



Statistical Significance in Biomarker Discovery

Karsten Borgwardt

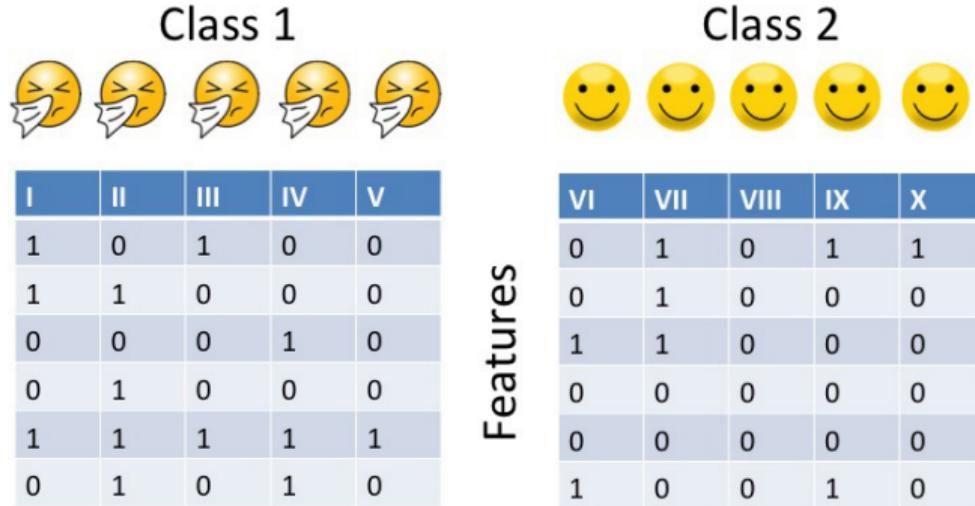
MLPM Summer School, Manchester, September 24, 2015

Biomarker Discovery as a Pattern Mining Problem

Finding groups of disease-related molecular factors

- Single genetic variants, gene expression levels, protein abundancies are often not sufficiently indicative of disease outbreak, progression or therapy outcome.
- Searching for combinations of these molecular factors creates an enormous search space, and two inherent problems:
 - 1 Computational level: How to efficiently search this large space?
 - 2 Statistical level: How to properly account for testing an enormous number of hypotheses?
- The vast majority of current work in this direction (e.g. Achlioptas et al., KDD 2011) focuses on Problem 1, the computational efficiency.
- **But Problem 2, multiple testing, is also of fundamental importance!**

Biomarker Discovery as a Pattern Mining Problem



- Feature Selection: Find features that distinguish classes of objects
- Pattern Mining: Find higher-order **combinations of binary features**, so-called *patterns*, to distinguish one class from another

Mining Significant Patterns

Fisher's exact test

■ Contingency Table

	$S = 1$	$S = 0$	
$y = 1$	a	$n_1 - a$	n_1
$y = 2$	$x - a$	$n - n_1 - x + a$	$n - n_1$
	x	$n - x$	n

- A popular choice is Fisher's exact test to test whether S is overrepresented in one of the two classes.
- The common way to compute p -values for Fisher's exact test is based on the hypergeometric distribution and assumes fixed total marginals (x, n_1, n) .

Mining Significant Patterns

Multiple Testing Problem

- Each S and contingency table corresponds to one hypothesis that is tested.
- To control the Family-Wise Error Rate (probability of detecting at least one false positive), we have to perform multiple testing correction.
- Without multiple testing correction, we will discover millions and billions of false positives in biomarker discovery.
- The classic approach is Bonferroni correction (1936), dividing the significance level α by the number of tests m , that is, $\frac{\alpha}{m}$.

Mining Significant Patterns

Tarone's approach (1990)

- For a discrete test statistics $T(S)$ for a pattern S , such as in Fisher's exact test, there is a minimum obtainable p-value, $p_{min}(S)$.
- For some S , $p_{min}(S) > \frac{\alpha}{m}$. Tarone refers to them as *untestable hypotheses* \bar{S} .
- **Tarone's strategy:** Ignore untestable hypotheses \bar{S} when counting the number of tests m for Bonferroni correction.
- If the p -values of the test are conditioned on the total marginals (as in Fisher's exact test), this does not affect the Family-Wise Error Rate.
- Difficulty: There is an interdependence between m and \bar{S} .

Mining Significant Patterns

Tarone's approach (1990)

- Assume k is the number of tests that we correct for.
- $m(k)$ is the number of testable hypotheses at significance level $\frac{\alpha}{k}$.
- Then the optimization problem is

$$\begin{aligned} \min k \\ \text{s. t. } k \geq m(k) \end{aligned}$$

Mining Significant Patterns

Tarone's approach (1990)

- Assume k is the number of tests that we correct for.
- $m(k)$ is the number of testable hypotheses at significance level $\frac{\alpha}{k}$.

procedure TARONE

$k := 1$;

while $k < m(k)$ **do**

$k := k + 1$;

return k

Mining Significant Patterns

Terada's link to frequent itemset mining (Terada et al., PNAS 2013)

- For $0 \leq x \leq n_1$, the minimum p-value $p_{min}(S)$ decreases monotonically with x .
- One can use *frequent itemset mining* to find all S that are testable at level α , with frequency $\psi^{-1}(\alpha)$.
- They propose to use a decremental search strategy:

procedure TERADA'S DECREMENTAL SEARCH (LAMP)

