# New Information-Theory-Based Methods in the Analysis of Childhood Development Data

Nikita A. Sakhanenko,[1] David J. Galas,[1] Representing the Healthy Birth, Growth, and Development knowledge integration (HBGDki) Community[2]

[1]Pacific Northwest Diabetes Research Institute, Seattle, WA, USA; [2]Bill & Melinda Gates Foundation, Seattle, WA, USA
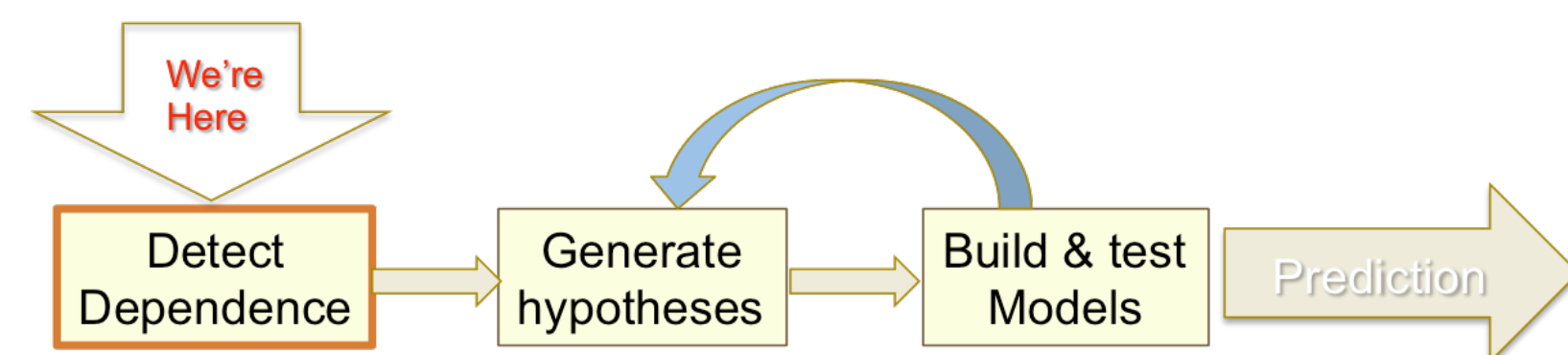
## Introduction

The complexity of infant growth and development processes and resultant data reflect the deep complexity of biological systems.

The Problem: Can we detect complex dependencies of biological variables?

Given data sets with many instances of many variables:

- Can we find a measure for a subset of variables that is significantly nonzero if and only if the variables are interdependent?
- Can we separate detection of the existence of dependence from any models of the nature of dependence?



We're Here → Detect Dependence → Generate hypotheses → Build & test Models → Prediction

### Dependency Measures
*Information about one variable is contained in information about another.*

Mutual information
$$I(X_i, X_j) = H_i + H_j - H_{ij}$$

"Interaction information"
$$I(X_i, X_j, X_k) = I(X_i, X_j) - I(X_i, X_j | X_k)$$
$$I(X_i, X_j, X_k) = H(X_i) + H(X_j) + H(X_k) - H(X_i, X_j) - H(X_i, X_k) - H(X_j, X_k) + H(X_i, X_j, X_k)$$

"Conditional mutual information, or differential interaction information"
$$-I(X_i, X_j | X_k) = I(X_i, X_j, X_k) - I(X_i, X_j)$$
$$\Delta_{ijk} = -I(X_i, X_j | X_k) = H_k - H_{ki} - H_{kj} + H_{ijk}$$

