# Deciphering splicing with sparse regression techniques in the era of high-throughput RNA sequencing

Elsa Bernard

CBIO Mines ParisTech, Institut Curie, INSERM U900

Journées MAS 2016

# Deciphering splicing with sparse regression techniques in the era of high-throughput RNA sequencing

- **(alternative) splicing.** Functional importance, human diseases, therapies.

- **RNA-seq.** Next generation sequencing of RNA molecules.

- **sparse regression.** Estimating splicing variants.

# Deciphering splicing with sparse regression techniques in the era of high-throughput RNA sequencing

- **(alternative) splicing.** Functional importance, human diseases, therapies.

- **RNA-seq.** Next generation sequencing of RNA molecules.

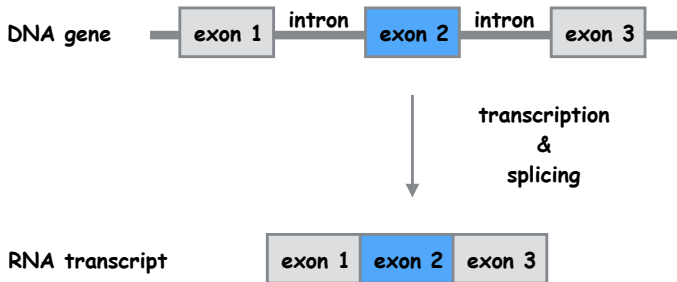- **sparse regression.** Estimating splicing variants.

# Deciphering splicing with sparse regression techniques in the era of high-throughput RNA sequencing

- **(alternative) splicing.** Functional importance, human diseases, therapies.

- **RNA-seq.** Next generation sequencing of RNA molecules.

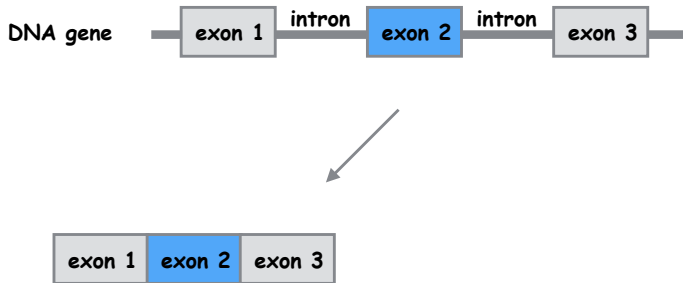- **sparse regression.** Estimating splicing variants.

## Deciphering splicing with sparse regression techniques in the era of high-throughput RNA sequencing

- **(alternative) splicing.** Functional importance, human diseases, therapies.

- **RNA-seq.** Next generation sequencing of RNA molecules.

- **sparse regression.** Estimating splicing variants.
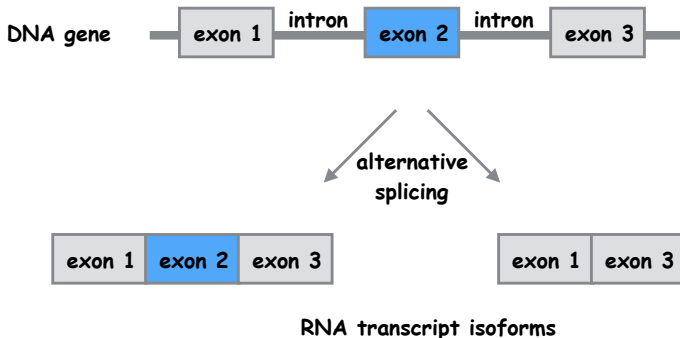
# Split genes and splicing of introns



"The discovery of split genes has been of fundamental importance for today's basic research in biology, as well as for more medically oriented research concerning the development of cancer and other diseases"

Nobel Prize Press Release, 1993.

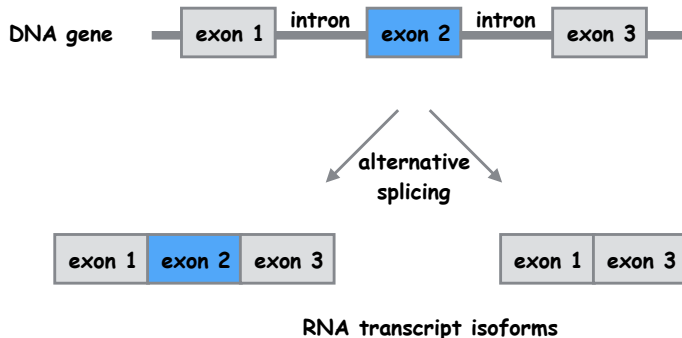# Alternative splicing produces transcript isoforms

# Alternative splicing produces transcript isoforms



- The splicing pattern determines the final genetic message.
- In human, 28k genes give 120k known transcript isoforms (Pal et al., 2012).
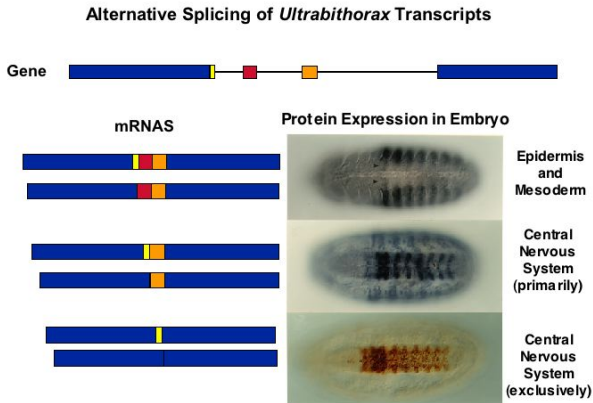
# The isoform identification and quantification problem



