# Alternative Model Forms for Multi-set, Multi-level and Multi-block Data

EIGENVECTOR RESEARCH INCORPORATED

# *Outline*

- Definitions
- Multi-level data
  - DOE, crossed and nested designs
- ASCA
  - ANOVA simultaneous component analysis
  - Example
- MLSCA
  - Multi-level simultaneous component analysis.
  - Example
- Multi-block Data
  - Levels of data fusion
  - Examples

**EIGENVECTOR**
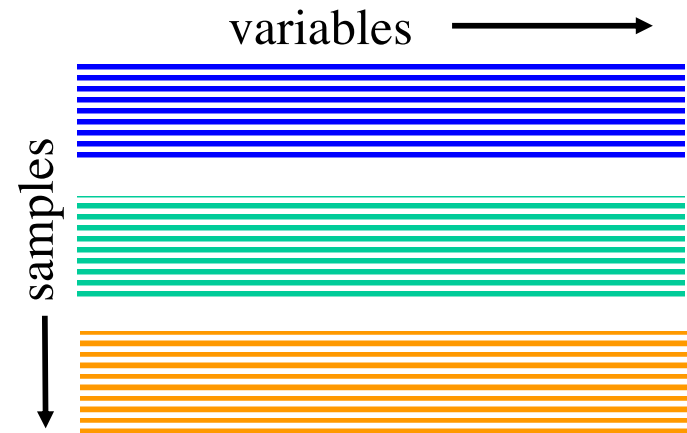RESEARCH INCORPORATED

# Definitions

- **Single-block**: data that is logically contained in a single matrix

- **Two-block**: two single block data sets that share a common mode (typically the sample mode)

- **Multi-block**: multiple single blocks that share a common mode

- **Multi-set**: groups of related samples that have the same variables, typically from designed experiments

- **Multi-level**: same as multi-set except typically from nested or happenstance designs

**EIGENVECTOR**
RESEARCH INCORPORATED

# *Definitions (cont.)*

- Multi-way: Data that is logically arranged in 3-way (or more) arrays
- Data fusion: the process of combining multiple sources of data to improve accuracy

**EIGENVECTOR RESEARCH INCORPORATED**

# *Multi-set Data*

variables ⟶

- Groups (sets) of related samples which have the same variables.

samples ↓

- Differences between groups may hide variability inherent to all samples.
- For samples grouped according to a DoE can separate variability
  - Due to each factor
  - Remaining systematic variability
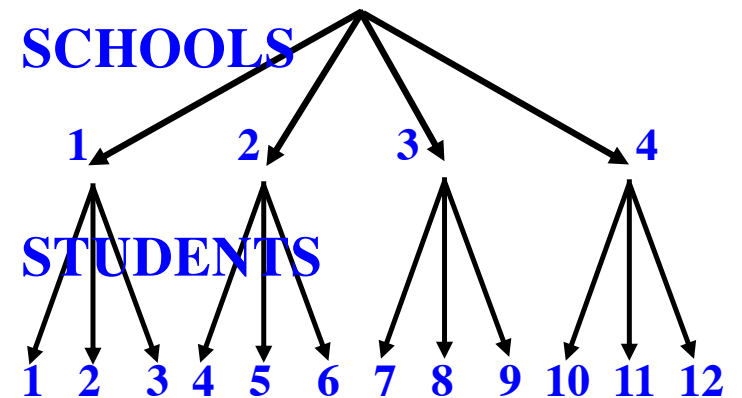- This is the purpose of ASCA and MLSCA

**EIGENVECTOR**
**RESEARCH INCORPORATED**

# Crossed and nested designs

- Crossed (factorial) designs: One or more factors with samples measured for every combination of factor levels.

| | | Treatment | | | |
|---|---|---|---|---|---|
| | | A | B | C | D |
| Dose | 1.1 | | | | |
| | 2.0 | | | | |
| | 3.5 | | | | |

- Nested designs: samples belong to groups which are organized hierarchically.



SCHOOLS
1    2    3    4

STUDENTS
1  2  3  4  5  6  7  8  9  10  11  12

These are both 2-factor designs

6

# *ASCA*
## *ANOVA Simultaneous Component Analysis*

For multivariate datasets based on crossed experimental designs, ASCA applies ANOVA decomposition and dimension reduction (PCA) to :

- Separate the variability associated with each factor.
- Estimate contribution of each factor to total variance.
- Test main factor and interaction effects for significance.
- View scores and loadings for these effects.

Especially useful for high-dimension datasets where traditional ANOVA is not possible.

**EIGENVECTOR**
**RESEARCH INCORPORATED**

# ASCA Method

- X data matrix, with 2 factors A and B.
- Decompose into DOE components

$$\mathbf{X} = \mathbf{X}_{avg} + \mathbf{X}_A + \mathbf{X}_B + \mathbf{X}_{AB} + \mathbf{E}$$

- Build PCA model for each main effect and interaction

$$\mathbf{X} = \mathbf{X}_{avg} + \mathbf{T}_A\mathbf{P}_A^T + \mathbf{T}_B\mathbf{P}_B^T + \mathbf{T}_{AB}\mathbf{P}_{AB}^T$$

- Calculate permutation P-value to estimate each factor's significance.
- Project residuals onto each PCA sub-model.

EIGENVECTOR RESEARCH INCORPORATED

# ASCA Demo data: asca_data

X: Measured glucosinolate levels in cabbage plants,

3 treatments, Control, Root, Shoot.

4 time points, Days 1, 3, 7, and 14.

5 replicates for each time-treatment.

11 measured concentrations.

X: (60, 11)

F: (60, 2) design matrix.

See X.description for details.

| | | Time (Day) | | | |
|---|---|---|---|---|---|
| | | 1 | 3 | 7 | 14 |
| Treatment | C | | | | |
| | R | 5 replicates each | | | |
| | S | | | | |

