

Consistency & Comprehensiveness of Pathway Databases

Limsoon Wong
23 February 2010

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



2

Plan



- **Past successes in gene expression analysis**
- **Towards more meaningful genes**
- **Issues on pathway sources**
 - Comprehensiveness
 - Consistency
 - Compatibility
 - Matching pathways in different sources

Past Success: An Example

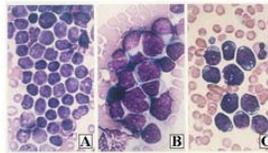


4

Childhood Acute Lymphoblastic Leukemia



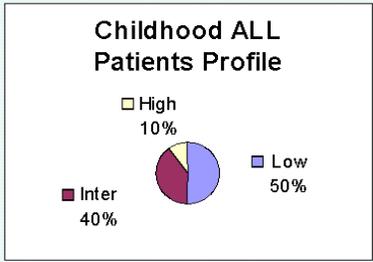
- Major subtypes: T-ALL, E2A-PBX, TEL-AML, BCR-ABL, MLL genome rearrangements, Hyperdiploid >50
- Diff subtypes respond differently to same Tx
- Over-intensive Tx
 - Development of secondary cancers
 - Reduction of IQ
- Under-intensive Tx
 - Relapse
- The subtypes look similar
- Conventional diagnosis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



5



Patient Profiles & Treatment Costs



Childhood ALL Patients Profile

- High 10%
- Inter 40%
- Low 50%

- **Treatment for childhood ALL over 2 yrs**
 - Intermediate intensity: US\$60k
 - Low intensity: US\$36k
 - High intensity: US\$72k
- **Treatment for relapse: US\$150k**
- **Cost for side-effects: Unquantified**

- **2000 new cases a year in ASEAN countries**

Talk at Korea-Singapore Workshop on Bioinformatics and NLP, KAIST, Feb 2010 Copyright 2010 © Limsoon Wong

6



Why not high/low intensity to everyone?

- **High-intensity Tx**
 - Over intensive for 90% of patients, thus a lot more side effects
 - US\$144m (US\$72k * 2000) for high-intensity tx

- **Low-intensity Tx**
 - Under intensive for 50% of patients, thus a lot more relapse
 - US\$72m (US\$36k * 2000) for low-intensity tx
 - US\$150m (US\$150k * 2000 * 50%) for relapse tx

⇒ **Total US\$144m/yr plus un-quantified costs for dealing with side effects**

⇒ **Total US\$222m/yr**

Talk at Korea-Singapore Workshop on Bioinformatics and NLP, KAIST, Feb 2010 Copyright 2010 © Limsoon Wong

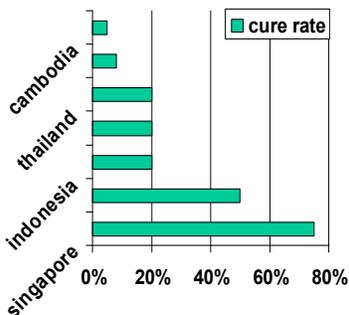
7



Current Situation

- **Intermediate intensity conventionally applied in ASEAN countries**

- **Over intensive for 50% of patients, thus more side effects**
- **Under intensive for 10% of patients, thus more relapse**
- **US\$120m (US\$60k * 2000) for intermediate intensity tx**
- **US\$30m (US\$150k * 2000 * 10%) for relapse tx**
- **Total US\$150m/yr plus unquantified costs for dealing with side effects**

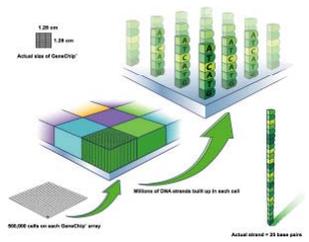


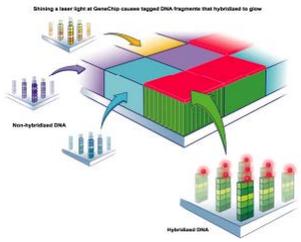
Talk at Korea-Singapore Workshop on Bioinformatics and NLP, KAIST, Feb 2010 Copyright 2010 © Limsoon Wong

