# SQANTI and TAPP**AS**:
# Making Sense of Iso-Seq Data
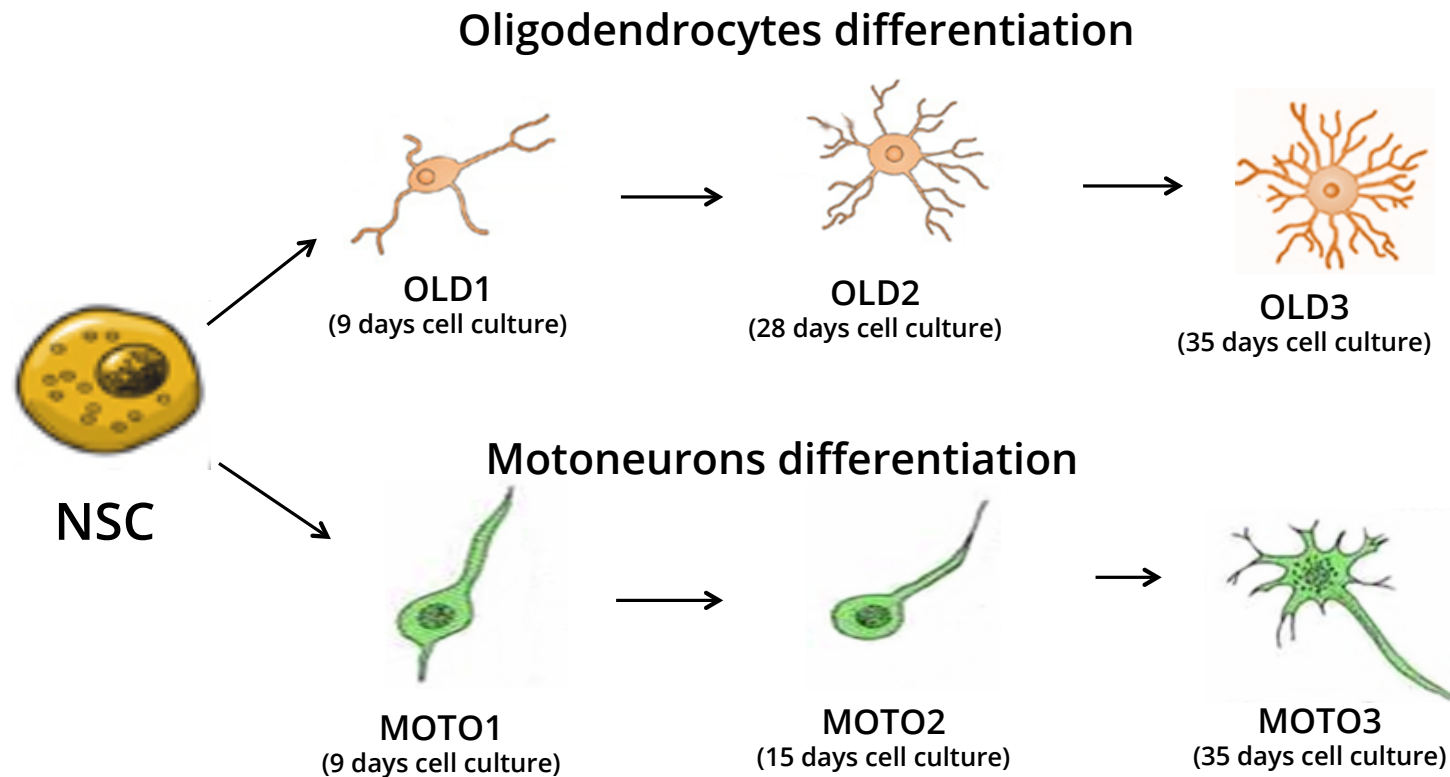
Ana Conesa, PhD
Genomics of Gene Expression Lab
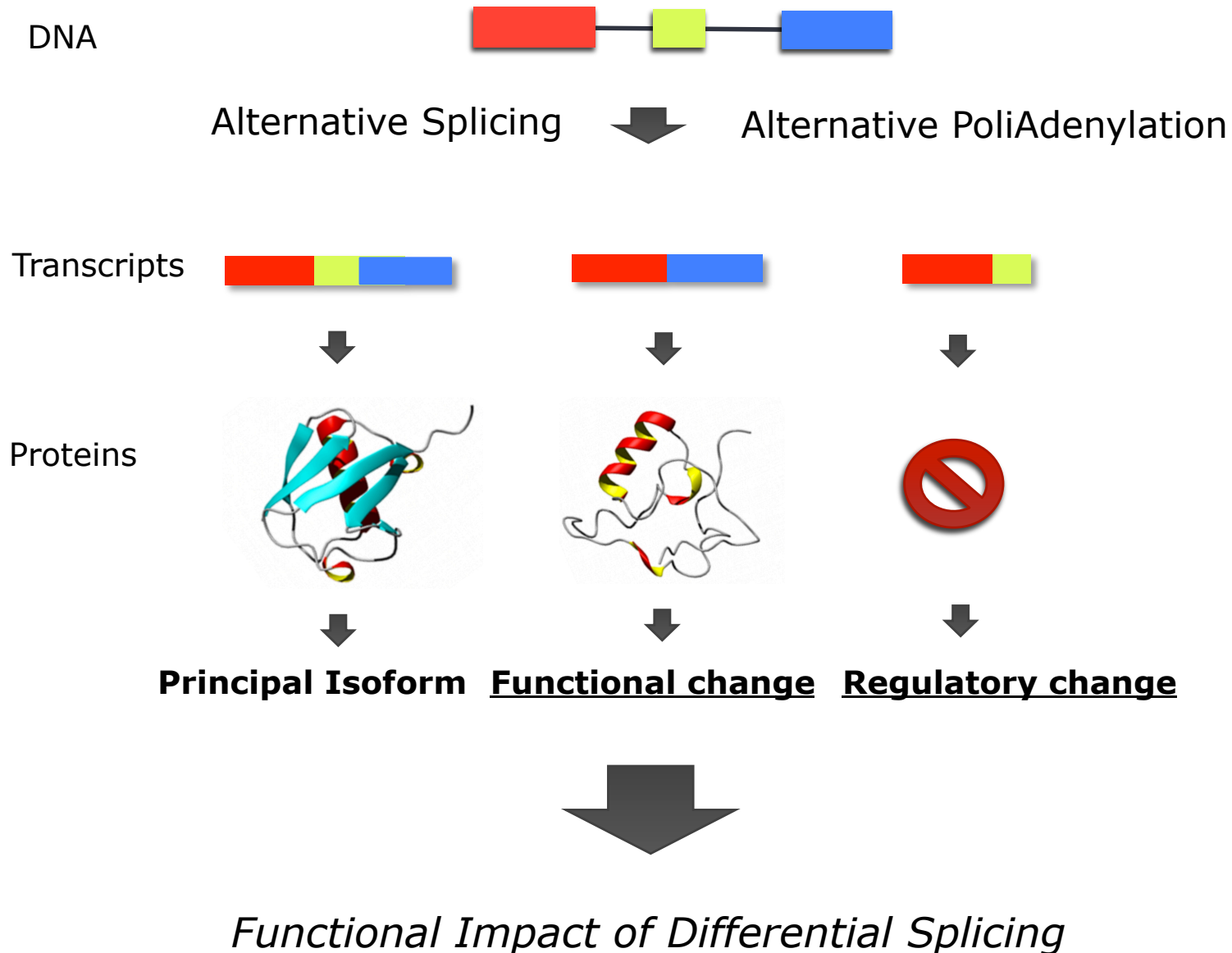CIPF/UF