### General dependency measures:
based on interaction information
**Symmetric deltas**

If $\{\tau_m\}$ are all subsets of $v_m \subseteq V_m$ that contain $X_m$
$$\Delta(V_{m-1}; X_m) = I(V_m) - I(V_{m-1}) = \sum_{\{\tau_m\}} (-1)^{|\tau_n|+1} H(\tau_m)$$
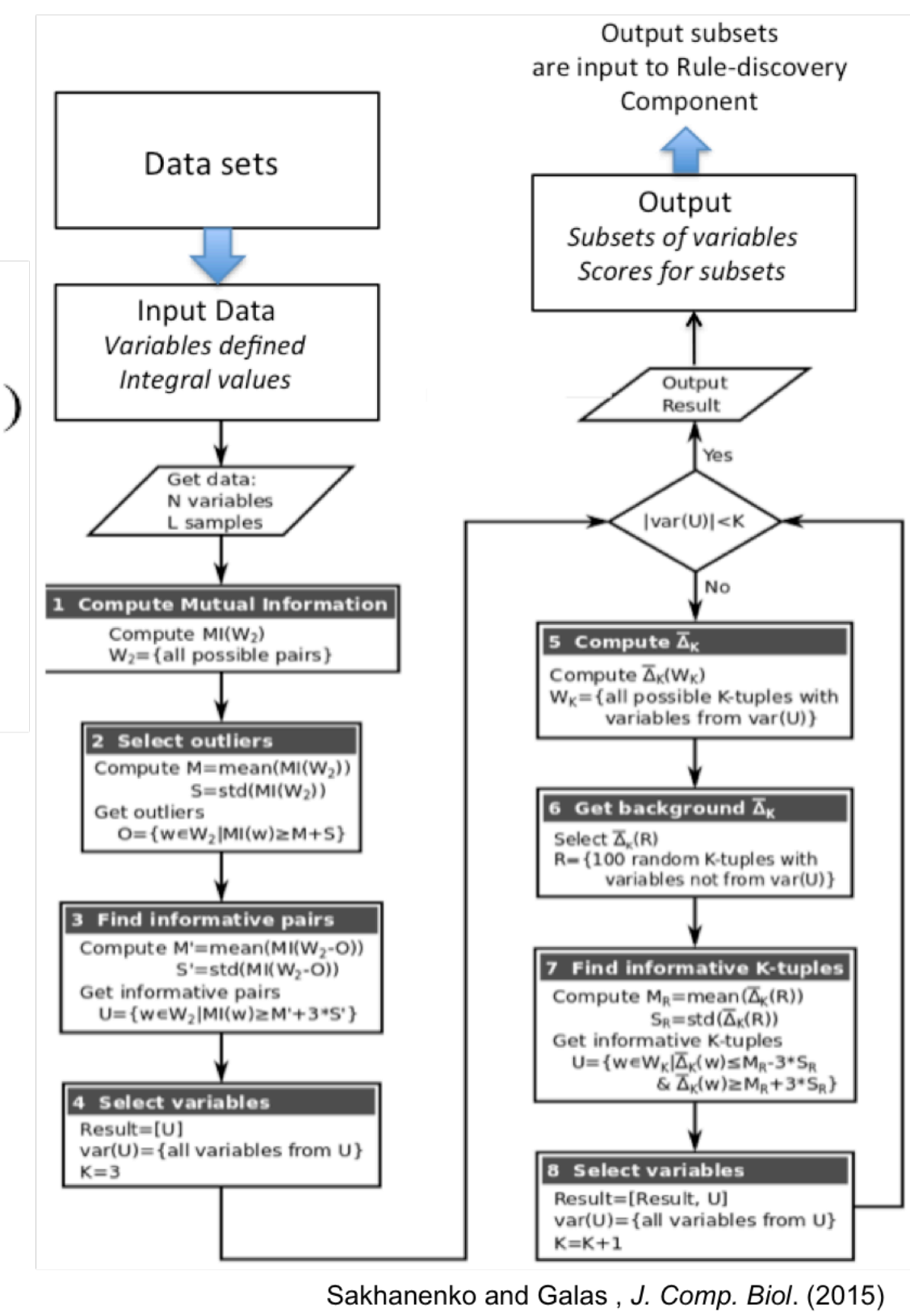
The general, symmetric measure is then
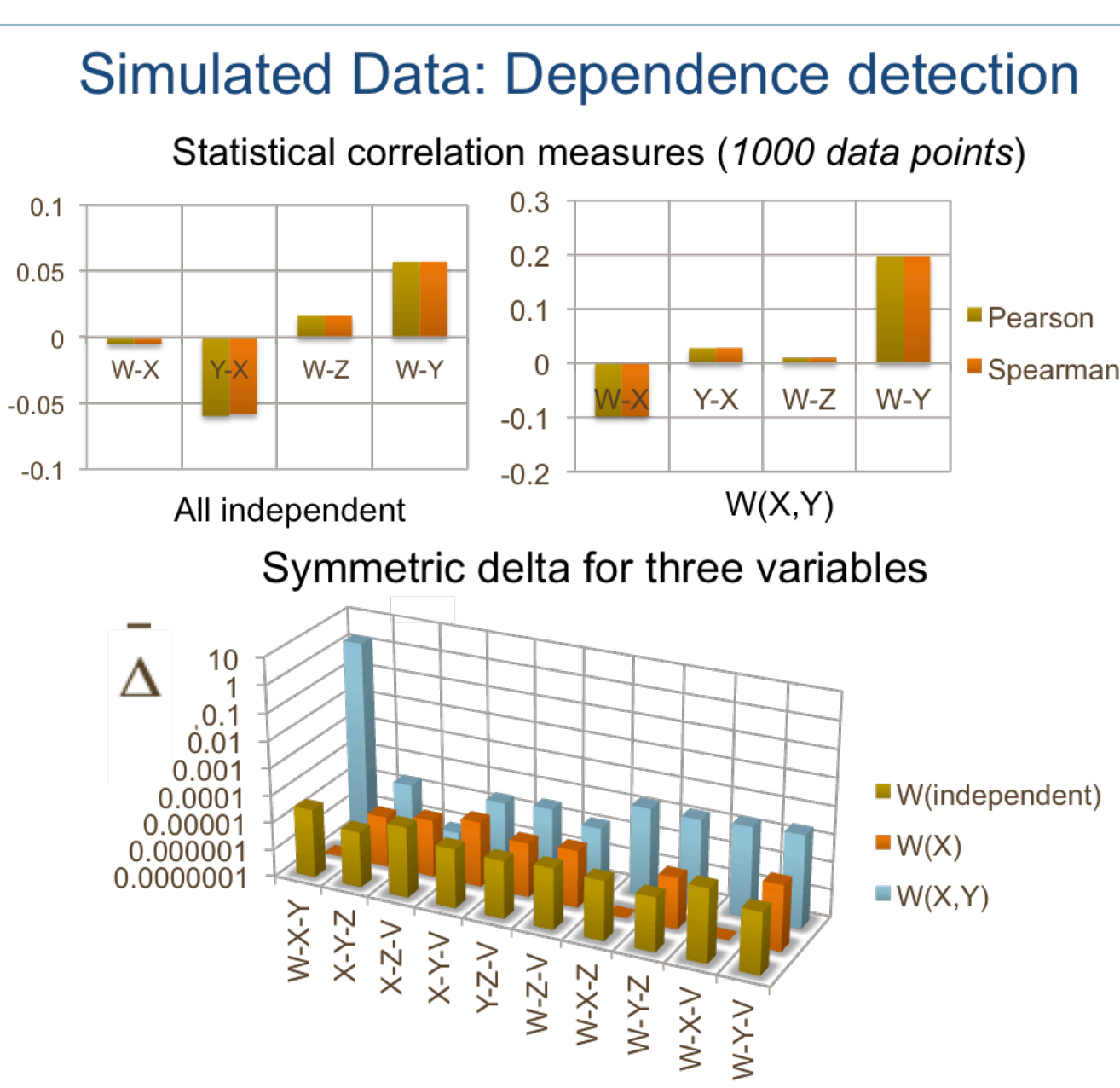$$\bar{\Delta}(V_m) \equiv \prod_{permutations} \Delta(V_{m-1}; X_m)$$