$k :=$ "very large";

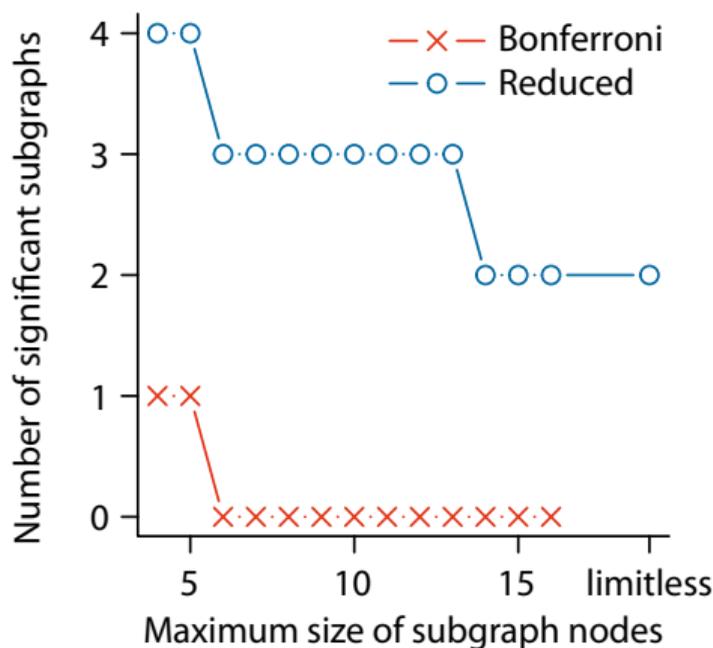
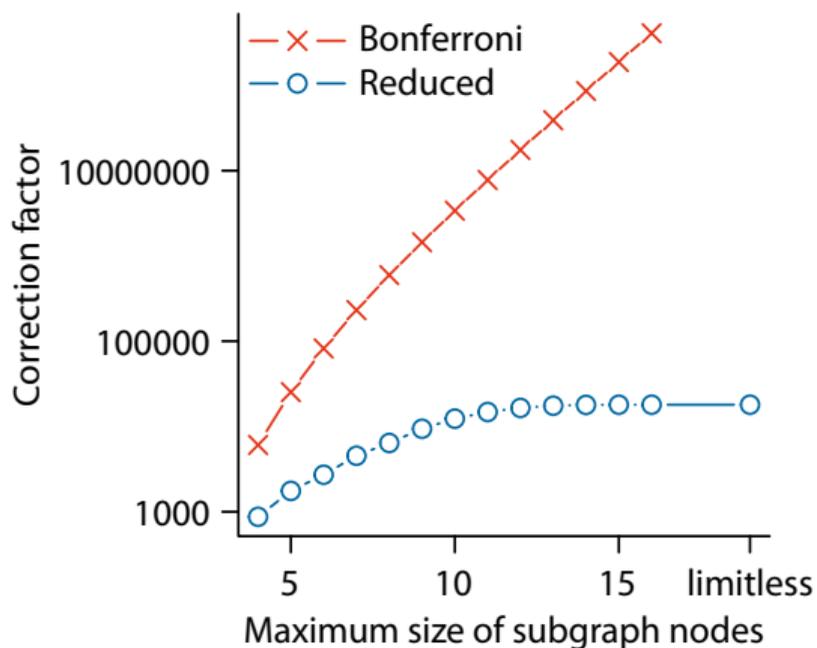
while $k > m(k)$ **do**

$k := k - 1$;

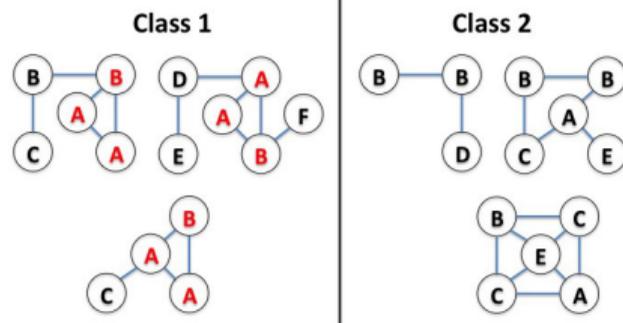
$m(k) :=$ frequent itemset mining($D, \psi^{-1}(\frac{\alpha}{k})$);

return $k + 1$

Example: PTC dataset (Helma et al., 2001)



Significant Subgraph Mining (Sugiyama et al., SDM 2015)



Significant Subgraph Mining

- Each object is a graph.
- A pattern is a subgraph in these graphs.
- Typical application in Drug Development: Find subgraphs that discriminate between molecules with and without drug effect.
- Counting all tests (= all patterns) requires exponential runtime in the number of nodes.