**EIGENVECTOR**
RESEARCH INCORPORATED

# ASCA Model

# *Time Model Scores and Loadings*
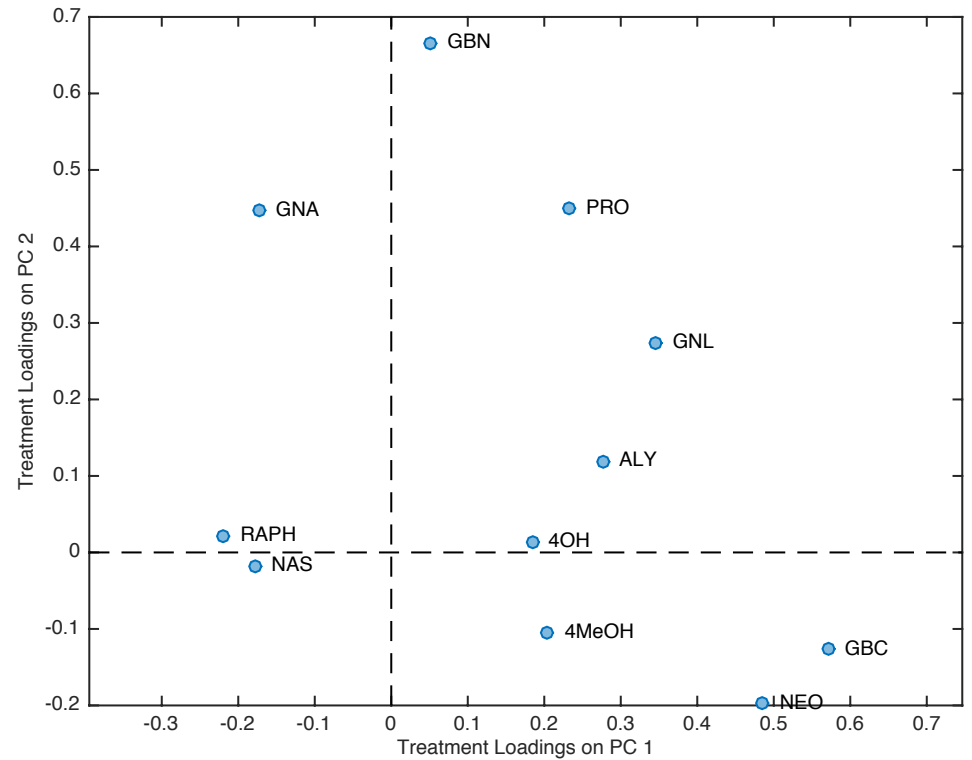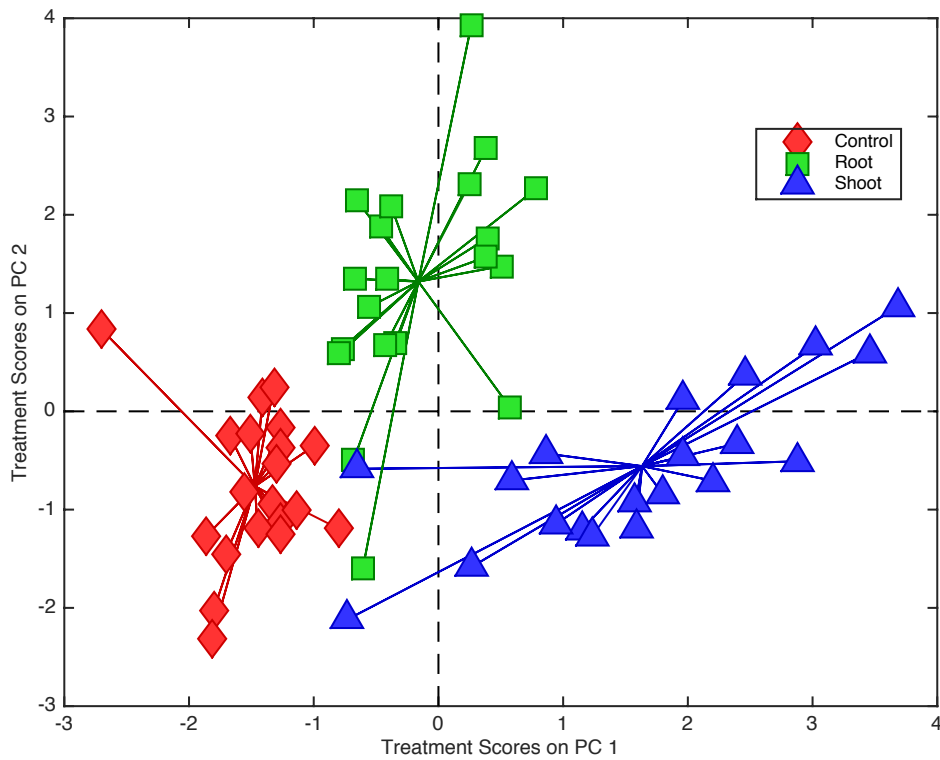
EIGENVECTOR RESEARCH INCORPORATED

# ASCA Scores Plot
## "Time" factor sub-model, PC 1



PC 1 of Time dependency common to all Treatments.

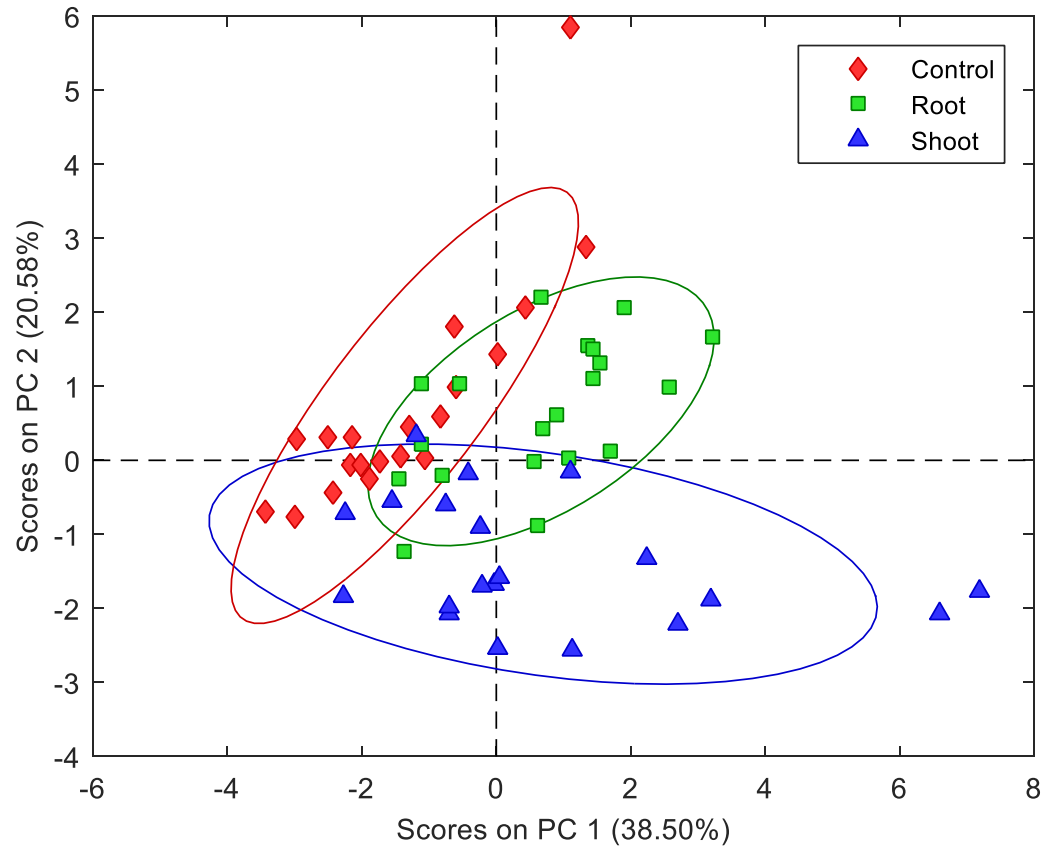Class = Treatment. Connect Classes = Mean at each X

EIGENVECTOR RESEARCH INCORPORATED

# *Treatment Model Scores and Loadings*

# ASCA *Treatment Scores Plot*



Separating out the Time and Time x Treatment effects highlights the Treatment effect

EIGENVECTOR
RESEARCH INCORPORATED

# PCA Scores Plot



…better than is seen by simply applying PCA to the data.

# *ASCA Conclusions*

- ASCA allows the variation associated with each factor to be resolved, and to see the main variables involved.

- For a perturbed biological system
    - Time factor scores reveal the common response independent of Treatment
    - Treatment factor scores show the Treatment effect independent of Time
    - Time x Treatment interaction scores show the additional time dependency at each Treatment level.

**EIGENVECTOR**
RESEARCH INCORPORATED

# *ASCA Conclusions, cont.*

- The % contribution of each factor or interaction to the total SSQ shows which effects are important.
- Perturbation P-values for each factor estimates the probability that there is no difference between the factor level averages for this effect.