Galas and Sakhanenko, J. Comp. Biol. (2014)

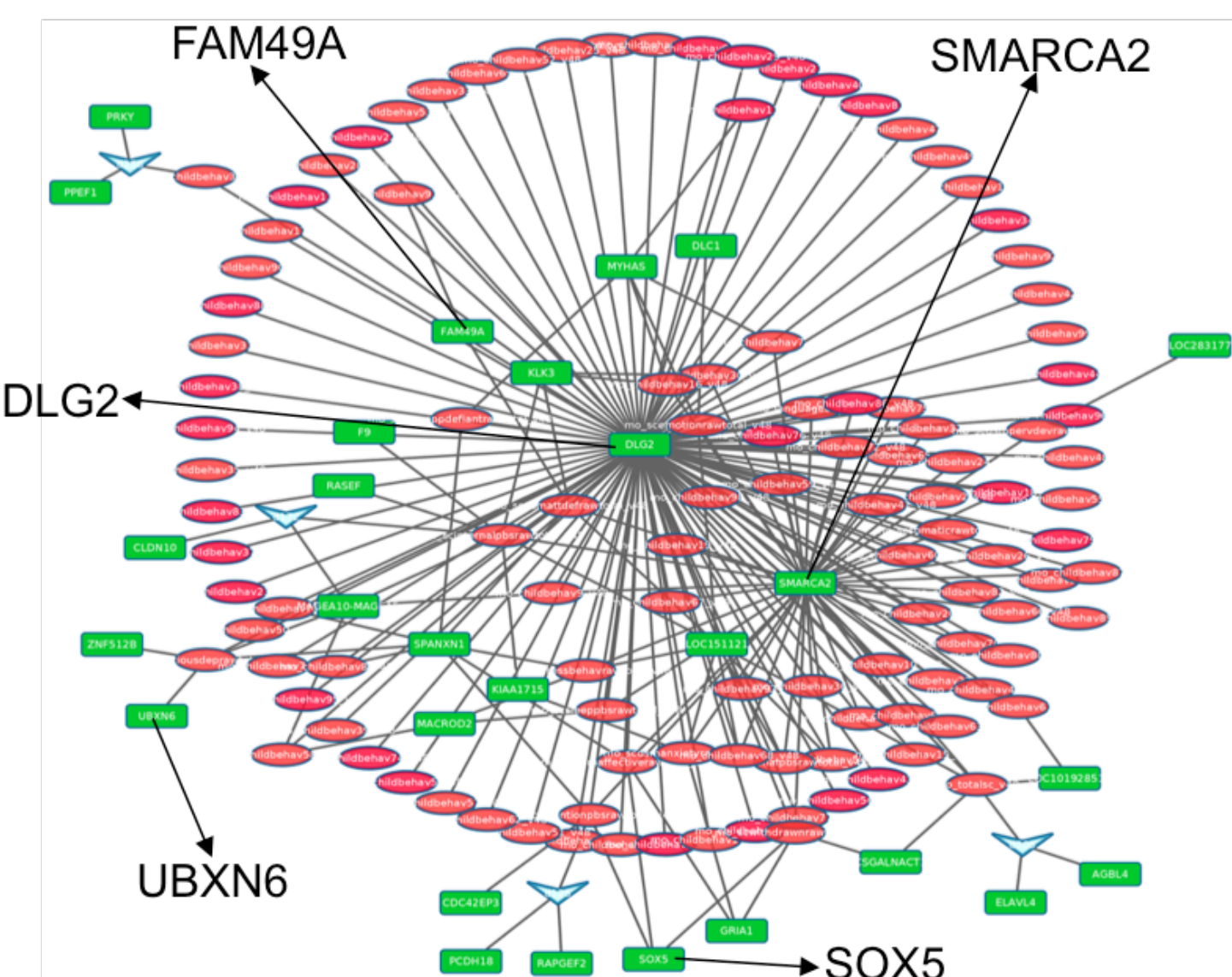### Avoiding the Combinatorial Explosion: the "Shadows Algorithm"

