

Figure 2: Bioset Ranking and Directionality

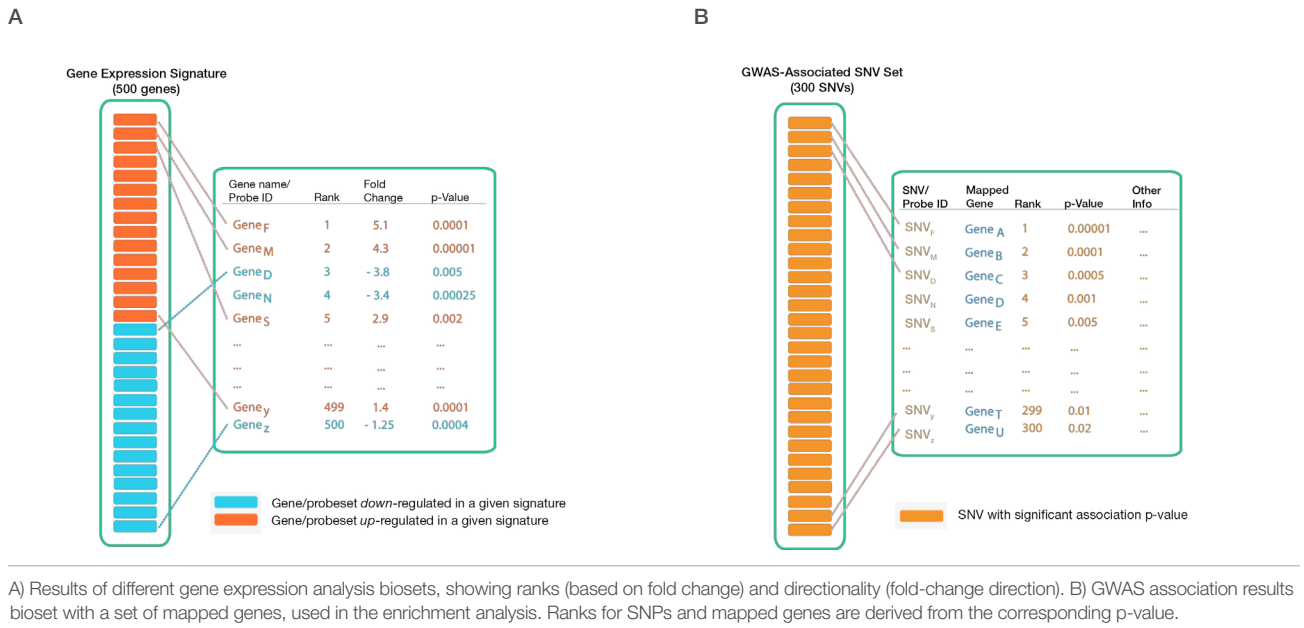
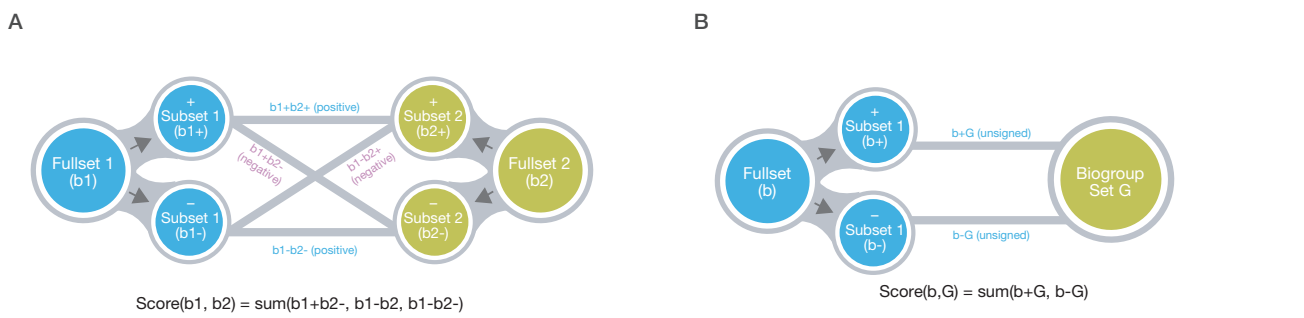


Figure 3: Rank-Based Directional Enrichment



The directional relationship between the 2 signatures is captured by the sign of the correlation score. Upregulated genes (*b+*) and downregulated genes (*b-*) are separated into directional subsets, and correlation scores are computed for each directional subset from one signature (*b1+*, *b1-*) against each subset from the other signature (*b2+*, *b2-*). A positive sign is given to a subset pair of the same direction (*b1+b2+*, *b1-b2-*), and a negative sign is given to a subset pair of opposite directions (*b1+b2-*, *b1-b2+*). The overall correlation

score is the sum of directional subset scores, and the sign of the sum determines whether the 2 signatures are positively or negatively correlated (Figure 3A).

