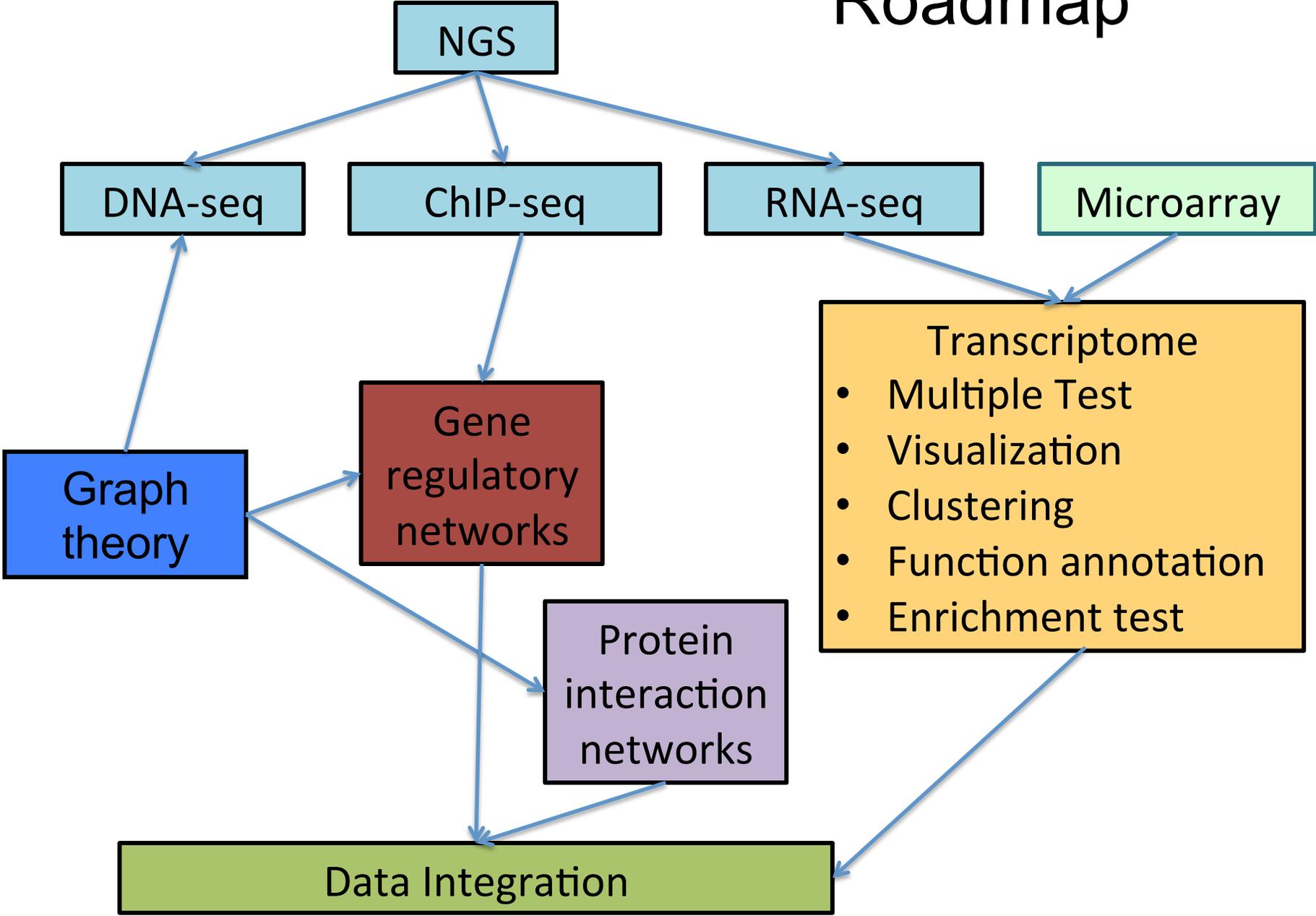
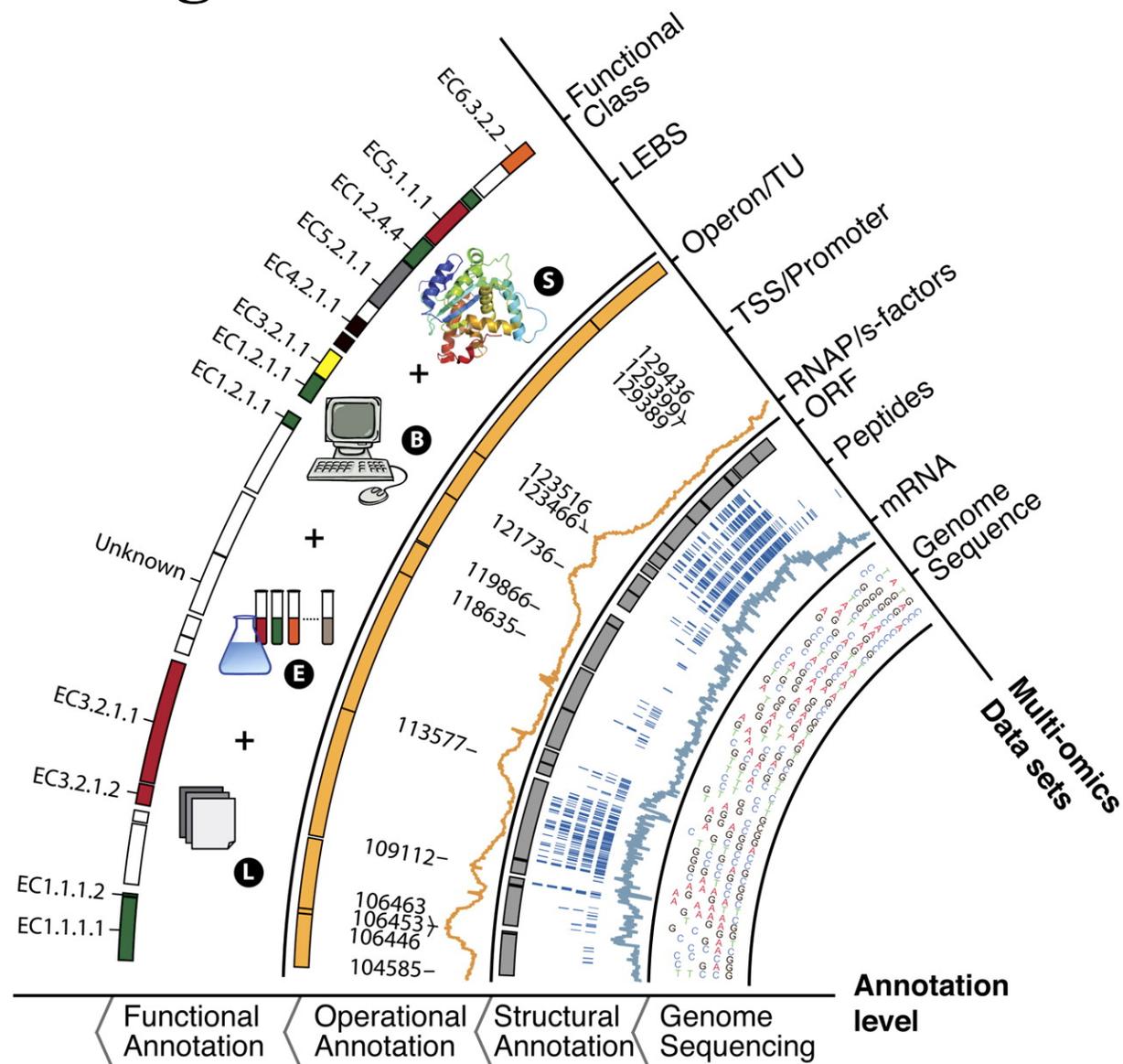


Data Integration

Roadmap

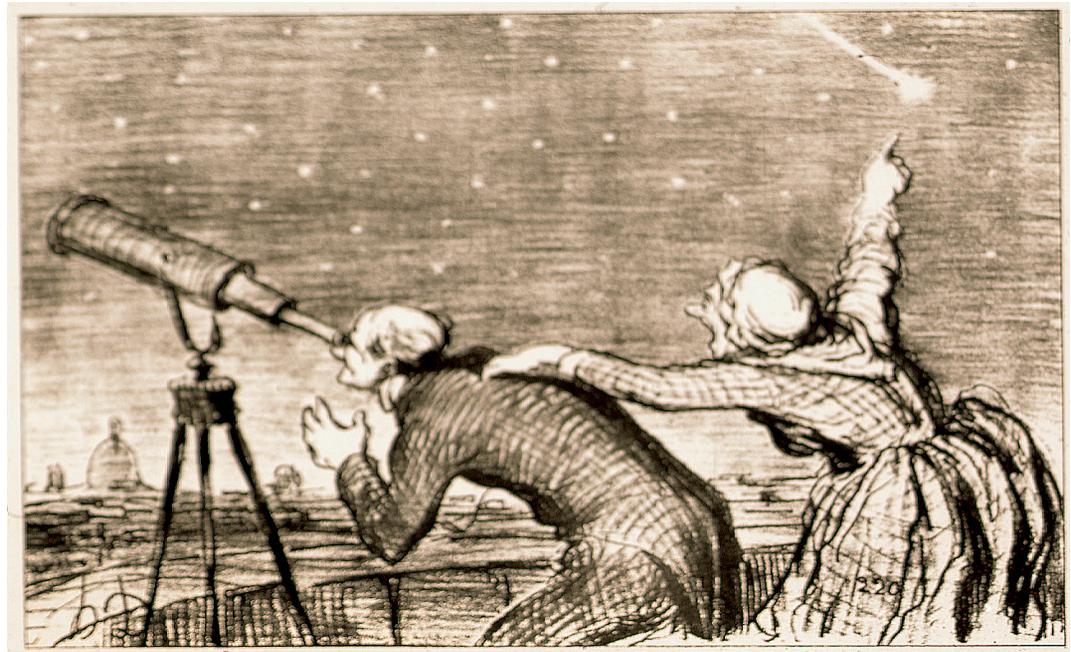


Metastructures: systems approach to determine genome annotation.



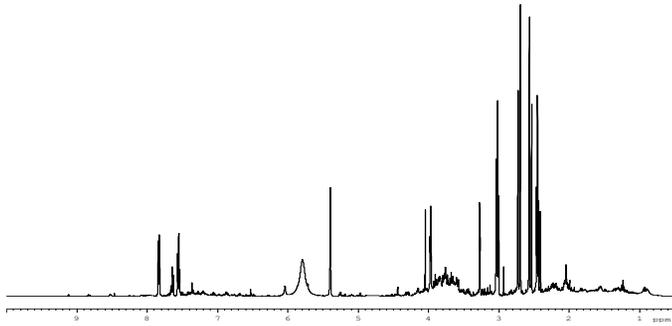
Why do we need to integrate various types of omic data?

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

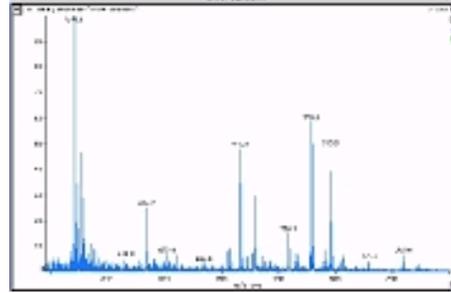


Experimental Platforms

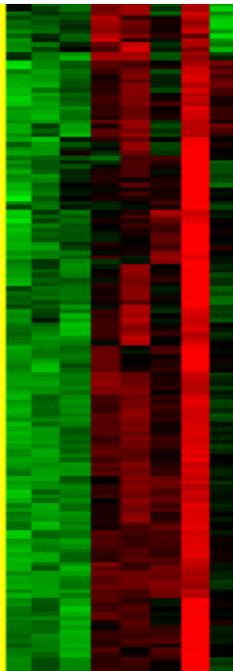
Non-omics and Omics, what are they?