Given a biological sample, can we:

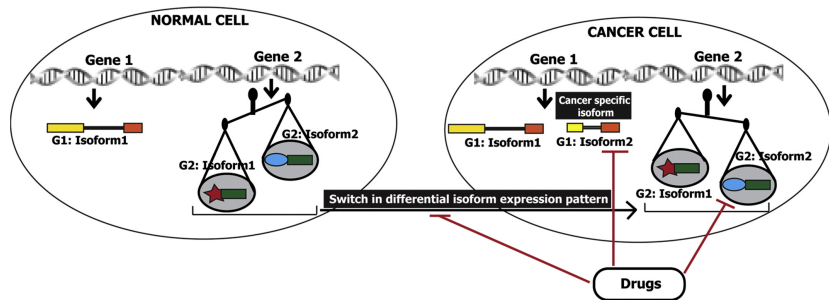1. identify the isoforms expressed by each gene?
2. quantify their abundances?

# Functional importance of alternative splicing

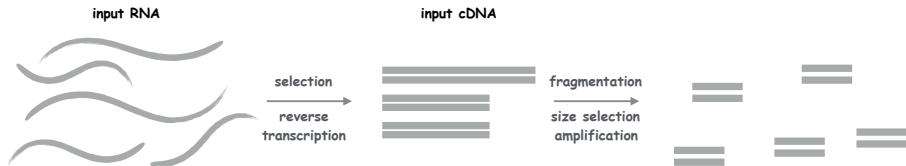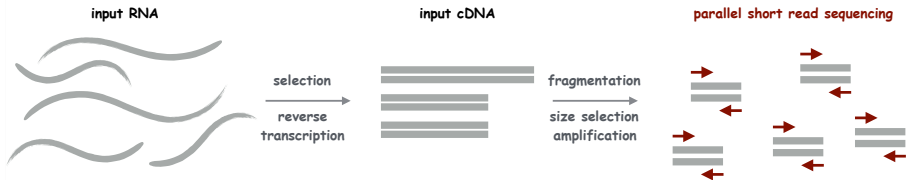- Developmental regulation of alternative splicing in Drosophila:



Alternative Splicing of *Ultrabithorax* Transcripts
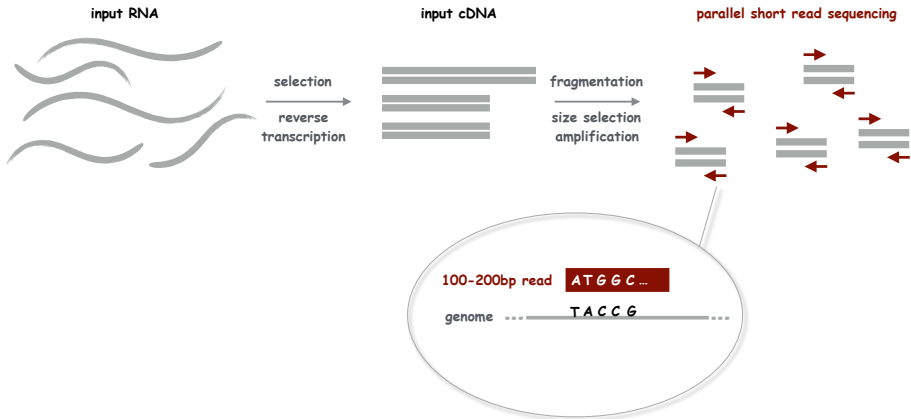
*http://orchid.bio.cmu.edu/research.html*

# Drug targets



*(Pal et al., 2012)*

# RNA-seq: shear RNA into pieces and sequence

Sample 1 ... Sample t ... Sample T

Isoforms 1? Isoforms t? Isoforms T?