BaseSpace Engine biogroups are genesets that are not directional and not ranked. The enrichment of a biogroup in a bioset is computed for each directional subset (*b+*, *b-*) comprising the bioset. Because the biogroups are not directional, the subset scores are unsigned and the overall enrichment score is the sum of the subset scores (Figure 3B).

Computing Pairwise Correlation Scores

The correlation scores between gene signature (bioset) pairs are computed with the Running Fisher algorithm, a nonparametric rank-based statistical approach. Detailed steps for 2 signature sets (*b1*, *b2*) are as follows:

First, each gene signature set is rank-ordered according to fold change, p-value, or a particular score. If appropriate metrics are not provided, then the gene signature set is unranked. The upregulated genes and downregulated genes are noted with positive and negative signs to imply directionality, respectively. A directional subset is generated for each direction, such as *b1+*, *b1-*, *b2+*, and *b2-* (Figure 3A). If no directional data are provided, then the gene signature set is not directional and only one subset is formed with the whole signature set, such as *b1o*, or *b2o*.

Second, all the subset pairs are identified: *b1Di*, *b2Dj*, where *Di* and *Dj* are the available directions (+, -, or o) in *b1* and *b2*, respectively. The Running Fisher algorithm is applied to each subset pair. The top 1% of genes in the platform in the first subset *b1Di* are collected as a group *G*. The second subset, *b2Dj*, is scanned top to bottom in the rank order to identify each rank with a gene matching a member in the group *G*. If the subset is unranked, all the genes in the subset are retrieved at the first scan.

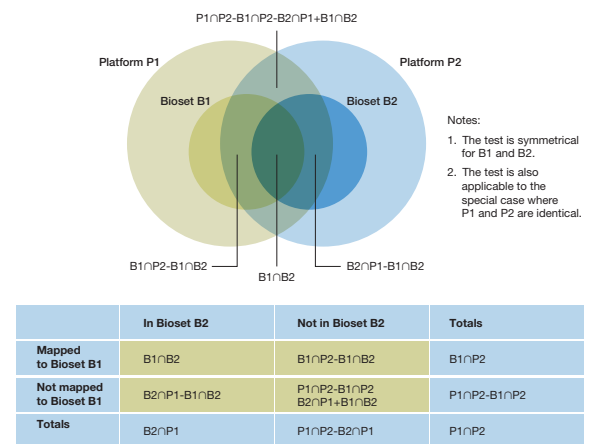
At each matching rank *K*, the scanned portion of the second subset *b2Dj* consists of *N* genes, and the overlap between group *G* and *N* genes is *M*. A Fisher's exact test is performed at rank *K*, to evaluate the statistical significance of observing *M* overlaps between a set of size *G* and a set of size *N*. In this analysis, the set of size *G* comes from platform *P1* and the set of size *N* comes from platform *P2*, given the sizes of *P1* and *P2* as well as the overlap between *P1* and *P2*. The parameters for Fisher's exact test are determined as shown in Figure 4.

At the end of the scan, the best p-value is retained, and the negative log of the best p-value is a score for the subset pair.

Next, the Running Fisher algorithm is performed in the reverse direction: the top-ranking genes in the second subset *b2Dj* are collected as a group *G*, and the first subset *b1Di* is scanned in the rank order. The same procedure in this reverse direction produces another score for the same subset pair. The 2 scores are averaged to represent the magnitude of the similarity between the 2 subsets. A positive sign is given to the final subset pair score if *Di* and *Dj* are the same. A negative sign is given if *Di* and *Dj* are opposite. The score is unsigned if any of *Di* and *Dj* are not directional.

Finally, the overall score is computed by summing up all directional subset pair scores (Figure 3A). The sign of the sum determines whether the 2 signatures are positively or negatively correlated. If both signatures are not directional, the overall score is the same as the only subset pair score that is unsigned. If one of the 2 signature sets is directional and the other is not directional, the overall score is represented by the larger of the 2 subset pair scores, annotated with the contributing direction from the directional signature. The matching genes between 2 typical gene signatures are depicted in Figure 5. The directionality and rank for each gene is shown.

Figure 4: Pairwise Correlation between Biosets



Fisher's exact test parameters for Bioset vs. Bioset—assessment of statistical significance on a bioset's enrichment in another bioset. When the bioset is directional, the direction-specific subset (*b1+* or *b1-* in place of *B1*; *b2+* or *b2-* in place of *B2*) is used. The scan of the bioset dynamically cuts off the top *K* rows of the gene set for evaluation at rank *K*.

