

# Tools and technologies to characterize isoforms at proteome-scale

Gloria Sheynkman

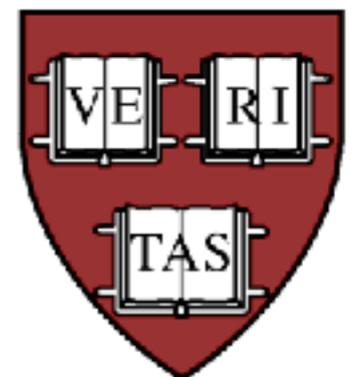
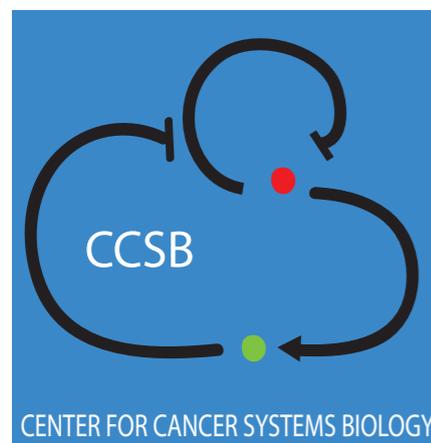
Marc Vidal Laboratory

Center for Cancer Systems Biology, Dana Faber Cancer Institute

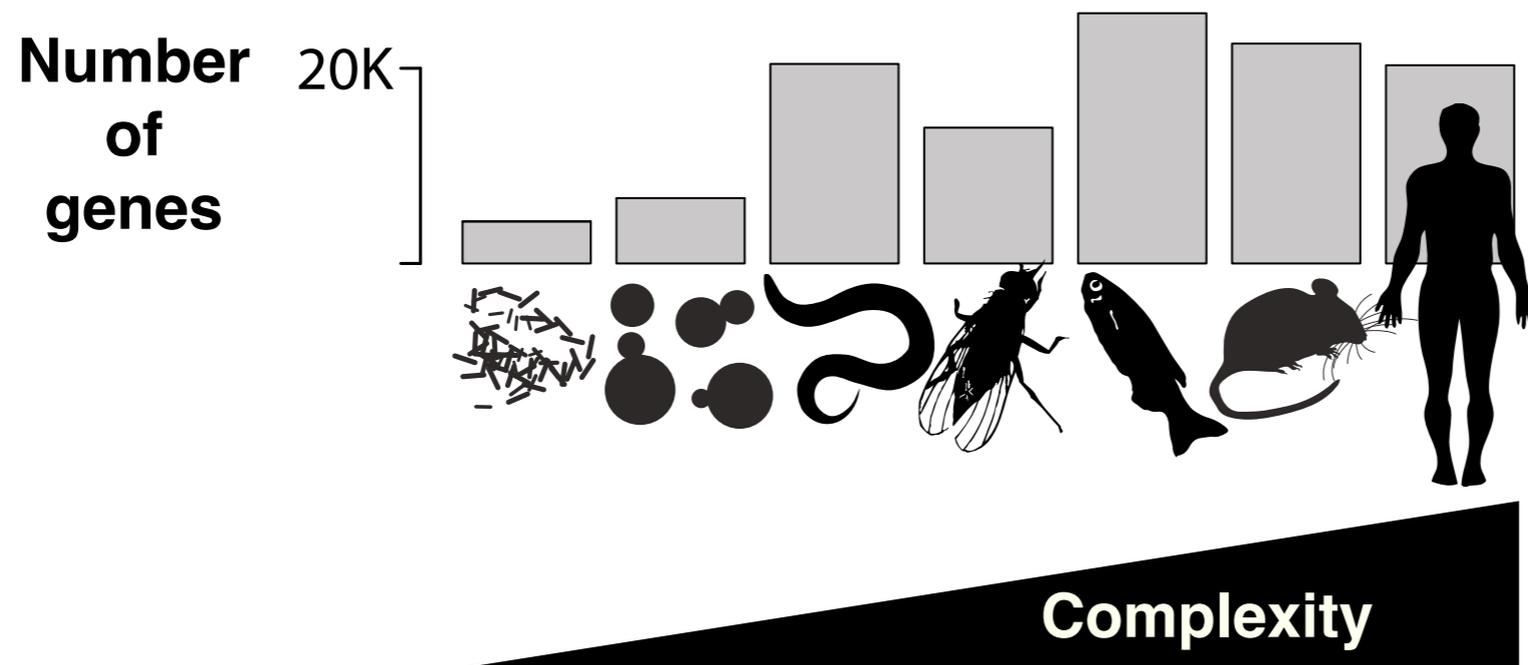
Department of Genetics, Harvard University

SMRTLeiden

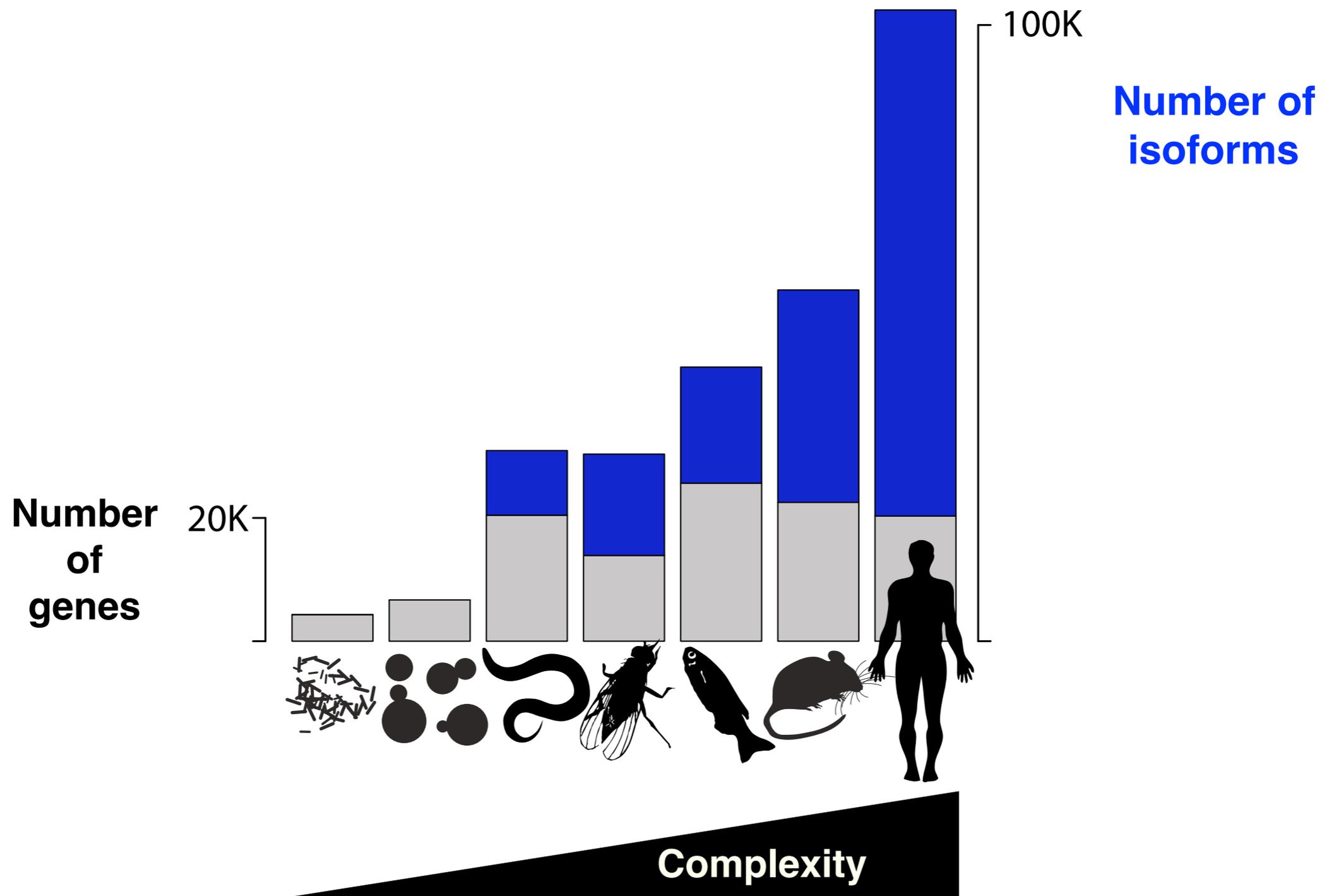
May 2nd, 2017



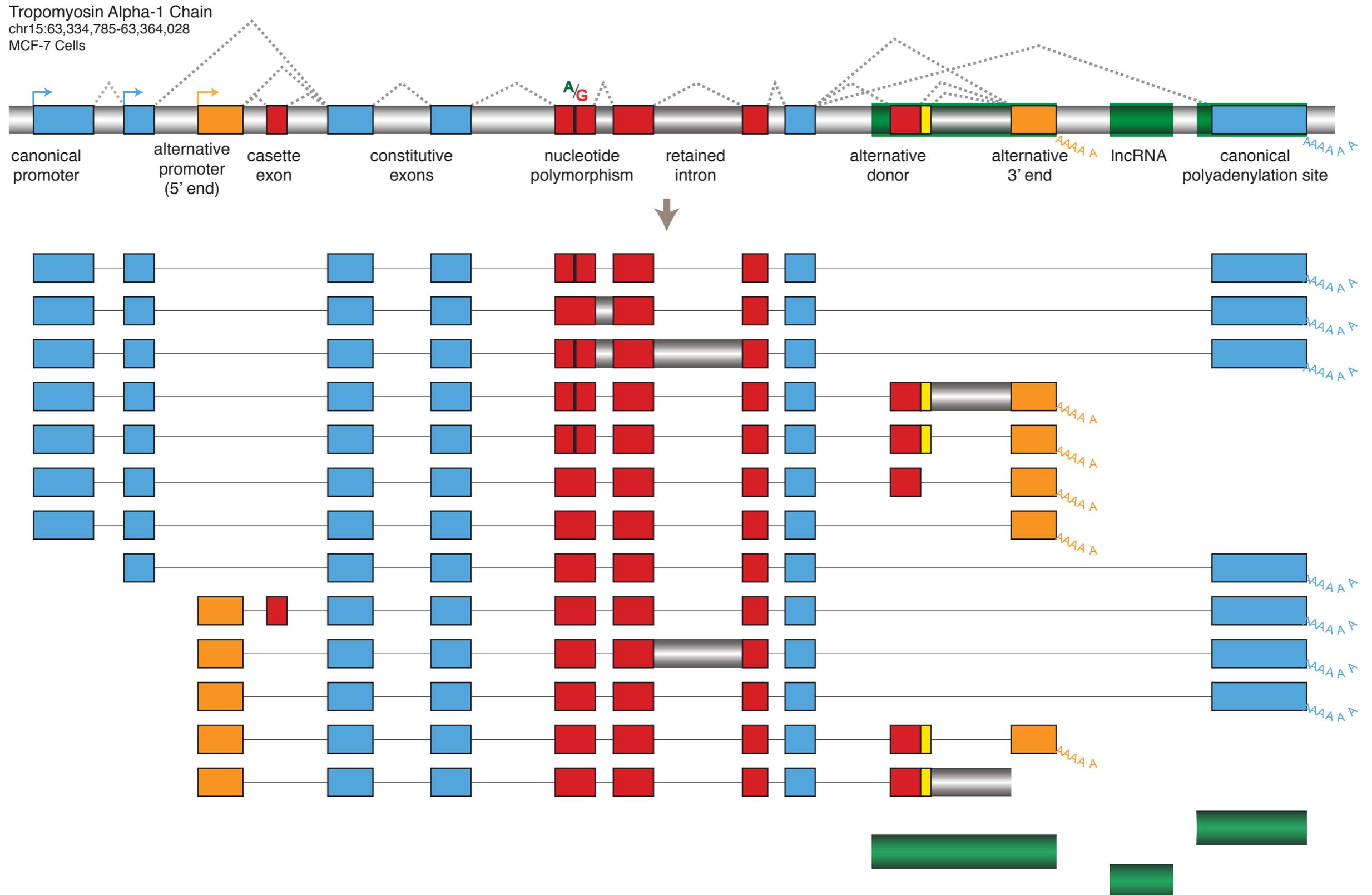
# Gene numbers



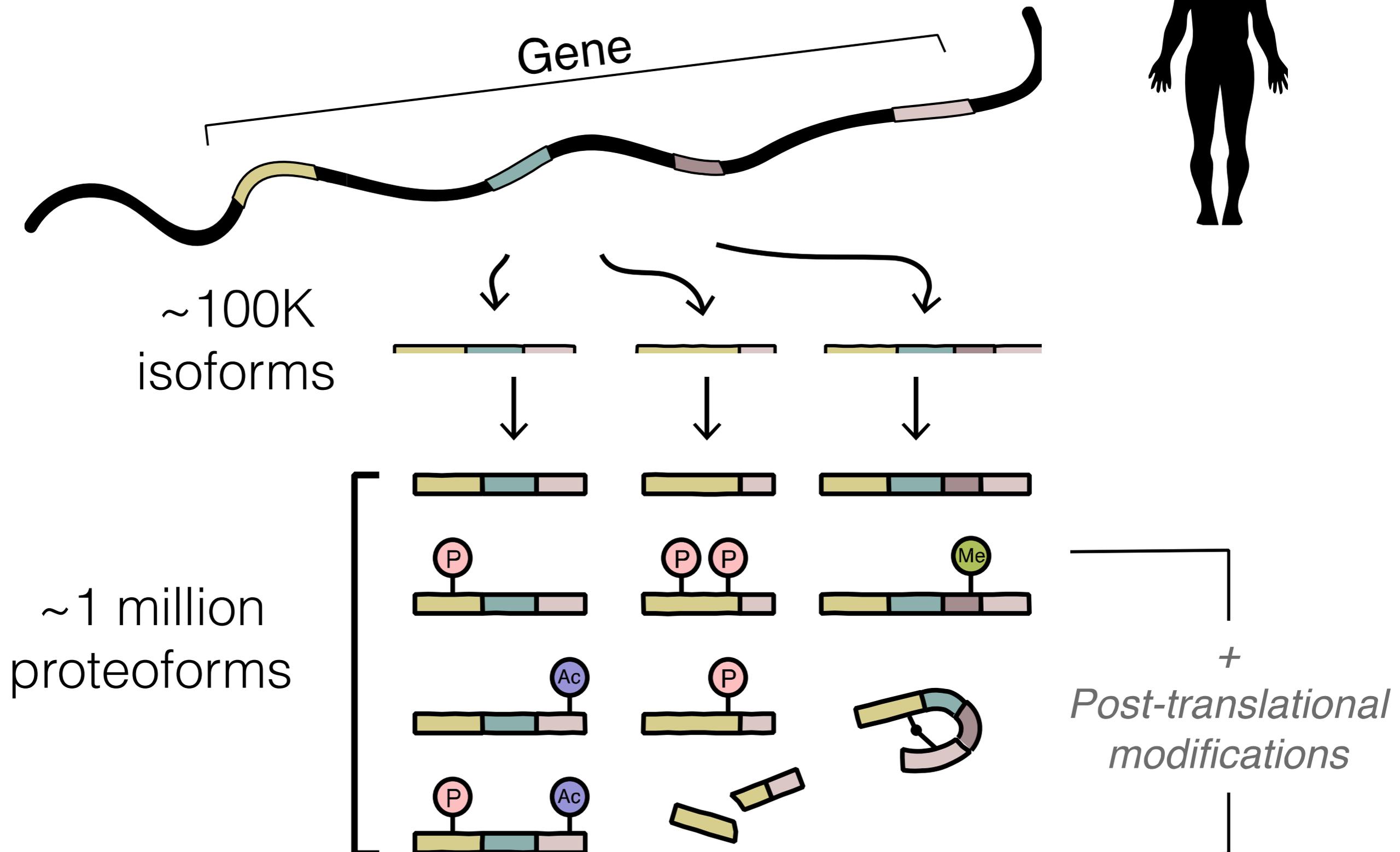
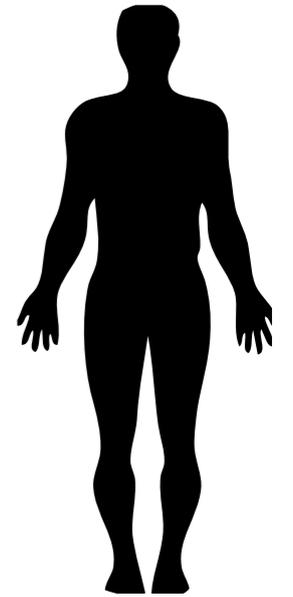
# Isoform numbers



# Combinations of splice sites produce diverse protein isoforms.



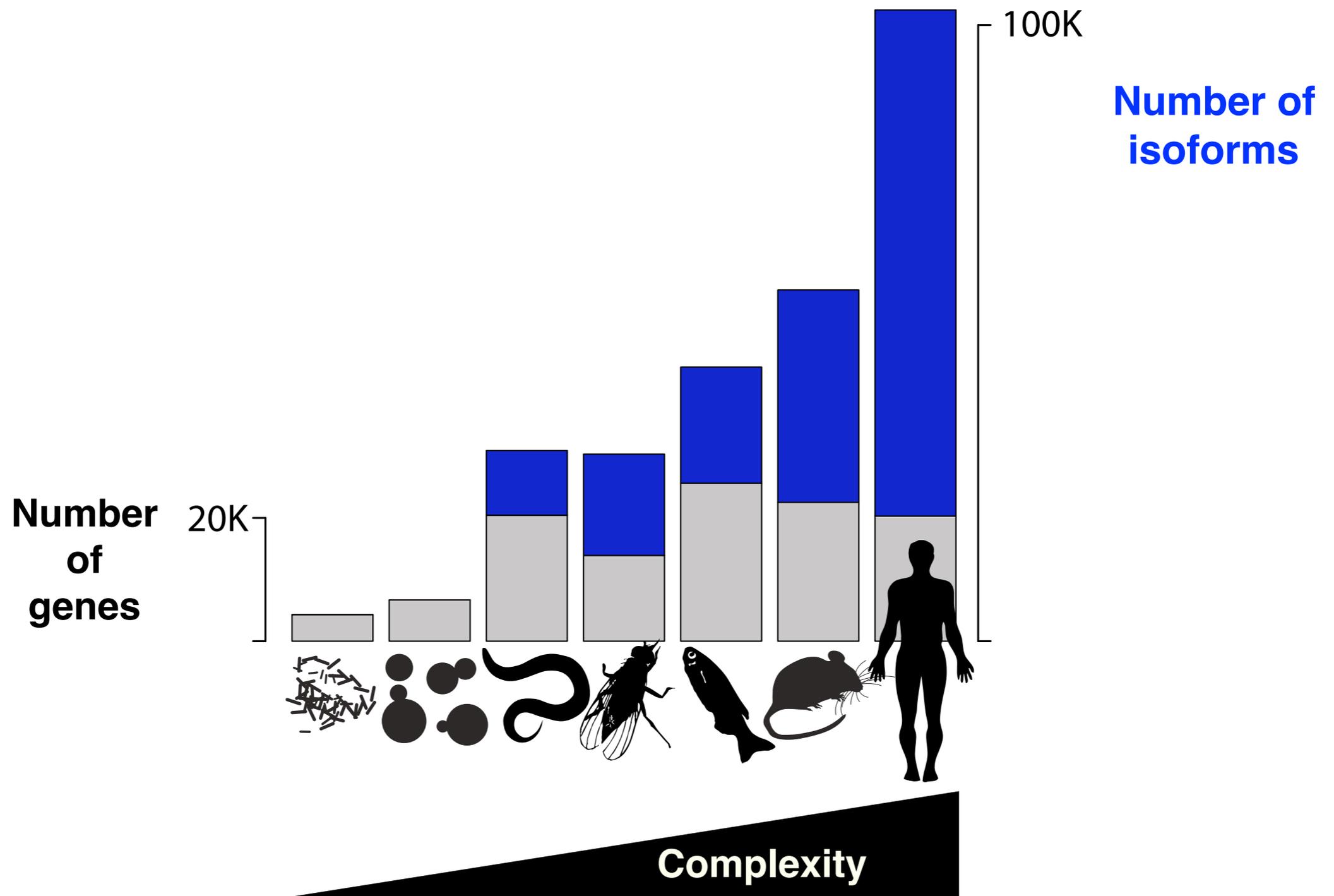
# The proteoform hypothesis



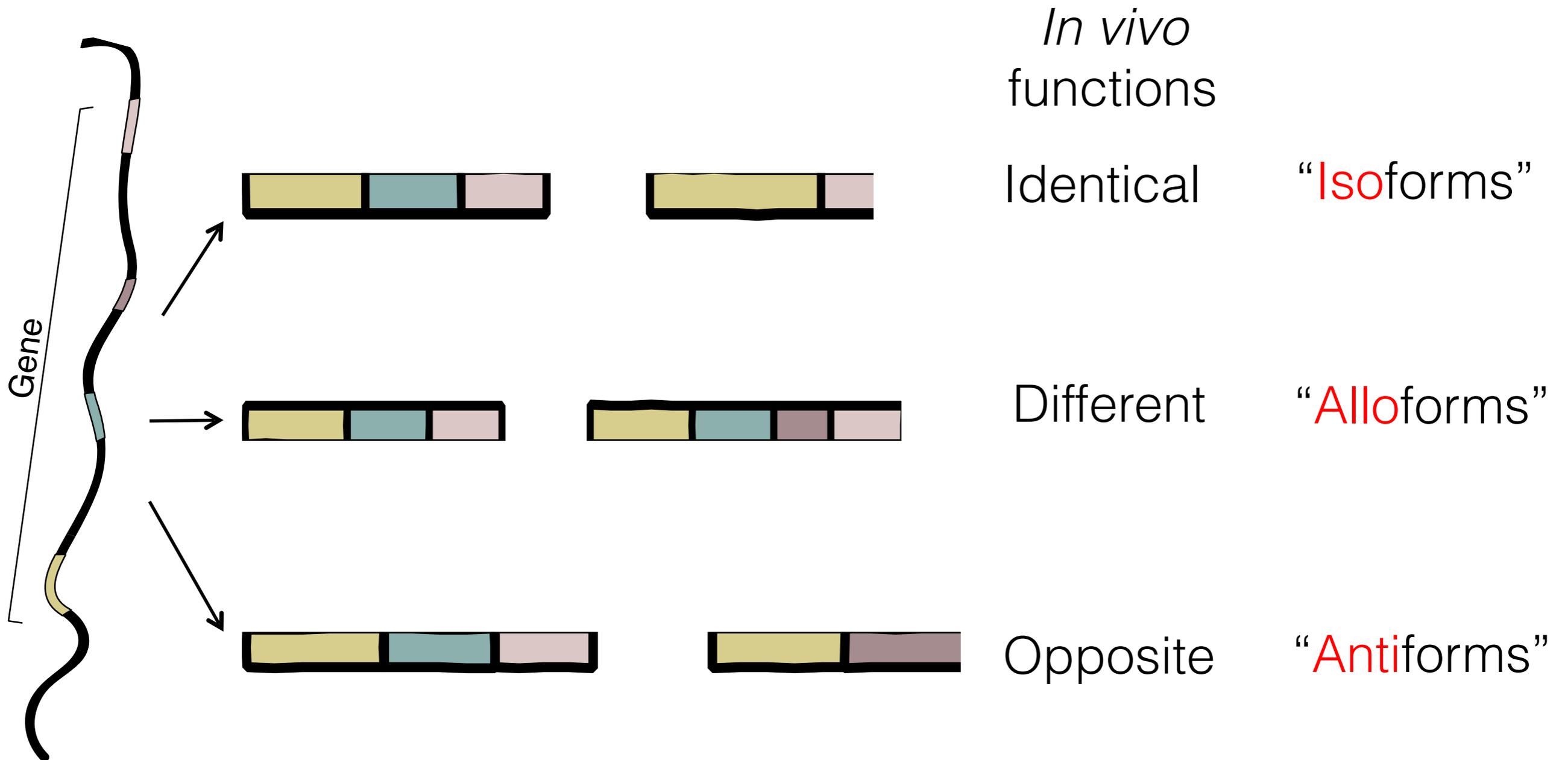
# Splicing regulation and disease

- splicing is pervasive, inherent to encoded products of the genome
- splicing is highly regulated in space and time
  - high tissue- and developmental- specificity
  - “splice code”
- splicing is dysregulated in many diseases, including cancer
  - estimates of 50% all disease variants affect splicing
  - splice-modulating therapies (e.g. antisense oligos)

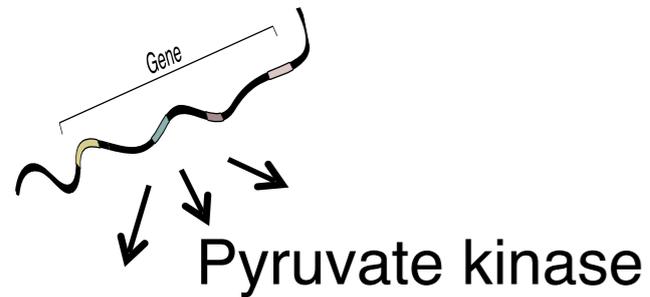
# Isoform function?



