# RNA-Seq methods and gene fusions: libraries, case reports, and algorithms

Clinical and Research Genomics

Spring 2018 Course

Andrea Sboner

03.21.2018

# Outline

- Background of transcriptome profiling

- Next Generation Sequencing: a revolution in molecular biology

- RNA-seq application: gene fusion detection

Images throughout the presentation from pixabai.com, commons.wikimedia.org
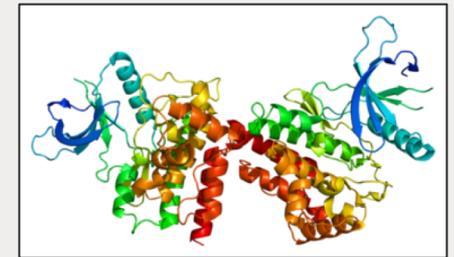
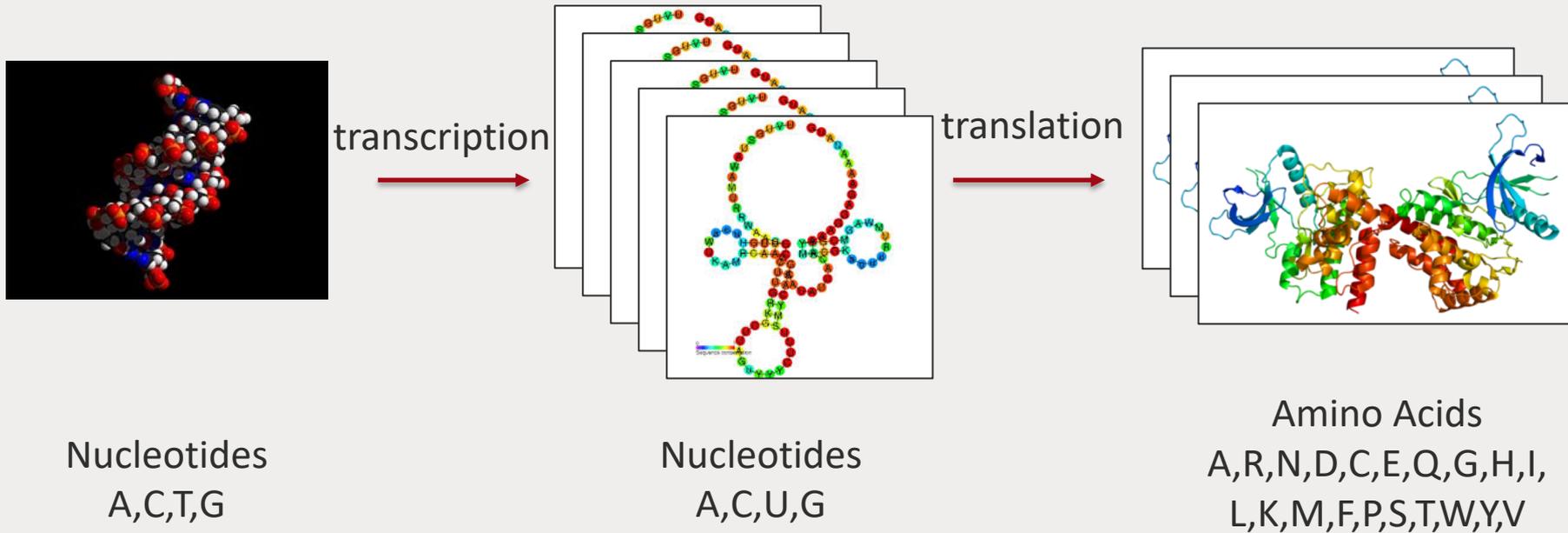# Central dogma of molecular biology

# Central dogma of molecular biology



transcription          translation

# Transcriptome profiling

transcription

translation

Nucleotides
A,C,T,G

Nucleotides
A,C,U,G

Amino Acids
A,R,N,D,C,E,Q,G,H,I,
L,K,M,F,P,S,T,W,Y,V

Transcriptome profiling goal is to characterize RNA in a tissue or cell.

The 'simpler' structure of RNA allows to employ most techniques used for DNA analysis – hybridization, polymerase chain reaction, etc.

# Pre-genome era (< 1990s)

~1970 <u>Reverse Transcriptase</u> → Allows the *reverse* transcription of RNA to DNA, generating cDNA

~ 1977 <u>Sanger Sequencing</u> → enables to 'read' the sequence of DNA

~ 1977 Northern blot → enable to measure the expression of RNA

~1983 <u>Polymerase Chain Reaction (PCR)</u> → allows the duplication/amplification of pieces of DNA

Formaldehyde gel (1%) with RNA samples run at 100V for 1 hour in 1x MOPS buffer.

Polymerase chain reaction - PCR

original DNA to be replicated

DNA primer

nucleotide

1 Denaturation at 94-96°C
2 Annealing at ~68°C
3 Elongation at ca. 72 °C

Weill Cornell Medicine        New York-Presbyterian

# Genome Era (1990s – 2000s)

~ 1991 Expressed Sequence Tags (ESTs) sequencing (500-800 nucleotides)

*Science* 21 Jun 1991; Vol. 252:Issue 5013: 1651-6

## Complementary DNA Sequencing: Expressed Sequence Tags and Human Genome Project

MARK D. ADAMS, JENNY M. KELLEY, JEANNINE D. GOCAYNE, MARK DUBNICK, MIHAEL H. POLYMEROPOULOS, HONG XIAO, CARL R. MERRIL, ANDREW WU, BJORN OLDE, RUBEN F. MORENO, ANTHONY R. KERLAVAGE, W. RICHARD MCCOMBIE, J. CRAIG VENTER*

~ 1995 Series Analysis of Gene Expression (SAGE) (9-12 nucleotides)

