



POLITECNICO  
DI TORINO

---

# Bioinformatics methods and tools for identifying disease-associated B-cell populations, miRNA patterns and gene fusions in RNA-Seq data

---

Torino, 1<sup>st</sup> June 2017

Candidate: Giulia Paciello  
Supervisors: Elisa Ficarra and Enrico Macii



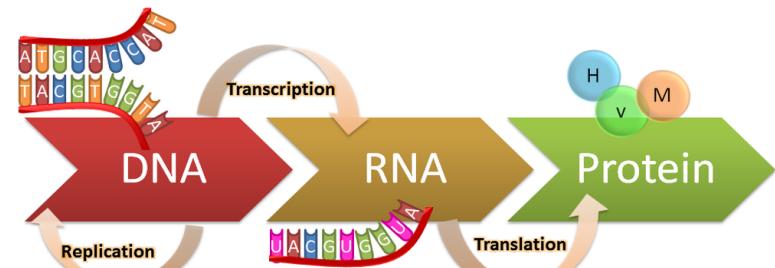
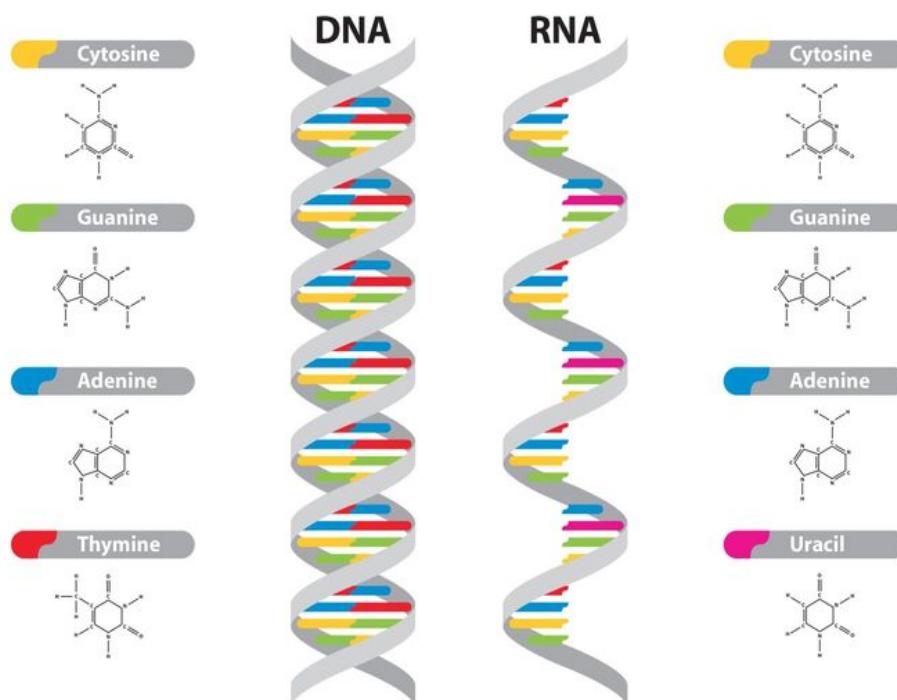
# Outline

- ✓ Introduction
- ✓ Objective
- ✓ Methodologies and tools: VDJSeq-Solver, isomiR-SEA, FuGePrior
  - Biological Background
  - Motivation
  - Algorithm
  - Results
  - Remarks and Perspectives
- ✓ Conclusions



# What does the term ‘sequencing’ mean?

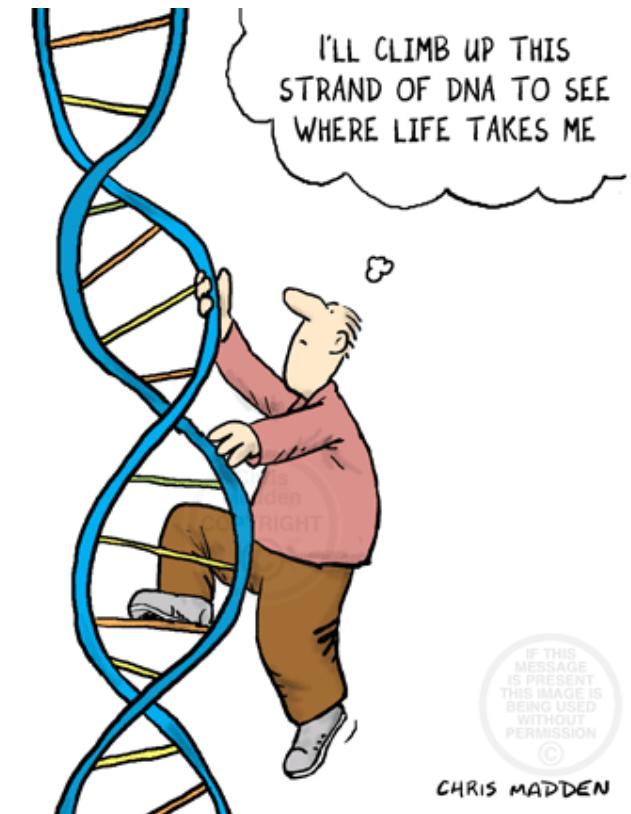
## The Central Dogma of Molecular Biology



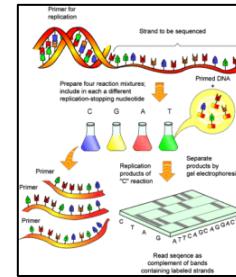
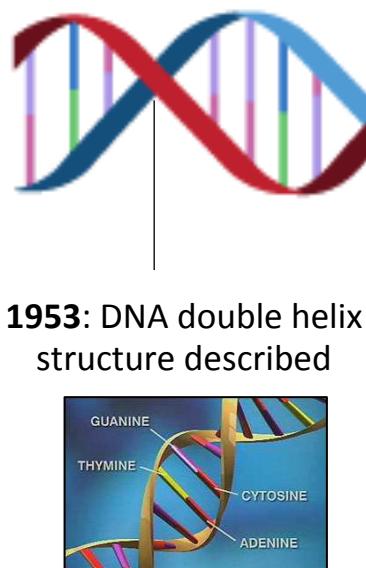
The **sequencing** is the process of determining the **sequence of nucleotide bases** in a nucleic acid molecule

# Why is sequencing important?

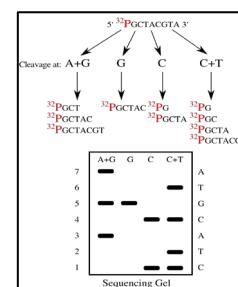
“...Decoding of the DNA that constitutes the human is a fundamental contribution toward understanding human **evolution**, the causation of **disease**, and the interplay between the environment and heredity in defining the human condition...”



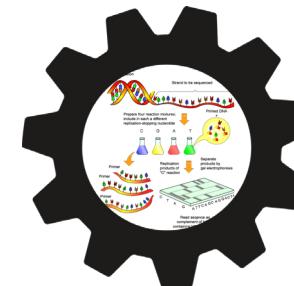
# The sequencing time line



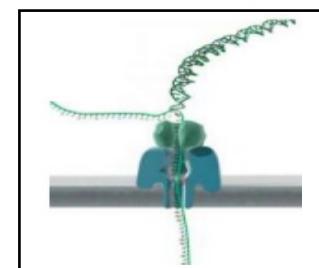
**1977:** Maxam-Gilbert sequencing technique



**1980-1990:** First Generation Sequencing



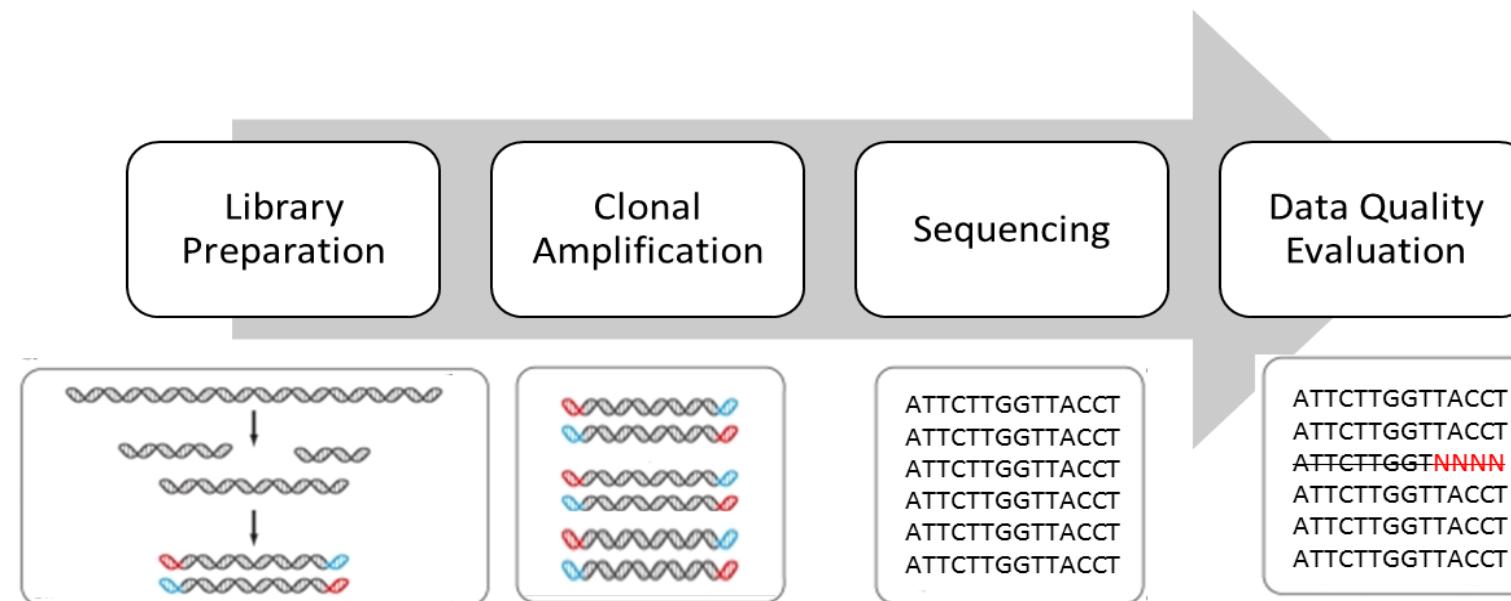
**2008:** Third Generation Sequencing



# The Next Generation Sequencing



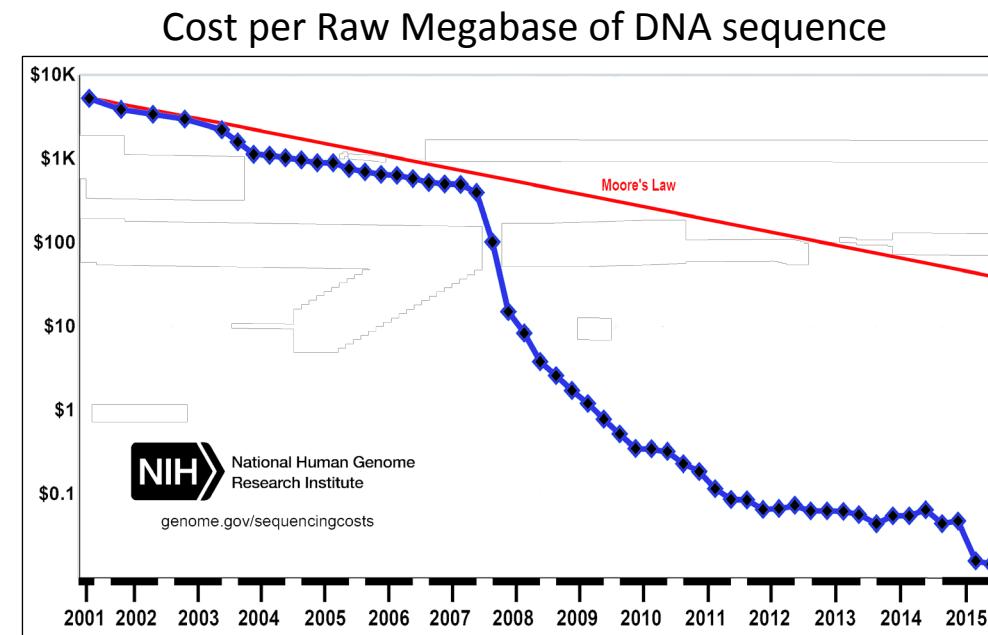
Next Generation Sequencing (NGS) technologies implement the concept of **cyclic-array sequencing**: Dense arrays of DNA are sequenced by repetitive cycles of enzymatic reaction and imaging-based data collection



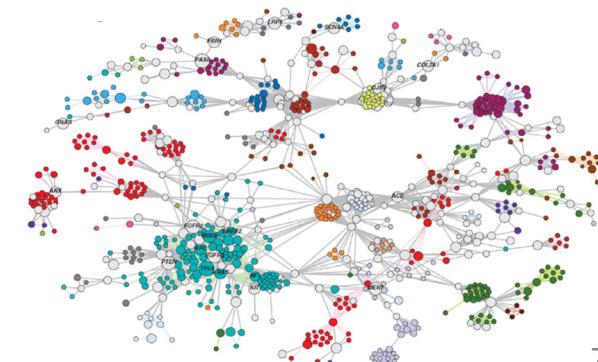


# NGS data: A twofold problem

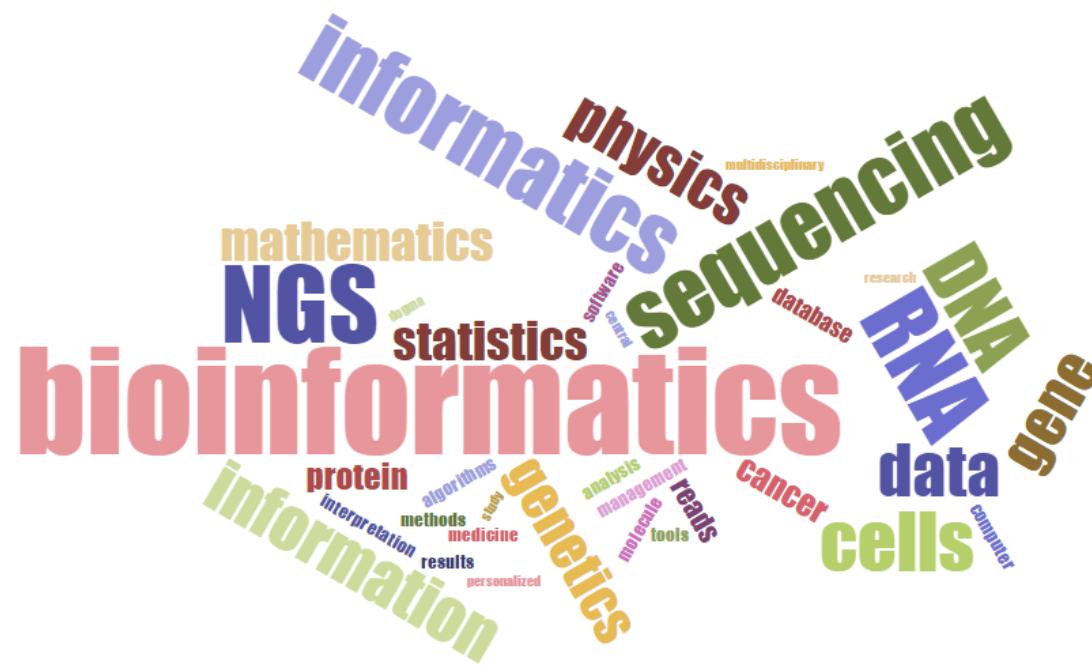
The **huge amounts** of NGS data produced worldwide need to be stored, organized and shared



The **analysis** of NGS data deals with the **complexity** of both structure and function of biological organisms

