

Finding Consistent Disease Subnetworks

Limsoon Wong
26 March 2010

(Joint work with Donny Soh, Difeng Dong, Yike Guo)



2

Plan



- **An issue in gene expression analysis**
- **Comparing pathway sources**
 - Comprehensiveness
 - Consistency
 - Compatibility
- **Matching pathways in different sources**
- **Finding more consistent disease subnetworks**

An Issue in Gene Expression Analysis



4



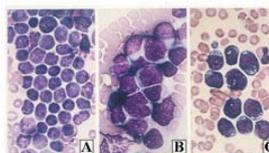
First, the good news..

5



Childhood Acute Lymphoblastic Leukemia

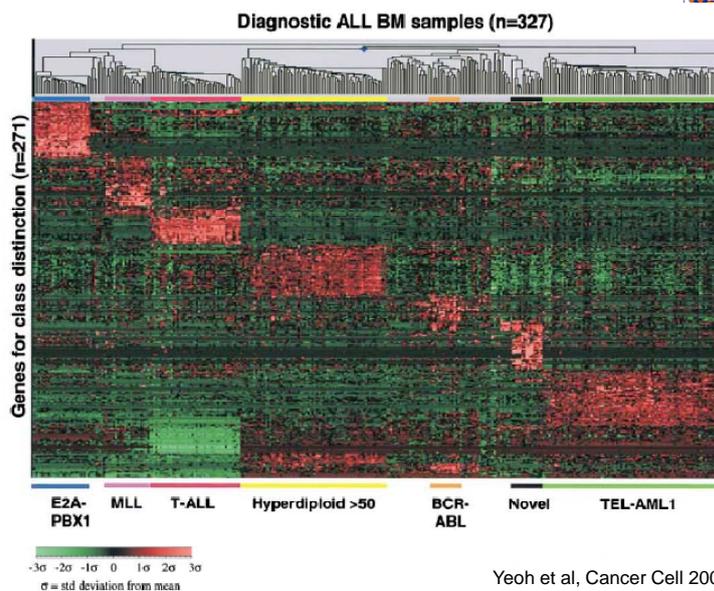
- Major subtypes: T-ALL, E2A-PBX1, TEL-AML, BCR-ABL, MLL genome rearrangements, Hyperdiploid >50
- Diff subtypes respond differently to same Tx
- Over-intensive Tx
 - Development of secondary cancers
 - Reduction of IQ
- Under-intensive Tx
 - Relapse
- The subtypes look similar
- Conventional diagnosis
 - Immunophenotyping
 - Cytogenetics
 - Molecular diagnostics
 - ⇒ Unavailable in developing countries



