#### **Data Integration**



## Metastructures: systems approach to determine genome annotation.



Qiu Y et al. Genome Res. 2010;20:1304-1311

# Why do we need to integrate various types of omic data?

- Get a consensus results (reduce false positive rate)
- Focusing on one type of data may miss an obvious signal





#### Integration of omics data sets

Genomics	Transcriptomics	Proteomics	Metabolomics	Protein–DNA interactions	Protein-protein interactions	Fluxomics	Phenomics
Genomics (sequence annotation)	<ul> <li>ORF validation</li> <li>Regulatory element identification<sup>74</sup></li> </ul>	• SNP effect on protein activity or abundance	• Enzyme annotation	• Binding-site identification <sup>75</sup>	• Functional annotation <sup>79</sup>	• Functional annotation	<ul> <li>Functional annotation<sup>71,103</sup></li> <li>Biomarkers<sup>125</sup></li> </ul>
	Transcriptomics (microarray, SAGE)	• Protein: transcript correlation <sup>20</sup>	• Enzyme annotation <sup>109</sup>	• Gene-regulatory networks <sup>76</sup>	<ul> <li>Functional annotation<sup>89</sup></li> <li>Protein complex identification<sup>82</sup></li> </ul>		• Functional annotation <sup>102</sup>
		Proteomics (abundance, post- translational	• Enzyme annotation <sup>99</sup>	• Regulatory complex identification	• Differential complex formation	• Enzyme capacity	• Functional annotation
		modification)	Metabolomics (metabolite abundance)	• Metabolic- transcriptional response		<ul> <li>Metabolic pathway bottlenecks</li> </ul>	<ul> <li>Metabolic flexibility</li> <li>Metabolic engineering<sup>109</sup></li> </ul>
				Protein–DNA interactions (ChlP–chip)	• Signalling cascades <sup>89,102</sup>		• Dynamic network responses <sup>84</sup>
					Protein–protein interactions (yeast 2H,		<ul> <li>Pathway identification activity<sup>89</sup></li> </ul>
					COAP-MS)	Fluxomics (isotopic tracing)	<ul> <li>Metabolic engineering</li> </ul>
							Phenomics (phenotype arrays, RNAi screens, synthetic lethals)

Nature Reviews Molecular Cell Biology, 7:198–210, 2006.

#### Multi-omic data integration



H. Latif et al. PLoS Genetics 2013

#### Challenges

- 1. Data Pre-processing
- 2. High Dimensionality
- 3. Multiple Testing for Marker Selection
- 4. Data Integration
- 5. Validation of the Prediction Model

#### Challenge #1: Data Pre-processing

 Peak Alignment for different platforms



- Normalization among various types of data
  - Why? Remove systematic bias in the data
  - Normalization within the platform makes data comparable across samples



#### Challenge # 2: High Dimensionality

# of subjects << # of variables</pre>



- Genes variants: 5000 peaks
- Gene Expression (RNA-seq): 22,000 probe sets
- Protein-DNA interactions (ChIP-seq): 2, 000 peaks
- Protein interactions: 100,000,000

#### Challenge #4: Data integration (How?)





Identification and quantitative comparison of genetic elements for transcription and translation initiation.

H. Latif et al. PLoS Genetics 2013

#### Example 2

#### $\sigma$ -factor network in E. coli.





hypB

hypC

hyp

B. Cho et al. BMC Biology 2014

#### Example 3

#### Integrative Modeling Defines the Nova Splicing-Regulatory Network and Its Combinatorial Controls

Chaolin Zhang et al. Science 2010

#### Data Integration

- More and more diverse "omics" data exist
- "It is essential to integrate various kinds of biological information and large-scale omics data sets through systematic analysis" with statistically rigorous and physically sound models
- Bayesian Network
- CLIP-seq (sequencing) + TF binding site (bioinformatics) + expression profile (microarry) +evolutionary signature.

## **Bayesian Network**

- A probabilistic graphical model that represents a set of random variables and their conditional dependencies via a directed acyclic graph.
- used to represent causal relationships.



