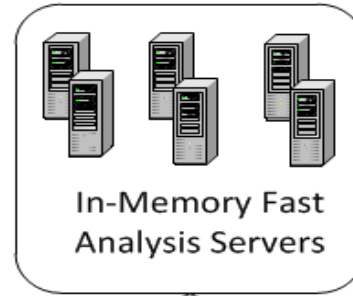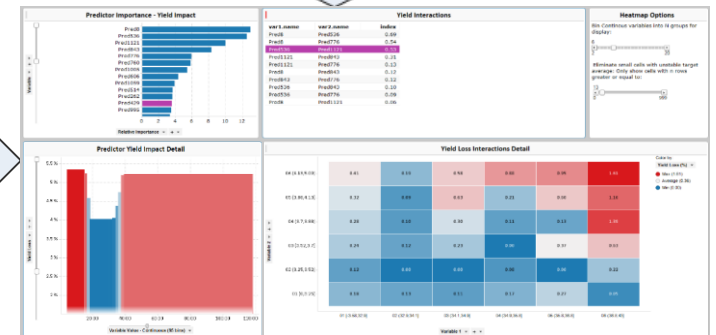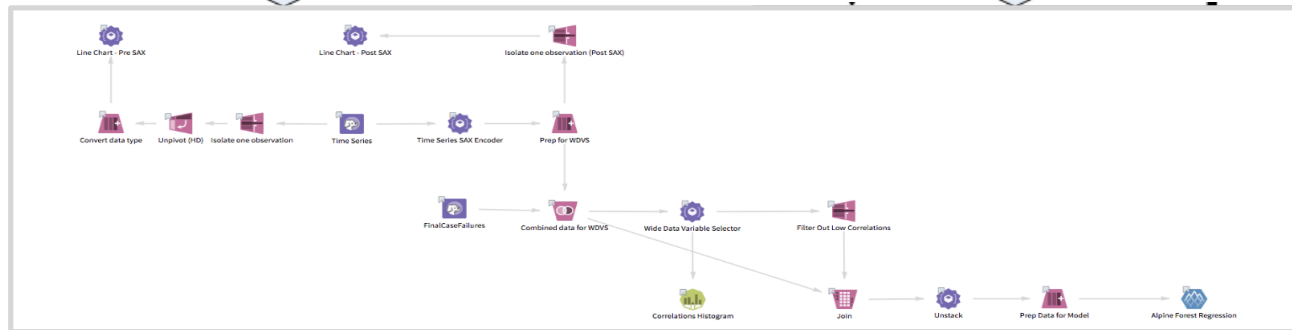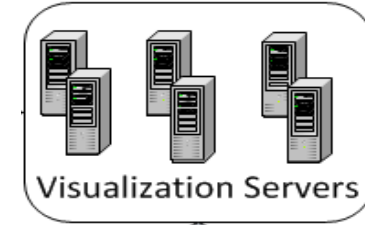- Demonstrated performance for time series data from 20,000 sensors, 10,000 wafers in **under 2 minutes**

- Current system scales to time series for 20,000 sensors, 100,000 wafers (**2.5 TB**) with results in **15 minutes**

  - More capacity and better performance can be achieved by adding nodes to the Spark cluster

- Working with top memory manufacturer to deploy production system

- System can provide automated real-time feedback on emerging equipment issues affecting yield

| Big Data Feature Selection Performance Benchmarks – Run Time[1] (minutes) | | | | | |
|---|---|---|---|---|---|
| | **20 Sensors** (1K variables) | **200 Sensors** (10K variables) | **2K Sensors** (100K variables) | **20K Sensors** (1M variables) | *Dataset Size* for *1M Variables* |
| **1K Wafers** | 0.47 | 0.48 | 0.72 | 1.0 | 25 GB |
| **10K Wafers** | 0.50 | 0.53 | 0.77 | 1.75 | 253 GB |
| **100K Wafers** | 0.53 | 0.67 | 1.25 | 15.15 | 2,530 GB |

