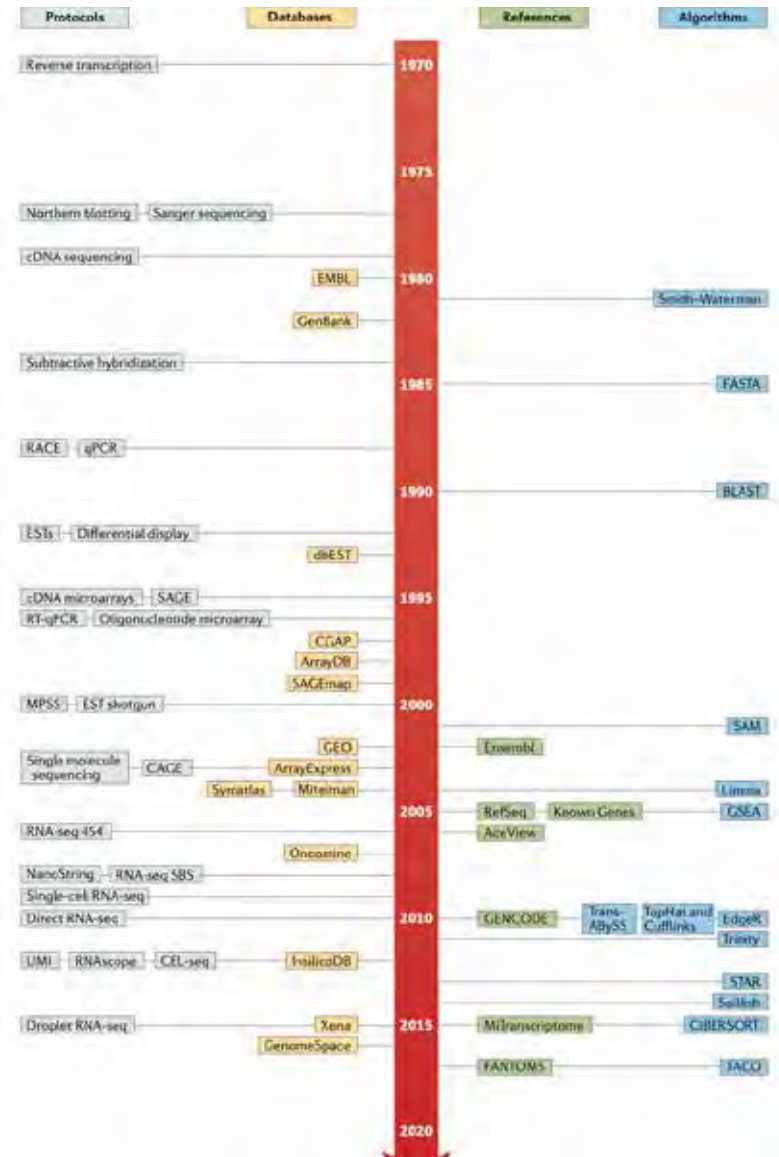


TRANSCRIPTOMICS

David W. Craig, Ph.D

BACKGROUND



Objective: Broad survey of RNA-Seq & Cancer

- └ Focus on breadth over depth
- └ Focus on methods focusing on tumor RNA-seq
 - ┆ Will not cover eQTL
- └ Focus on human applications and RNA
 - ┆ Epigenetics/mouse largely not covered

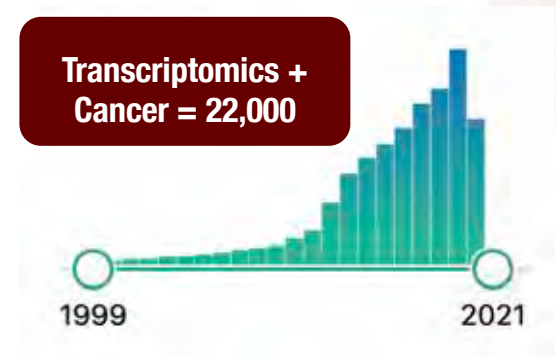
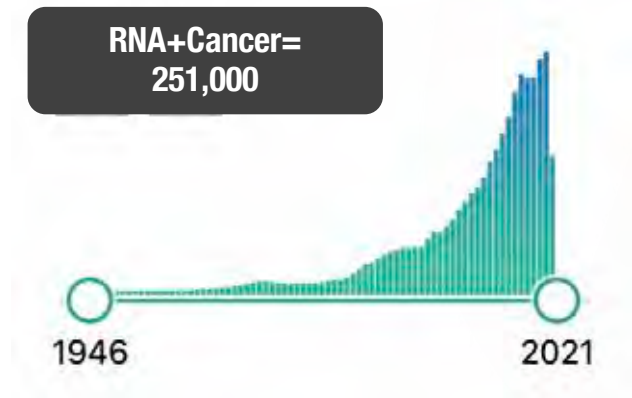
Transcriptomics:

- └ Background & Core concepts w/ NGS

Applications

- └ Bulk Applications
- └ Emerging Methods

Basics of Analysis



Bound to disappoint those dedicated those who live transcriptomics

REVIEW PAPERS

APPLICATIONS OF NEXT-GENERATION SEQUENCING

Translating RNA sequencing into clinical diagnostics: opportunities and challenges

Sara A. Byron¹, Kendall R. Van Keuren-Jensen², David M. Engelthaler¹, John D. Corbett¹ and David W. Craig¹

Abstract | With the emergence of RNA sequencing (RNA-seq) technologies, RNA-based biomarkers look expanded promise for their diagnostic, prognostic, and therapeutic applicability in various diseases, including cancers and infectious diseases. Detection of gene fusions and differential expression of lncRNA disease-causing transcripts by RNA-seq represent some of the most immediate opportunities. However, it is the diversity of RNA species detected through RNA-seq that holds new promise for the multi-focused clinical applicability of RNA-based measures, including the potential of extracellular RNAs as non-invasive diagnostic indicators of disease. Ongoing efforts towards the establishment of bench-to-bed standards, assay optimization for clinical conditions and demonstration of assay reproducibility are required to expand the clinical utility of RNA-seq.

RNA sequencing: the teenage years

Rory Stark¹, Marta Grzelik² and James Hastfield³*

Abstract | Over the past decade, RNA sequencing (RNA-seq) has become an indispensable tool for transcriptome-wide analysis of differential gene expression and differential splicing of mRNAs. However, as next-generation sequencing technologies have developed, so too has RNA-seq. Now, RNA-seq methods are available for studying many different aspects of RNA biology, including single-cell gene expression, translation (the translatome) and RNA structure (the structurome). Exciting new applications are being explored, such as spatial transcriptomics (spatialomics). Together with new long-read and direct RNA-seq technologies and better computational tools for data analysis, innovations in RNA-seq are contributing to a fuller understanding of RNA biology, from questions such as when and where transcription occurs to the folding and intermolecular interactions that govern RNA function.

APPLICATIONS OF NEXT-GENERATION SEQUENCING

Cancer transcriptome profiling at the juncture of clinical translation

Marcin Cieřlik^{1,2} and Arul M. Chinnaiyan^{1,3}*

Abstract | Methodological breakthroughs over the past four decades have repeatedly revolutionized transcriptome profiling. Using RNA sequencing (RNA-seq), it has now become possible to sequence and quantify the transcriptional outputs of individual cells or thousands of samples. These transcriptomes provide a link between cellular phenotypes and their molecular underpinnings, such as mutations. In the context of cancer, this link represents an opportunity to dissect the complexity and heterogeneity of tumours and to discover new biomarkers or therapeutic strategies. Here, we review the rationale, methodology and translational impact of transcriptome profiling in cancer.