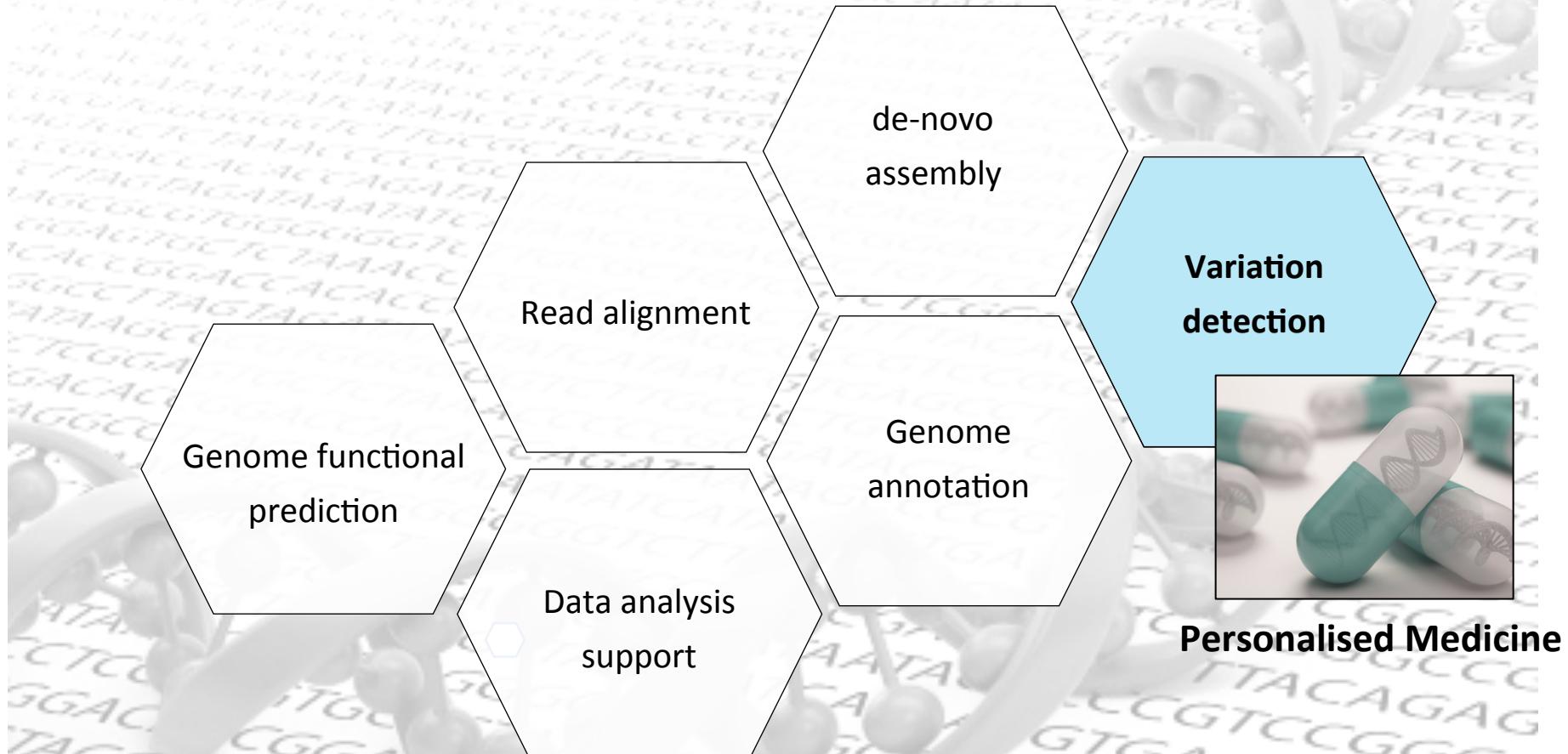


# How to deal with NGS data?



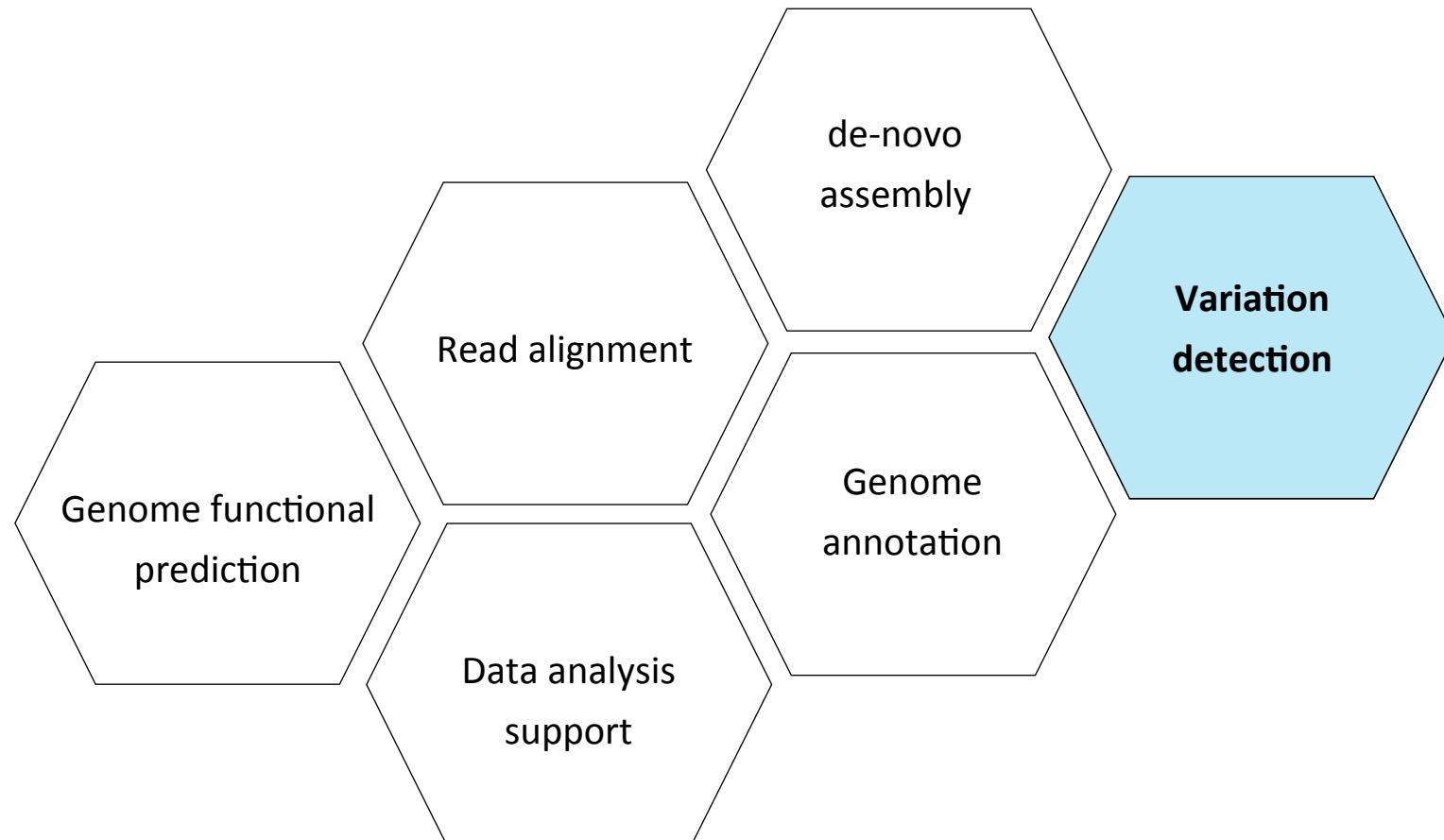
The term **Bioinformatics** refers to the application of **computational techniques** to the **understanding** and **organisation** of the information associated to biological macromolecules on a large scale and comprises research fields at the interface between computer and biological sciences

# Bioinformatics revolutionised genomics

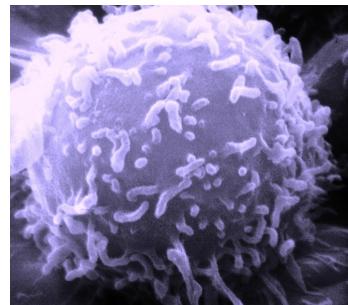


Zhang, Jun, et al. "The impact of next-generation sequencing on genomics." Journal of genetics and genomics 38.3 (2011): 95-109

# My contribution in this field

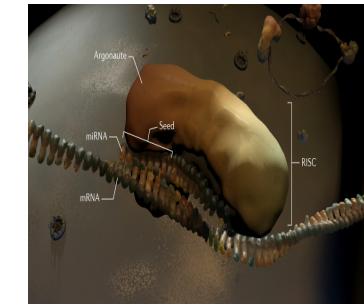


# Which variations have I investigated?

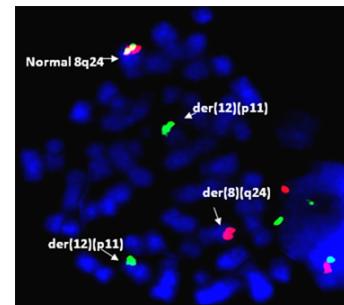


B cell clonal populations

Variation  
detection

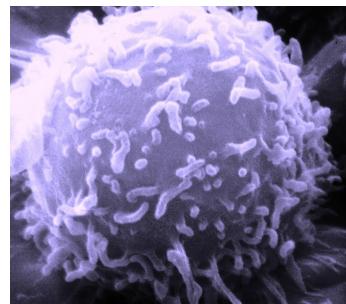


miRNA/isomiR patterns

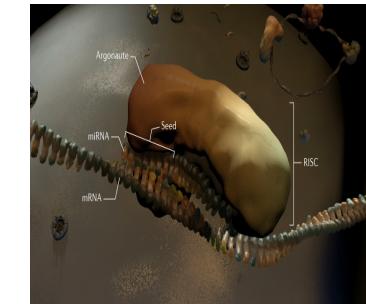
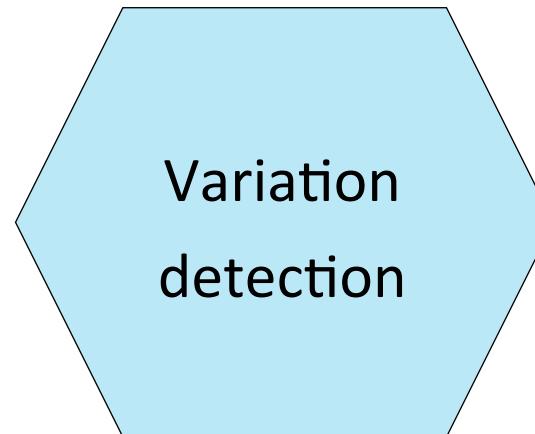


Gene Fusions

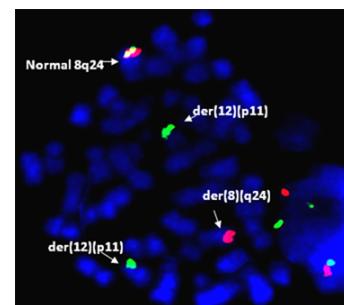
# How did I investigate these variations?



VDJSeq-Solver tool



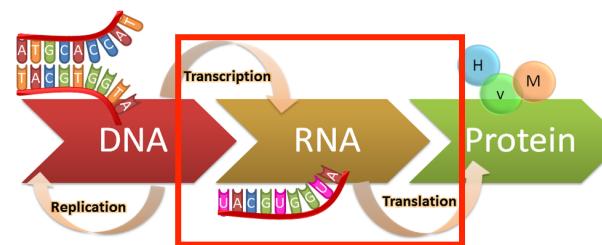
isomiR-SEA

UNIVERSITÀ  
di VERONA

FuGePrior tool



# Which features do these tools have in common?



- ✓ RNA-Seq data is used as input
- ✓ Propose novel **in silico** approaches for **disease** characterization
- ✓ Designed considering latest **biological knowledge**
- ✓ Built-on-top of state-of-the-art **bioinformatics methods** and tools
- ✓ Results **confirm** previous studies or **introduce** novel genomic knowledge

Introduction

Objective

**Methodologies&Tools**

Conclusions

**VDJSeq-Solver**

isomiR-SEA

FuGePrior

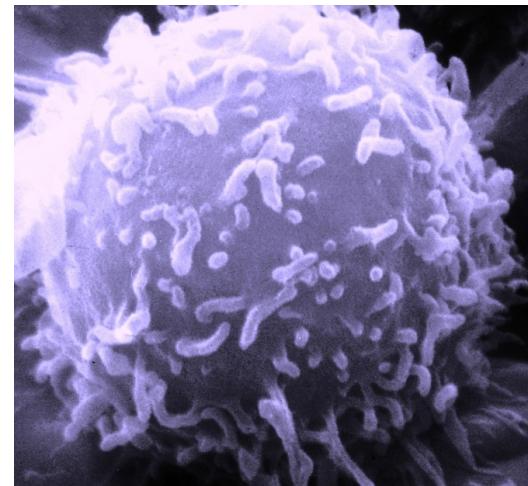


**POLITECNICO  
DI TORINO**

# VDJSeq-Solver

Biological Background

## The B cells

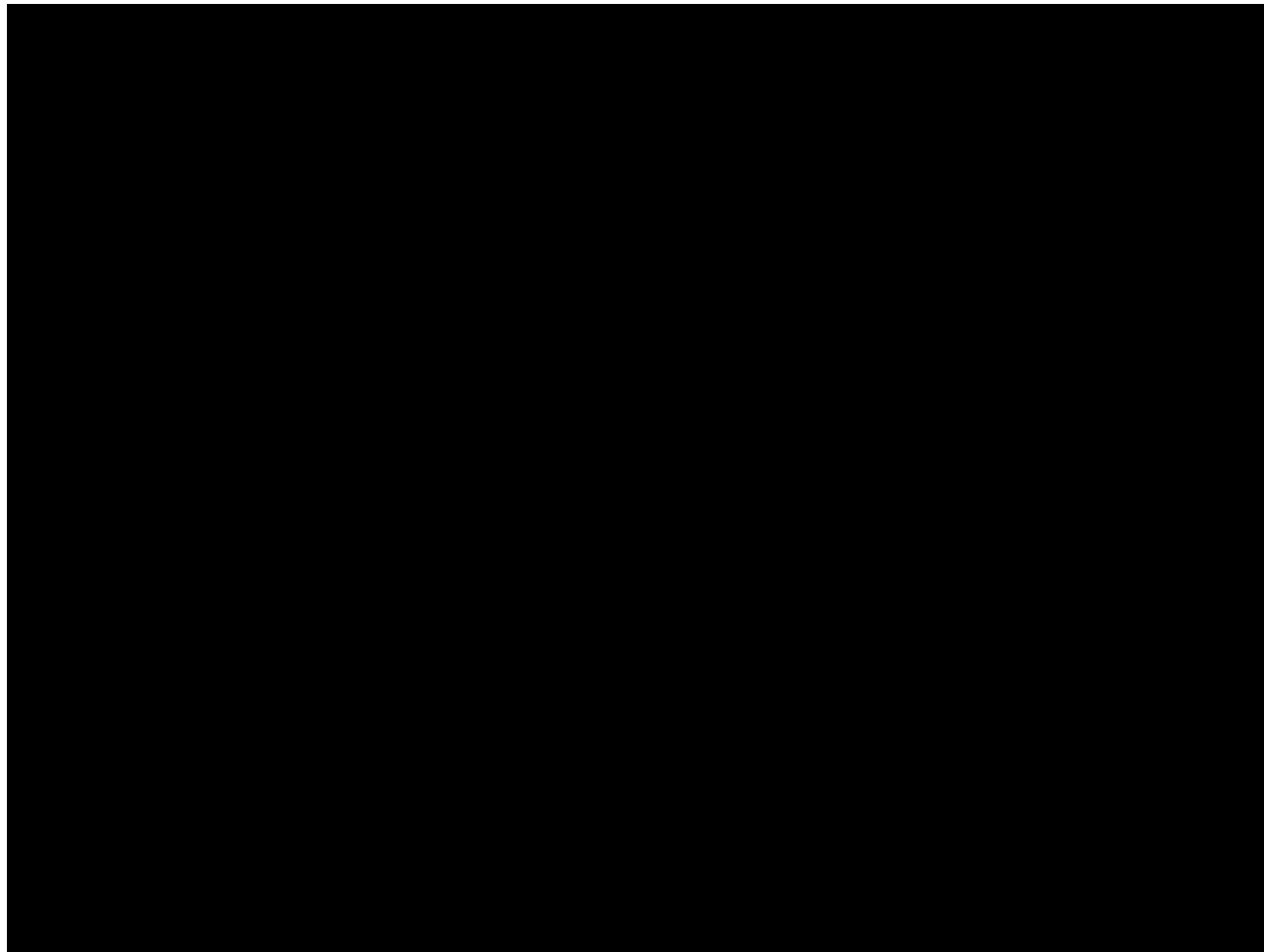


- ✓ One of the major effector molecules involved in **Adaptive Immune Responses**
- ✓ Originate in the **bone marrow** and undergo maturation in the **lymphoid tissues**
- ✓ Produce **antibodies**
- ✓ More than  **$10^{16}$**  different types of **B cells** ensure protection against infectious diseases