# Isoforms and functional divergence



# Examples of functionally divergent isoforms



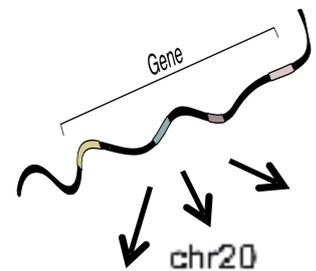
## Alloforms

chr15:72,200,000-72,230,000 (introns not to scale)

M1 isoform (↑adult tissues)



M2 isoform (↑embryo, ↑tumor)



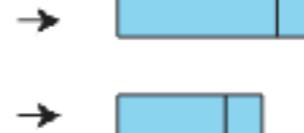
## Antiforms

31,723,000 bp    31,722,000 bp    31,666,000 bp    31,665,000 bp

Bcl2L1

Bcl-X(L)

Bcl-X(S)



Bcl-X **anti**-apoptotic

Bcl-X **pro**-apoptotic

# Divergent functional capabilities described in literature

Isoforms for a  
few hundred genes



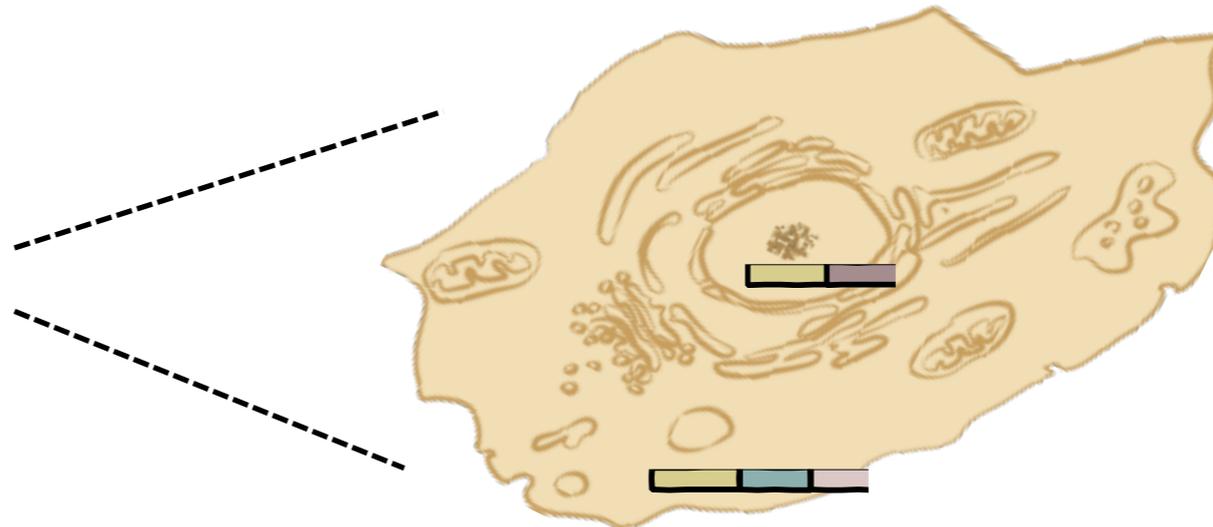
Physical interactions

Cellular localization

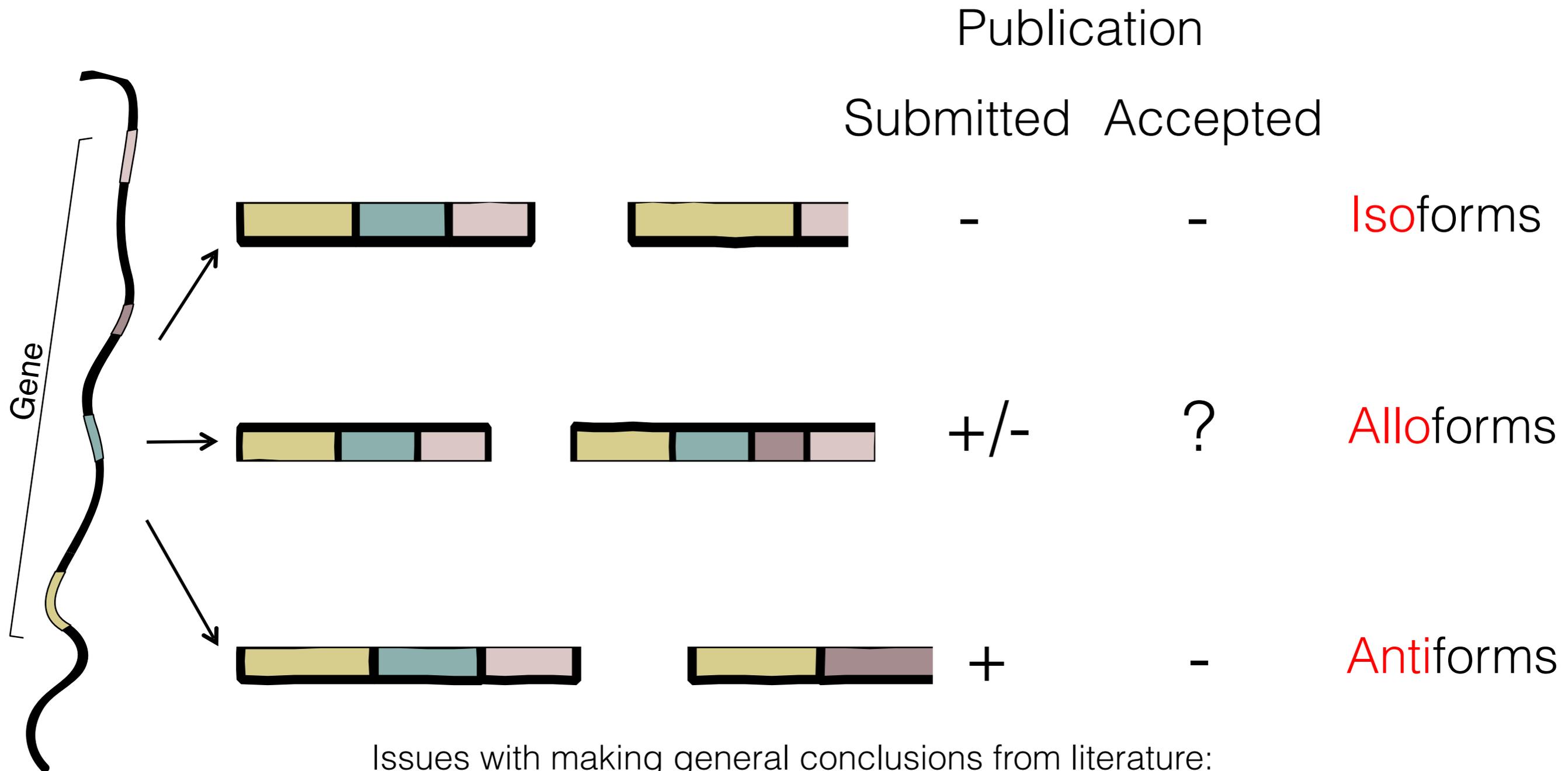
Enzymatic activities

Stability

.....



# Sociological biases in literature



Issues with making general conclusions from literature:  
-confirmation bias (sampling not random)  
-experimental approaches for characterization heterogenous  
-isoform identity unknown

# How widespread is isoform functional divergence in the whole proteome?

Systematic identification  
of large numbers of  
isoform pairs  
for large numbers  
of human genes



Physical interactions

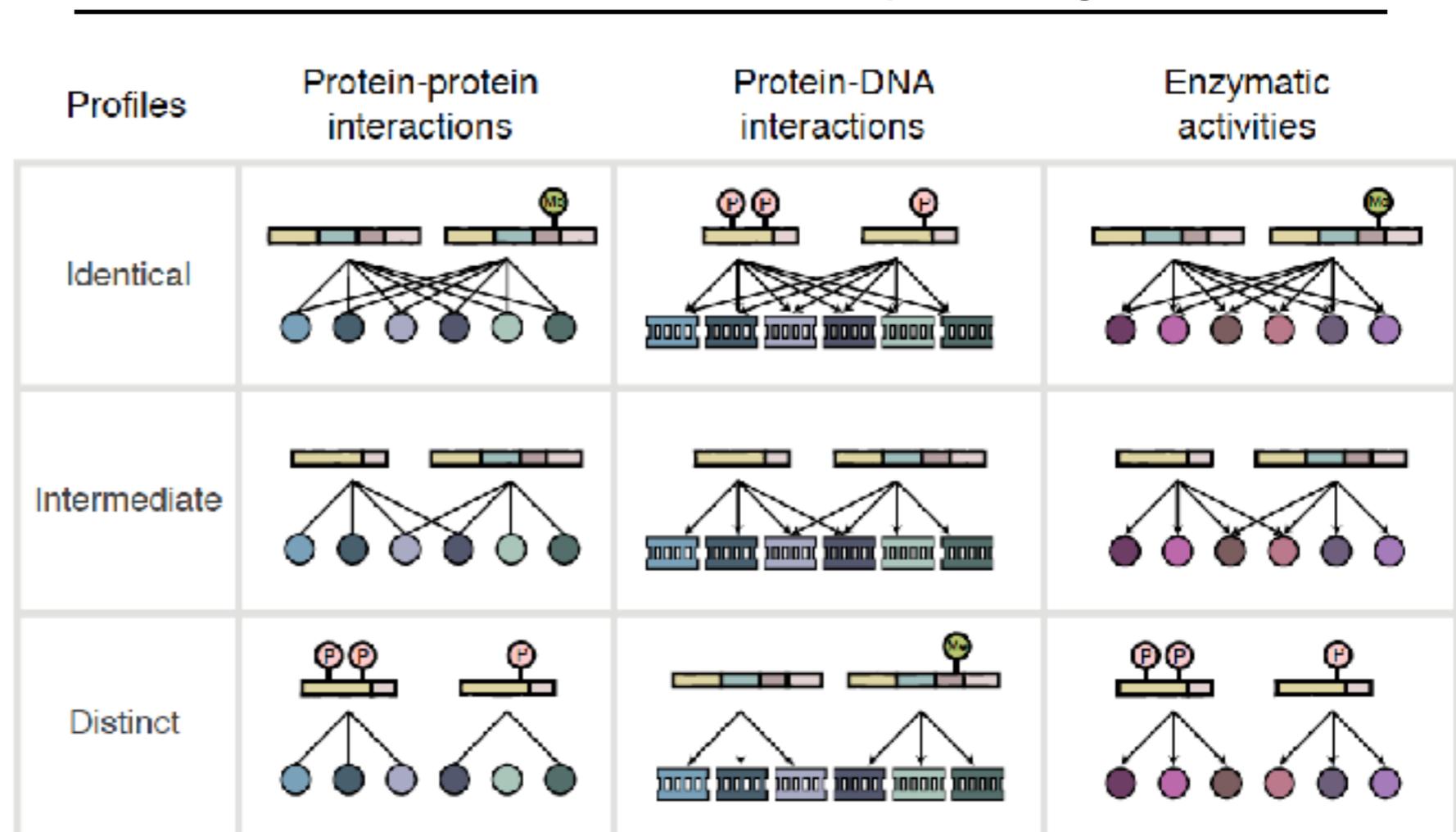
Enzymatic activities

Cellular localization

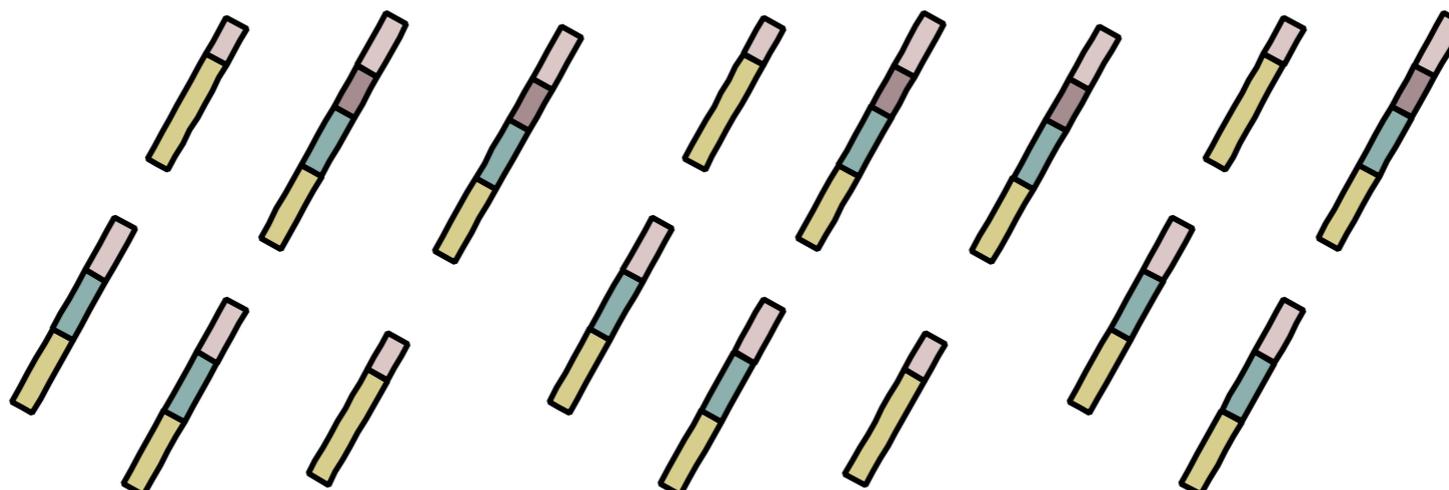
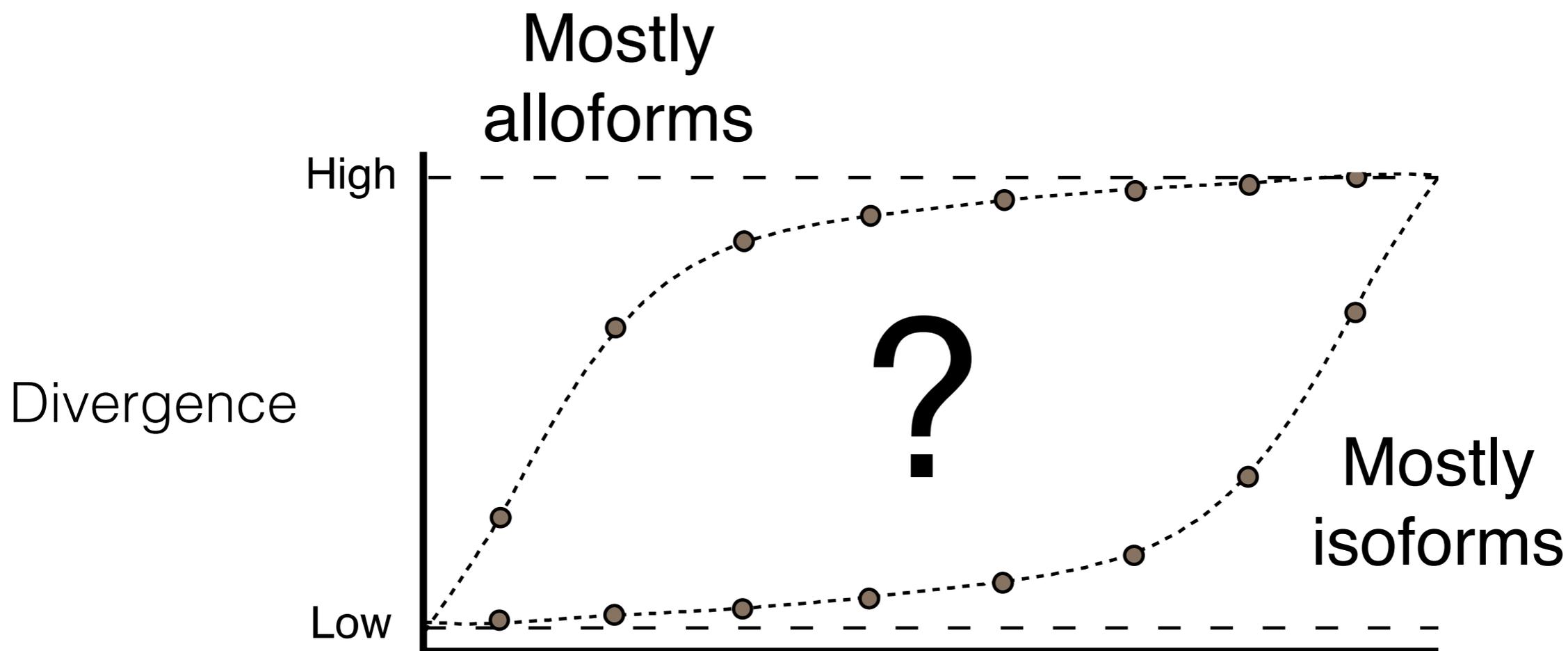
Stability

.....

## Unbiased functional profiling



# Landscape of protein isoform functional divergence



Large numbers of pairs of isoforms encoded by common genes

# How widespread is isoform functional divergence in the whole proteome?

Systematic identification of large numbers of isoform pairs for large numbers of human genes



Physical interactions

Enzymatic activities

Cellular localization

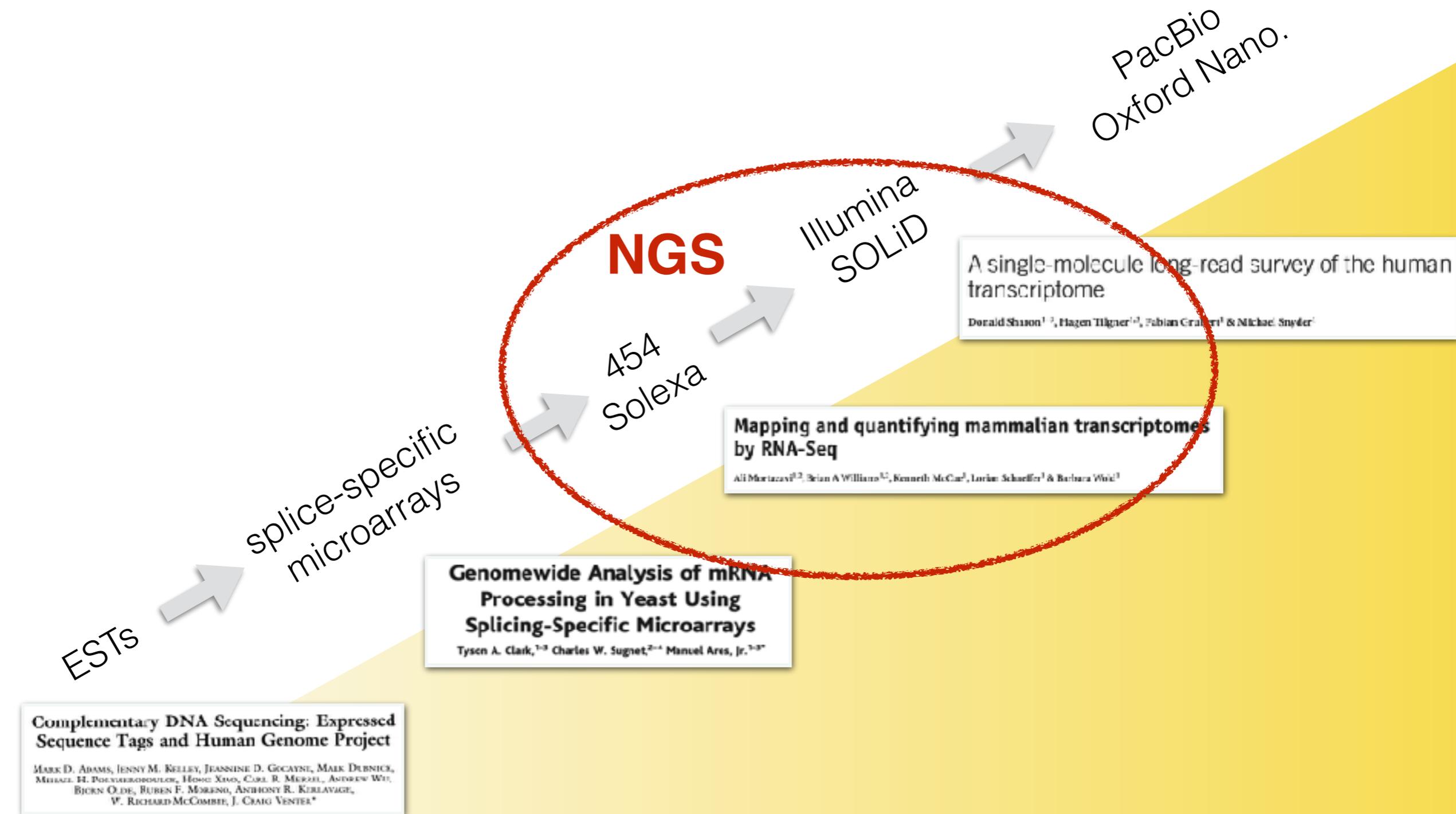
Stability

.....

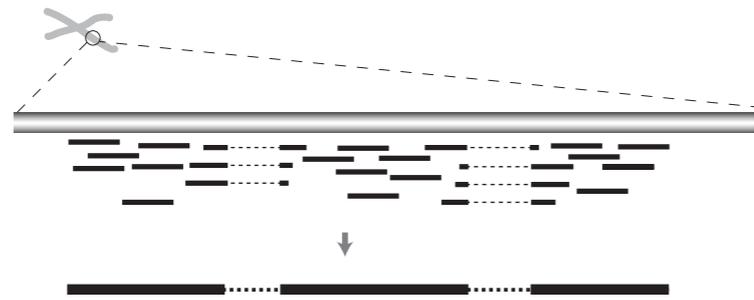
## Unbiased functional profiling

Profiles	Protein-protein interactions	Protein-DNA interactions	Enzymatic activities
Identical			
Intermediate			
Distinct			

# RNA sequencing data has been the primary means to characterize isoforms



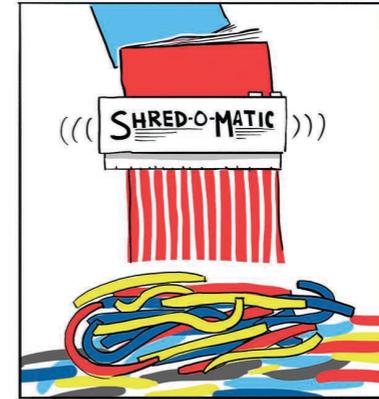
NGS (RNA-Seq) data can reveal the presence of exons and junctions, but fails to accurately reconstruct full-length isoforms.



transcript reconstruction



transcript assembly



Korf *Nature Methods* 2013

## Assessment of transcript reconstruction methods for RNA-seq

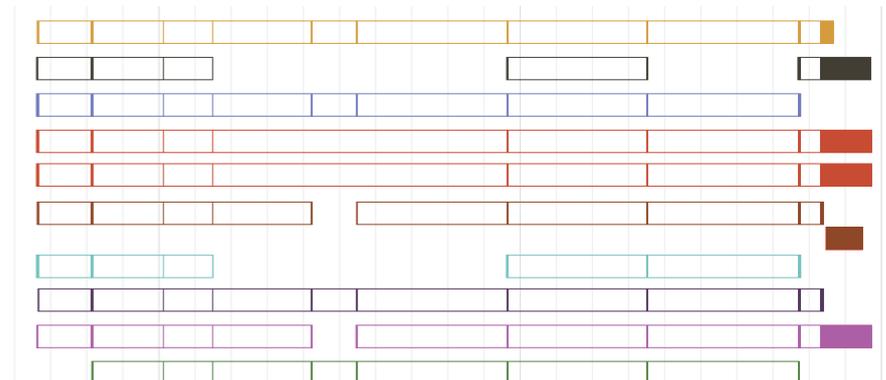
Steijger et al *Nature Methods* 2013

- 14 reconstruction/assembly methods evaluated
- high method variability
- poor performance on simulated data

Augustus all  
Cufflinks  
iReckon full  
mGene  
mGene graph  
mTim  
SLIDE all  
Transomics all  
Trembly all  
Tromer

RPKM (scaled)

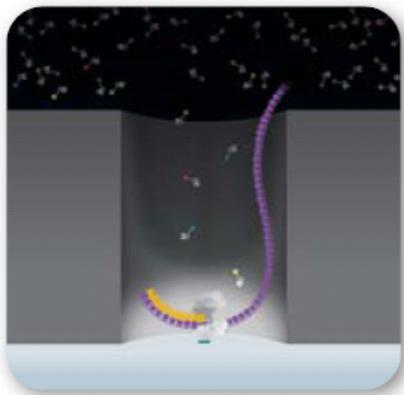
10.55  
42.92  
13.23  
7.66  
6.54  
4.42  
30.70  
7.05  
6.34  
4.21



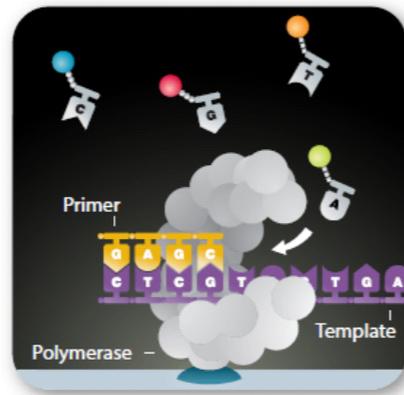
# Iso-Seq enables direct sequencing of full-length isoforms and thus characterize transcriptome complexity

PacBio  
detection: fluorescence

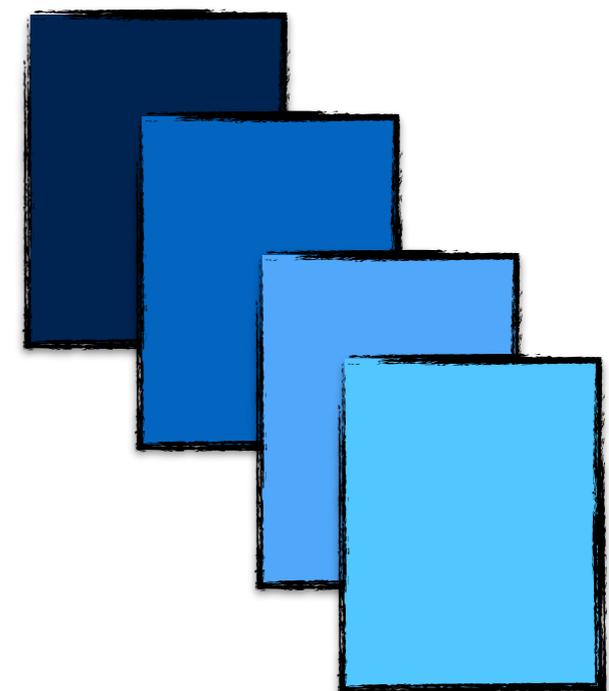
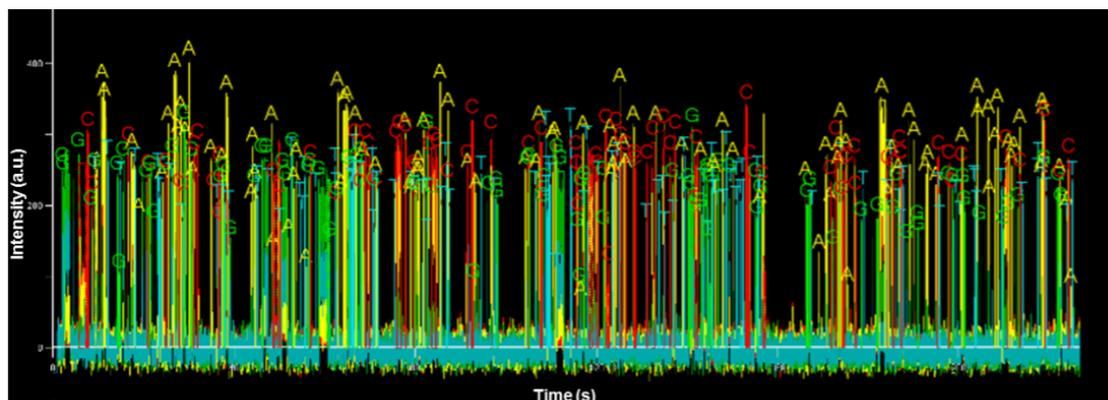
Zero-Mode Waveguides