NMR protein structures and proteomics

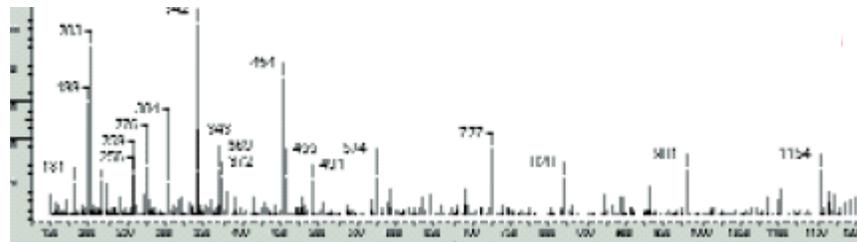
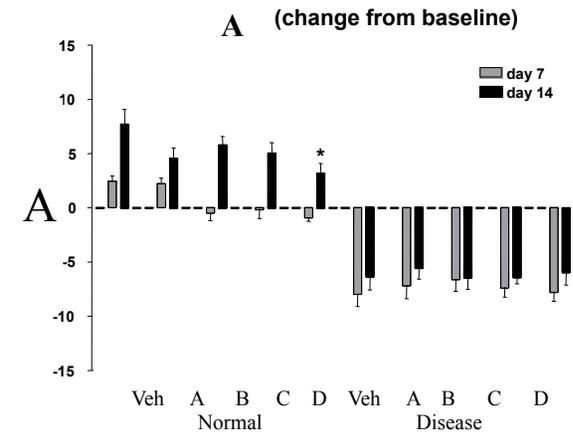


DNA sequences



Transcriptome

“Non-omic”
markers



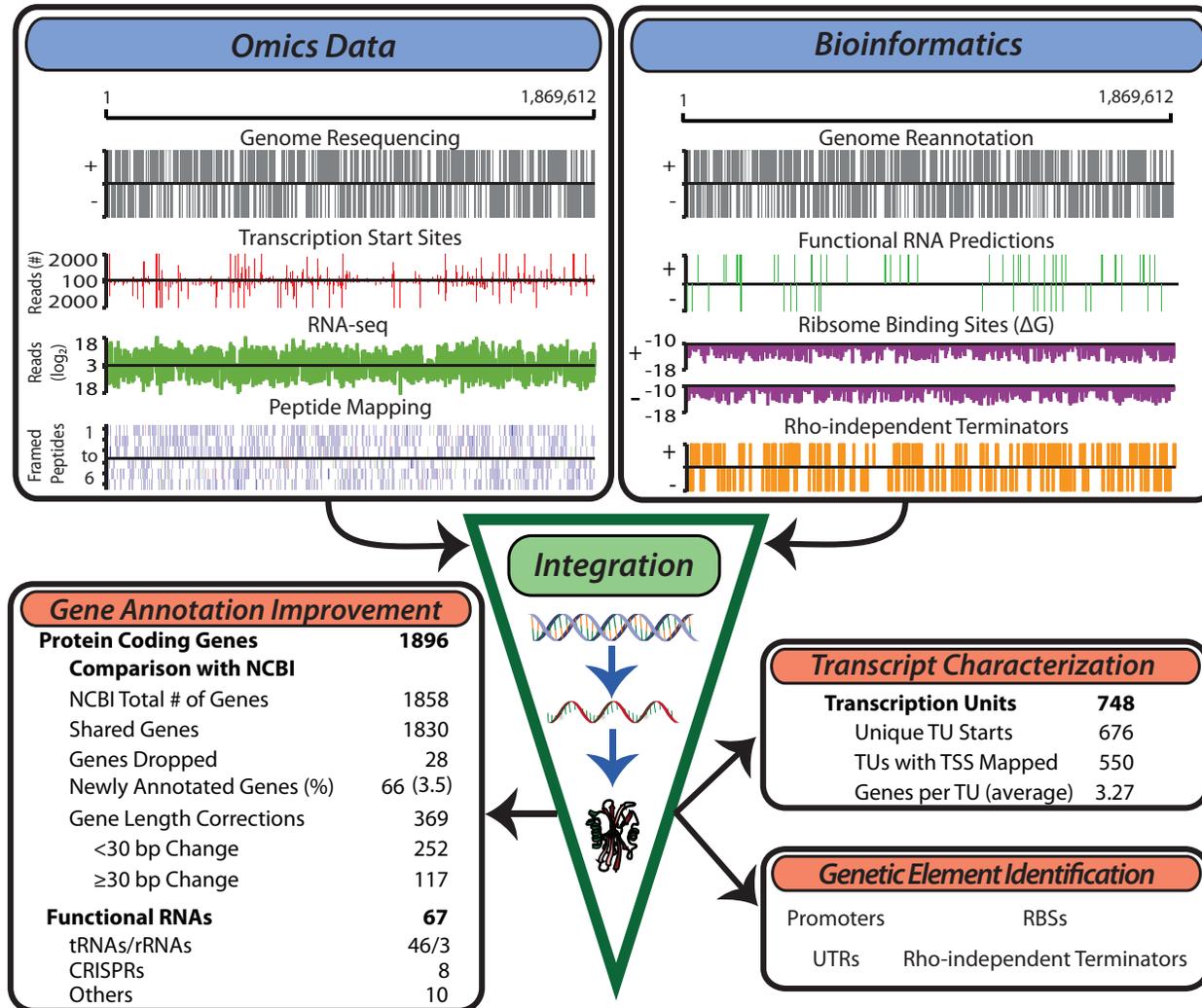
LC-MS protein sequences and interactions

Integration of omics data sets

Genomics	Transcriptomics	Proteomics	Metabolomics	Protein–DNA interactions	Protein–protein interactions	Fluxomics	Phenomics
Genomics (sequence annotation)	<ul style="list-style-type: none"> • ORF validation • Regulatory element identification⁷⁴ 	<ul style="list-style-type: none"> • SNP effect on protein activity or abundance 	<ul style="list-style-type: none"> • Enzyme annotation 	<ul style="list-style-type: none"> • Binding-site identification⁷⁵ 	<ul style="list-style-type: none"> • Functional annotation⁷⁹ 	<ul style="list-style-type: none"> • Functional annotation 	<ul style="list-style-type: none"> • Functional annotation^{71,103} • Biomarkers¹²⁵
	Transcriptomics (microarray, SAGE)	<ul style="list-style-type: none"> • Protein: transcript correlation²⁰ 	<ul style="list-style-type: none"> • Enzyme annotation¹⁰⁹ 	<ul style="list-style-type: none"> • Gene-regulatory networks⁷⁶ 	<ul style="list-style-type: none"> • Functional annotation⁸⁹ • Protein complex identification⁸² 		<ul style="list-style-type: none"> • Functional annotation¹⁰²
		Proteomics (abundance, post-translational modification)	<ul style="list-style-type: none"> • Enzyme annotation⁹⁹ 	<ul style="list-style-type: none"> • Regulatory complex identification 	<ul style="list-style-type: none"> • Differential complex formation 	<ul style="list-style-type: none"> • Enzyme capacity 	<ul style="list-style-type: none"> • Functional annotation
			Metabolomics (metabolite abundance)	<ul style="list-style-type: none"> • Metabolic-transcriptional response 		<ul style="list-style-type: none"> • Metabolic pathway bottlenecks 	<ul style="list-style-type: none"> • Metabolic flexibility • Metabolic engineering¹⁰⁹
				Protein–DNA interactions (ChIP–chip)	<ul style="list-style-type: none"> • Signalling cascades^{89,102} 		<ul style="list-style-type: none"> • Dynamic network responses⁸⁴
					Protein–protein interactions (yeast 2H, coAP–MS)		<ul style="list-style-type: none"> • Pathway identification activity⁸⁹
						Fluxomics (isotopic tracing)	<ul style="list-style-type: none"> • Metabolic engineering
							Phenomics (phenotype arrays, RNAi screens, synthetic lethals)

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

Multi-omic data integration

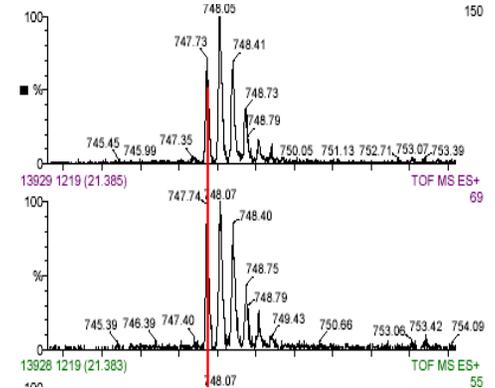
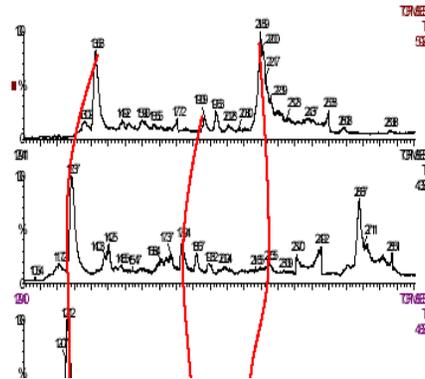


Challenges

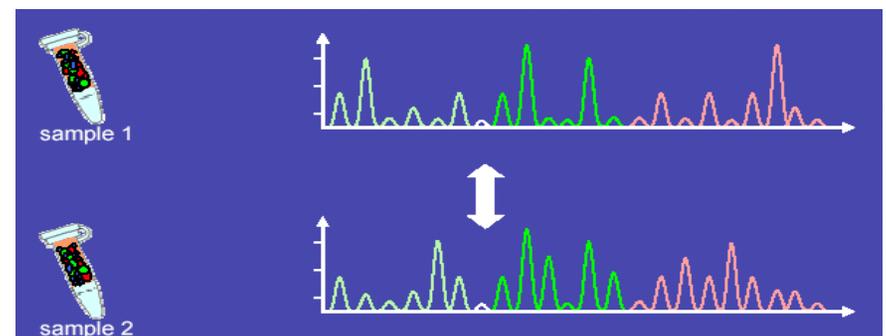
1. Data Pre-processing
2. High Dimensionality
3. Multiple Testing for Marker Selection
4. Data Integration
5. Validation of the Prediction Model

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

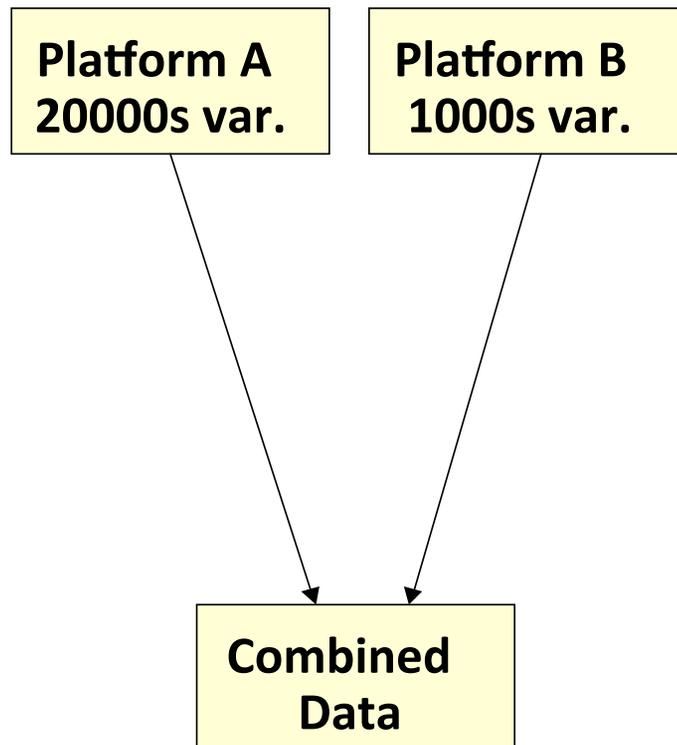
of subjects \ll # of variables

	Choles, Trig,... ...	probe set 1 22,000	Lipid 1 2,000	Metabolite 1 ... 3,000	NMR 1 500
Animal 1					
Animal 2					
.					
.					
.					
.					
.					
.					
.					
Animal 100					

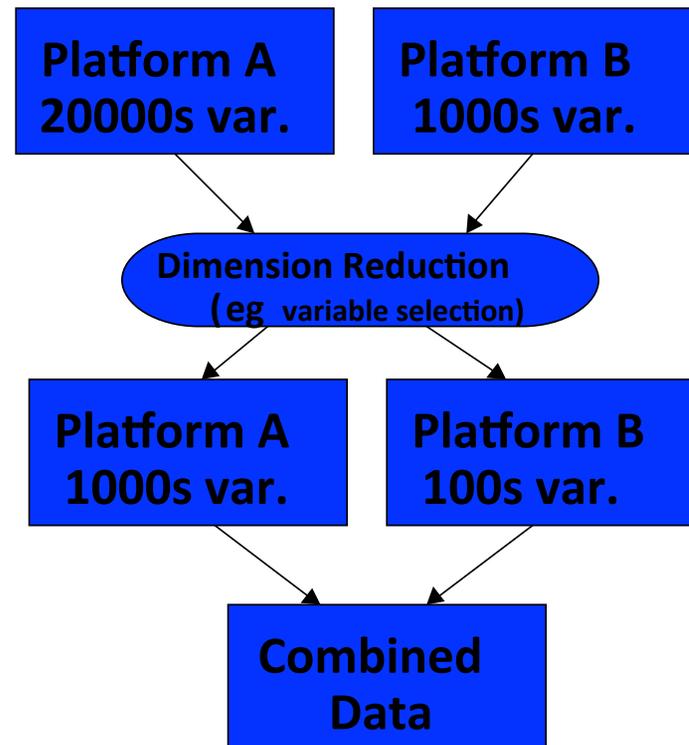
- Genes variants: 5000 peaks
- Gene Expression (RNA-seq): 22,000 probe sets
- Protein-DNA interactions (ChIP-seq): 2, 000 peaks
- Protein interactions: 100,000,000

Challenge #4: Data integration (How?)