Significant Subgraph Mining (Sugiyama et al., SDM 2015)

Incremental search with early stopping

- **procedure** INCREMENTAL SEARCH WITH EARLY STOPPING

$\theta := 0$

repeat

$\theta := \theta + 1; FS_{\theta} := 0;$

repeat

find next frequent subgraph at frequency θ

$FS_{\theta} := FS_{\theta} + 1$

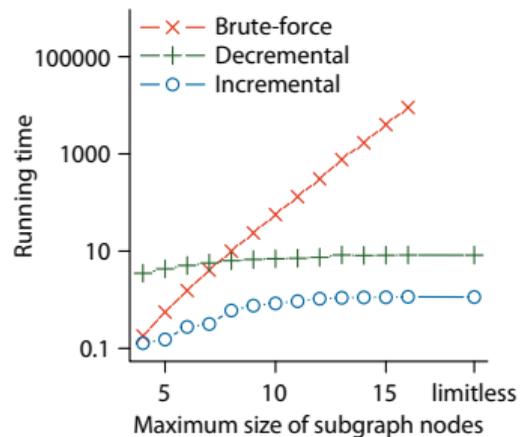
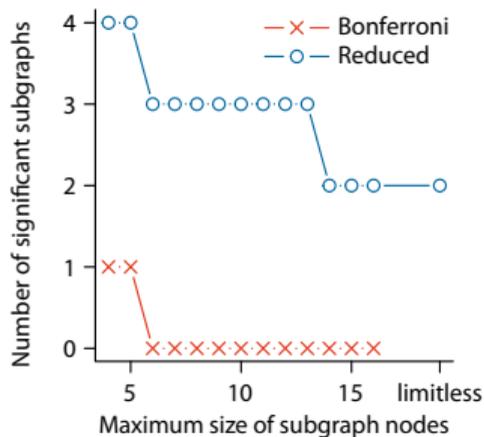
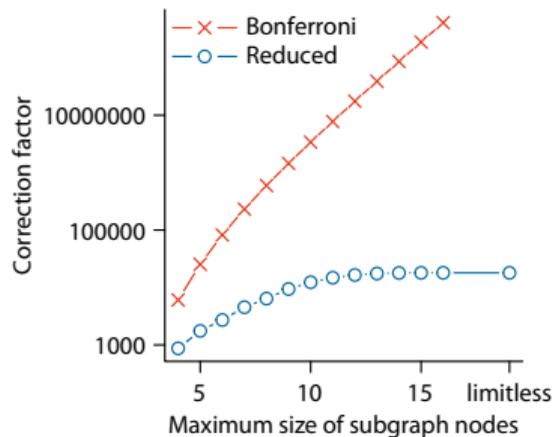
until (no more frequent subgraph found) or $(FS_{\theta} > \frac{\alpha}{\psi(\theta)})$

until $FS_{\theta} \leq \frac{\alpha}{\psi(\theta)}$

return $\psi(\theta)$

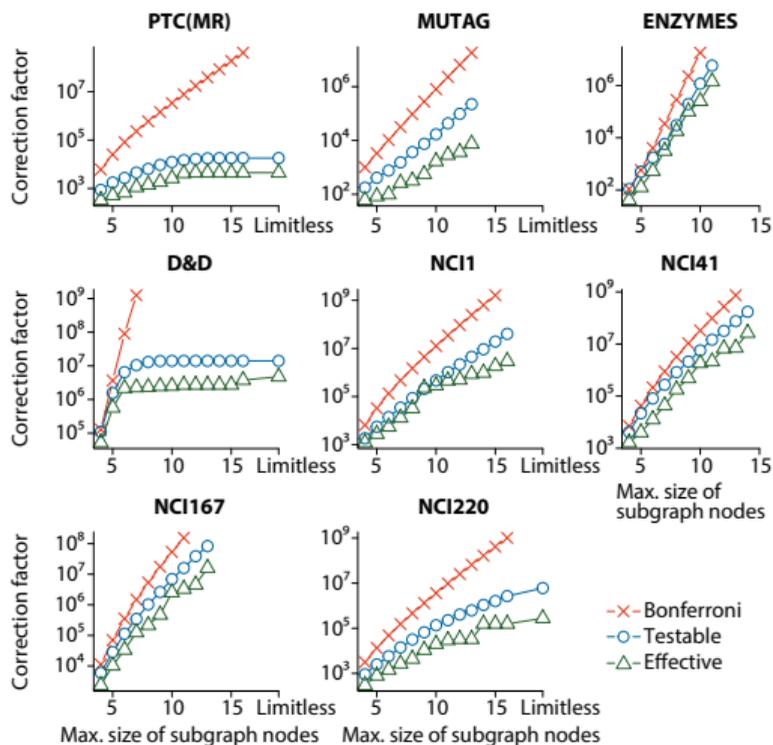
- $\frac{\alpha}{\psi(\theta)}$ is the maximum correction factor, such that subgraphs with frequency θ can be significant at level $\psi(\theta)$.