PRINCIPE FELIPE
CENTRO DE INVESTIGACION
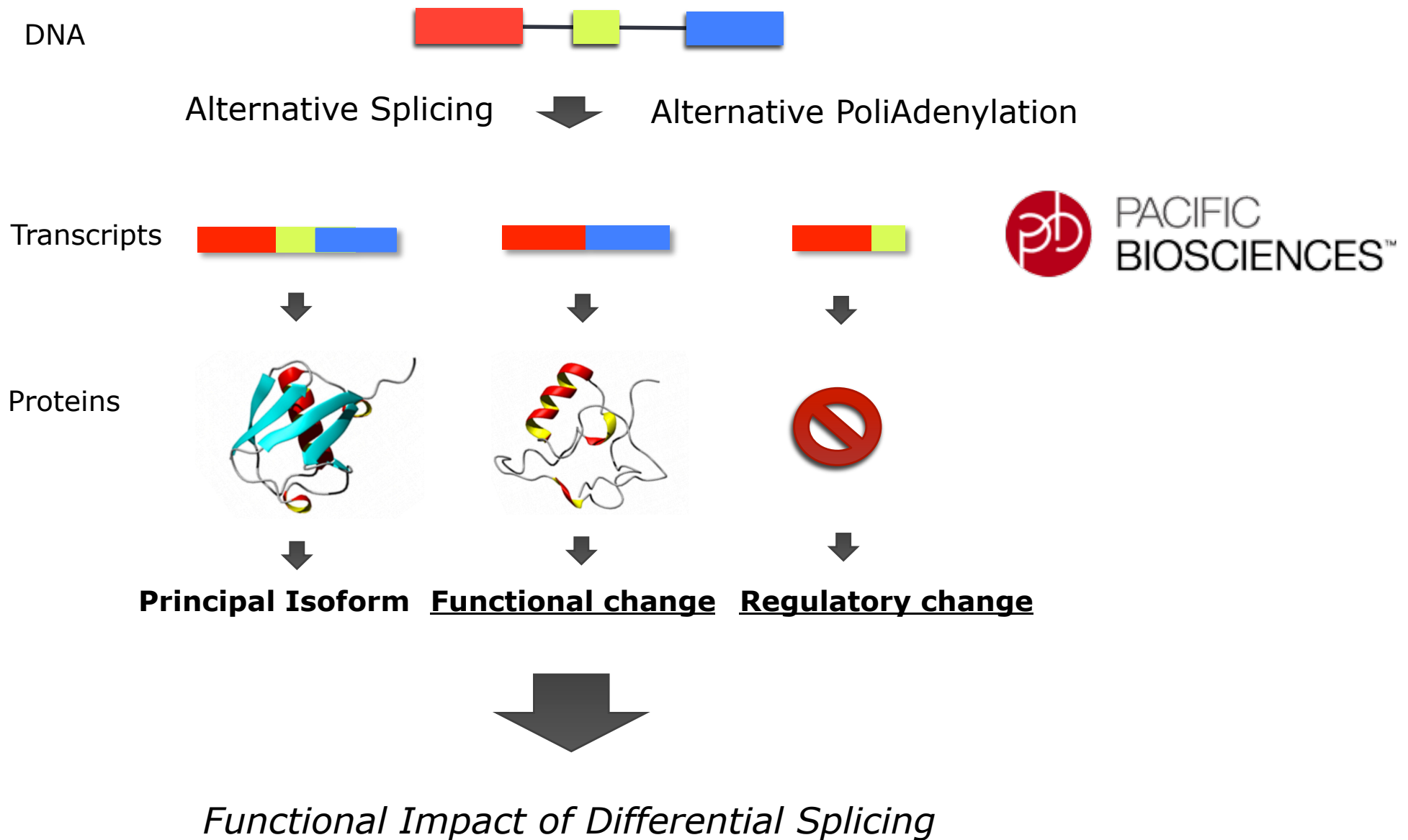