Phospholinked Nucleotides



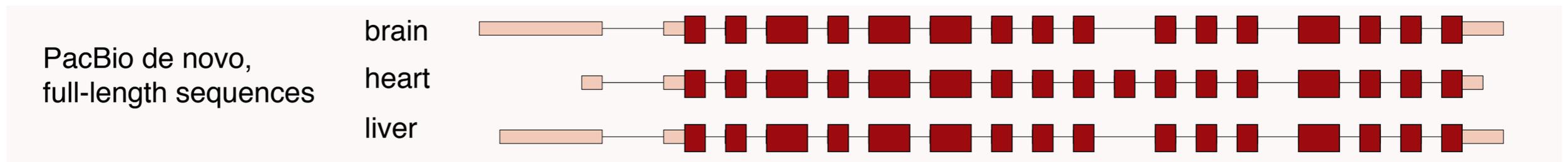
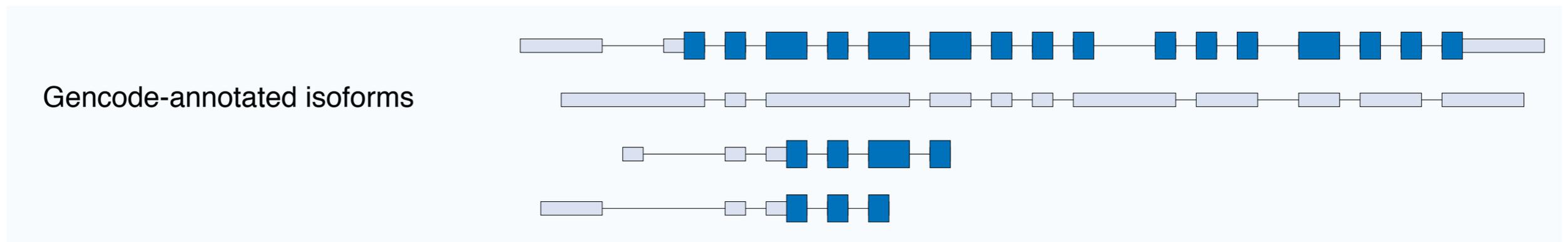
Trace



Given a base mRNA sequence from Iso-Seq, one can predict open reading frames or “ORFs”

These become candidate protein isoforms for downstream functional analysis

Gene:SSRP1, Strand: - chr11:57335895-57327743



□ □ UTRs    ■ ■ ■ CDS'

# The challenge of detecting low abundance isoforms

Wide dynamic range of human transcriptomes

Limited knowledge of isoform complexity for low-abundance transcripts

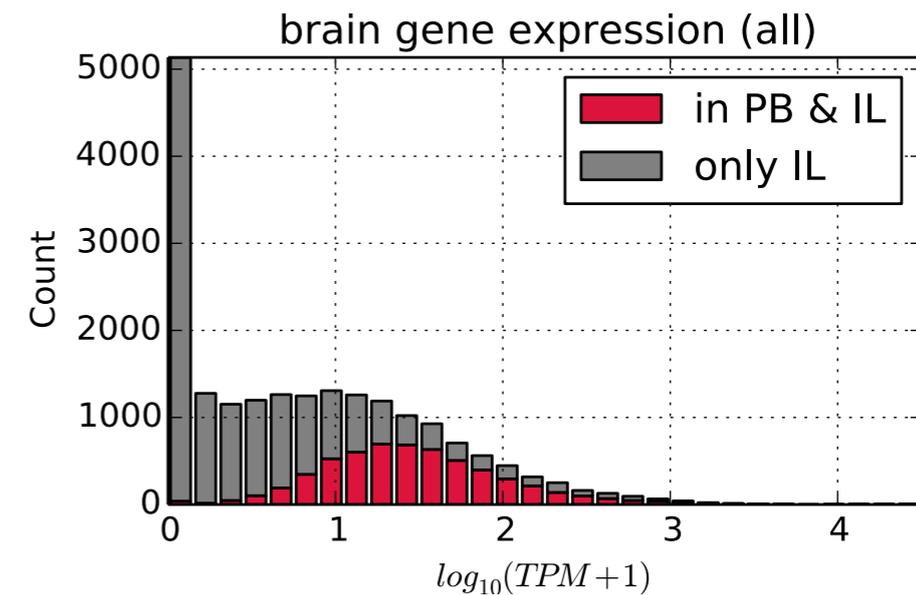
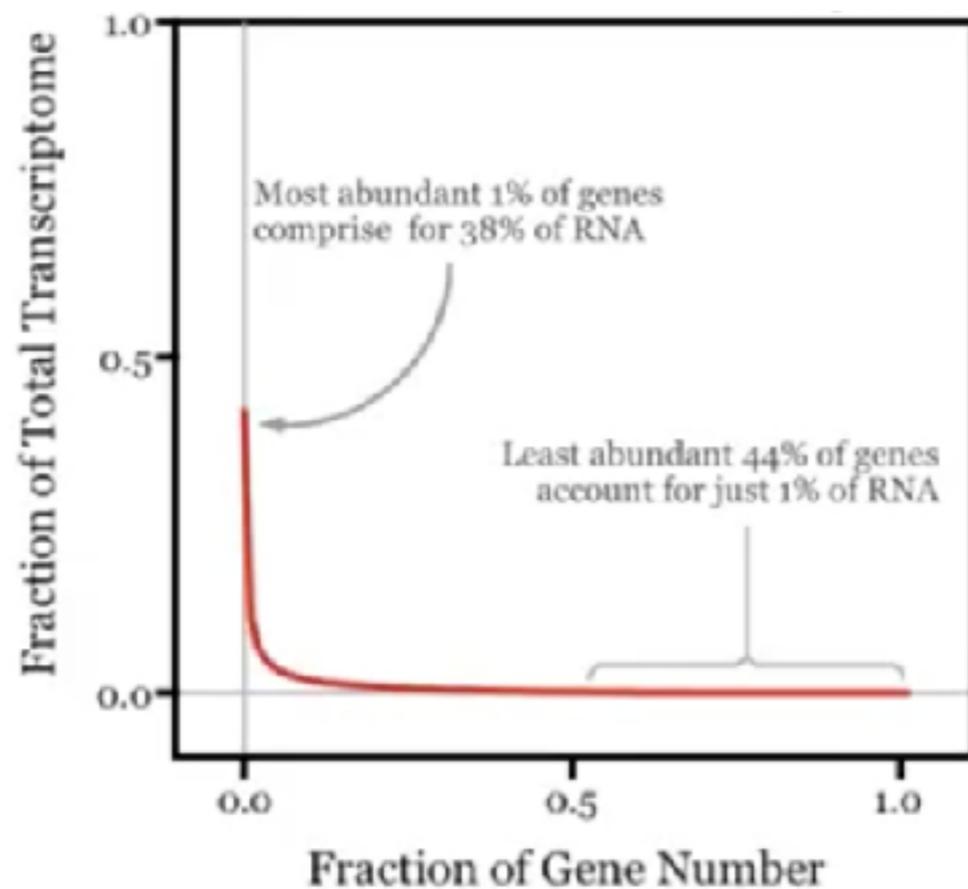
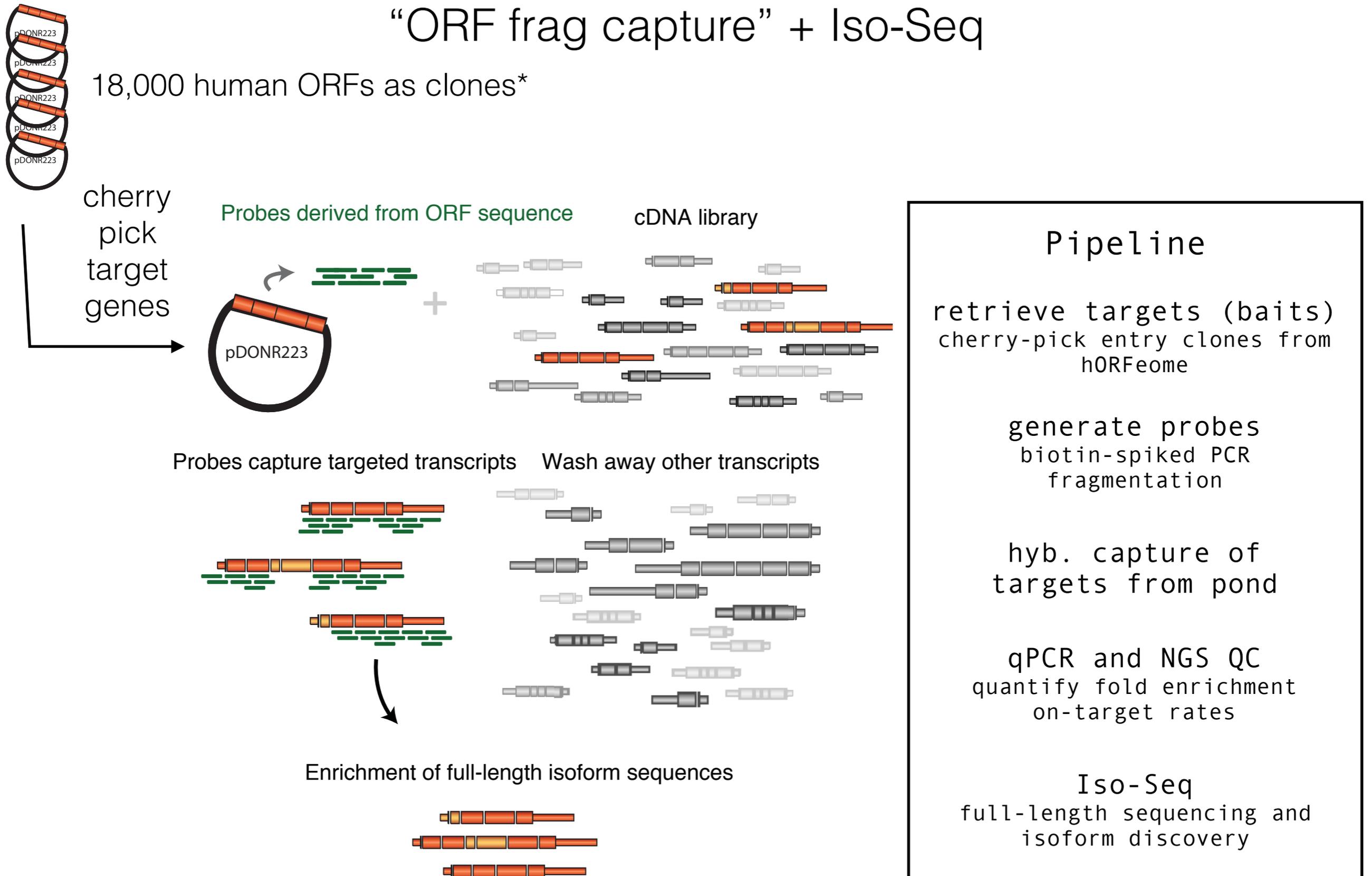


Figure from Tim Mercer (Capture-Seq)

IL = Illumina data  
PB = PacBio data

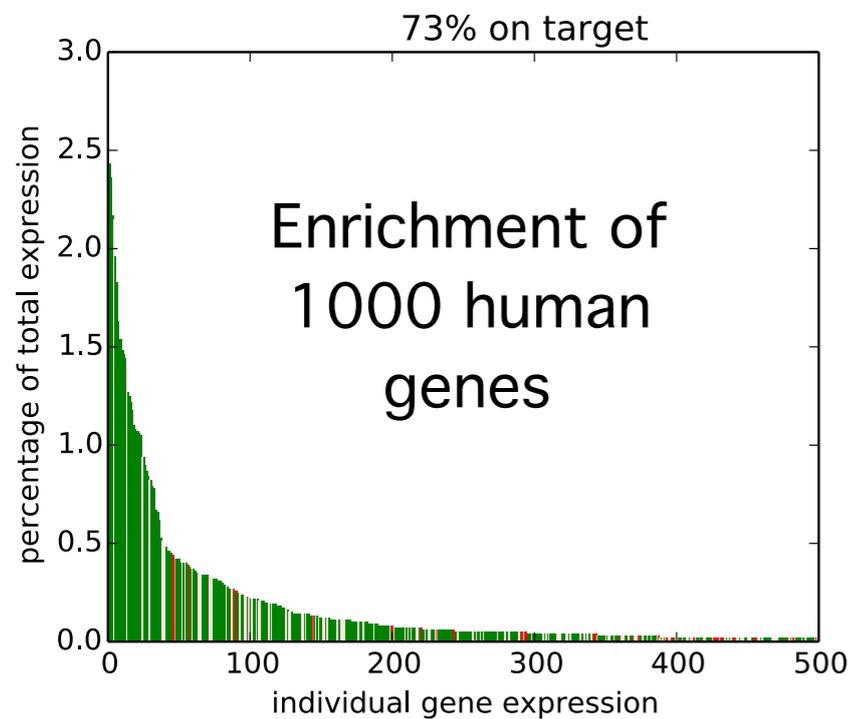
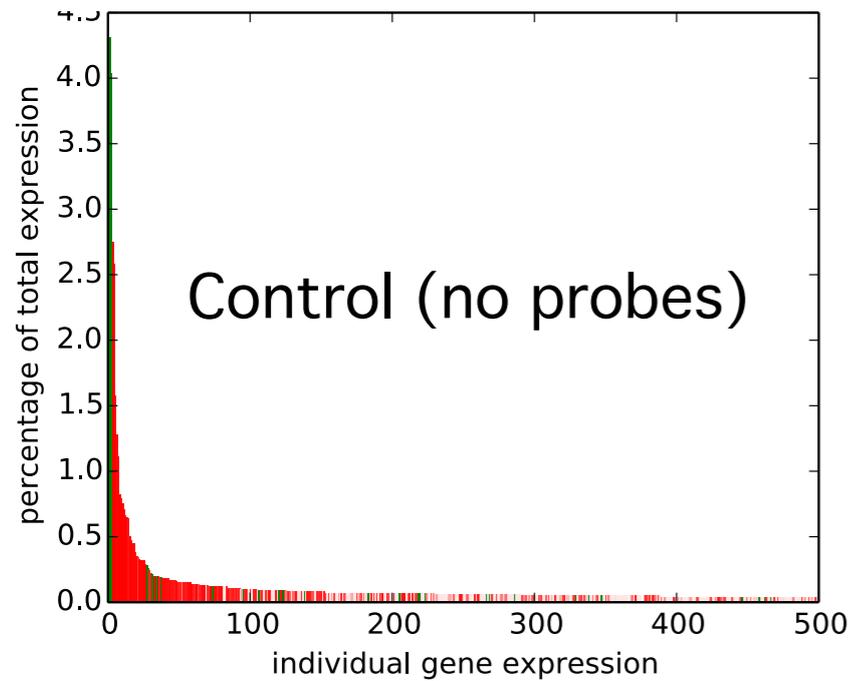
# We developed a new method to generate probes en masse for use in hybridization-based target capture: “ORF frag capture” + Iso-Seq



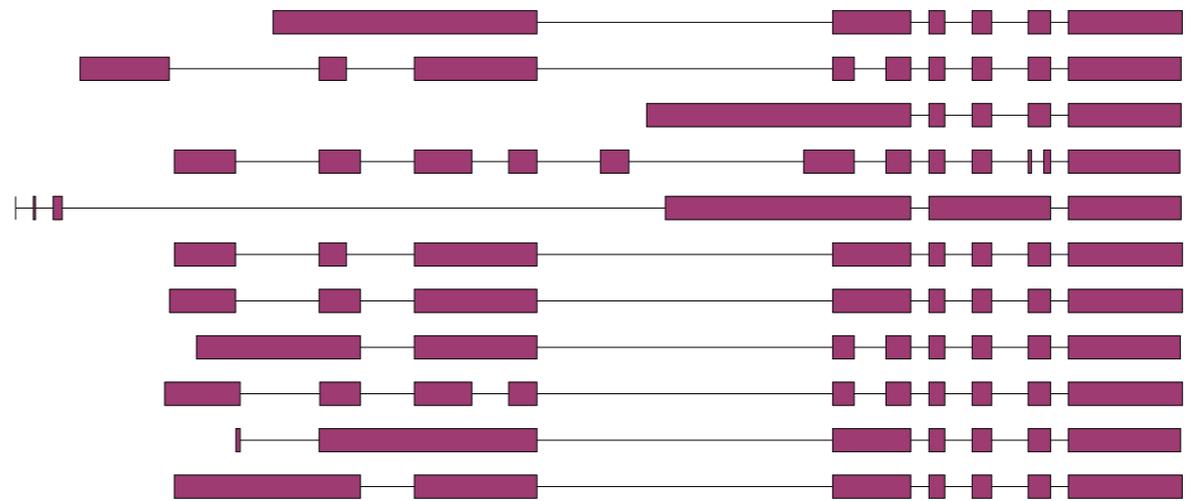
\*The ORFeome Collaboration *Nature Methods* (2016)

w Jason Underwood, Tyson Clark

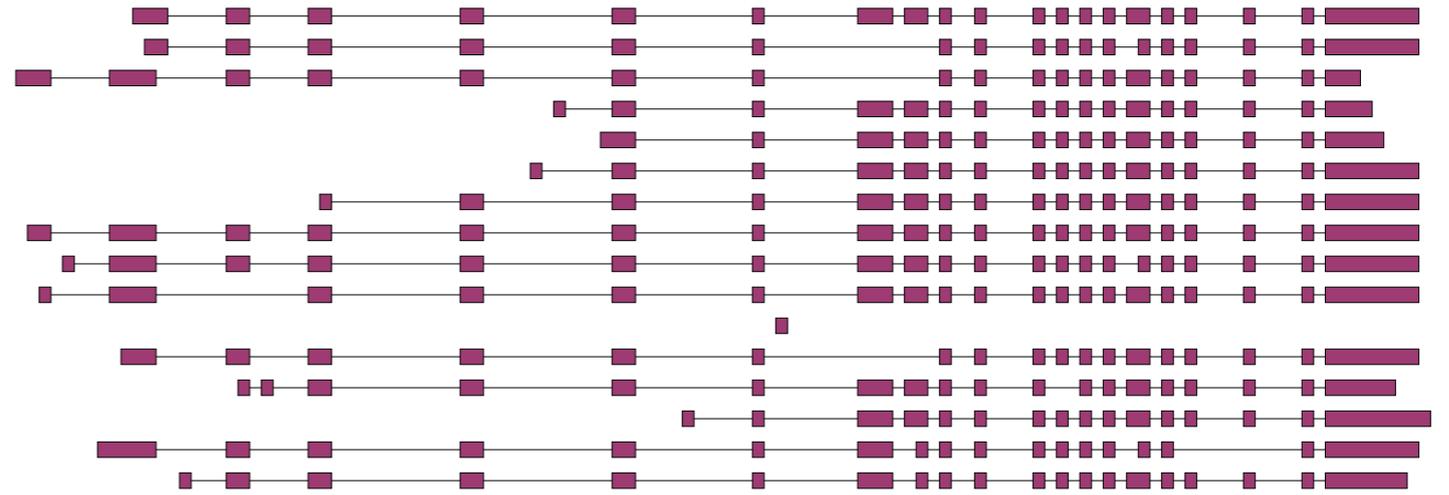
# “ORF frag capture” + Iso-Seq enables of high sensitivity discovery of full-length isoforms



## CREB3L4 isoforms detected\*



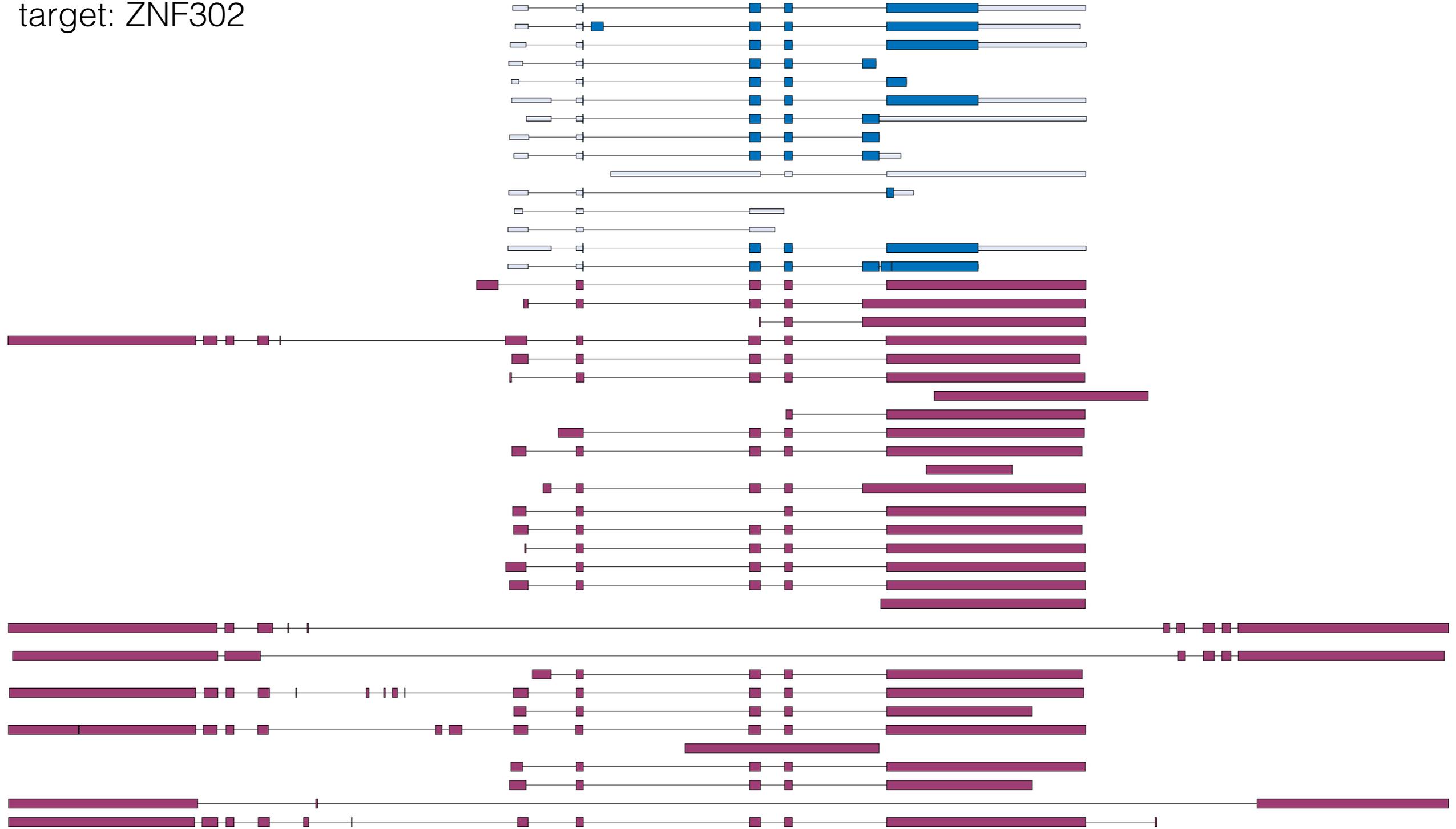
## FOXP1 isoforms detected\*



\*No isoforms detected in ultra-deep-coverage PacBio dataset of the same brain cDNA library.

# “ORF frag-based” probe capture, sequencing, and discovery of full-length isoforms

target: ZNF302



# How widespread is isoform functional divergence in the whole proteome?

Systematic identification of large numbers of isoform pairs for large numbers of human genes



Physical interactions

Enzymatic activities

Cellular localization

Stability

.....

## Unbiased functional profiling

Profiles	Protein-protein interactions	Protein-DNA interactions	Enzymatic activities
Identical			
Intermediate			
Distinct			

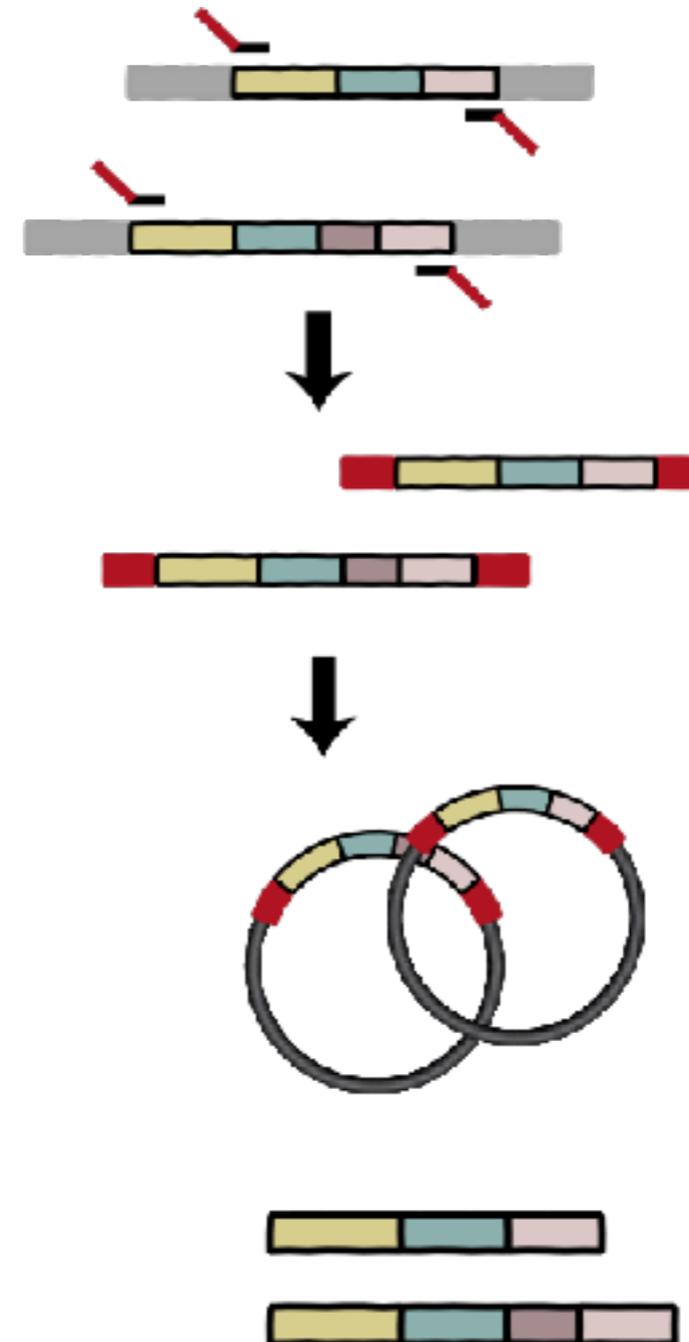
# ORF-Seq:

## Systematic full-length isoform sequencing and cloning

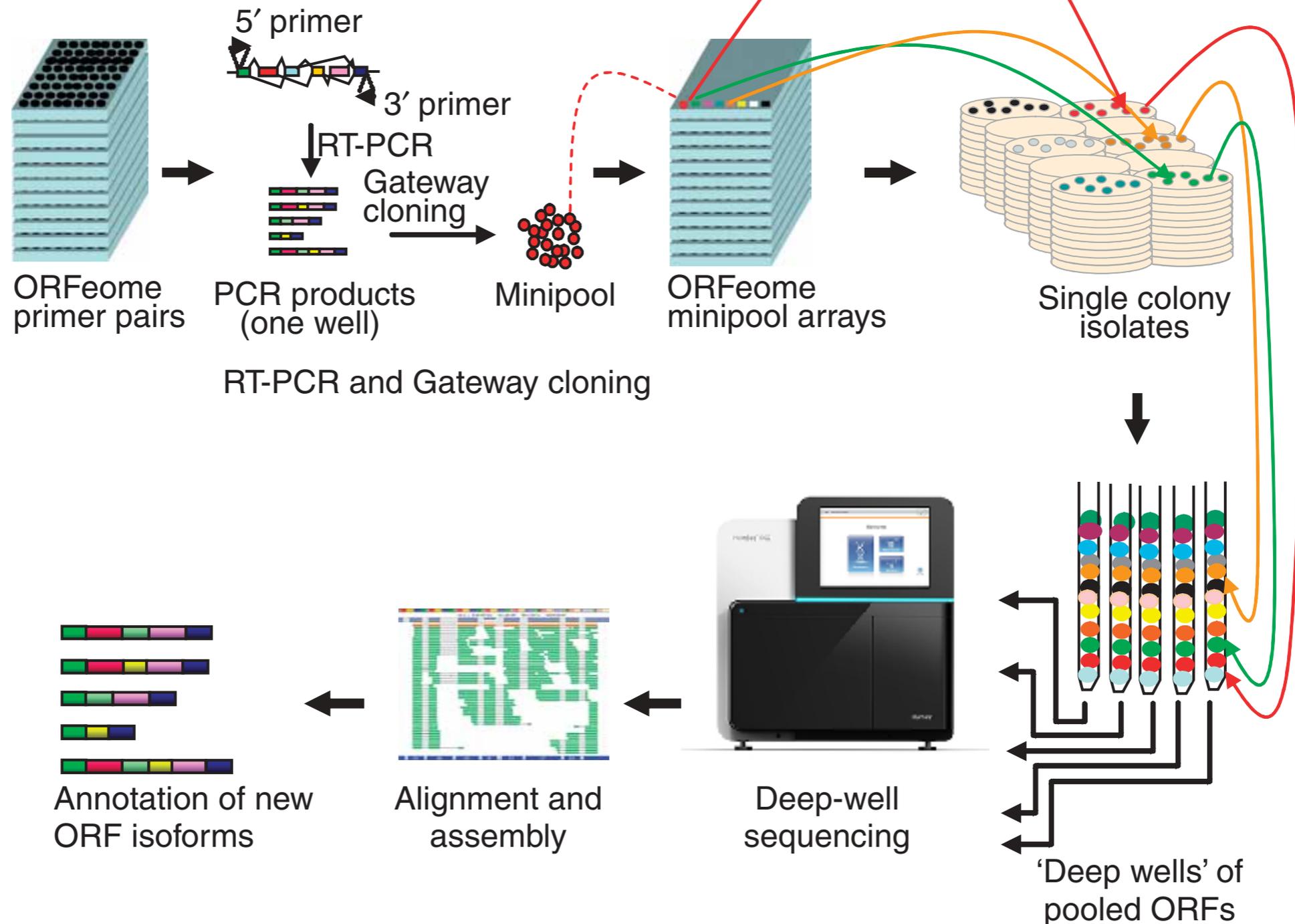
### Isoform discovery by targeted cloning, 'deep-well' pooling and parallel sequencing

Kourosh Salehi-Ashtiani<sup>1,2,5</sup>, Xinpeng Yang<sup>1,2,5</sup>,  
Adnan Derti<sup>1,3,5</sup>, Weidong Tian<sup>1,3,5</sup>, Tong Hao<sup>1,2,5</sup>,  
Chenwei Lin<sup>1,2</sup>, Kathryn Makowski<sup>4</sup>, Lei Shen<sup>4</sup>,  
Ryan R Murray<sup>1,2</sup>, David Szeto<sup>1,2</sup>, Nadeem Tusneem<sup>4</sup>,  
Douglas R Smith<sup>4</sup>, Michael E Cusick<sup>1,2</sup>, David E Hill<sup>1,2</sup>,  
Frederick P Roth<sup>1,3</sup> & Marc Vidal<sup>1,2</sup>

**NATURE METHODS** | VOL.5 NO.7 | JULY 2008

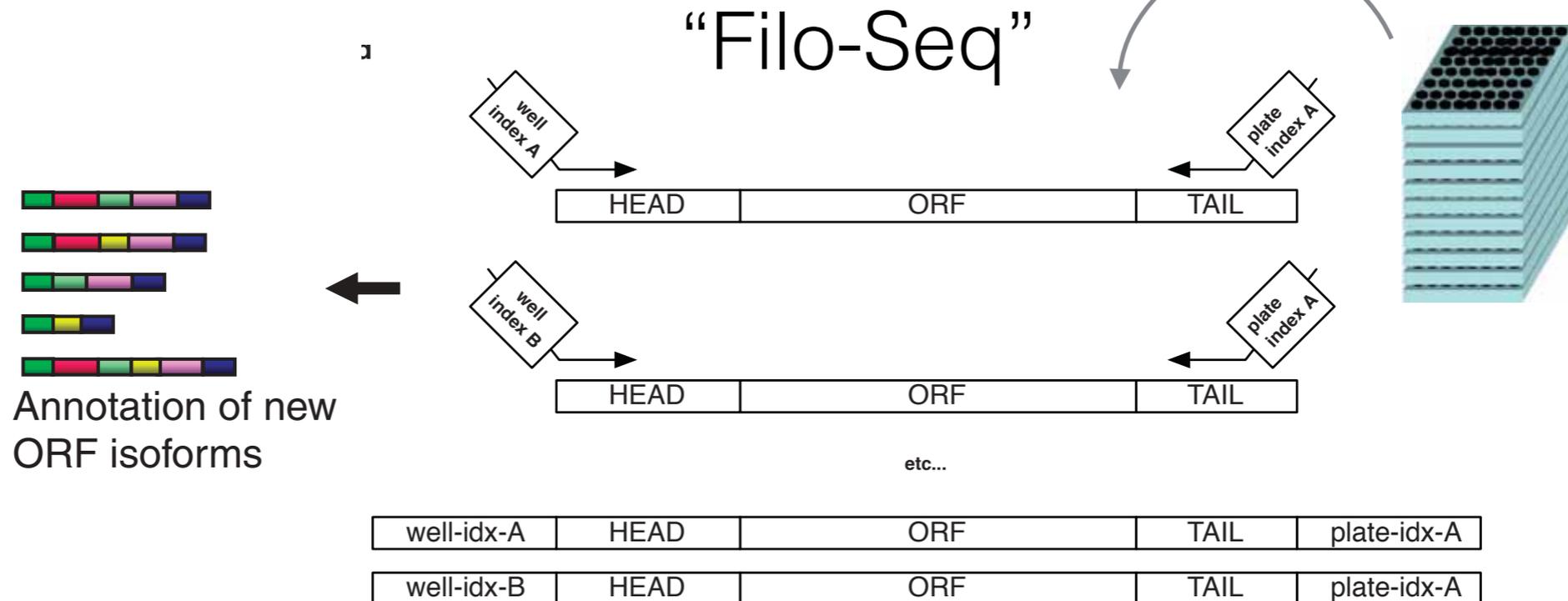
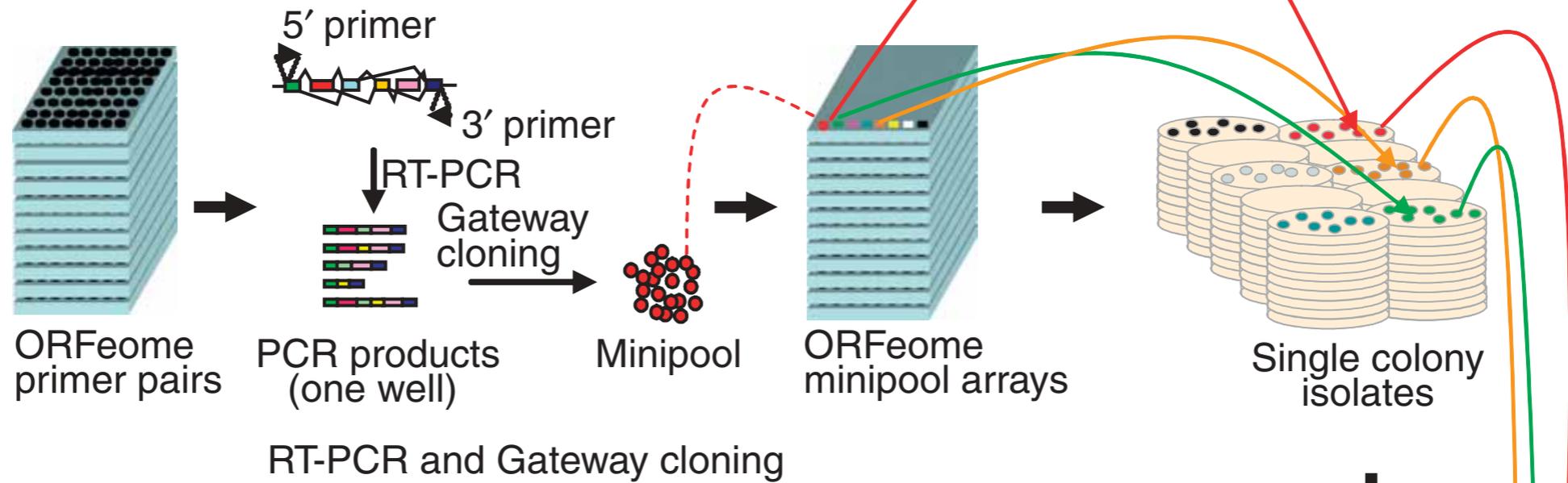


# ORF-Seq: Systematic full-length isoform sequencing and cloning

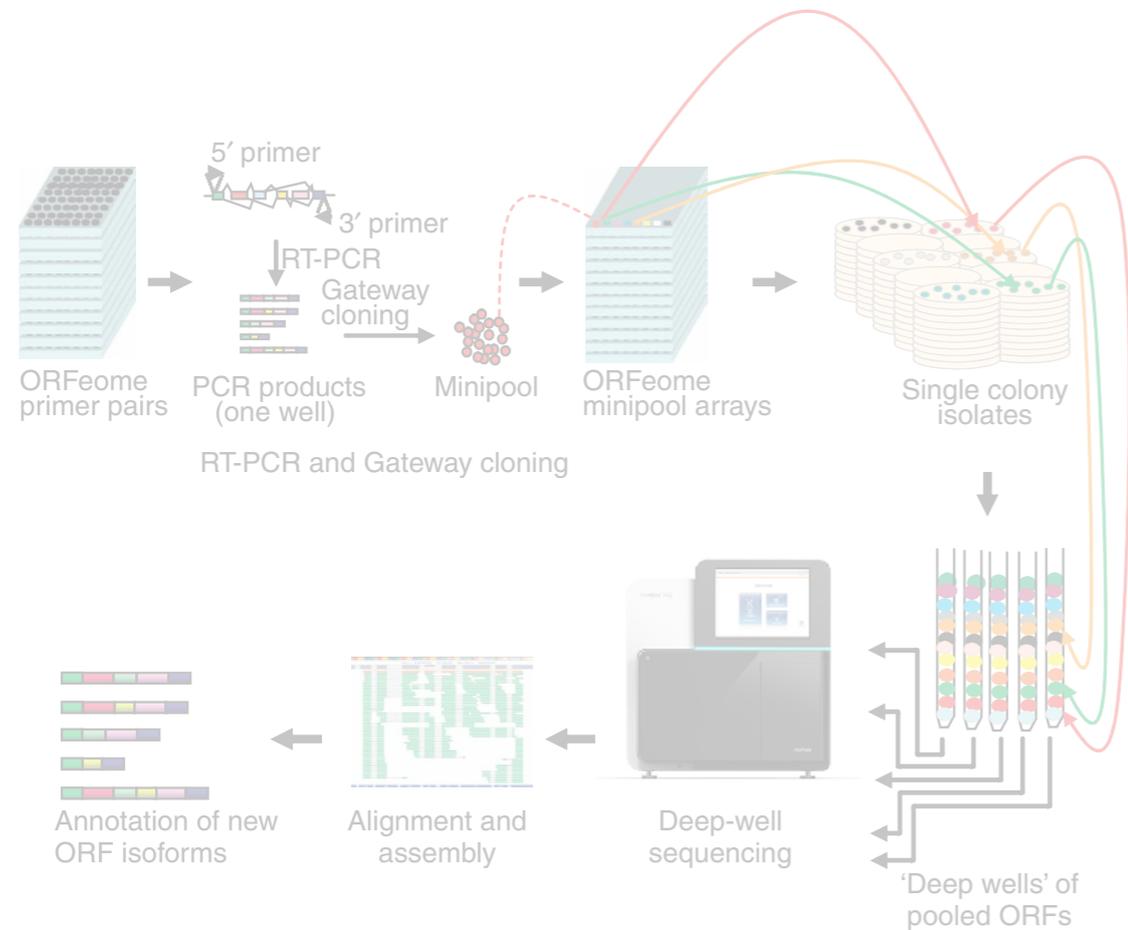


# ORF-Seq:

## Systematic full-length isoform sequencing and cloning



# ORF-Seq: Systematic full-length isoform sequencing and cloning



~1,500 genes, 5 human tissues

1,423 full-length isoforms

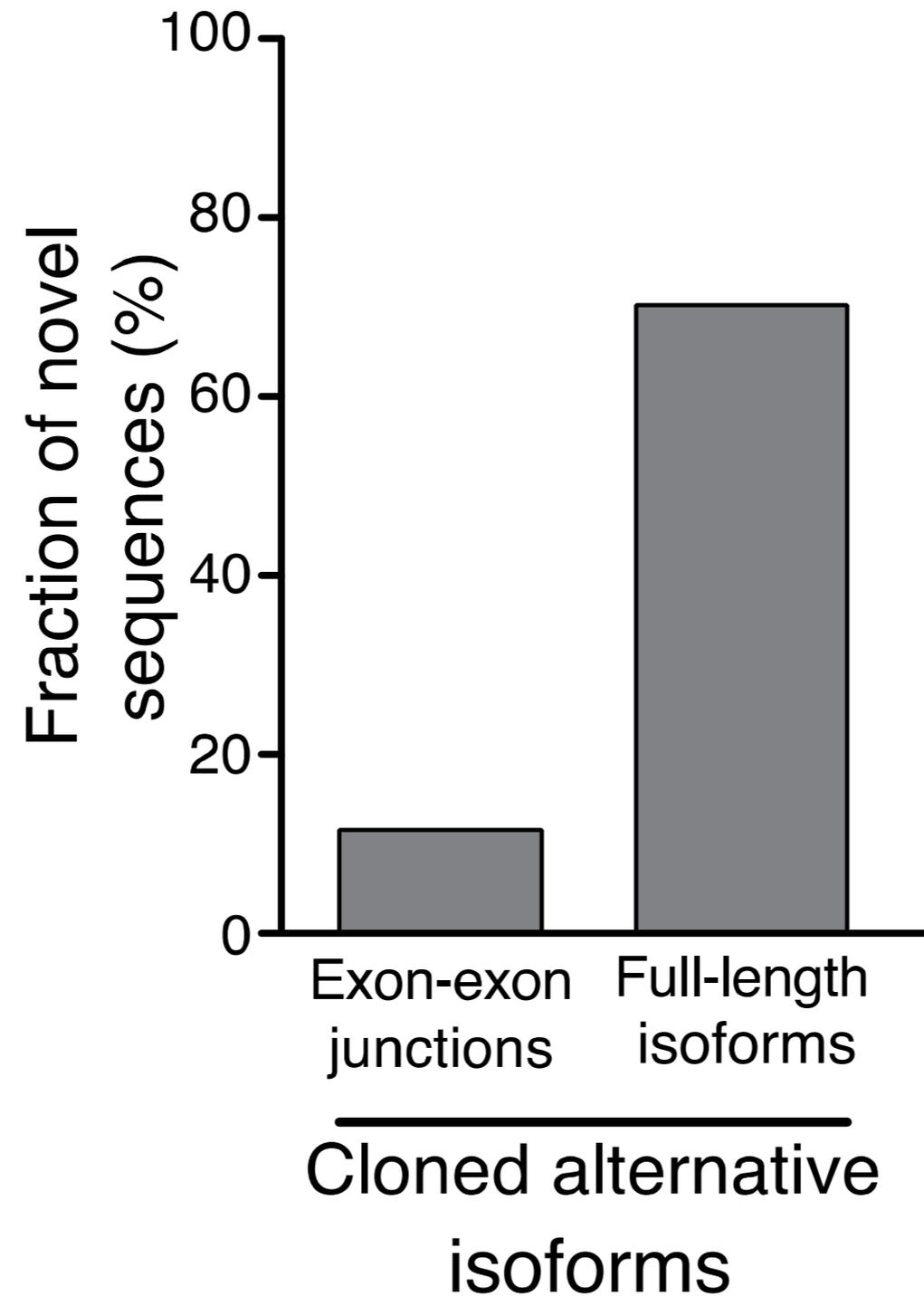
506 genes with  $\geq 2$  isoforms

1677 pairs of alternative isoforms,  
each encoded by a common gene

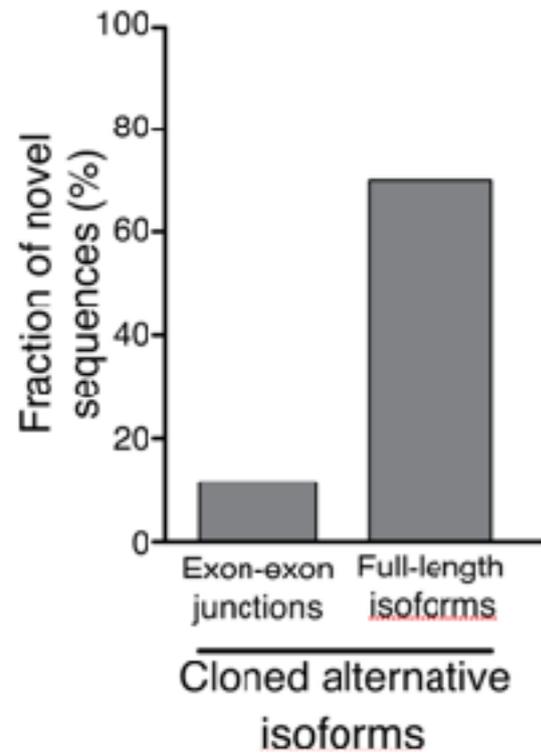


# Novel full-length isoforms

506 genes with  $\geq 2$  isoforms

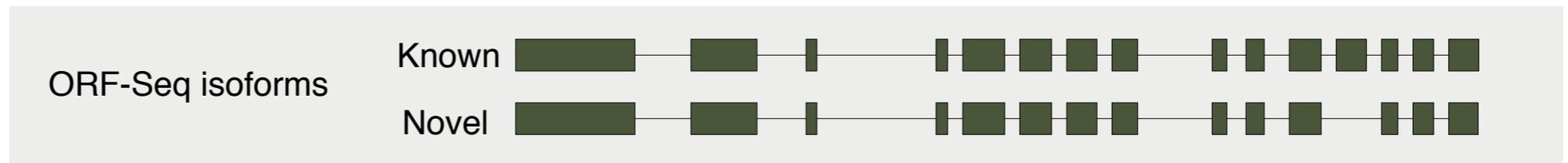
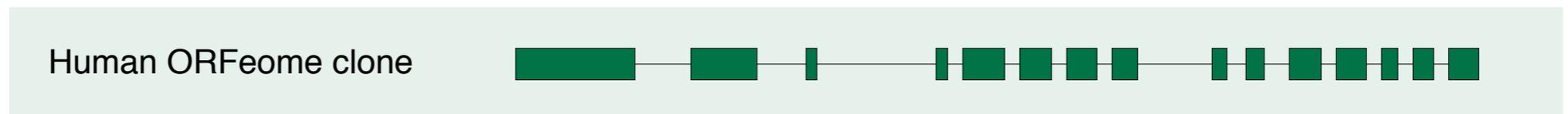
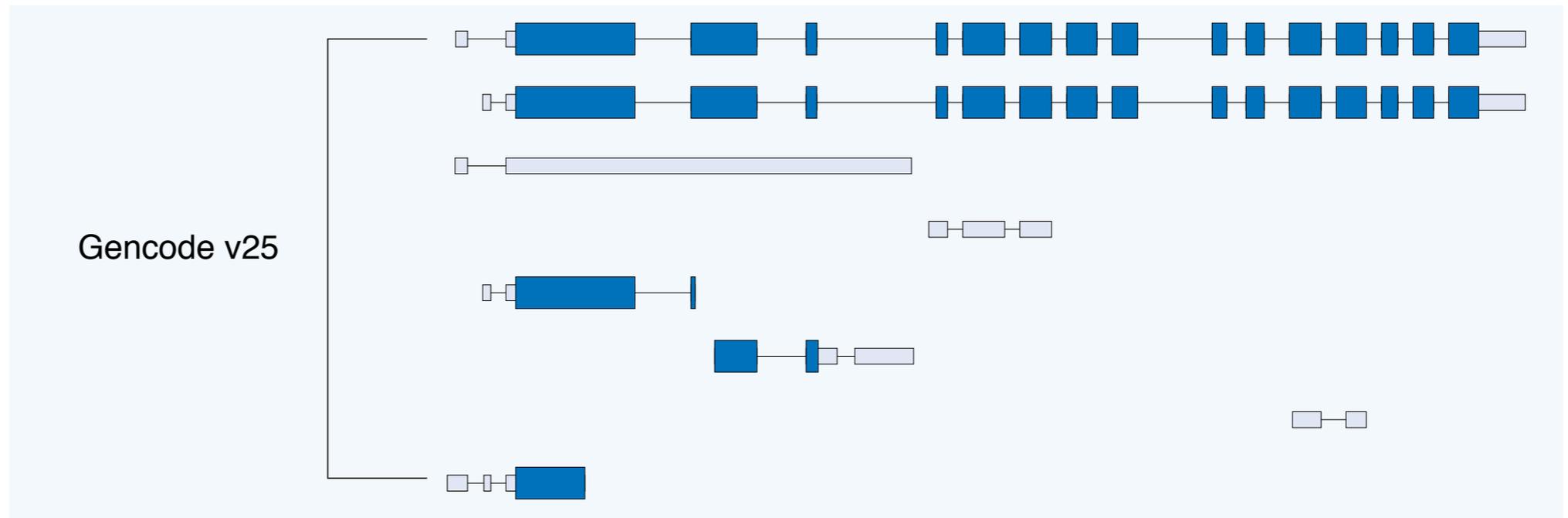


# Novel full-length isoforms



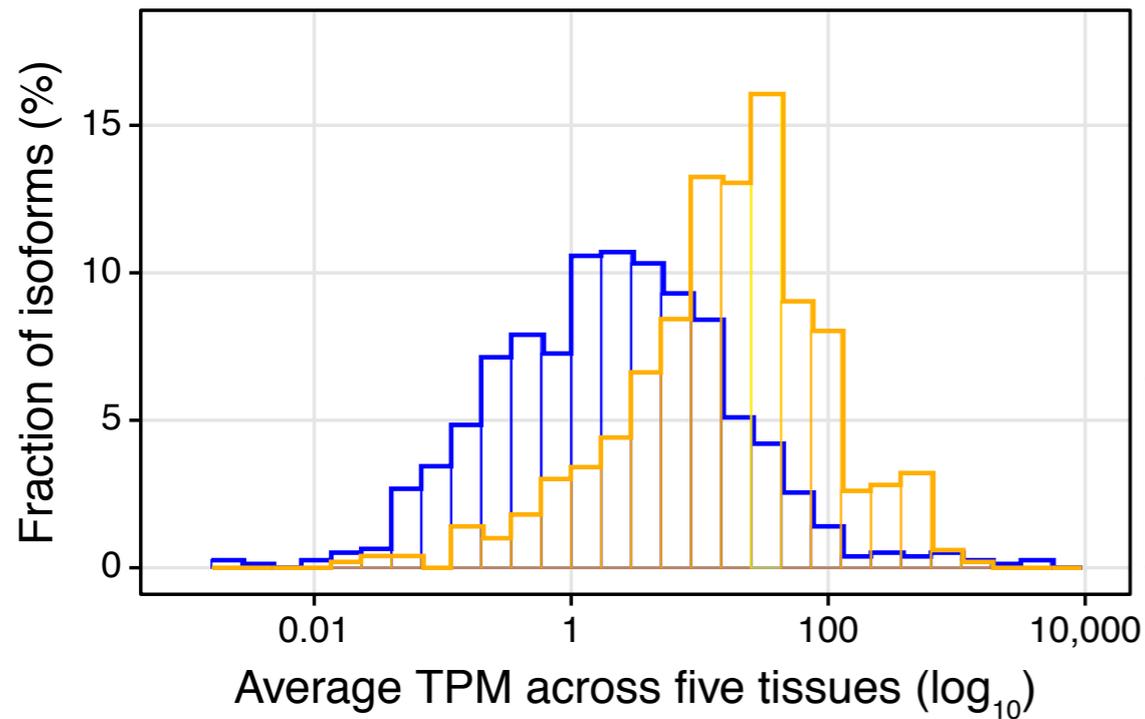
Gene: ARHGEF15

Strand: + chr17:8322516 - 8310241



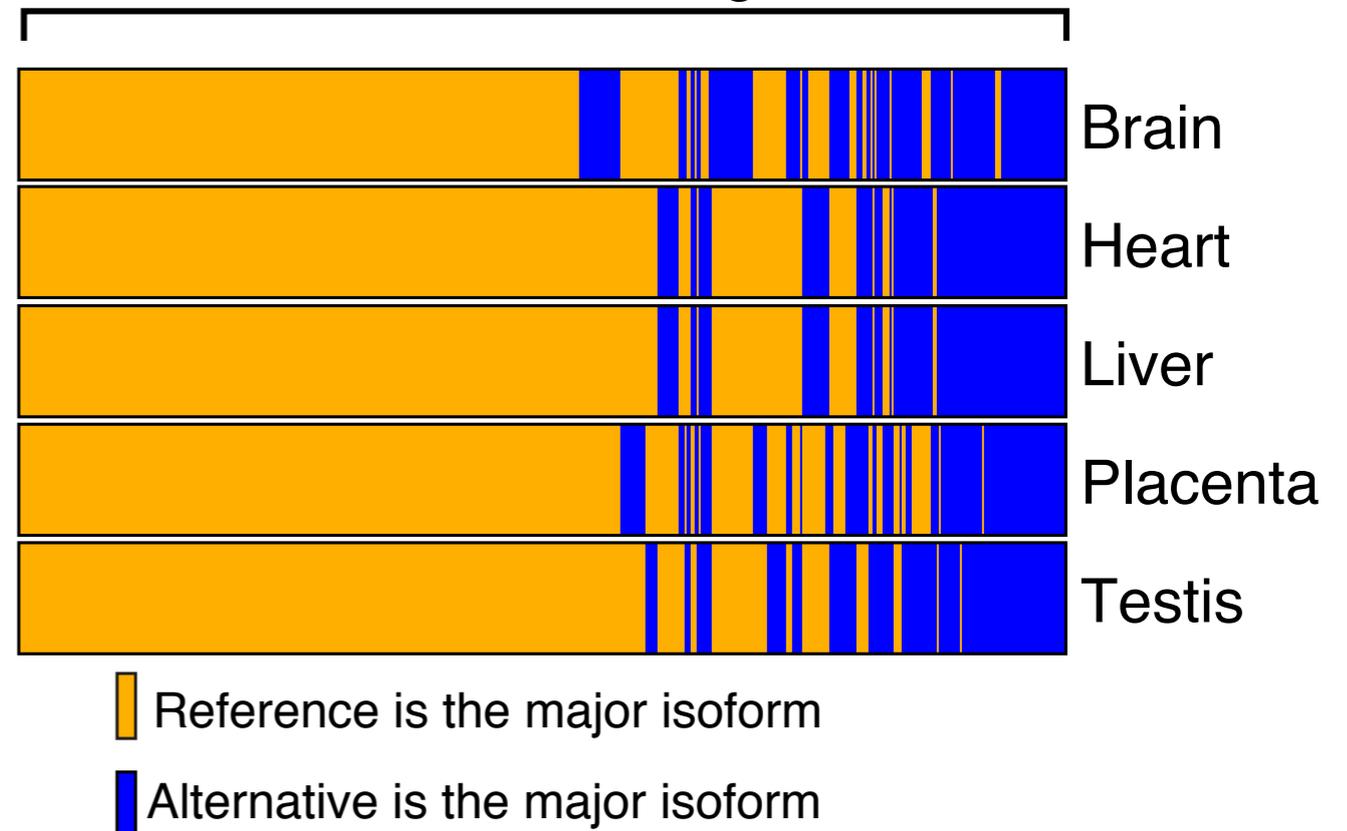
Intron lengths are not to scale.