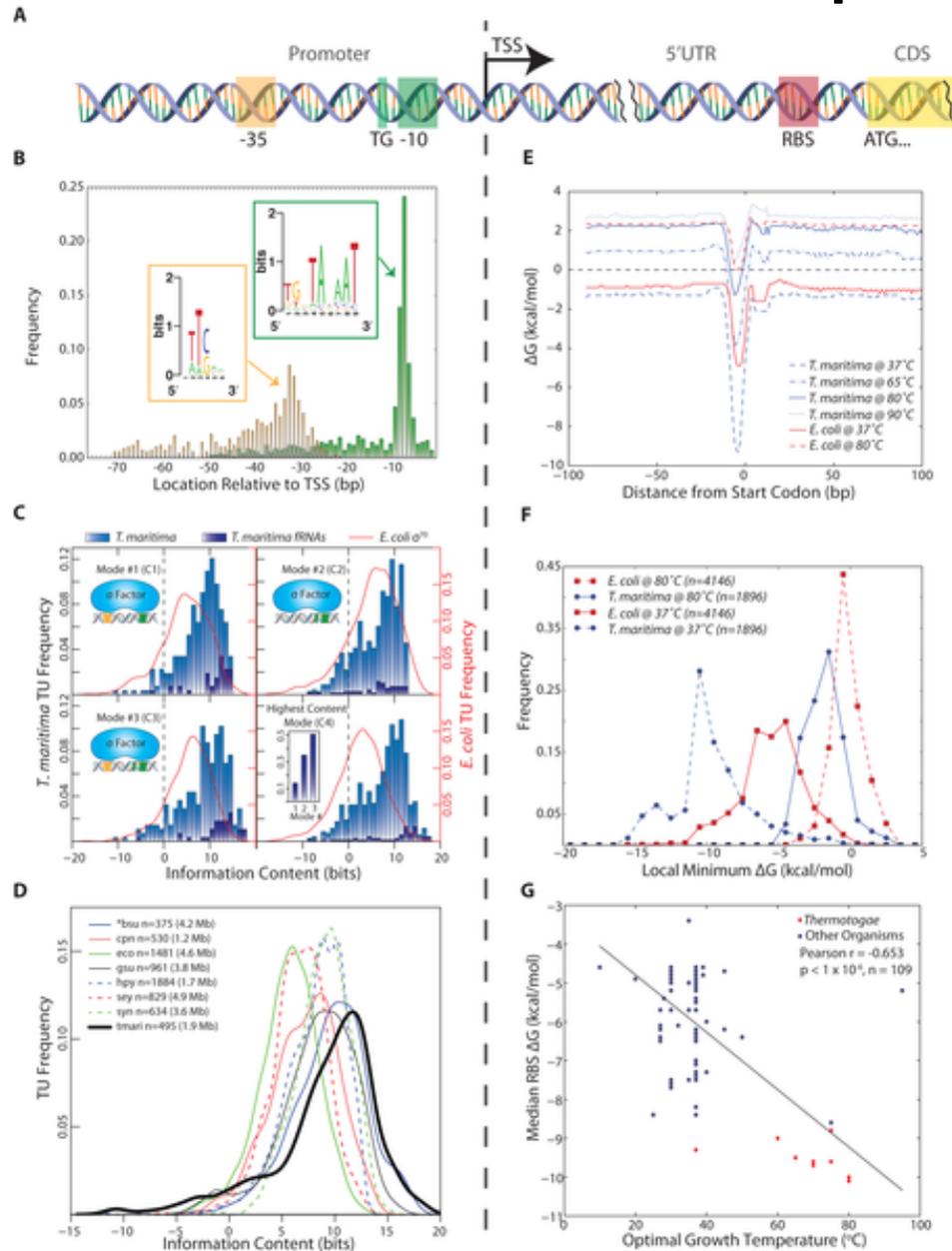
Integration Approach 1:



Integration Approach 2:



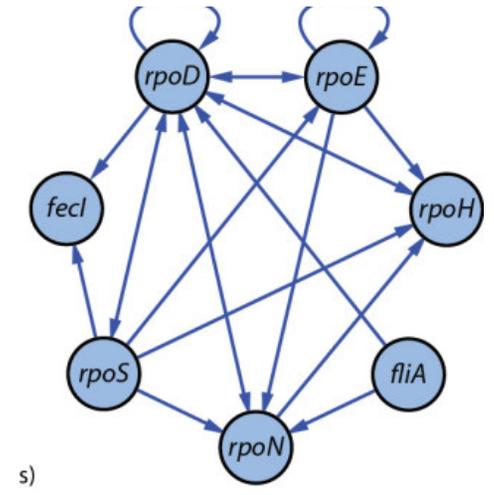
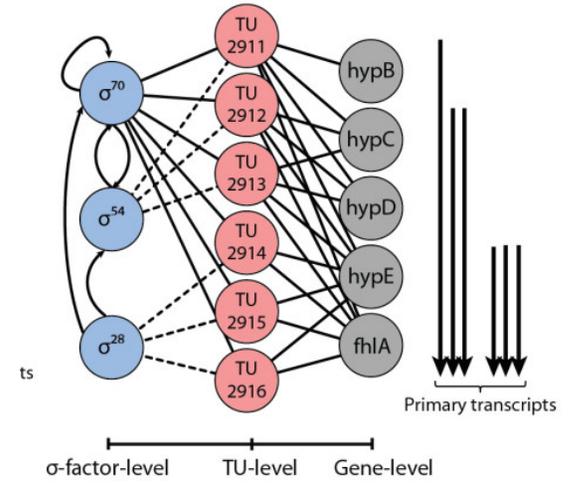
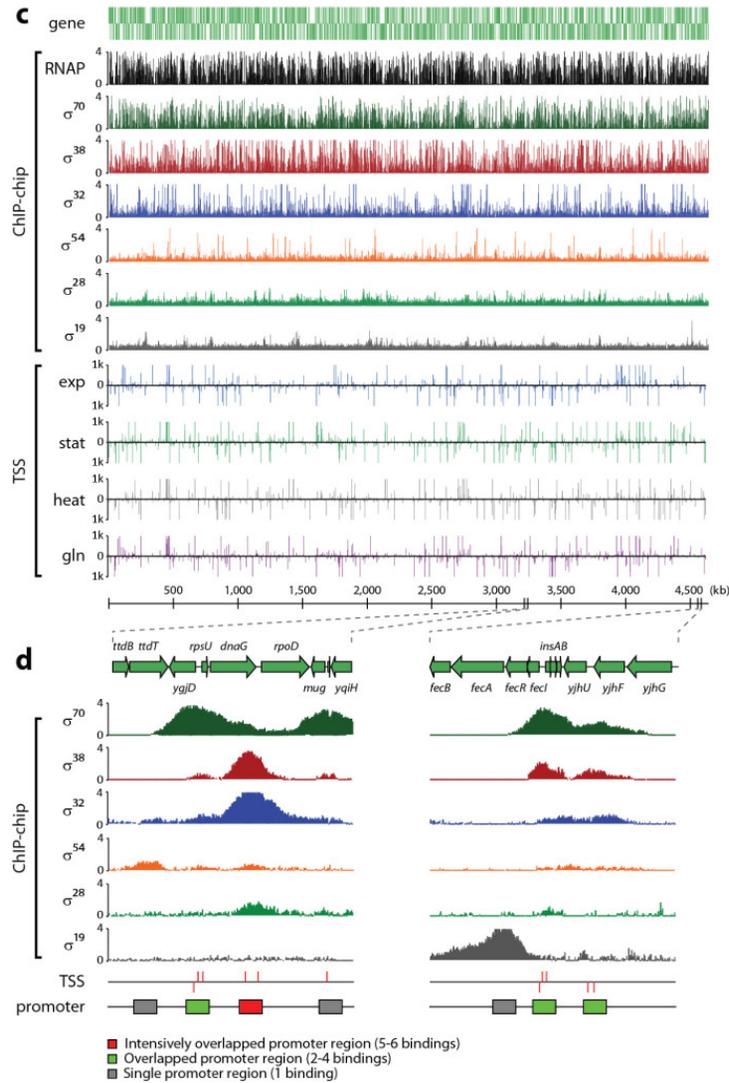
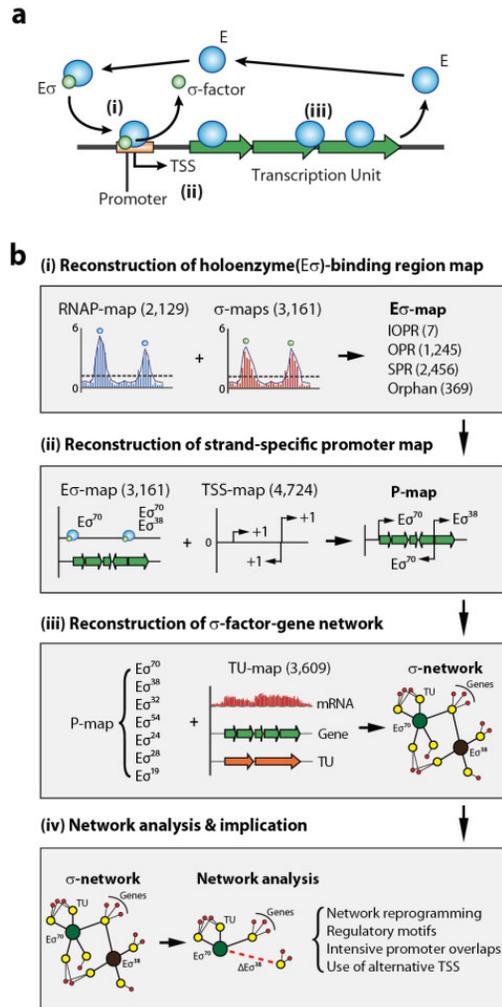
Example 1



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

Example 2

σ -factor network in *E. coli*.



Example 3

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

Chaolin Zhang et al. Science 2010

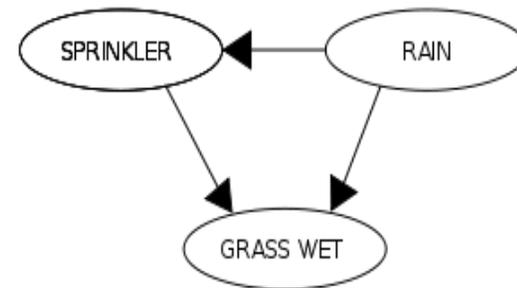
Data Integration

- 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 (microarray) + 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.

		SPRINKLER	
RAIN		T	F
F		0.4	0.6
T		0.01	0.99



		RAIN	
		T	F
		0.2	0.8

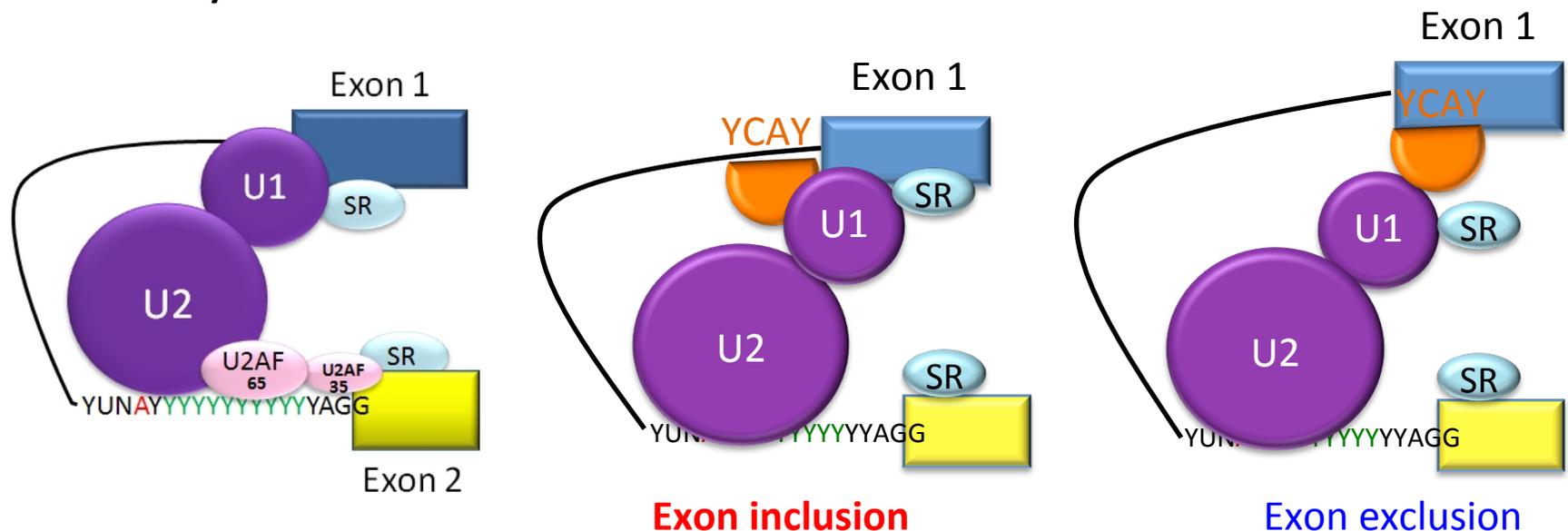
Applications:

1. Inferring unobserved variables
2. Parameter learning

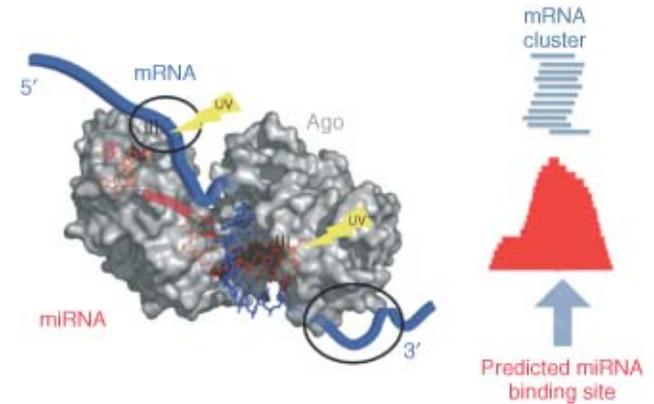
		GRASS WET	
SPRINKLER	RAIN	T	F
F	F	0.0	1.0
F	T	0.8	0.2
T	F	0.9	0.1
T	T	0.99	0.01

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 “YCAAY” clusters.
- Nova binding to an **exonic YCAAY** cluster changed the protein complexes assembled on pre-mRNA, blocking U1 snRNP binding and exon inclusion.
- Nova binding to an **intronic YCAAY** cluster enhanced spliceosome assembly and exon inclusion.



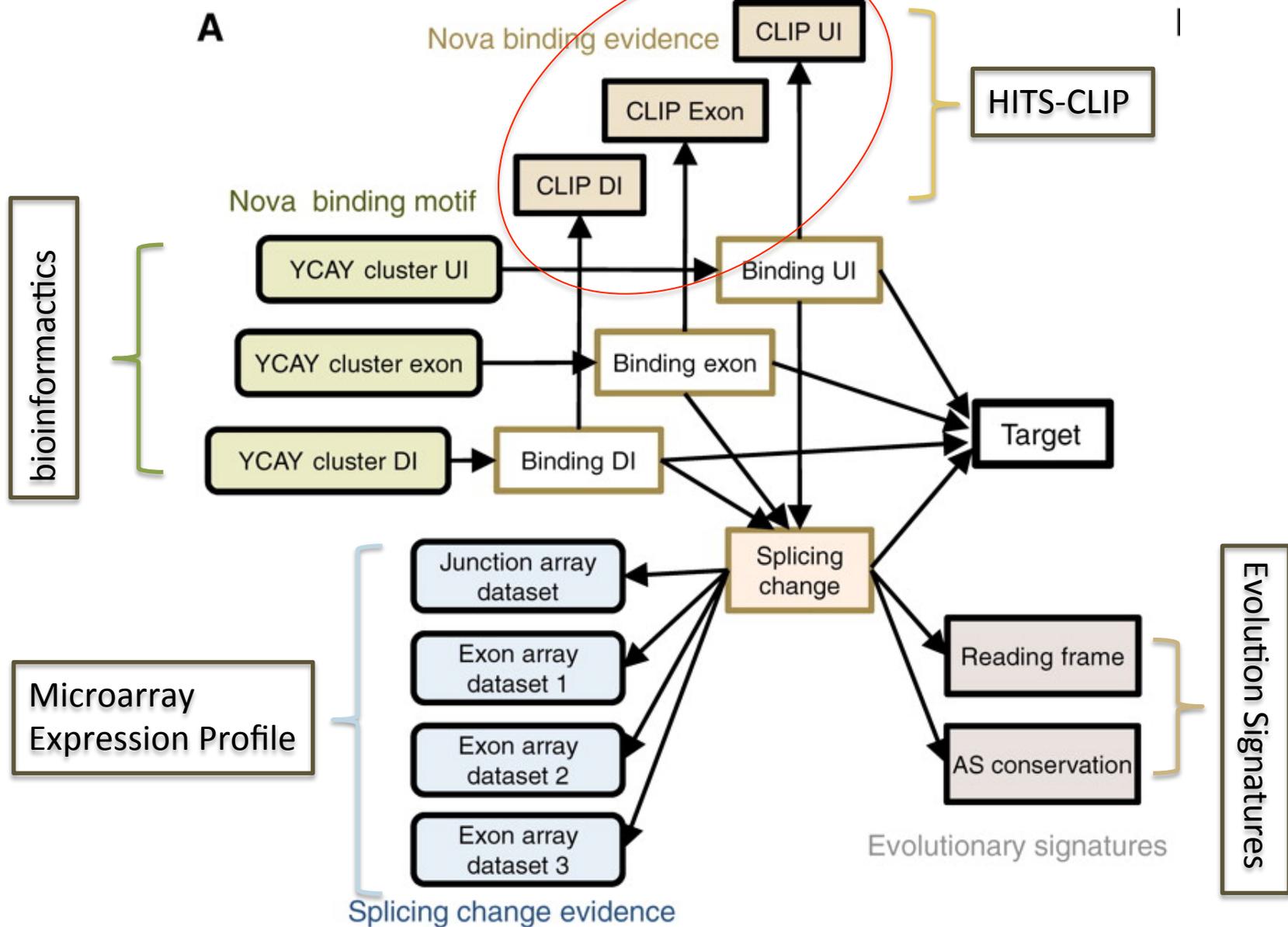
Data sources



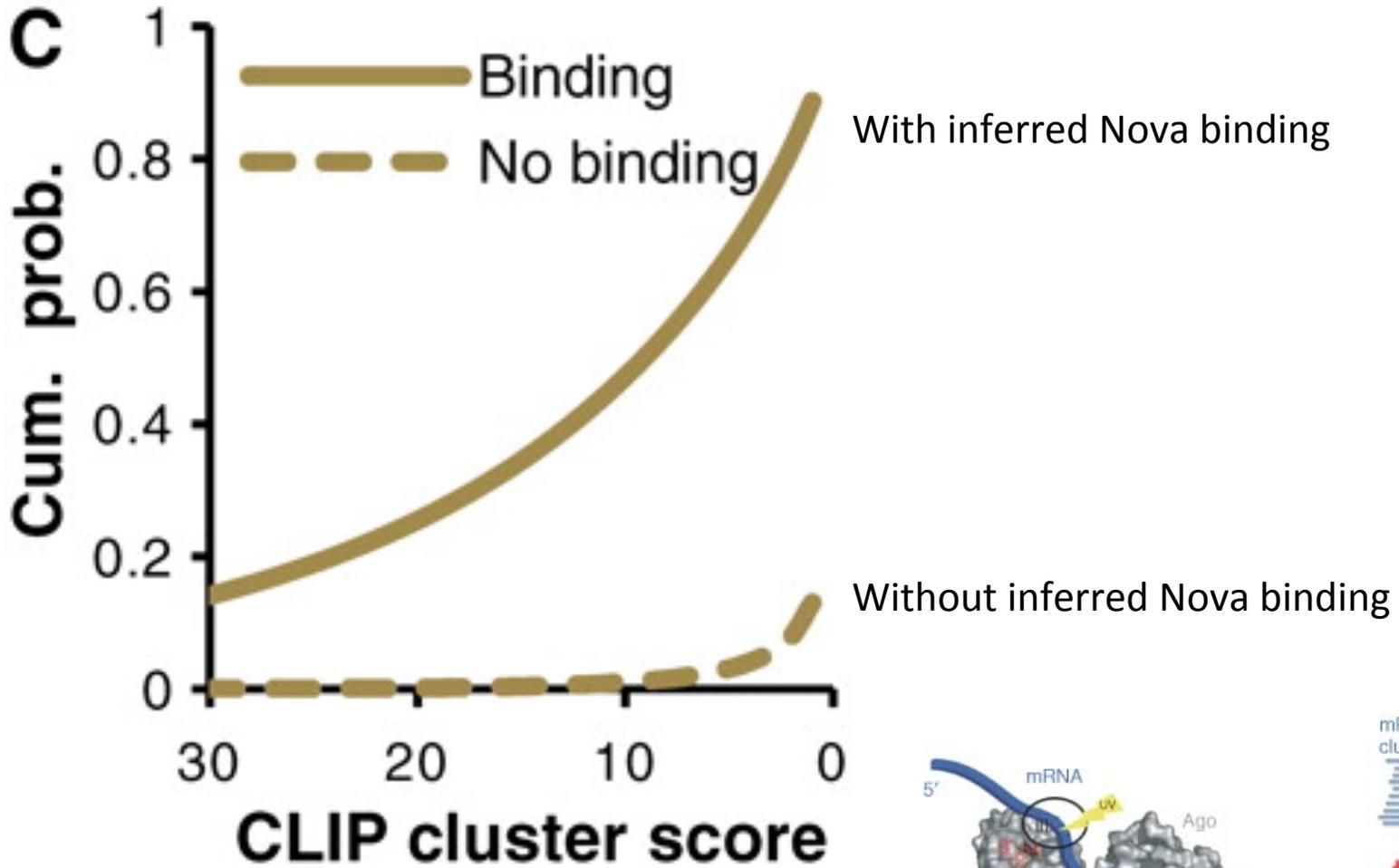
- **HITS-CLIP** (CLIP-seq). Study Protein-RNA binding by crosslinking between RNA and the protein, followed by immunoprecipitation and high-throughput sequencing.
- **Genome-wide searching of YCAY motif.** Bioinformatics approach.
- **Microarray data compared WT and Nova knockout.**
- **Evolution signature.**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

Bayesian Model

E: Exon
 U: Upstream of Intron
 D: Downstream of Intron

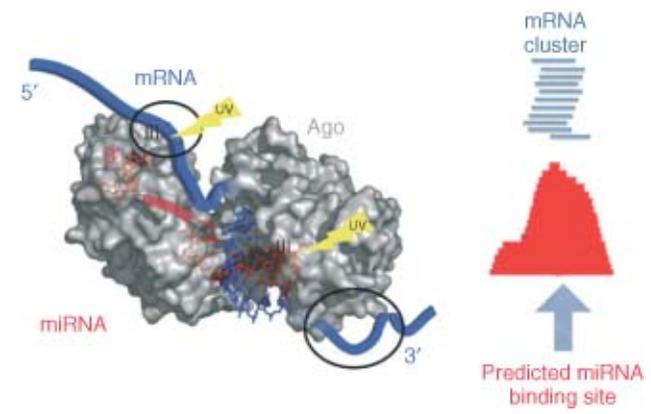


Estimated Conditional prob.