Talk at HKUST, 26 March 2010

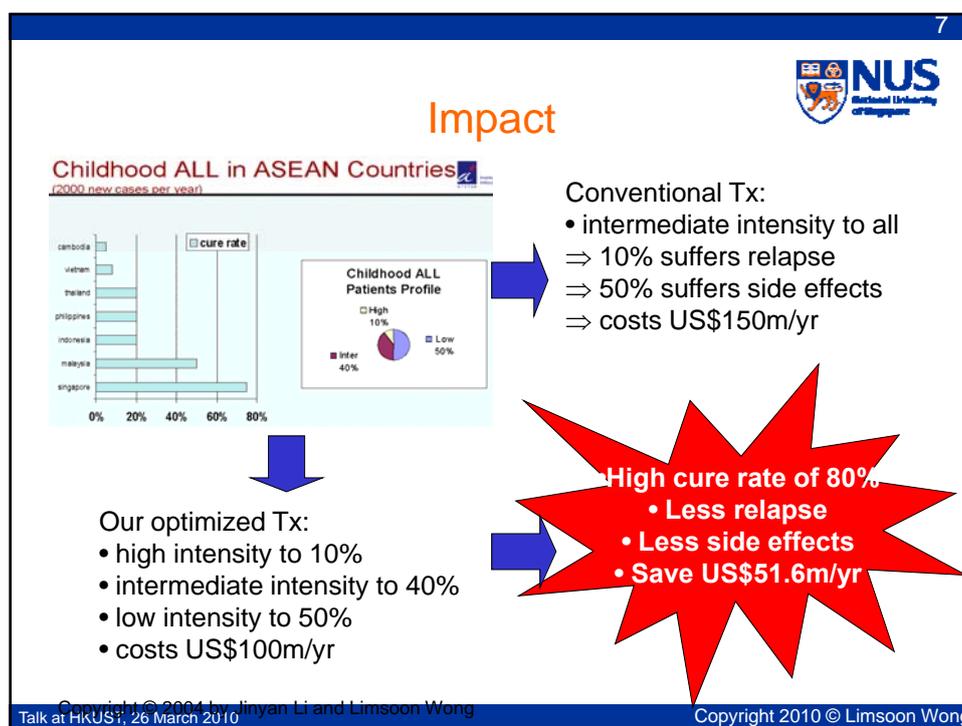
Copyright 2010 © Limsoon Wong

6



Talk at HKUST, 26 March 2010

Copyright 2010 © Limsoon Wong



8



Now, the bad news..

Talk at HKUST, 26 March 2010

Copyright 2010 © Limsoon Wong

9

 NUS
National University of Singapore

Percentage of Overlapping Genes

- **Low % of overlapping genes from diff expt in general**
 - Prostate cancer
 - Lapointe et al, 2004
 - Singh et al, 2002
 - Lung cancer
 - Garber et al, 2001
 - Bhattacharjee et al, 2001
 - DMD
 - Haslett et al, 2002
 - Pescatori et al, 2007

Datasets	DEG	POG
Prostate Cancer	Top 10	0.30
	Top 50	0.14
	Top100	0.15
Lung Cancer	Top 10	0.00
	Top 50	0.20
	Top100	0.31
DMD	Top 10	0.20
	Top 50	0.42
	Top100	0.54

Zhang et al, Bioinformatics, 2009

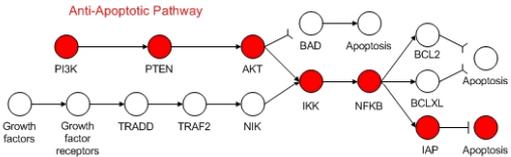
Talk at HKUST, 26 March 2010 Copyright 2010 © Limsoon Wong

10

 NUS
National University of Singapore

Gene Regulatory Circuits

Anti-Apoptotic Pathway



- Each disease subtype has underlying cause
- There is a unifying biological theme for genes that are truly associated with a disease subtype

- **Uncertainty in selected genes can be reduced by considering biological processes of the genes**
- **The unifying biological theme is basis for inferring the underlying cause of disease subtype**

Talk at HKUST, 26 March 2010 Copyright 2010 © Limsoon Wong

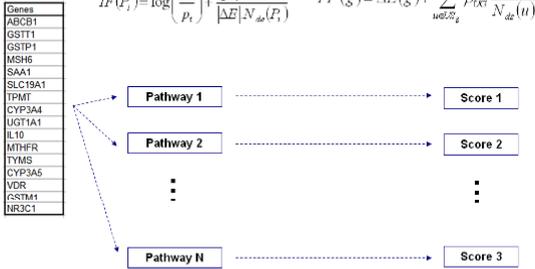
11

 **NUS**
National University of Singapore

Towards More Meaningful Genes

- **ORA**
 - Khatri et al
 - *Genomics*, 2002
- **FCS**
 - Pavlidis & Noble
 - *PSB* 2002
- **GSEA**
 - Subramanian et al
 - *PNAS*, 2005
- **Pathway Express**
 - Draghici et al
 - *Genome Res*, 2007

Gene Class Testing: Pathway Express

$$IF(P_i) = \log\left(\frac{1}{p_i}\right) + \frac{\sum PF(g)}{|\Delta E|N_{\Delta}(P_i)} \quad PF(g) = \Delta E(g) + \sum_{w \in \Delta(g)} \beta_{(w)} \frac{PF(w)}{N_{\Delta}(w)}$$


Draghici et al. *Genome Res*. 2007

Talk at HKUST, 26 March 2010 Copyright 2010 © Limsoon Wong

12

 **NUS**
National University of Singapore

All of these newer methods rely on gene group or pathway information.