Significant Subgraph Mining on PTC Dataset

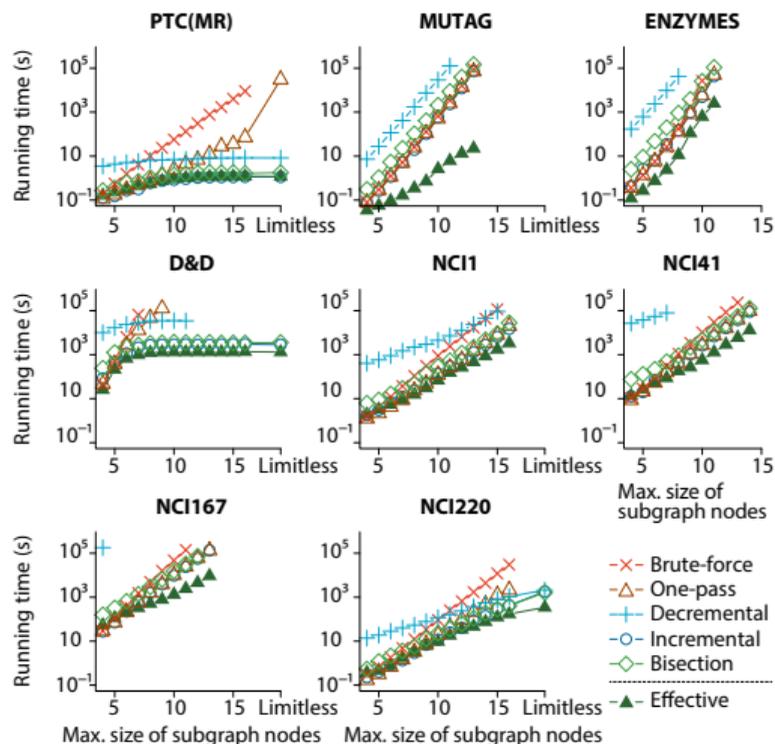


Dataset from Helma et al. (2001)

Significant Subgraph Mining: Correction Factor



Significant Subgraph Mining: Runtime



Westfall-Young light (Llinares-Lopez et al., KDD 2015)

Dependence between hypotheses

- As patterns are often in sub-/superpattern-relationships, they do not constitute independent hypotheses.
- Informally: The underlying number of hypotheses may be much lower than the raw count.
- Westfall-Young-Permutation tests (Westfall and Young, 1993), in which the class labels are repeatedly permuted to approximate the null distribution, are one strategy to take this dependence into account.
- **Computational problem: How to efficiently perform these thousands of permutations?**
- There is one existing approach, FastWY (Terada et al., ICBB 2013), which suffers from either memory or runtime problems.

Westfall-Young light (Llinares-Lopez et al., KDD 2015)

The Algorithm

- 1 Input:** Transactions D , class labels \mathbf{y} , target FWER α , number of permutations j_p .
- 2** Perform j_p permutations of the class label \mathbf{y} and store each permutation as \mathbf{c}_j .
- 3** Initialize $\theta := 1$ and $\delta^* := \psi(\theta)$ and $p_{min}^{(j)} := 1$.
- 4** Perform a depth first search on the patterns:
 - Compute the p -value of pattern S across all permutations, update $p_{min}^{(j)}$ if necessary.
 - Update δ^* by α -quantile of $p_{min}^{(j)}$, and increase θ accordingly.
 - Process all children of S with frequency $\geq \psi^{-1}(\delta^*)$.
- 5 Output:** Corrected significance threshold δ^* .

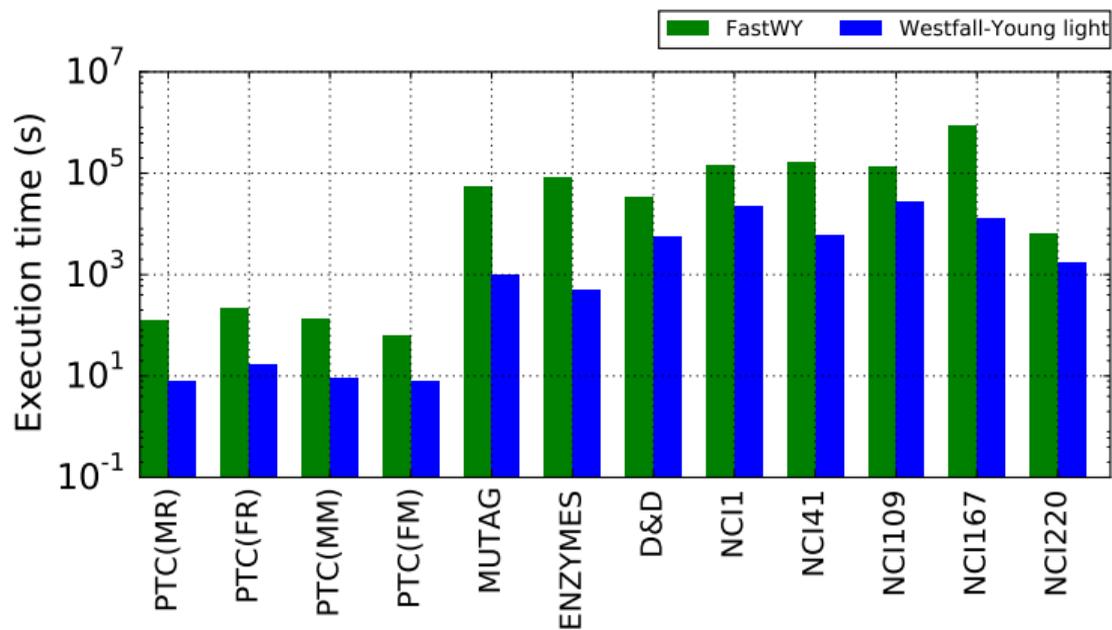
Westfall-Young light (Llinares-Lopez et al., KDD 2015)

Speed-up tricks of Westfall-Young light

- Follows incremental search strategy rather than decremental search strategy of FastWY
- Performs only one iteration of frequent pattern mining
- Does not store the occurrence list of patterns
- Does not compute the upper $1 - \alpha$ quantile of minimum p-values exactly.
- Reduces the number of cell counts that have to be evaluated
- Shares the computation of p-values across permutations