**One-sample:** can we perform accurate de novo isoform reconstruction for one given RNA-seq sample?

## 1) the one-sample case

**FlipFlop** Fast Lasso based Isoform Prediction as a FLOw Problem

## 2) the multi-sample case

Isoform detection from multiple RNA-seq samples

## 3) clinical application

Quantify abnormal splicing from targeted RNA-seq

**1) the one-sample case**

**FlipFlop** Fast Lasso based Isoform Prediction as a FLOw Problem

**2) the multi-sample case**

**Isoform detection from multiple RNA-seq samples**

**3) clinical application**

**Quantify abnormal splicing from targeted RNA-seq**

**1) the one-sample case**

**FlipFlop** Fast Lasso based Isoform Prediction as a FLOw Problem

**2) the multi-sample case**

**Isoform detection from multiple RNA-seq samples**

**3) clinical application**

**Quantify abnormal splicing from targeted RNA-seq**

# Outline

### 1) the one-sample case

**FlipFlop** Fast Lasso based Isoform Prediction as a FLOw Problem
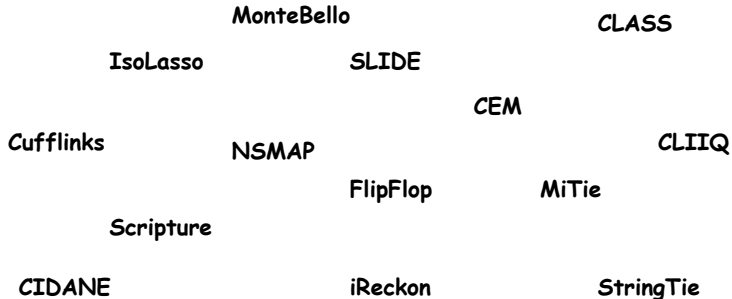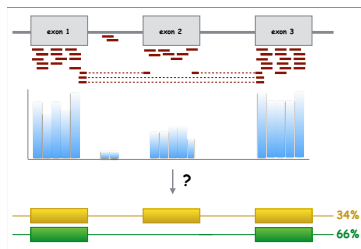
### 2) the multi-sample case

Isoform detection from multiple RNA-seq samples

### 3) clinical application

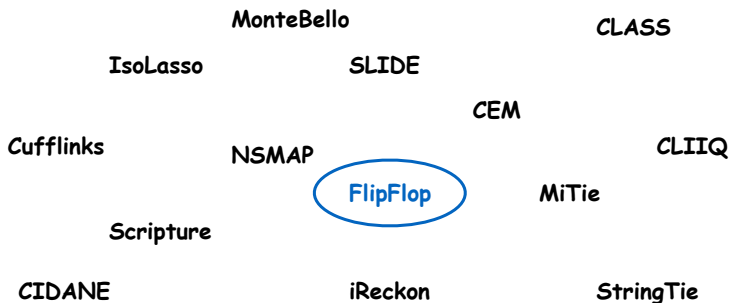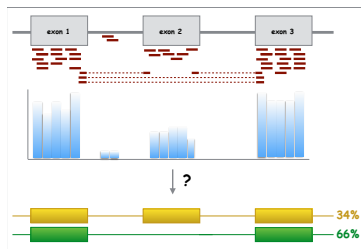Quantify abnormal splicing from targeted RNA-seq

# Genome-guided isoform reconstruction

- Input: spliced alignment of reads against reference genome
- Goal: reconstruct transcripts (multi-assembly problem)



MonteBello

CLASS

IsoLasso

SLIDE

CEM

Cufflinks

NSMAP

CLIIQ

FlipFlop

MiTie

Scripture

CIDANE

iReckon

StringTie

# What's new?

- Input: spliced alignment of reads against reference genome
- Goal: reconstruct transcripts (multi-assembly problem)



MonteBello

CLASS

IsoLasso

SLIDE

CEM

Cufflinks

NSMAP

CLIIQ

FlipFlop

MiTie

Scripture

CIDANE

iReckon

StringTie

# Contributions

1. **No need to filter** candidate transcript isoforms

2. **Faster** than existing methods that solve the same problem

   *Flow methods*

3. Adapted to **long reads**   *Particular splicing graph*

4. R package (open-access, maintained, parallelizable)   *Bioconductor*

# Contributions

1. **No need to filter** candidate transcript isoforms

2. **Faster** than existing methods that solve the same problem

   **Flow methods**

3. Adapted to **long reads** | Particular splicing graph

4. R package (open-access, maintained, parallelizable) | Bioconductor

# Contributions

1. **No need to filter** candidate transcript isoforms
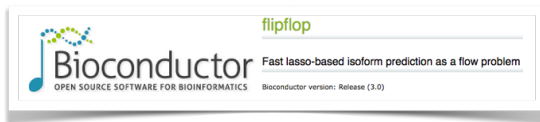2. **Faster** than existing methods that solve the same problem

Flow methods

3. Adapted to **long reads** **Particular splicing graph**

4. R package (open-access, maintained, parallelizable) Bioconductor

# Contributions

1. **No need to filter** candidate transcript isoforms

2. **Faster** than existing methods that solve the same problem

} Flow methods

3. Adapted to **long reads** } Particular splicing graph

4. R package (open-access, maintained, parallelizable) } **Bioconductor**
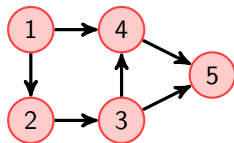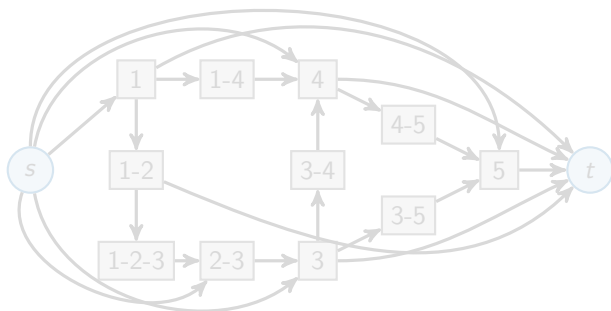
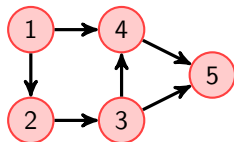# Isoforms are paths in a graph