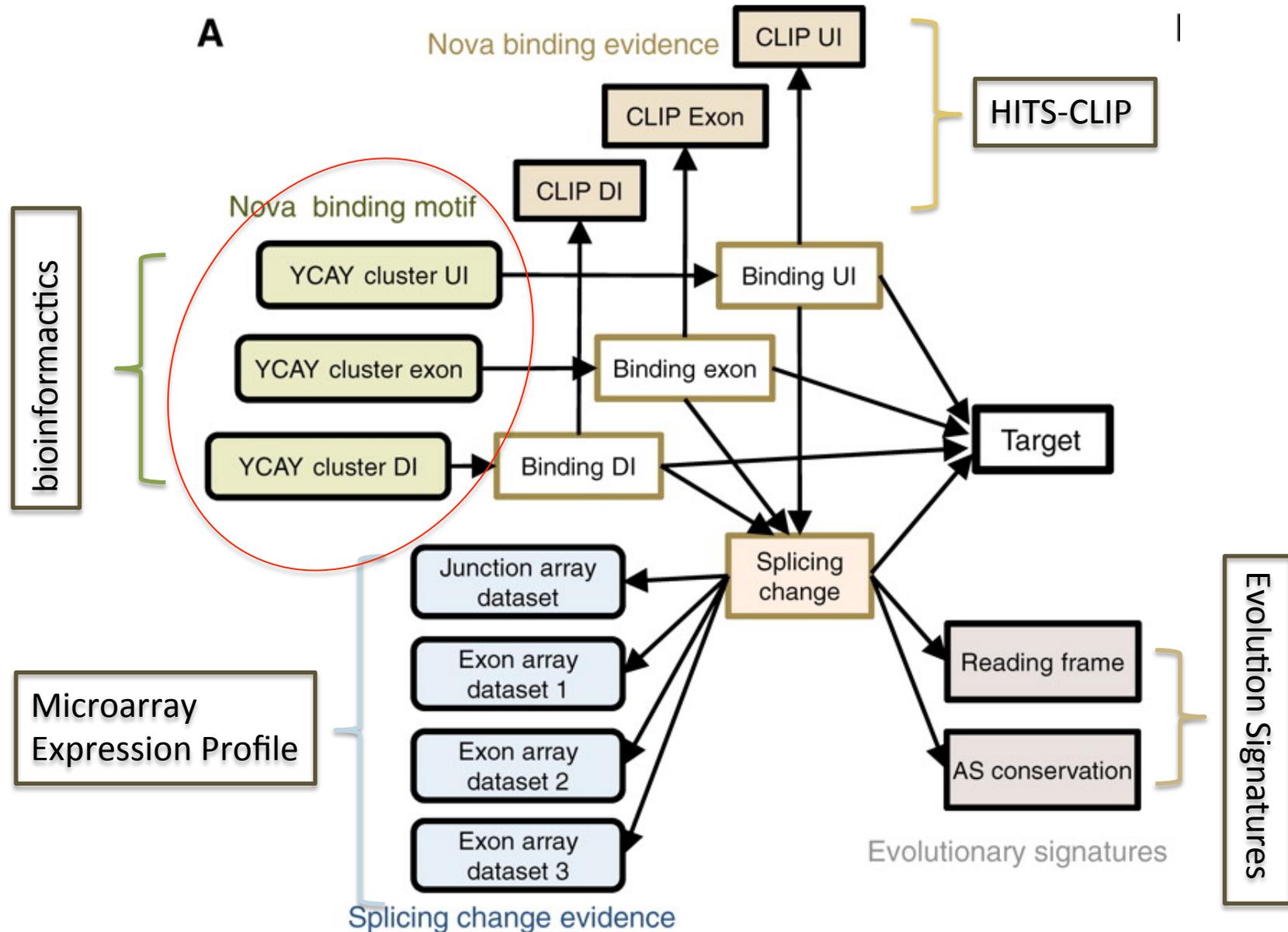


With inferred Nova binding

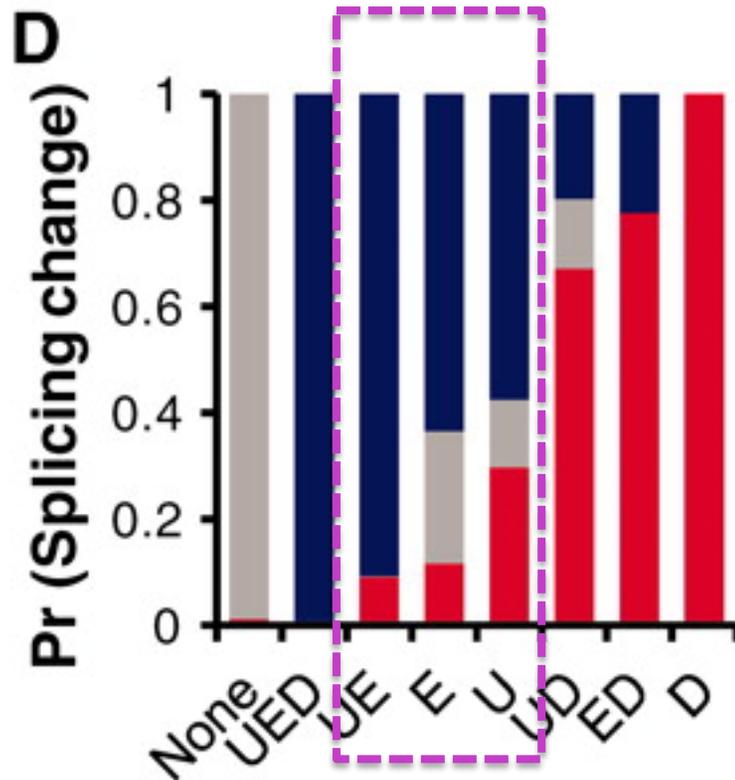
Without inferred Nova binding



Bayesian Model

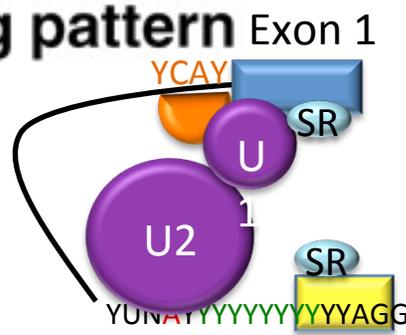


Regulation of Nova binding position

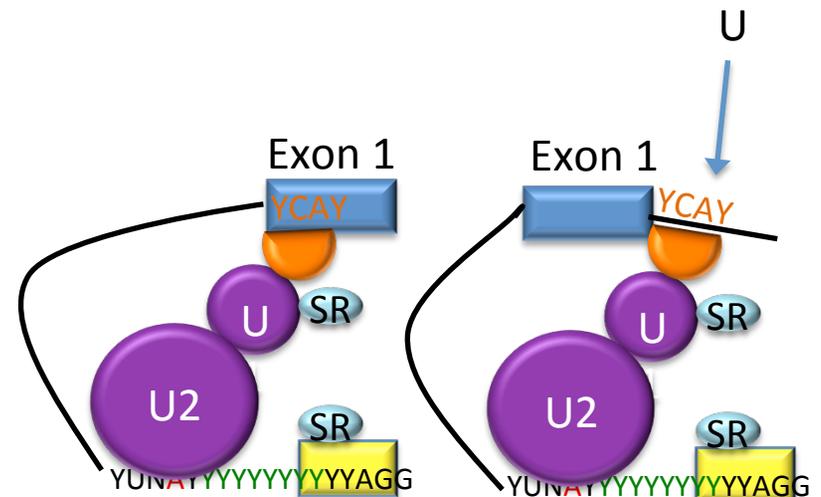


E: Exon
 U: Upstream of Intron
 D: Downstream of Intron

Binding pattern

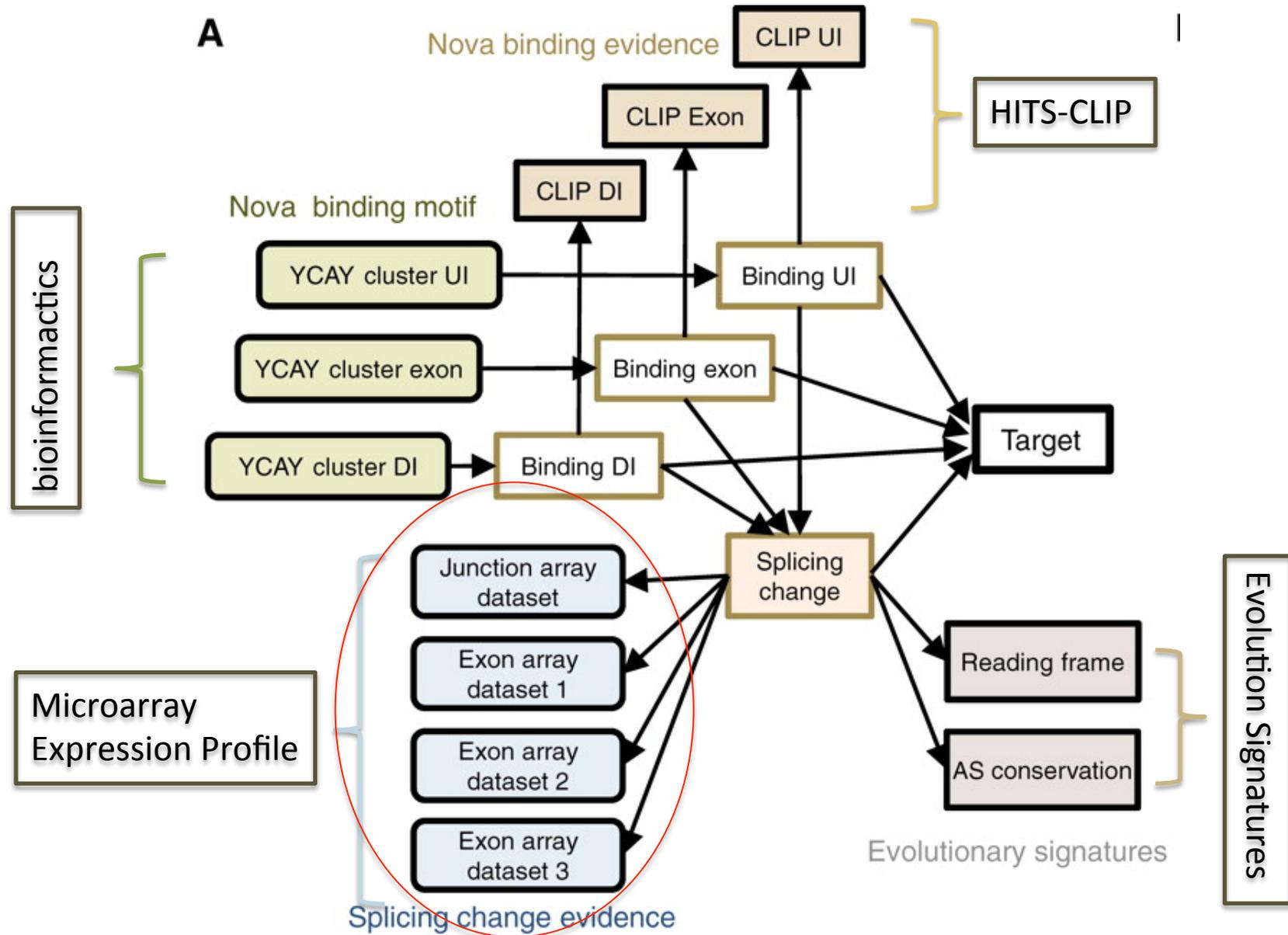


Exon inclusion

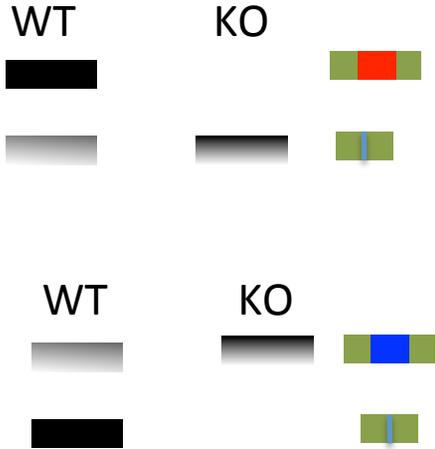
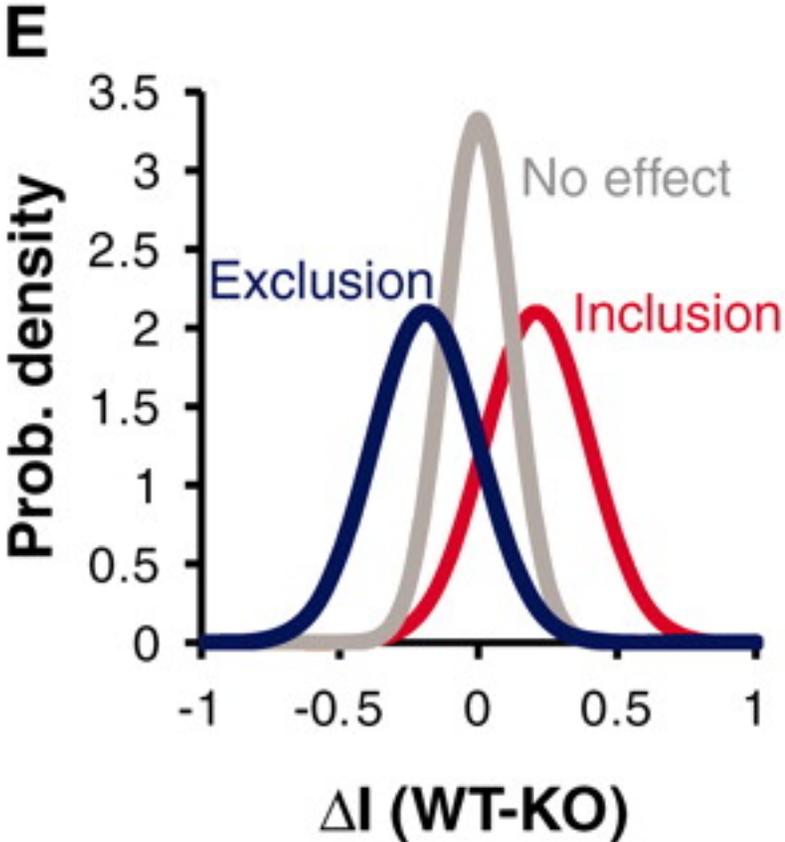