But how good are the available sources of pathway information?

Talk at HKUST, 26 March 2010 Copyright 2010 © Limsoon Wong

Comparing Pathway Sources

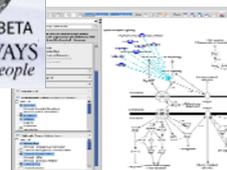


14

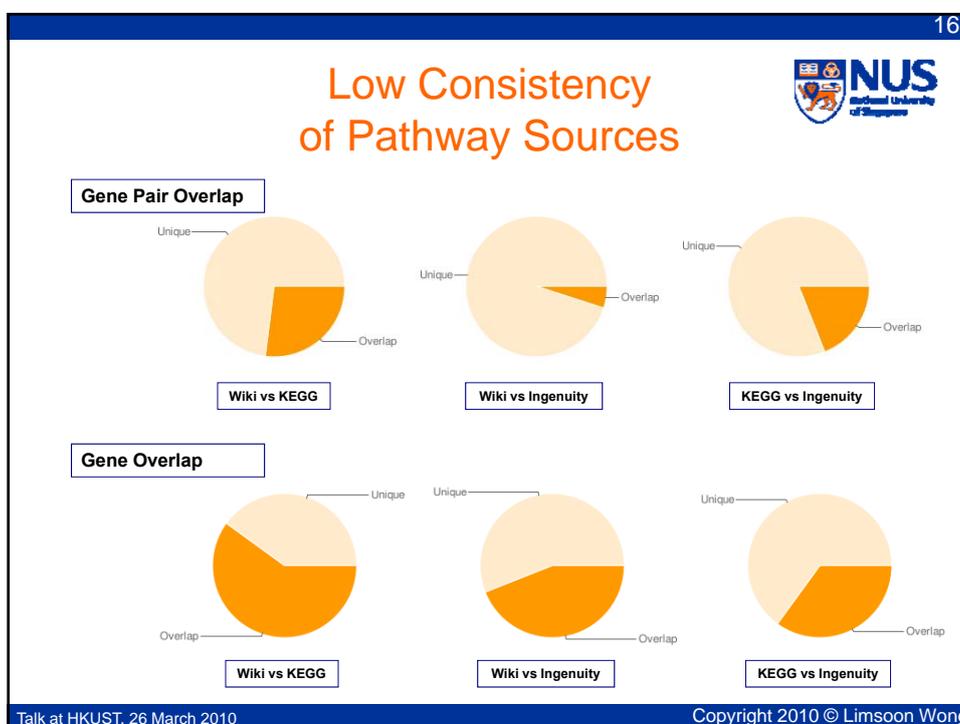
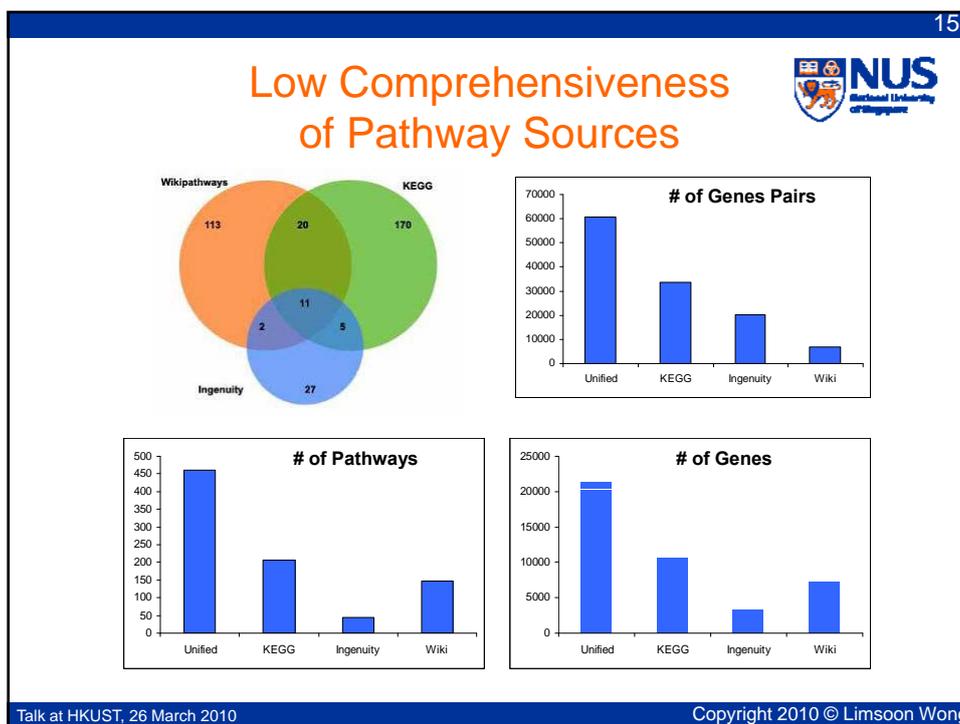
Data Sources