*Science* 20 Oct 1995:Vol. 270, Issue 5235, pp. 484-487

## Serial Analysis of Gene Expression

Victor E. Velculescu, Lin Zhang, Bert Vogelstein, Kenneth W. Kinzler*

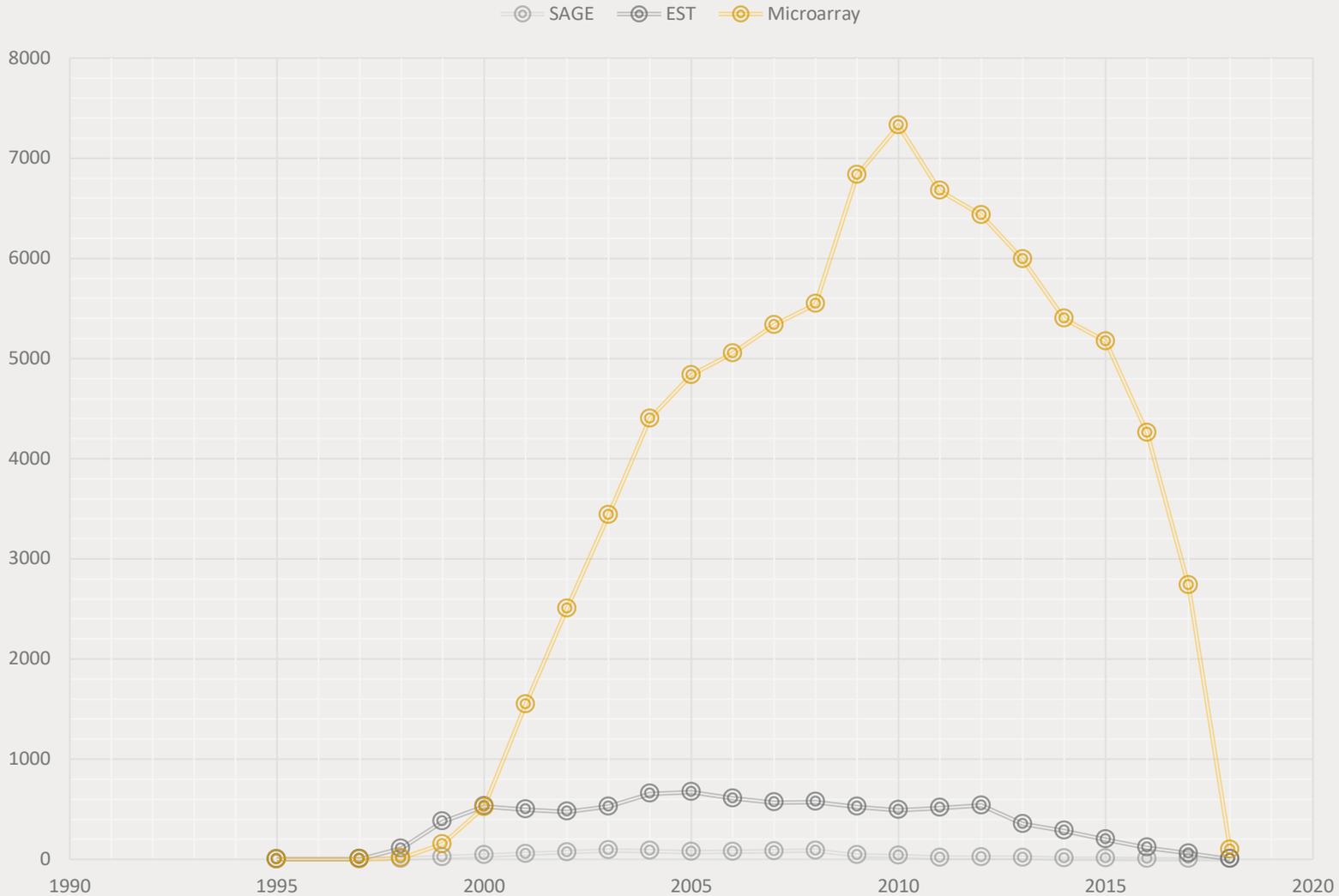Weill Cornell Medicine     NewYork-Presbyterian

# Genome Era (1990s – 2000s)

~ 1991 Expressed Sequence Tags (ESTs) sequencing (50-800 nucleotides)

~ 1995 Series Analysis of Gene Expression (SAGE) (9-12 nucleotides)

~ **1999 Microarray**



*Science* **286**, 531 (1999);