Genomic basis for RNA alterations in cancer

<https://doi.org/10.1038/s41566-020-1970-0>

Received: 29 March 2018

Accepted: 11 December 2019

Published online: 5 February 2020

Open access

PCAWG Transcriptome Core Group^{1,2}, Claudia Calabrese^{3,4}, Natalie R. Davidson^{1,5,6,7,8,9,10}, Deniz Demirciođlu^{11,12}, Nuno A. Fonseca^{1,13}, Yao He^{14,15}, André Kahles^{1,16,17,18}, Kjong Van Lehmann^{1,19,20,21}, Fenglin Liu^{22,23}, Yuichi Shirahishi^{24,25}, Cameron M. Soulette^{26,27}, Lara Urban²⁸, Liliana Greger²⁹, Siqiang Liu^{30,31}, Dongbing Liu^{32,33}, Marc D. Perry^{34,35}, Qian Xiang³⁶, Fan Zhang³⁷, Junjun Zhang³⁸, Peter Bailey³⁹, Serap Erkek⁴⁰, Katherine A. Hoodley⁴¹, Yong Hou^{42,43}, Matthew R. Huska⁴⁴, Helena Kilpinen⁴⁵, Jan O. Korbel⁴⁶, Maximilian G. Marin⁴⁷, Julia Markowski⁴⁸, Tannistha Nandi⁴⁹, Qiang Pan-Hammarström^{50,51}, Chandra Sekhar Pedamallu^{52,53,54}, Reiner Siebert⁵⁵, Stefan G. Stark^{56,57}, Hong Su^{58,59}, Patrick Tan⁶⁰, Sebastian M. Waszak⁶¹, Christina Yang⁶², Shida Zhu^{63,64}, Philip Awadalla^{65,66}, Chad J. Creighton⁶⁷, Matthew Meyerson^{68,69,70}, B. F. Francis Ouellette⁷¹, Kui Wu^{72,73}, Huanming Yang⁷⁴, PCAWG Transcriptome Working Group⁷⁵, Alvis Brazma^{76,77}, Angela N. Brooks^{78,79,80,81}, Jonathan Cooke^{82,83}, Gunnar Rätsch^{84,85,86,87}, Roland F. Schwarz^{88,89,90,91}, Oliver Stegle^{92,93,94}, Zemin Zhang^{95,96} & PCAWG Consortium⁹⁷

Annual Review of Cancer Biology

Deciphering Human Tumor Biology by Single-Cell Expression Profiling

Itzy Tirosh¹ and Mario L. Suvà^{2,1}

¹Department of Molecular Cell Biology, Weizmann Institute, 23102 Rehovot, ISRAEL; itzy.tirosh@weizmann.ac.il

²Broad Institute of Harvard and MIT, Cancer Program, Massachusetts 02142, USA

³Department of Pathology and Center for Cancer Research, Massachusetts General Hospital and Harvard Medical School, Boston, Massachusetts 02114, USA; mario.suva@broadinstitute.org

Next-generation computational tools for interrogating cancer immunity

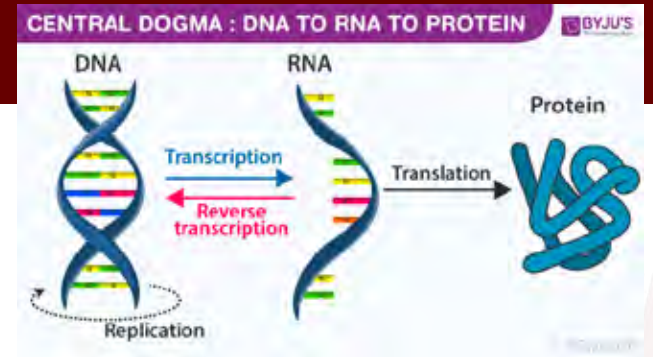
Francesca Finotello¹, Dietmar Rieder, Hubert Hochl² and Tatjana Trajanoski³*

Abstract | The remarkable success of cancer therapies with immune checkpoint blockers is revolutionizing oncology and has sparked an intensive basic and translational research into the mechanisms of cancer-immune cell interactions. In parallel, numerous novel cutting-edge technologies for comprehensive molecular and cellular characterization of cancer immunity have been developed, including single-cell sequencing, mass cytometry and multiplexed spatial cellular phenotyping. In order to process, analyze and visualize multidimensional data sets generated by these technologies, computational methods and software tools are required. Here, we review computational tools for interrogating cancer immunity, discuss advantages and limitations of the various methods and provide guidelines to assist in method selection.

WITH IMPORTANT EXCEPTIONS...

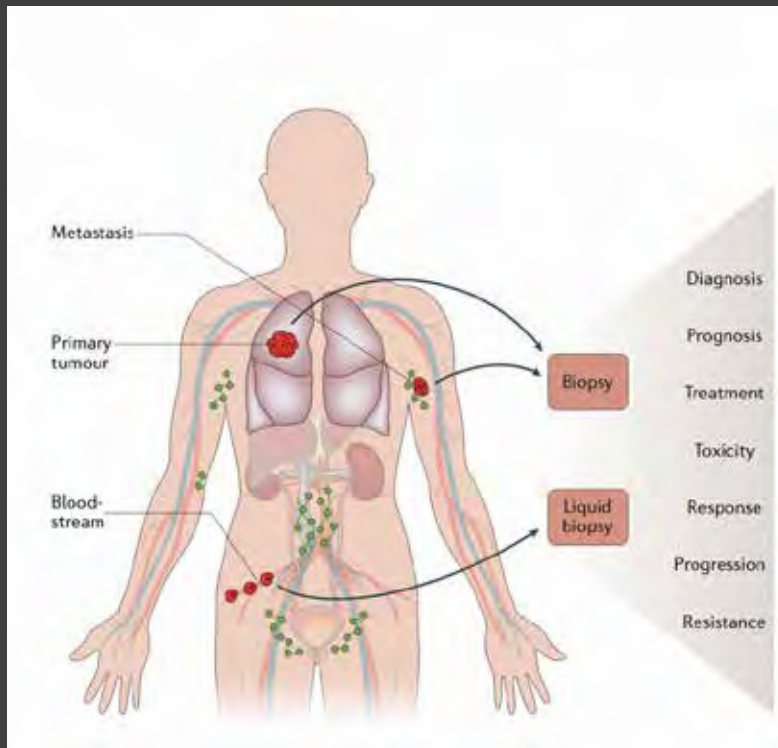
- ... you are diploid with a maternal and paternal copy
- ... you have two copies of 22 chromosomes plus X and sometimes Y
- ... there are four nucleotides (A, T, C, G), about 3 billion bases long (ATTATA..)
- ... a copy of your genome is every cell.
- ... there are 4 millions genetic variants between two people (3 billion)
- ... Variants are of different types. Single nucleotide substitutions (SNVs), Insertions/Deletions (indels), Structural Variants (inversions, duplications, translocations)
- ... most genetic variants are not functional. There are many variants in genes like BRCA1, having a variant does not mean you carry the BRCA1 gene.
- ... changes occurring in a specific tissue or cell during our life are called somatic events
- ... SNPs \in SNVs. SNPs are inherited. Polymorphisms are common in population & you weren't born with cancer.
- ... 1% of your genome is coded in genes, sometimes this is called your exome
- ... in genes, DNA is transcribed to RNA, RNA is translated to proteins
- ... genes are frequently transcribed as exons broken by introns, where the introns are spliced out of mRNA
- ... a considerable number of modifications can occur to proteins (e.g. phosphorylation)
- ... 99% of your genome we don't understand, but we all recognize its important.
- ... two identically cloned calico cats look nothing alike because epigenetics matters

The exceptions are often the most important aspects of understanding and treating diseases.