- **KEGG**
 - Curated by a single lab
 - Long famous history
 - Used by many people
- **Wikipathways**
 - Community effort
 - new curation model
- **Ingenuity**
 - Commercial effort
 - Used by many biopharma's



INGENUITY
S Y S T E M S



17

 NUS
National University of Singapore

Example: Apoptosis Pathway

Apoptosis Pathway			
	Wiki x KEGG	Wiki x Ingenuity	KEGG x Ingenuity
Gene Pair Count:	144 vs 172	144 vs 3557	172 vs 3557
Gene Count:	85 vs 80	85 vs 176	80 vs 176
Gene Overlap:	38	28	30
Gene % Overlap:	48%	33%	38%
Gene Pair Overlap:	23	14	24
Gene Pair % Overlap:	16%	10%	14%









Talk at HKUST, 26 March 2010
Copyright 2010 © Limsoon Wong

18

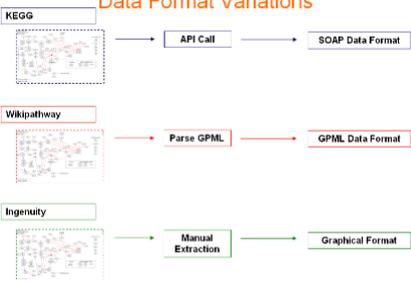
 NUS
National University of Singapore

Would Unifying Pathway Sources Help?

- **Incompatibility Issues!**

- Data extraction method variations
- Format variations
- Data differences
- Pathway name differences
- Gene/GenelD name differences

Data Format Variations



```

graph LR
    KEGG[KEGG] --> API[API Call] --> SOAP[SOAP Data Format]
    Wikipathway[Wikipathway] --> Parse[Parse GPML] --> GPML[GPML Data Format]
    Ingenuity[Ingenuity] --> Manual[Manual Extraction] --> Graphical[Graphical Format]
  
```

Talk at HKUST, 26 March 2010
Copyright 2010 © Limsoon Wong



The preceding analyses hide an intricate issue...

The same pathways in the different sources are often given different names.

So how do we even know two pathways are the same and should be compared / merged?

Intricacy of Pathway Matching



Possible Ways to Match Pathways



- **Match based on name**
 - Pathways w/ similar name should be the same pathway
 - But annotations are very noisy
 - ⇒ Likely to mismatch pathways?
 - ⇒ Likely to match too many pathways?