Extending the Bioset-Bioset Scoring Algorithm

The bioset-bioset scoring algorithm can be extended to compute the correlation score between 2 biosets of other data types. For example, the algorithm can be used to compute correlation scores between 2 exon data sets, after the following requirements are satisfied:

1. Identification of exons from 2 different platforms can be mapped to a common standard set of exon IDs.
2. All exons in each platform are stored so that the overlap between one exon set with the platform for the other exon set can be computed, as well as the overlap between the 2 platforms.
3. Rank and directionality for each exon are defined and retrievable.

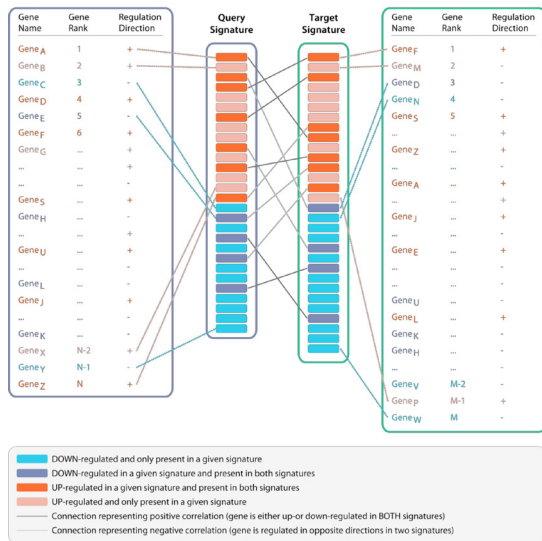
Similarly, other data types—such as isoform-specific transcript data sets or microRNA (miRNA) data sets—can be compared against any other data set within the space of same data types. Also, the mapping from these nongene data types to genes generates a gene-level representation of the data set, so that all the data sets can be compared at the gene level.

Computing the Enrichment Score for a Biogroup in a Bioset

The enrichment score for a biogroup in a bioset is computed with the Running Fisher algorithm. The detailed steps given a signature set (*b*) and a geneset in a biogroup (*C*) are as follows:

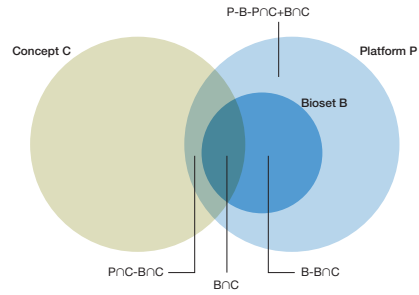
First, the gene signature set is rank-ordered and assigned directionality as previously described for biosets. Directional subsets are also generated as described (Figure 3B).

Figure 5: Pairwise Signature Correlation Scores



An outline of pairwise directional enrichment analysis between 2 signatures.

Figure 6: Bioset vs. Biogroup



	In Bioset B	Not in Bioset B	Totals
Mapped to C	BnC	PnC-BnC	PnC
Not mapped to C	B-BnC	P-B-PnC+BnC	P-PnC
Totals	B	P-B	P

Fisher's Exact Test parameters for Bioset vs. Biogroup—assessment of statistical significance on a biogroup's enrichment in a bioset. When the bioset is directional, the direction-specific subset ($b+$ or $b-$ in place of B) is used. The scan of the bioset dynamically cuts off the top K rows of the gene set for evaluation at rank K.

Second, the Running Fisher algorithm is applied to compute the enrichment of the biogroup genes in each subset. The subset $b+$ or $b-$ is scanned top to bottom in the rank order to identify each rank with a gene matching a member in the biogroup C. If the subset is unranked, all the genes in the subset are retrieved at the first scan.

At each matching rank K, the scanned portion of the subset $b+$ or $b-$ consists of N genes, and the overlap between group C and N genes is M. A Fisher's exact test is performed at rank K, to evaluate the statistical significance of observing M overlaps between a set of size C and a set of size N, where the set of size N comes from platform P, given the size of P and the overlap between P and C. The parameters for Fisher's exact test are determined as in Figure 6.

At the end of the scan, the best p-value is retained. The negative natural log of the best p-value is the enrichment score for the biogroup in the subset.

Finally, the overall score is computed by summing up the subset pair scores (Figure 3B). The sum score is unsigned, as are the subset scores.

References

1. Kupersmidt I, Su QJ, Grewal A, Sundaresh S, Halperin I, et al. (2010). Ontology-based metaanalysis of global collections of high-throughput public data. PLoS ONE 5(9): e13066.
2. Lamb J, Crawford ED, Peck D, Modell JW, Blat IC, et al. (2006) The Connectivity Map: using gene-expression signatures to connect small molecules, genes, and disease. Science 313(5795): 1929–1935.
3. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, et al. (2005) Gene set enrichment analysis: a knowledge- based approach for interpreting genome-wide expression profiles. Proc Natl Acad Sci USA 102(43): 15545–15550.
4. Newton MA, Quintana FA, den Boon JA, Sengupta S, Ahlquist P (2007) Random-set methods identify distinct aspects of the enrichment signal in gene-set analysis. Ann Appl Stat 1(1): 85–106.