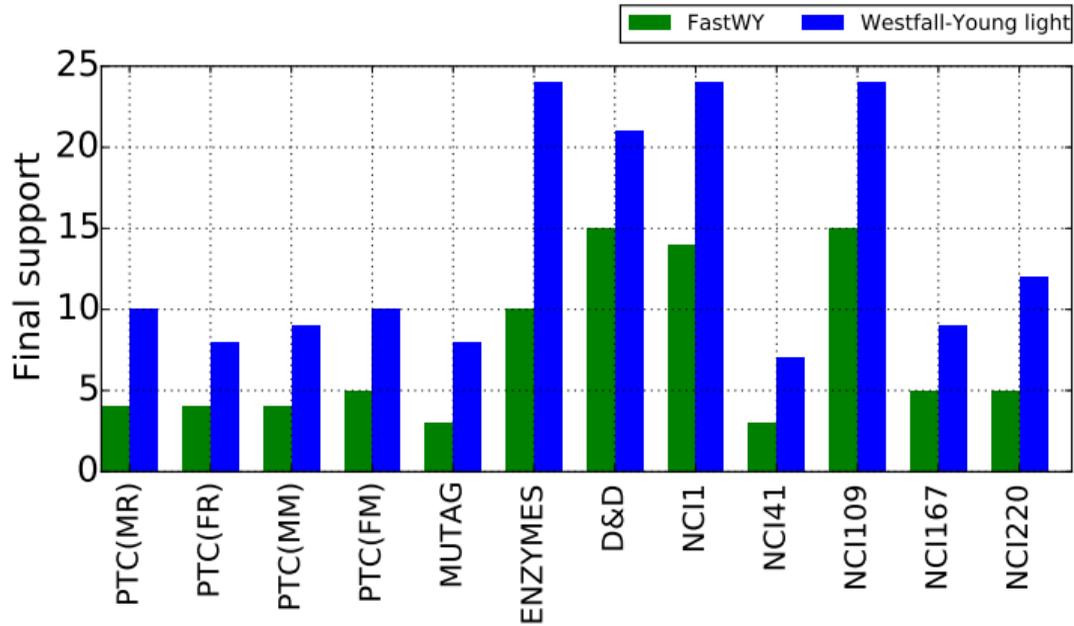
Westfall-Young light

■ Runtime



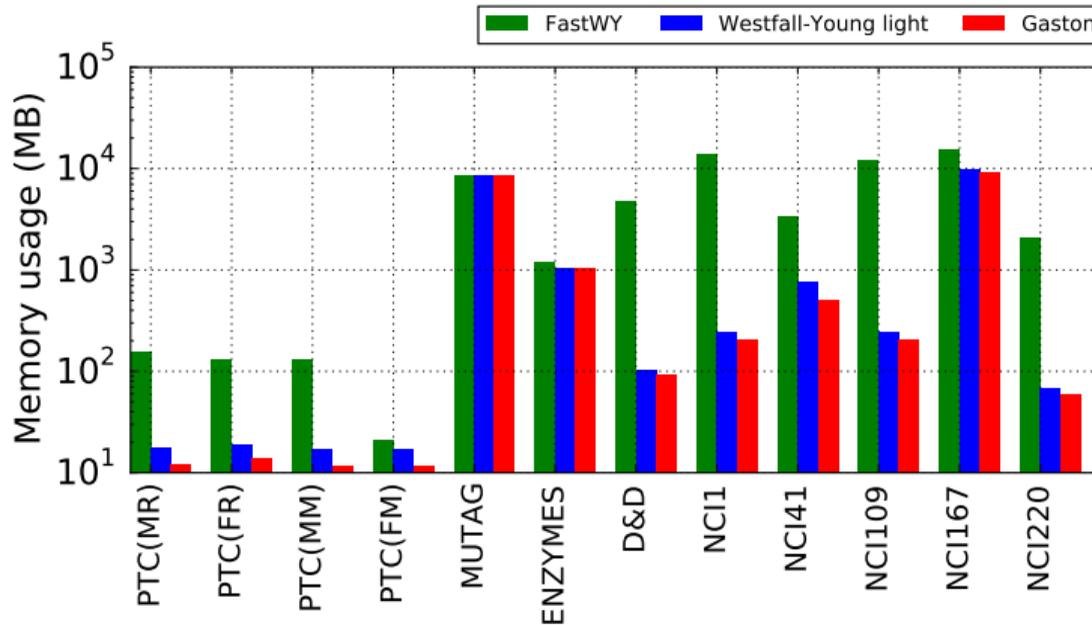
Westfall-Young light

- Final frequency threshold (support)