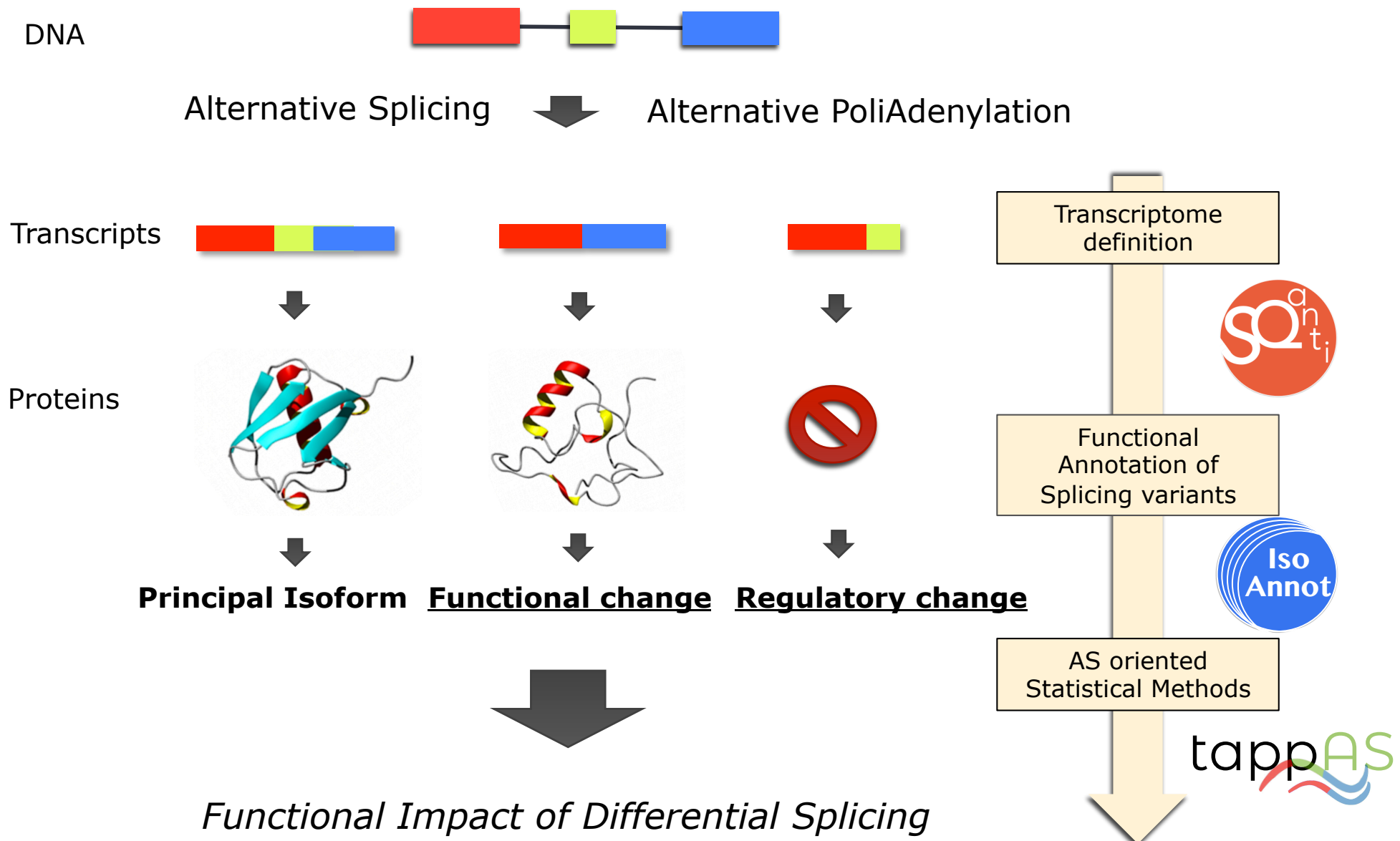
UF
UNIVERSITY of
FLORIDA