EIGENVECTOR
RESEARCH INCORPORATED

# MLSCA
## Multi-level Simultaneous Component Analysis

MLSCA is a special case of ASCA applied to data from designed experiments with nested factors.

- Separates variability associated with each factor and residual.

- Estimate contribution of each factor to total sum of squares.

- View scores and loadings for these effects.

- Also builds PCA model on the residuals, or "within" variability. "Within" is often the focus of the analysis.

- Note that "Class Center" pre-processing can achieve same result if there is a single nesting factor.

EIGENVECTOR
RESEARCH INCORPORATED

# *MLSCA Method*

- X data matrix, with 2 nested factors A and B.

- Decompose into DOE components

$$\mathbf{X} = \mathbf{X}_{avg} + \mathbf{X}_A + \mathbf{X}_{B(A)} + \mathbf{E}$$

$\mathbf{X}_A$ contains factor A level averages

$\mathbf{X}_{B(A)}$ contains factor B level averages for each level A

$\mathbf{E}$ are the residuals, "within" component

- Build PCA model for each effect and residual

$$\mathbf{X} = \mathbf{X}_{avg} + \mathbf{T}_A\mathbf{P}_A^T + \mathbf{T}_{B(A)}\mathbf{P}_{B(A)}^T + \mathbf{T}_E\mathbf{P}_E^T$$

*constant   between A     between B         within*

**EIGENVECTOR**
**RESEARCH INCORPORATED**

# *MLSCA: simple example*

MLSCA can be used to reveal systematic variability within grouped samples which can be obscured by inter-group differences.
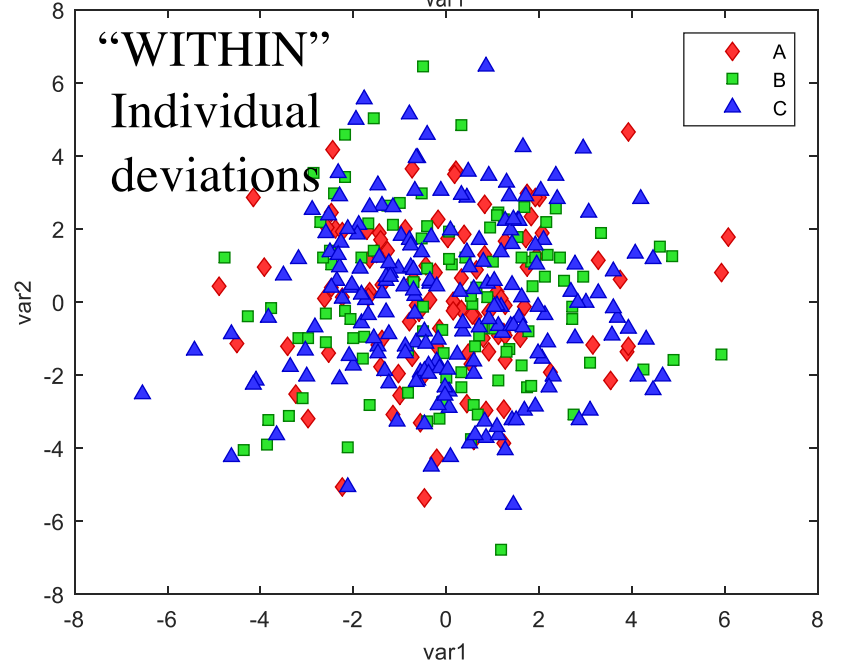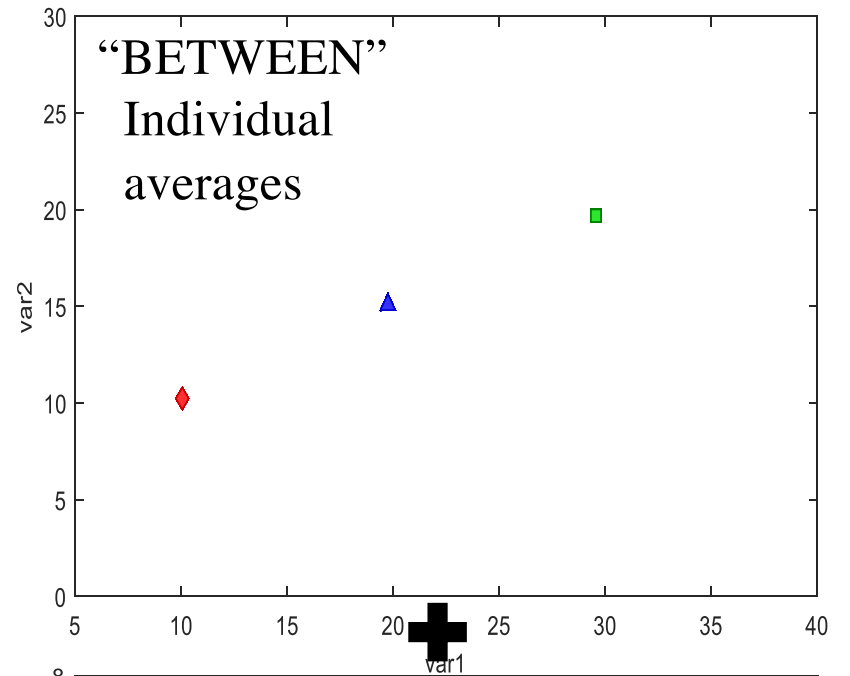
Example:   X: (400,2)
400 samples from 3 individuals, A, B, and C.

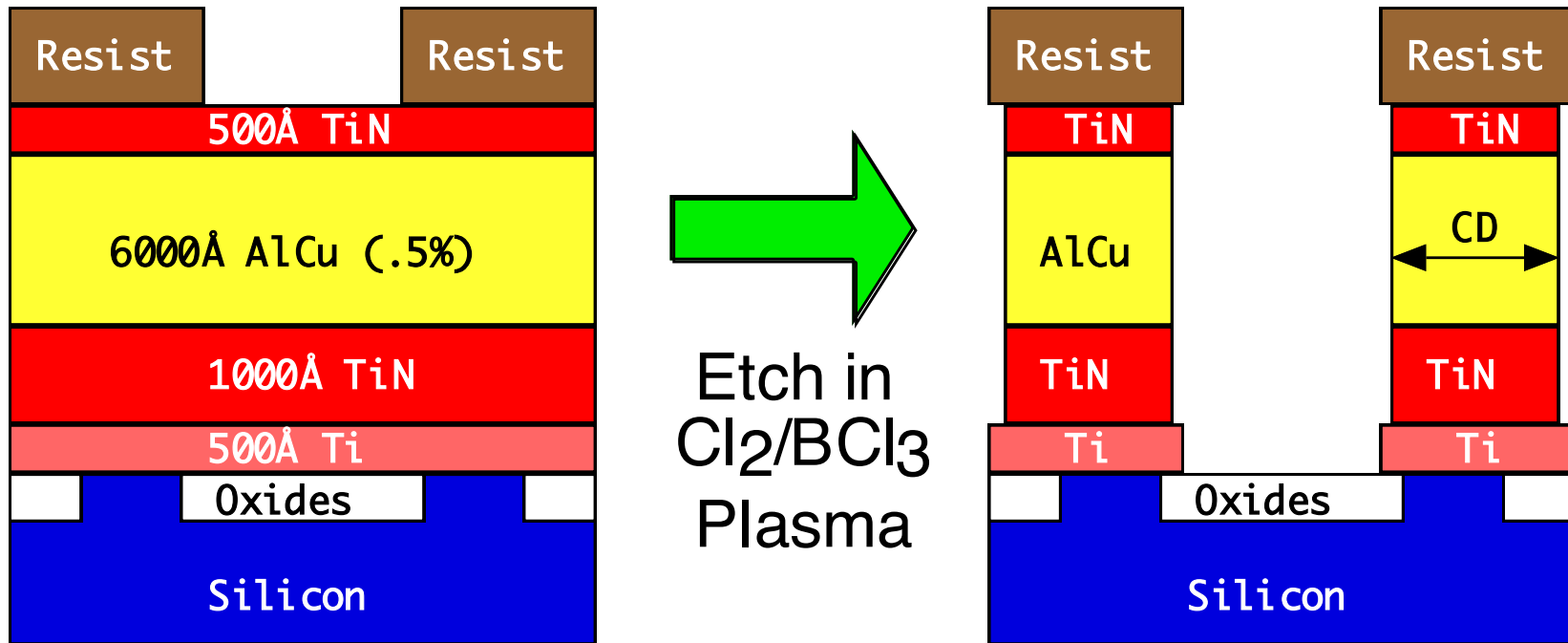Need to remove offsets for each individual to see the internal, "within" individual variation.