**[1]Test Conditions:**
- Data stored in Hadoop data source
- 25 node Spark cluster – 16 cores, 32 GB for each node
- Each sensor time series compressed to 50 variables with SAX encoder prior to feature selection step
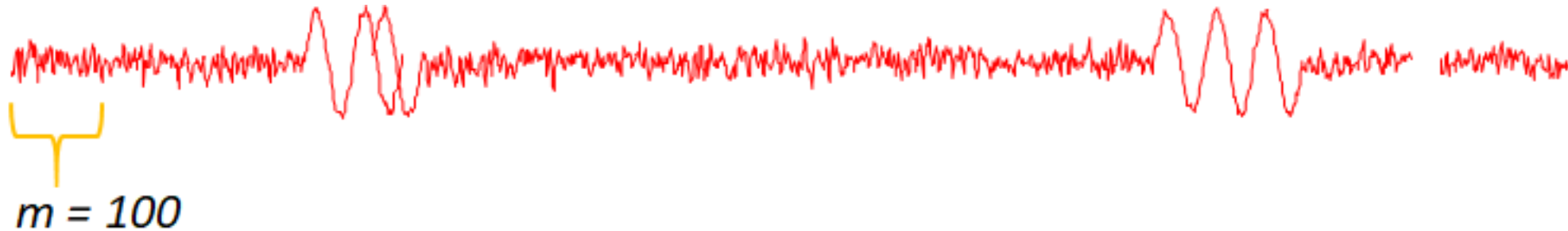
TIBCO®

# Longitudinal Anomaly Analysis

## *Subsequence Search*

A New Method for Identifying Anomalous Patterns in Time Series (Trace Analytics)

# Mueen's Algorithm for Similarity Search

**Mueen's Algorithm for Similarity Search (MASS)** is specialized for finding anomalous (versus typical) subsequences of time series



$m = 100$

*Extremely fast algorithm for this use case*

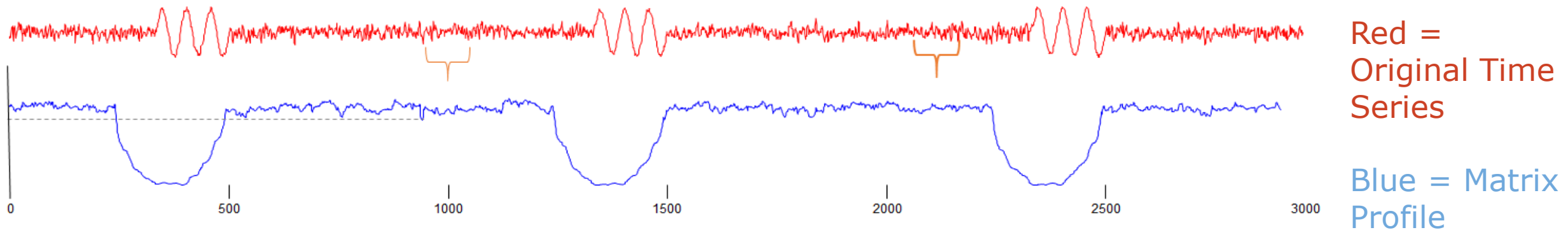Suitable for further acceleration using GPU

Material adapted from:
https://www.cs.ucr.edu/~eamonn/matrix_profile_i.pptx

TIBC○®

# Mueen's Algorithm for Similarity Search

Quickly create a matrix profile = a partial distance matrix

This uses a sliding window to define a series of subsequences

The Matrix Profile plots the distance of each subsequence to its nearest match, with the time sequence of the start of each subsequence on the x-axis
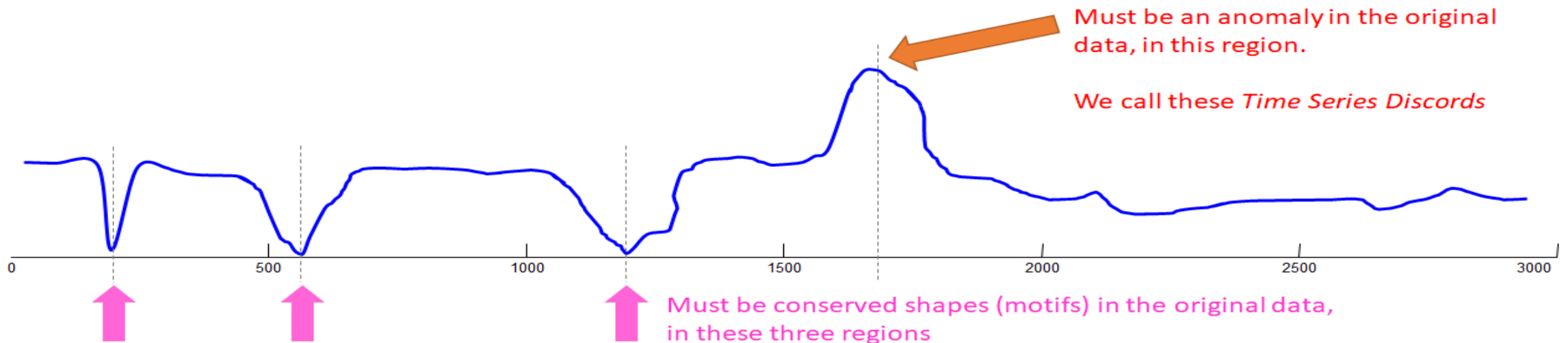