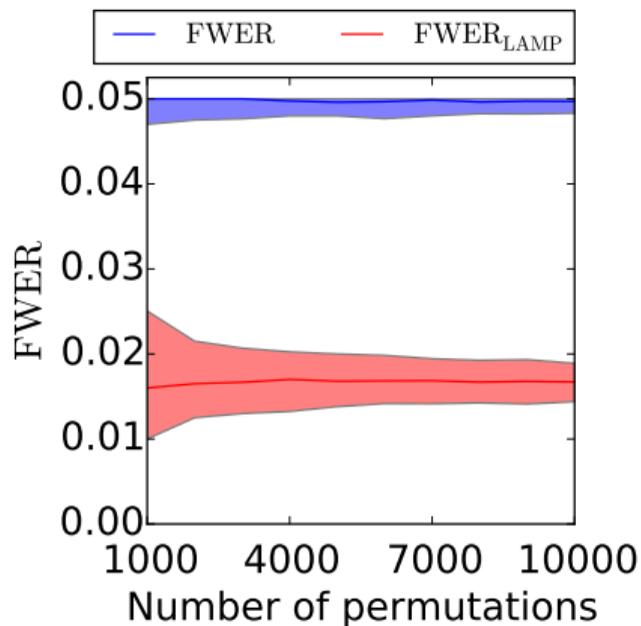
Westfall-Young light

■ Peak memory usage



Westfall-Young light

- Better control of the Family-wise error rate (Enzymes)



FAIS: Finding intervals that exhibit genetic heterogeneity

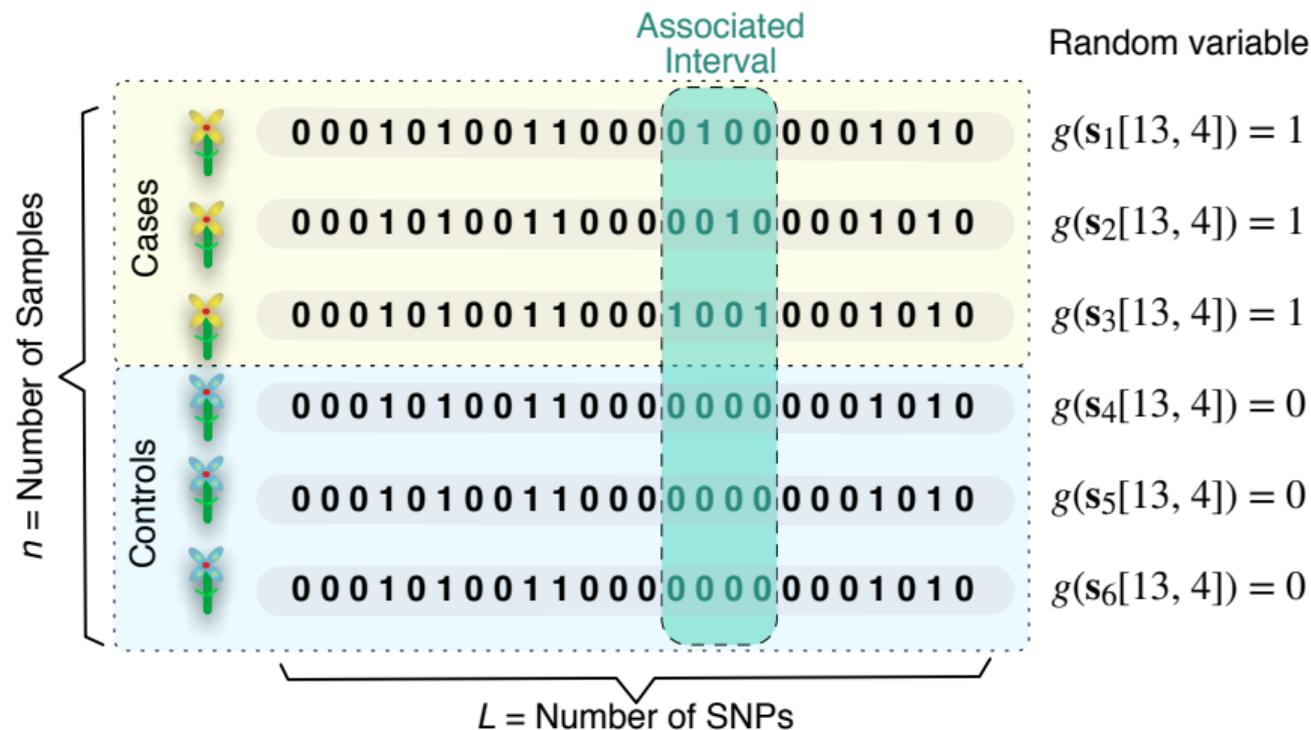
Genetic heterogeneity

- Genetic heterogeneity refers to the phenomenon that several different genes or sequence variants may give rise to the same phenotype.
- The correlation between each individual gene or variant and the phenotype may be too weak to be detected, but the group may have a strong correlation.
- The only current way to consider genetic heterogeneity is to consider fixed groups of variants. Genome-wide scans cause tremendous computational and statistical problems.

Fast Automatic Interval Search (Llinares-Lopez et al., ISMB 2015)

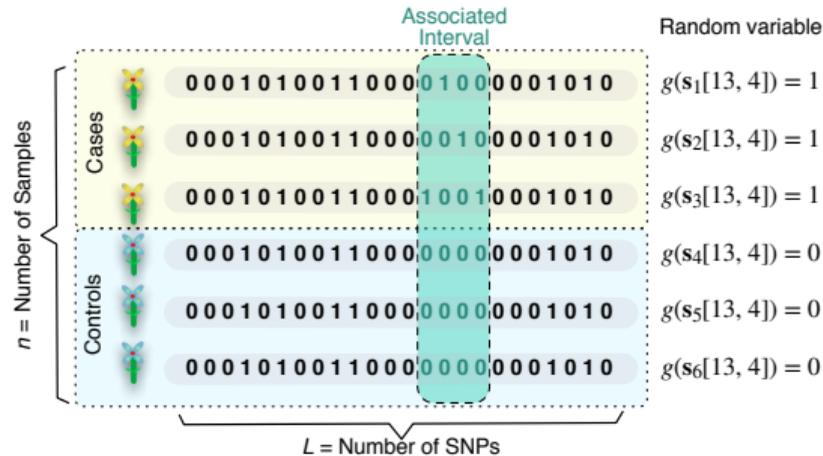
- FAIS finds all contiguous sets of variants that are significantly associated with a given phenotype under a model of genetic heterogeneity.