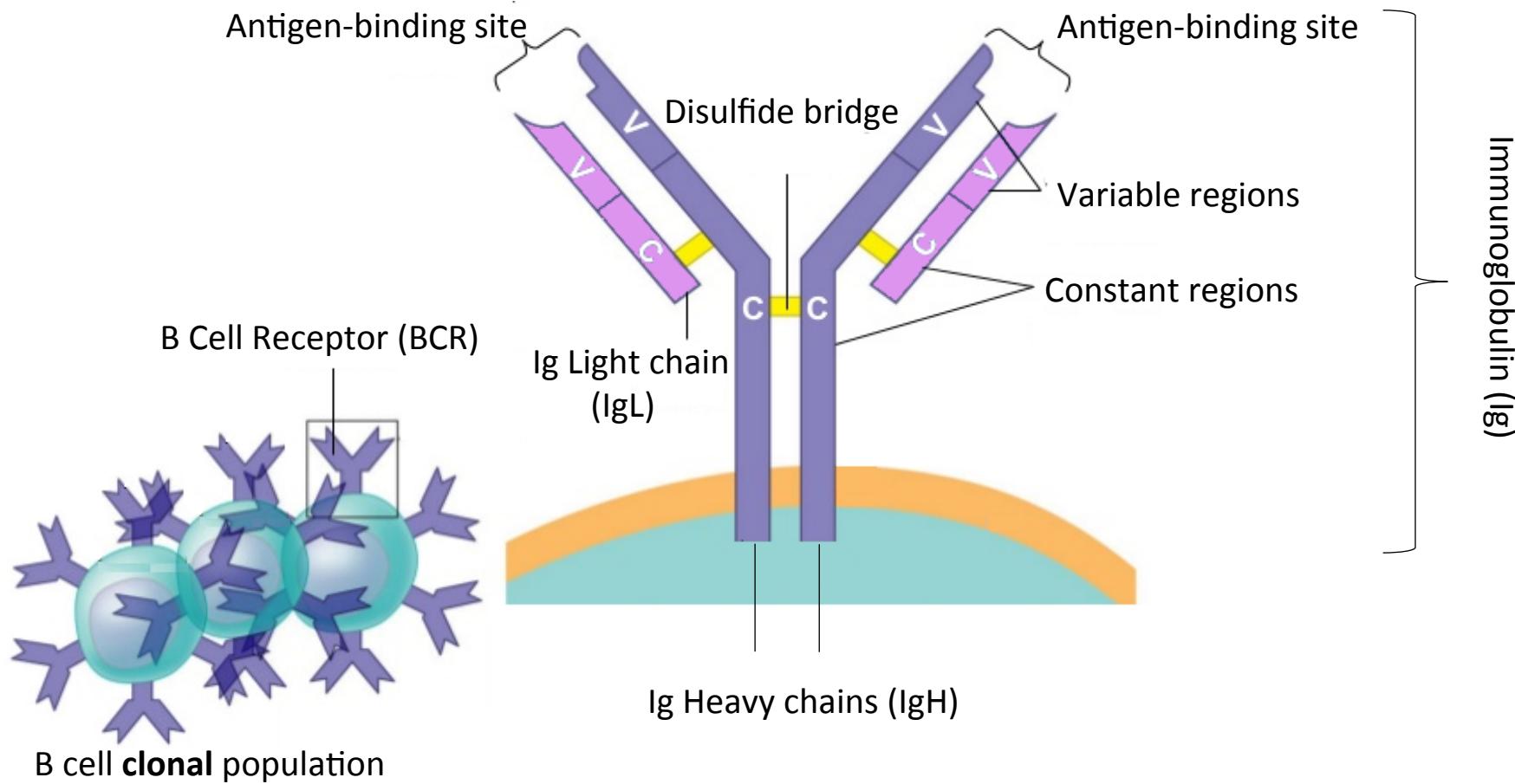
Biological Background

# How do B cells work?



Biological Background

# How can B cells recognize different antigens?



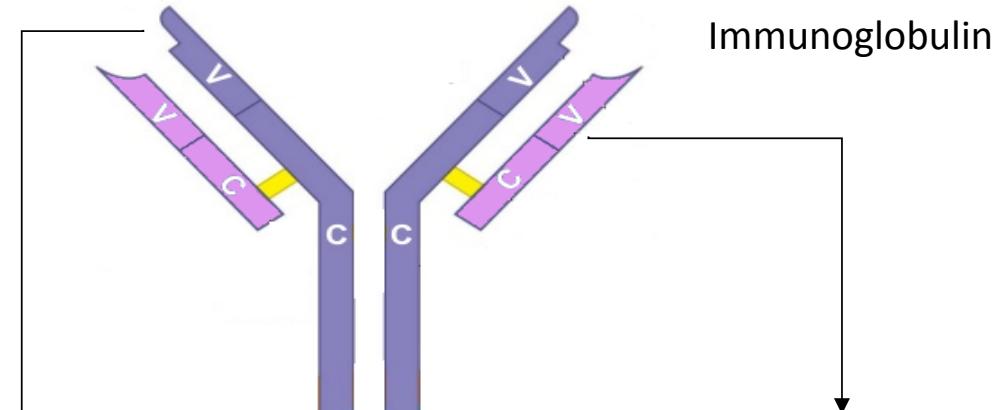
B cell clonal population

Different B Cell Receptors bind different antigens

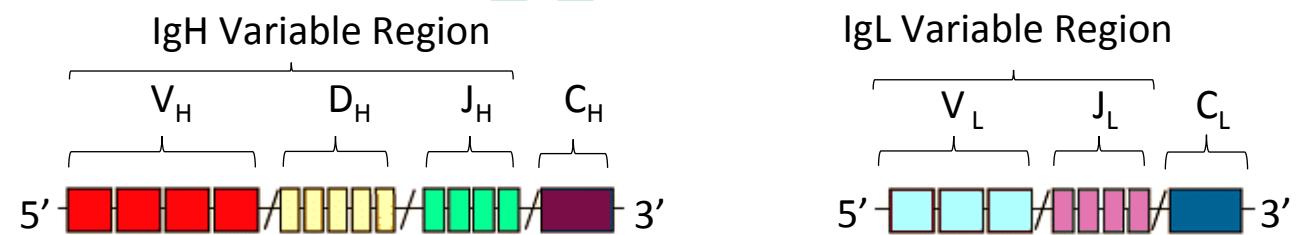
Biological Background

# V(D)J Recombination accounts for BCR diversity

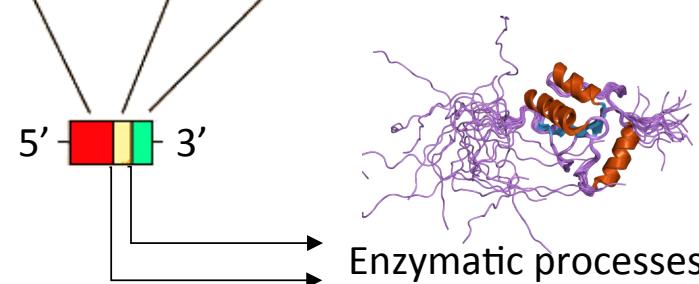
Protein Level



DNA Level



RNA Level



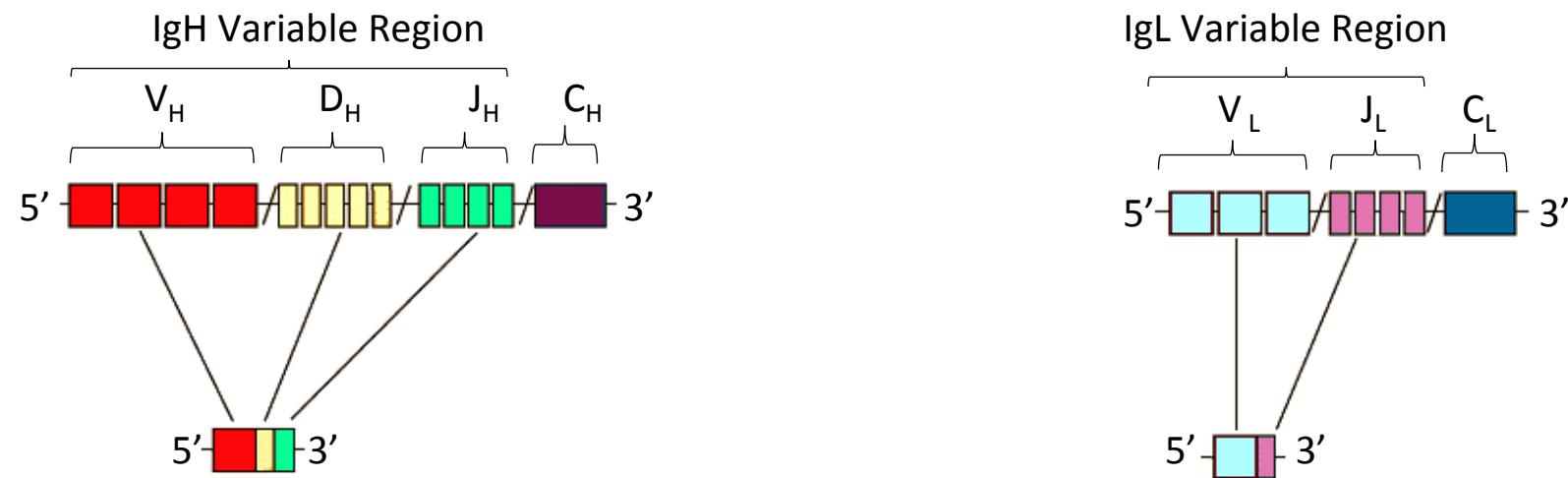
Enzymatic processes



Motivation

# Why is B cell clone identification important?

- ✓ Several B cell diseases are characterized by **massive expansion** of a single B cell clone
- ✓ B cell clonal expansion is generally interpreted by clinicians as **evidence of cancer**
- ✓ There are **biases** in IgH and IgL **recombined genes** in different pathologies and **correlations** between the clinical course of the disease and Ig gene usage

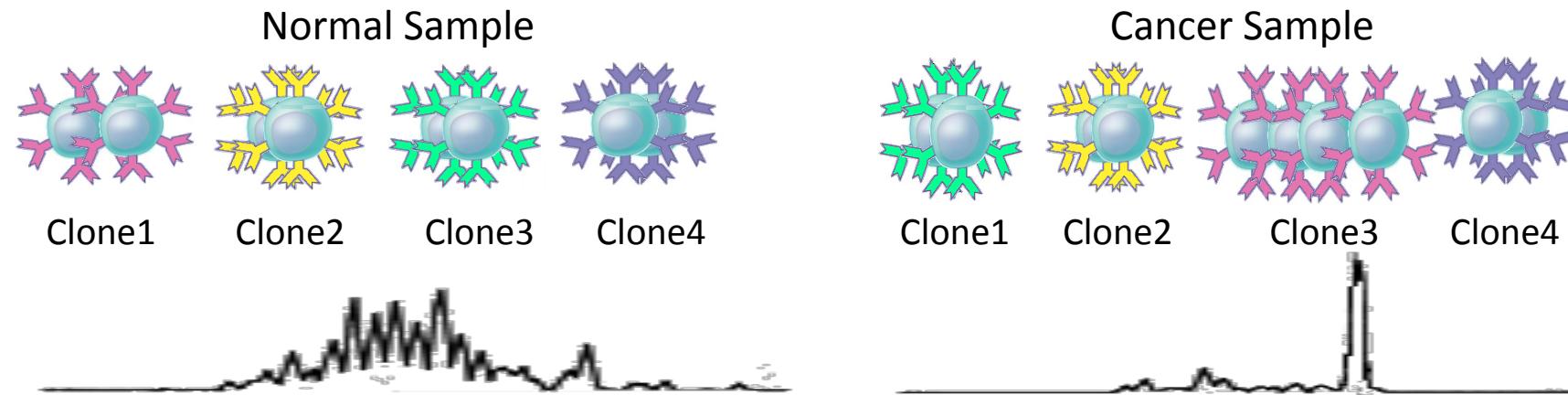


Capello, D., et al. "Evidence of biased immunoglobulin variable gene usage in highly stable B-cell chronic lymphocytic leukemia." Leukemia 18.12 (2004): 1941-1947



## Motivation

# B cell clone identification approaches

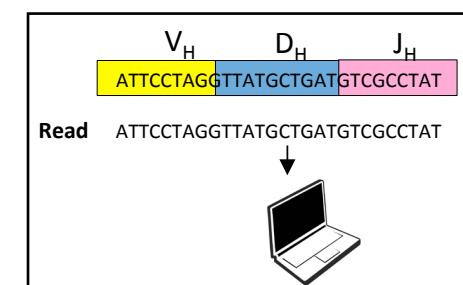


PCR Clonality tests

- ✓ **Limited sensitivity** associated to the normal polyclonal background

NGS analyses

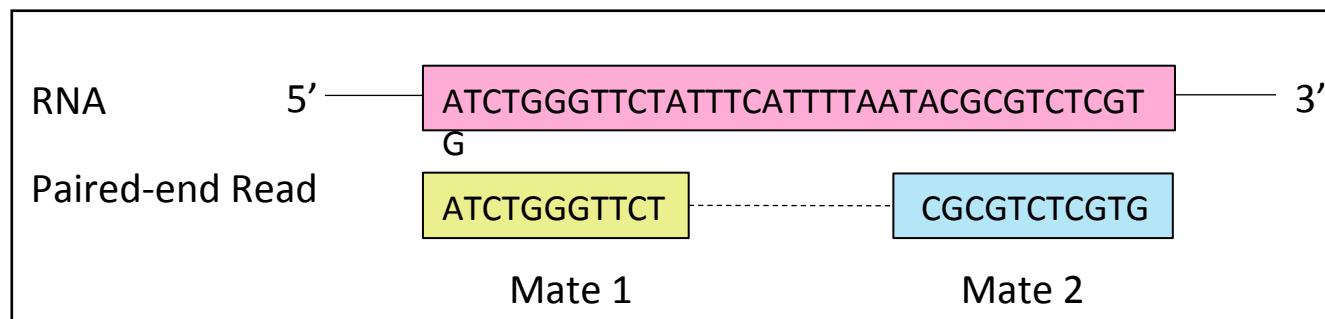
- ✓ Rely on the usage of **Ig amplified DNA/RNA molecules**



Gazzola, Anna, et al. "The evolution of clonality testing in the diagnosis and monitoring of hematological malignancies." Therapeutic advances in hematology 5.2 (2014): 35-47.