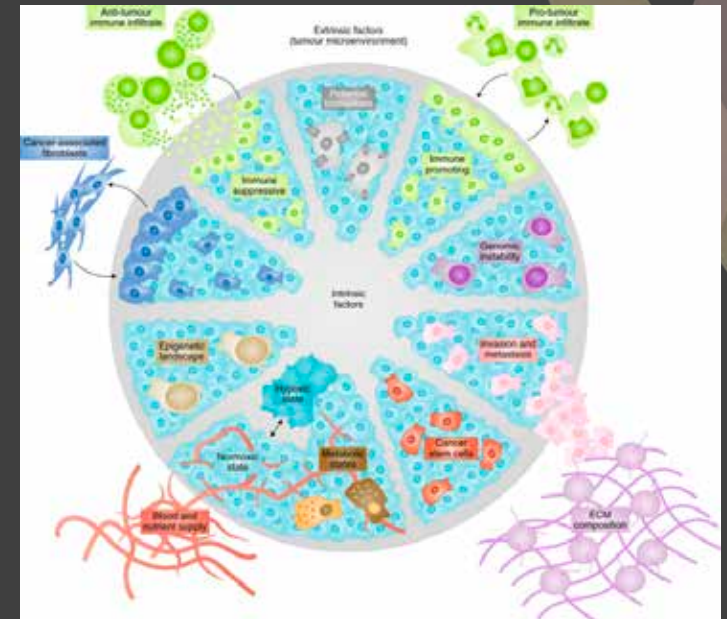
WHY TRANSCRIPTOMICS MATTERS IN CANCER

Interpretation of functional impact of DNA variation
Provide possible biomarkers for diagnosis or progression
Give insight into biological drivers, response, therapies



Tumour heterogeneity and metastasis at single-cell resolution

Shen R, Liaw D, et al. Cell. 2015;161(6):1338-1347.



HUMAN GENOME PROJECT

articles

Initial sequencing and analysis of the human genome

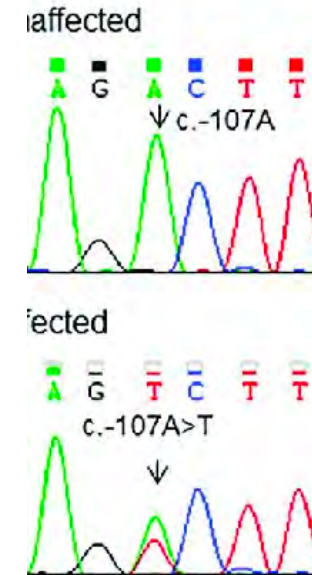
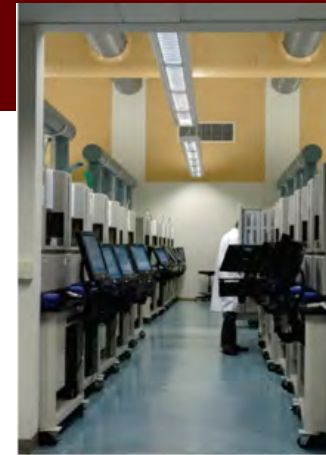
International Human Genome Sequencing Consortium*

* A partial list of authors appears on the opposite page. Affiliations are listed at the end of the paper.

The human genome holds an extraordinary trove of information about human development, physiology, medicine and evolution. Here we report the results of an international collaboration to produce and make freely available a draft sequence of the human genome. We

Table 8 Chromosome size estimates

Chromosome*	Sequenced bases† (Mb)	FCC gaps‡		SCC gaps‡		Sequence gaps‡		Heterochromatin and short arm adjustments** (Mb)	Total estimated chromosome size (including artefactual duplication in draft genome sequence)†† (Mb)	Previously estimated chromosome size†† (Mb)
		Number	Total bases in gaps§ (Mb)	Number	Total bases in gaps§ (Mb)	Number	Total bases in gaps§ (Mb)			
All	2,892.9	897	152.0	4,076	142.7	145,514	80.6	212	3,289	3,286
1	212.2	104	17.7	347	12.1	11,803	6.5	30	279	263
2	221.6	50	8.5	296	10.4	12,660	7.1	3	251	255
3	186.2	71	12.1	336	11.8	14,669	8.1	3	221	214
4	168.1	39	6.6	343	12.0	12,768	7.1	3	197	203
5	169.7	46	7.8	337	11.8	10,304	5.7	3	198	194
6	158.1	15	2.6	275	9.6	5,225	2.9	3	176	183
7	146.2	27	4.6	195	6.8	4,338	2.4	3	163	171
8	124.3	41	7.0	249	8.7	8,692	4.6	3	148	155
9	106.9	19	3.2	122	4.3	6,063	3.4	22	140	145
10	127.1	14	2.4	163	5.7	8,947	5.0	3	143	144
11	128.8	29	4.9	183	6.8	8,279	4.6	3	148	144
12	124.5	26	4.4	168	5.9	8,226	4.6	3	142	143
13	92.9	12	2.0	115	4.0	5,065	2.8	16	118	114
14	86.9	13	2.2	40	1.4	775	0.4	16	107	109
15	73.4	18	3.1	104	3.6	5,717	3.2	17	100	106
16	73.1	56	9.4	102	3.6	4,757	2.6	15	104	98
17	72.8	41	7.0	95	3.3	4,201	2.4	3	88	92
18	72.9	22	3.7	113	4.0	4,324	2.4	3	86	85
19	55.4	49	8.3	108	3.8	2,344	1.3	3	72	67
20	60.5	7	1.2	33	1.2	499	0.3	3	56	72
21	33.8	4	0.1	0	0.0	0	0.0	11	45	50
22	33.8	10	1.0	0	0.0	0	0.0	13	48	56
X	127.7	141	24.0	182	6.4	4,282	2.4	3	163	164
Y	21.8	6	1.0	19	0.7	113	0.1	27	51	59
NA	5.1	0	0	134	0.0	577	0.3	0	0	0
UL	9.3	38	0	7	0.0	566	0.3	0	0	0