FAIS: Finding intervals that exhibit genetic heterogeneity



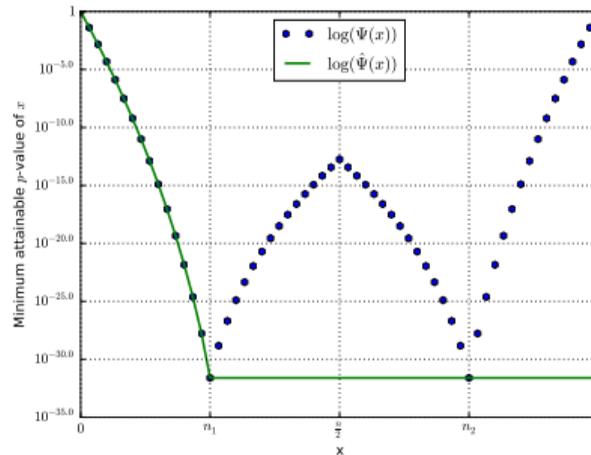
FAIS: Finding intervals that exhibit genetic heterogeneity

Finding trait-associated genome **segments** with at least one minor allele



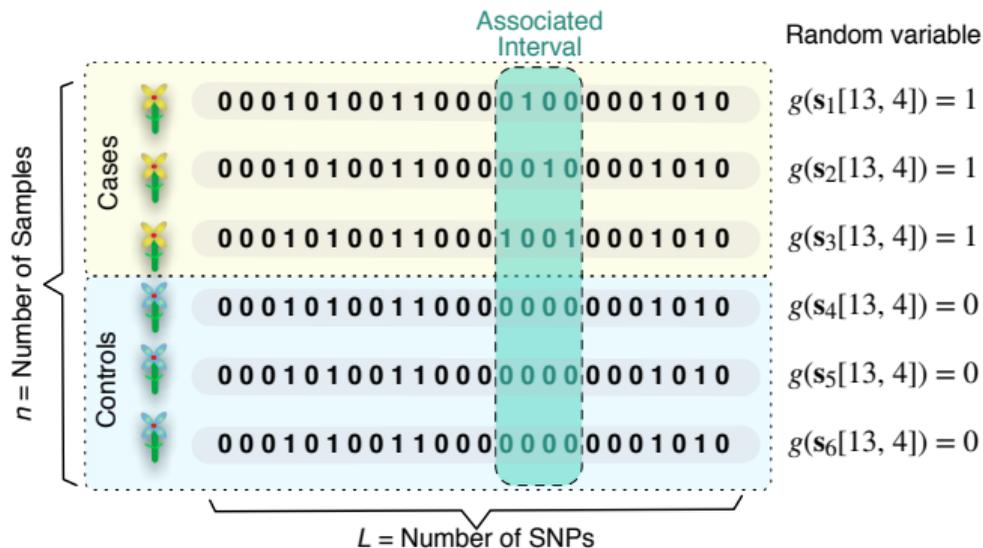
- An interval is represented by its maximum value. The longer an interval, the more likely it is that this maximum is 1.

FAIS: Finding intervals that exhibit genetic heterogeneity



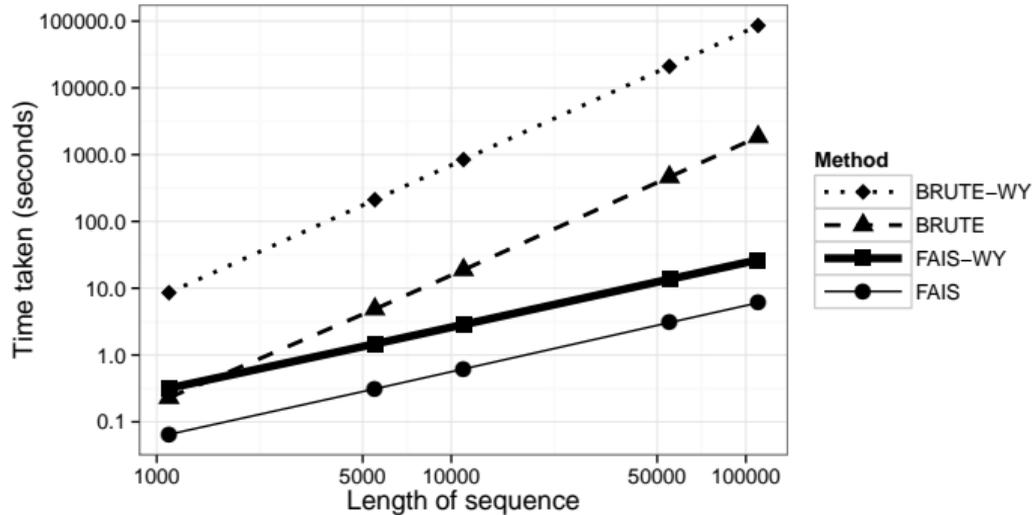
- **Pruning criterion 1:** If an interval is represented by 1 for too many individuals, the interval is not testable.