- Splicing graph for a gene with 5 exons:



- FlipFlop graph: **1 type of read ↔ 1 node**

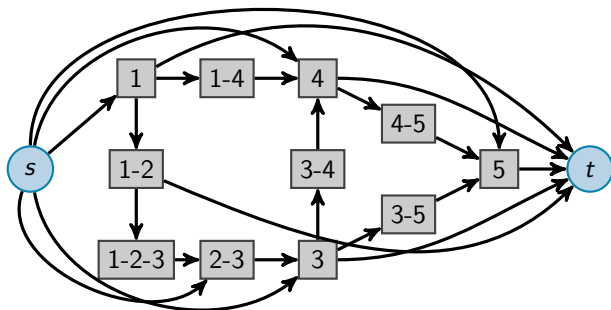# Graph adapted to long reads

- Splicing graph for a gene with 5 exons:



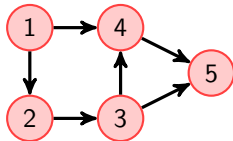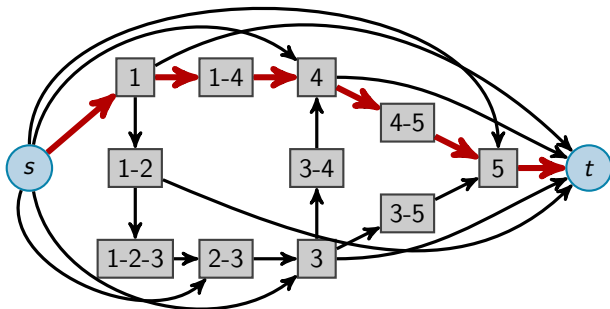- FlipFlop graph: **1 type of read ↔ 1 node**

# Isoforms are paths in a graph

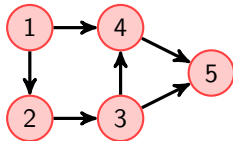- Splicing graph for a gene with 5 exons:
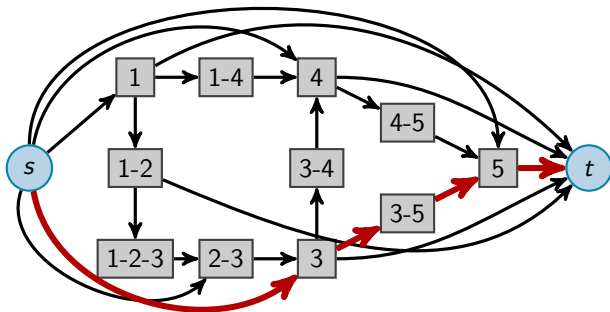


- FlipFlop graph: **one path with abundance $\theta_1$**

# Isoforms are paths in a graph

- Splicing graph for a gene with 5 exons:



- FlipFlop graph: **another path with abundance $\theta_2$ ...**

$n$ **exons** $\rightarrow$ $\sim 2^n$ **paths/candidate isoforms**

feature selection problem with $\sim 10^3$ candidates for 10 exons
and $\sim 10^6$ for 20 exons

**Minimum path cover**
- Cufflinks, CLASS
- ✗ **do not use read counts**

**Sparse regression**
- IsoLasso, NSMAP, SLIDE, CEM, iReckon, MiTie, **FlipFlop**, CIDANE
- ✓ **use read counts**

- Estimate $\theta$ sparse by solving:

$$\min_{\theta} \quad \underbrace{\mathcal{L}(\theta)}_{} \quad + \quad \underbrace{\lambda\|\theta\|_1}_{}$$

big vector!

**fit to the data**
do you well explain
read counts with the
selected isoforms?
e.g: minus log–likelihood

**sparsity–inducing effect**
you select a few isoforms
among many candidates

- **Computationally challenging**
  $\rightarrow$ IsoLasso: strong filtering
  $\rightarrow$ NSMAP, SLIDE: number of exons cut-off

- **FlipFlop**
  $\rightarrow$ no filtering
  $\rightarrow$ no exon restrictions

# Isoform deconvolution with the $\ell_1$-norm penalization

- Estimate $\theta$ sparse by solving:

$$\min_{\theta} \quad \underbrace{\mathcal{L}(\theta)}_{\substack{\textbf{fit to the data} \\ \text{do you well explain} \\ \text{read counts with the} \\ \text{selected isoforms?} \\ \text{e.g: minus log-likelihood}}} + \underbrace{\lambda\|\theta\|_1}_{\substack{\textbf{sparsity-inducing effect} \\ \text{you select a few isoforms} \\ \text{among many candidates}}}$$

**big vector!** ↑

- **Computationally challenging**
  - → IsoLasso: strong filtering
  - → NSMAP, SLIDE: number of exons cut-off

- **FlipFlop**
  - → no filtering
  - → no exon restrictions

# Fast isoform deconvolution

The isoform deconvolution problem

$$\min_{\theta}\ \mathcal{L}(\theta) + \lambda\|\theta\|_1\ ,$$

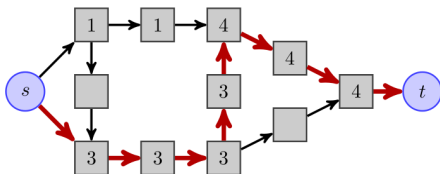is solvable in **polynomial time** with the number of nodes of the splicing graph.