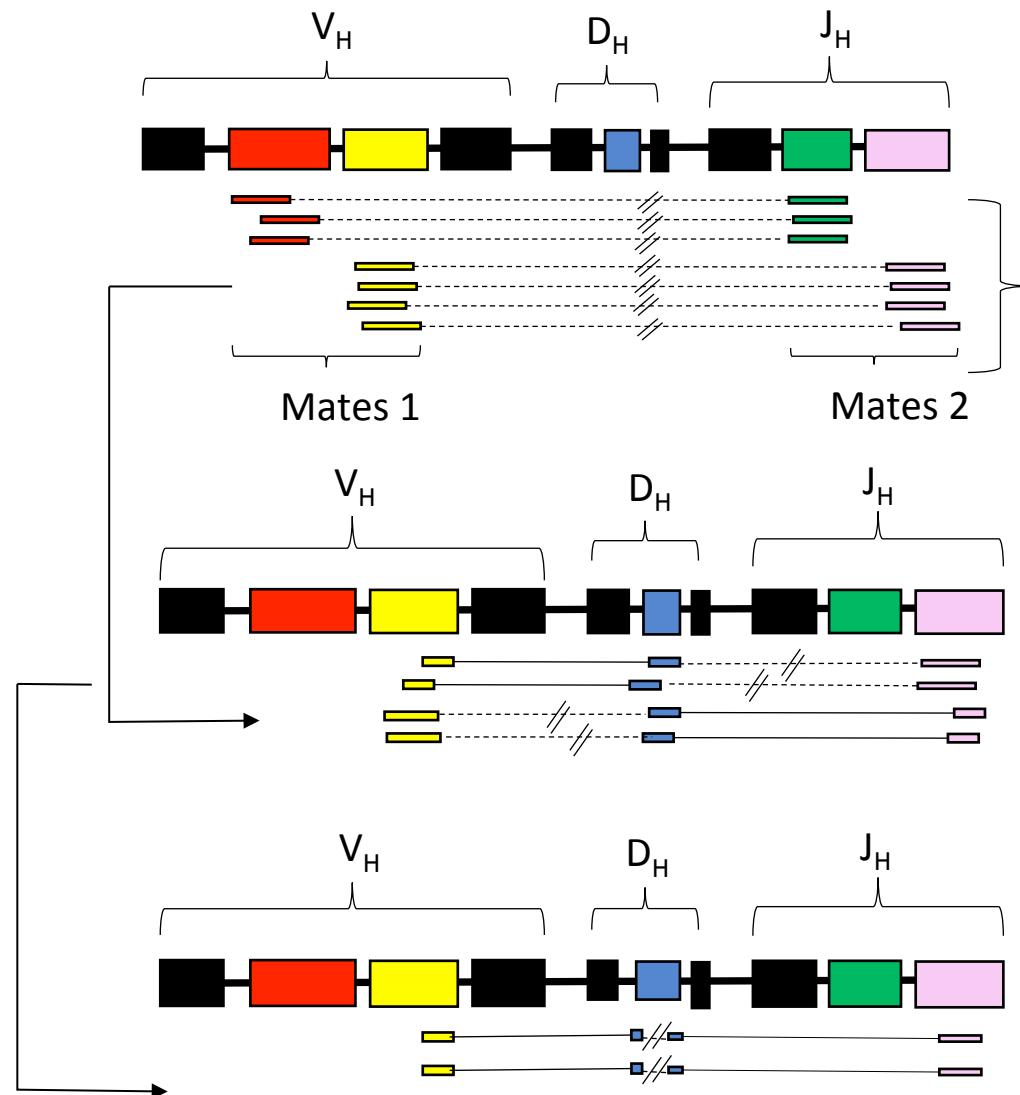
Motivation

# VDJSeq-Solver



Algorithm

# B cell clone identification



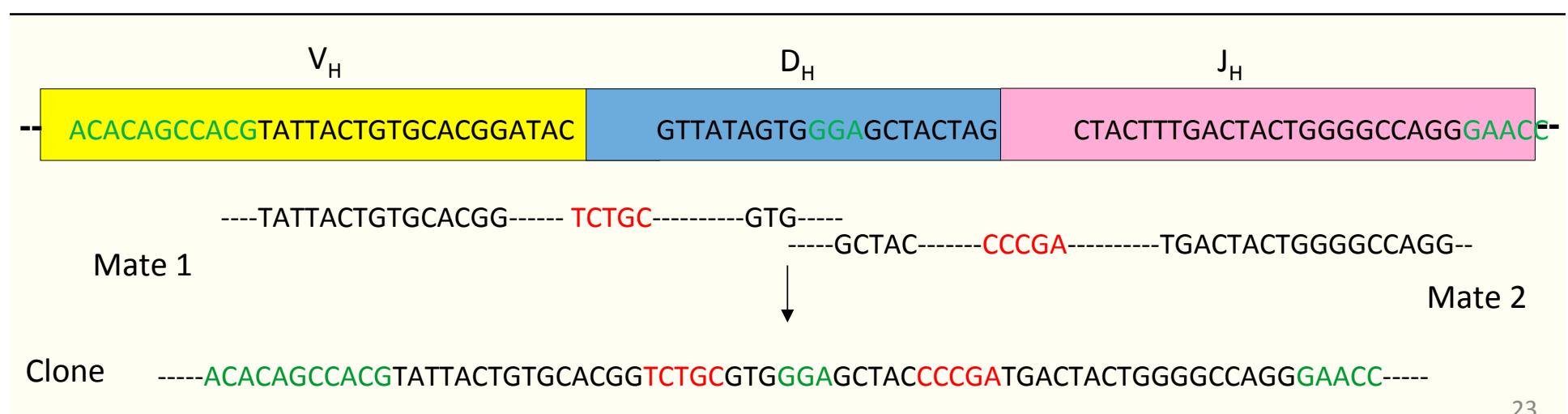
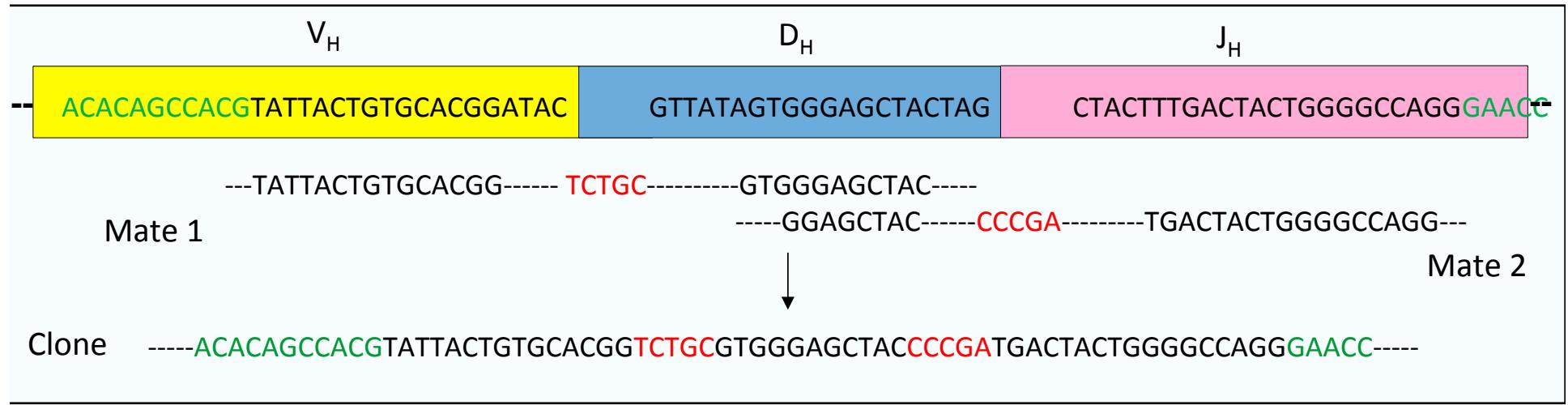
VJ Encompassing reads allow to identify VJ clones in both IgHs and IgLs

VJ Encompassing reads having **one or both** mates partially mapped onto a D gene identify IgH V(D)J clones

VJ Encompassing reads having **both** mates mapped onto a D gene allow to reconstruct IgH V(D)J clone sequences

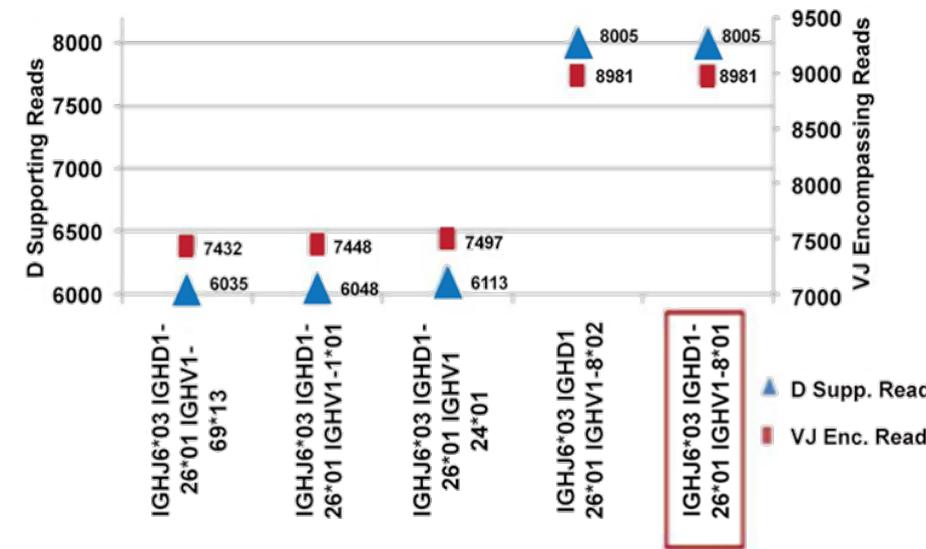
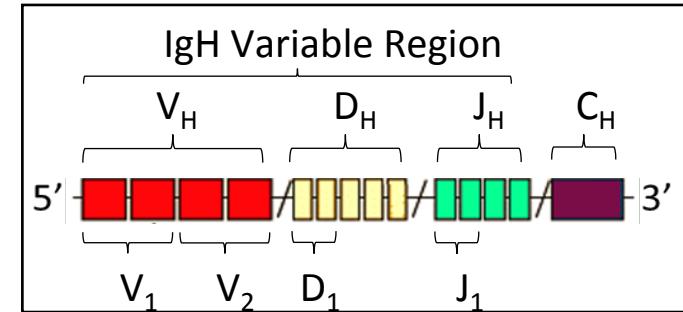
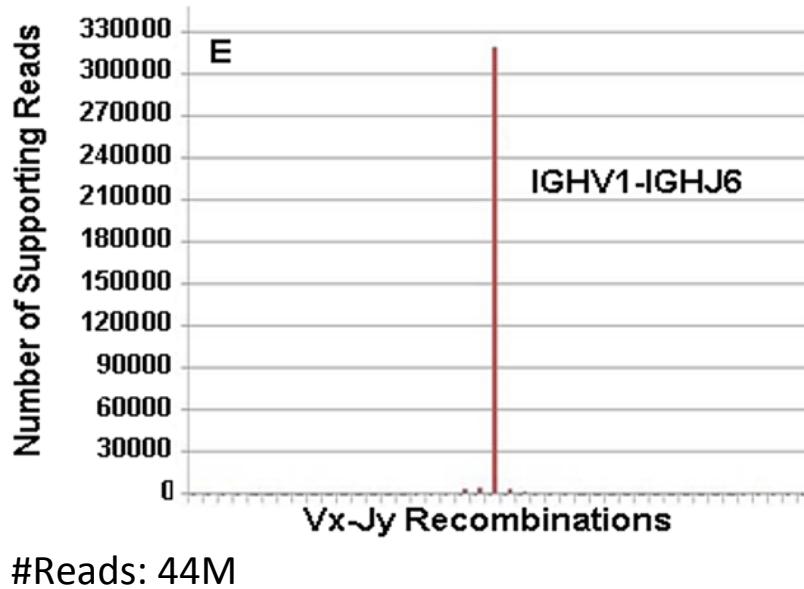
## Algorithm

# B cell clone sequence reconstruction



Results

# IgH main clone identification in MCL samples



VDJSeq-Solver correctly identified IgH main clone in MCL and DLBCL samples



## Results

# IgH main clone sequence in MCL samples

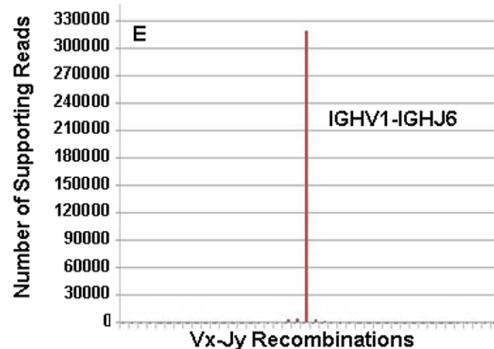
VDJSeq-Solver IgH main clone sequences are **highly similar** to those obtained via wet-lab experiments, maximum 1.56% error

The neglectable impact of the identified mismatches is confirmed by the analysis performed adopting **state-of-the-art tools** for Ig gene assignment

Tool	$V_H$		$D_H$		$J_H$	
	Sanger	VDJSeq-Solver	Sanger	VDJSeq-Solver	Sanger	VDJSeq-Solver
IMGT/V-QUEST	IGHV1-8*01	IGHV1-8*01	IGHD1-26*01	IGHD1-26*01	IGHJ6*03	IGHJ6*03
JOINSOLVER	IGHV1-8*01	IGHV1-8*01	IGHD1-26*01	IGHD1-26*01	IGHJ6*03	IGHJ6*03
VDJsolver	IGHV1-8*01	IGHV1-8*01	IGHD1-26*01	IGHD1-26*01	IGHJ6*03	IGHJ6*03
SoDA	IGHV1-8*01	IGHV1-8*01	IGHD1-26*01	IGHD1-26*01	IGHJ6*03	IGHJ6*03
iHMMune-align	IGHV1-8*01	IGHV1-8*01	IGHD1-26*01	IGHD1-26*01	IGHJ6*03	IGHJ6*03

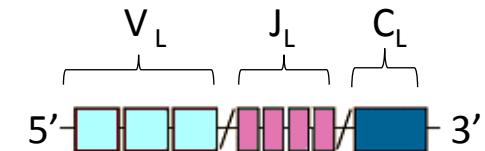
Remarks and Perspectives

# To summarise



VDJSeq-Solver **identified** the IgH main clonal population in MCL and DLBCL samples and **reconstructed** its sequence in MCL samples