"BETWEEN"
Individual
averages

"WITHIN"
Individual
deviations

X = average for each individual
+ deviations from that

24

# *Example: Plasma Metal Etch*



- Linewidth (Critical Dimension) Control
  - Constant linewidth reduction run to run and across wafer
  - Constant linewidth reduction for every material in stack
- Minimal damage to oxide

EIGENVECTOR
RESEARCH INCORPORATED

# *Available Measurements*

- Machine State Data: Equipment has SECS-II Port
  - Provides traces with time stamp and step number
  - Regulatory controller setpoints & controlled variable measured values
    - gas flows, pressure, plasma powers
  - Regulatory controller manipulated variables
    - exhaust throttle valve, capacitors
    - mass flow controller do not provide valve position
  - Additional process measurements
    - broadband plasma emission (often used for endpoint)
    - impedance measurements
- Optical Emission Spectroscopy (OES)
- RF Data

EIGENVECTOR RESEARCH INCORPORATED

# Nested dataset "mlsca_data"

12 engineering variables from a LAM 9600 Metal Etcher over the course of etching 107 wafers.

- Three experiments were run at different times.
- Experiment have 34, 36 and 37 wafers each, for 107 unique wafers.
- 80 samples (replicates) measured for each wafer during etching.
- X is (8560, 12)

| | EXPERIMENT | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | | | | 2 | | | | 3 | | | | |
| WAFER | 1 | 2 | ··· | 34 | 35 | 36 | ··· | 70 | 71 | 72 | ··· | 107 | |
| *80* REPLI-CATES | x x x . . . . x | x x x . . . . x | | x x x . . . . x | x x x . . . . x | x x x . . . . x | | x x x . . . . x | x x x . . . . x | x x x . . . . x | | x x x . . . . x | |

Nested factors are not crossed.