Ideas:

1. the sum of isoform abundances corresponds to a **flow** on the graph
2. reformulation as a **convex cost flow problem** (Mairal and Yu, 2012)
3. recover isoforms by flow decomposition algorithm

# Combinations of isoforms are flows



(a) Reads at every node corresponding to one isoform.

(b) Reads at every node after adding another isoform.

- Linear combinations of isoforms $\Rightarrow$     Flow value on every edges
- Flow value on every edges $\Rightarrow$     Paths with given value/abundance

**Flow Decomposition**
**(linear time algorithm)**

A Novel Min-Cost Flow Method for Estimating Transcript Expression with RNA-Seq. RECOMB-2013.

# Equivalent flow problem (simpler!)



(a) Reads at every node corresponding to one isoform.    (b) Reads at every node after adding another isoform.

- $\mathcal{L}(\theta)$ depends only on the values of the flow on the vertices

- $\|\theta\|_1 = \sum_{\text{path } p} \theta_p = f_t$

- Therefore,

$$\min_\theta \mathcal{L}(\theta) + \lambda\|\theta\|_1 \quad \text{is equivalent to} \quad \min_{f \text{ flow}} \tilde{\mathcal{L}}(f) + \lambda f_t$$

# FlipFlop Summary

## Isoform detection = Path selection problem

$\sim 2^n$ variables (all paths in the splicing graph)

$\Updownarrow$

## Equivalent network flow problem

$\sim \frac{n^2}{2}$ variables (all nodes of the splicing graph)

$\downarrow$

## Network flow algorithms

Efficient algorithms. Polynomial time.

# Human Simulation: precision / recall

hg19, 1137 genes on chr1, 1million 200 bp single-end reads by transcript levels.
Simulator: http://alumni.cs.ucr.edu/~liw/rnaseqreadsimulator.html

# Speed Trial

hg19, 1137 genes on chr1, 1million reads by exon levels.

Simulator: `http://alumni.cs.ucr.edu/~liw/rnaseqreadsimulator.html`

**FlipFlop** $\rightarrow$ **transcripts reconstruction over an exponential number of candidates in polynomial time**

- http://cbio.ensmp.fr/flipflop/
- http://cbio.ensmp.fr/flipflop/experiments.html
- R package
  ```
  > source("http://bioconductor.org/biocLite.R")
  > biocLite("flipflop")
  ```

📄 E. Bernard, L. Jacob, J. Mairal and J.-P. Vert. **Efficient RNA isoform identification and quantification from RNA-seq data with network flows**. *Bioinformatics*, 2014.

# Multi-dimensional case



Sample 1 · · · ·  Sample t · · · ·  Sample T

**Multi-dimensional splicing graph**

Can we find a sparse set of paths that explains
the multi-dimensional read counts?

# Multi-dimensional case



Sample 1  ····  Sample t  ····  Sample T

**Multi-dimensional splicing graph**

**Can we find a sparse set of paths that explains the multi-dimensional read counts?**

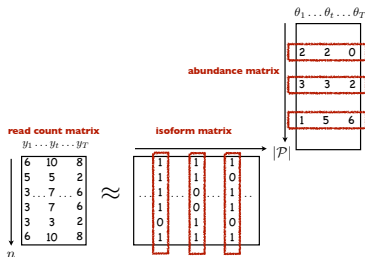# Group-Lasso strategy

# Group-Lasso strategy

# More formally



- each isoform defines a **group** $\boldsymbol{\theta}_p = \{\theta_p^t, t \in [\![1, T]\!]\}$
- the multi-sample loss is the sum of the independent losses

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{t=1}^{T} \text{loss}(y_t, \theta_t)$$

- ideally we want to solve the NP-hard $\ell_0$ problem

$$\min_{\{\boldsymbol{\theta}_p\}} \mathcal{L}(\boldsymbol{\theta}) + \lambda \sum_{p \in \mathcal{P}} \mathbf{1}_{\{\boldsymbol{\theta}_p \neq \mathbf{0}\}}$$

# More formally



- each isoform defines a **group** $\boldsymbol{\theta}_p = \{\theta_p^t, t \in [\![1, T]\!]\}$
- the multi-sample loss is the sum of the independent losses

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{t=1}^{T} \text{loss}(y_t, \theta_t)$$

- instead we solve the **group-lasso convex relaxation**

$$\min_{\{\boldsymbol{\theta}_p\}} \mathcal{L}(\boldsymbol{\theta}) + \lambda \sum_{p \in \mathcal{P}} \|\boldsymbol{\theta}_p\|_2$$

# Simulation: GroupLasso vs Merging



$$\forall t \in \{1, \ldots, T\}, \mathrm{supp}\,\theta_t = \mathrm{supp}\,\theta_o$$

# modENCODE data
## Time course development of D.melanogaster

# Multi-sample case summary

**FlipFlop** $\rightarrow$ **transcript reconstruction using several samples simultaneously leads to more statistical power**

- `http://cbio.ensmp.fr/flipflop/details.html`

📄 E. Bernard, L. Jacob, J. Mairal, E. Viara and J.-P. Vert. **A convex formulation for joint RNA isoform detection and quantification from multiple RNA-seq samples.** *BMC Bioinformatics*, 2015.