Oligodendrocytes differentiation

OLD1
(9 days cell culture)

OLD2
(28 days cell culture)

OLD3
(35 days cell culture)

NSC

Motoneurons differentiation

MOTO1
(9 days cell culture)

MOTO2
(15 days cell culture)

MOTO3
(35 days cell culture)

DNA

Alternative Splicing　　➜　　Alternative PoliAdenylation

Transcripts

Proteins

**Principal Isoform**　**Functional change**　**Regulatory change**

*Functional Impact of Differential Splicing*

DNA

Alternative Splicing → Alternative PoliAdenylation

Transcripts

PACIFIC BIOSCIENCES™

Proteins

**Principal Isoform**    **Functional change**    **Regulatory change**

*Functional Impact of Differential Splicing*

# Functional Implications of Differential Splicing

DNA

Alternative Splicing ➡ Alternative PoliAdenylation

Transcripts

Transcriptome
definition

Proteins

Functional
Annotation of
Splicing variants

**Principal Isoform**  **Functional change**  **Regulatory change**

AS oriented
Statistical Methods

*Functional Impact of Differential Splicing*

# Structural and Quality Annotation of Transcript Isoforms

https://bitbucket.org/ConesaLab/sqanti

## 1. Classification
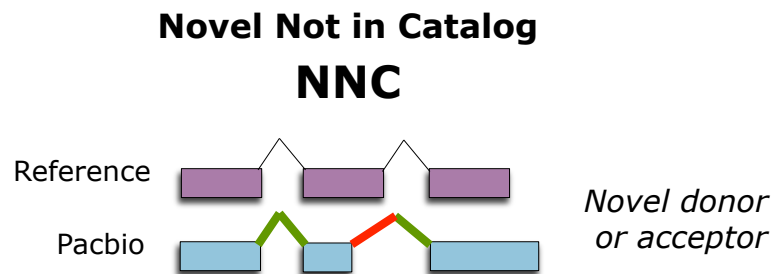
**Known Isoforms**

**Full-Splice Match**
**FSM**



**Incomplete-Splice Match**
**ISM**

## 1. Classification

### Novel Isoforms – Known genes

**Novel In catalog**
**NIC**

Reference

Pacbio

Pacbio

*Known donors and acceptors*

**Novel Not in Catalog**
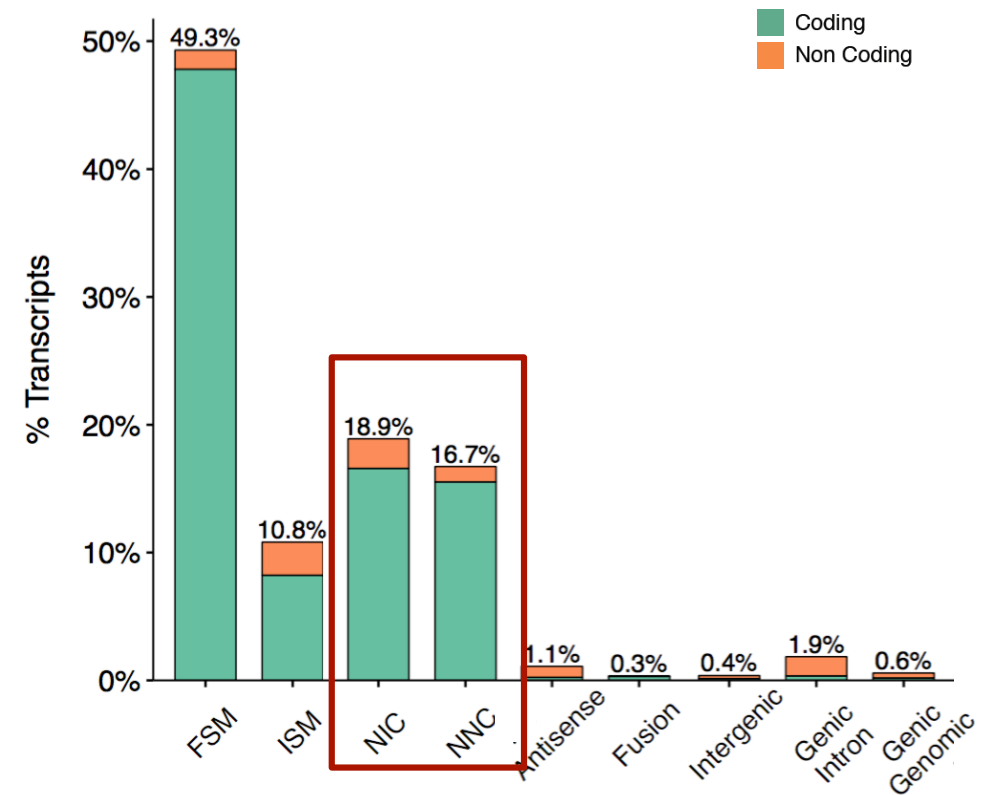**NNC**

Reference

Pacbio

*Novel donor or acceptor*

# 1. Classification

## 1. Classification

35 % of novel isoforms in mouse…

Are all of them real?