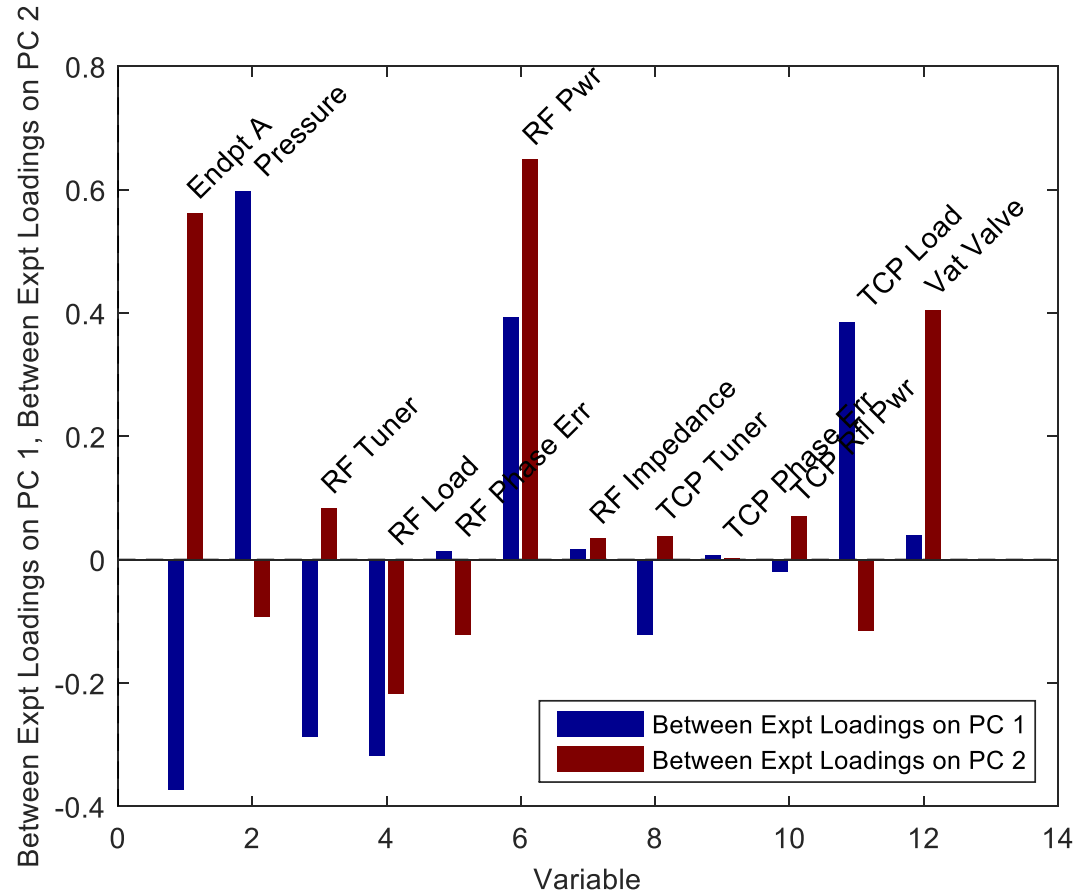**EIGENVECTOR RESEARCH INCORPORATED**

# *MLSCA Model*

# MLSCA Scores Plot
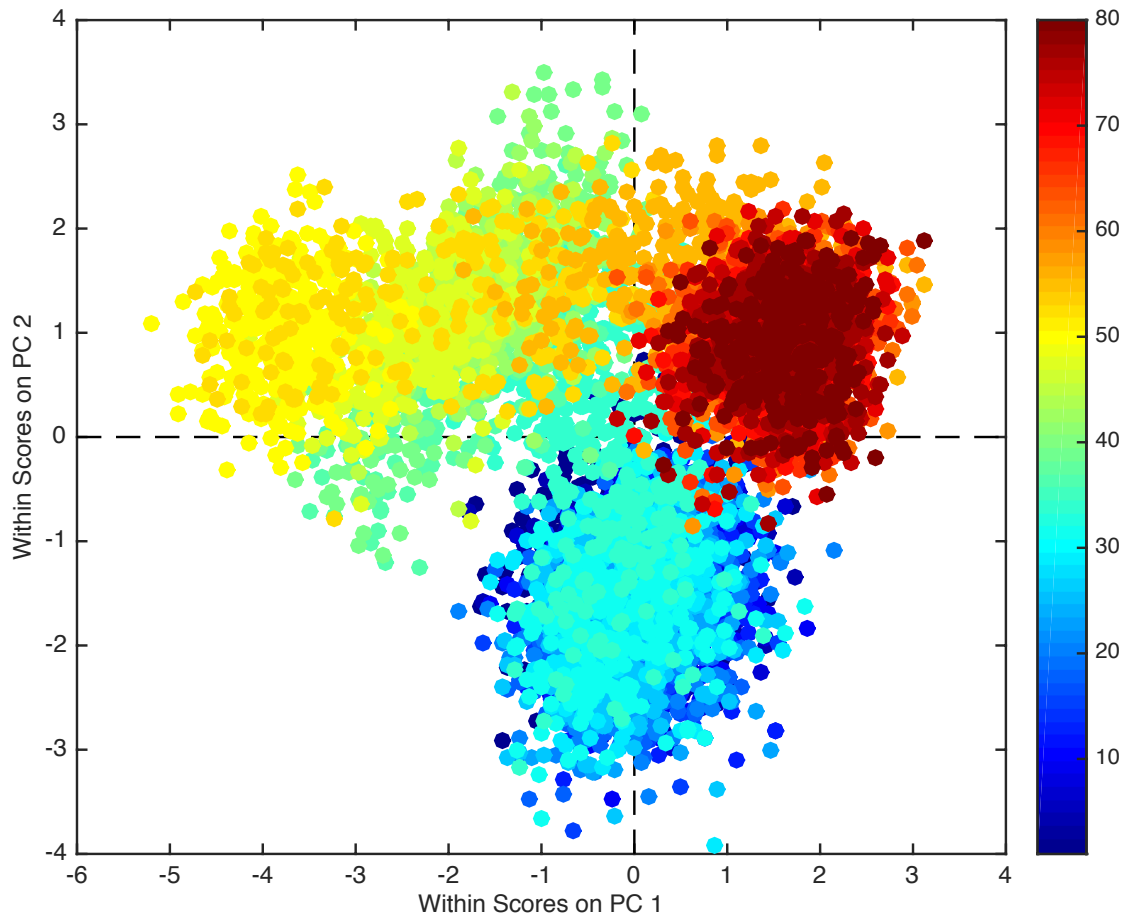## "Experiment" factor sub-model, PC 1 vs 2

# MLSCA Loadings Plot
## "Experiment" factor sub-model, PC 1 and 2

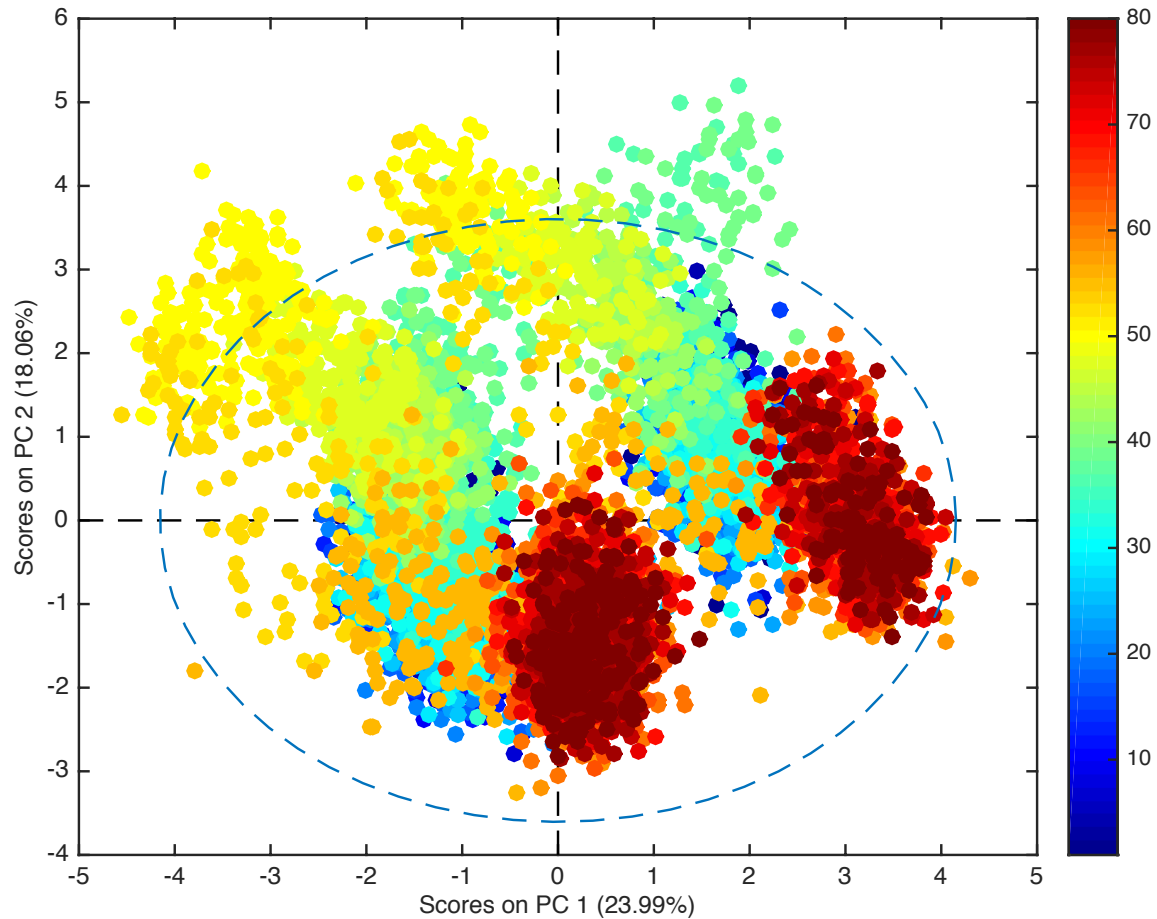# MLSCA Scores Plot
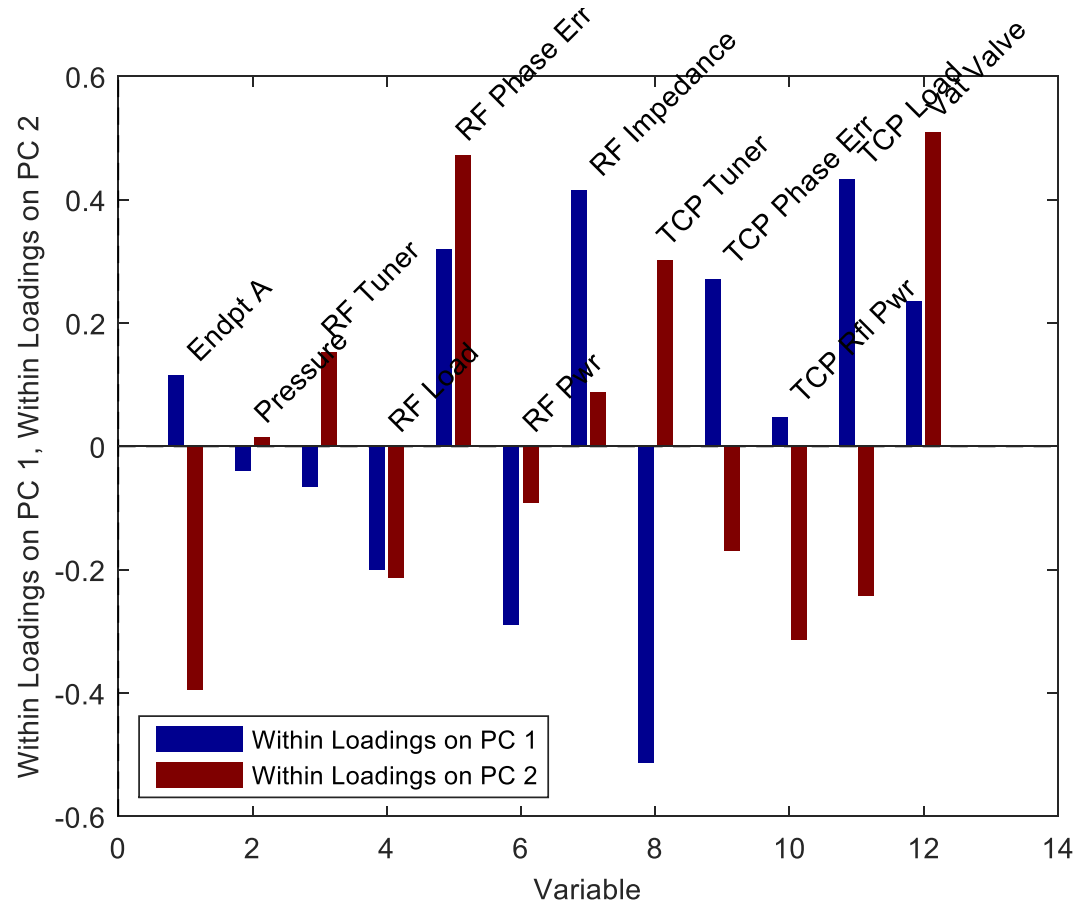## ''Within'' sub-model, PC 1 vs 2, colored by time