# Outline

**1) the one-sample case**

FlipFlop Fast Lasso based Isoform Prediction as a FLOw Problem

**2) the multi-sample case**

Isoform detection from multiple RNA-seq samples

**3) clinical application**

**Quantify abnormal splicing from targeted RNA-seq**

# Molecular diagnosis and splicing

- Various splicing enhancing and silencing motifs:



- Variants disrupting/creating these consensus sequences can affect normal splicing

$\Rightarrow$ **molecular diagnosis:** correct interpretation of these variants on splicing is imperative for genetic counseling

# Molecular diagnosis and splicing

- Various splicing enhancing and silencing motifs:



- Variants disrupting/creating these consensus sequences can affect normal splicing

## Development of a new diagnostic tool

- time and cost-effective identification and quantification of transcripts using targeted high-throughput RNA-seq
- extension of sparse regression techniques to a new experimental design

# Promising results on BRCA1

- BRCA1: Breast Cancer susceptibility gene
- Involved in DNA repair pathway and cell cycle
- High number of splicing events (regulated in a cell-cycle- and cell-type-specific manner)

**Combined genetic and splicing analysis of BRCA1 c.[594-2A>C; 641A>G] highlights the relevance of naturally occurring in-frame transcripts for developing disease gene variant classification algorithms**

Miguel de la Hoya[1,*], Omar Soukarieh[2], Irene López-Perolio[1], Ana Vega[3],

# Promising results on BRCA1

- BRCA1: Breast Cancer susceptibility gene
- Involved in DNA repair pathway and cell cycle
- High number of splicing events (regulated in a cell-cycle- and cell-type-specific manner)

ORIGINAL ARTICLE    *Human Molecular Genetics, 2016, Vol. 0, No. 0    1–13*

**Combined genetic and splicing analysis of BRCA1 c.[594-2A>C; 641A>G] highlights the relevance of naturally occurring in-frame transcripts for developing disease gene variant classification algorithms**

Miguel de la Hoya[1,*], Omar Soukarieh[2], Irene López-Perolio[1], Ana Vega[3],

Accurate quantification of overlapping splicing events:

31

# Thanks

Laurent Jacob

JP Vert

Julien Mairal

Eric Viara

Elodie Girard

# Supplementary Slides

## Part 1: one-sample approach

**FlipFlop** Fast Lasso based Isoform Prediction as a FLOw Problem

## Technical details

Poisson Loss:

$$\mathcal{L}(\theta) = \sum_{u \in V} \left[ Nl_u \left( \sum_{\text{path } p \ni u} \theta_p \right) - \mathbf{y}_u \log \left( Nl_u \sum_{\text{path } p \ni u} \theta_p \right) \right]$$

Flow Decomposition:

$$f_{uv} = \sum_{\text{path } p \ni (u,v)} \theta_p$$

$$\Rightarrow f_v = \sum_{u \in V} f_{uv} = \sum_{\text{path } p \ni v} \theta_p$$

Convex Cost Flow:

$$\min_{f \text{flow}} \sum_{u \in V} [Nl_u f_u - \mathbf{y}_u \log(f_u)] + \lambda f_t$$

Solved using $\epsilon$-relaxation method (Bertsekas 1998)

# Effective length



1)  $l_{\text{left}} \geq L, \quad l_{\text{right}} \geq L$  ⬛🟥🟥🟥⬜⬜⬜⬜⬜  $l_i = L - 1$

2)  $l_{\text{left}} < L, \quad l_{\text{right}} \geq L$  ⬛🟥⬜⬜⬜⬜⬜  $l_i = l_{\text{left}}$

3)  $l_{\text{left}} \geq L, \quad l_{\text{right}} < L$  ⬛🟥⬜⬜⬜  $l_i = l_{\text{right}}$

4)  $l_{\text{left}} < L, \quad l_{\text{right}} < L$  ⬛🟥⬜⬜⬜  $l_i = l_{\text{left}} + l_{\text{right}} - L + 1$

| 1 | 2 | 3 |

- FlipFlop graph:

- FlipFlop graph:

- FlipFlop graph:

- FlipFlop graph:

- FlipFlop graph:

# Graph adapted to long reads

- FlipFlop graph:

- FlipFlop graph:

- FlipFlop graph:

# Performance increases with coverage

# Extension to paired-end reads OK

# Real Data

Human: 50 million 75bp reads.