# Estimated isoform clone expression in the original, endogenous tissues



Isoform category  
Reference  
Alternative

~500 multi-isoform genes



# How widespread is isoform functional divergence in the whole proteome?

Systematic identification  
of large numbers of  
isoform pairs  
for large numbers  
of human genes



Physical interactions

Enzymatic activities

Cellular localization

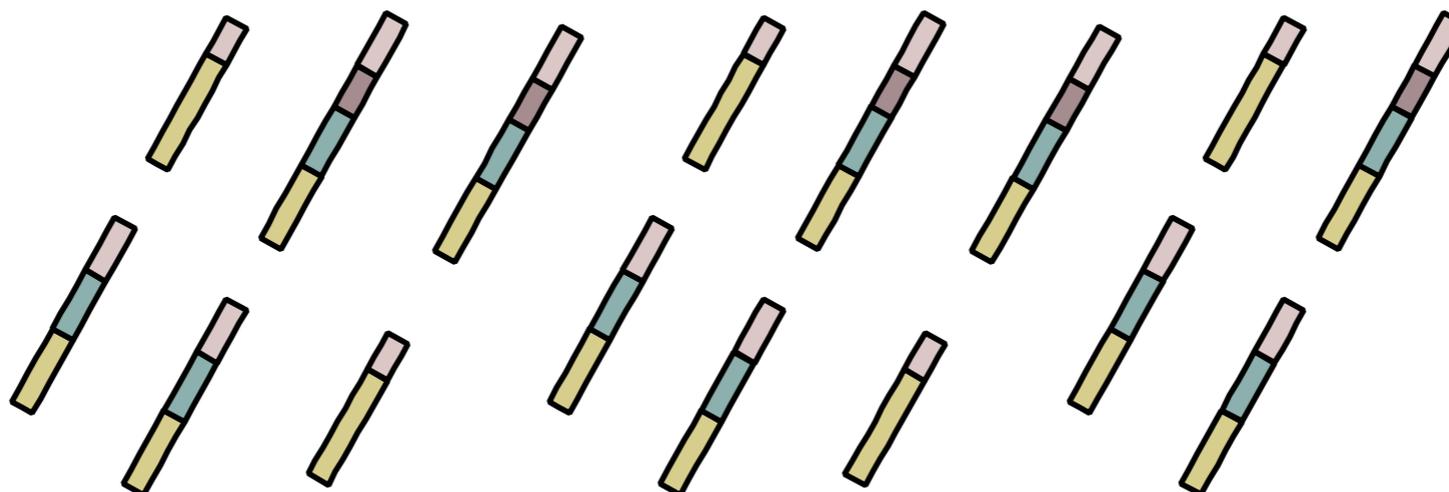
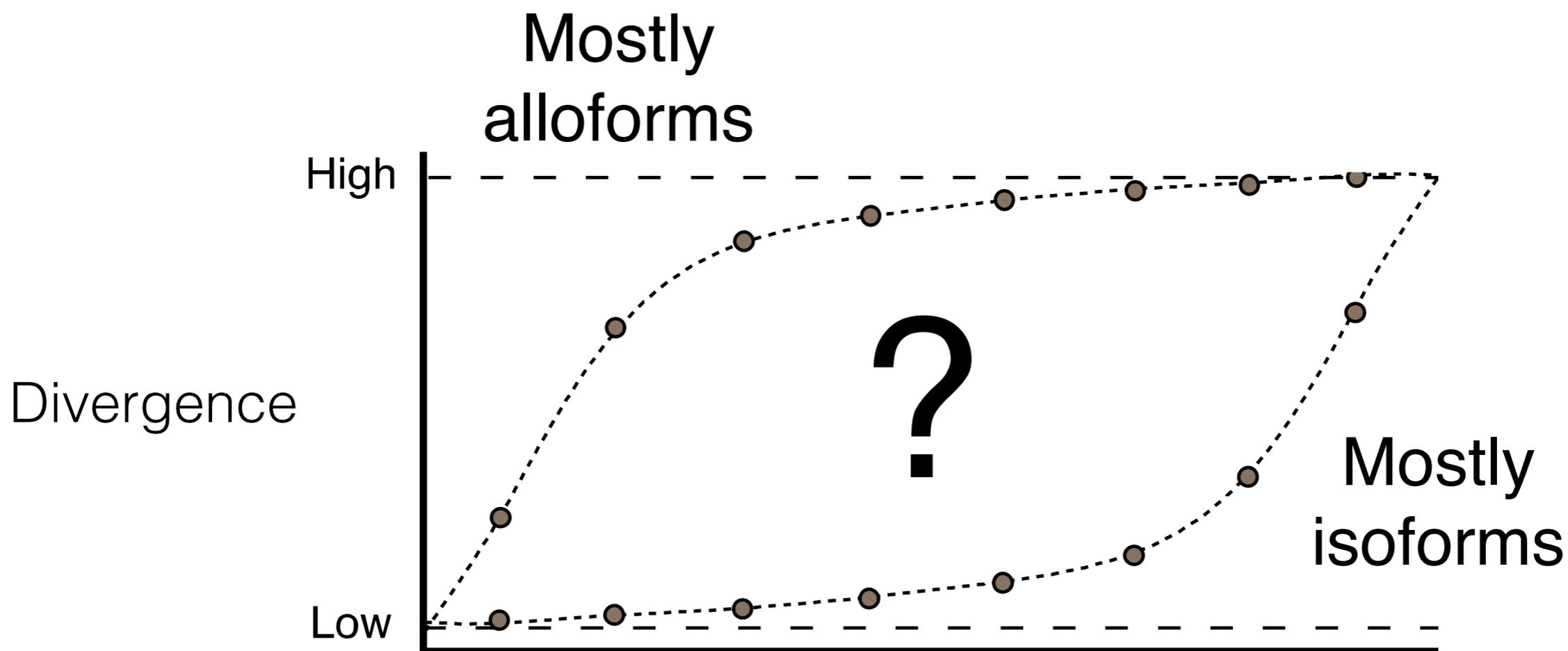
Stability

.....

## Unbiased functional profiling

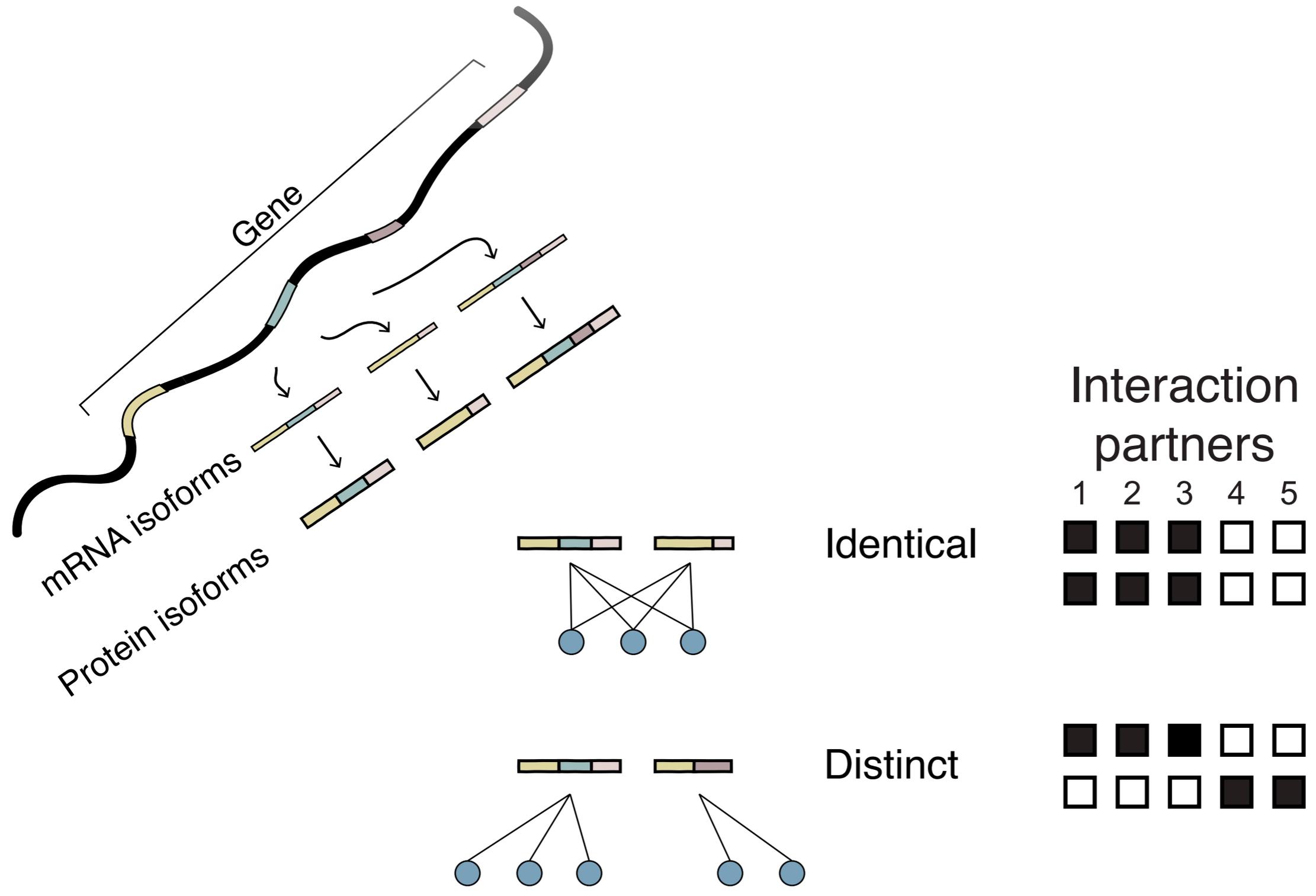
Profiles	Protein-protein interactions	Protein-DNA interactions	Enzymatic activities
Identical			
Intermediate			
Distinct			

# Landscape of protein isoform functional divergence



Large numbers of pairs of isoforms encoded by common genes

# Comparative protein-protein interaction profiling for large numbers of isoform pairs



# Comparative protein-protein interaction profiling for large numbers of isoform pairs

**Primary screen using all isoforms**



**PPI matrix profiling**



**Verification and  
sequence confirmation**

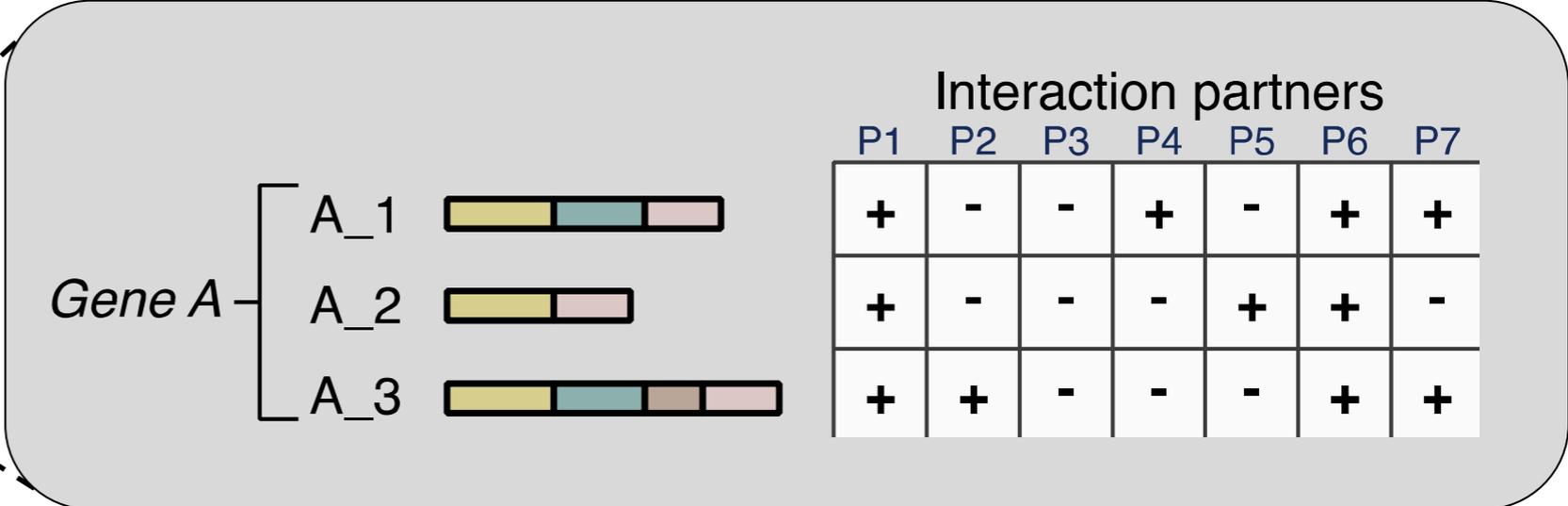
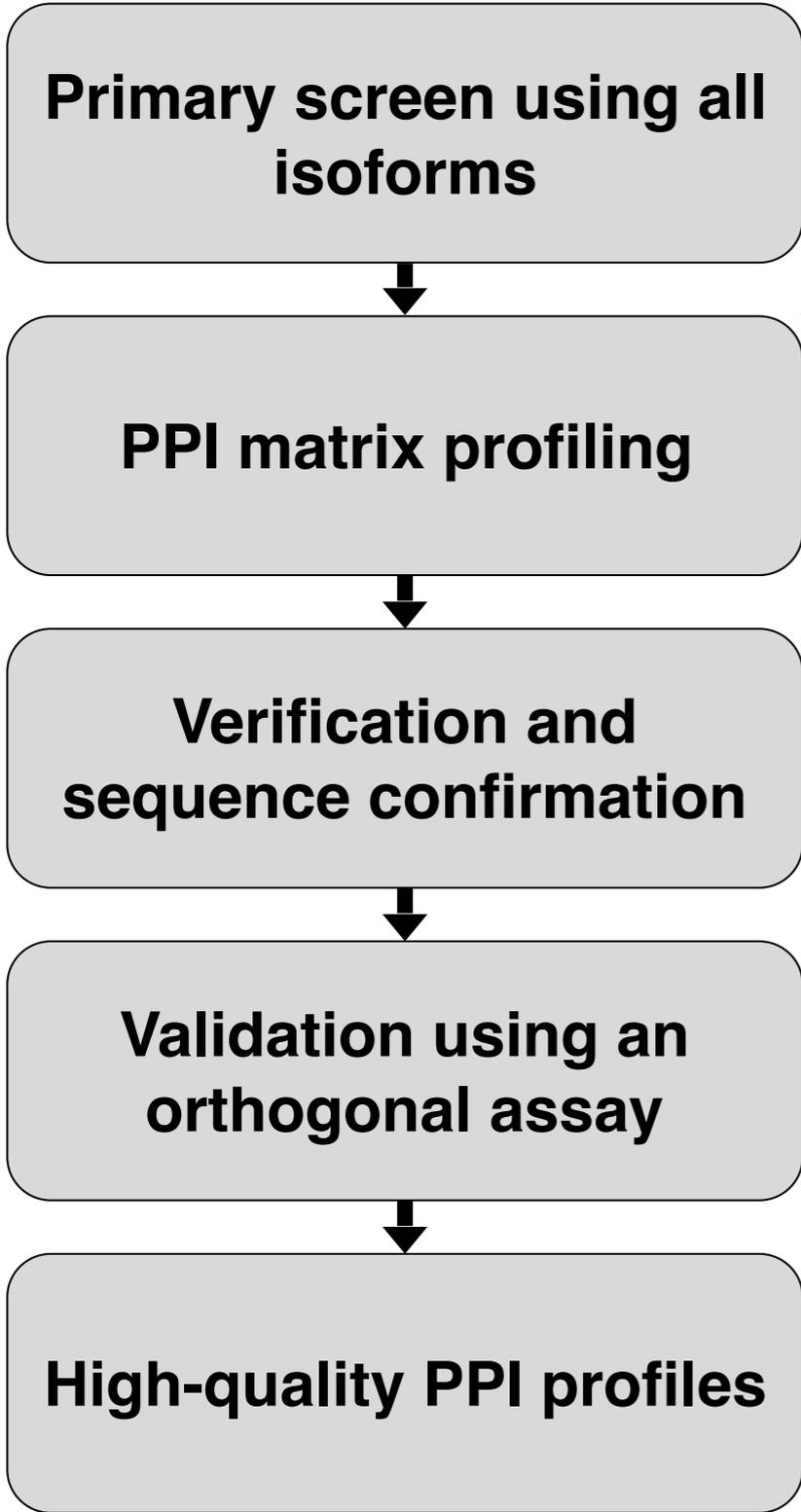


**Validation using an  
orthogonal assay**



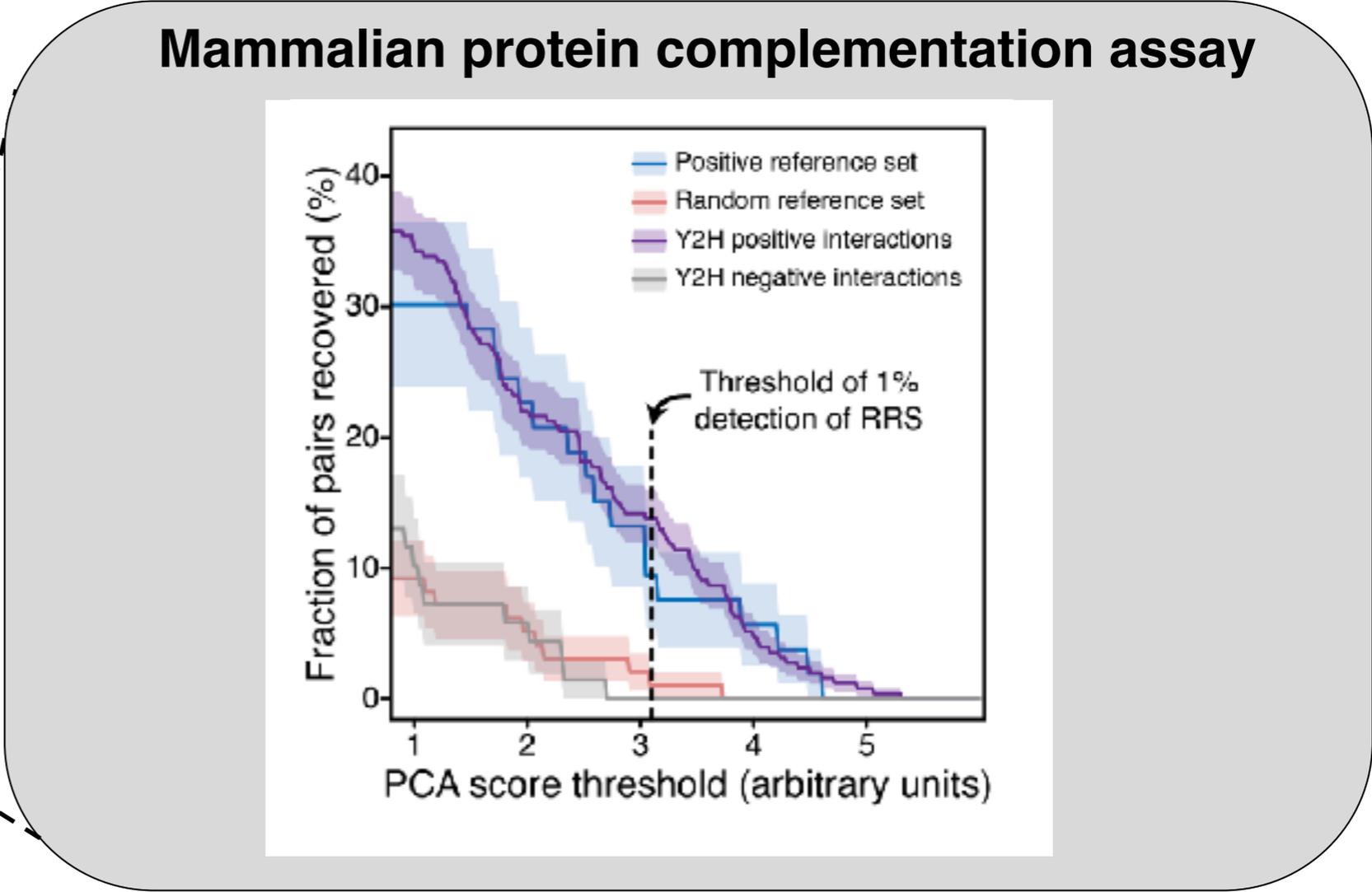
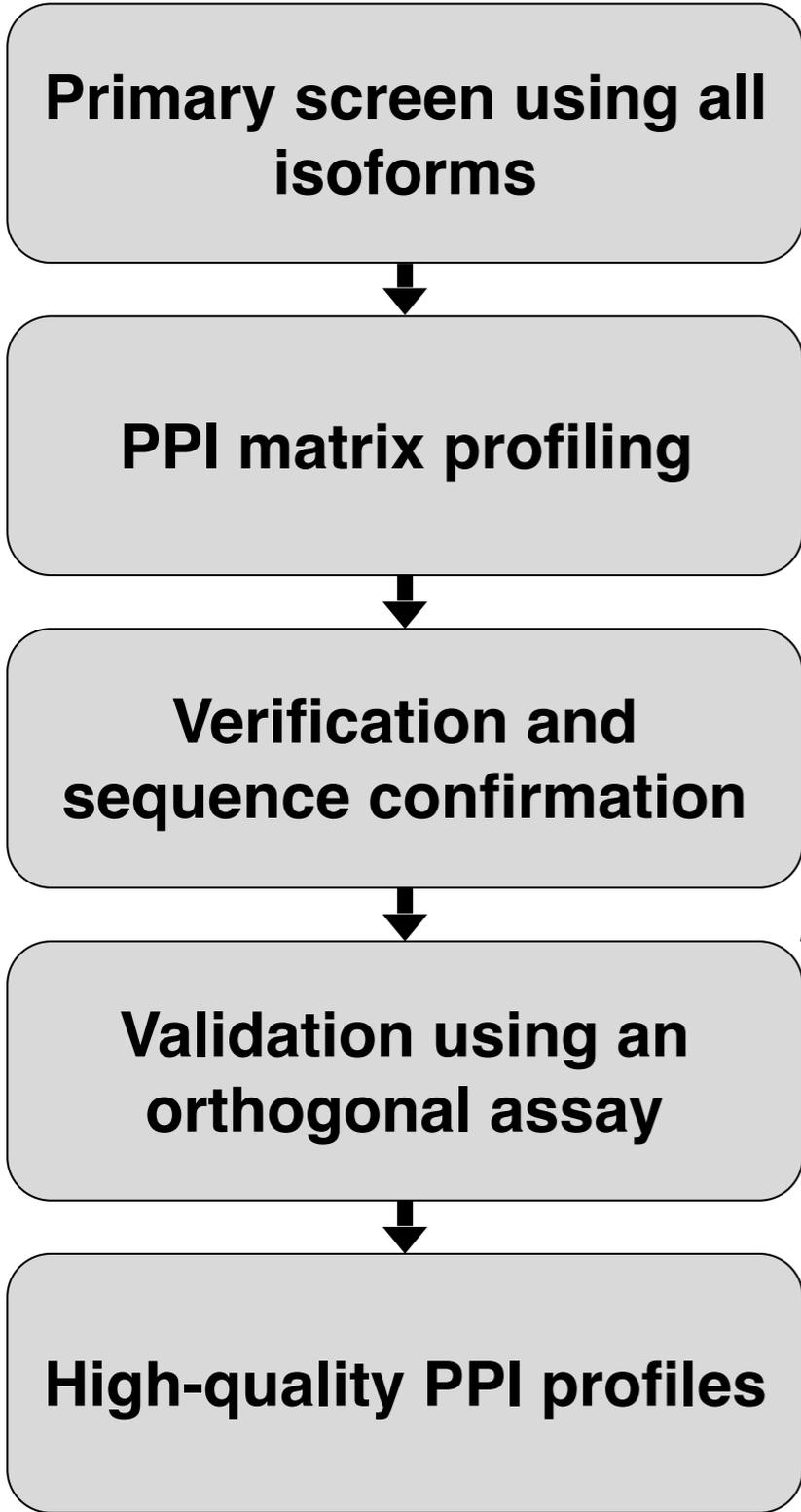
**High-quality PPI profiles**

# Comparative protein-protein interaction profiling for large numbers of isoform pairs



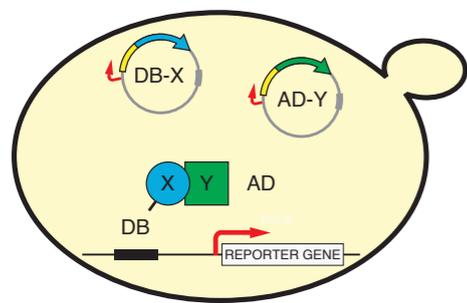
All first-pass pairs identified for any isoform were tested against all isoforms from a common gene.