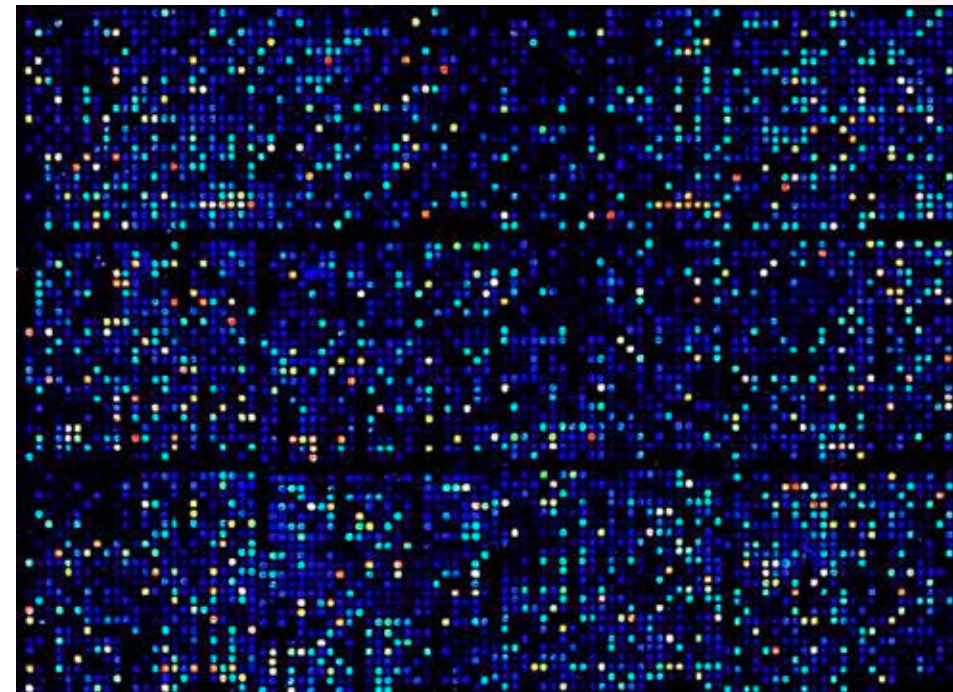
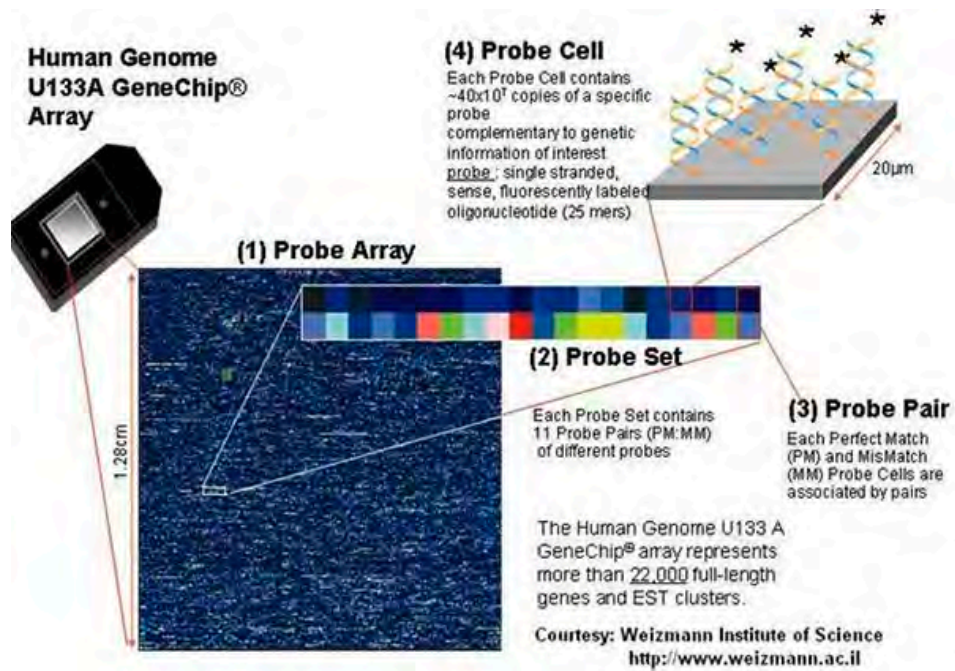


C
DN/
Non

Sec

ANOTHER FORM OF ANALOG MEASUREMENTS: ARRAYS

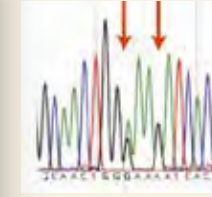
Quantification of complementary transcripts (25-120bp)



ENGINEERING NEW MEASUREMENTS

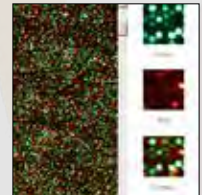
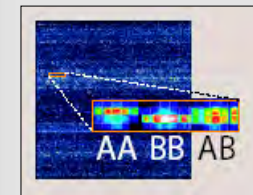
Traditional Sanger Sequencing: 1985+

- Engineering improvements (capillaries/dyes)
- Consensus of billions of molecules



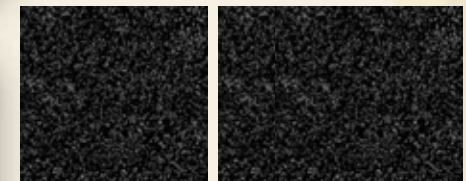
Microarrays: 1995+

- Detection of gene variation by hybridization
- Consensus of billions of molecules



Single molecule sequencing: 2007+

- Each read is a single molecule
- Millions of reads (lawn-sequencing)



Single cell sequencing: 2012+

- Full RNA-seq on single cells in solution
- Emerging spatial genomic



BEYOND TRANSCRIPT ABUNDANCE BY INTEGRATION OF DNA

Allele specific expression, non-sense mediated
decay, PSI, intron inclusion, and more

KEY PRINCIPLES THAT YOU MUST KNOW

Pseudo-single molecule reads

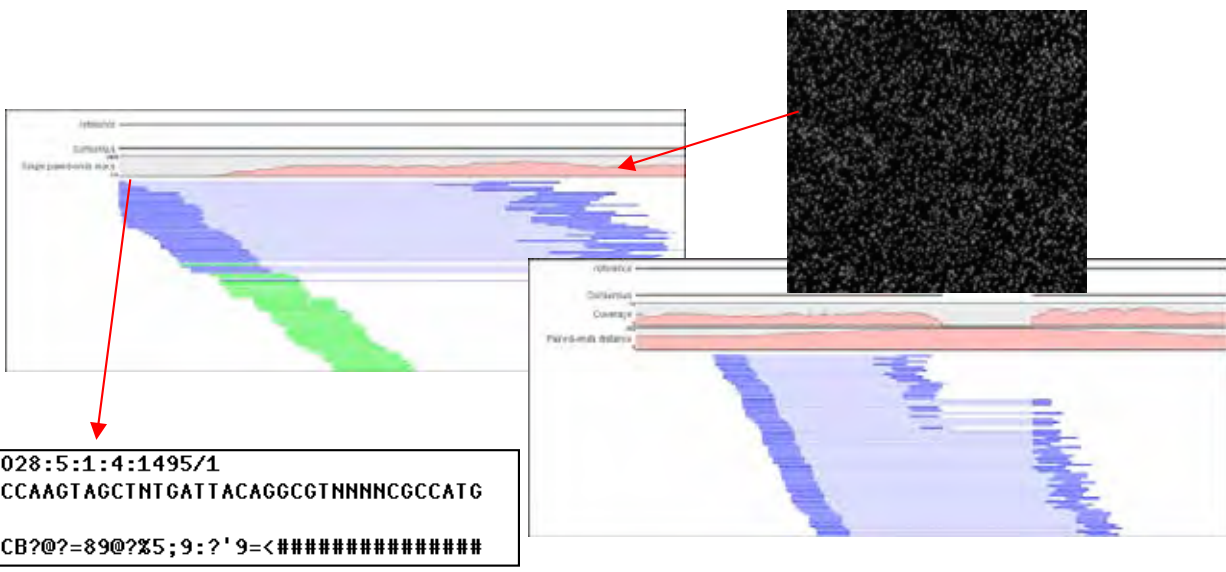
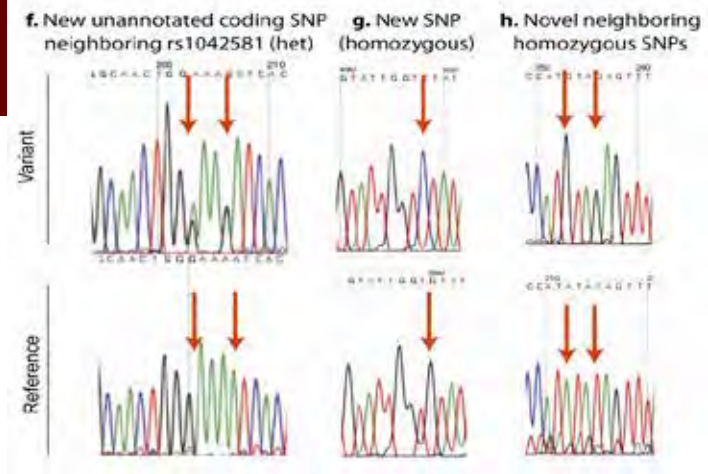
- A heterozygous SNP will give the paternal or maternal allele in a single read, not both