Preliminary results on **IgLs** in MCL samples are promising



**Future work:**

- { T Cell Receptor (TCR) analysis
- Clone Collapsing
- Parallelisation

Introduction

Objective

**Methodologies&Tools**

Conclusions

VDJSeq-Solver

**isomiR-SEA**

FuGePrior

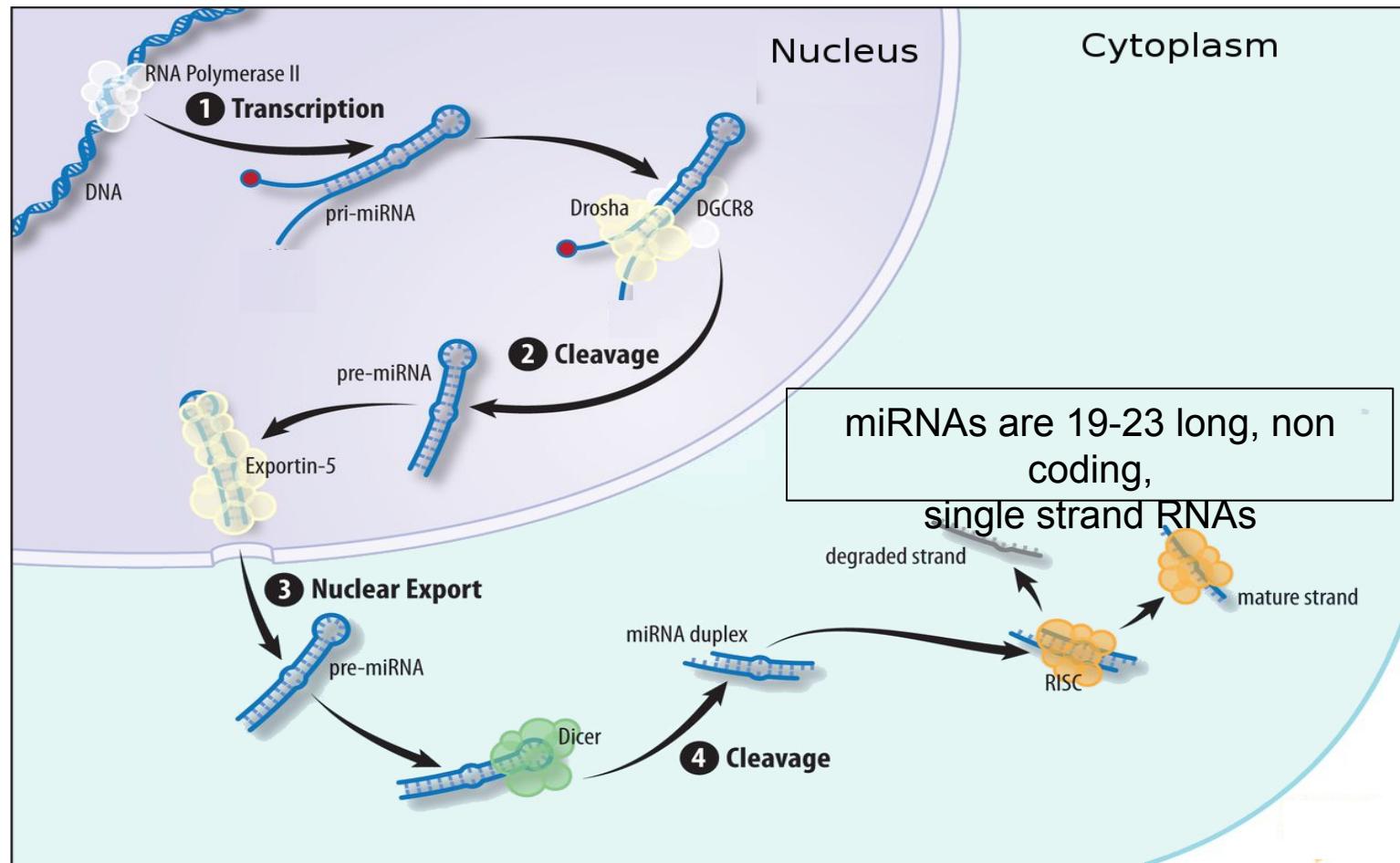


POLITECNICO  
DI TORINO

# isomiR-SEA

## Biological Background

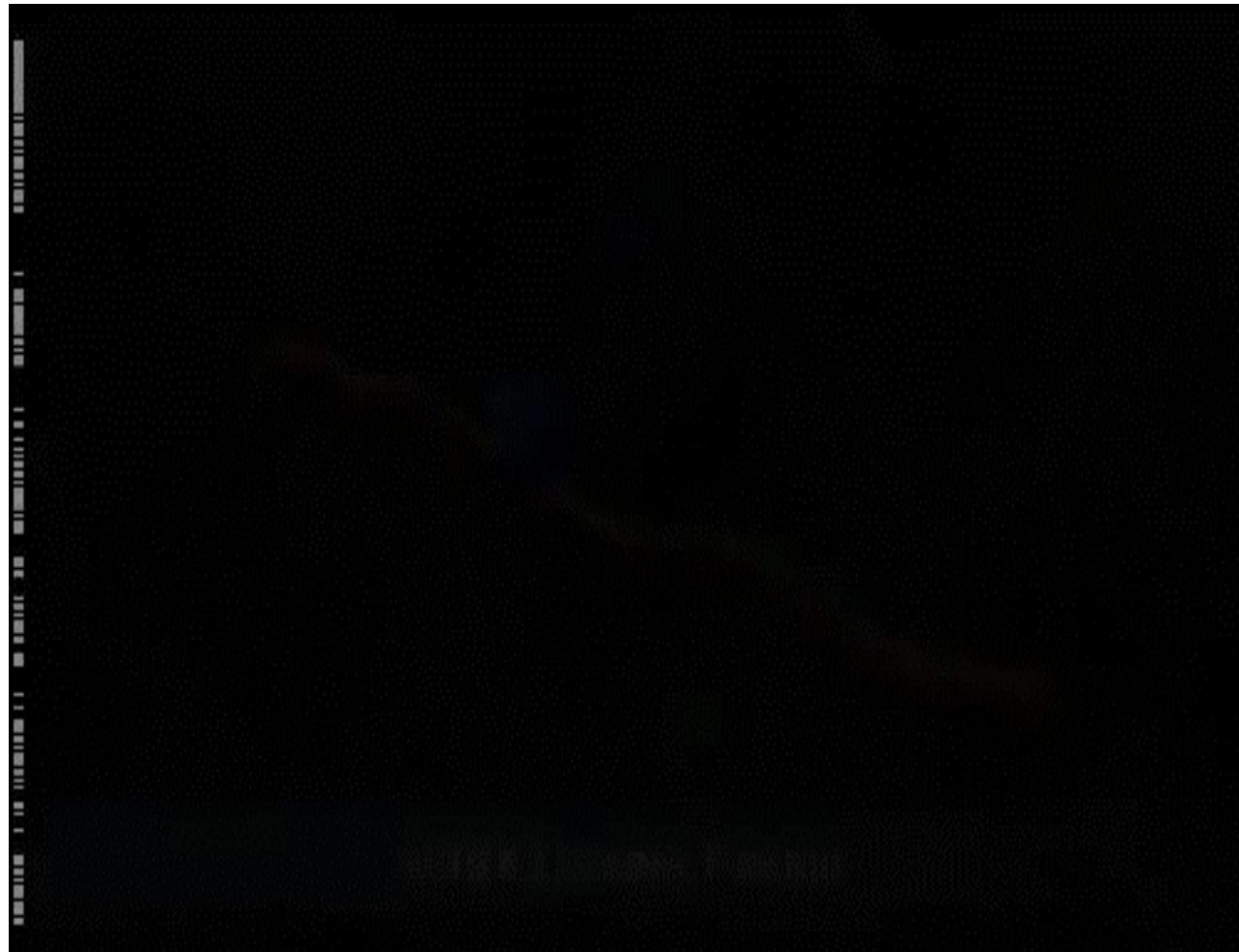
# The miRNAs





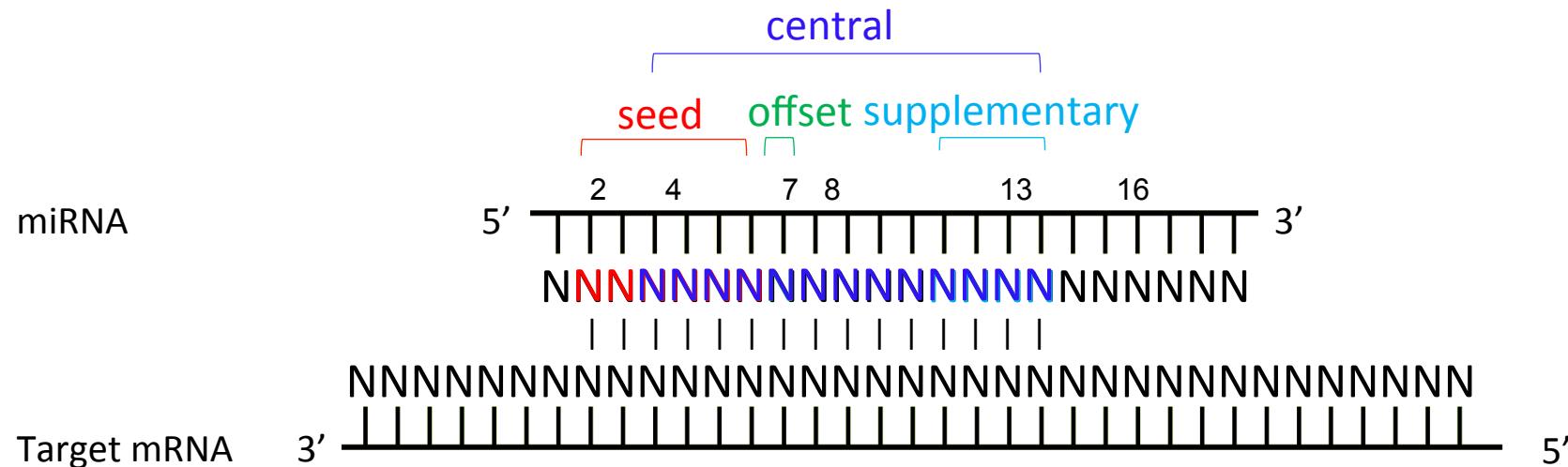
Biological Background

# How do miRNAs work?



Biological Background

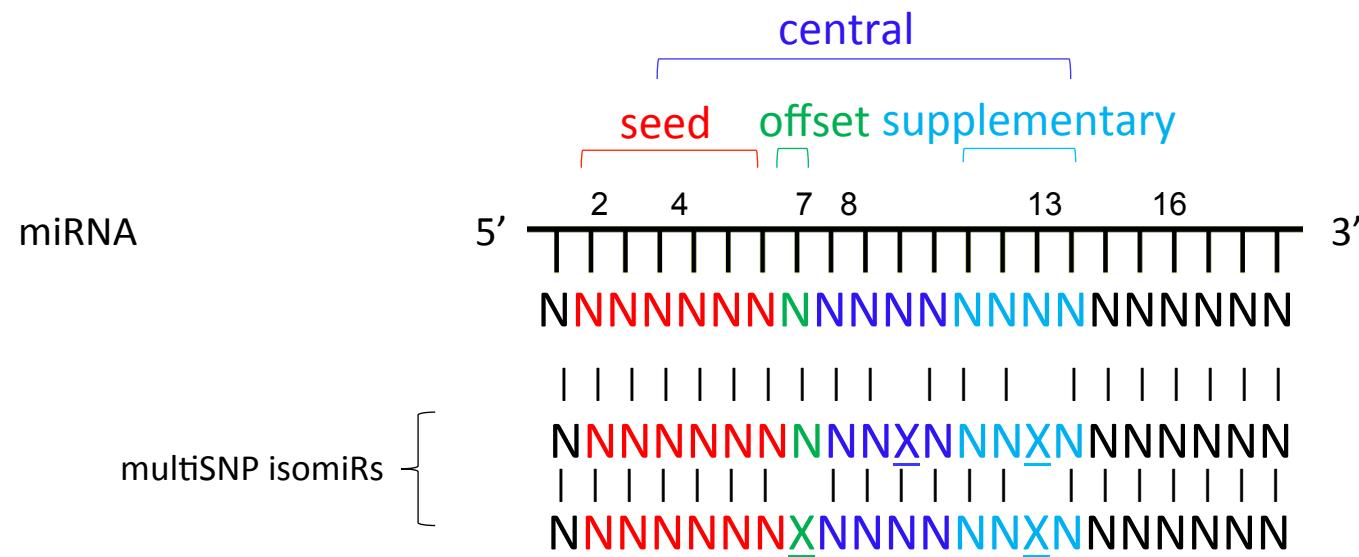
# miRNA:mRNA interaction



miRNAs belonging to the same **family** have identical **seed** sequences

## Biological Background

# miRNA variants: The isomiRs



Motivation

# Why is miRNA/isomiR identification important?

Metastatic  
miRNAs

Play a role in the  
Epithelial to Mesenchymal  
Transition (EMT)

Examples:  
miR-7, miR-22, miR-30

Oncogenic  
miRNAs

Promote tumor  
development by inhibiting  
tumor suppressor genes

Examples:  
miR-17, miR-21, miR-155

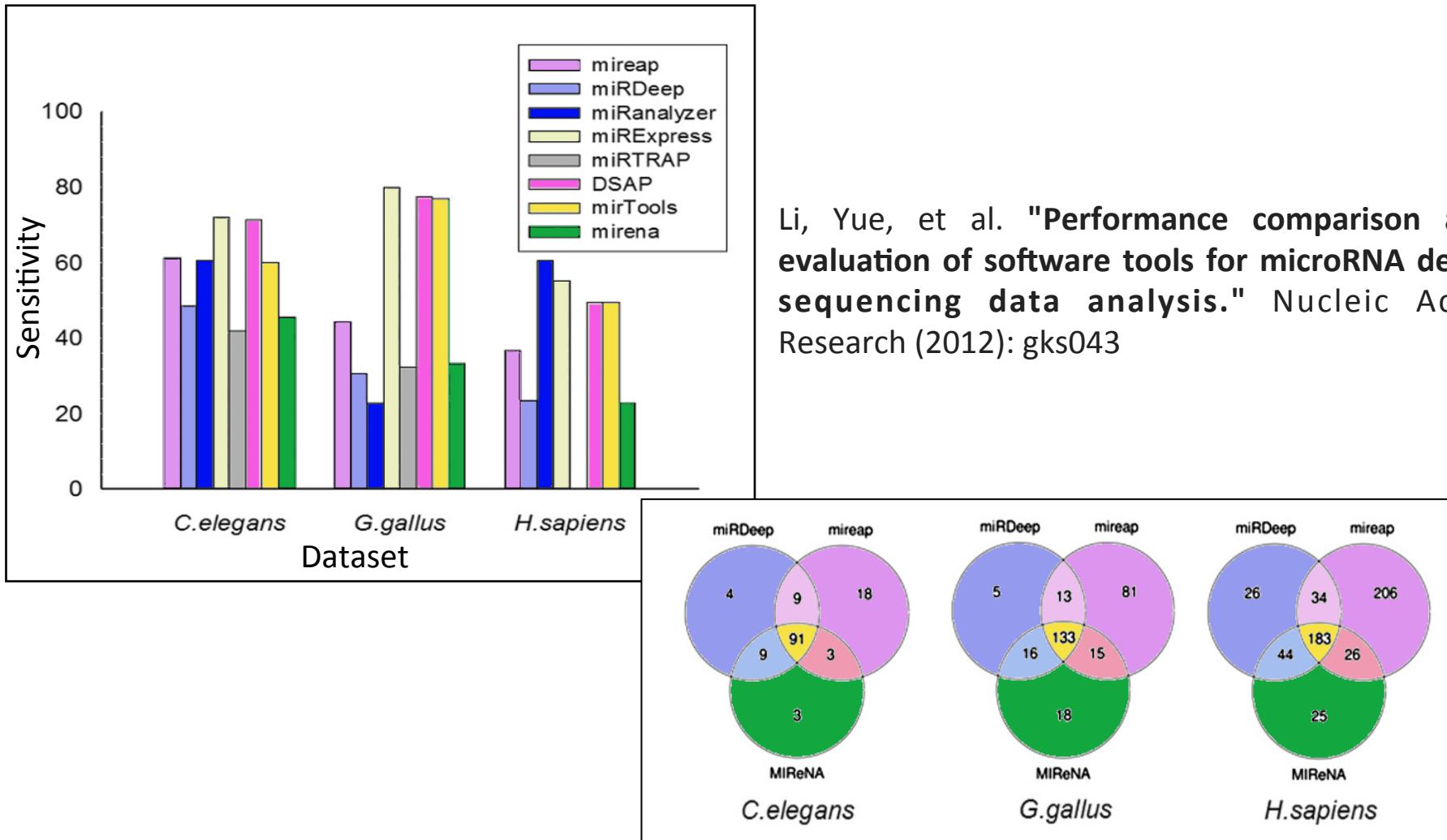
Tumor  
suppressor  
miRNAs

Prevent tumor  
development by inhibiting  
oncogenes

Examples:  
miR-145, miR-200, miR-205

## Motivation

# State-of-the-art miRNA detection tools



Li, Yue, et al. "Performance comparison and evaluation of software tools for microRNA deep-sequencing data analysis." Nucleic Acids Research (2012): gks043

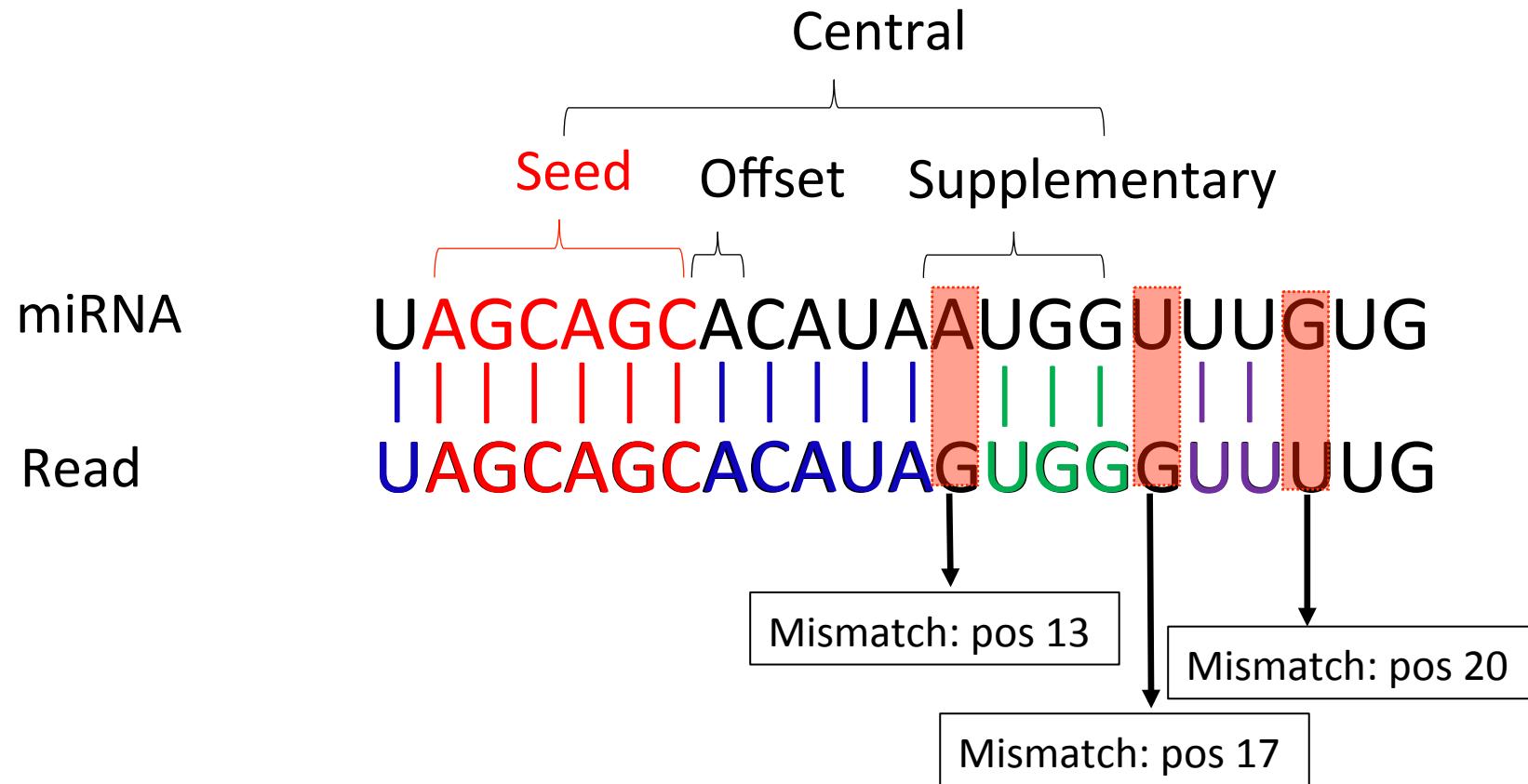
Motivation

# isomiR-SEA

- ✓ Implements a **miRNA-specific** alignment procedure
- ✓ Identifies **miRNA**, **isomiR** and **interaction site** expression profiles
- ✓ Uses **Seqan** library
- ✓ **Freely available** → <http://eda.polito.it/isomir-sea/>
- ✓ Users can easily **configure** the run