FAIS: Finding intervals that exhibit genetic heterogeneity



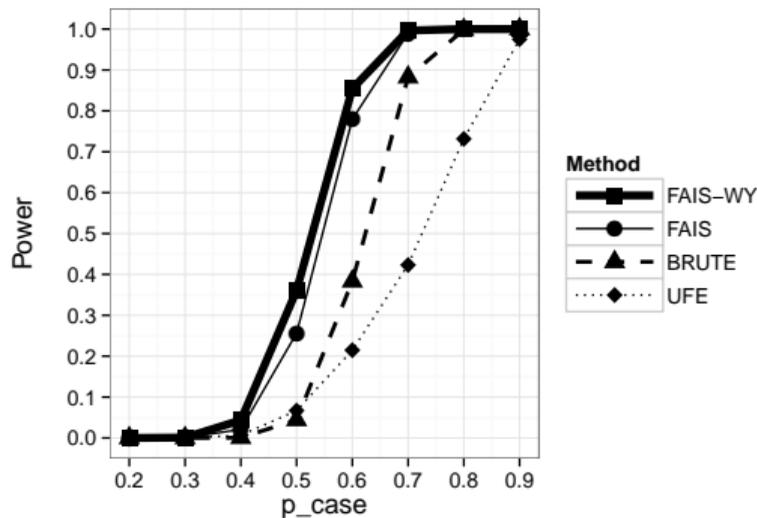
- Pruning criterion 2:** If an interval is too frequent to be testable, then none of its superintervals is testable.

FAIS: Finding intervals that exhibit genetic heterogeneity



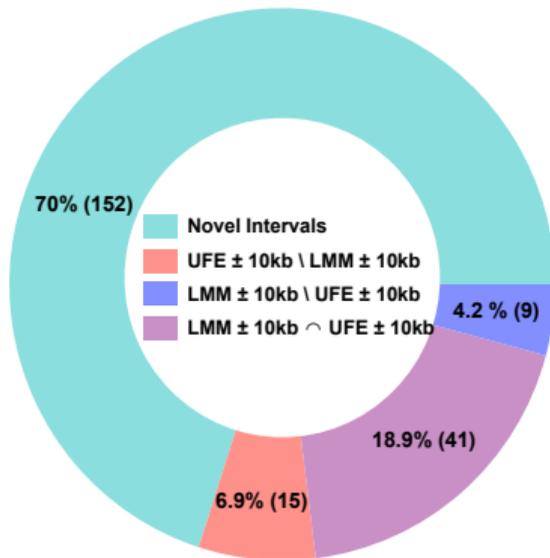
- Our method FAIS (Fast Automatic Interval Search) improves over the brute-force interval search in terms of runtime in simulations.

FAIS: Finding intervals that exhibit genetic heterogeneity



- Our method FAIS (Fast Automatic Interval Search) improves over brute-force interval search and univariate approaches in terms of power in simulations.

FAIS: Finding intervals that exhibit genetic heterogeneity



- Most significant intervals would have been missed by univariate approaches (UFE and LMM) on 21 binary phenotypes from *Arabidopsis thaliana* (Atwell et al., Nature 2010).

Outlook

Current and future topics

- Better empirical understanding of the impact of considering testability
- Pattern summarization
- Conditioning on covariates, e.g. to model population structure: Recent arxiv paper (Llinares-Lopez et al., 2015) which ignores untestable patterns in the Cochran-Mantel-Haenszel test on $K \times 2 \times 2$ contingency tables.
- Two postdoc positions and one PhD student position are available within this Starting Grant project 'Significant Pattern Mining'.

Thank You

- Felipe Llinares Lopez
- Menno Witteveen
- Dean Bodenham
- Udo Gieraths
- Dominik Grimm
- Elisabetta Ghisu
- Anja Gumpinger
- Xiao He
- Laetitia Papaxanthos
- Damian Roqueiro
- Birgit Knapp



Sponsors:

- Krupp-Stiftung
- Marie-Curie-FP 7
- SNSF Starting Grant (ERC backup)
- Horizon 2020

References: <http://www.bsse.ethz.ch/mlcb>

References I

-  R. E. Tarone, *Biometrics* **46**, 515 (1990).
-  P. H. Westfall, S. S. Young, *Statistics in Medicine* **13**, 1084 (1993).
-  D. R. Nyholt, *American Journal of Human Genetics* **74**, 765 (2004).
-  A. Terada, M. Okada-Hatakeyama, K. Tsuda, J. Sese, *Proceedings of the National Academy of Sciences* **110**, 12996 (2013).
-  A. Terada, K. Tsuda, J. Sese, *IEEE International Conference on Bioinformatics and Biomedicine* (2013), pp. 153–158.
-  M. Sugiyama, F. Llinares-López, N. Kasenburg, K. M. Borgwardt, *SIAM Data Mining* (2015).

References II

-  F. Llinares-López, M. Sugiyama, L. Papaxanthos, K. M. Borgwardt, *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Sydney, NSW, Australia, August 10-13, 2015*, L. Cao, et al., eds. (ACM, 2015), pp. 725–734.
-  F. Llinares-López, et al., *Bioinformatics* **31**, 240 (2015).