# Number or PubMed articles

# The human genome reference sequence is completed in 2003

# Cost of sequencing decreased faster than Moore's Law



Cost per Raw Megabase of DNA Sequence

Wetterstrand KA. DNA Sequencing Costs: Data from the NHGRI Genome Sequencing Program (GSP) Available at: www.genome.gov/sequencingcosts. Accessed 3.28.2017
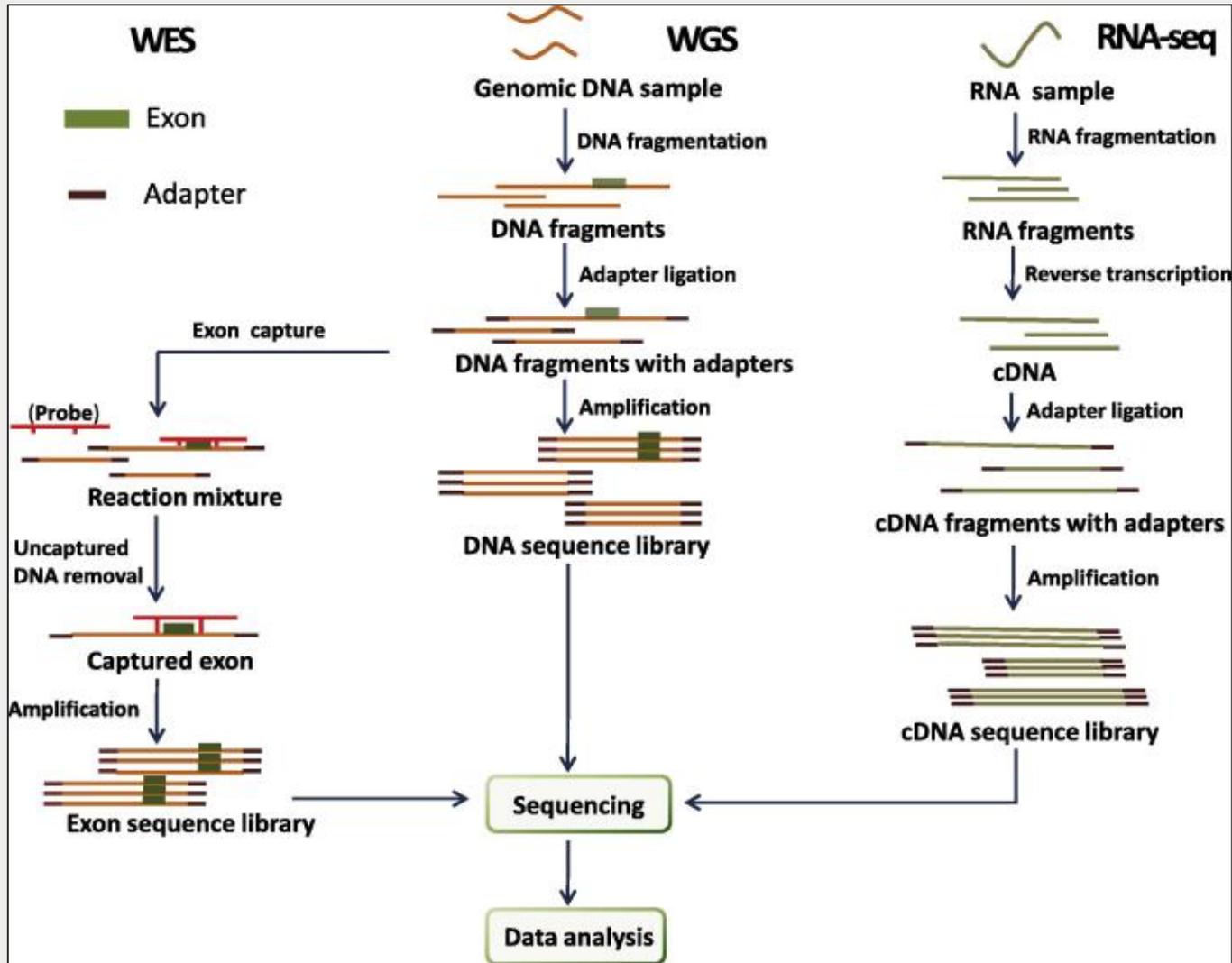
# NGS allows the rapid sequencing of millions of "short" DNA/cDNA fragments

Many applications of NGS have been developed

DNA/RNA sequencing are the most common applications of NGS



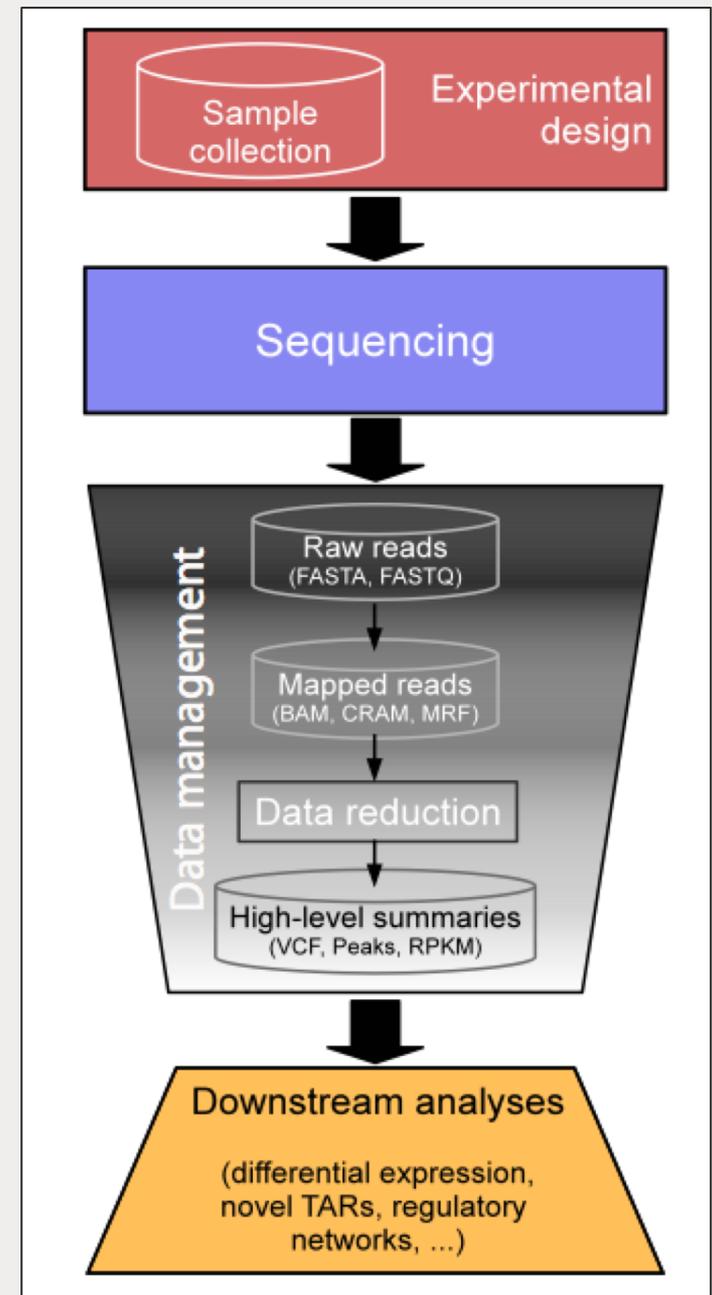For all you seq...

# Common NGS approaches

# RNA-Seq Experiment

**Data management**:

Mapping the reads
Creating summaries

**Downstream analysis**: *the interesting stuff*
Differential expression, chimeric transcripts, novel
transcribed regions, etc.