# Comparative protein-protein interaction profiling for large numbers of isoform pairs

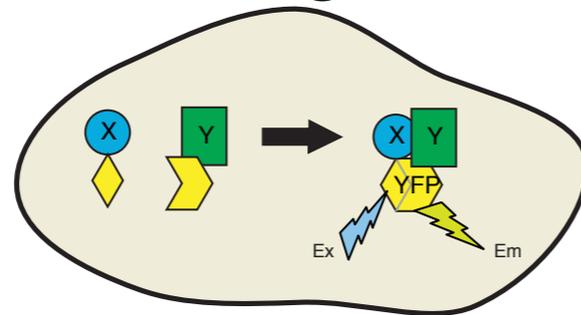


# Validation

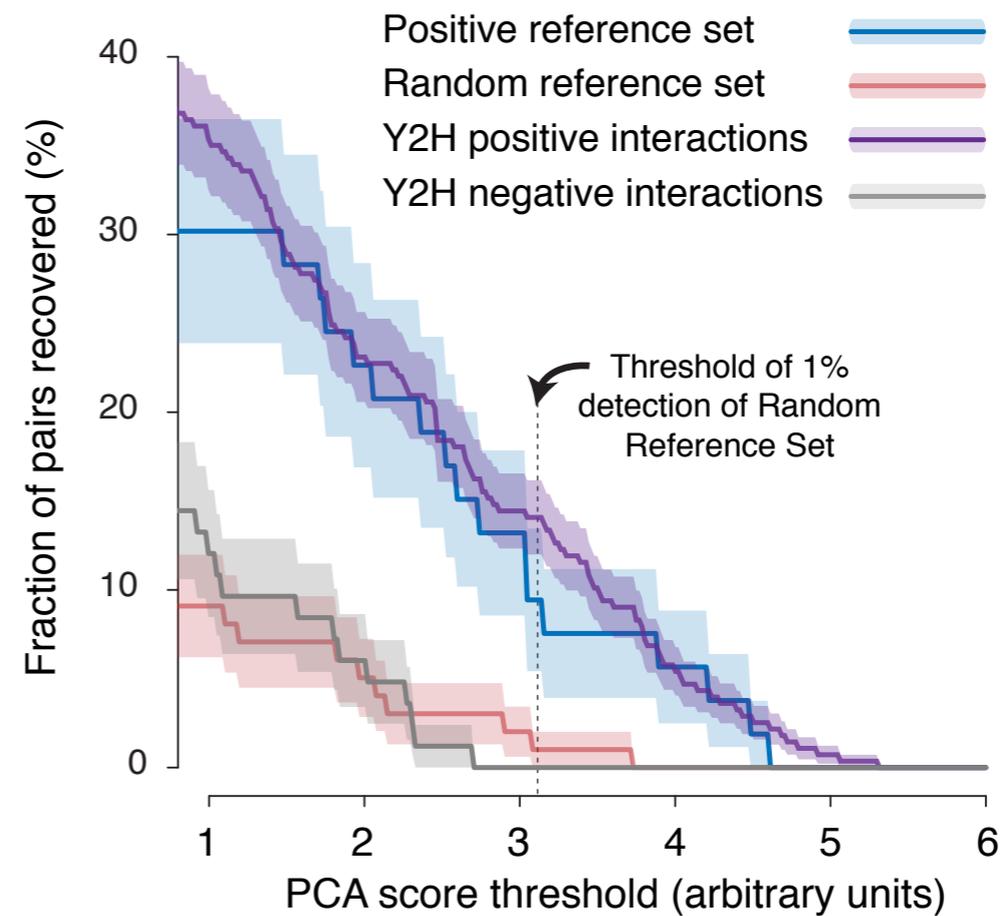
## Validation using an orthogonal assay (PCA)



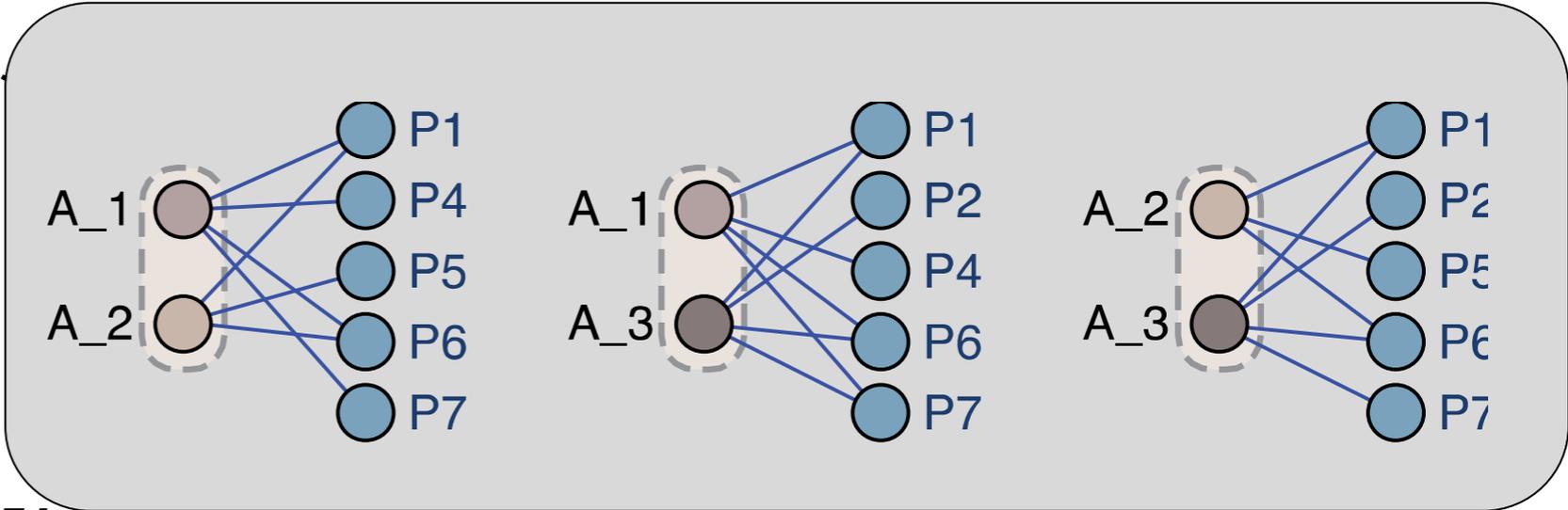
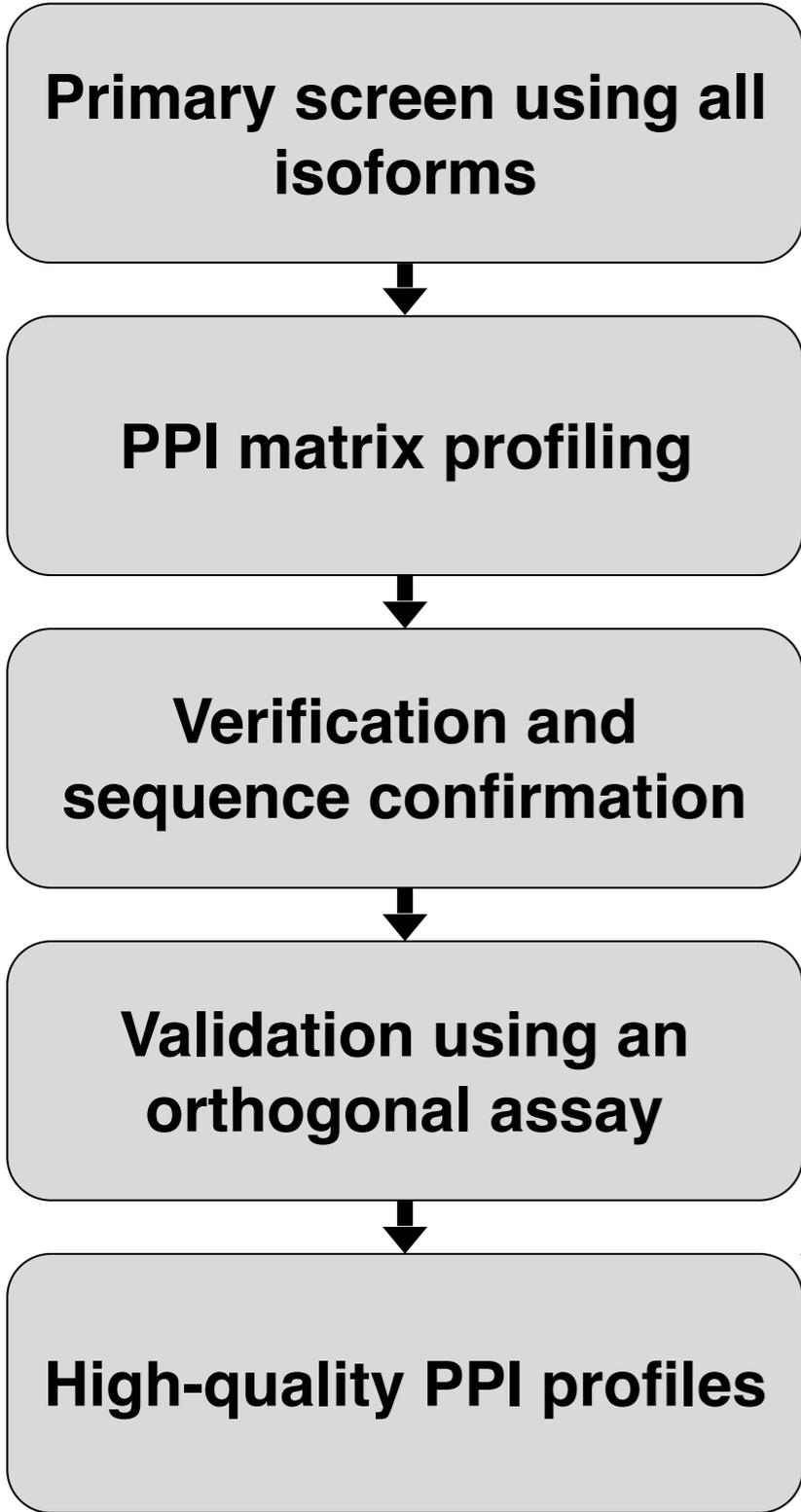
Two-hybrid  
In yeast cells



PCA  
In mammalian cells



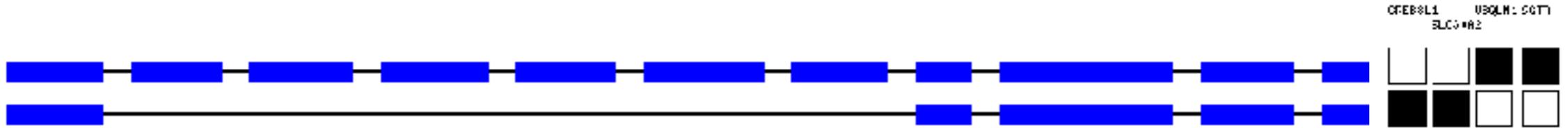
# Comparative protein-protein interaction profiling for large numbers of isoform pairs



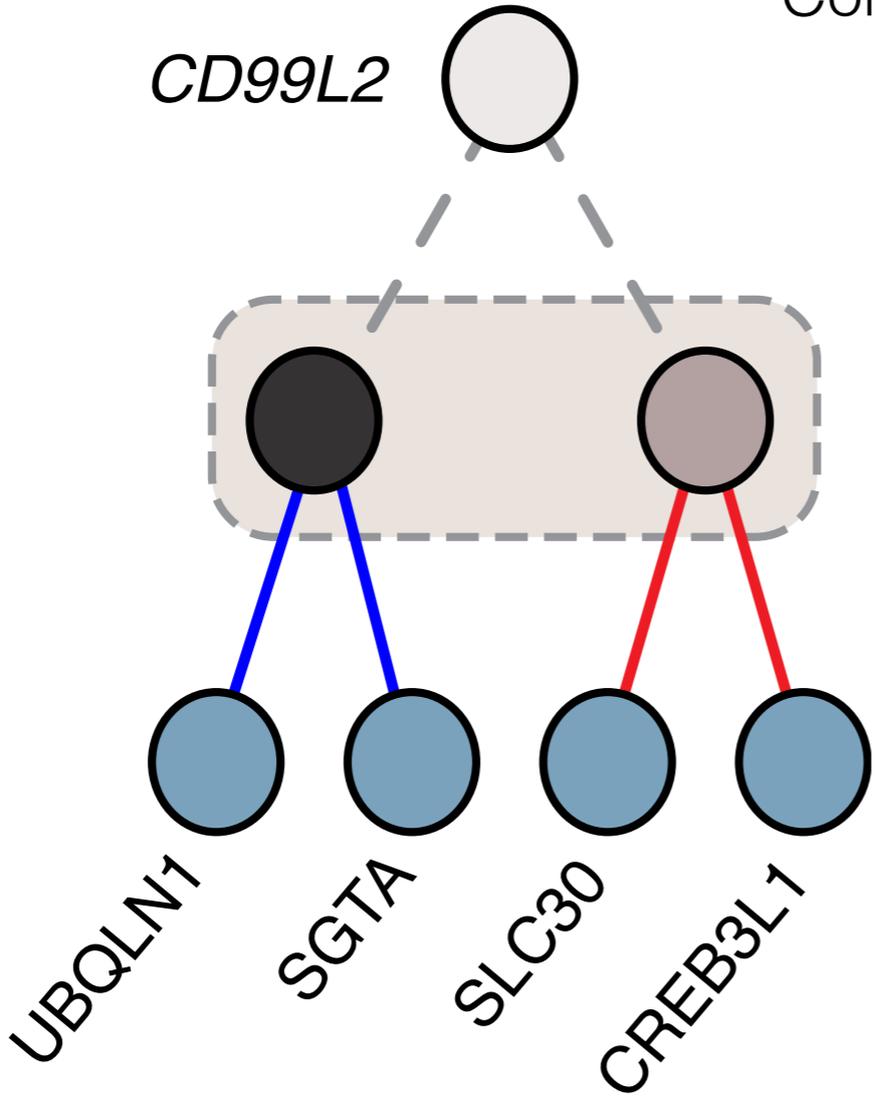


# Isoform-specific interaction network

CD99L2 isoforms

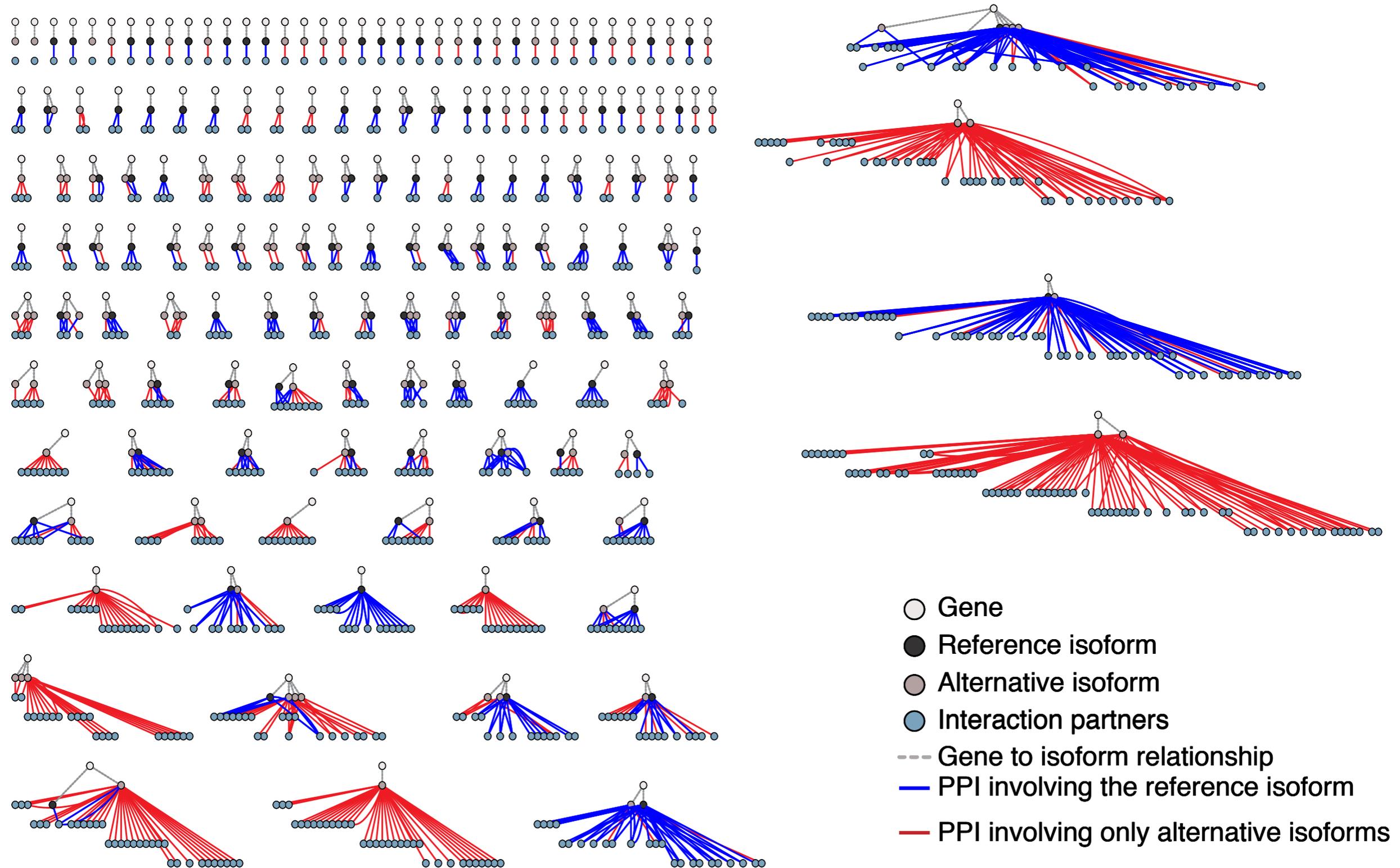


Convert into nodes and edges



- Gene
- Reference isoform
- Alternative isoform
- Interaction partners
- Gene to isoform relationship
- PPI involving the reference isoform
- PPI involving only alternative isoforms

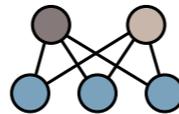
# Isoform-specific interaction network



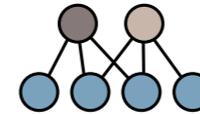
# Interaction profile dissimilarity

$$d_J(A, B) = 1 - J(A, B) = \frac{|A \cup B| - |A \cap B|}{|A \cup B|}$$

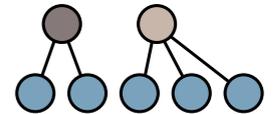
Identical



Intermediate



Distinct



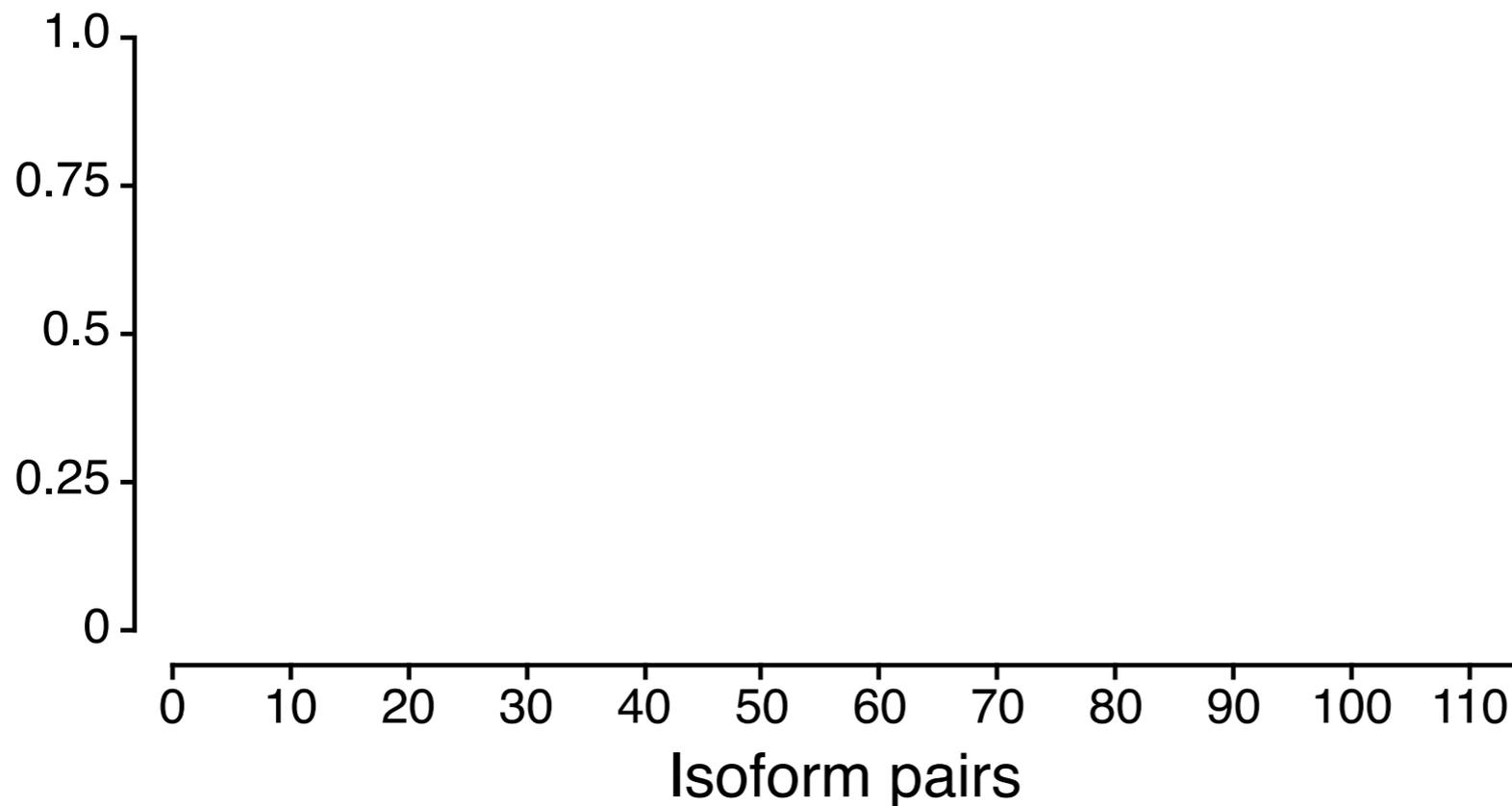
Jaccard distance:

0

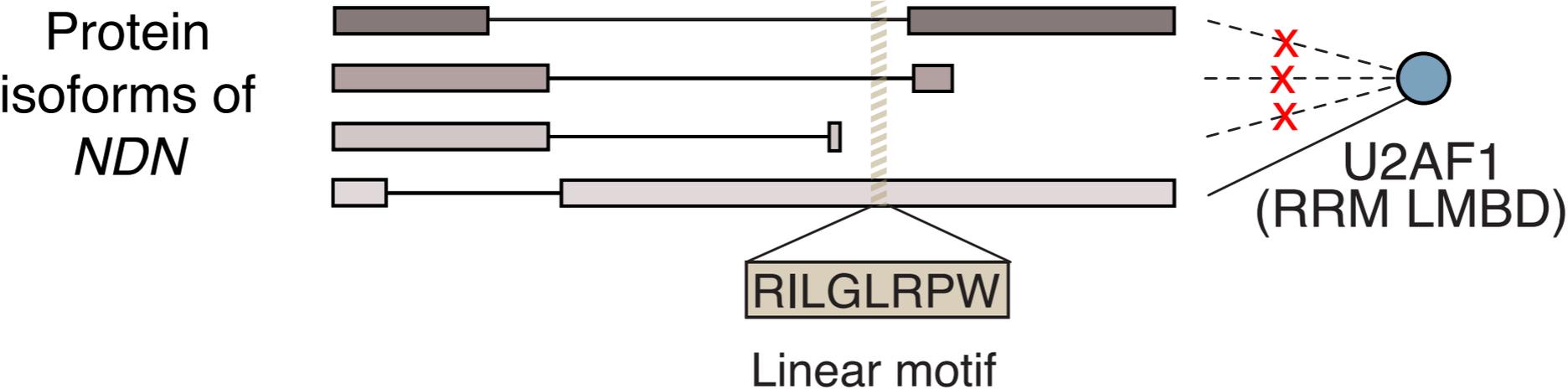
0.5

1.0

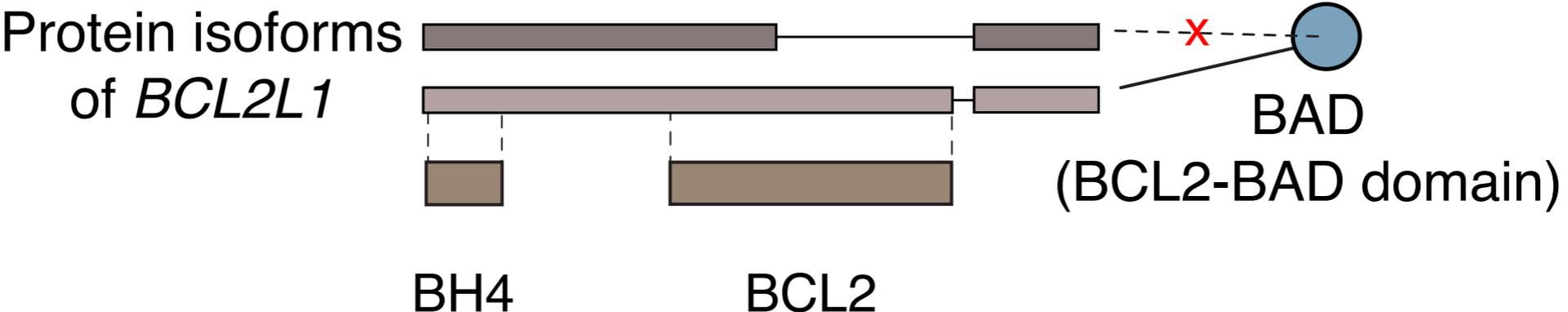
Interaction profile dissimilarity  
(Jaccard distance)



# Sequence features in isoforms underlie interaction perturbations



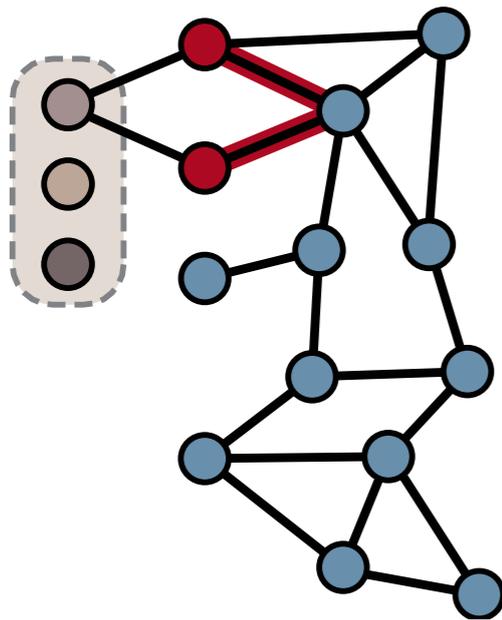
Isoforms with a high density of linear motifs tend to promote interactions (p-value =  $5.7 \times 10^{-4}$ )



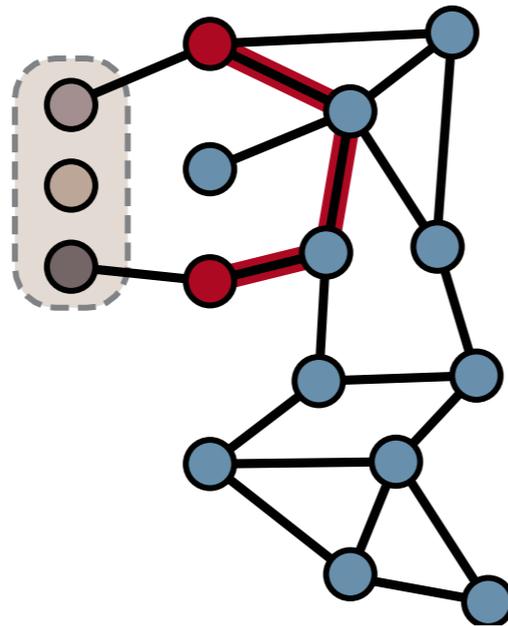
Truncated domains tend to correspond to a loss of interaction (p-value =  $6.4 \times 10^{-5}$ )

# How functionally divergent are isoform-specific interactors?

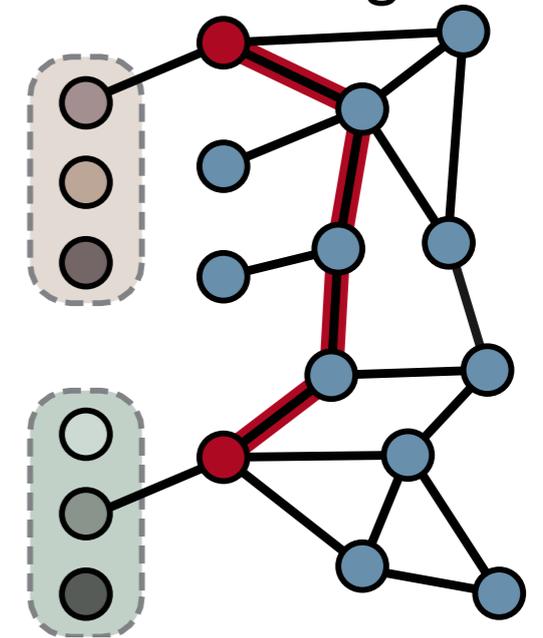
Single protein



Alternative isoforms

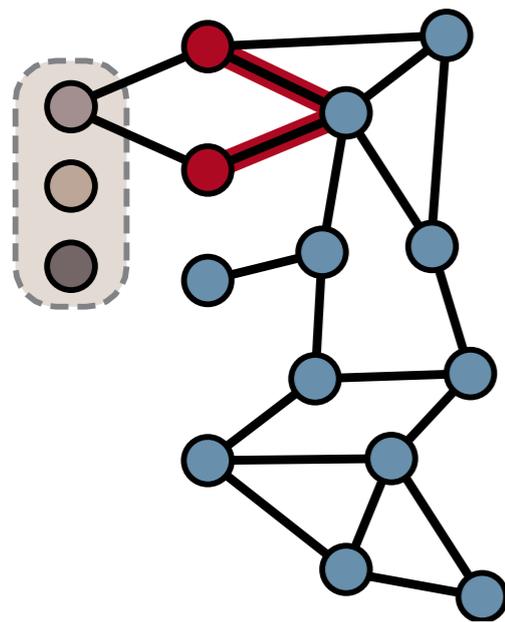


Proteins from different genes

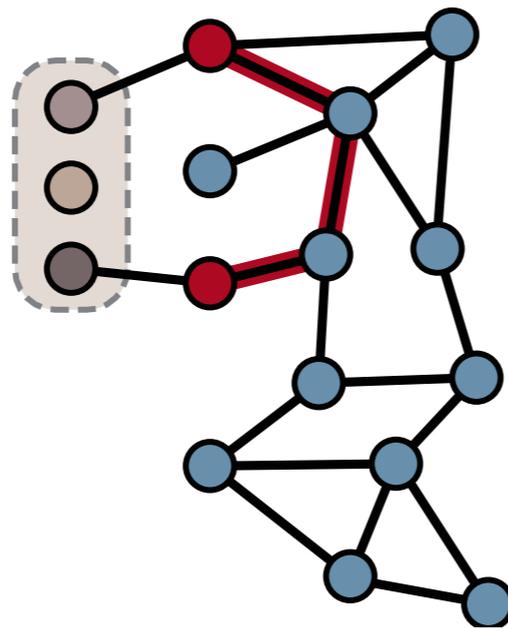


How functionally divergent  
are isoform-specific interactors?