8

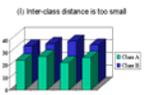


Single-Test Platform of Microarray & Machine Learning

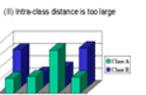




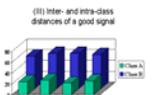
(i) Inter-class distance is too small



(ii) Intra-class distance is too large



(iii) Inter- and intra-class distances of a good signal



	Pos/neg	Negative	Pairs	In/Avg	Diff	Ass	Cell
00-0688-U00-C586-U00-0588-U00-C586-U00-0588-U0-Descriptors							
APF&Mur1	5	2	19	207.5	A	M16762	Mouse int
APF&Mur2	5	2	19	554.2	A	M37827	Mouse int
APF&Mur3	4	2	19	308.6	A	M29522	Mus mus
APF&Mur4	1	5	19	141	A	M95849	Mus mus
APF&BioE	13	1	19	6540.6	P	J34423	E coli bioE
APF&BioF	15	0	19	12662.4	P	J34423	E coli bioF
APF&BioG	12	0	19	9716.5	P	J34423	E coli bioG
APF&BioH	17	0	19	25942.5	P	J34423	E coli bioH
APF&BioI	16	0	20	29826.5	P	J34423	E coli bioI
APF&BioJ	17	0	19	25785.2	P	J34423	E coli bioJ
APF&BioK	19	0	20	143173.2	P	J34423	E coli bioK
APF&BioL	20	0	20	260326.6	P	M35455	Bacterioph
APF&BioM	20	0	20	401741.8	P	M35455	Bacterioph
APF&BioN	7	4	18	485	A	J34423	E coli bioN
APF&BioO	5	4	18	313.7	A	J34423	E coli bioO
APF&BioP	7	6	20	-1016.2	A	J34423	E coli bioP

Talk at Korea-Singapore Workshop on Bioinformatics and NLP, KAIST, Feb 2010 Copyright 2010 © Limsoon Wong

9



Individual Gene Testing

Fold Change

$$FC_{ratioi} = \frac{x_i}{y_i}$$

$$FC_{diffi} = x_i - y_i$$

x – Microarray value after drug
y – Microarray value before drug
i – Gene

T-test

$$T_i = \frac{\hat{x}_i - \hat{y}_i}{s_i}$$

$$T_i = \frac{\hat{x}_i - \hat{y}_i}{s_i + s_o}$$

$$T_i = \frac{\hat{x}_i - \hat{y}_i}{\sqrt{Bs^2 + (1 - B)s_i^2}}$$

x – Log2 value of treatment
y – Log2 value of control
s – Standard error
i – Gene

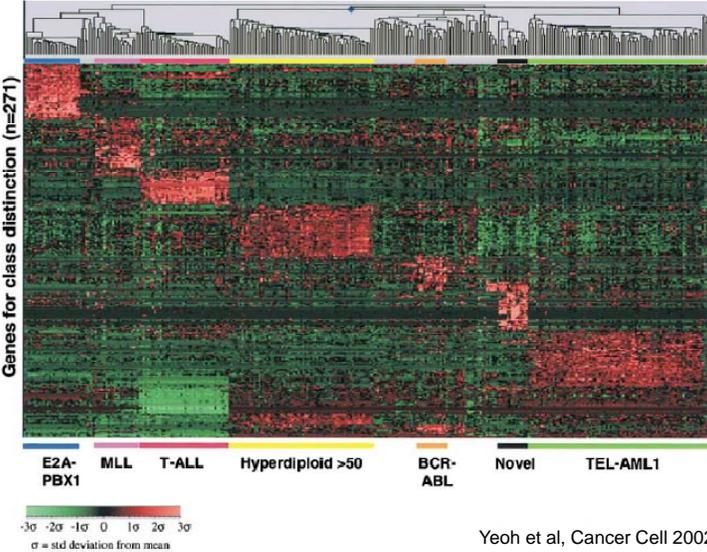
Golub et al, Science, 1999

Talk at Korea-Singapore Workshop on Bioinformatics and NLP, KAIST, Feb 2010
Copyright 2010 © Limsoon Wong

10



Diagnostic ALL BM samples (n=327)



Yeoh et al, Cancer Cell 2002

Talk at Korea-Singapore Workshop on Bioinformatics and NLP, KAIST, Feb 2010
Copyright 2010 © Limsoon Wong



Exploit Invariant Gene Expr Profiles

- Low intensity applied to 50% of patients
- Intermediate intensity to 40% of patients
- High intensity to 10% of patients
- US\$36m (US\$36k * 2000 * 50%) for low intensity
- US\$48m (US\$60k * 2000 * 40%) for intermediate intensity
- US\$14.4m (US\$72k * 2000 * 10%) for high intensity
- ⇒ **Reduced side effects**
- ⇒ **Reduced relapse**
- ⇒ **75-80% cure rates**
- **Total US\$98.4m/yr**
- ⇒ **Save US\$51.6m/yr**