# Transcriptome characterization

1. Classification

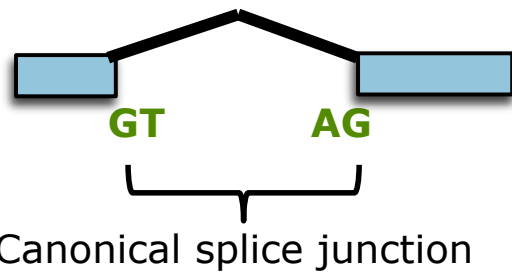2. **QC descriptors**

1. Classification

2. **QC descriptors: SJ canonical status**


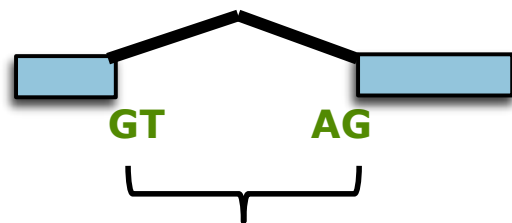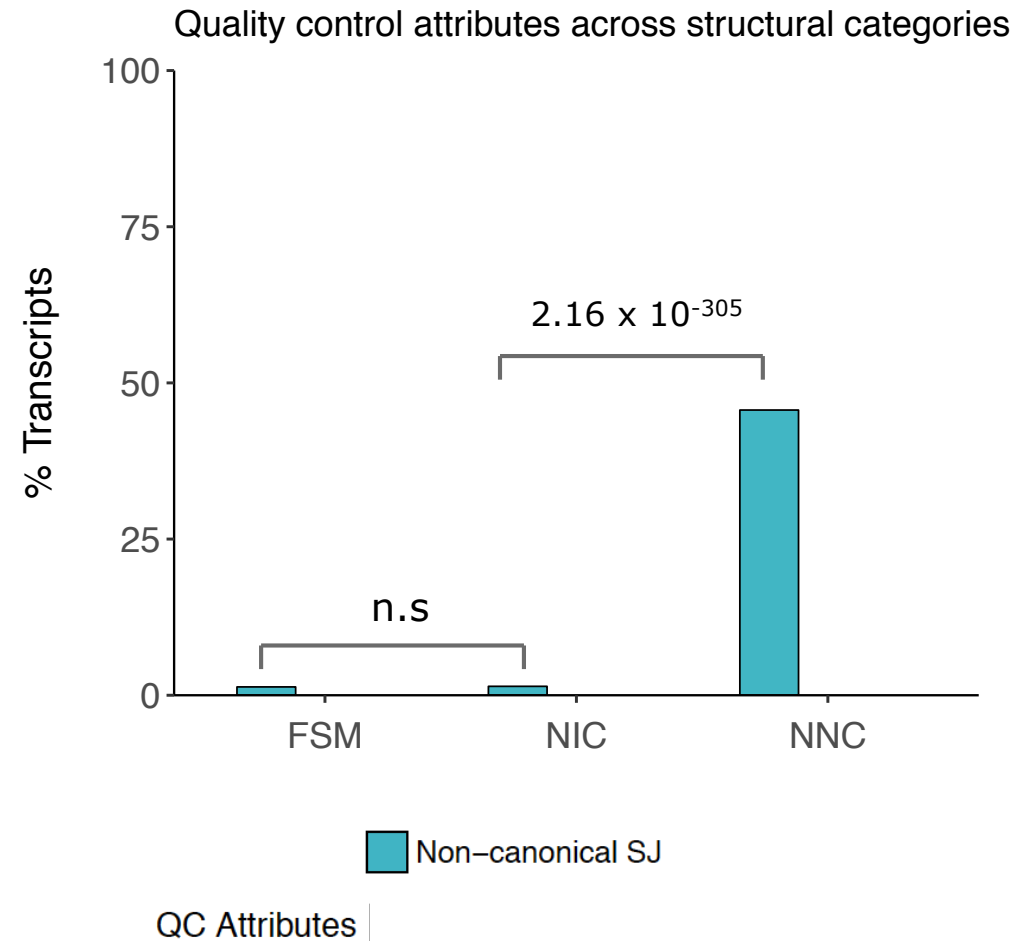
Canonical splice junction