# Compare to PCA
## Convolves between and within factors

# MLSCA Loadings Plot
## ''Within'' Residual sub-model, PC 1 and 2

EIGENVECTOR RESEARCH INCORPORATED

# *MLSCA Conclusions*

MLSCA allows the variation associated with each nested factor to be resolved, and to see the main variables involved.

- Often used to reveal the inherent "within" group variability of samples after factor effects are removed. For process data this allows separation of within-run variation from between-run variation.

- SSQ contributions show which nested factors are important.

**EIGENVECTOR**
RESEARCH INCORPORATED

# *Multi-block Data Fusion*

- Data fusion can be done at three levels
  - Low level: single model of combined data blocks appropriately scaled/preprocessed
  - Mid level: combining scores from individual data blocks into a consensus model
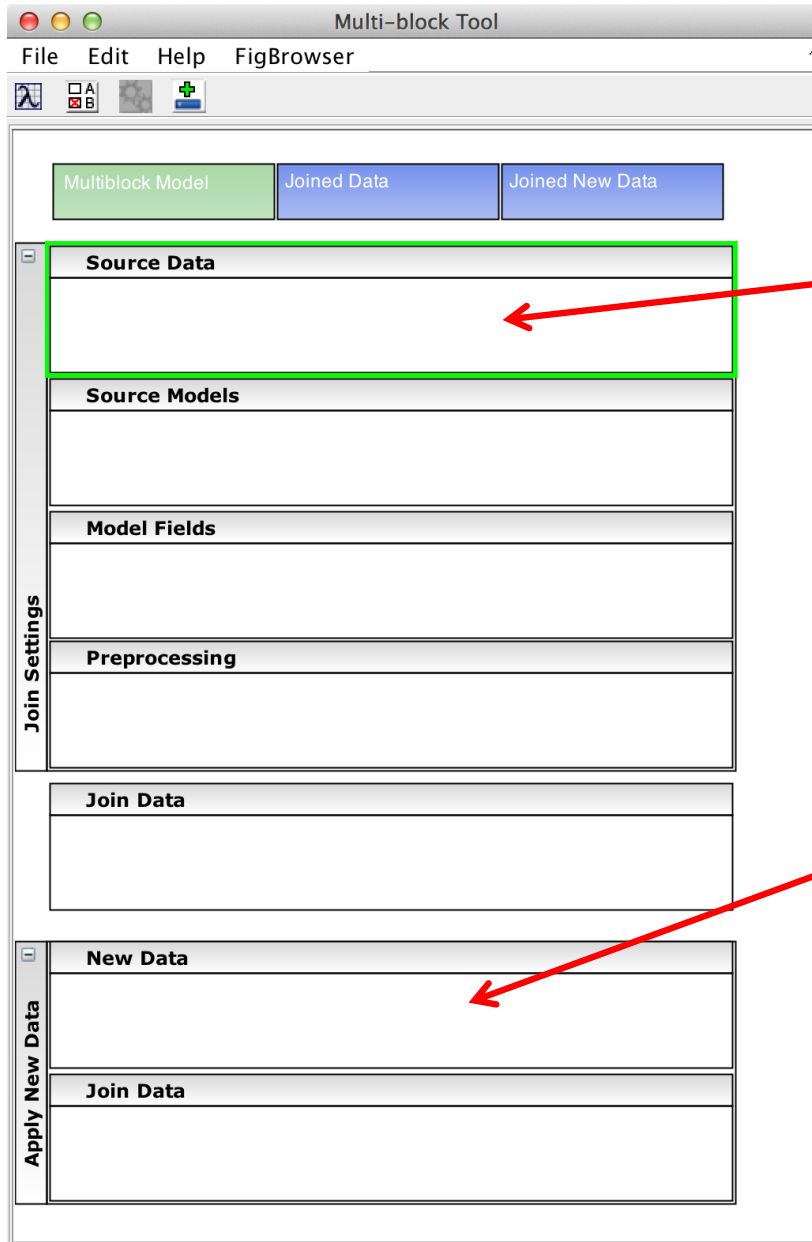  - High level: combining predictions from individual models in some sort of voting scheme

**EIGENVECTOR**
**RESEARCH INCORPORATED**

# *Sensitivity of MSPC Models*

- Three experiments performed with 21 "induced" faults on:
  - TCP top power
  - RF bottom power
  - Cl2 flow
  - BCl3 flow
  - Chamber pressure
  - Helium chuck pressure
- Data available for Machine State, RF and OES
- Goal: Compare ability of models considered for detecting faults: best case and for routine data
- Generated realistic faults to test models