- **Are the followings good alternative approaches?**
 - Match based on overlap of genes
 - Match based on overlap of gene pairs

Matching Pathways by Name



- | | |
|--|--|
| <ul style="list-style-type: none"> • LCS procedure <ul style="list-style-type: none"> – Given pathway X in db A – Sort pathways in db B by “longest common substring” with X – Manually scan the ranked list to choose closest nomenclatural match | <ul style="list-style-type: none"> • Issue: Accuracy <ul style="list-style-type: none"> – When LCS says two pathways are the same one, are they really the same?
 • Issue: Completeness <ul style="list-style-type: none"> – When LCS says two pathways are different, are they really different? |
|--|--|

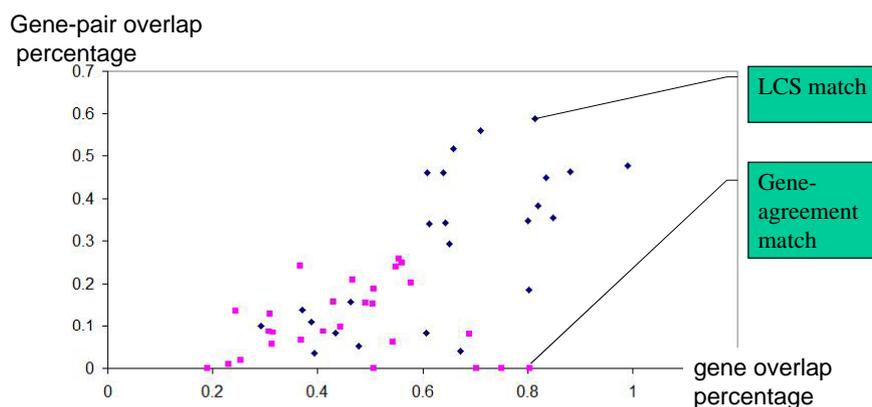
LCS vs Gene-Agreement Matching



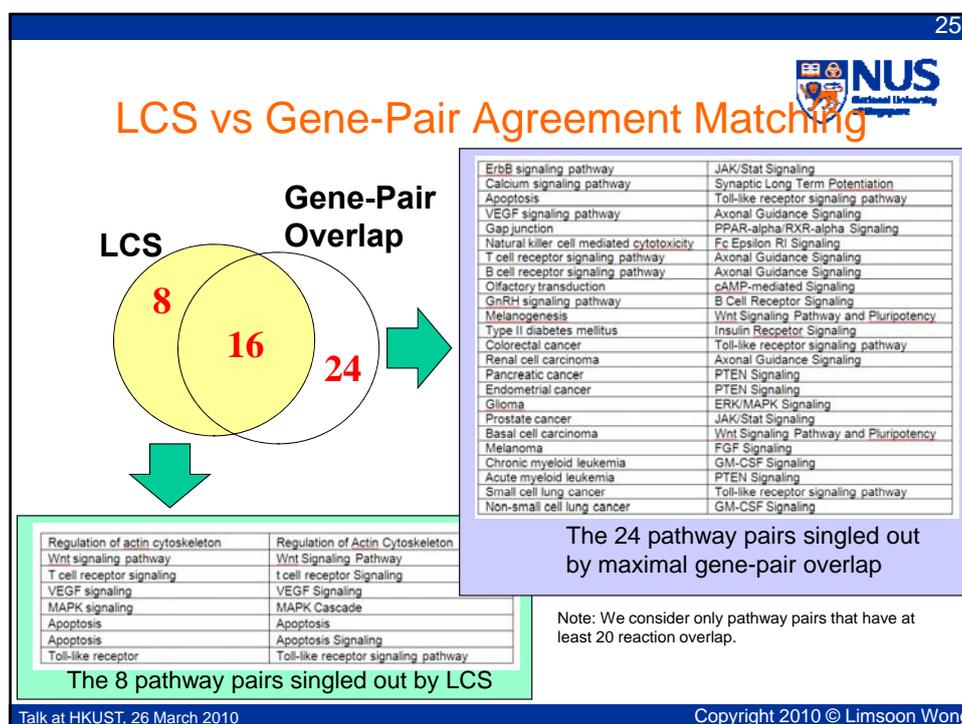
- **Accuracy**
 - 94% of LCS matches are in top 3 gene agreement matches
 - 6% of LCS matches not in top 3 of gene agreement matches; but their gene-pair agreement levels are higher
- **Completeness**
 - Let P_i be pathway in db A LCS cannot find match in db B
 - Let Q_i be pathway in db B with highest gene agreement to P_i
 - Gene-pair agreement of P_i - Q_i is much lower than pathway pairs matched by LCS