# Chimeric Transcripts
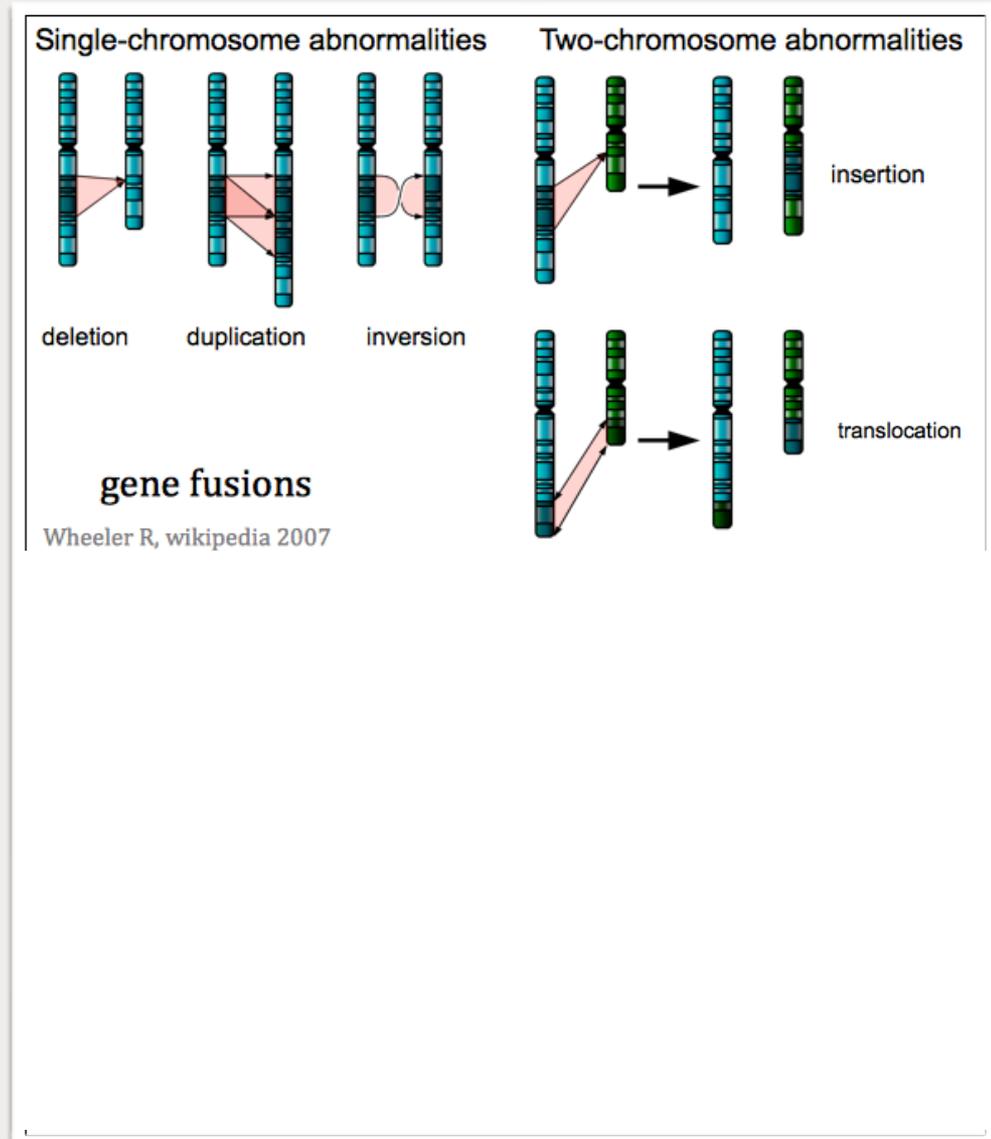
Shedding light on gene fusions

# What are chimeric transcripts?

Transcripts that are *not co-linear* in the genome space

They can arise from:

**genomic rearrangements, i.e.** *gene fusions*

**post-transcriptional events, i.e.** *trans-splicing or cis-splicing*



Single-chromosome abnormalities
Two-chromosome abnormalities

deletion   duplication   inversion

insertion

translocation

gene fusions

Wheeler R, wikipedia 2007

**Weill Cornell Medicine**    **New York-Presbyterian**

# Why are they (gene fusions) important?

Fusion genes are often *oncogenes*

**Ex: BCR-ABL1 (Philadelphia chromosome) in Chronic myelogenous leukemia (CML) and Acute Lymphoblastic leukemia (ALL) t(9;22)(q34;q11)**

Fusion involving a proto-oncogene with a strong promoter resulting in *upregulation* (lymphomas)

**Ex: (IgH locus)-MYC in Burkitt's lymphoma (cMYC over-expressed)**





Hampton OA et al. Genome Res 2009

# Why are they (trans-splicing events) important?

Trans(cis)-splicing was initially found in lower eukariotes, such as trypanosomes and worms
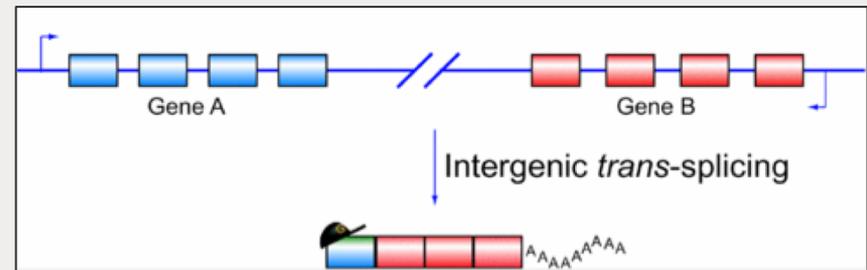
**Short sequences of nucleotides are trans-spliced to distant 5' of many protein coding genes**

Recently, they were found in mammalian cells:

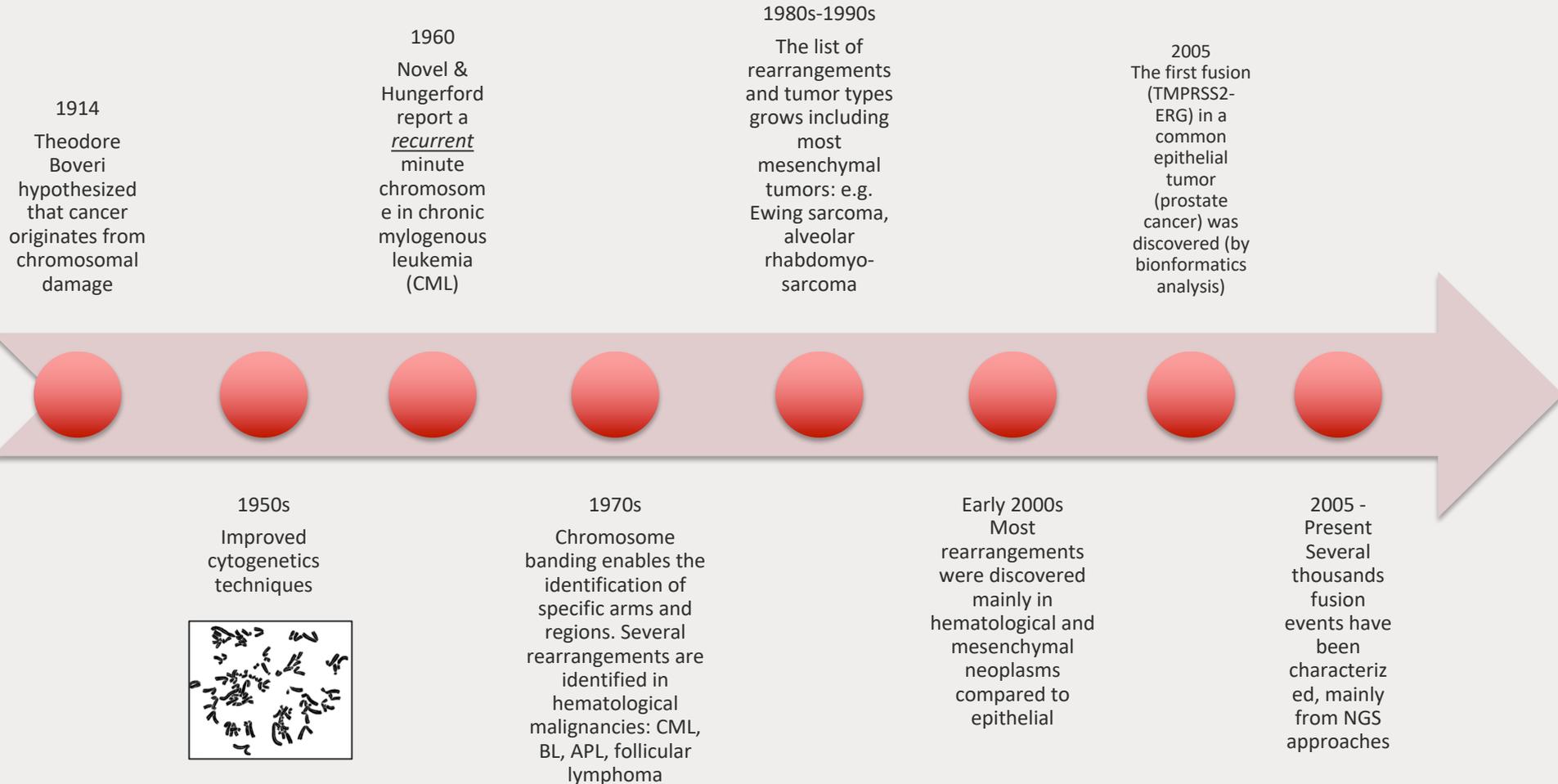**JAZF1-SUZ12 in endometrial stroma cells (Li et al. Science 2008)**

**SLC45A3-ELK4 in prostate tissues (Rickman et al. Cancer Res 2009)**

65% of protein-coding genes have distal 5' transcription start sites (ENCODE pilot) --> revised to ~50% the ENCODE 2012
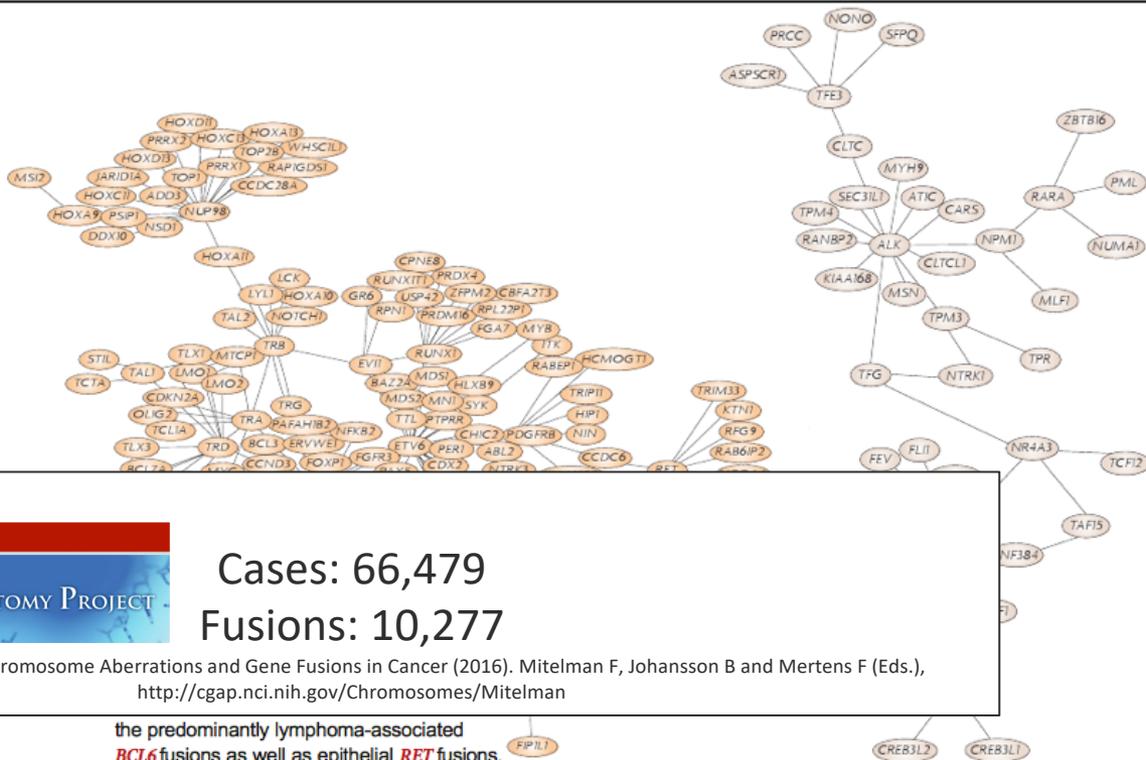


Horiuchi, Takayuki, and Toshiro Aigaki. Biology of the Cell 98, no. 2 (January 9, 2012): 135–140.