Paired-Reads

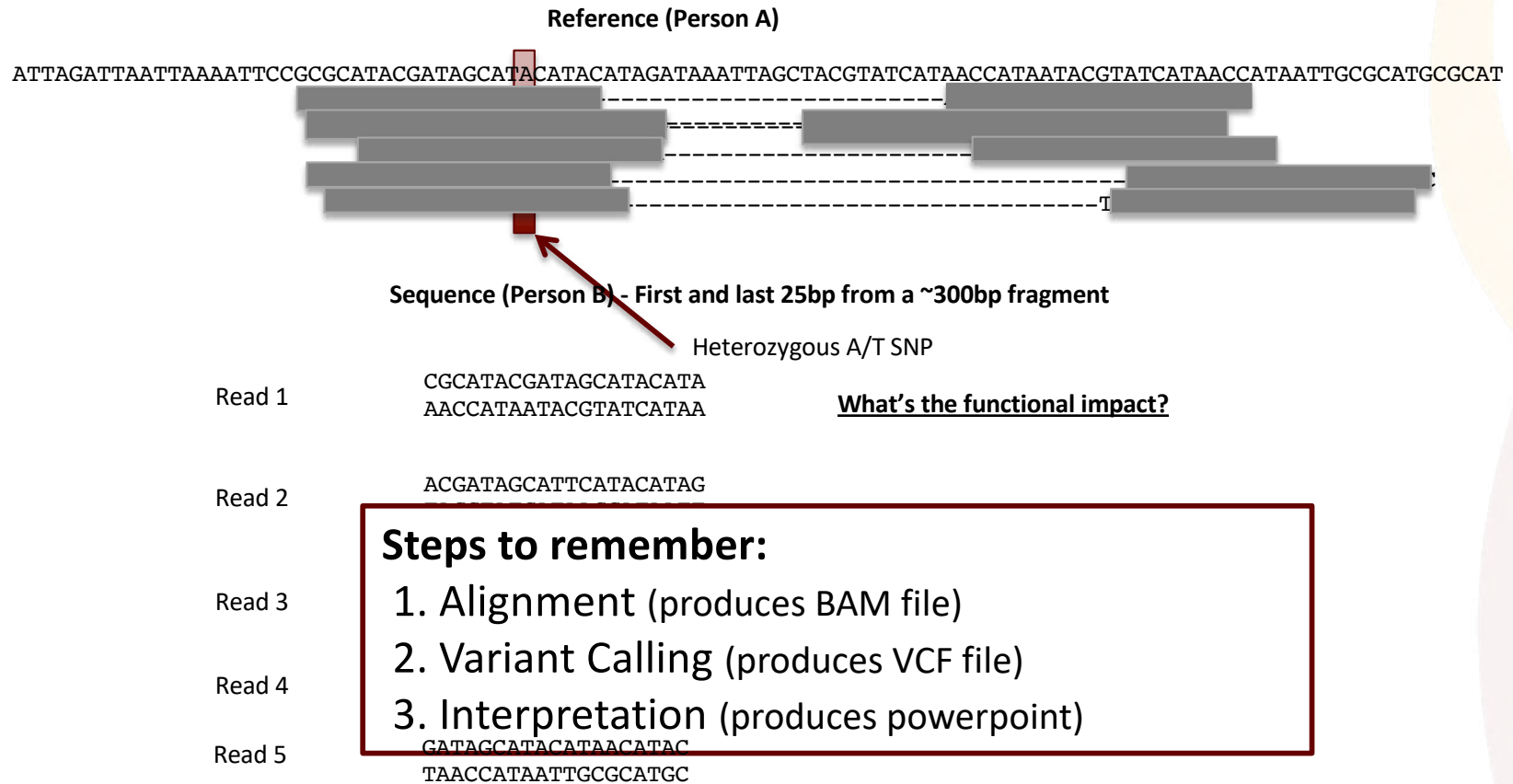
First 100 bases and last 100 bases of a ~500bp DNA molecule

Billions of reads in a sequencing run

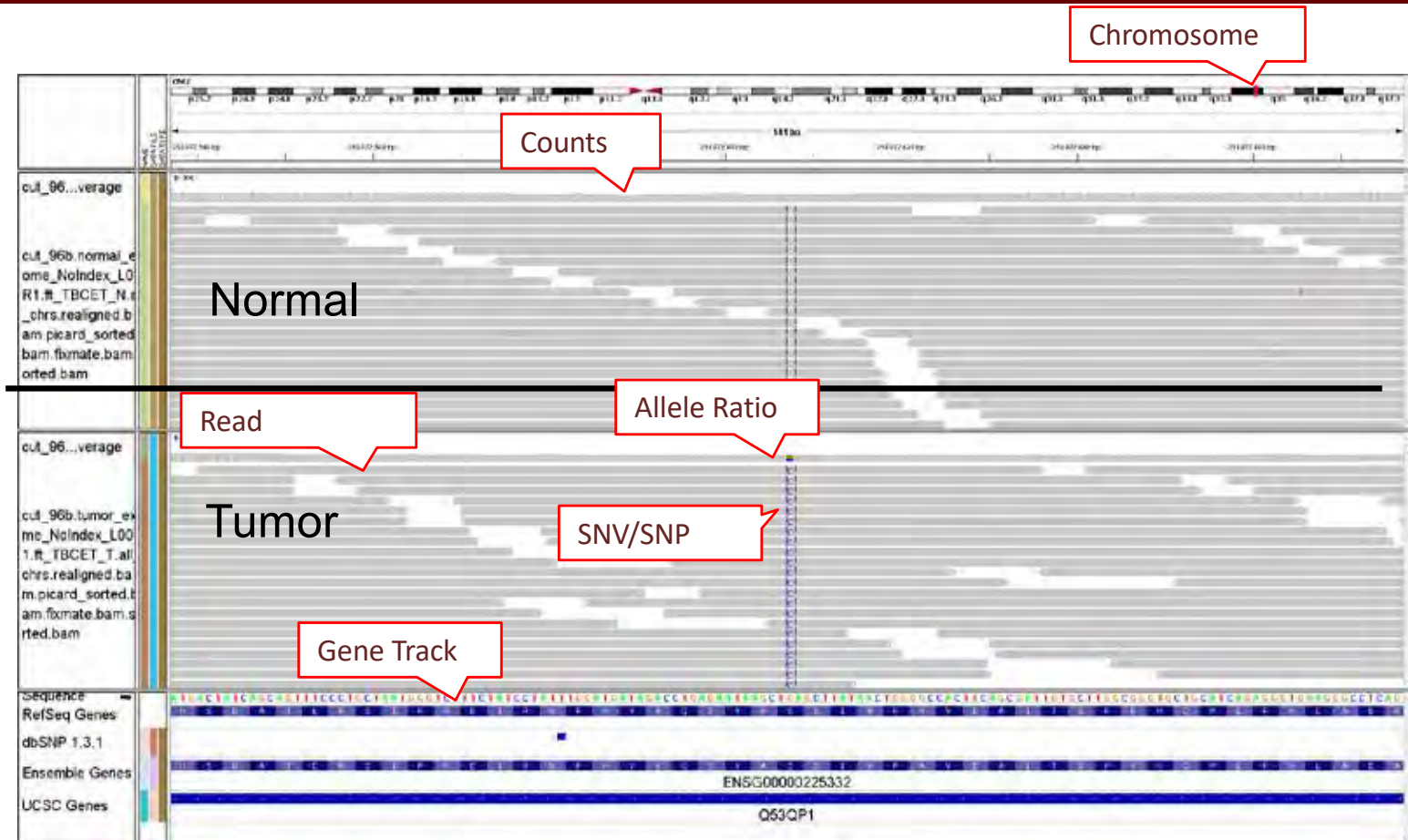
- Sampling matters and is how we control error