LCS is better than gene-agreement based matching!

LCS vs Gene-Agreement Matching



- **LCS consistently has higher gene-pair agreement**
 ⇒ **LCS is better than gene-agreement based matching!**



26



LCS vs Gene-Pair Agreement Matching

- **Gene-pair agreement match will miss when**
 - Pathway P in db A has few overlap with pathway P in db B due to incompleteness of db, even if pathway name matches perfectly!
 - Example: wnt signaling pathway, VEGF signaling pathway, MAPK signaling pathway, etc. in KEGG don't have largest gene-pair overlap w/ corresponding pathways in Wikipathways & Ingenuity

⇒ **Bad for getting a more complete unified pathway P**

Talk at HKUST, 26 March 2010 Copyright 2010 © Limsoon Wong

LCS vs Gene-Pair Agreement Matching



- Pathways having large gene-pair overlap are not necessarily the same pathways
 - Examples
 - “Synaptic Long Term Potentiation” in Ingenuity vs “calcium signalling” in KEGG
 - “PPAR-alpha/RXR-alpha Signaling” in Ingenuity vs “TGF-beta signaling pathway” in KEGG
- ⇒ Difficult to set correct gene-pair overlap threshold to balance against false positive matches

Bad Gene-Pair Agreement Matches



- | | |
|--|--|
| <ul style="list-style-type: none"> • “Synaptic Long Term Potentiation” in Ingenuity vs “calcium signalling” in KEGG <div style="background-color: #ffffcc; padding: 5px;"> <ul style="list-style-type: none"> • Calcium signaling pathway in KEGG describes general mechanism of external calcium signal transduction into cells • Calcium signal transduction can activate multiple downstream pathways, LTP is one of them <p>⇒ LTP in Ingenuity is only a downstream event of the calcium pathway in KEGG</p> </div> | <ul style="list-style-type: none"> • “PPAR-α/RXR-α signaling” in Ingenuity vs “TGF-β signaling” in KEGG <div style="background-color: #ffffcc; padding: 5px;"> <ul style="list-style-type: none"> • PPARα/RXRα plays essential roles in the regulation of cellular differentiation, development, metabolism, and tumorigenesis • TGF-β acts as antiproliferative factor in normal cells at early stages of oncogenesis <p>⇒ They are independent. The reason they are paired is that they have a mutual inhibition</p> </div> |
|--|--|



- Having found a good way to match up pathways in different datasources, we proceeded to build a big unified pathway db....

PathwayAPI
= KEGG
+ Wikipathways
+ Ingenuity

More Consistent Disease Subnetworks



31



But these methods still don't return the precise parts of a pathway that are significant...

- **ORA**
 - Khatri et al
 - Genomics, 2002
- **FCS**
 - Pavlidis & Noble
 - PSB 2002
- **GSEA**
 - Subramanian et al
 - PNAS, 2005
- **Pathway Express**
 - Draghici et al
 - Genome Res, 2007

Test whole gene group at a time

Test a node and its immediate neighbours

Talk at HKUST, 26 March 2010 Copyright 2010 © Limsoon Wong

32



The SNet Method

- **Group samples into type D and \neg D**
- **Extract & score subnetworks for type D**
 - Get list of genes highly expressed in most D samples
 - **These genes need not be differentially expressed!**
 - Segregate these genes into pathways
 - Locate largest connected components (ie., candidate subnetworks) from these pathway graphs
 - Score each subnetwork
- **Repeat the same on \neg D samples**
- **T-test on the two sets of scores to get significant subnetworks for D**

Talk at HKUST, 26 March 2010 Copyright 2010 © Limsoon Wong

33



SNet: Extract Subnetworks

List of Genes, GL

ARF
MDM2
MLC
Rac
Rho
P53
ATM
GSN
PAK
Rb
...