≈ 98,7 % of canoncal SJ in mammalian*

**97,7 % of total splice junctions in our neural transcriptome are canonical**
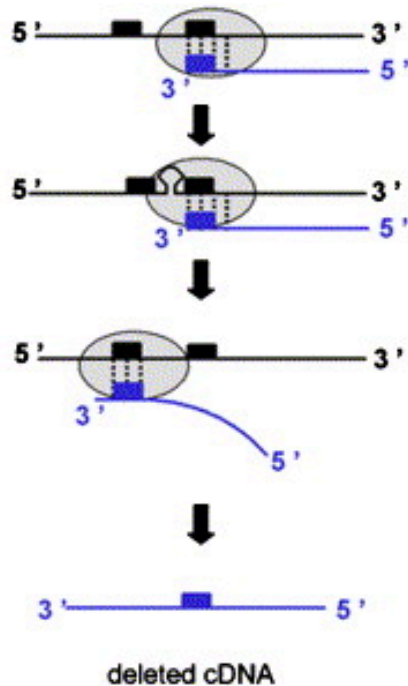
*Burset et al, 2000

1. Classification

2. **QC descriptors: SJ canonical status**

GT      AG

Canonical splice junction

≈ 98,7 % of canoncal SJ in mammalian*

**97,7 % of total splice junctions in our neural transcriptome are canonical**

Quality control attributes across structural categories

% Transcripts

100

75

$2.16 \times 10^{-305}$

50

25

n.s

0

FSM      NIC      NNC

☐ Non–canonical SJ

QC Attributes

*Burset et al, 2000

1.  Classification

**2.  QC descriptors:  <span style="color:darkred">RT-switching</span>**



deleted cDNA

- Reverse transcriptase template switching
- Caused by RNA secondary structure and repeated regions.
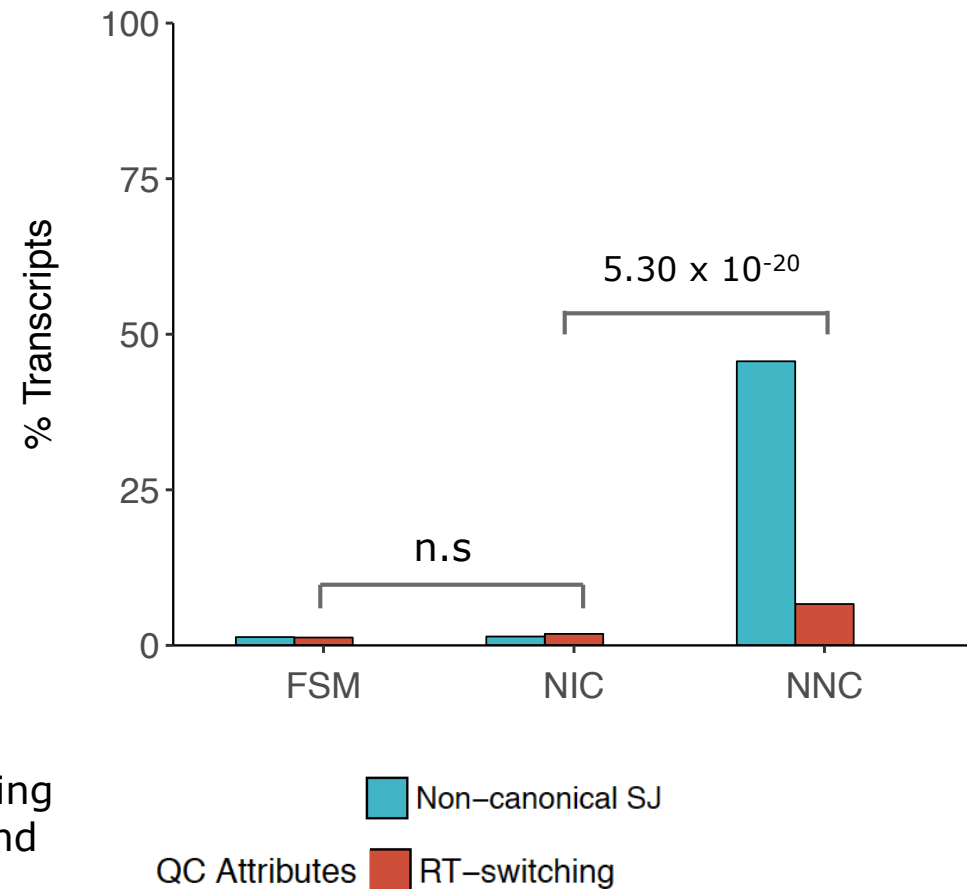- Appears as novel splice junctions

1. Classification

2. **QC descriptors: RT-switching**



Quality control attributes across structural categories

$5.30 \times 10^{-20}$

n.s

% Transcripts

FSM    NIC    NNC

Non-canonical SJ

QC Attributes    RT-switching

deleted cDNA

- Reverse transcriptase template switching
- Caused by RNA secondary structure and repeated regions.
- Appears as novel splice junctions

1. Classification

2. **QC descriptors: PolyA intra-priming**

preRNA



Oligo dT

- oligodT can prime outside polyA tail in A rich regions inside transcribed regions.