- For $N$ variables, and subsets of size $m$,
- Number of combinations increases like $\sim N^m$
- For $\sim100,000$ variables $m = 3$ & $4$ gives $10^{15}$ & $10^{20}$
- "Shadows algorithm" tracks *shadows* of the high degree dependencies in lower degree calculations.



Sakhanenko and Galas, J. Comp. Biol. (2015)



Simulated Data: Dependence detection
Statistical correlation measures (*1000 data points*)
Pearson / Spearman
All independent / W(X,Y)
Symmetric delta for three variables

**Example**: 3-variable dependence (2 SNPs and 1 phenotype), illustrated by a hypergraph
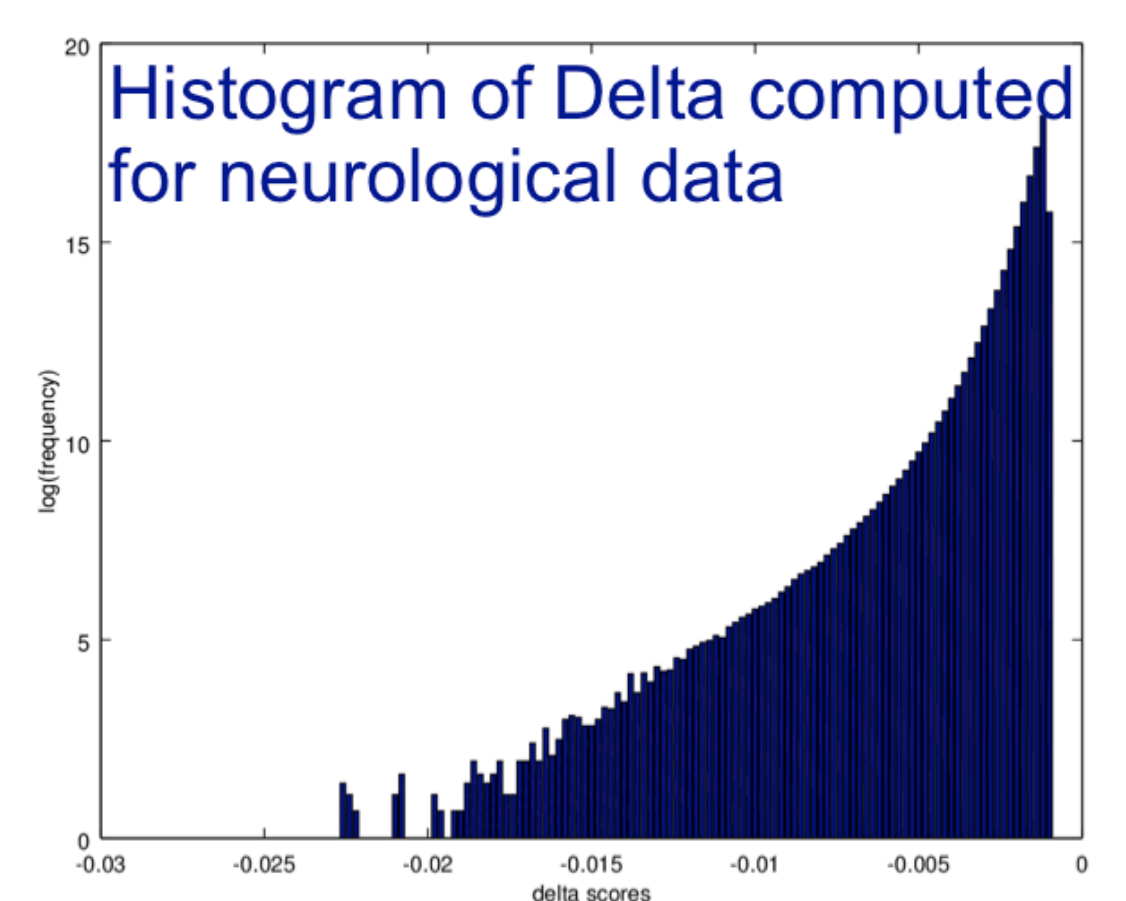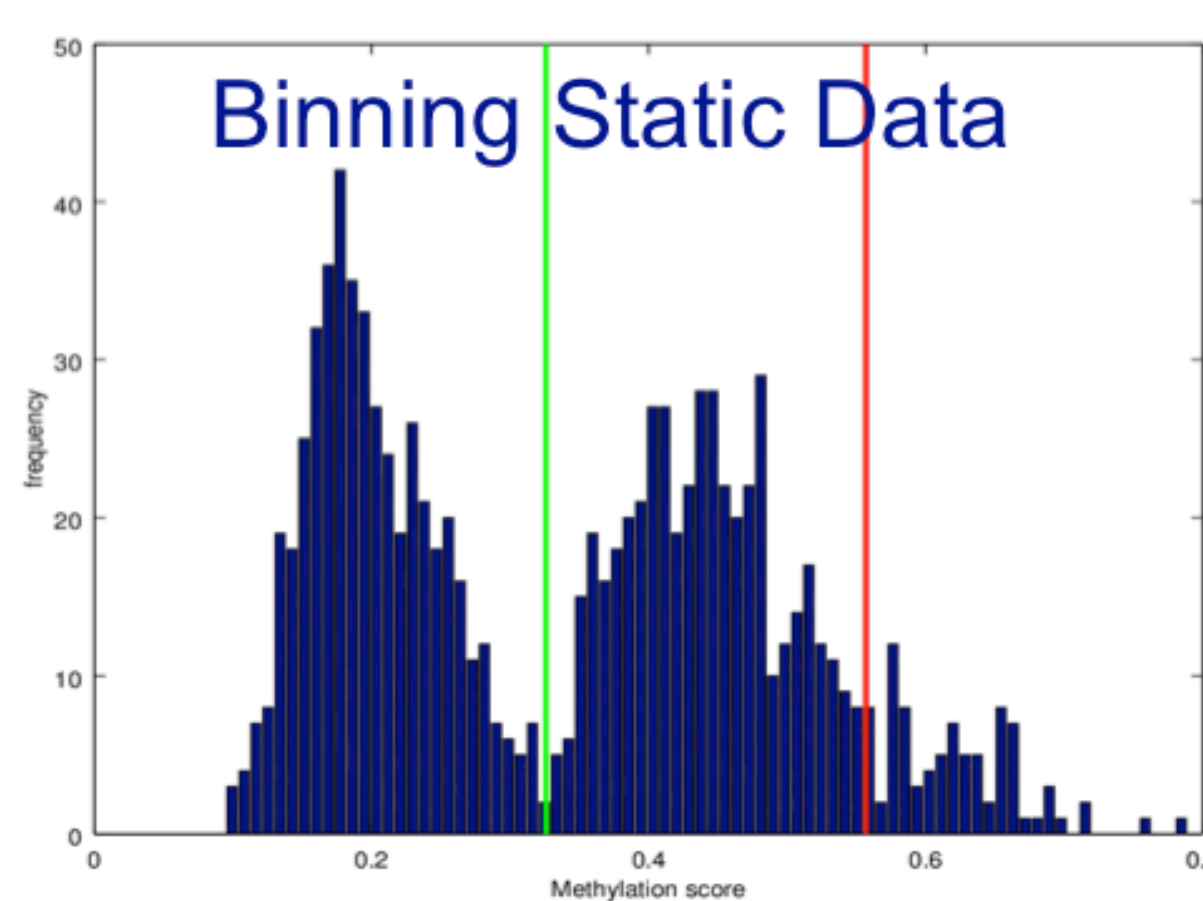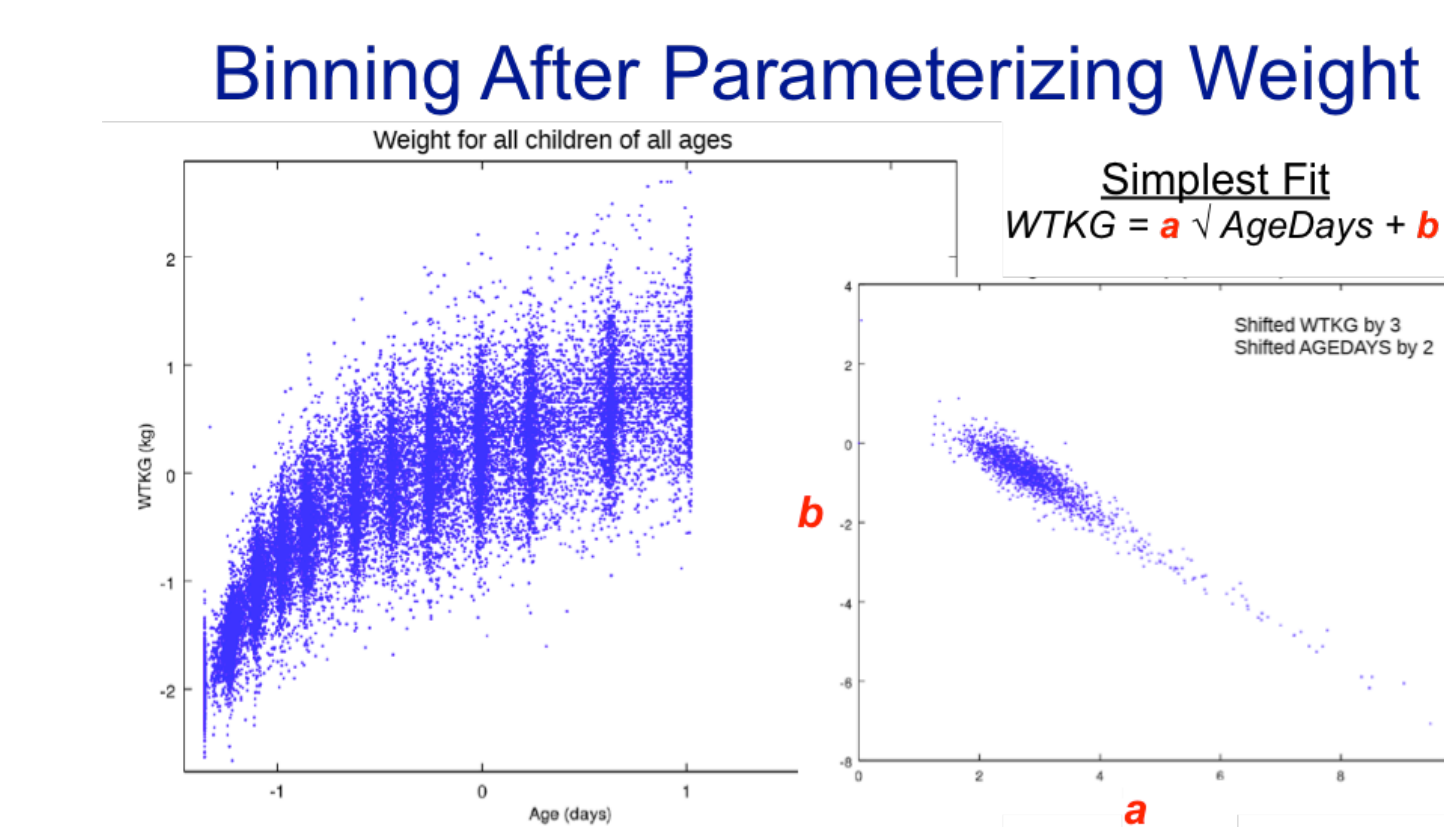