# An historical perspective of gene fusions

**1914**
Theodore Boveri hypothesized that cancer originates from chromosomal damage

**1960**
Novel & Hungerford report a _recurrent_ minute chromosome in chronic mylogenous leukemia (CML)

**1980s-1990s**
The list of rearrangements and tumor types grows including most mesenchymal tumors: e.g. Ewing sarcoma, alveolar rhabdomyo-sarcoma

**2005**
The first fusion (TMPRSS2-ERG) in a common epithelial tumor (prostate cancer) was discovered (by bionformatics analysis)

**1950s**
Improved cytogenetics techniques

**1970s**
Chromosome banding enables the identification of specific arms and regions. Several rearrangements are identified in hematological malignancies: CML, BL, APL, follicular lymphoma

**Early 2000s**
Most rearrangements were discovered mainly in hematological and mesenchymal neoplasms compared to epithelial

**2005 - Present**
Several thousands fusion events have been characterized, mainly from NGS approaches

# How many different gene fusions do we know?



- 358 gene fusion
- 337 different genes
- ~90% form three clusters

National Cancer Institute
**CANCER GENOME ANATOMY PROJECT**

Cases: 66,479
Fusions: 10,277

Mitelman Database of Chromosome Aberrations and Gene Fusions in Cancer (2016). Mitelman F, Johansson B and Mertens F (Eds.), http://cgap.nci.nih.gov/Chromosomes/Mitelman

the predominantly lymphoma-associated *BCL6* fusions as well as epithelial *RET* fusions.

haematological *MLL* fusions connected to the *HMGA2* fusions typically found in soft tissue tumours.

lymphoma-associated *ALK* fusions, the carcinoma-associated transcription factor for IGHM enhancer 3 (*TFE3*) fusions, and the sarcoma-associated *EWSR1* fusions.

Mitelman F et al, Nature Rev Cancer 2007

# Gene fusions are important for clinical treatment…

# … and diagnostic/prognostic purposes

Exclusively present in *epitheliod hemangioendothelioma*

| G | WWTR1 | | CAMTA1 | |
|---|---|---|---|---|
| | Positive /total | % | Positive /total | % |
| Epithelioid hemangioendothelioma | 42/47 | 89% | 39/45 | 87% |
| Angiosarcoma, NOS | 0/42 | 0% | 0/39 | 0% |
| Epithelioid angiosarcoma | 0/7 | 0% | 0/7 | 0% |
| Intimal sarcoma | 0/5 | 0% | 0/3 | 0% |
| Kaposi's sarcoma | 0/4 | 0% | 0/4 | 0% |
| Malignant hemangioendothelioma, NOS | 0/1 | 0% | 0/1 | 0% |
| Retiform hemangioendothelioma | 0/1 | 0% | 0/1 | 0% |
| Kaposiform hemangioendothelioma | 0/3 | 0% | 0/2 | 0% |
| Epithelioid hemangioma | 0/5 | 0% | 0/4 | 0% |
| Arteriovenous malformation | 0/2 | 0% | 0/2 | 0% |
| Angiomatosis | 0/1 | 0% | 0/1 | 0% |
| Hemangioma, NOS | 0/3 | 0% | 0/3 | 0% |
| Capillary/pyogenic hemangioma | 0/5 | 0% | 0/5 | 0% |
| Cavernous hemangioma | 0/5 | 0% | 0/5 | 0% |
| Juvenile hemangioma | 0/1 | 0% | 0/1 | 0% |
| Spindle cell hemangioma | 0/4 | 0% | 0/4 | 0% |
| Synovial hemangioma | 0/1 | 0% | 0/1 | 0% |
| Intramuscular hemangioma | 0/6 | 0% | 0/5 | 0% |
| Littoral cell hemangioma | 0/6 | 0% | 0/2 | 0% |
| Malignant hemangiopericytoma | 0/1 | 0% | 0/1 | 0% |
| Hemangiopericytoma, NOS | 0/1 | 0% | 0/1 | 0% |
| Sinonasal hemangiopericytoma | 0/1 | 0% | 0/1 | 0% |
| Glomus tumor | 0/1 | 0% | 0/1 | 0% |
| Atypical glomus tumor | 0/2 | 0% | 0/2 | 0% |
| Lymphangioma | 0/7 | 0% | 0/7 | 0% |
| Lymphangioleiomyomatosis | 0/1 | 0% | 0/1 | 0% |
| Papillary endothelial hyperplasia | 0/2 | 0% | 0/2 | 0% |
| Total cases | 165 | | 151 | |

# How to identify fusion transcripts from paired-end RNA-seq?



RNA-seq
RNA sample
RNA fragmentation
RNA fragments
Reverse transcription
cDNA
Adapter ligation
cDNA fragments with adapters
Amplification
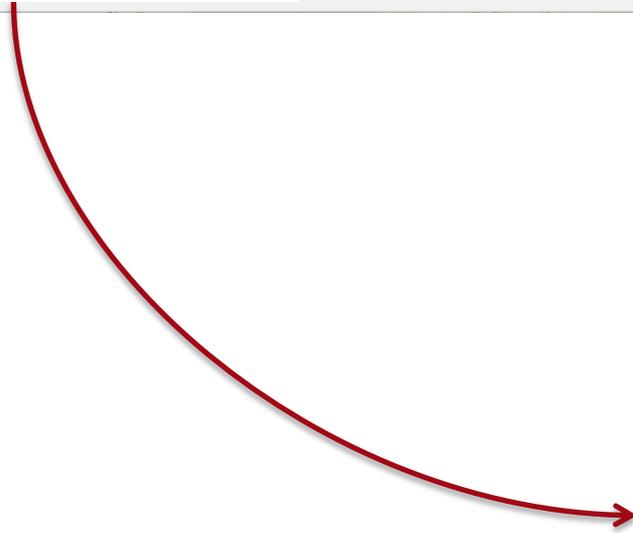cDNA sequence library
Sequencing