Algorithm

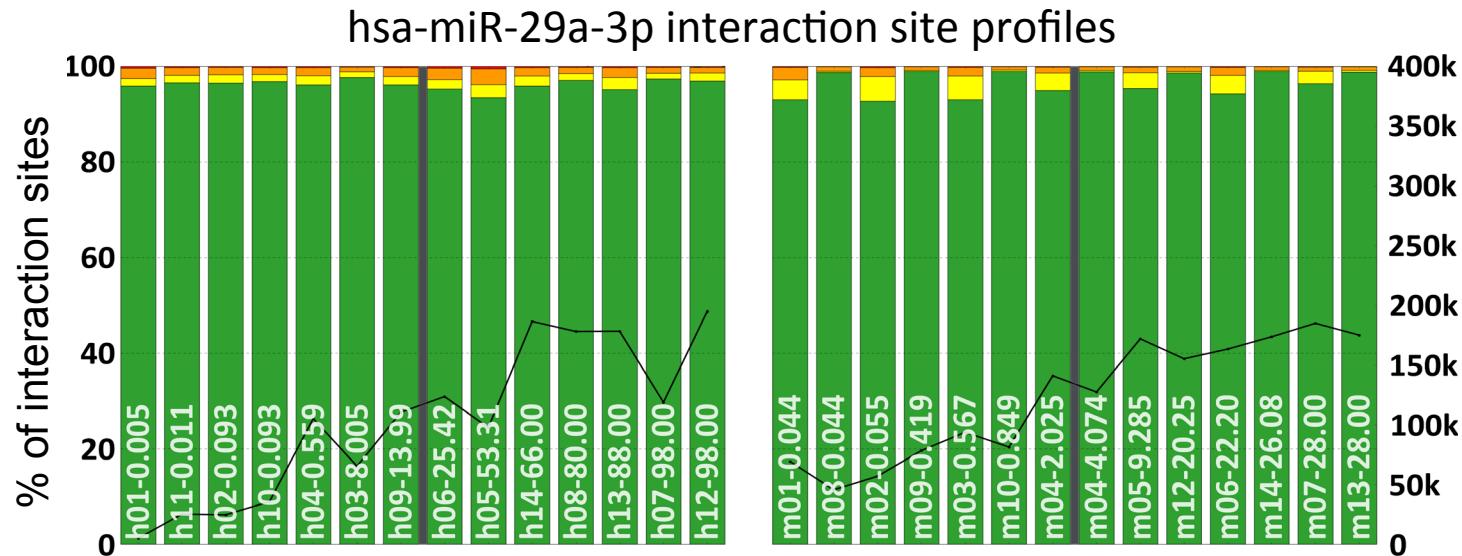
# The alignment procedure



By evaluating the **positions** of the encountered **mismatches**, isomiR-SEA can distinguish among miRNAs/isomiRs and assess the conservation of miRNA interaction sites

## Results

# miRNA/isomiR profiles in Somel dataset

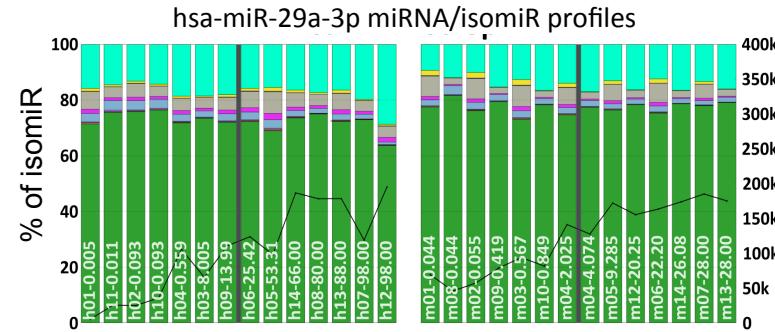


Avg#Reads: 5M

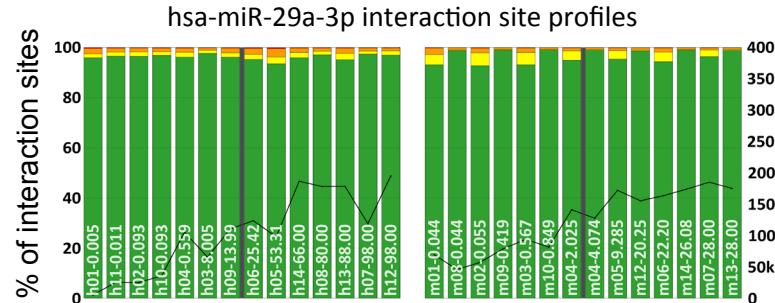
Col.	Interaction site Type	Conserved nt		
		nt 8	nt 9-11	nt 12-16
green	offset-suppl-central	X	X	X
yellow	offset-suppl	X		X
orange	offset-only	X		
red	suppl-only			X

Remarks and Perspectives

# To summarise



isomiR-SEA identified **miRNA/isomiR and conserved interaction site** expression levels in normal and cancer samples



isomiR-SEA results allow to explain miRNA/isomiR targeting activity

**Future work:**

pre-miRNA evaluation

read quality value evaluation

long non-coding RNA identification

Introduction

VDJSeq-Solver

Objective

isomiR-SEA

**Methodologies&Tools**

**FuGePrior**

Conclusions

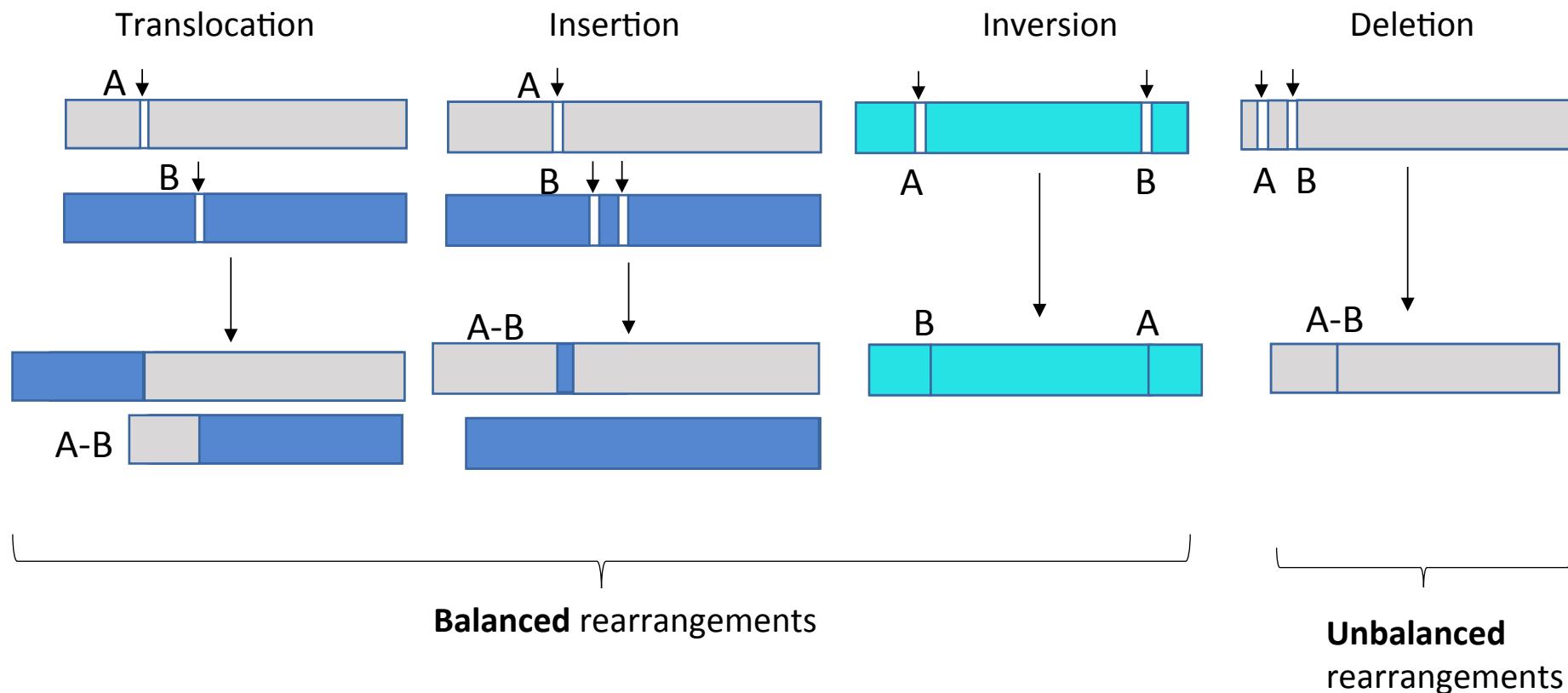


**POLITECNICO  
DI TORINO**

# FuGePrior

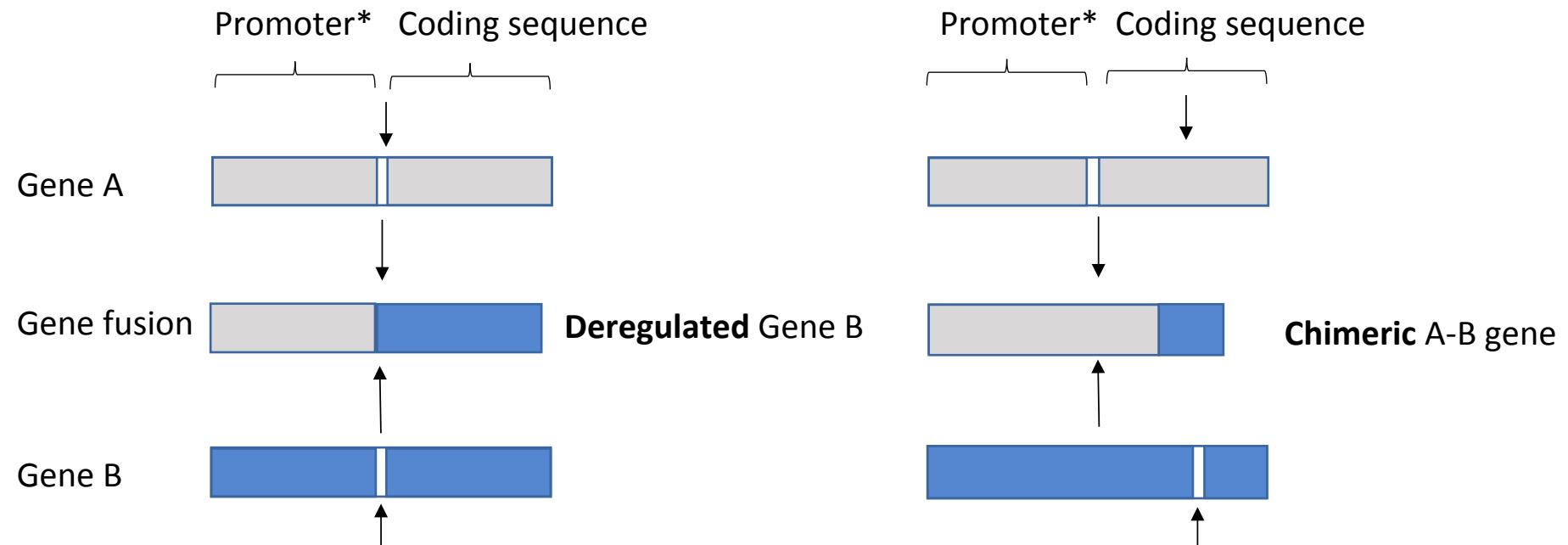
Biological Background

# The gene fusions



Biological Background

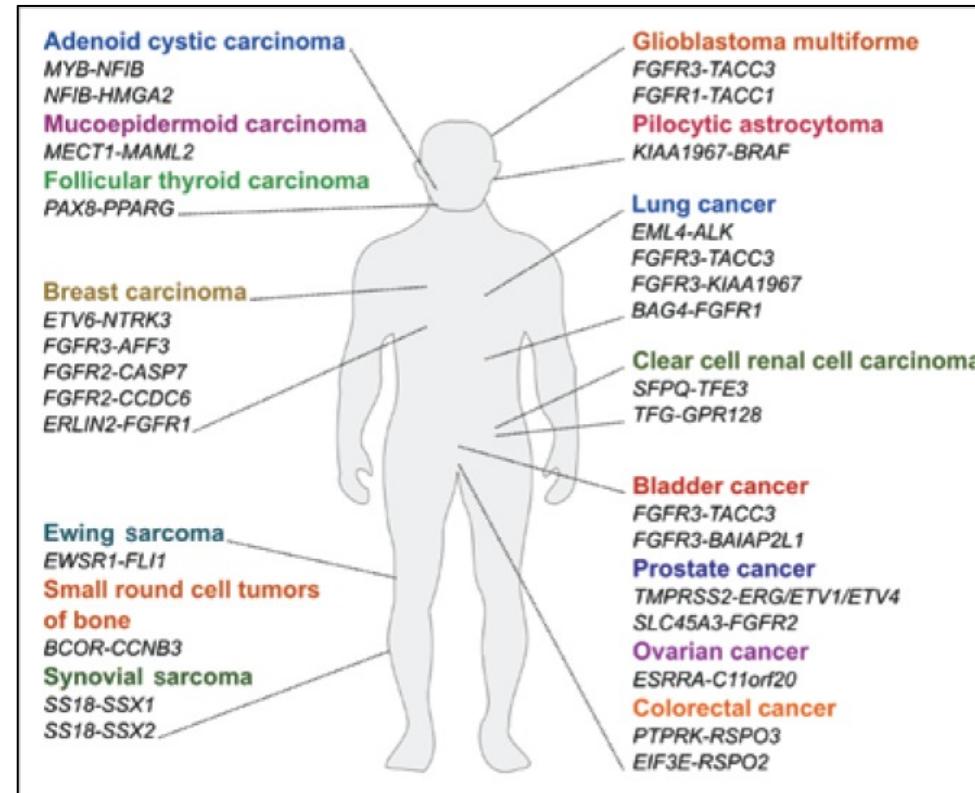
# Which is the effect of gene fusions?