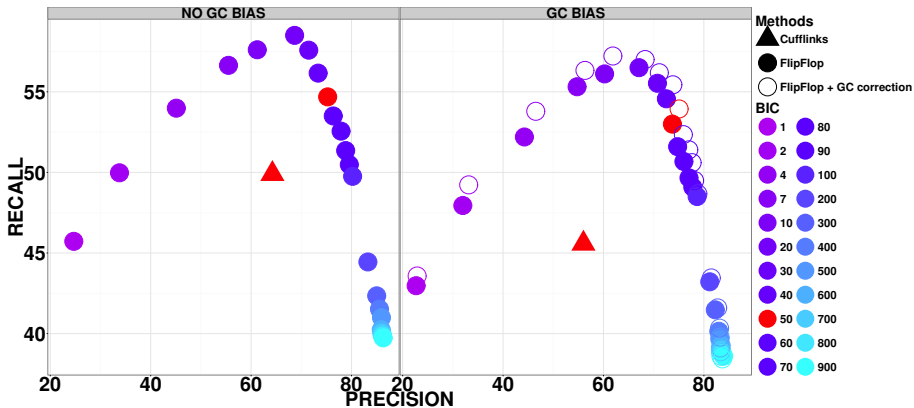# Precision-Recall curves on real data

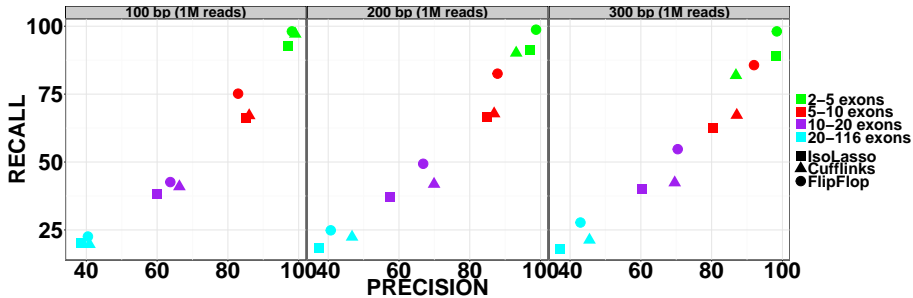# Speed comparison on real data

# GC bias - Precision-Recall curve
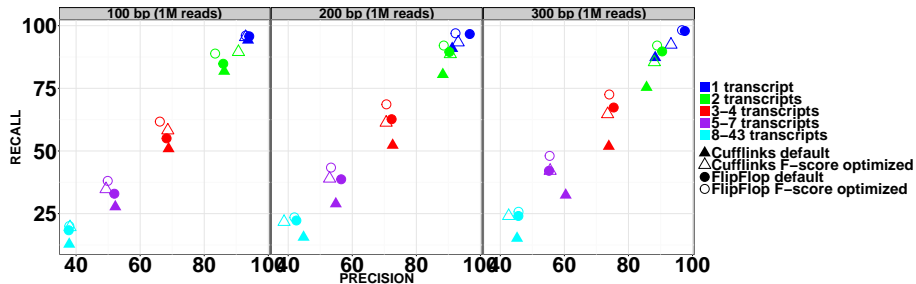
hg19, chr1, 4140 transcripts, 2million 150bp single-end reads

Simulator: FluxSimulator `http://sammeth.net/confluence/display/SIM/Home`

**Model selection**: set of solutions minimizing $\mathcal{L}(\theta) + \lambda\|\theta\|_1$ for different values of $\lambda \rightarrow$ BIC criteria

# Exon stratification
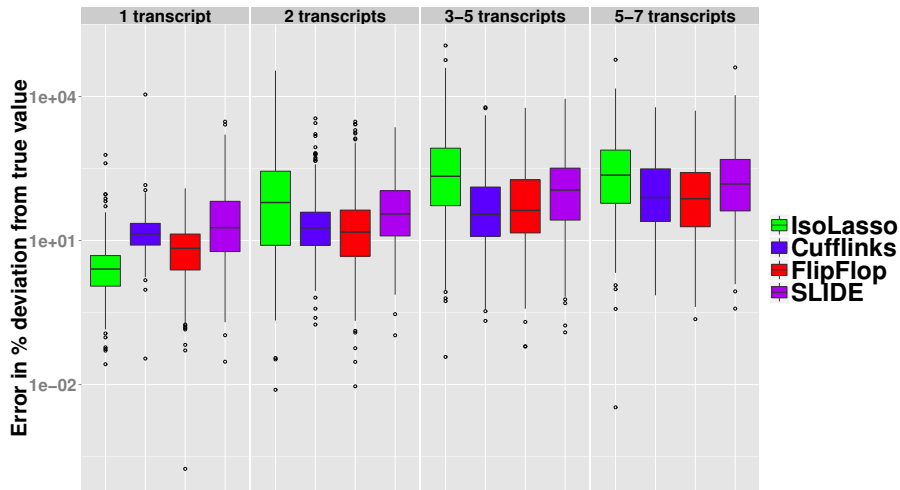
# Tuning

# Stability study

# Human Simulation: Abundances

hg19, 1137 genes on chr1, 1million 75 bp single-end reads by transcript levels.
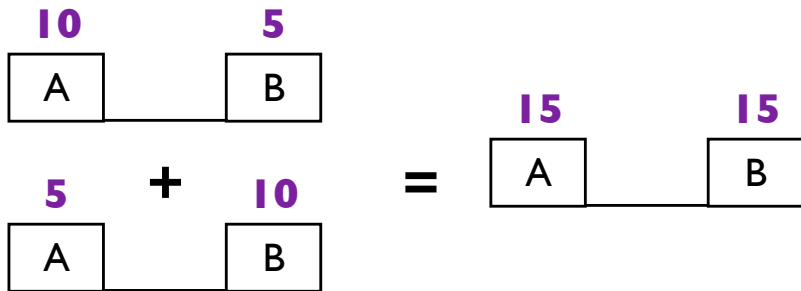
# Simulation: Deviation