Paired-end sequencing means that we know the sequence of the two ends of a fragment

# Mapping

Google   Institute for Computational Biomedicine   🔍

# Mapping

# How to identify fusion transcripts from paired-end RNA-seq?



RNA-seq
RNA sample
↓ RNA fragmentation
RNA fragments
↓ Reverse transcription
cDNA
↓ Adapter ligation
cDNA fragments with adapters
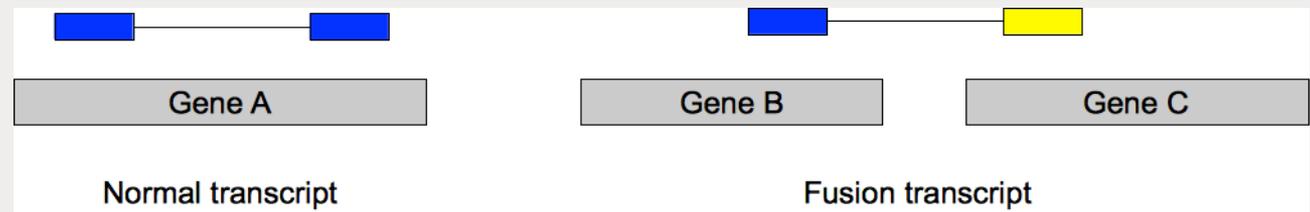↓ Amplification
cDNA sequence library
↓
Sequencing

Paired-end sequencing means that we know the sequence of the two ends of a fragment

*Straightforward*:

**If the two ends map to different genes, then we have a potential fusion transcript**

Gene A — Normal transcript

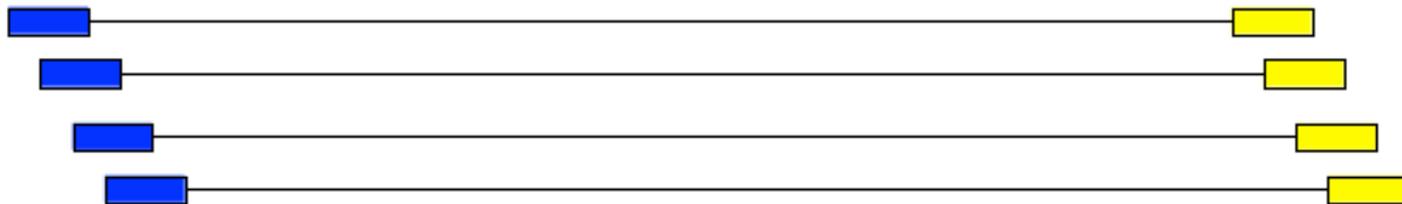Gene B — Gene C — Fusion transcript

# What about different isoforms?

# Composite model



- Each PE read can be assigned to one "gene"