\*Promoter: a region of DNA that initiates transcription of a particular gene

Motivation

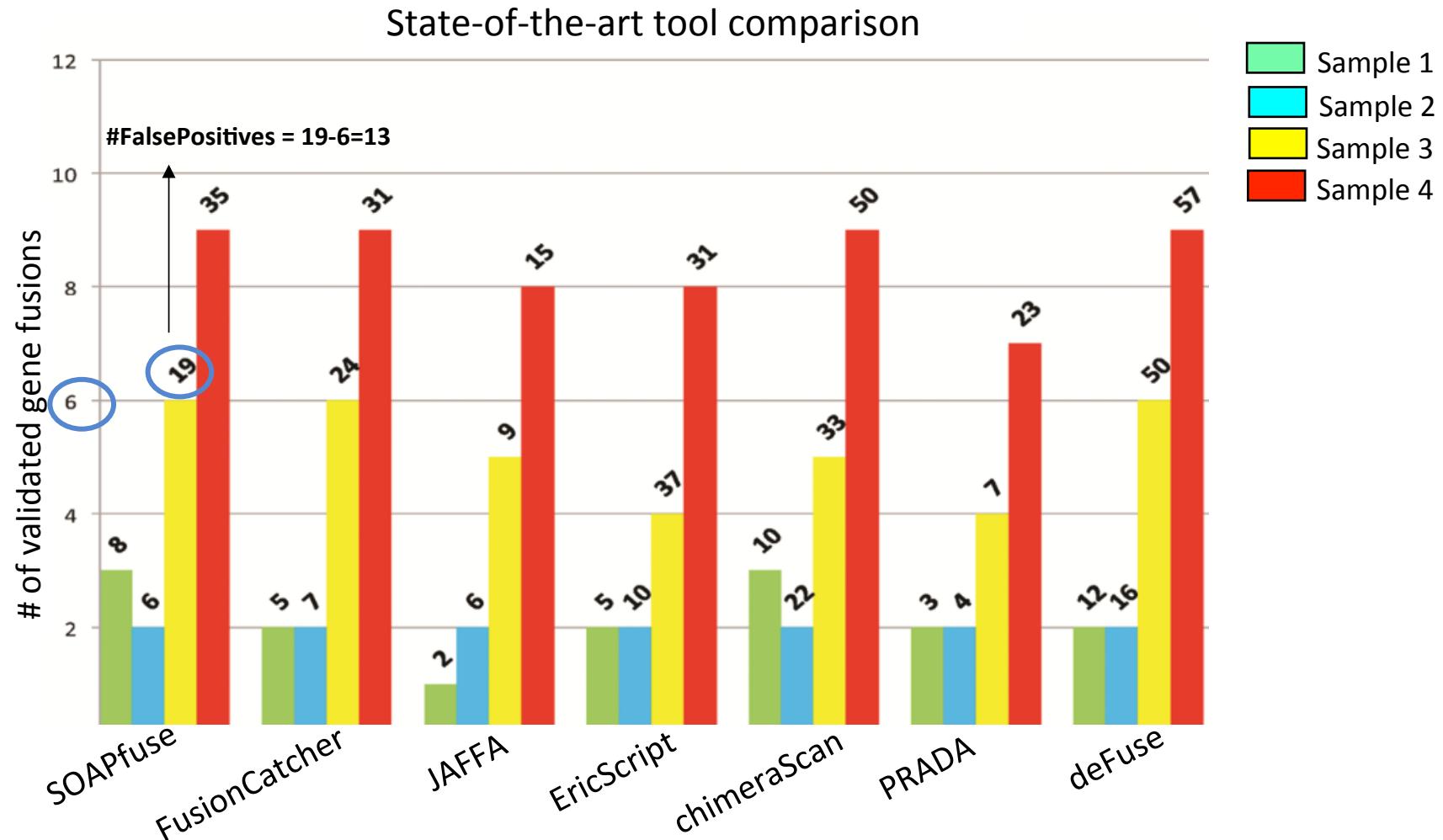
# Why is gene fusion identification important?



Gene fusions are **cancer-specific** and their identification is considered fundamental for the **diagnosis** and **prognosis** of several cancers

Motivation

# State-of-the-art gene fusion detection tools



Kumar, Shailesh, et al. "Comparative assessment of methods for the fusion transcripts detection from RNA-Seq data." *Scientific Reports* 6 (2016)



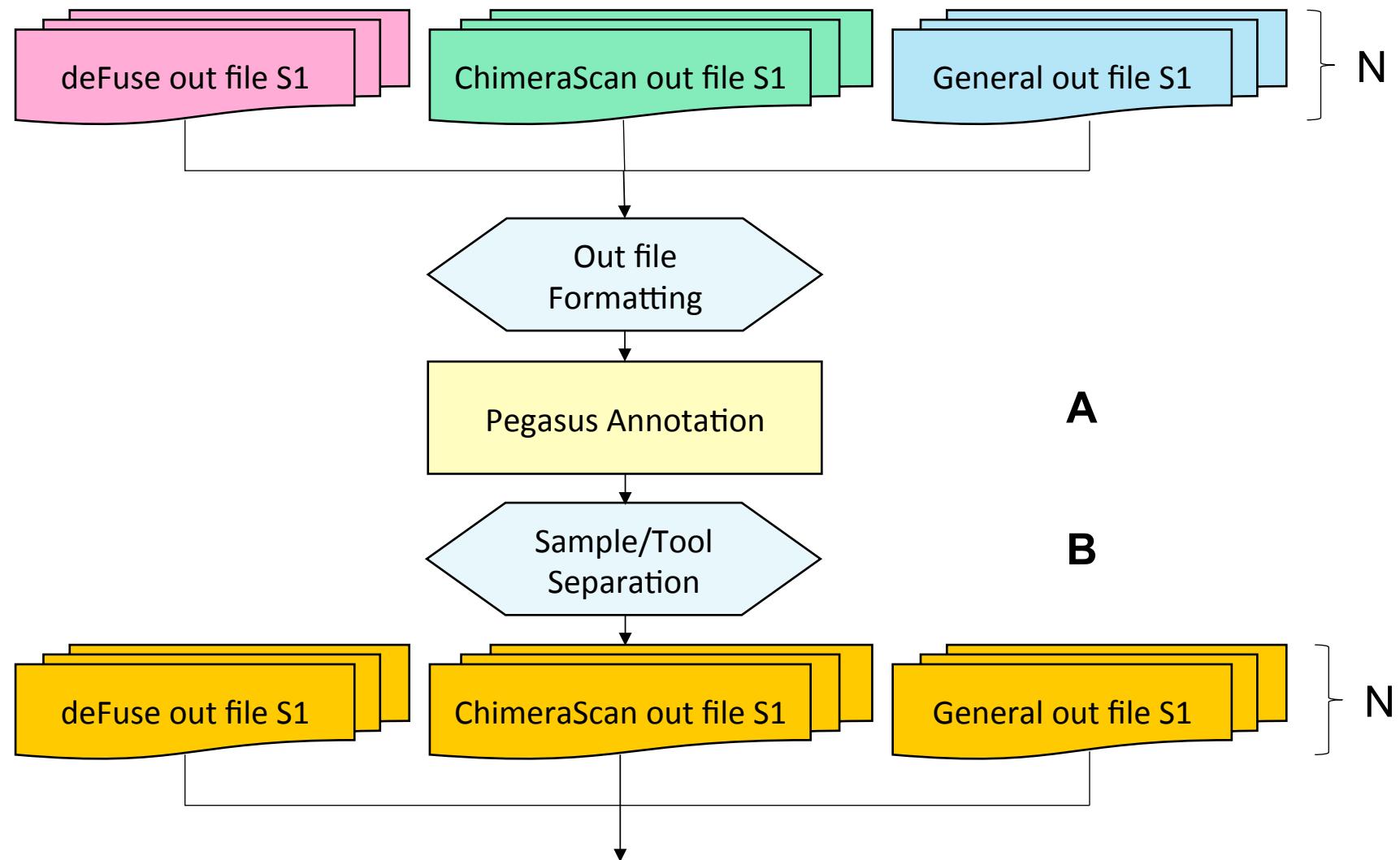
Motivation

## FuGePrior

- ✓ **Prioritizes** fusions from gene fusion discovery tools
- ✓ Implements a set of **filtering** and **processing** stages
- ✓ Integrates the driver scores from 2 **Machine Learning** algorithms
- ✓ **Freely available** → <https://philae.polito.it/paciello/FuGePrior/>
- ✓ Users can easily **configure** the run