Gene List is split into its relevant pathways

Actin Cytoskeleton

VCL	MLC	PI4P5
Rho	PIX	MLCK
PAK	OSN	ROCK

Cell Cycle Regulation

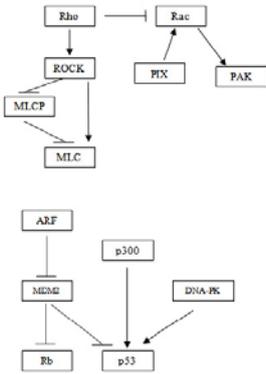
ARF	p53
MDM2	Rb
p300	ATM

Gene List is split into its relevant pathways

ROCK	→	MLC
Rho	→	Rac
PAK	→	PIX

ARF	→	MDM2
MDM2	→	p53
p300	→	p53

Matched genes are connected to form



Genes highly expressed in many type-D samples

Talk at HKUST, 26 March 2010
Copyright 2010 © Limsoon Wong

34



SNet: Score Subnetworks

Step 2: Subnetwork Scoring We assign a score vector $SN_{sn,d}^{u.score}$ with respect to phenotype d to each subnetwork sn within SN_{List} according to Equation 1.

$$SN_{sn,d}^{u.score} = \{SN_{sn,1,d}^{i.score}, SN_{sn,2,d}^{i.score}, \dots, SN_{sn,n,d}^{i.score}\} \quad (1)$$

Where n is the number of patients in phenotype d . The formula $SN_{sn,i,d}^{i.score}$ for the i^{th} patient (also the i^{th} element of this vector) is given by:

$$SN_{sn,i,d}^{i.score} = \sum_{j=1}^g G_{sn,i,d}^{j.score} \quad (2)$$

$G_{sn,i,d}^{j.score}$ refers to the score of the j^{th} gene (say, gene x) in the subnetwork sn for phenotype d . (This score $G_{sn,i,d}^{j.score}$ is given by Equation 3) and is simply given by:

$$G_{sn,i,d}^{j.score} = k/n \quad (3)$$

Where k is the number of patients of phenotype d who has gene x highly expressed (top $\alpha\%$) and n is the total number of patients of phenotype d . The entire Step 2 is repeated for the other disease phenotype $-d$, giving us the score vectors, $SN_{sn,d}^{u.score}$ and $SN_{sn,-d}^{u.score}$ for the same set of connected components. The t-test is finally calculated between these two vectors, creating a final t-score for each subnetwork sn within SN_{List} .

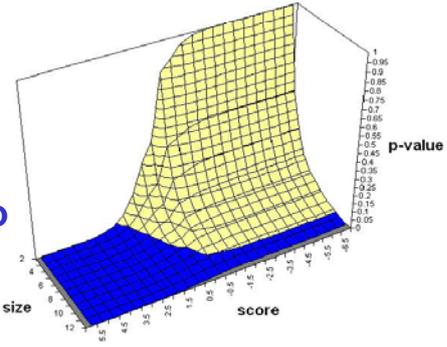
Talk at HKUST, 26 March 2010
Copyright 2010 © Limsoon Wong

35

 NUS
National University
of Singapore

SNet: Significant Subnetworks

- Randomize patient samples many times
- Get t-score for subnetworks from the randomizations
- Use these t-scores to establish null distribution
- Filter for significant subnetworks from real samples