- *Potential Fusion Transcripts*: if pair belongs to different genes

# Not an ideal word: sources of errors

*Mis-alignments*

**Base caller error**

**SNPs**

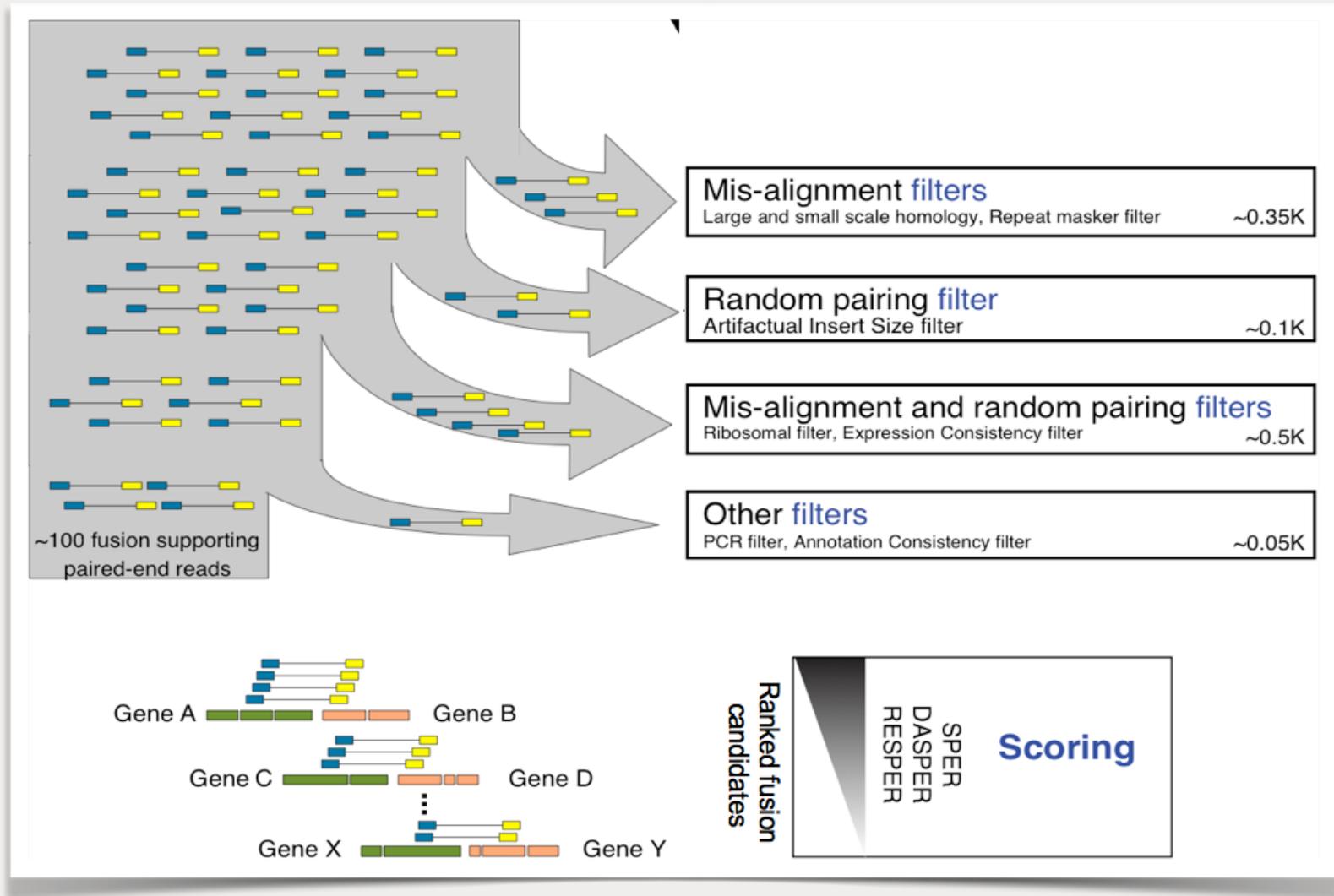**RNA editing**

**Sequence similarity (paralogs, pseudogenes)**

*Random pairing of transcript fragments*

**Library preparation**

*Combination of mis-alignment and random pairing*

*PCR amplification, gene annotation inconsistencies/incompleteness*
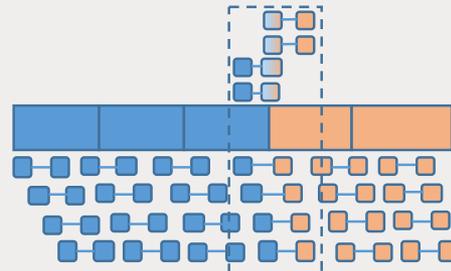
# Filtration Cascade Module



Mis-alignment **filters**
Large and small scale homology, Repeat masker filter — ~0.35K

Random pairing **filter**
Artifactual Insert Size filter — ~0.1K

Mis-alignment and random pairing **filters**
Ribosomal filter, Expression Consistency filter — ~0.5K

Other **filters**
PCR filter, Annotation Consistency filter — ~0.05K

~100 fusion supporting paired-end reads

Gene A — Gene B
Gene C — Gene D
Gene X — Gene Y

Ranked fusion candidates

SPER DASPER RESPER

**Scoring**

# Augmenting the support for fusion: fusion junction reads

# Tools for detecting fusion transcripts

From sequencing data

http://omictools.com/gene-fusion-detection-category

http://omictools.com/transcriptome-assembly-category

**RNA-seq short-reads "only"**
Bellerophontes
BreakFusion
chimeraScan
CRAC
deFuse
EricScript
FusionAnalyser
FusionCatcher
FusionFinder
FusionHunter
FusionQ
FusionSeq
Jaffa
MapSplice
PRADA
shortFuse
SnowShoes-FTD
SOAPFuse/Fusion
TopHat-Fusion
STAR-fusion

**RNA-seq & DNA-seq**
BreakTrans
Comrad
nFuse

**Gene fusion annotation**
Chimera
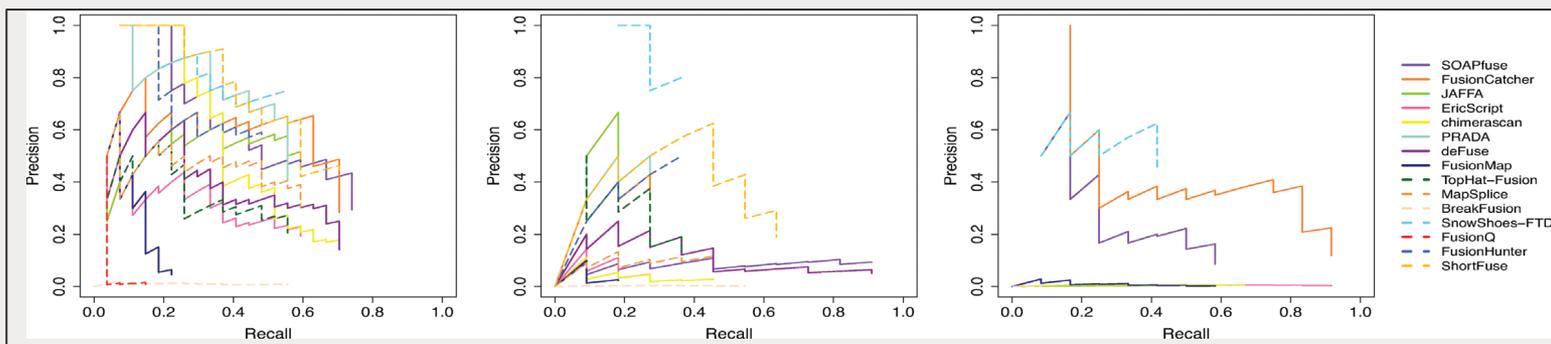Pegasus

**Transcript Assembly**
CuffLinks
Scripture
Trinity
Trans-Abyss

**Comprehensive evaluation of fusion transcript detection algorithms and a meta-caller to combine top performing methods in paired-end RNA-seq data**

**Silvia Liu[1,2,†], Wei-Hsiang Tsai[3,†], Ying Ding[1,2,†], Rui Chen[1], Zhou Fang[1], Zhiguang Huo[1], SungHwan Kim[1], Tianzhou Ma[1], Ting-Yu Chang[4], Nolan Michael Priedigkeit[5], Adrian V. Lee[6], Jianhua Luo[7], Hsei-Wei Wang[3,4,8,*], I-Fang Chung[3,8,*] and George C. Tseng[1,2,*]**

# Summary and Future directions

- Massively Parallel Sequencing has enabled the discovery of fusion transcripts

- Specificity is the main challenge: too many false positives!

- <u>Longer reads</u>: could help overcome the limitations of short reads

- <u>Combination of tools</u> may help further improve on the reduction of FP

- "For the large bioinformatics community, development of a high-performing (accurate and fast) fusion detection tool or methods to combine top- performing tools remains an important and open question"

<u>ans2077@med.cornell.edu</u>