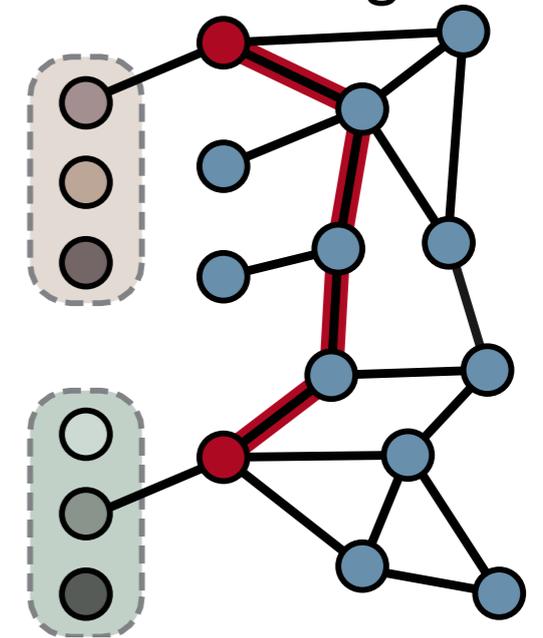
Single protein



Alternative isoforms



Proteins from  
different genes

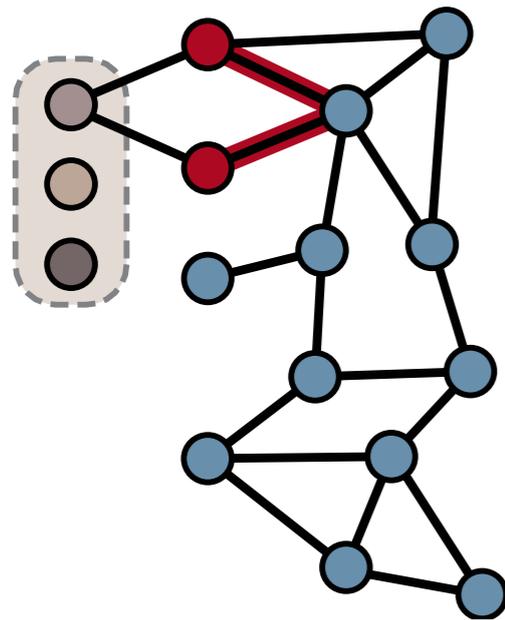


# A Proteome-Scale Map of the Human Interactome Network

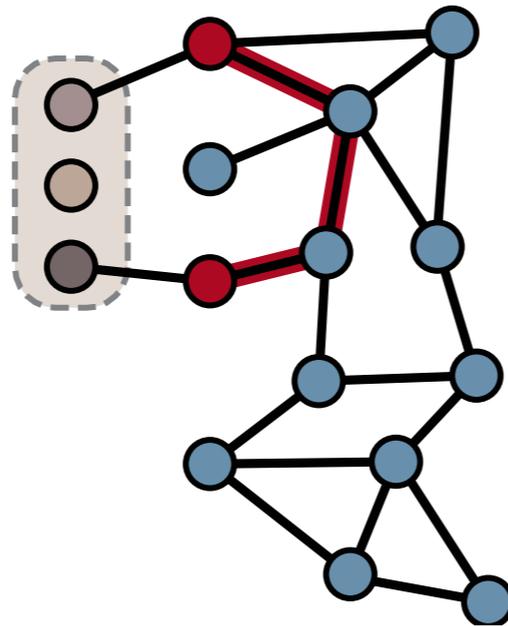


How functionally divergent  
are isoform-specific interactors?

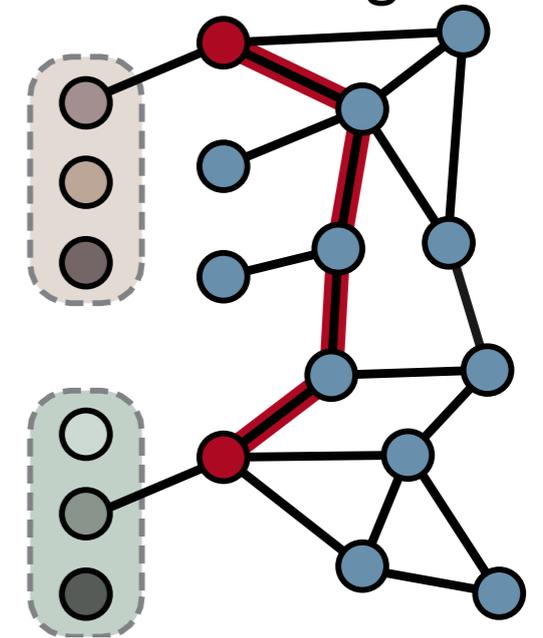
Single protein



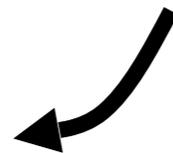
Alternative isoforms



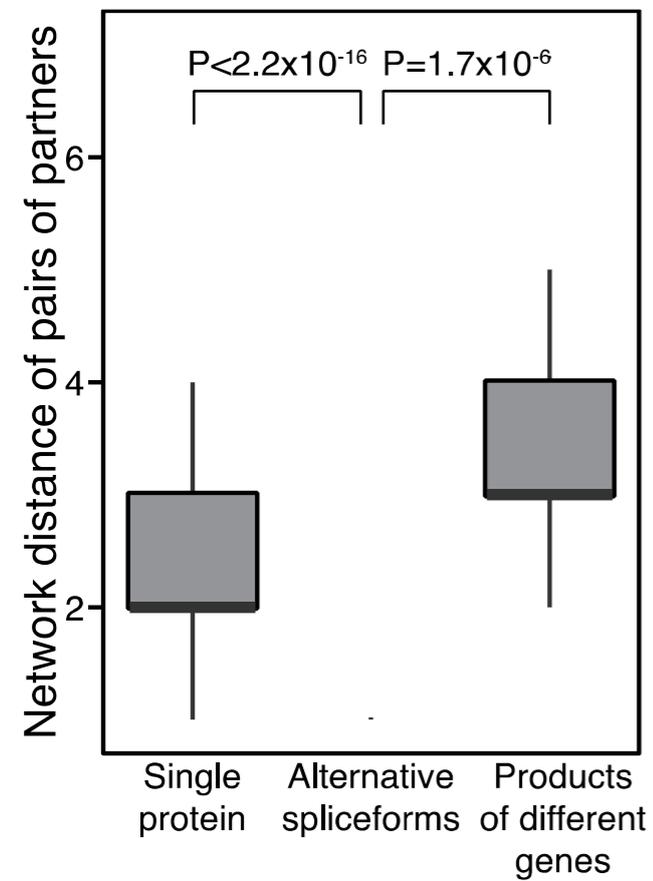
Proteins from  
different genes



# A Proteome-Scale Map of the Human Interactome Network

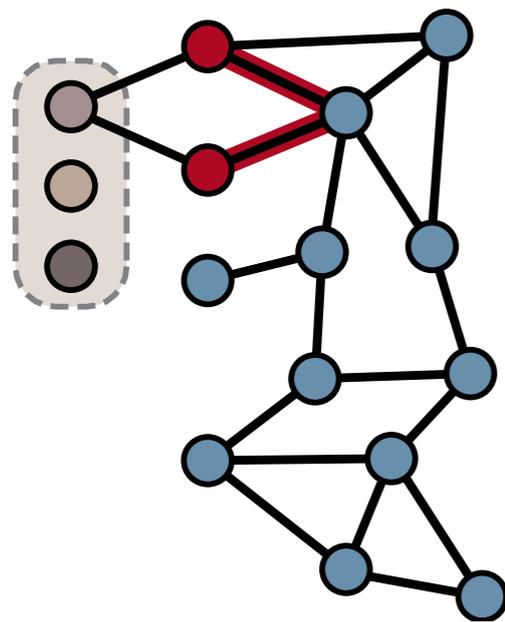


# How functionally divergent are isoform-specific interactors?

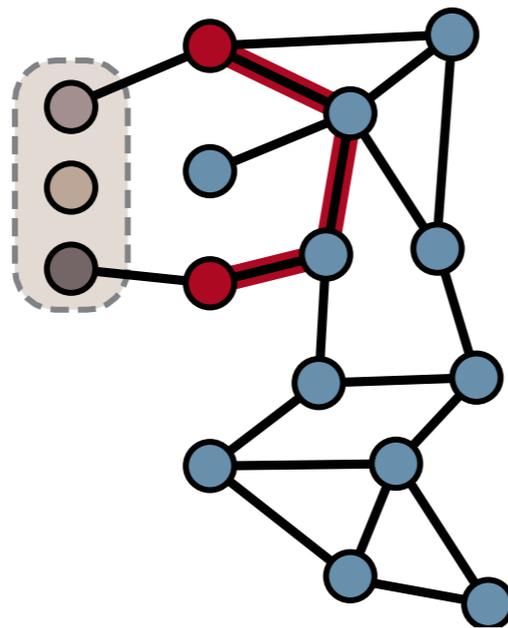


# How functionally divergent are isoform-specific interactors?

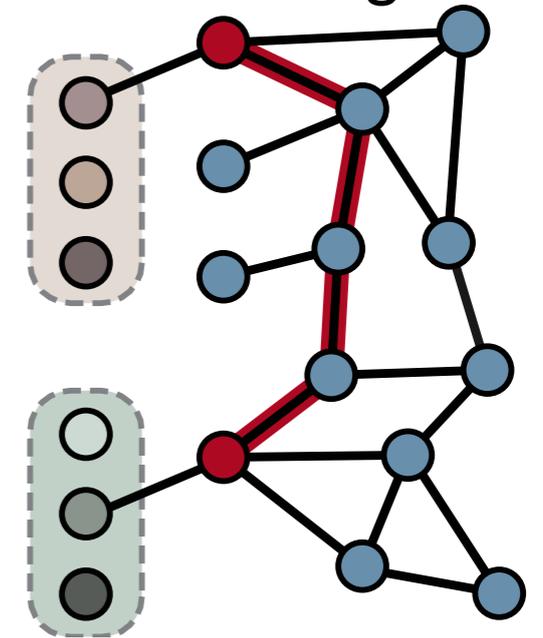
Single protein



Alternative isoforms

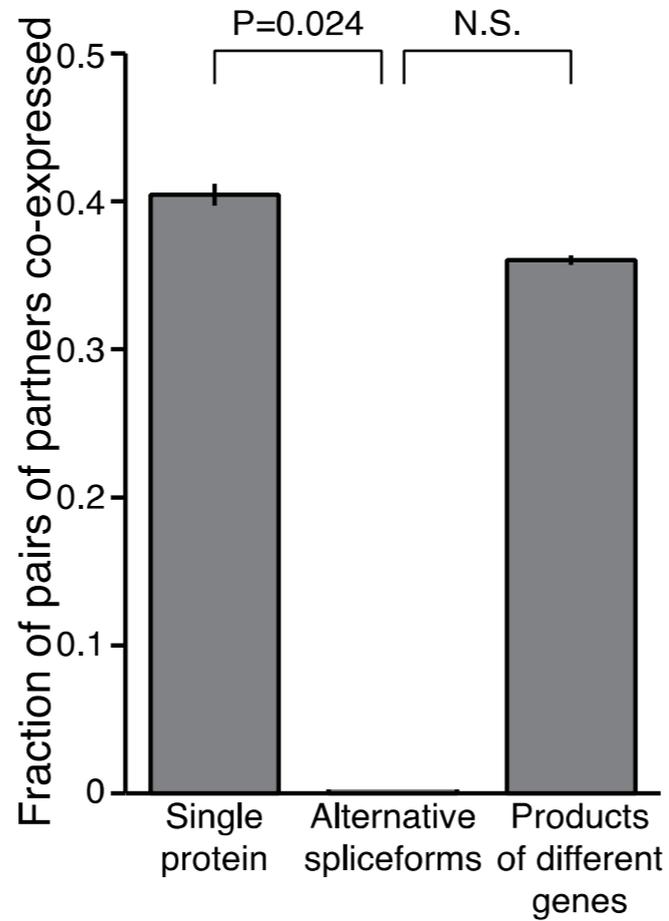
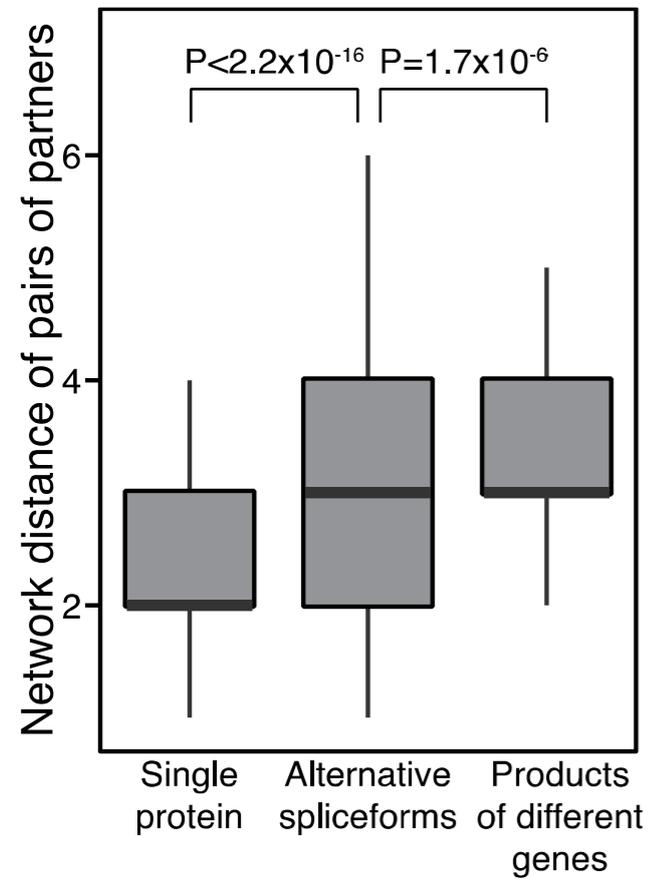


Proteins from  
different genes



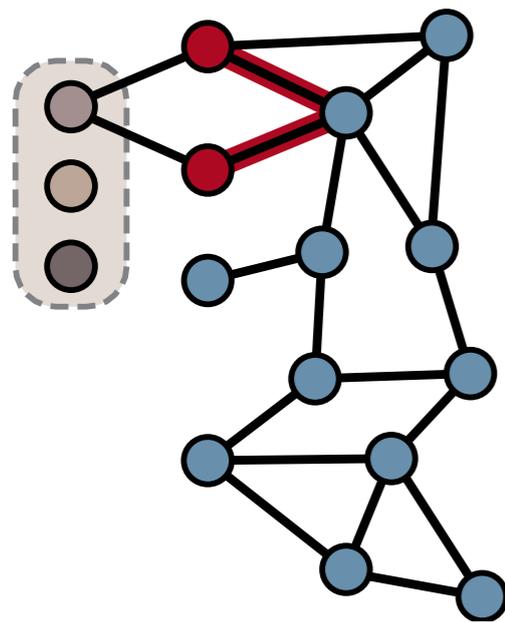
Illumina BodyMap 2.0  
Gene expression for 16 human tissues

# How functionally divergent are isoform-specific interactors?

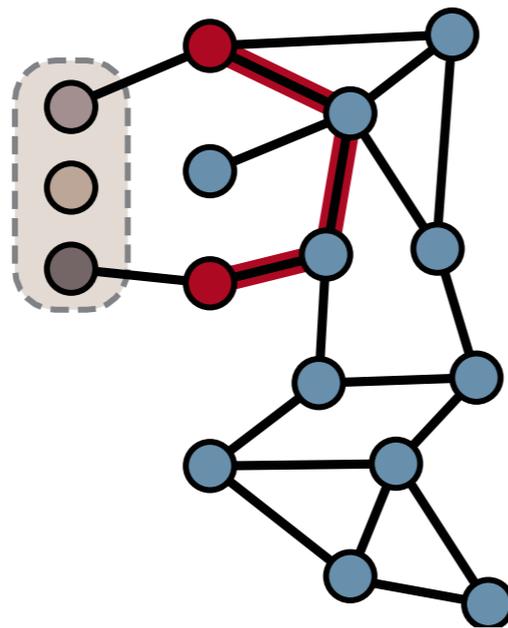


# How functionally divergent are isoform-specific interactors?

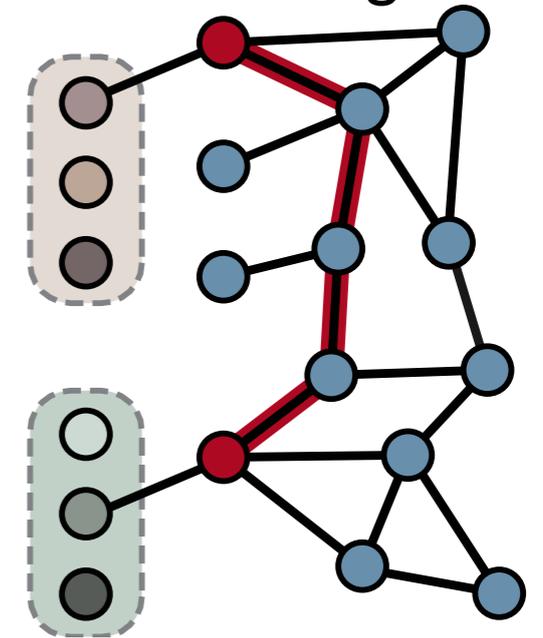
Single protein



Alternative isoforms

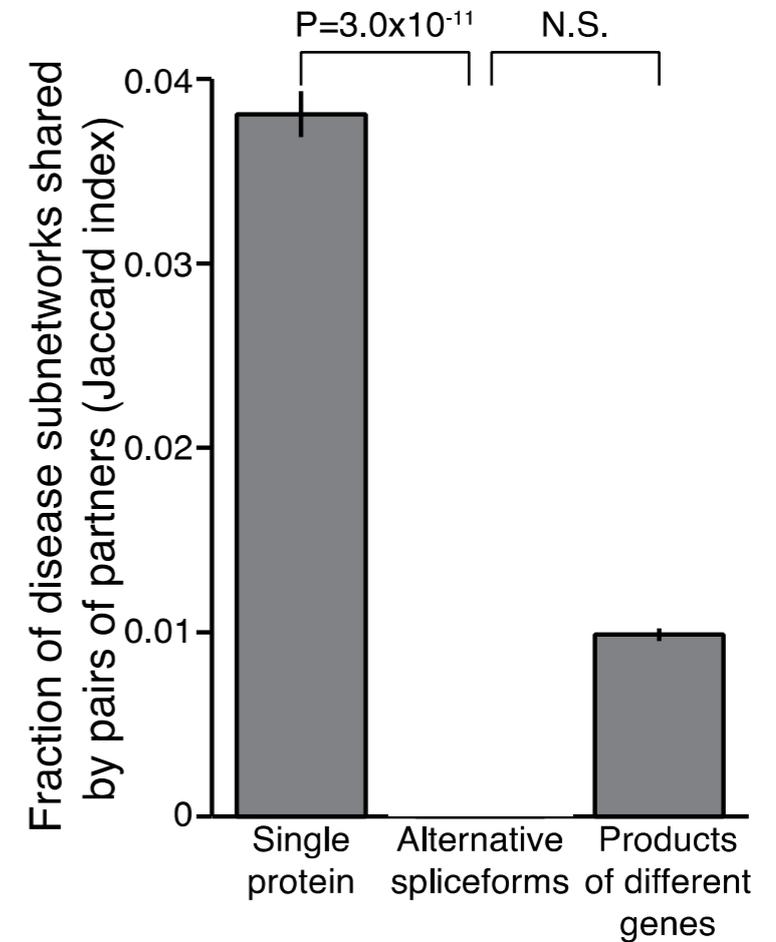
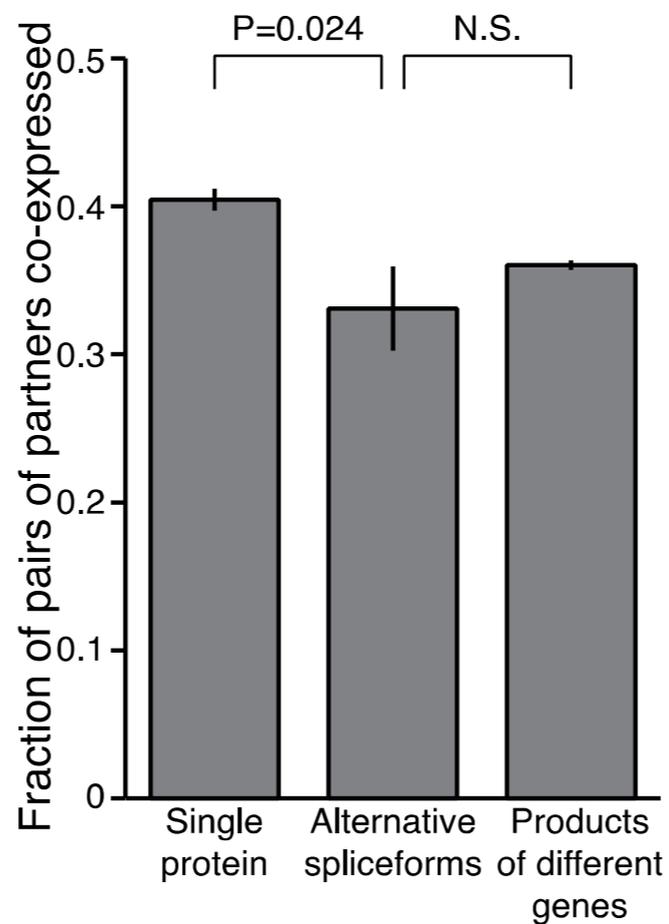
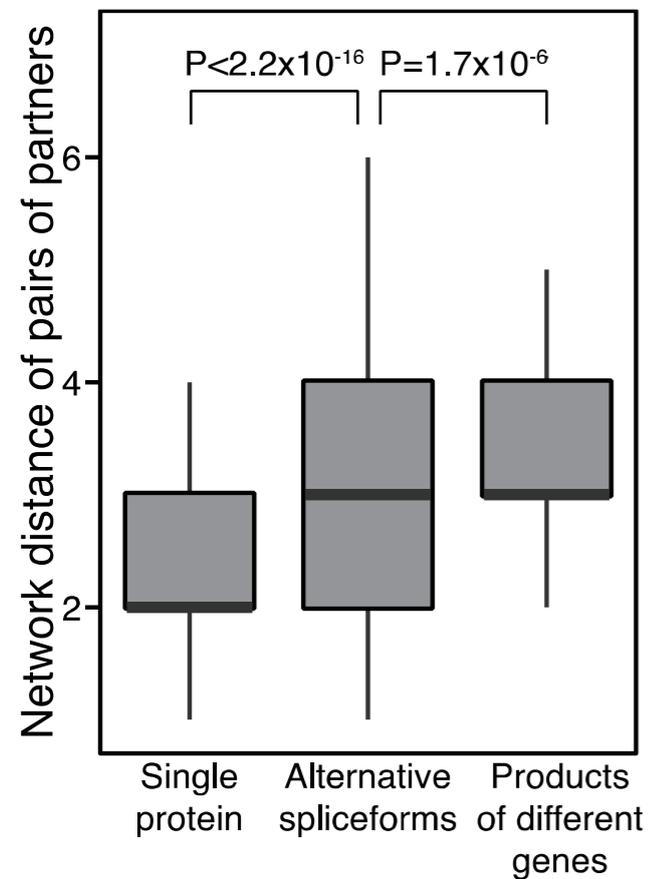