Yeoh et al, Cancer Cell 2002

Toward More Meaningful Genes





Percentage of Overlapping Genes

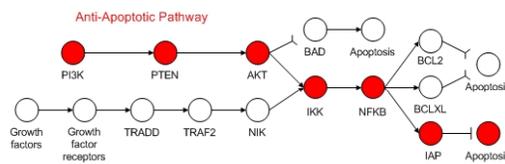
- Low % of overlapping genes from diff expt in general
 - Prostate cancer
 - Lapointe et al, 2004
 - Singh et al, 2002
 - Lung cancer
 - Garber et al, 2001
 - Bhattacharjee et al, 2001
 - DMD
 - Haslett et al, 2002
 - Pescatori et al, 2007

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



Gene Regulatory Circuits



- 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

17



Towards More Meaningful Genes

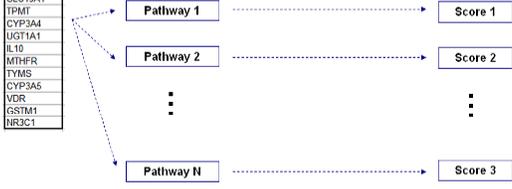
- **ORA**
 - Khatri et al
 - Genomics, 2002
- **FCS**
 - Pavlidis & Noble
 - PSB 2002
- **GSEA**
 - Subramanian et al
 - PNAS, 2005
- **Pathway Express**
 - Draghici et al
 - Genome Res, 2007

Gene Class Testing: Pathway Express

Gene
ABCB1
GSTT1
GSTP1
MSH6
SAA1
SLC19A1
TPMT
CYP3A4
UGT1A1
IL10
MTHFR
TYMS
CYP3A5
VDR
GSTM1
NR3C1

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

$$PF(g) = \Delta E(g) + \sum_{u \in \mathcal{G}_g} \beta_{u, \mathcal{G}_g} \frac{PF(u)}{N_{gs}(u)}$$



Draghici et al. Genome Res. 2007

Talk at Korea-Singapore Workshop on Bioinformatics and NLP, KAIST, Feb 2010
Copyright 2010 © Limsoon Wong

18



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

But how good are the available sources of pathway information?

Talk at Korea-Singapore Workshop on Bioinformatics and NLP, KAIST, Feb 2010
Copyright 2010 © Limsoon Wong