Exon exclusion

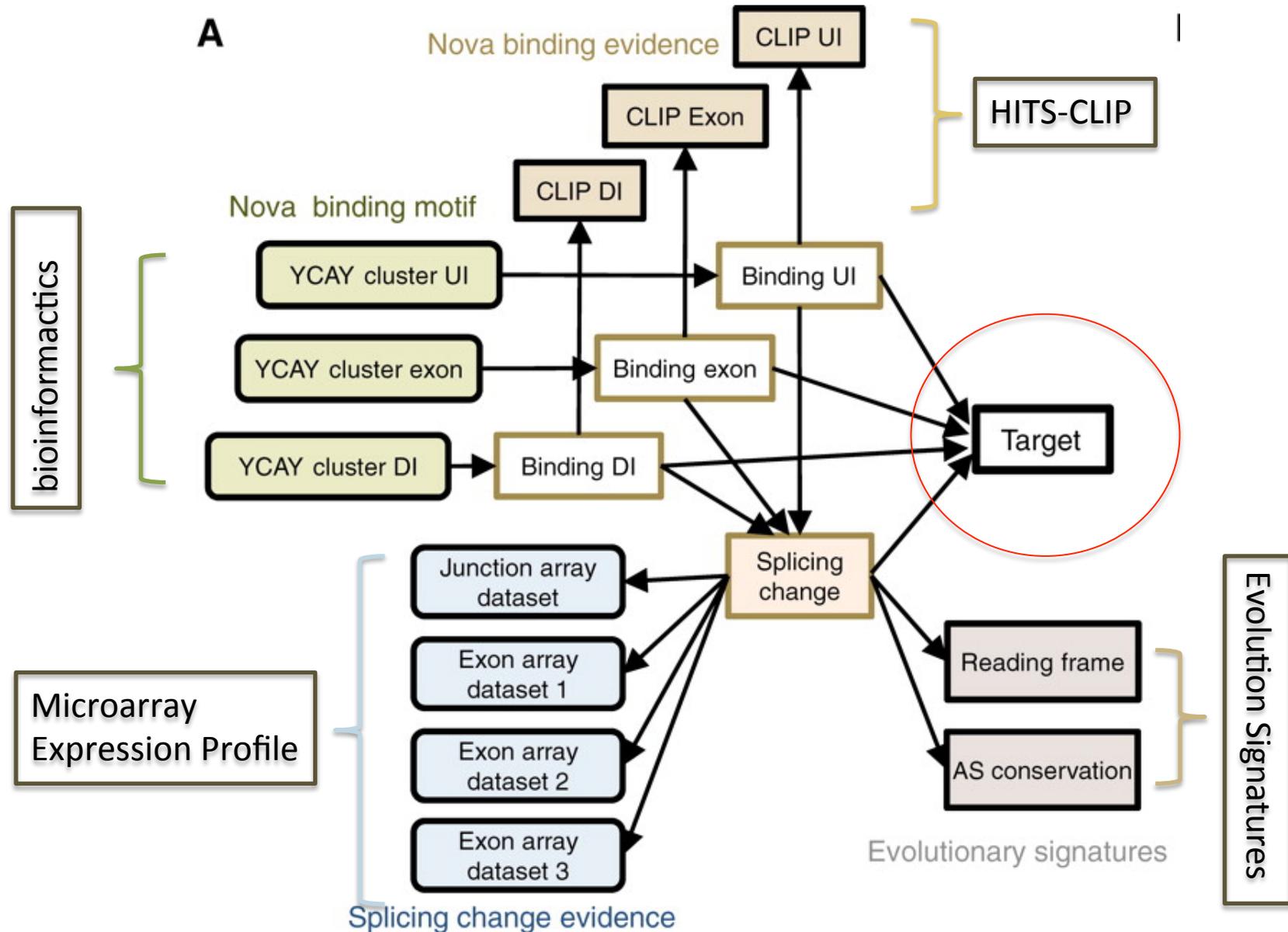
Bayesian Model



Splicing changes

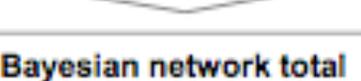


Bayesian Model



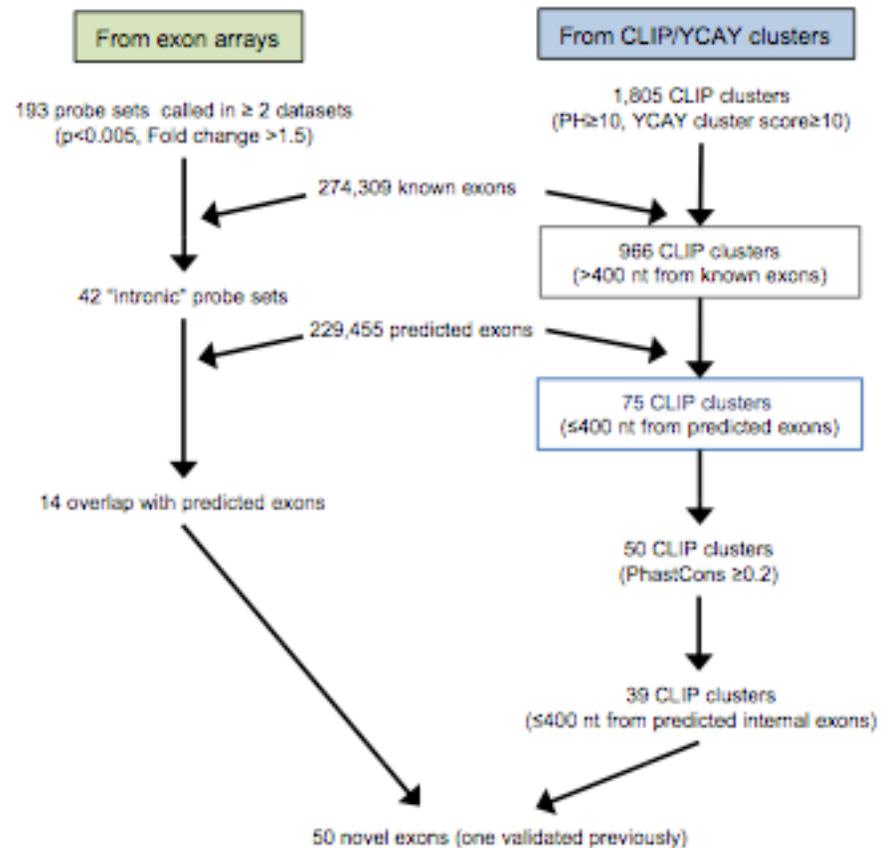
Predicted Nova-regulated targets

- 13,357 annotated cassette exons

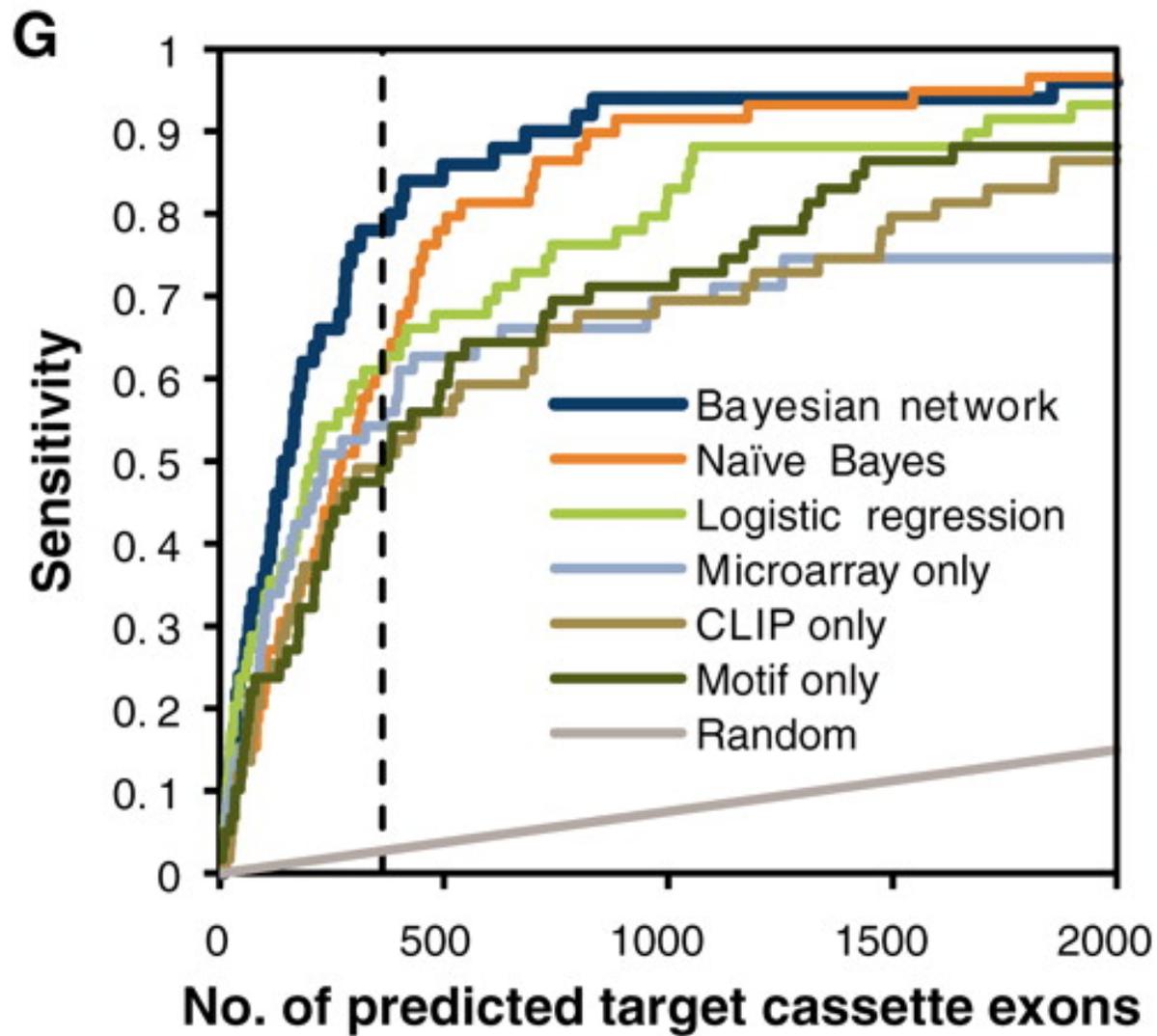
AS type	AS diagram	No. AS events
Bayesian network predictions:		
CASS		363
TACA		141
MUTX		37
ALT5		9
ALT3		9
APA5		13
APA3		16
Bayesian network total		588

Novel Nova-regulated targets

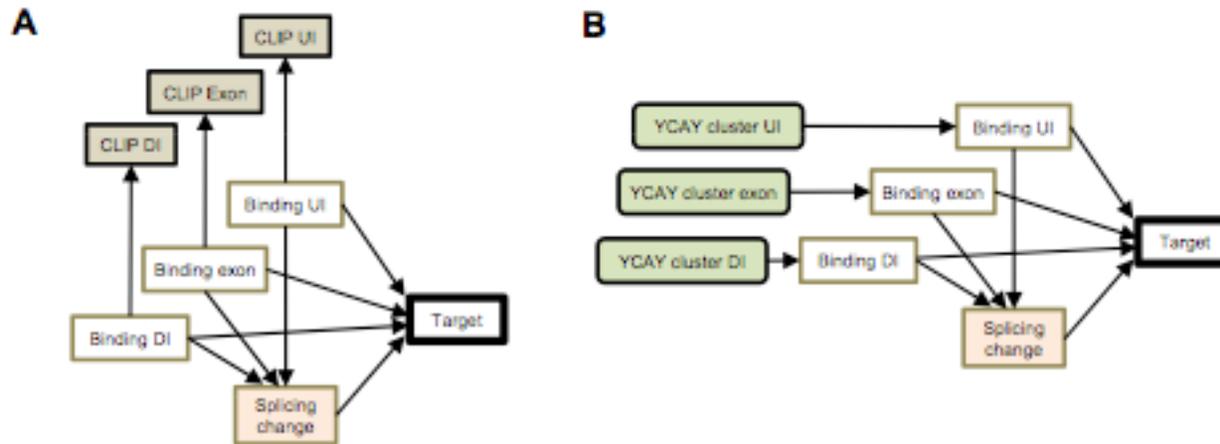
- Besides AS from database, searched novel exons with high sequence conservations.
- Additional 76 novel exons as Nova targets