hg19, 1137 genes on chr1, 1million 75 bp single-end reads by transcript levels.
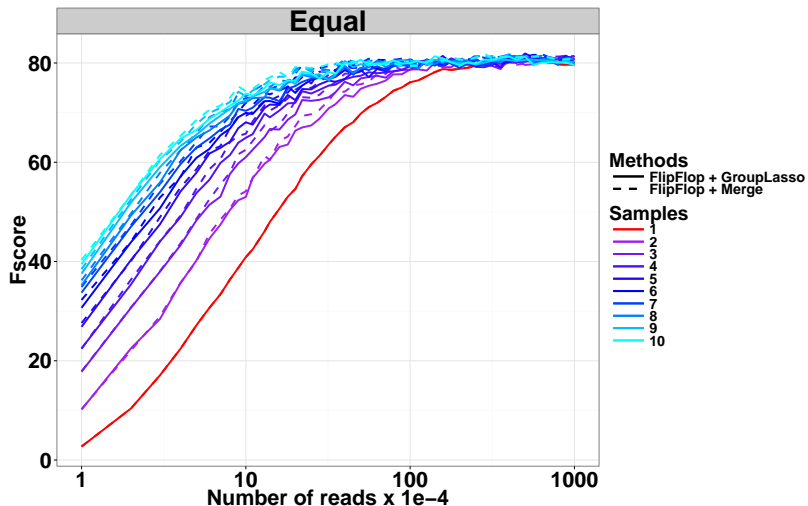
## Part 2: multi-sample approach

**Isoform detection from multiple RNA-seq sample**
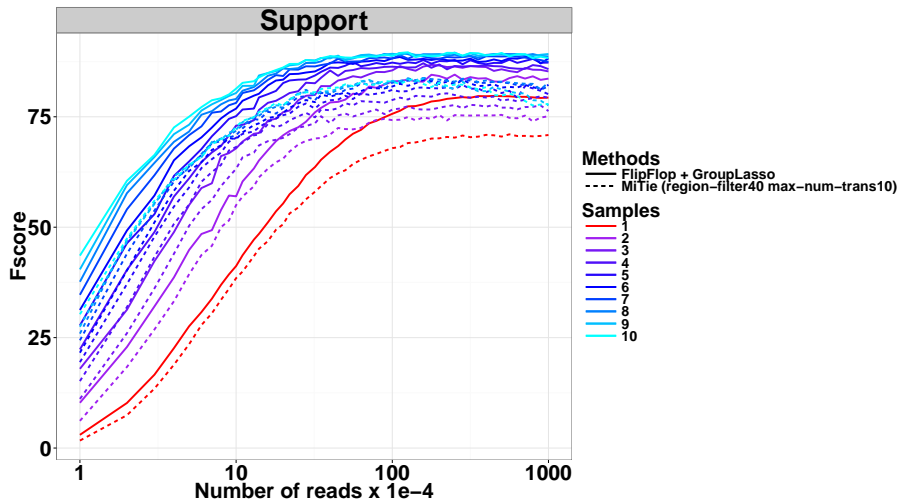
# Why Aggregating can be bad

# Toy simulation



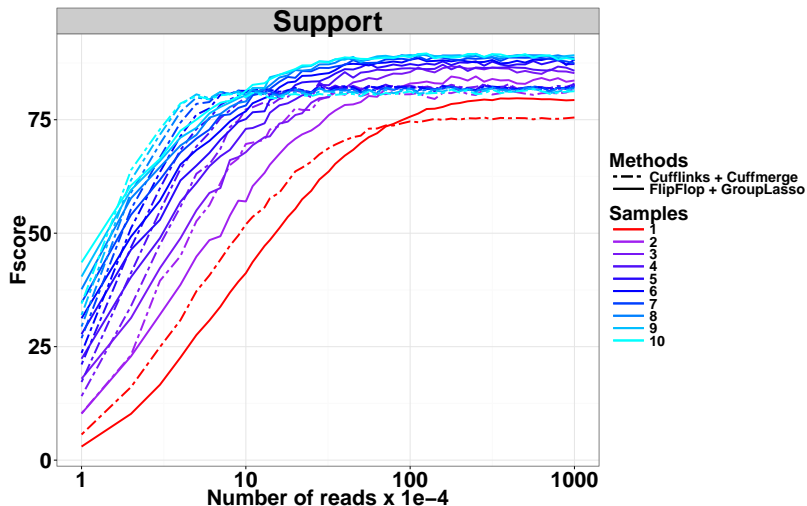$$\forall t \in \{1, \ldots, T\}, \theta_t = \theta_o$$

$$\forall t \in \{1, \dots, T\}, \operatorname{supp}\theta_t = \operatorname{supp}\theta_o$$

$\forall t \in \{1, \ldots, T\}, \text{supp}\theta_t = \text{supp}\theta_o$

# Simulation: read length