#### Nova-regulated Alternative Splicing

- Nova proteins are a family of neuron-specific alternative splicing factors. (Ule *et al. Nature* 2006)
- The Nova protein binds to pre-message RNA at a binding site of "YCAY" clusters.
- Nova binding to an exonic YCAY cluster changed the protein complexes assembled on pre-mRNA, blocking U1 snRNP binding and exon inclusion.
- Nova binding to an intronic YCAY cluster enhanced spliceosome assembly and exon inclusion.



#### Data sources



- **HITS-CLIP** (CLIP-seq). Study Protein-RNA binding by crosslinking between RNA and the protein, followed by immunoprecipitation and high-throughput sequencing.
- Genome-wide searching of YCAY motif. Bioinformatics approach.
- Microarray data compared WT and Nova knockout.
- Evolution signature. Conserved Alternating Splicing between human and rats.



#### Estimated Conditional prob.



Predicted miRNA binding site

#### **Bayesian Model**



## Regulation of Nova binding position



#### **Bayesian Model**



#### Splicing changes





#### **Bayesian Model**

![](_page_24_Figure_1.jpeg)

#### Predicted Nova-regulated targets

 13,357 annotated cassette exons

![](_page_25_Figure_2.jpeg)

#### Novel Nova-regulated targets

- Besides AS from database, searched novel exons with high sequence conservations.
- Additional 76 novel exons as Nova targets

![](_page_26_Figure_3.jpeg)

#### **Prediction Performance**

![](_page_27_Figure_1.jpeg)

#### **Reduced Bayesian Network**

в

![](_page_28_Figure_1.jpeg)

![](_page_28_Figure_2.jpeg)

• Clip Data Only

![](_page_28_Figure_4.jpeg)

![](_page_28_Figure_5.jpeg)

#### **Experimental Validation**

![](_page_29_Figure_1.jpeg)

#### **Exon Inclusion**

#### Two more Casset Exon Cases

![](_page_30_Figure_1.jpeg)

**Exon Exclusion** 

#### Other examples

![](_page_31_Figure_1.jpeg)

TACA Exon Inclusion

![](_page_31_Figure_3.jpeg)

#### **Conservation regions**

![](_page_32_Figure_1.jpeg)

#### Functions of Nova targets

- Nova regulates alternative splicing of transcripts encoding synaptic proteins.
- Go-term analysis and KEGG metabolic pathway analysis confirmed this.
- It is unclear how Nova-regulated AS might effect the interactions between those synaptic proteins.
- Protein annotations revealed that about half
   Nova target transcripts encoded
   phosphoproteins.

#### GO-term enrichment of Nova targets

GO Term	Gene	%	P-Value	Fold Enrichment	Benjamini FDR
Biological process					
GO:0016043~cellular component organization	93	26.05	4.02E-12	2.02	6.93E-09
GO:0007399~nervous system development	53	14.85	1.01E-09	2.48	8.72E-07
GO:0032989~cellular component morphogenesis	31	8.68	2.32E-09	3.56	1.33E-06
GO:0030030~cell projection organization	30	8.40	3.55E-09	3.60	1.53E-06
GO:0007268~synaptic transmission GO:0048667~cell morphogenesis involved in	22	6.16	8.39E-09	4.64	2.90E-06
neuron differentiation	22	6.16	1.18E-08	4.55	3.40E-06
GO:0051179-localization	104	29.13	1.63E-08	1.66	4.03E-06
GO:0000902~cell morphogenesis	28	7.84	1.64E-08	3.57	3.53E-06
GO:0019226~transmission of nerve impulse	24	6.72	2.47E-08	4.01	4.73E-06
GO:0007154~cell communication	34	9.52	4.41E-08	2.93	7.60E-06
Cellular component					
GO:0045202~synapse	38	10.64	1.64E-15	4.85	5.00E-13
GO:0044459~plasma membrane part	82	22.97	2.06E-15	2.51	3.16E-13
GO:0042995~cell projection	47	13.17	3.24E-14	3.62	3.24E-12
GO:0030054~cell junction	42	11.76	3.53E-14	4.00	2.65E-12
GO:0005856~cytoskeleton	61	17.09	6.36E-12	2.60	3.81E-10
GO:0005886~plasma membrane	104	29.13	6.82E-11	1.84	3.41E-09
GO:0042734~presynaptic membrane	11	3.08	1.07E-09	14.80	4.61E-08
GO:0016323~basolateral plasma membrane	19	5.32	3.11E-09	5.83	1.17E-07
GO:0044456~synapse part	23	6.44	5.65E-09	4.55	1.88E-07
GO:0043005~neuron projection	25	7.00	8.93E-09	4.09	2.68E-07
Molecular function					
GO:0005515~protein binding	198	55.46	9.43E-15	1.49	4.50E-12
GO:0008092~cytoskeletal protein binding	35	9.80	2.57E-10	3.51	6.14E-08
GO:0003779~actin binding	23	6.44	9.83E-07	3.40	1.56E-04
GO:0030695~GTPase regulator activity GO:0060589~nucleoside-triphosphatase regulator	26	7.28	1.17E-06	3.06	1.39E-04
activity	26	7.28	1.61E-06	3.01	1.54E-04