- We looked for transcripts showing >= 80% Adenines in the 20 nts downtream "detected" 3' end



**Novel Genes**

1. Classification

2. **QC descriptors: SJ support**



Supported splice junction

***Illumina Reads from same cDNA sequenced by PacBio***

Quality control attributes across structural categories

$2.22 \times 10^{-193}$

n.s

% Transcripts

FSM    NIC    NNC

QC Attributes
- Non-canonical SJ
- RT-switching
- Not coverage SJ

## Transcript level attributes

1. Transcript Classification
    1. Reference Gene match
    2. Reference Transcript match
    3. Structural Category

2. Structural characteristics
    1. Detected/Reference Length
    2. Detected/Reference number of exons
    3. Distance to nearest annotated TSS
    4. Distance to nearest annotated TTS
    5. Bite

3. Quality Control attributes
    1. RT-switching
    2. PolyA Intrapriming
    3. Canonical status
    4. Indels near SJ

4. Support
    1. Minimum splice junction coverage
    2. Minimum sample coverage
    3. Minimum coverage position
    4.. Number of Full-length reads supporting the transcript

5. Expression levels:
    1. Transcript level
    2. Gene level

6. Coding potential
    1. Coding/non coding
    2. ORF/CDS length
    3. CDS start and end positions

...

## Junction level attributes

1. Junction Classification
    1. Novel/Known
    2. Splice site motif

2. Structural characteristics
    1. Diffterence to nearest ref. donor
    2. Diffterence to nearest ref. acceptor
    3. Bite

3. Quality Control attributes
    1. Canonical
    2. Rts_junction
    3. Indel near junc

4. Support
    1. Samples with cov
    2. Total coverage
    3. Coverage per sample

# Functional Annotation of Splicing Variants

Transcriptome definition

**Functional Annotation of Splicing variants**
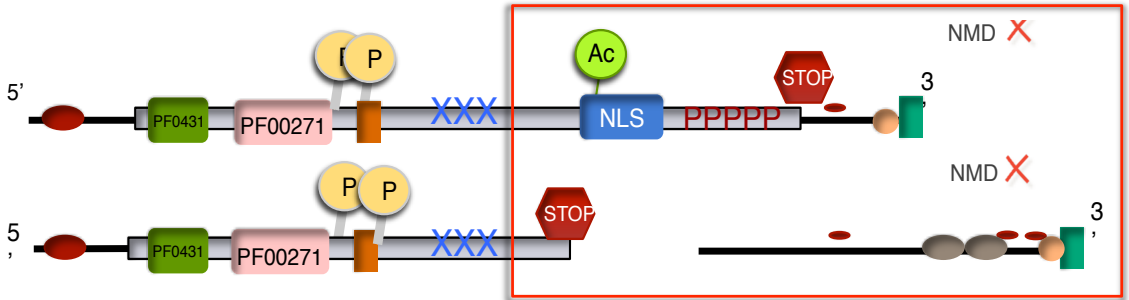
Functional Profiling at isoform level

Splicing Variants (known and novel)

Sequenced-based predictors

Coordinate-based transference algorithms

TMHMM SignalP HMMER-Pfam

UTRscan

RepeatMasker

NMD algorithm

Uniprot

Phospho sitePlus

TargetScan

CLIPdb

PFAM Domains

NMD prediction

Post-translational modifications

Motifs

UTR motifs

Transmembranes regions

Binding CA-bind Crosslink Np-bind

miRNA binding

Active Sites

Compositional bias

Repeats

Signal Peptides

RNA binding proteins

AU rich regions



Domain and Motif annotation at Isoform Resolution

# Functional Profiling at Isoform Level

## Structural Annotation and Functional Annotation



Reference Annotation provided

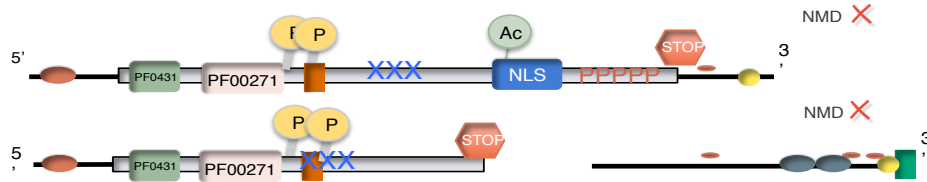OPTIONAL: User-defined

## Isoform quantification



User-defined

INPUT

tappAS

## Structural Annotation and Functional Annotation

## Isoform quantification



Reference Annotation provided

OPTIONAL: User-defined

**Module 1**

Visualization Interface

tappAS

INPUT

User-defined

## Structural Annotation and Functional Annotation

## Isoform quantification



Reference Annotation provided

OPTIONAL: User-defined

**Module 1**

Visualization Interface

**Module 2**

Motif Diversity
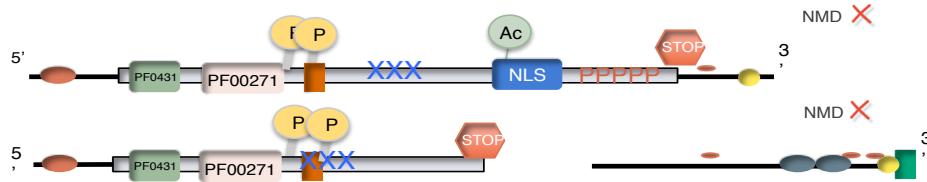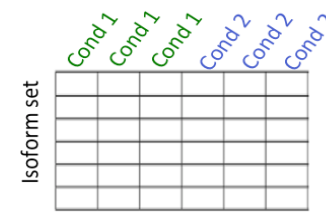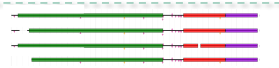
*Functional Diversity Analysis (FDA)*