Proteins from  
different genes

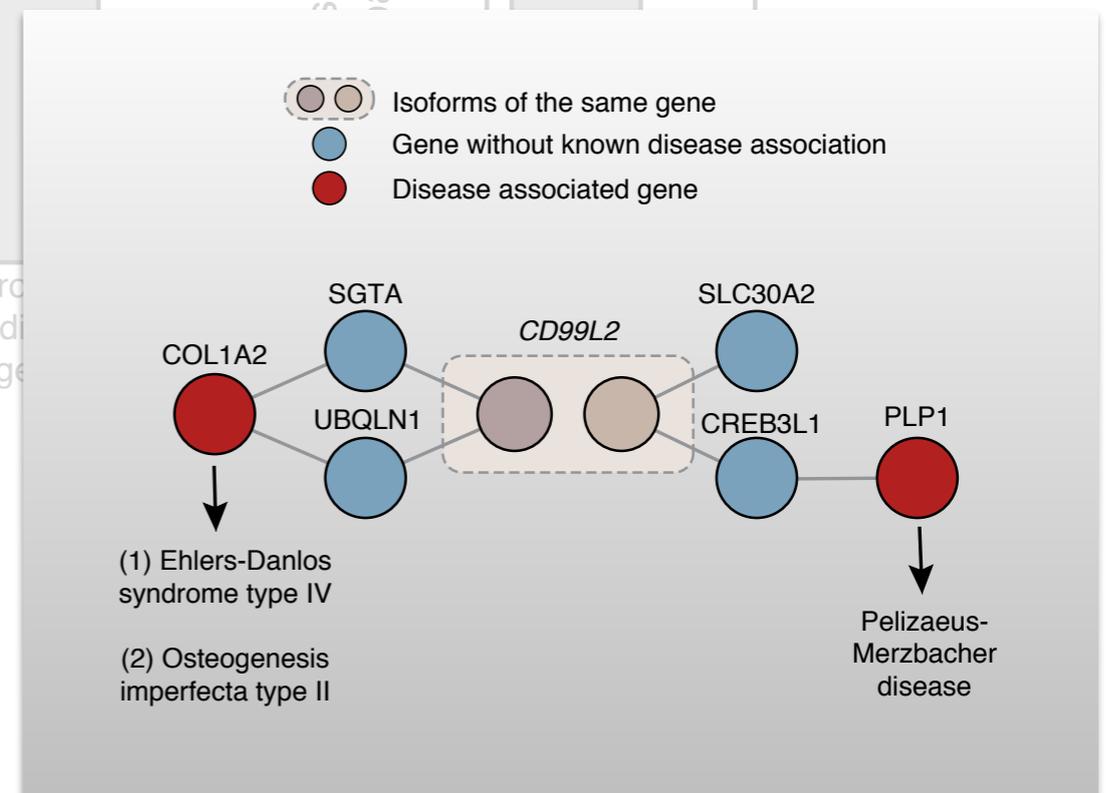
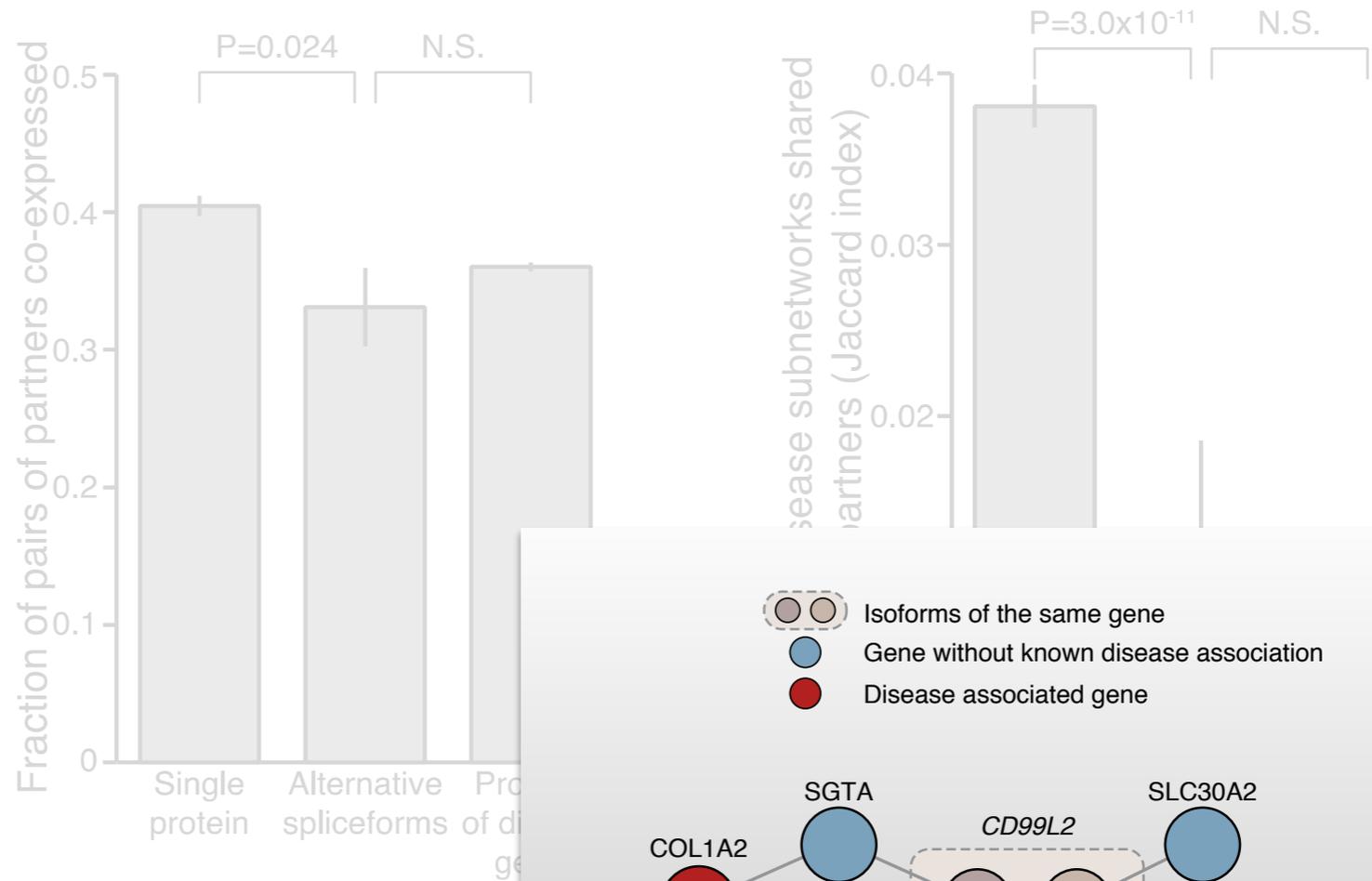
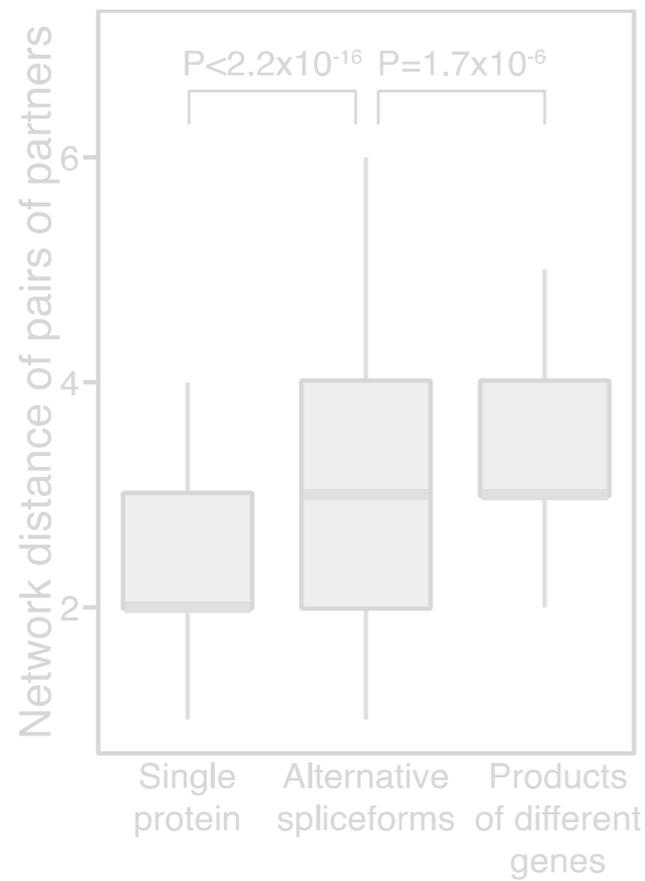


Disease subnetworks derived from  
GeneCards database

# How functionally divergent are isoform-specific interactors?



# How functionally divergent are isoform-specific interactors?



# How widespread is isoform functional divergence in the whole proteome?

Systematic identification  
of large numbers of  
isoforms pairs  
for large numbers  
of human genes



Physical interactions

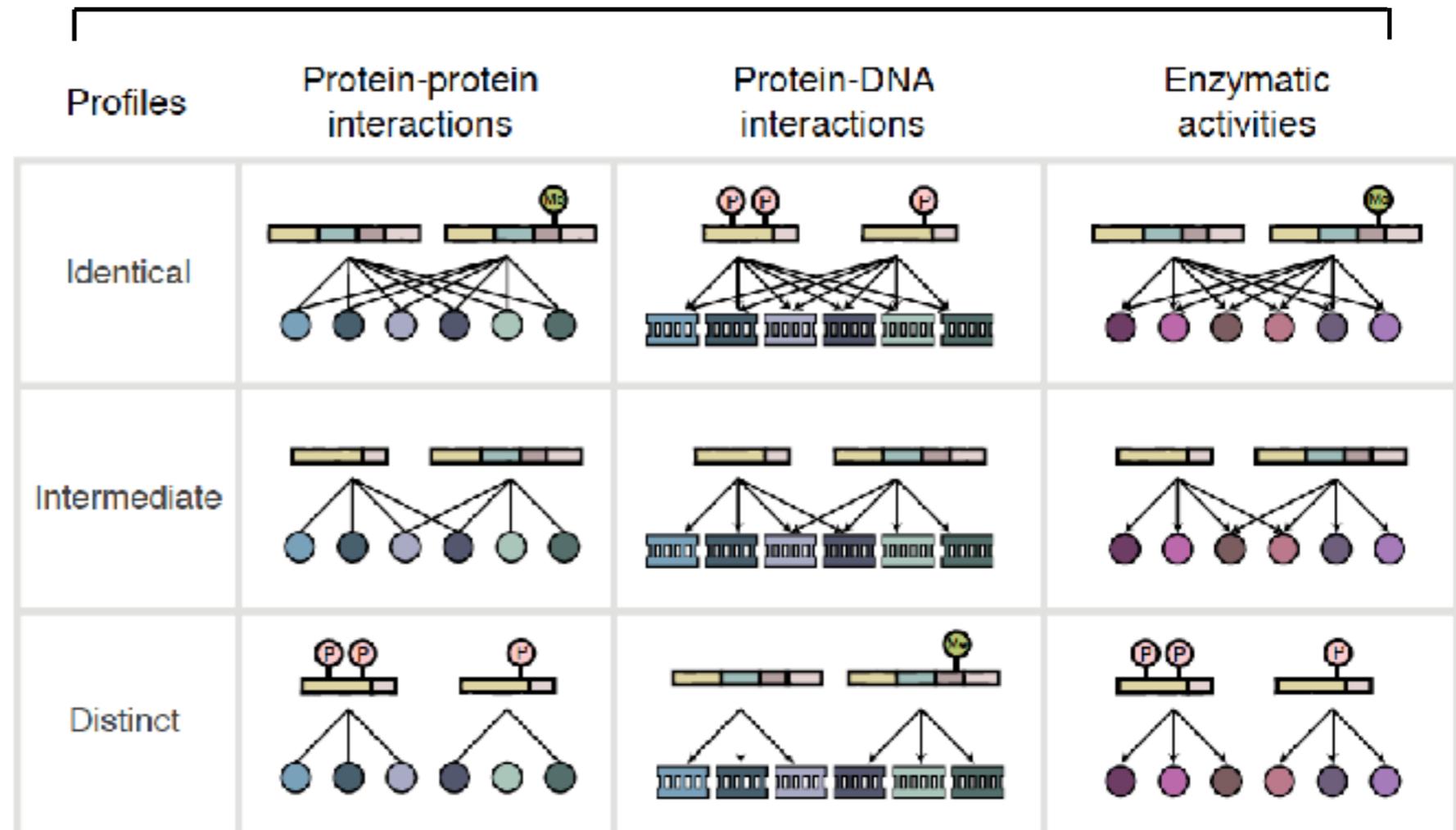
Enzymatic activities

Cellular localization

Stability

.....

## Unbiased functional profiling

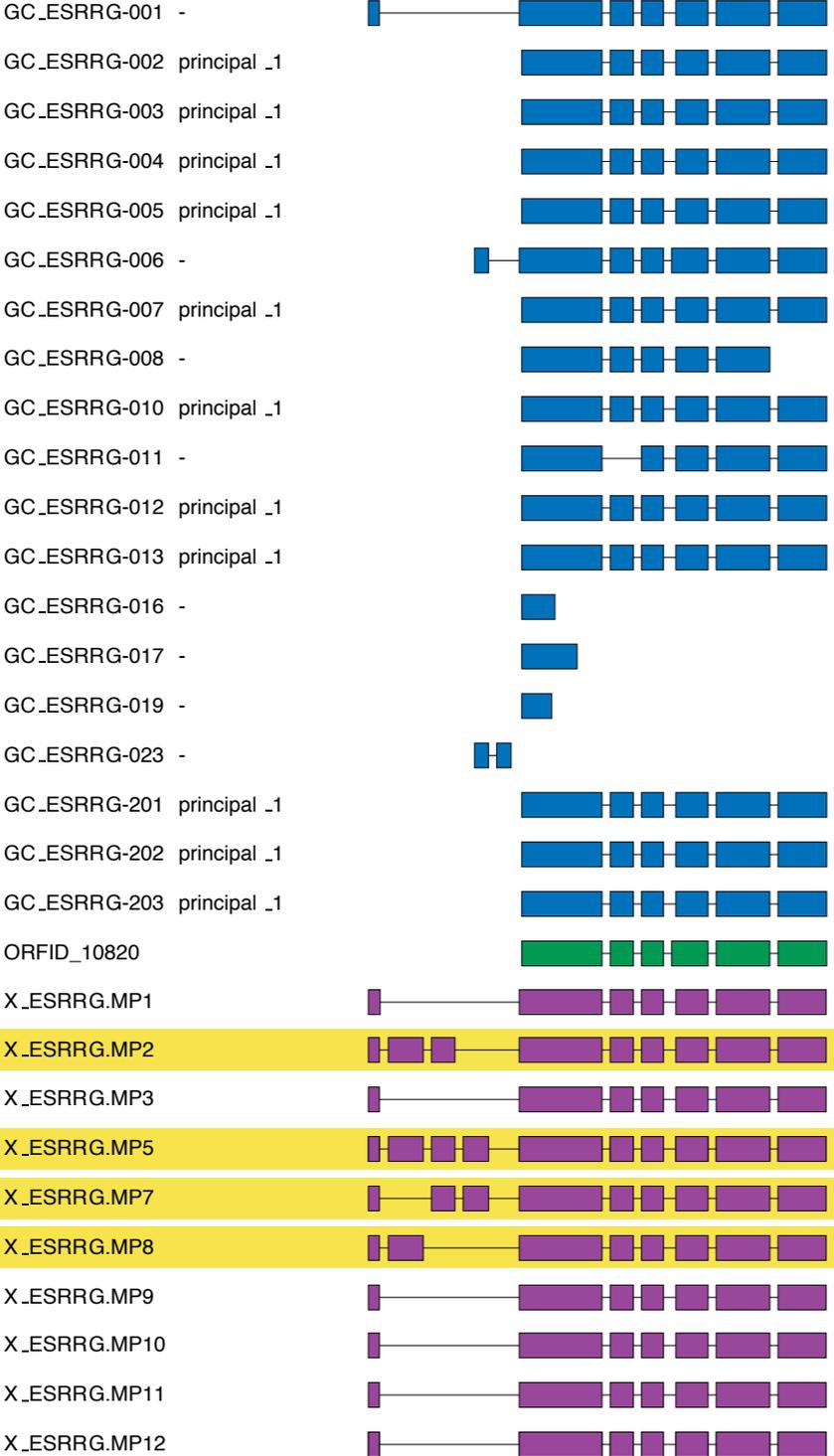


# Integration of ORF-Seq with third generation sequencing technologies

## Pipeline in progress for 800 human TFs

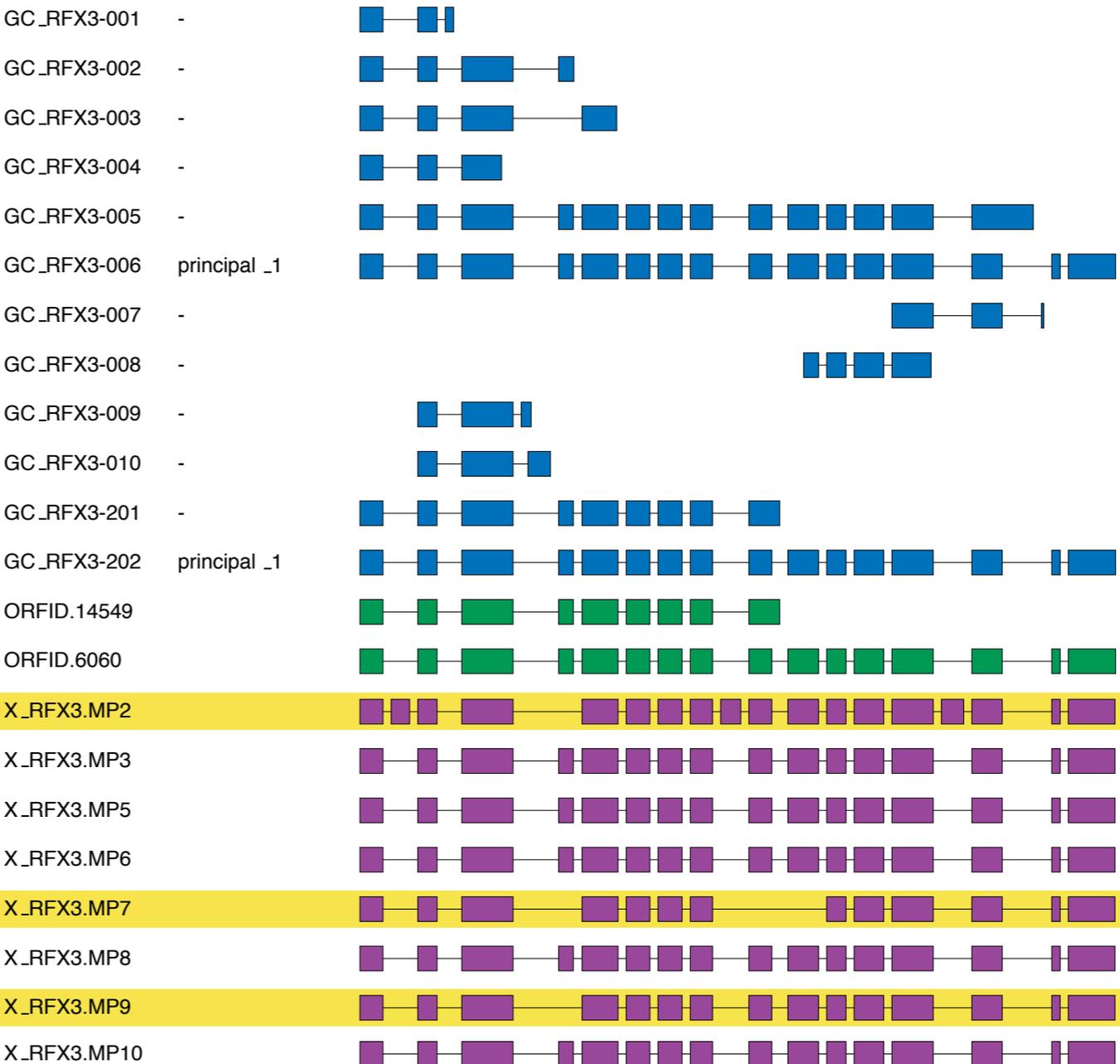
Gene:ESRRG, Strand: - chr1:216505473- 216723335

Isoform Name Appris Status



Gene:RFX3, Strand: - chr9:3225042- 3395588

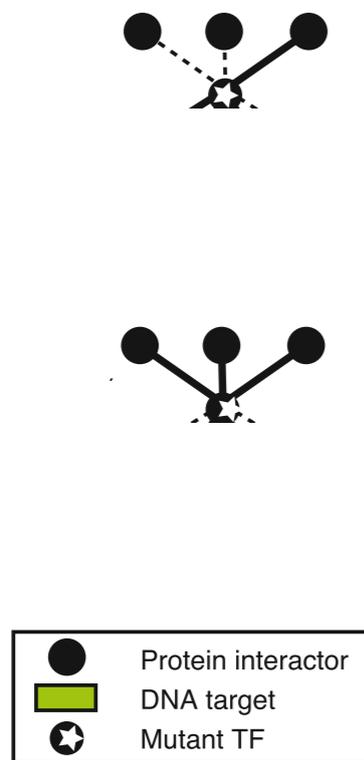
Isoform Name Appris Status



Isoform clones that are not present in Gencode v25, and thus represent novel forms.

# Multi-parameter comparative profiling of isoforms

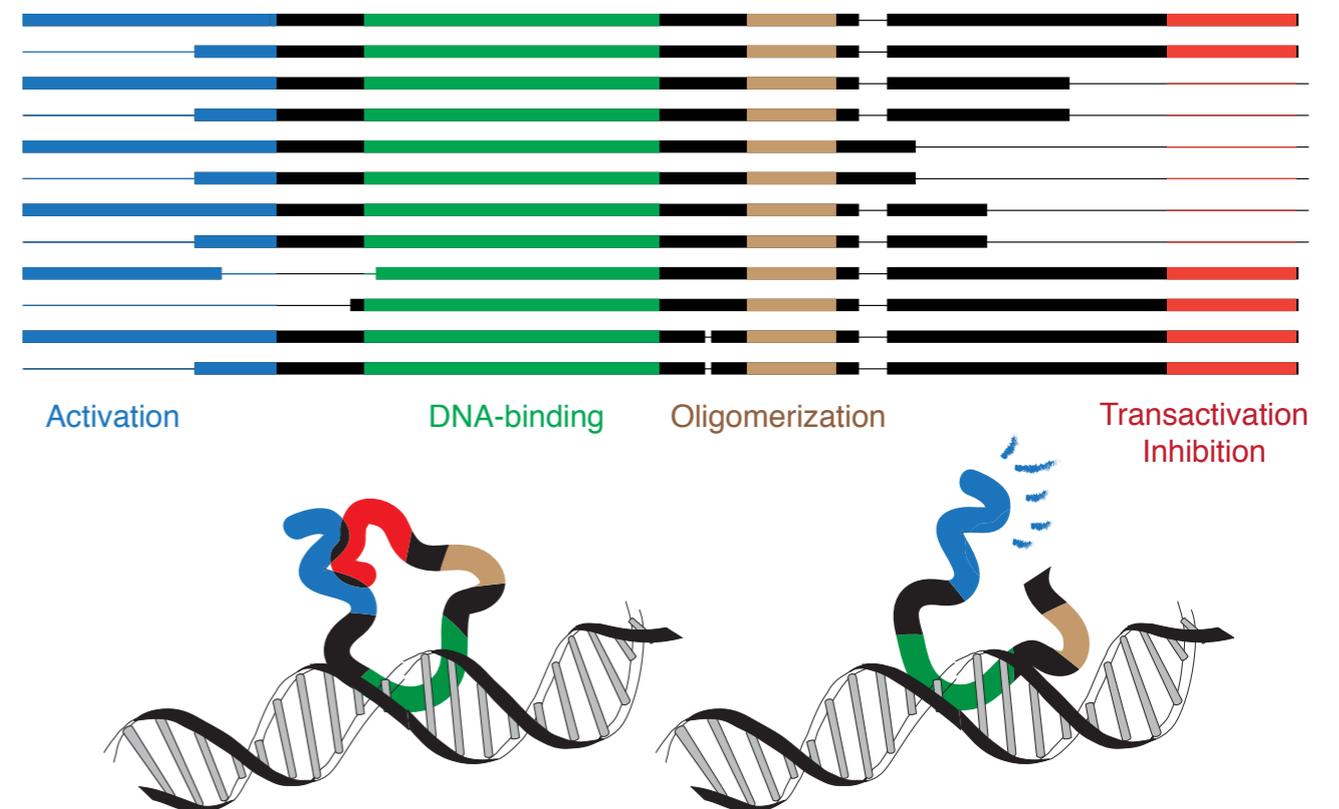
Variations in TF sequence perturb protein and/or DNA interactions



Sahni *et al* Cell 2016

Similar perturbations likely for isoforms

TP63 isoforms

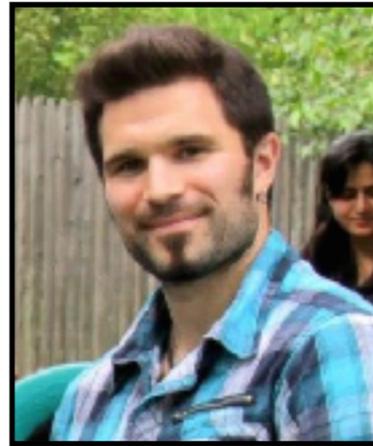


Will conduct screens for differential protein-protein and protein-DNA interactions for TF isoforms

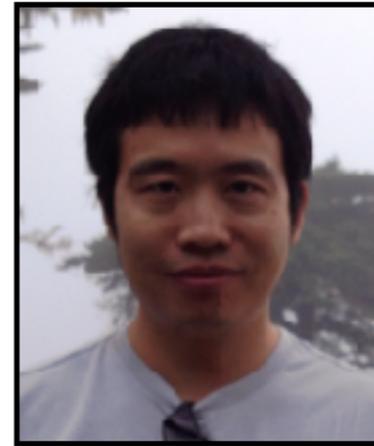
# Acknowledgements



Xinping Yang



Jasmin  
Coulombe-Huntington



Shuli Kang



Tong Hao

**DFCI, Boston**  
Marc Vidal  
Dawit Balcha  
Wenting Bian  
Tiziana Cafarelli

Michael Calderwood  
Soon Gang Choi  
Meaghan Daley  
David DeRidder  
Alice Desbuleux

Tong Hao  
David Hill  
Katja Luck  
Dylan Markey  
Julien Olivet

Carl Pollis  
Aaron Richardson  
Sadie Schlabach  
Kerstin Spirohn  
Yang Wang

**Donnelly Center, Toronto**  
Fritz Roth



**UCSD, San Diego**  
Lilia Iakoucheva  
Shuli Kang

**McGill, Montreal**  
Yu Xia  
Jasmin Coulombe-Huntington