Talk at HKUST, 26 March 2010

Copyright 2010 © Limsoon Wong

36

 NUS
National University
of Singapore

Let's see whether SNet gives us subnetworks that are

- (i) more consistent between datasets of the same types of disease samples
- (ii) larger and more meaningful

Talk at HKUST, 26 March 2010

Copyright 2010 © Limsoon Wong

Recall Examples from “Bad News”



- **Low % of overlapping genes from diff expt in general**
 - Prostate cancer
 - Lapointe et al, 2004
 - Singh et al, 2002
 - Lung cancer
 - Garber et al, 2001
 - Bhattacharjee et al, 2001
 - DMD
 - Haslett et al, 2002
 - Pescatori et al, 2007

Datasets	DEG	POG
Prostate Cancer	Top 10	0.30
	Top 50	0.14
	Top100	0.15
Lung Cancer	Top 10	0.00
	Top 50	0.20
	Top100	0.31
DMD	Top 10	0.20
	Top 50	0.42
	Top100	0.54

Zhang et al, Bioinformatics, 2009

Better Subnetwork Overlap



Table 1. Table showing the percentage overlap significant subnetworks between the datasets. Each row refers to a separate disease (as indicated in the first column). Each disease is tested against two datasets depicted in the second and third column. The overlap percentages refer to the pathway overlaps obtained from running SNet (column 4) and GSEA (column 5) The actual number of overlaps are parenthesized in the same columns.

Disease	Dataset 1	Dataset 2	SNet	GSEA
Leuk	Golub	Armstrong	83.3% (20)	0.0% (0)
Subtype	Ross	Yeoh	47.6% (10)	23.1% (6)
DMD	Haslett	Pescatori	58.3% (7)	55.6% (10)
Lung	Bhatt	Garber	90.9% (9)	0.0% (0)

- **For each disease, take significant subnetworks from one dataset and see if it is also significant in the other dataset**



Better Gene Overlaps

Table 2. Table showing the number and percentage of significant overlapping genes. γ refers to the number of genes compared against and is the number of unique genes within all the significant subnetworks of the disease datasets. The percentages refer to the percentage gene overlap for the corresponding algorithms.

Disease	γ	SNet	GSEA	SAM	t-test
Leuk	84	91.3%	2.4%	22.6%	14.3%
Subtype	75	93.0%	4.0%	49.3%	57.3%
DMD	45	69.2%	28.9%	42.2%	20.0%
Lung	65	51.2%	4.0%	24.6%	26.2%

- For each disease, take significant subnetworks extracted independently from both datasets and see how much their genes overlap



Larger Subnetworks

Table 3. Table comparing the size of the subnetworks obtained from the t-test and from SNet. The first column shows the disease and the second column shows the number of genes which comprised of the subnetworks. The third and fourth column depicts the number of genes present within each subnetwork for the t-test and SNet respectively. So for instance in the leukemia dataset, we have 8 subnetworks with size 2 genes, 1 subnetwork with size 3 genes for the t-test. For SNet, we have 2 subnetworks with size 5 genes, 3 subnetworks with size 6 genes, 2 subnetworks with size 7 genes and 1 subnetwork with a size of ≥ 8 genes

Disease	γ	Num Genes (t-test)				Num Genes (SNet)			
		2	3	4	5	5	6	7	≥ 8
Leuk	84	8	1	0	0	2	3	2	1
Subtype	75	5	1	1	1	1	0	1	6
DMD	45	3	1	0	0	1	0	0	5
Lung	65	3	2	1	0	5	3	0	1

Remarks



42

What have we learned?



- **Significant lack of concordance betw db's**
 - Level of consistency for genes is 0% to 88%
 - Level of consistency for genes pairs is 0%-61%
 - Most db contains less than half of the pathways in other db's
- **Matching pathways by name is better than matching by gene overlap or gene-pair overlap**
- **SNet method yields more consistent and larger disease subnetworks**

Acknowledgements



Donny Soh



Difeng Dong



Yike Guo

- **A*STAR AIP scholarship**
- **A*STAR SERC PSF grant**