Prediction Performance

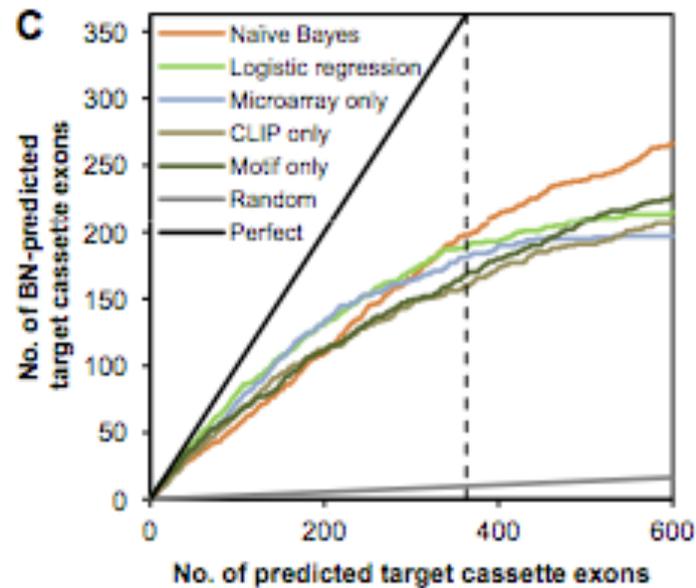


Reduced Bayesian Network



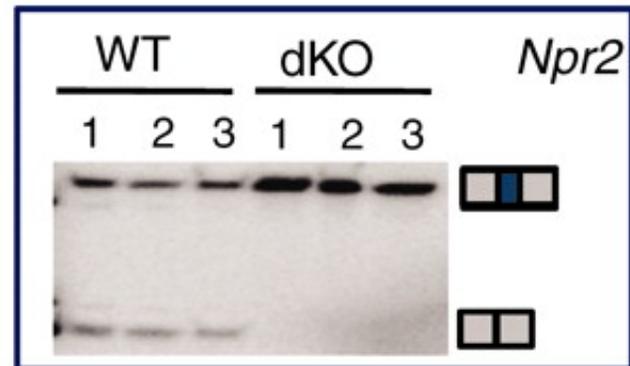
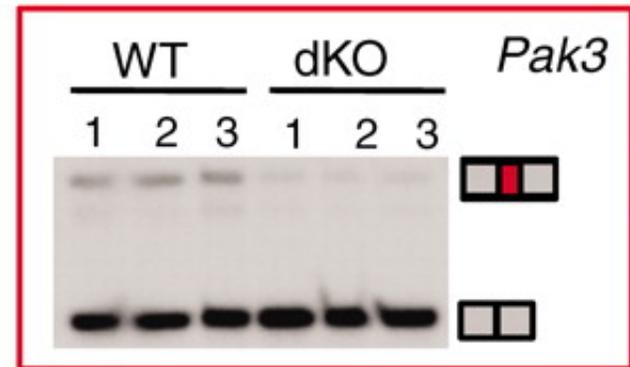
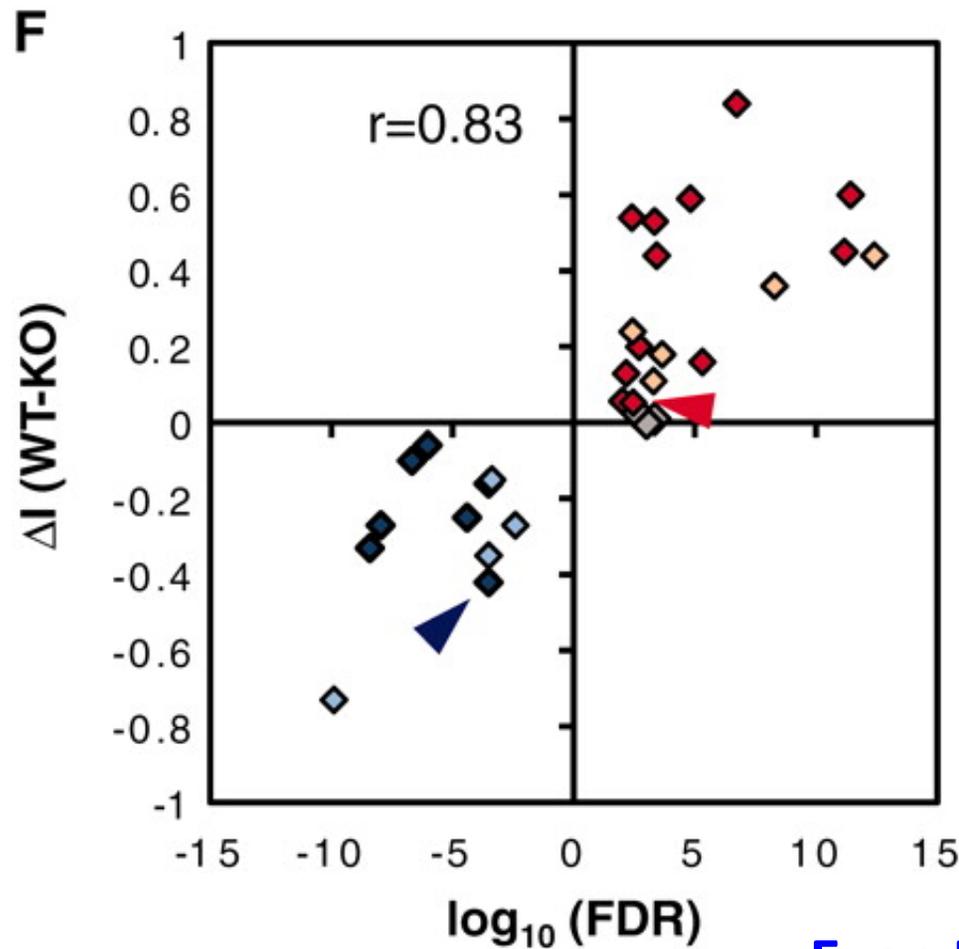
- Clip Data Only