**EIGENVECTOR**
**RESEARCH INCORPORATED**

# *Example with Etch Data*

- Available data: Machine, OES and RFM data for 104 normal wafers and 20 induced faults
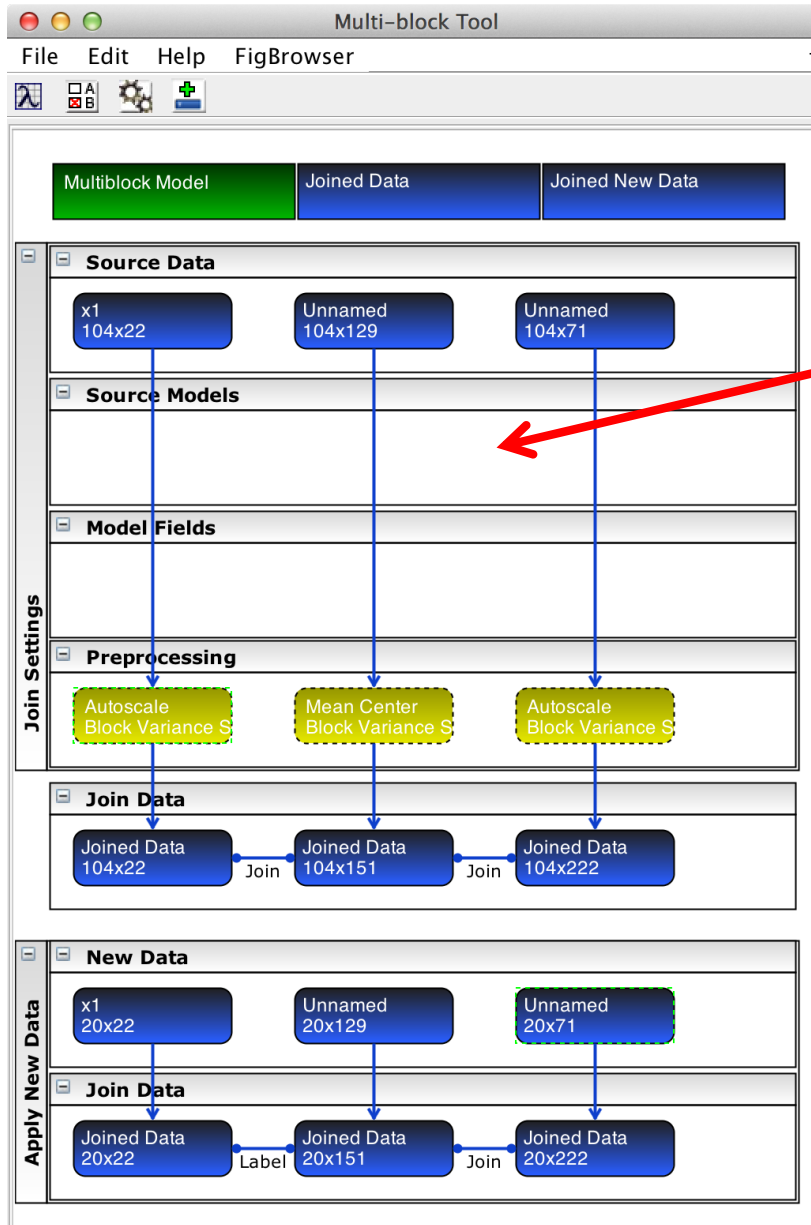
- Data reduced just to mean over each batch