Issues on Pathway Sources

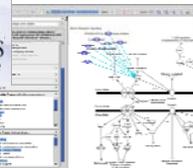


20

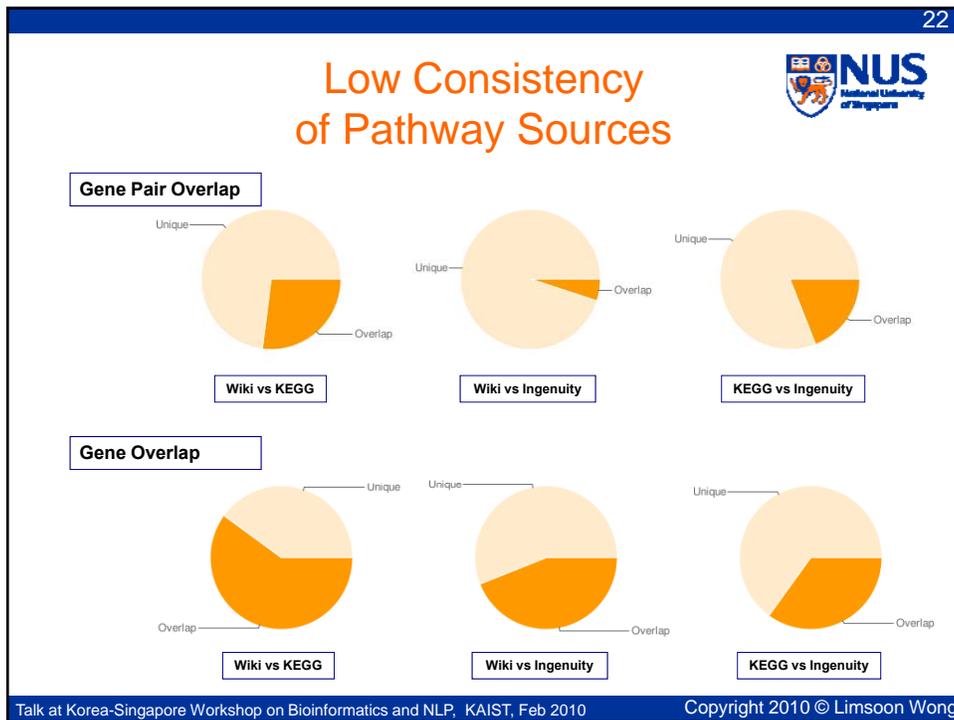
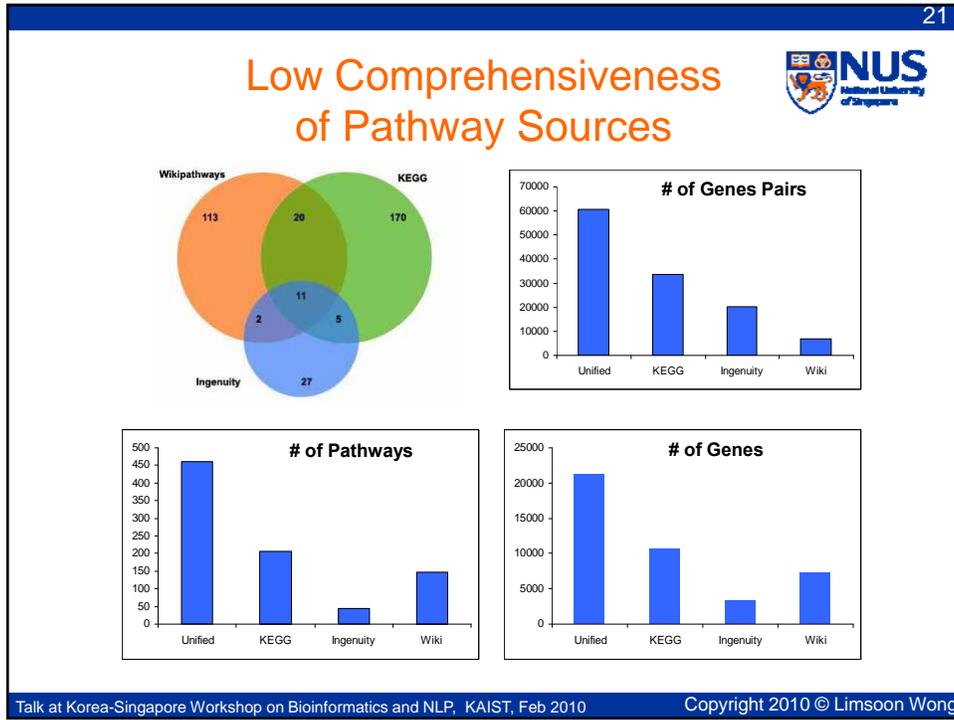
Data Sources



- **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





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%



GenMAPP



PharmGKB

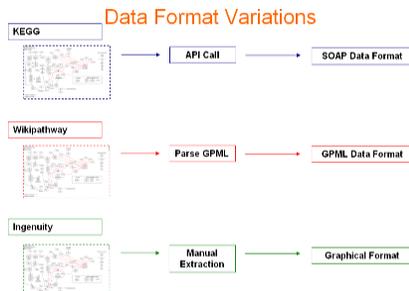
INGENUITY SYSTEMS



Would Unifying Pathway Sources Help?

- Incompatibility Issues!**

- Data extraction method variations



- Format variations
- Data differences
- Pathway name differences
- Gene/GeneID name differences