- Motif Data Only



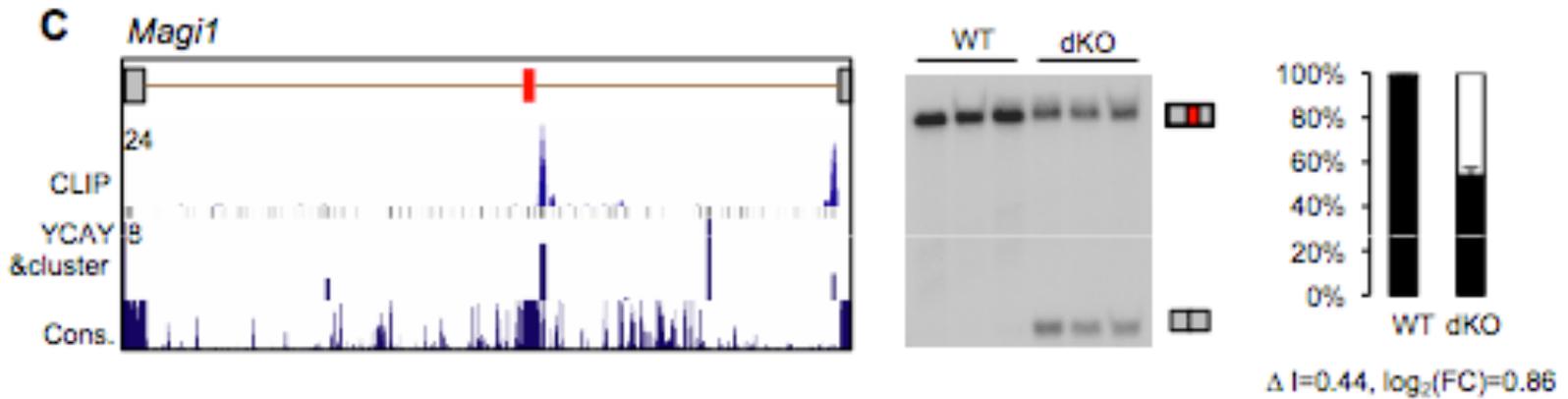
Experimental Validation

Exon Inclusion

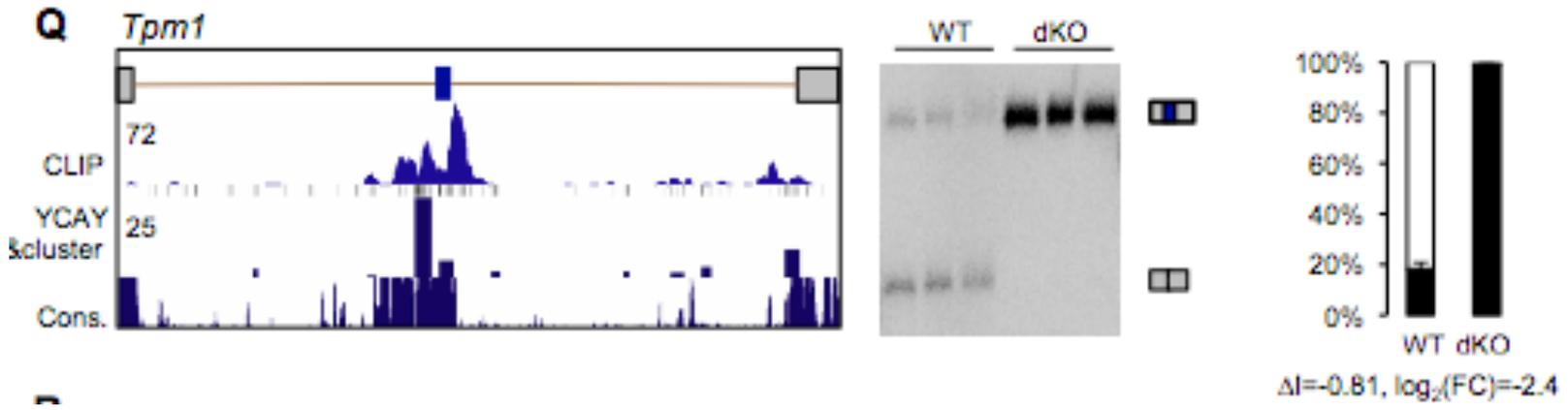


Exon Exclusion

Two more Casset Exon Cases

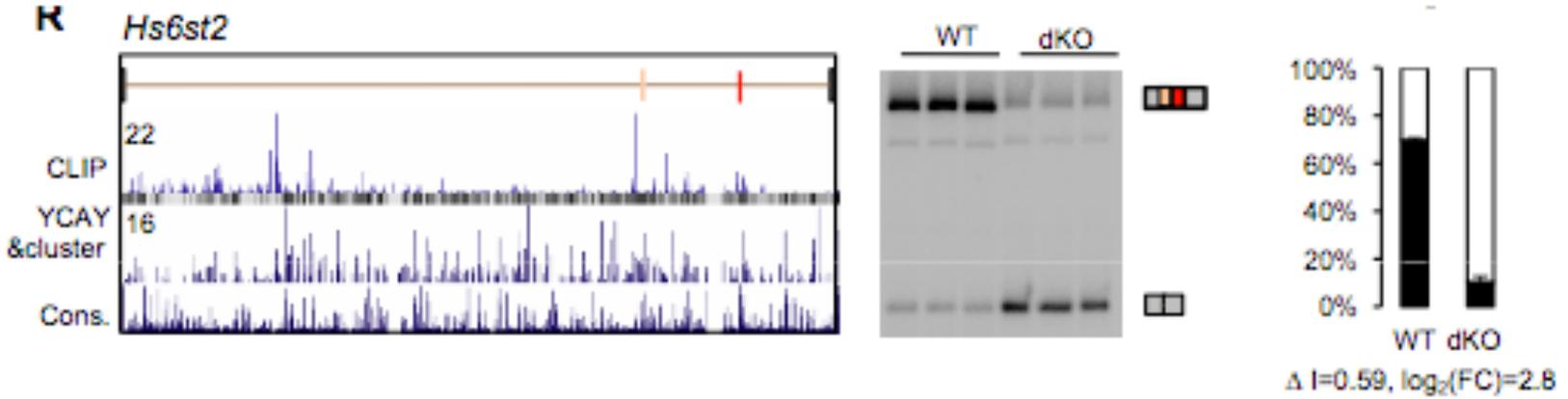


Exon Inclusion



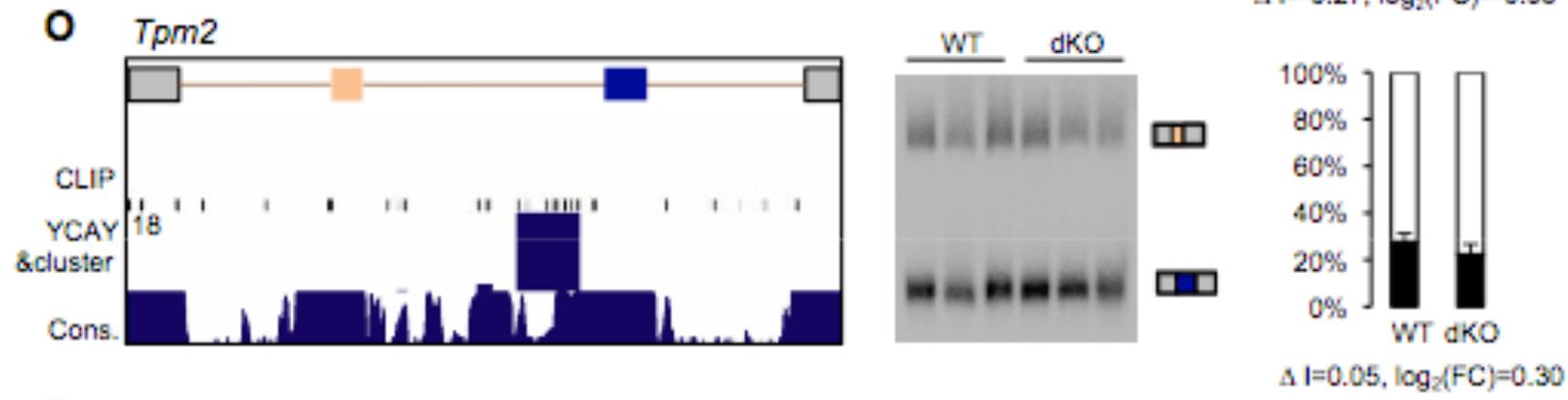
Exon Exclusion

Other examples



TACA

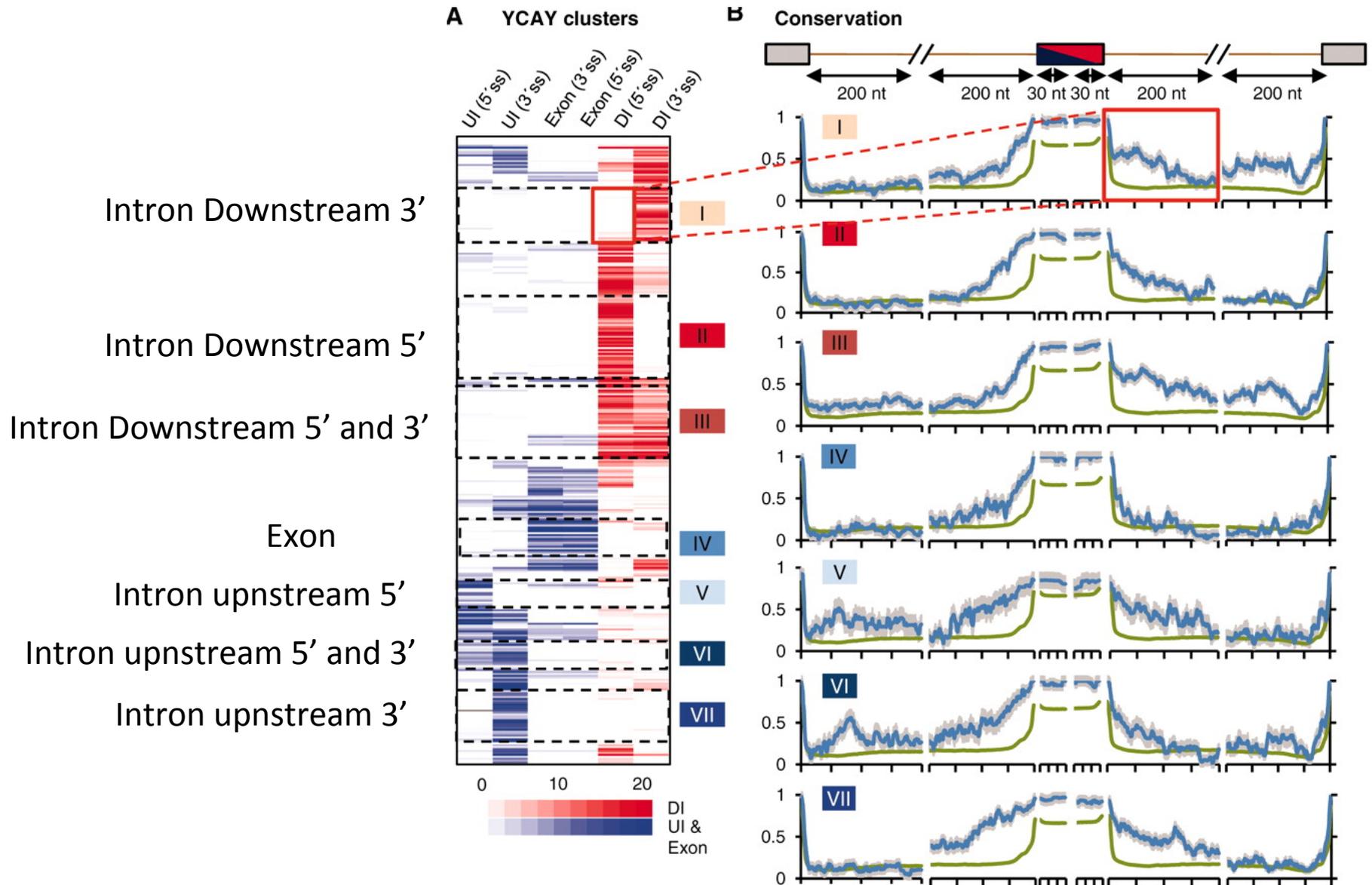
Exon Inclusion



MUTX

Exon Exclusion

Conservation regions



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.

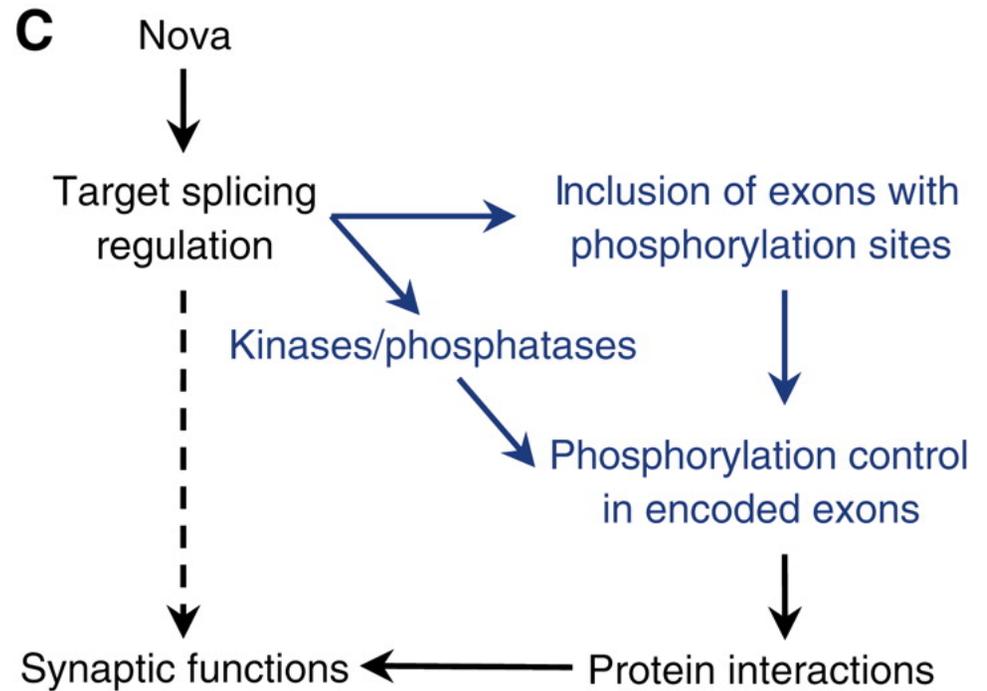
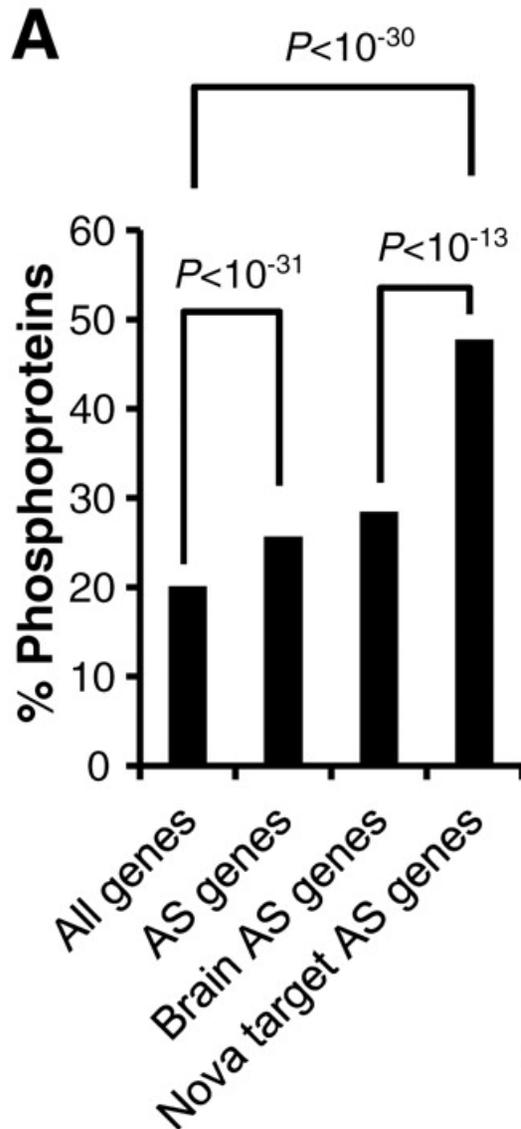
GO-term enrichment of Nova targets

GO Term	Gene count	%	P-Value	Fold Enrichment	Benjamini FDR
<i>Biological process</i>					
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	22	6.16	8.39E-09	4.64	2.90E-06
GO:0048667-cell morphogenesis involved in 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
<i>Cellular component</i>					
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
<i>Molecular function</i>					
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	26	7.28	1.17E-06	3.06	1.39E-04
GO:0060589-nucleoside-triphosphatase regulator activity	26	7.28	1.61E-06	3.01	1.54E-04
-----	---	---	---	---	---

Pathway enrichment of Nova targets

KEGG pathway	Gene count	Fold Enrichment	Benjamini FDR	Genes
mmu04020 Calcium signaling pathway	17	3.5	0.001	<i>Atp2b1, Atp2b2, Cacna1c, Cacna1d, Cacna1b, Cacna1g, Camk2a, Camk2g, Camk2b, Grin1, Gnas, Plcb4, Ppp3cb, Ppp3cc, Ryr2, Slc8a1, Erbb4</i>
mmu04720 Long-term potentiation	10	5.3	0.003	<i>Cacna1c, Camk2a, Camk2g, Camk2b, Gria2, Grin1, Plcb4, Ppp1r12a, Ppp3cb, Ppp3cc</i>
mmu04514 Cell adhesion molecules (CAMs)	12	4.1	0.003	<i>Alcam, Cadm1, Cadm3, Mpz1, Neo1, Nrxa3, Nfasc, Nfasc, Ptpfr, Ptpm, Nlgn1, Nrcam, Nrxa1</i>
mmu04520 Adherens junction	10	4.4	0.006	<i>Actn4, Baiap2, Ctnna2, Ctnnd1, Pard3, Smad2, Smad4, Ptpfr, Ptpm, Sorbs1</i>
mmu04360 Axon guidance	13	3.3	0.006	<i>Ablim1, Cxcl12, Dcc, EphA5, Efna5, Ablim2, Ntng1, Pak3, Ppp3cb, Ppp3cc, Arhgef12, Robo2, Unc5c</i>
mmu04912 GnRH signaling pathway	10	3.6	0.017	<i>Cacna1c, Cacna1d, Camk2a, Camk2g, Camk2b, Gnas, Mapk8, Mapk9, Map2k4, Plcb4</i>
mmu04310 Wnt signaling pathway	12	2.8	0.032	<i>Apc, Camk2a, Camk2g, Camk2b, Smad2, Smad4, Mapk8, Mapk9, Plcb4, Porcn, Ppp3cb, Ppp3cc</i>
mmu04930 Type II diabetes mellitus	6	5.2	0.048	<i>Cacna1c, Cacna1d, Cacna1b, Cacna1g, Mapk8, Mapk9</i>
mmu04260 Cardiac muscle contraction	7	4.2	0.049	<i>Tpm2, Cacna1d, Cacna1c, Tpm1, Ryr2, Slc8a1, Tpm3</i>
mmu04012 ErbB signaling pathway	8	3.3	0.074	<i>Camk2a, Camk2g, Camk2b, Mapk8, Mapk9, Map2k4, Pak3, Erbb4</i>
mmu04530 Tight junction	9	2.7	0.01	<i>Actn4, Cask, Ctnna2, Pard3, Epb4.1, Epb4.1f1, Epb4.1f2, Epb4.1f3, Magi1</i>

Nova targets - phosphoproteins



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.