FAM49A, SMARCA2, DLG2, UBXN6, SOX5

Neuro data: top CBCL phenotypes & SNP connections
All point to genes with known neurological effects.

| SNP1 (Gene1) | SNP2 (Gene2) | Phenotype | Description of the phenotype |
|---|---|---|---|
| Rs896996 (DLG2) | Rs1022308 (DLG2) | mo_totalsc_v48 | Total score summary |
| Rs6475635 (SMARCA2) | Rs1555646 (SMARCA2) | mo_childbehav18_v48 | Destroys things belonging to family |
| Rs751192 (FAM49A) | Rs888575 (FAM49A) | mo_scdsmattdefrawtotal_v48 | Attention Deficit/ Hyperactivity total score |
| Rs741923 (UBXN6) | Rs4807584 (MPND) | mo_scanxiousdeprawtrtotal_v48 | Anxious/Depressed: raw total score |



Binning After Parameterizing Weight
Weight for all children of all ages
Simplest Fit
$WTKG = a \sqrt{AgeDays} + b$
Shifted WTKG by 3
Shifted AGEDAYS by 2



Binning Static Data



Histogram of Delta computed for neurological data

## Methods

- This method was model-free and insensitive to under sampling.
- The method was used to detect multivariable dependencies among variables.
- The measures were significantly nonzero only when the subset of variables had an essential, collective dependency [1].
- We used our approach to detect multivariable dependencies in childhood data in a large cohort of children with known genotype, environmental, and phenotype information (GUSTO data from Singapore).
- The approach was used to develop new hypotheses about causal relations.
- We calculated dependency values for variable sets of large degree; this enabled us to identify dependent subsets, but was limited by the combinatorial explosion:

  *Calculations $\sim$ (Number of variables)$^{degree}$*

- We used properties of the measure to avoid the combinatorial explosion by following the "shadows" that the multivariable dependency cast onto smaller subsets [2].
- The Shadows Algorithm enabled us to calculate measures of any degree of dependency.

## Results

- We analyzed a large, high-dimensional data set about the development of Singapore children (GUSTO).
- The GUSTO study collected a diverse range of information about children to capture a full view of child development.
- We considered 3 categories of phenotypes – anthropometric, neurological, and asthma/eczema – and their dependence on genetic variation.
- We identified a small set of strong 2- and 3-variable collective dependencies among phenotypes and SNPs.
- These dependencies formed interconnected networks of variables and enabled us to seek biological relations in these dependencies and form new hypotheses.
- Our method returned a set of candidate multivariable dependencies, which were input to functional analysis.
- The SNP-phenotype dependencies and their networks suggested several involved biological pathways – essential for precise models.
- Genetics can help stratify populations to better detect environmental effect signals.

## Conclusions

**A. The method works.**

- The application of our method to the Singapore data (GUSTO) showed promising initial results.
- We identified complex dependencies in very large and heterogeneous data sets.
- We are adding other types of data to the analysis and integrating them into a more unified network.

**B. Preprocessing of data is extremely important.**

- Missing data and other noise can strongly affect our ability to detect dependencies.
- Binning variable quantities also is key.
- These issues are being actively investigated.

### References

1. Galas DJ, Sakhanenko NA, Skupin A, Ignac T. Describing the complexity of systems: multivariable "set complexity" and the information basis of systems biology. *J Comput Biol*. 2014;21(2):118-140.
2. Sakhanenko NA, Galas DJ. Biological data analysis as an information theory problem: multivariable dependence measures and the shadows algorithm. *J Comput Biol*. 2015;22(11):1005-1024.