Intricacy of Pathway Matching



26

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

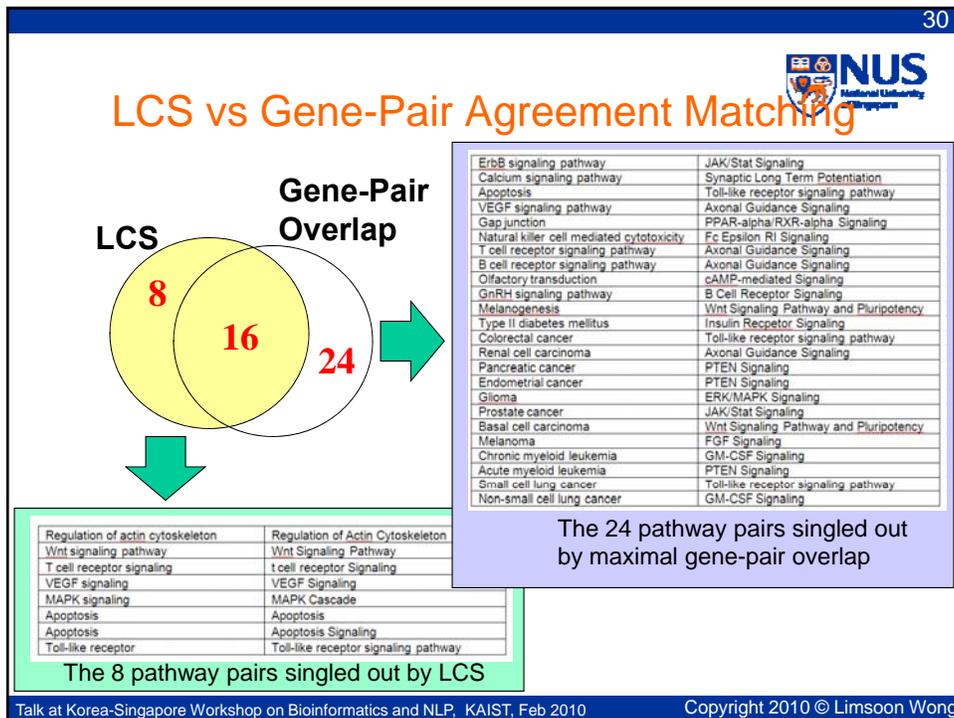
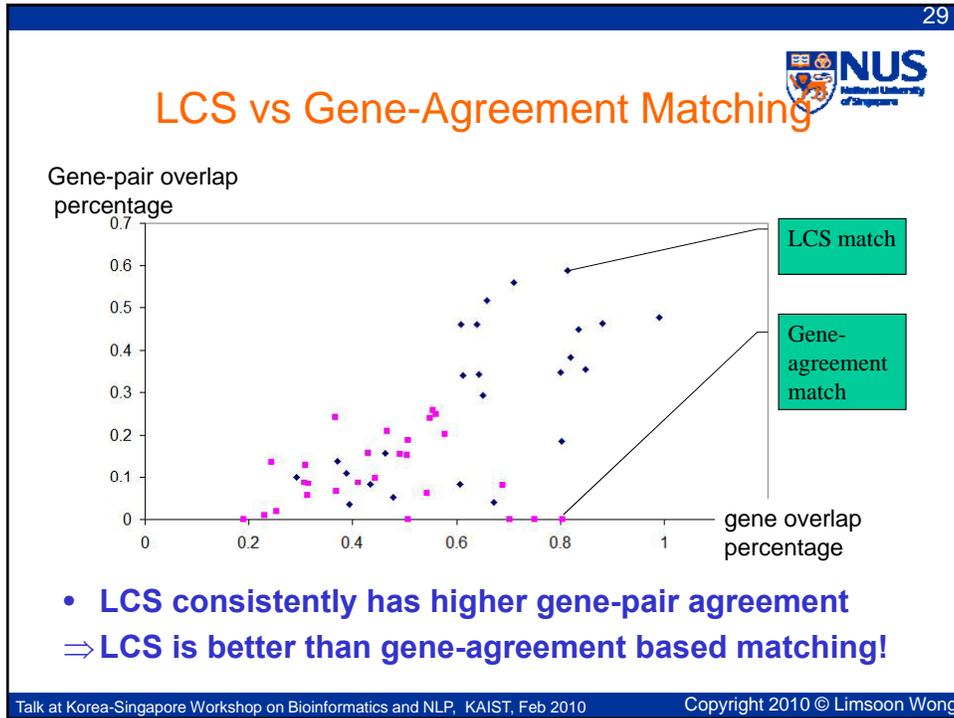
- **LCS procedure**
 - Given pathway X in db A
 - Sort pathways in db B by “longest common substring” with X
 - Manually scan the ranked list to choose closest nomenclatural match
- **Issue: Accuracy**
 - When LCS says two pathways are the same one, are they really the same?
- **Issue: Completeness**
 - 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-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**

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

- “Synaptic Long Term Potentiation” in Ingenuity vs “calcium signalling” in KEGG
- “PPAR- α /RXR- α signaling” in Ingenuity vs “TGF- β signaling” in KEGG

- 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
- ⇒ LTP in Ingenuity is only a downstream event of the calcium pathway in KEGG

- PPAR α /RXR α 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
- ⇒ They are independent. The reason they are paired is that they have a mutual inhibition

Remarks



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%
- **Significant lack of comprehensiveness**
 - 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**



Acknowledgements



Donny Soh



Difeng Dong



Yike Guo

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