**EIGENVECTOR**
**RESEARCH INCORPORATED**

# *Multi-block Tool Interface*
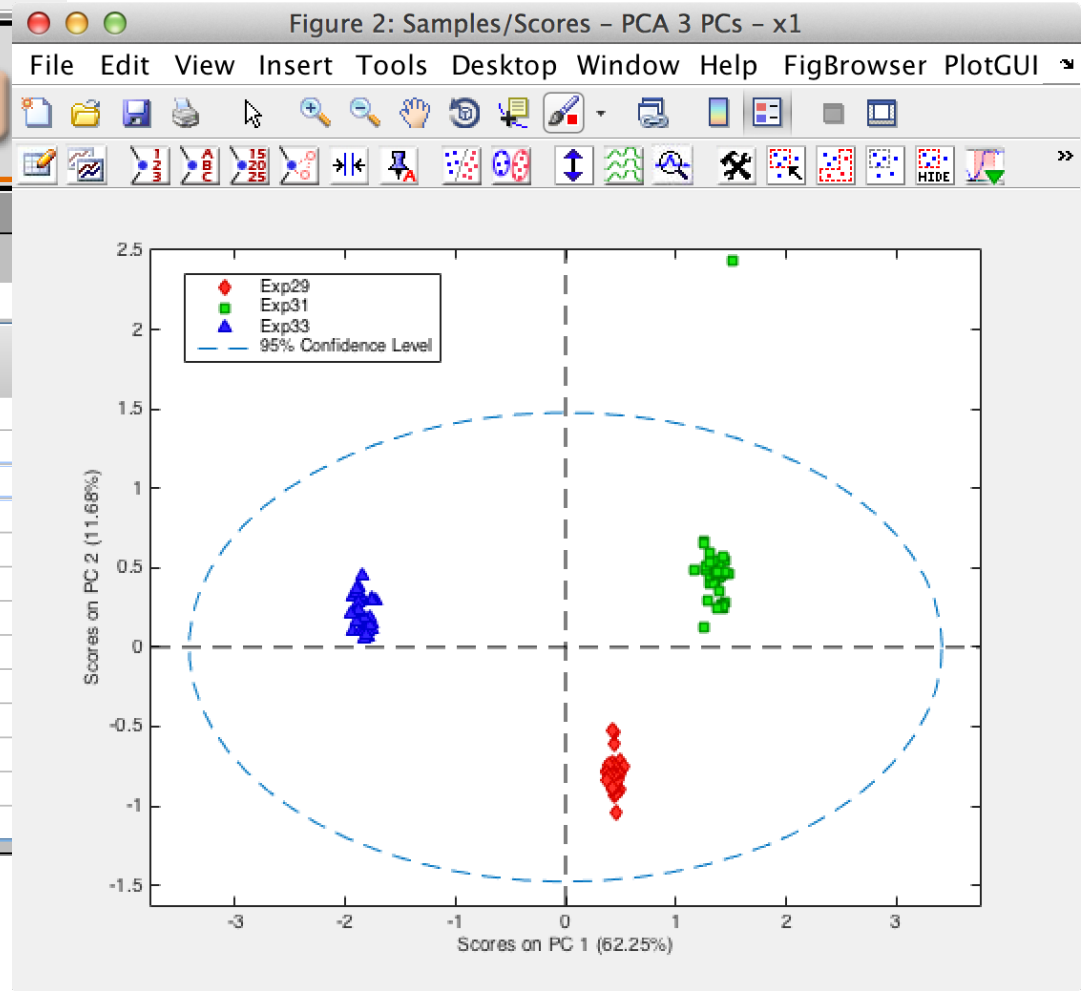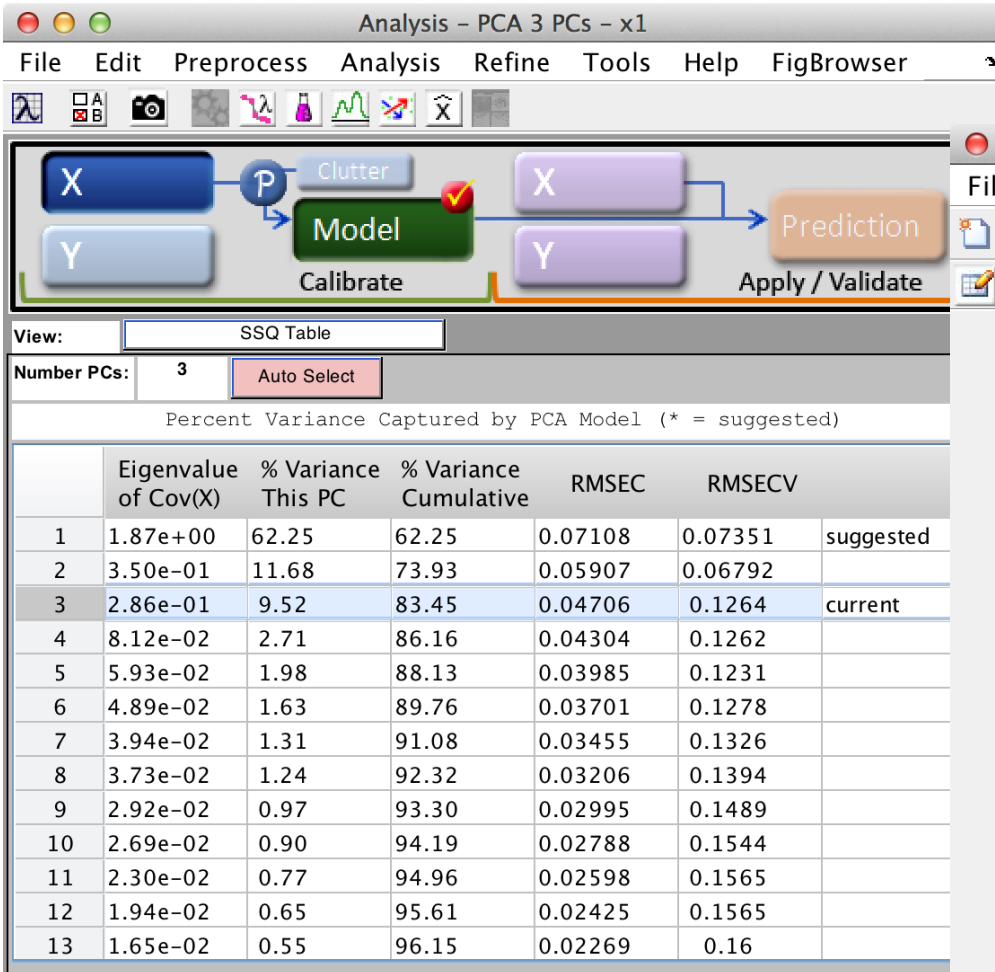


Drag calibration data sets here

Drag test data sets here

# *Separately Preprocessed Then Joined Data*



Or put models here for mid level fusion
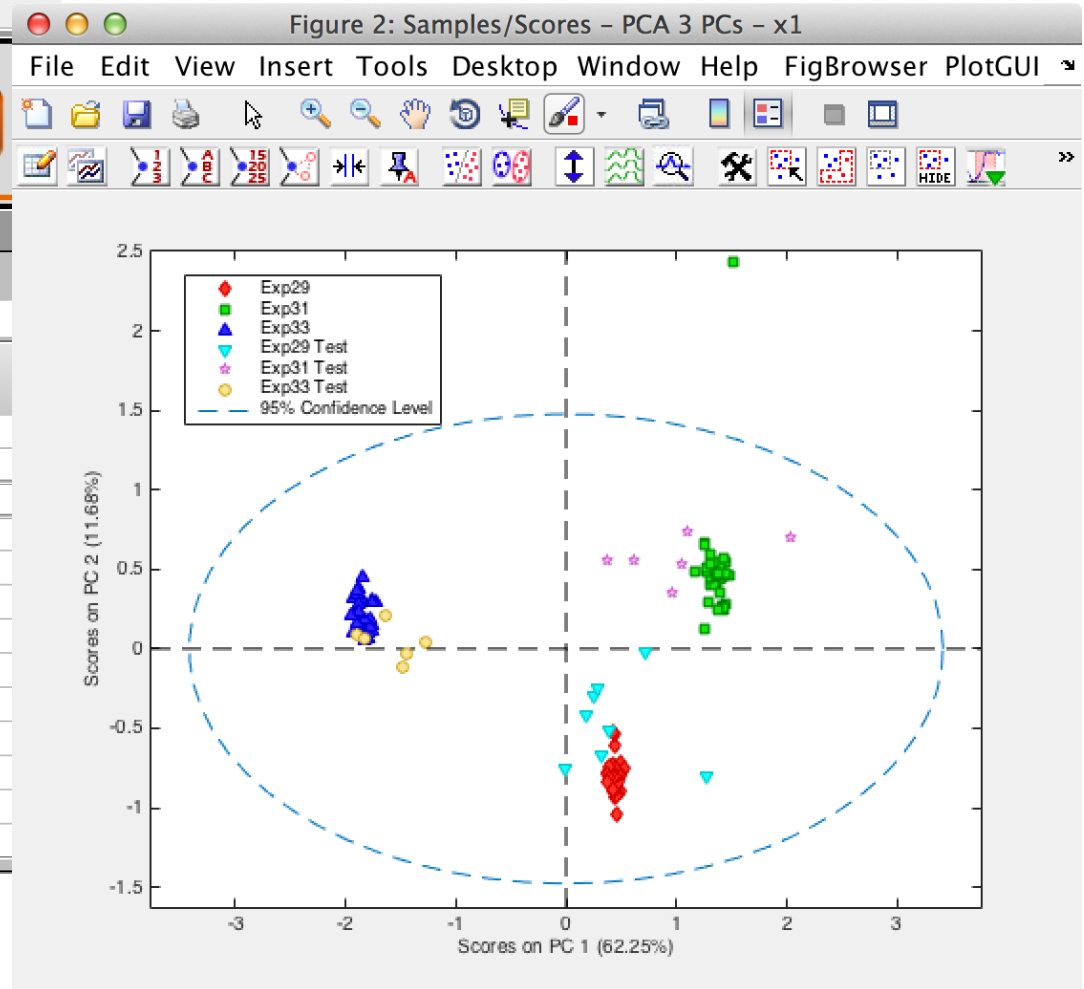
# *Data pushed into PCA*

# With Test Data Loaded

# *Redo at Mid-level*

- Develop individual PCA models of data blocks
- Load models into Multi-block tool
- Choose model outputs
- Join and push into PCA
- Results similar

**EIGENVECTOR**
RESEARCH INCORPORATED

# *Conclusions I*

- ASCA
  - for multi-set data typically from designed experiments
- MLASCA
  - for multi-level data typically from happenstance data (often semi-batch)
- ASCA and MLASCA allow new ways to partition and understand variance

**EIGENVECTOR**
RESEARCH INCORPORATED

# *Conclusions II*

- Data Fusion methods combine multi-block data that share a common mode

- Data Fusion can be done at three levels
  - Low Level: joining blocks after preprocessing
  - Mid Level: joining model outputs such as scores
  - High Level: Combine predictions from multiple models in some sort of voting scheme

- Often brings out aspects of data that aren't obvious in blocks analyzed separately

EIGENVECTOR
RESEARCH INCORPORATED

# *References*

## ASCA:

- Smilde, A.K., J.J. Jansen, H.C.J. Hoefsloot, R-J.A.N. Lamars, J. van der Greef, M.E. Timmerman, "ANOVA-simultaneous component analysis (ASCA): a new tool for analyzing designed metabolomics data", Bioinformatics, 2005, 21, 3043-3048.
- Zwanenburg, G., H.C.J. Hoefsloot, J.A. Westerhuis, J.J. Jansen, and A.K. Smilde, "ANOVA-principal component analysis and ANOVA-simultaneous component analysis: a comparison". J. Chemometrics, 2011.

## MLSCA:

- de Noord, O.E., and E.H. Theobald, Multilevel component analysis and multilevel PLS of chemical process data. J. Chemometrics 2005; 301–307
- Timmerman, M.E., Multilevel Component Analysis. Brit. J. Mathemat. Statist. Psychol. 2006, 59, 301-320.
- Jansen, J.J., H.C.J. Hoefsloot, J. van der Greef, M.E. Timmerman and A.K. Smilde, Multilevel component analysis of time-resolved metabolic fingerprinting data. Analytica Chimica Acta, 530, (2005), 173–183.

**EIGENVECTOR**
**RESEARCH INCORPORATED**