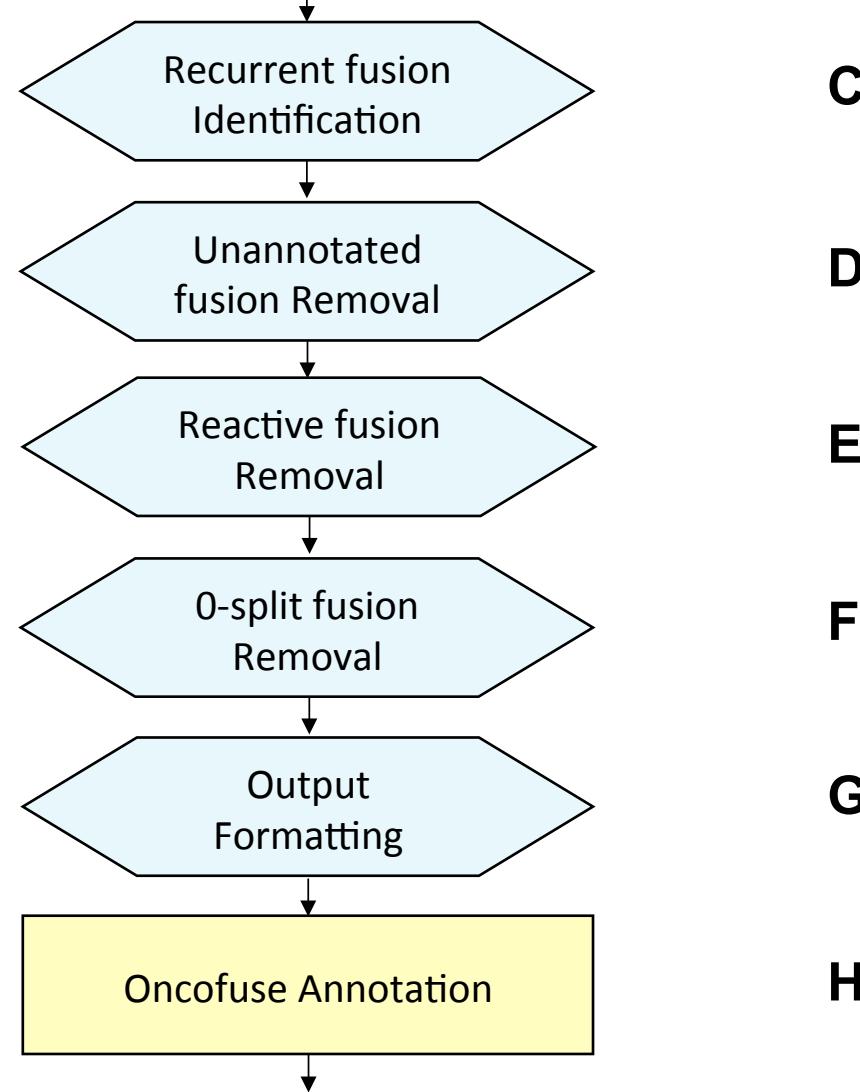
Algorithm

# How are fusions prioritized? (1)



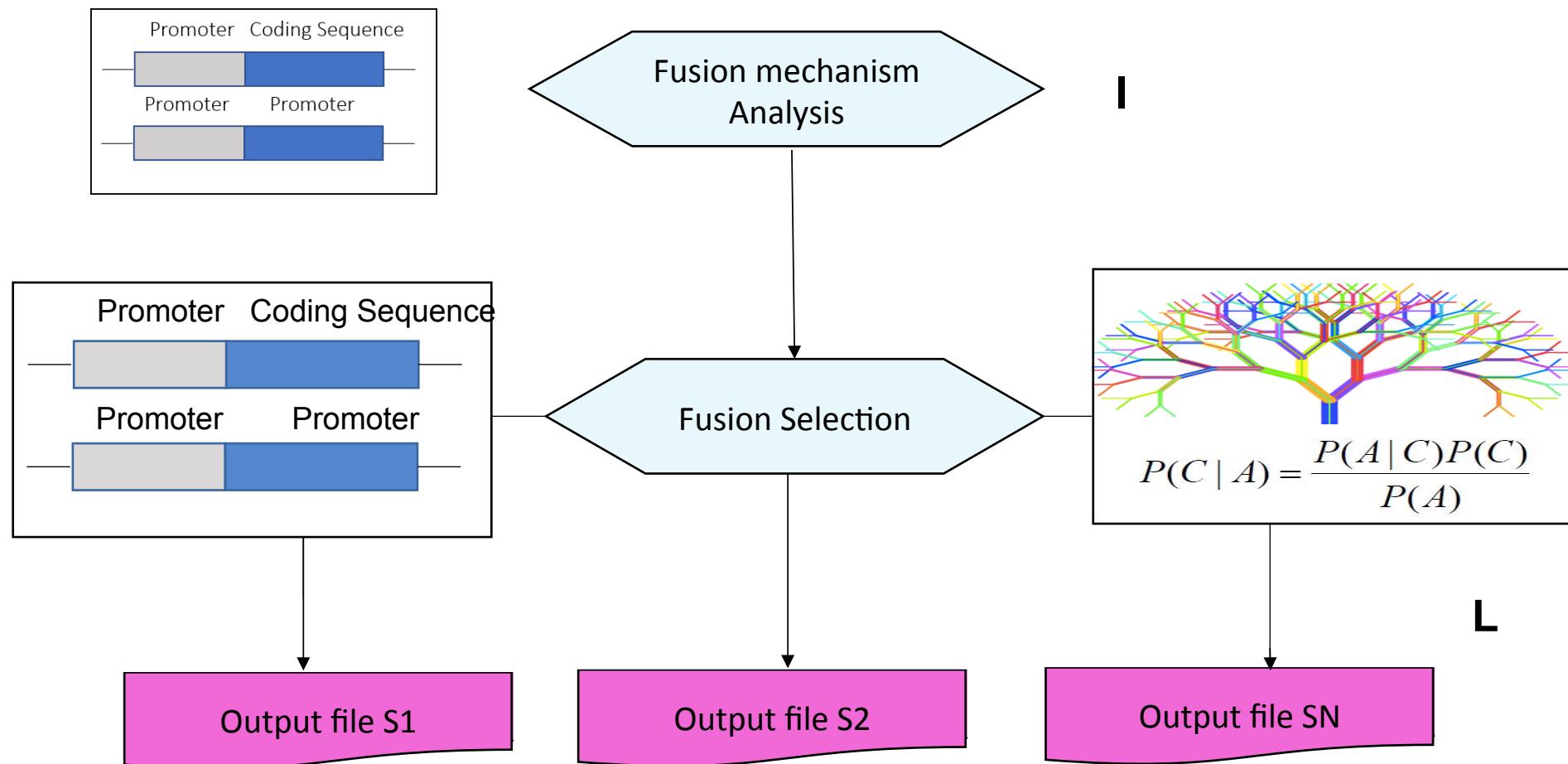
Algorithm

## How are fusions prioritized? (2)

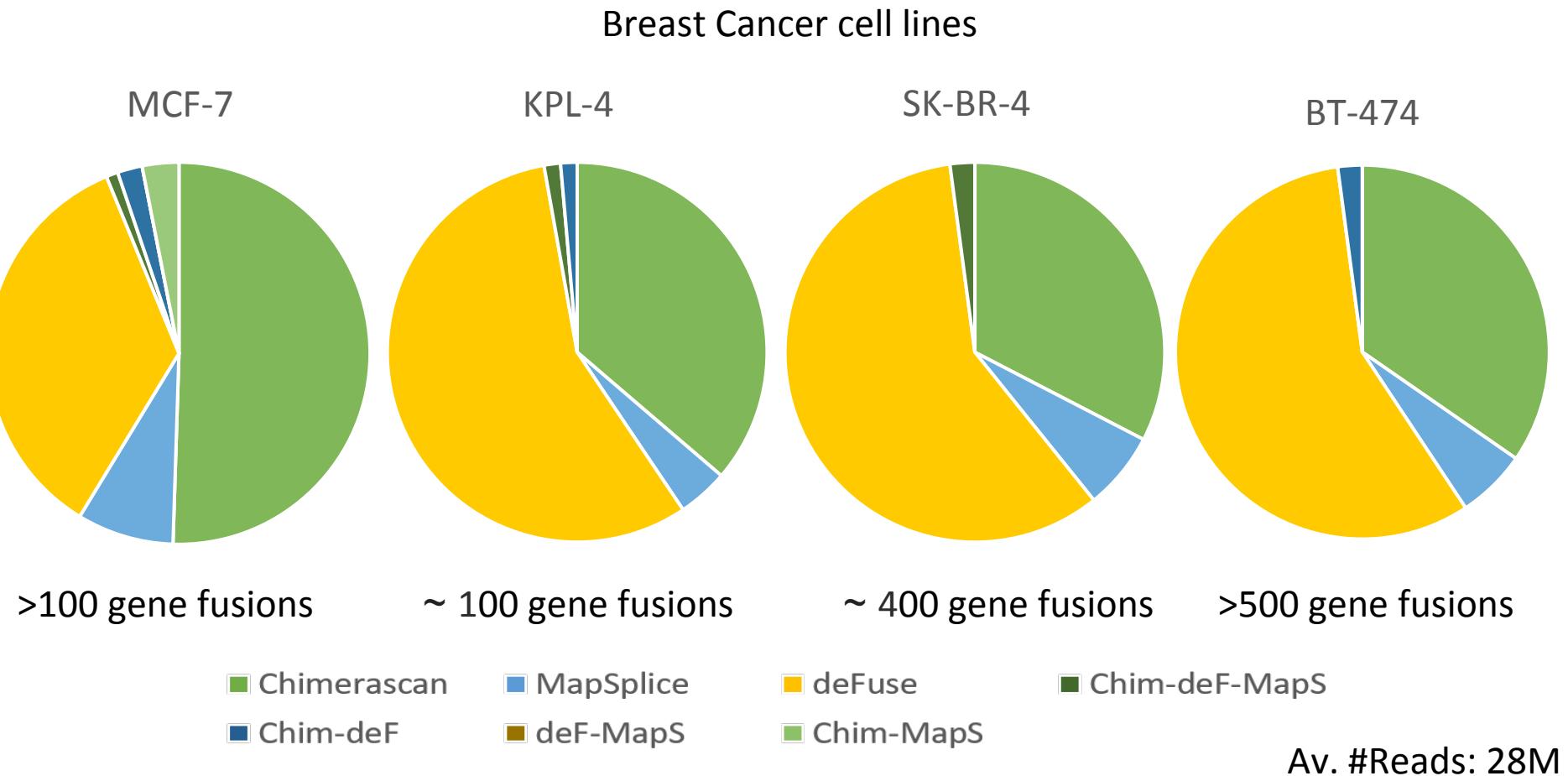


Algorithm

# How are fusions prioritized? (3)

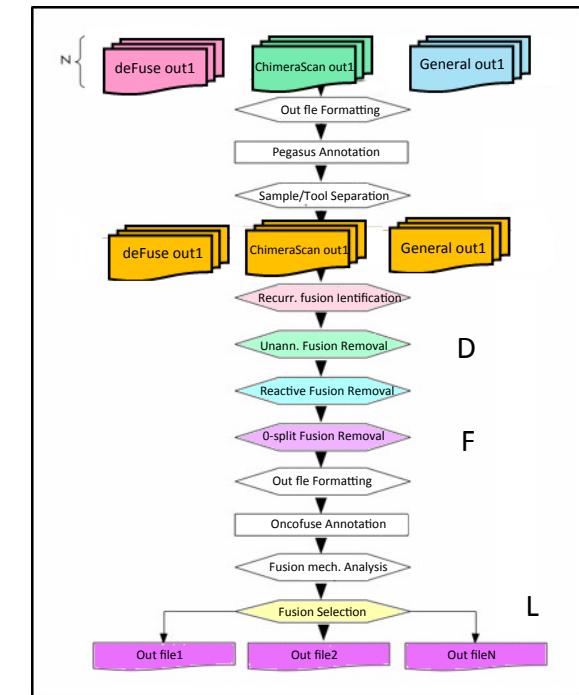
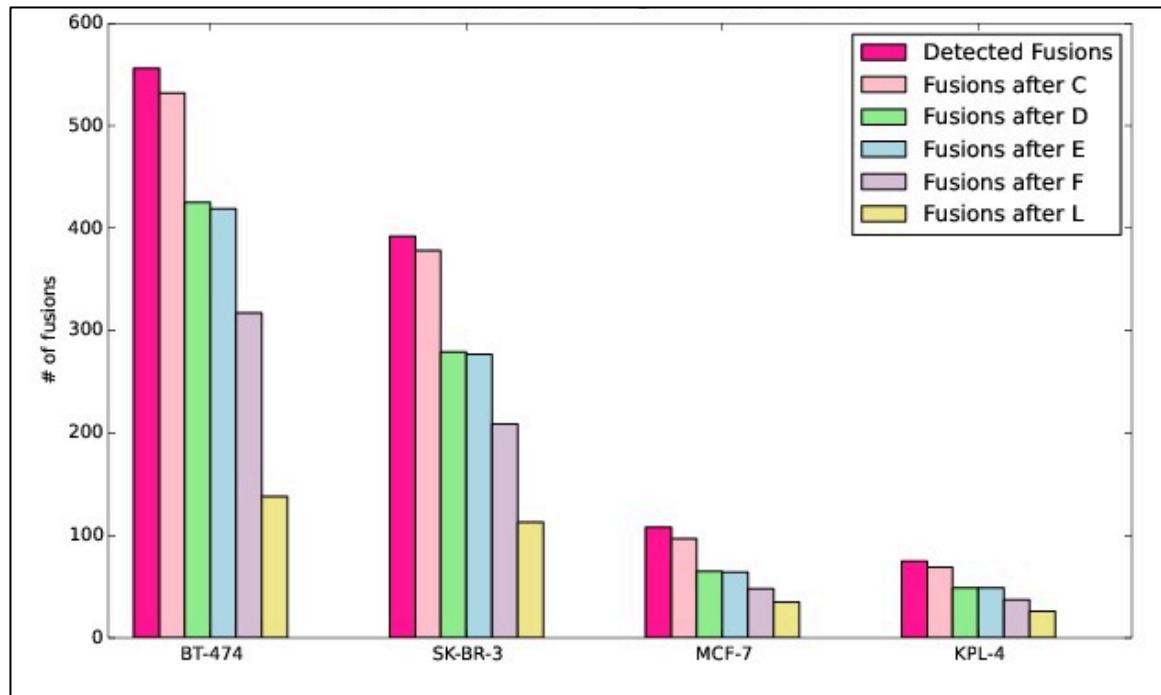


# Do we really need to prioritize fusions?



## Results

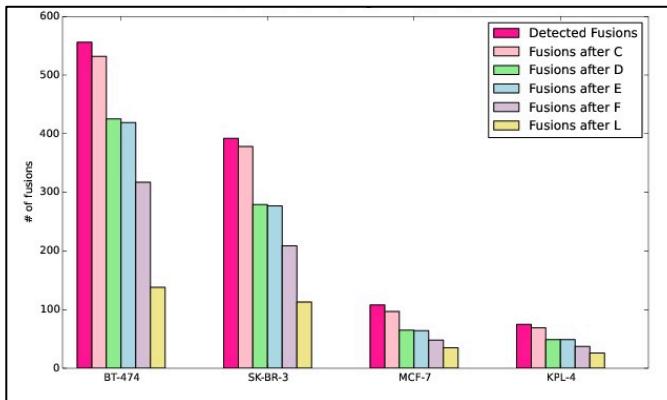
# Priority fusions in breast cancer cell lines



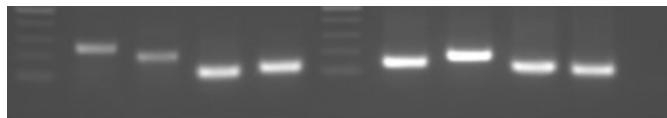
- ✓ **Reduction** in the number of fusions ranging from 65% to 75%
- ✓ 25 out of 26 **validated fusions** reported in output

Remarks and Perspectives

# To summarise



FuGePrior **prioritized** fusions from state-of-the-art gene fusion discovery tools in breast, prostate, AML and Burkitt's lymphoma cancer samples



FuGePrior **priority** fusions were confirmed by both previous researches or by PCR experiments

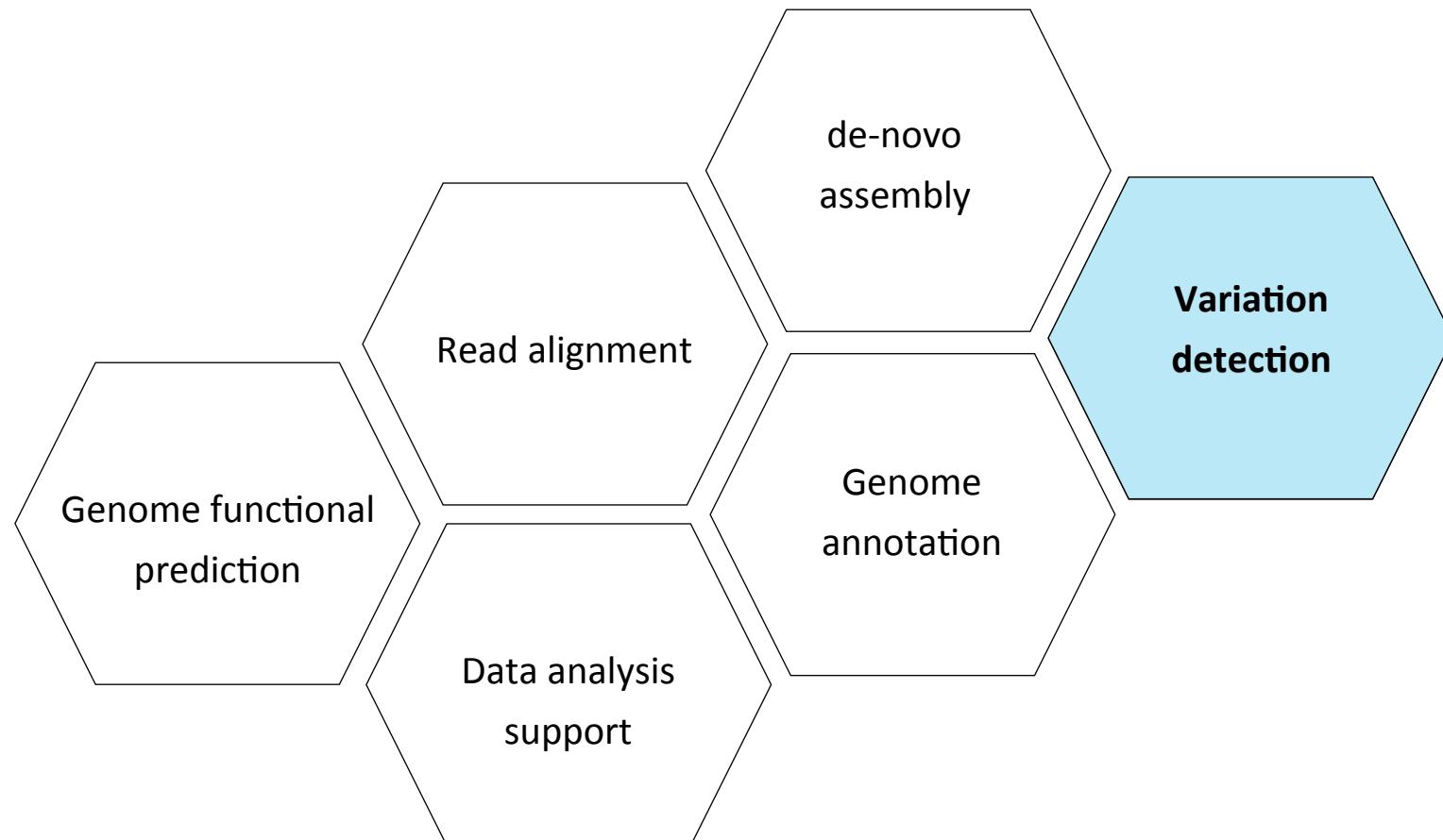
**Future work:**

additional gene fusion discovery tools

additional filtering stages

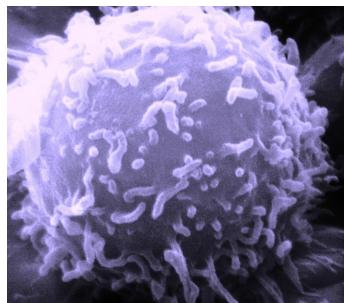
extensive ALL sample analysis

# Conclusions

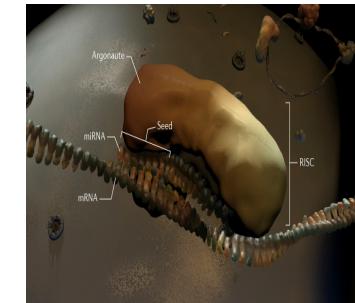
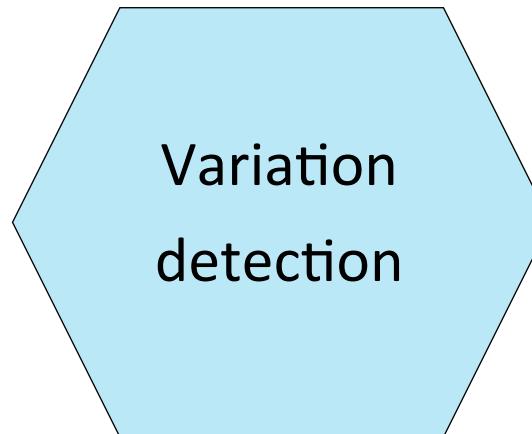




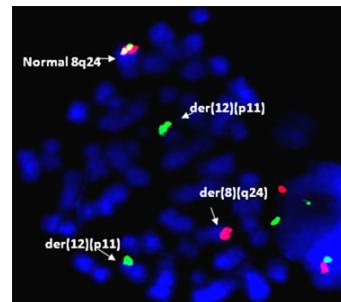
# Conclusions



VDJSeq-Solver tool

UNIVERSITÀ  
di VERONA

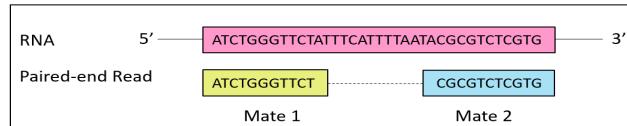
isomiR-SEA



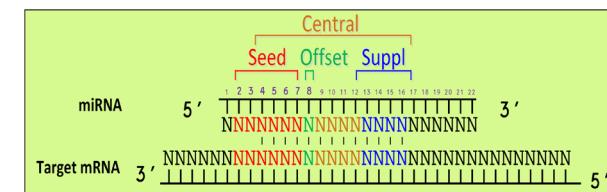
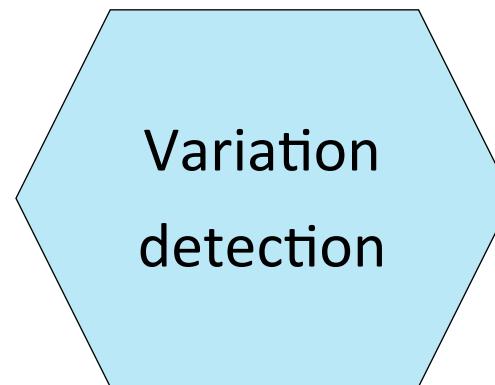
FuGePrior tool



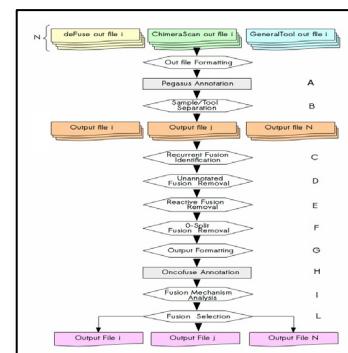
# Conclusions



- ✓ No Ig-Seq experiments
- ✓ Multiple analyses on the same datasets



- ✓ miRNA reads have a **complete seed**
- ✓ **isomiR** profiles
- ✓ Evaluation of **miRNA/isomiR targeting** activity



- ✓ **Reduction** of false positive predictions
- ✓ FuGePrior fusions are highly **reliable**



# Publications

Paciello, Giulia, et al. "A novel pipeline for V (D) J junction identification using RNA-Seq paired-end reads." In BIOSTEC 2013: , International Conference on Bioinformatics Models, Methods and Algorithms, Barcelona, 11-14 February 2013

Paciello, Giulia, et al. "VDJSeq-Solver: in silico V (D) J recombination detection tool." PloS one 10.3 (2015): e0118192

Paciello, Giulia, et al. "**Detection of rearranged light chain sequences by RNA-Seq in B-cell lymphomas and reactive lymphadenopathies.**" In: EAHP 2016: 18th Meeting of the European Association for Hematopathology, Basilea, 3-8 Sept. 2016

Urgese, Gianvito, Paciello, Giulia, et al. "miR-SEA: miRNA Seed Extension based Aligner Pipeline for NGS Expression Level Extraction." In IWBBIO 2014: 2nd International Work-Conference on Bioinformatics and Biomedical Engineering, Granada, 7-9 Apr, 2014

Urgese, Gianvito, Paciello, Giulia, et al. "On the relevance of a complete characterisation of miRNAs, isomiRs and miRNA-mRNA interaction sites through miRNA-specific alignment tools." In RNAi 2015: Short & Long Non-coding RNAs, 10th International Conference & Exhibition, Oxford, 24-26 March, 2015

Urgese, Gianvito, Paciello, Giulia, et al. "isomiR-SEA: an RNA-Seq analysis tool for miRNAs/isomiRs expression level profiling and miRNA-mRNA interaction sites evaluation." BMC bioinformatics 17.1 (2016): 148

Padella, Antonella, Paciello, Giulia, et al. "Next-Generation Sequencing Analysis Revealed That BCL11B Chromosomal Translocation Cooperates with Point Mutations in the Pathogenesis of Acute Myeloid Leukemia." Blood 124.21 (2014): 2352-2352. Proceedings of AHS Annual Meeting

Padella, Antonella, Paciello Giulia, et al. "RNA Sequencing Reveals Novel and Rare Fusion Transcripts in Acute Myeloid Leukemia." Blood 126.23 (2015): 3627-3627. Proceedings of AHS Annual Meeting

Padella, Antonella, Paciello, Giulia, et al. "Novel fusion transcripts identified by RNAseq cooperate with somatic mutations in the pathogenesis of acute myeloid leukemia." Cancer Res (76) (2016). Proceedings of AACR 107th Annual Meeting