Red = Original Time Series

Blue = Matrix Profile

Material adapted from:
https://www.cs.ucr.edu/~eamonn/matrix_profile_i.pptx
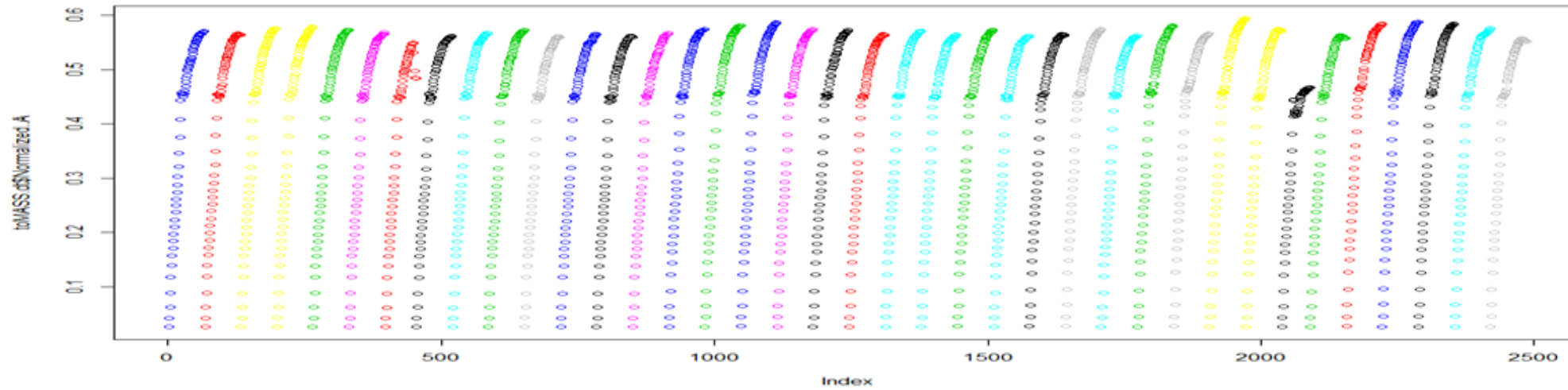
TIBCO®

# How to "Read" a Matrix Profile

Where you see relatively low values, you know that the subsequence in the original time series must have (at least one) relatively similar subsequence elsewhere in the data (such regions are "motifs" or reoccurring patterns)

Where you see relatively high values, you know that the subsequence in the original time series must be unique in its shape (such areas are "discords" or anomalies)



Must be an anomaly in the original data, in this region.

We call these *Time Series Discords*

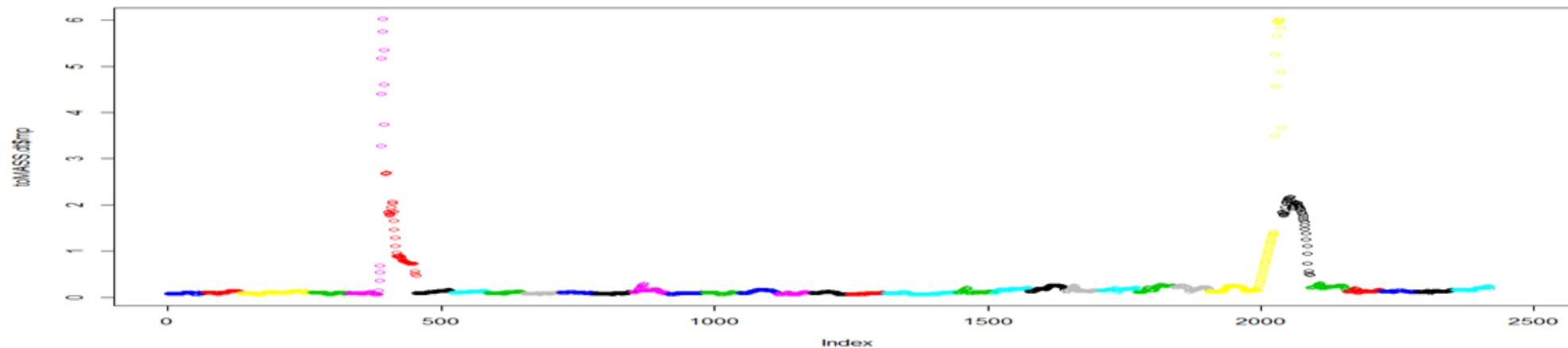Must be conserved shapes (motifs) in the original data, in these three regions