Concept of NGS Sequence Analysis



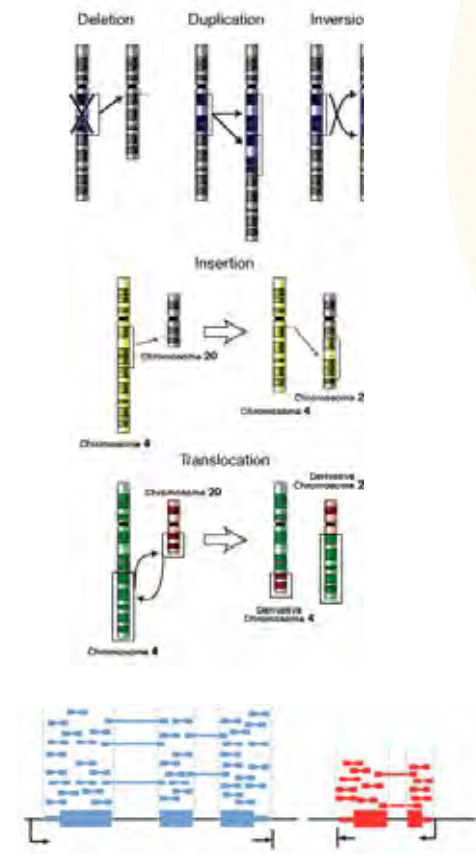
VARIANTS: EXAMPLE



NEXT GENERATION SEQUENCING

Quantum Measurement of Molecular Variation

- ┆ **Point Mutations** – Single Nucleotide Variation (SNVs) & Small Insertions/mutations (Indels)
- ┆ **Copy Number** – Changes in abundance – both DNA/RNA
- ┆ **Rearrangements** – Translocations and Structural variants via read mapping
- ┆ **Transcriptional Profiling** – Abundance, exon level, isoform level, and study splicing defects



VISUALIZING TUMOR / REFERENCE IN IGV



QUANTIFICATION/ABUNDANCE/DIFFERENTIAL EXPRESSION

Abundance:

- FPKM: Fragments Per Exon Kilobase of Sequence Per Million Reads
 - Some genes are longer than other genes and they get counted more
- TPM: For every 1,000,000 **RNA** molecules in the **RNA-seq** sample, x came from this gene/transcript
- Transcripts? Why the use of genes....

Differential Expression Should Not be done using Abundance

- Let's say you have erythrocytes higher in 1 sample, adding lots of globin
 - HBB: TPM in sample A is 600K
 - HBB: TPM in sample B is 300K.
 - Because TPM is fractional, all TPMs are lower in Sample B.
 - You need to normalize before differential expression!
 - You need the count level data to address these types of issues.

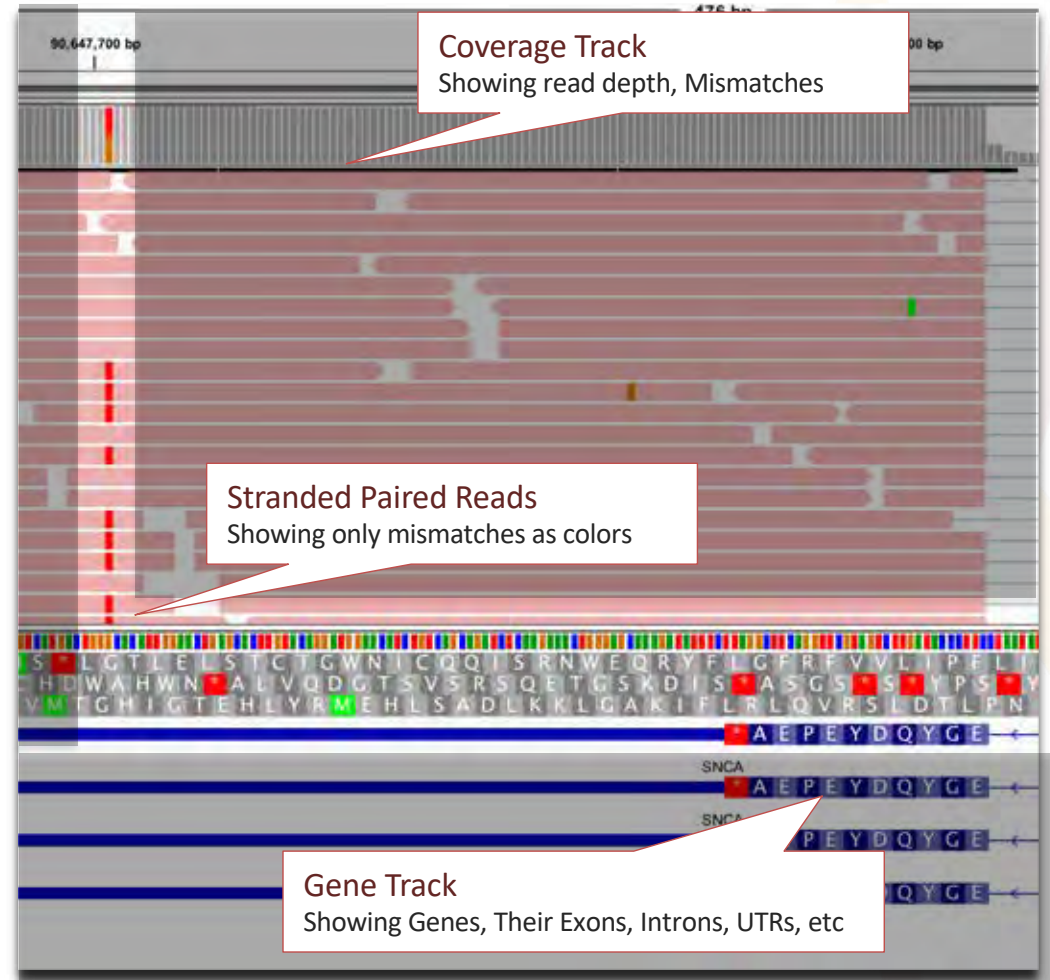
RNA



TRANSCRIPTOME: BEYOND QUANTIFICATION

- Raw Data View:

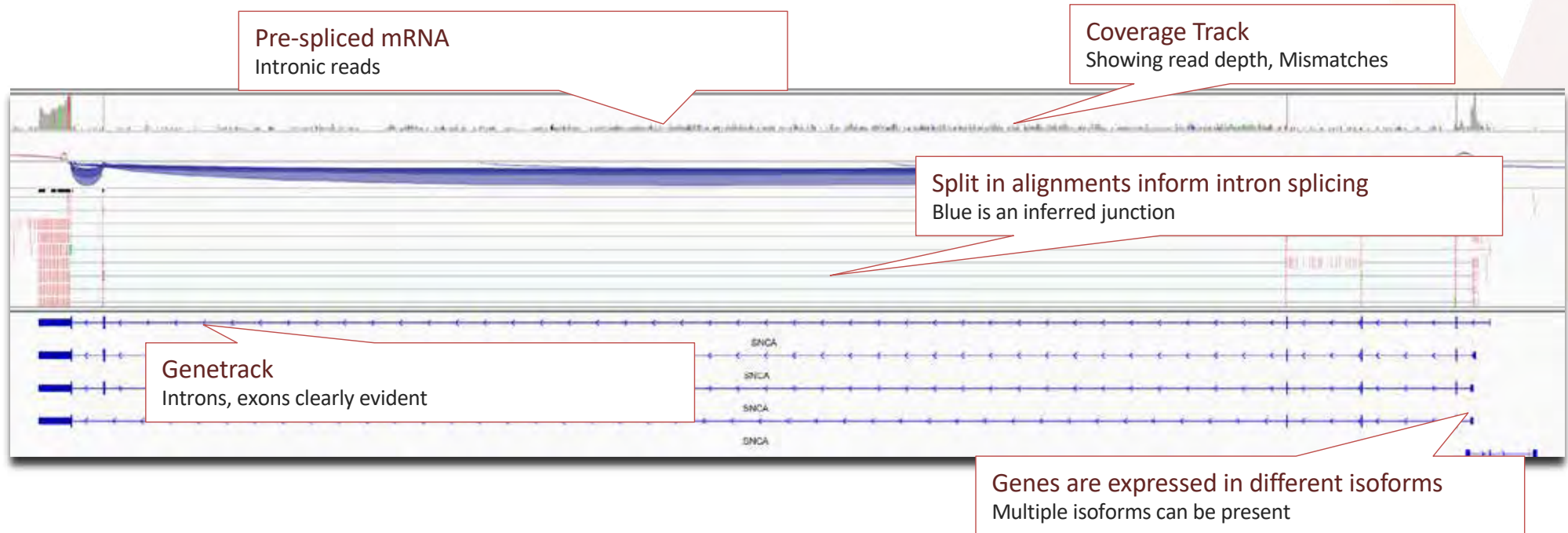
- Pre-spliced, Spliced, Strandedness
- Allele counts, etc.



TRANSCRIPTOME: BEYOND QUANTIFICATION

- Transcript Quantification (ZOOM OUT)

- Pre-spliced, Spliced, Strandedness



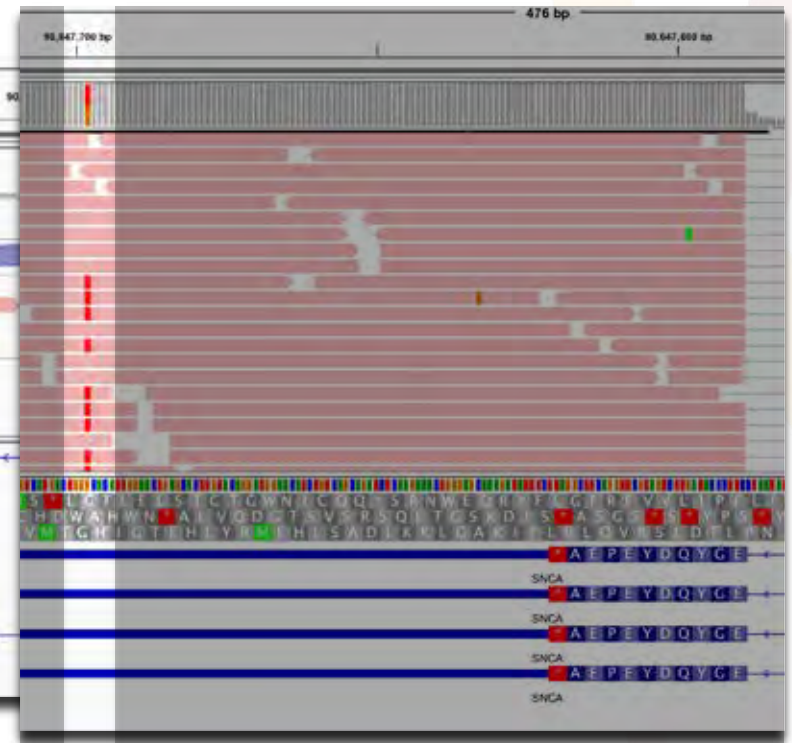
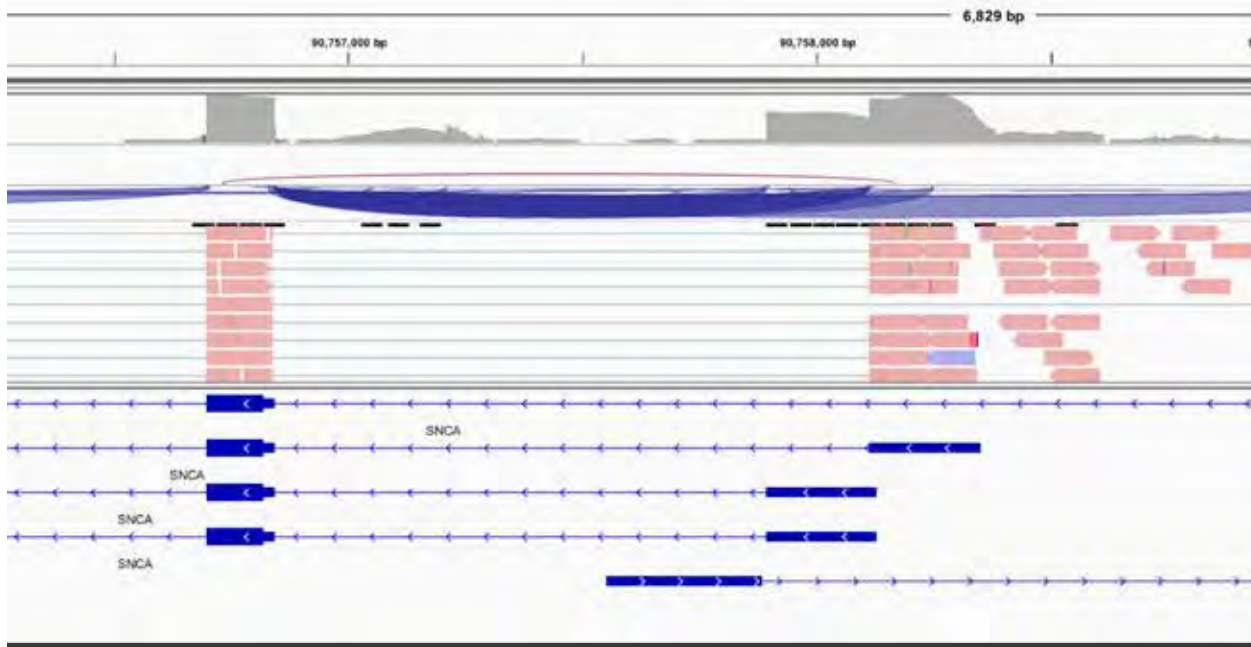
TRANSCRIPTOME: BEYOND QUANTIFICATION