CDS/UTRs Diversity

INPUT

User-defined

# Functional Profiling at Isoform Level

## Structural Annotation and  Functional Annotation

Reference Annotation provided
OPTIONAL: User-defined

**Module 1**
Visualization Interface

**Module 2**

Motif Diversity

*Functional Diversity Analysis (FDA)*

CDS/UTRs Diversity
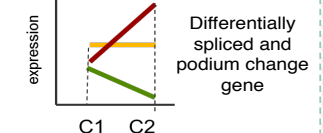
## Isoform quantification

INPUT

User-defined

**Module 3**

Differentially Spliced genes

Mayor Isoform Switching

Differentially expressed genes/transcripts/ORFs

*Differential Analysis (DSA/DEA)*

Differentially spliced and podium change gene

# Functional Profiling at Isoform Level

## Structural Annotation and Functional Annotation



NMD ✗

NMD ✗

Reference Annotation provided
OPTIONAL: User-defined

## Isoform quantification



INPUT

User-defined

**Module 1**
Visualization Interface

**Module 2**

Motif Diversity

CDS/UTRs Diversity

*Functional Diversity Analysis (FDA)*

miRNA-A

tappAS

**Module 3**

Differentially Spliced genes

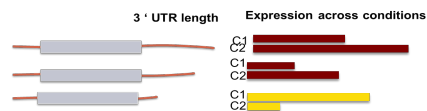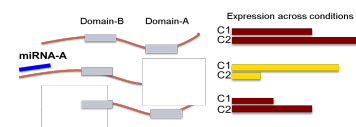Mayor Isoform Switching

Differentially expressed genes/transcripts/ORFs

*Differential Analysis (DSA/DEA)*

expression

Differentially spliced and podium change gene

C1   C2

## Integrative methods

**Module 4**

UTR length regulation

3 ' UTR length   Expression across conditions

C1
C2
C1
C2
C1
C2

*UTR Differential Splicing (UtrDS)*

Regulation of annotated elements
Motifs, Domains, miRNA binding sites

Domain-B  Domain-A

miRNA-A

Expression across conditions

C1
C2
C1
C2
C1
C2

*Motif Differential Splicing (mDS)*

# Functional Profiling at Isoform Level

**Structural Annotation and Functional Annotation**

Reference Annotation provided
OPTIONAL: User-defined

**Isoform quantification**

INPUT

User-defined

**Module 1**
Visualization Interface

**Module 2**
Motif Diversity

CDS/UTRs Diversity

*Functional Diversity Analysis (FDA)*

tappAS

**Module 3**
Differentially Spliced genes

Mayor Isoform Switching

Differentially expressed genes/transcripts/ORFs

*Differential Analysis (DSA/DEA)*

Differentially spliced and podium change gene

Integrative methods

**Module 4**
UTR length regulation

Regulation of annotated elements
Motifs, Domains, miRNA binding sites

*UTR Differential Splicing (UtrDS)*

*Motif Differential Splicing (mDS)*

**Module 5**
Functional enrichment over any annotated category

*Functional Enrichment and Gene Set Analysis (FEA/GSA)*

26

**Gene A**

Domain-B    Domain-A

Expression across conditions

AU-rich

Significant differential usage of AU-rich motif in Gene A?

AU-rich associated expression

Not AU-rich associated expression

GLM

**DS AU-rich region**

**AU-rich element favored in condition 1 by DS**

*Regulation of protein motifs by differential splicing*

_Experimental validation_



Nuclear Localization Signal

4

# Rufy3
## Generation of neuronal polarity formation and axon growth

*Regulation of UTR motifs by differential splicing*

| # | Gene | Feature | Feature Id | Position | FDSA Result | Q-Value | Favored Condi... | PodiumChg | TotalChg + |
|---|------|---------|-----------|----------|-------------|---------|------------------|-----------|-----------|
| 1 | Rufy3 | miRNA | mmu-miR-590-3p | T88641599-88... | DS | 5.2611E-9 | NSC | NO | 48.79 |



Rufy3 - RUN and FYVE domain containing 3

Transcripts View - Aligned for Project 'NSCtoOLD'

**mir590 binding site in splicing variant with longer UTR**

*Experimental validation **ongoing**:*

- *Analysis of miRNA 590 expression.*
- *Validation of the mirna binding site in Isoform 1 by miRNA pull-down assays.*

DE genes    DS genes

mRNA processing

Neurological system process

Cell-cell contact zone

40 genes

80 genes

FDR<0.05

Neuron projection

# http://tappas.org

# http://tappas.org

# http://tappas.org

# Acknowledgements





**UF**
William Farmerie
Eric Triplett
Lauren McIntyre

**UCI**
Ali Mortazavi

**Pacbio**
Liz Tseng

**CIPF**
Victoria Moreno
Susana Rodriguez