TIBCO®

# Manufacturing Batches

Raw Amperage - Each color delimits a batch



Matrix Profile highlights anomalies - set sliding window close to batch size

# Community

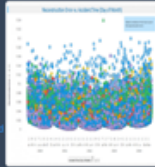## Anomaly Detection Template for TIBCO Spotfire®

This template detects anomalous data points in a dataset using an autoencoder algorithm.

## Data Function for TIBCO® Data Science - Team Studio in TIBCO Spotfire®

This data function enables users to execute a TIBCO® Data Science - Team Studio workflow from Spotfire.

## IoT Accelerator

Capture and analyze sensor data in real-time from your Internet of Things devices with TIBCO's IoT Accelerator. Integrate through industry-standard protocols like OSI PI, MQTT, and Web Services. Alternatively, implement custom adapters for your own protocols, all the way down to baseline serial port integration. Apply custom validations, cleansing policies, rules, and feature statistics on data feeds to identify trends and gain insight.

## Random Forest - Data Function for TIBCO Spotfire®

**Random forests** are an ensemble decision tree machine learning method for classification and regression.

## Statistical Process Control Template for TIBCO Spotfire® using TIBCO® Data Science - Statistica

This TIBCO Spotfire template is designed to enable user to build wide range of quality control charts with possibility do define charts specifications interactively according to user's needs. This comprehensive template is constructed based on Statistica data function and utilizing wide range of parameter settings already implemented in TIBCO Data Science - Statistica. It is an example of no-code data function.

## Risk Management Accelerator

Identify potentially risky activities in a high-frequency event stream using machine learning in TIBCO's Risk Management Accelerator. Build supervised and/or unsupervised models and hot deploy these to the streaming event processing platform, then score events in real-time. Raise alerts when potentially risky behaviour is detected.

## Gradient Boosting Machine Analysis Template for TIBCO Spotfire®

This template is used to create a GBM machine learning model to understand the effects of predictor variables on a single response.
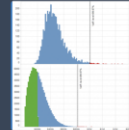
## Loss Distribution Approach to Operational Risk - Analysis Template for TIBCO Spotfire®

This analysis implements simple frequency-severity models for Operational Risk event types. This forms the basis of the Loss Distribution Approach alternative in the Basel regulations.

## Dynamic Pricing Accelerator

Take control of your pricing platform with TIBCO's Dynamic Pricing Accelerator. Applicable to insurance, retail, travel, or any industry where personalized pricing would be an advantage. Transform into an algorithmic business by deploying personalized pricing and propensity models that you build and manage to gain advantage over competitors while using industry-standard modelling languages. Hot deploy these models and watch the results in real-time with the TIBCO Insight Platform.

**TIBCO® Spotfire®**     **TIBCO® Data Science**     **TIBCO® Streaming**

**TIBCO®**

# AI in Operations
*Cloud Starters, Accelerators, Analytic Apps*

*Thoughtleader-Led Solutions*



**AI on Demand**

Data Science in Operations



Driving Customer Engagement

Click Here for Demo



Fraud and Risk Managment

Click Here for Demo



Insurance Dynamic Pricing

Click Here for Demo



Digital Twins in Manufacturing

Click Here for Demo



The Industrial Internet - Production Surveillance

Click Here for Demo

TIBCO®

# Questions & Contact

Steven Hillion

Sr. Director, Data Science

shillion@tibco.com

@StevenHillion

Michael O'Connell

Chief Analytics Officer

moconnell@tibco.com

@MichOConnell

## TIBCO Community

community.tibco.com

## TIBCO Exchange

community.tibco.com/exchange

## TIBCO Tech Blog

community.tibco.com/blog

*Acknowledgements*

Dr. Tom Hill, Mike Alperin, Nico Rode, David Katz, Jagrata Minardi, Siva Ramalingam

TIBCO®

!

Please complete the session survey in the mobile app.

aws