#### Pathway enrichment of Nova targets

KEGG pa	thway	Gene count	Fold Enrichment	Benjamini FDR	Genes
mmu0402 Calcium s	0 ignaling pathway	17	3.5	0.001	Atp2b1, Atp2b2, Cacna1c, Cacna1d, Cacna1b, Cacna1g, Camk2a, Camk2g, Camk2b, Grin1, Gnas, Picb4, Ppp3cb, Ppp3cc, Ryr2, Sic8a1, Erbb4
mmu0472 Long-term	0 potentiation	10	5.3	0.003	Cacna1c, Camk2a, Camk2g, Camk2b, Gria2, Grin1, Plcb4, Ppp1r12a, Ppp3cb, Ppp3cc
mmu0451 Cell adhe:	4 sion molecules	12	4.1	0.003	Alcam, Cadm1, Cadm3, Mpzl1, Neo1, Nrxn3, Nfasc, Nfasc, Ptprf, Ptprm, Nign1, Nrcam,Nrxn1
(CAMS) mmu0452 Adherens	0 junction	10	4.4	0.006	Actn4, Baiap2, Ctnna2, Ctnnd1, Pard3, Smad2, Smad4, Ptprf, Ptprm,Sorbs1
mmu0436 Axon guid	0 ance	13	3.3	0.006	Ablim1, Cxcl12, Dcc, Epha5, Efna5, Ablim2, Ntng1, Pak3, Ppp3cb, Ppp3cc, Arhgef12, Robo2, Unc5c
mmu0491 GnRH sig	2 naling pathway	10	3.6	0.017	Cacna1c, Cacna1d, Camk2a, Camk2g, Camk2b, Gnas, Mapk8, Mapk9, Map2k4, Picb4
mmu0431 Wnt signa	0 ling pathway	12	2.8	0.032	Apc, Camk2a, Camk2g, Camk2b, Smad2, Smad4, Mapk8, Mapk9, Plcb4, Porcn, Ppp3cb, Ppp3cc
mmu0493 Type II dia	0 abetes mellitus	6	5.2	0.048	Cacna1c, Cacna1d, Cacna1b, Cacna1g, Mapk8, Mapk9
mmu0426 Cardiac m	0 ouscle contraction	7	4.2	0.049	Tpm2, Cacna1d, Cacna1c, Tpm1, Ryr2, Slc8a1, Tpm3
mmu0401 ErbB sign	2 aling pathway	8	3.3	0.074	Camk2a, Camk2g, Camk2b, Mapk8, Mapk9, Map2k4, Pak3, Erbb4
mmu0453 Tight junc	0 tion	9	2.7	0.01	Actn4, Cask, Ctnna2, Pard3, Epb4.1, Epb4.111, Epb4.112, Epb4.113, Magi1

#### Nova targets - phosphoproteins

![](_page_36_Figure_1.jpeg)

## **Applications of Bayesian Network**

Can we apply Bayesian Network into our research?

- Next generation sequencing data, such as RNA-seq, Chip-seq etc.
- Microarray data
- Motif data, for example, TF binding sites, miRNA sequences etc.
- Genome sequence data, Ath, Maize, Rice, Soybean etc.

#### Summary

- Recent technological advances present challenging and interesting biological data at molecular level.
- Statistics and multivariate analysis play an important role in understanding and extracting knowledge from these type of data.
- Integrative analysis is even more challenging and we presented some solutions to these challenges. There is plenty of room for improvement.