- Junctions/Isoforms

- Alternative start-sites
- Integrated to individual w/ DNA

SNPs

- Allele specific expression
- Non-sense mediated decay, eQTLs



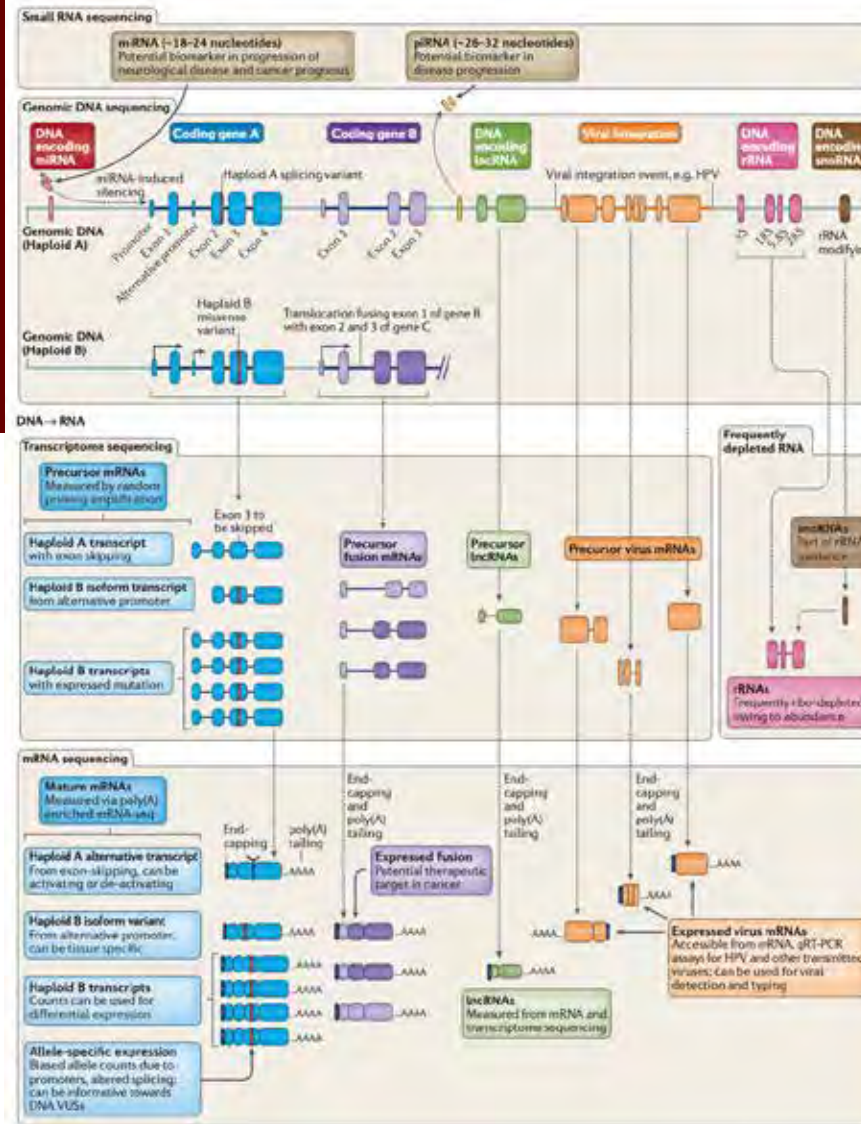
NGS VARIANTS



BEYOND THE CODING



HUMAN VARIATION AT THE RNA LEVEL

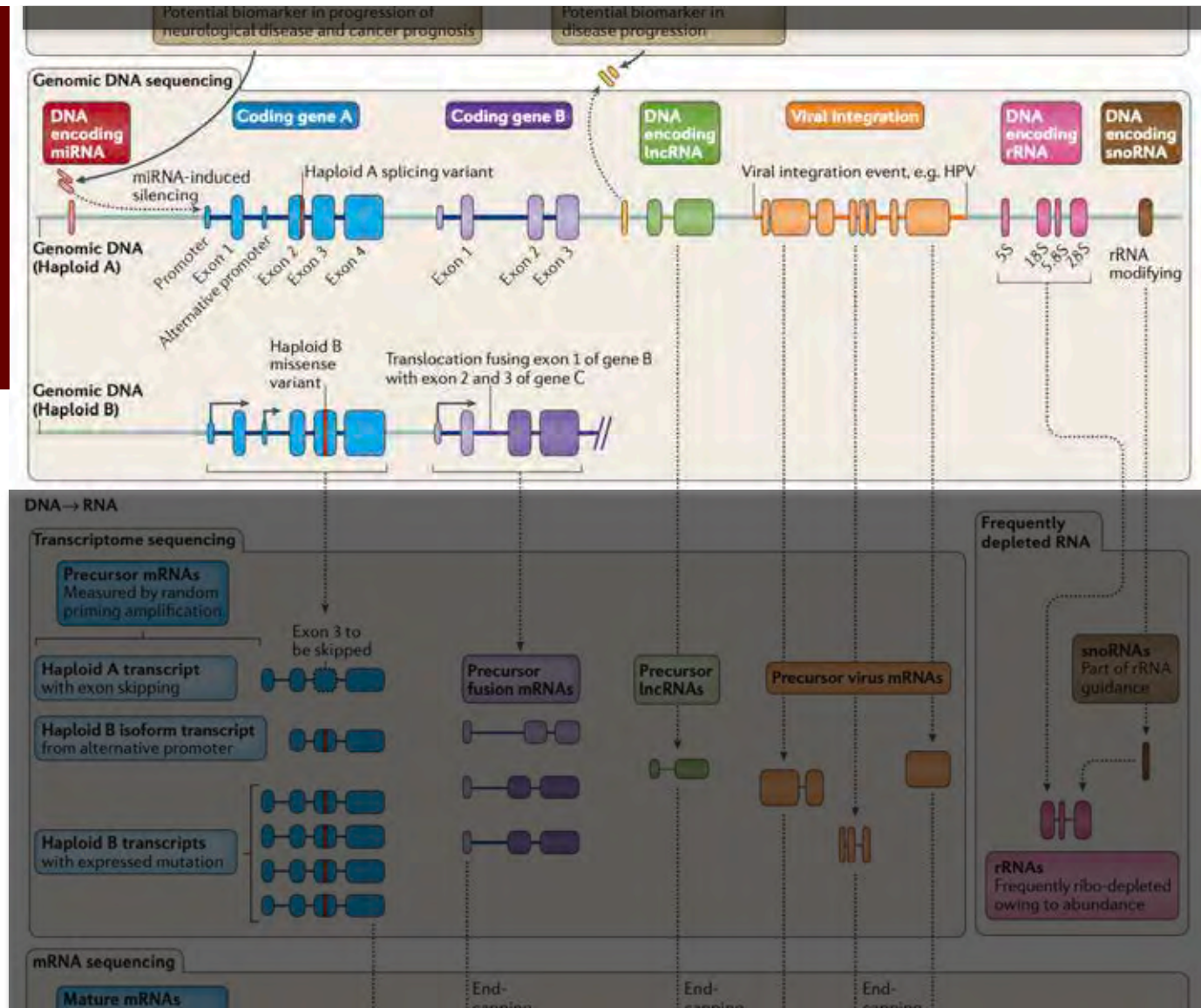


APPLICATIONS OF NEXT-GENERATION SEQUENCING

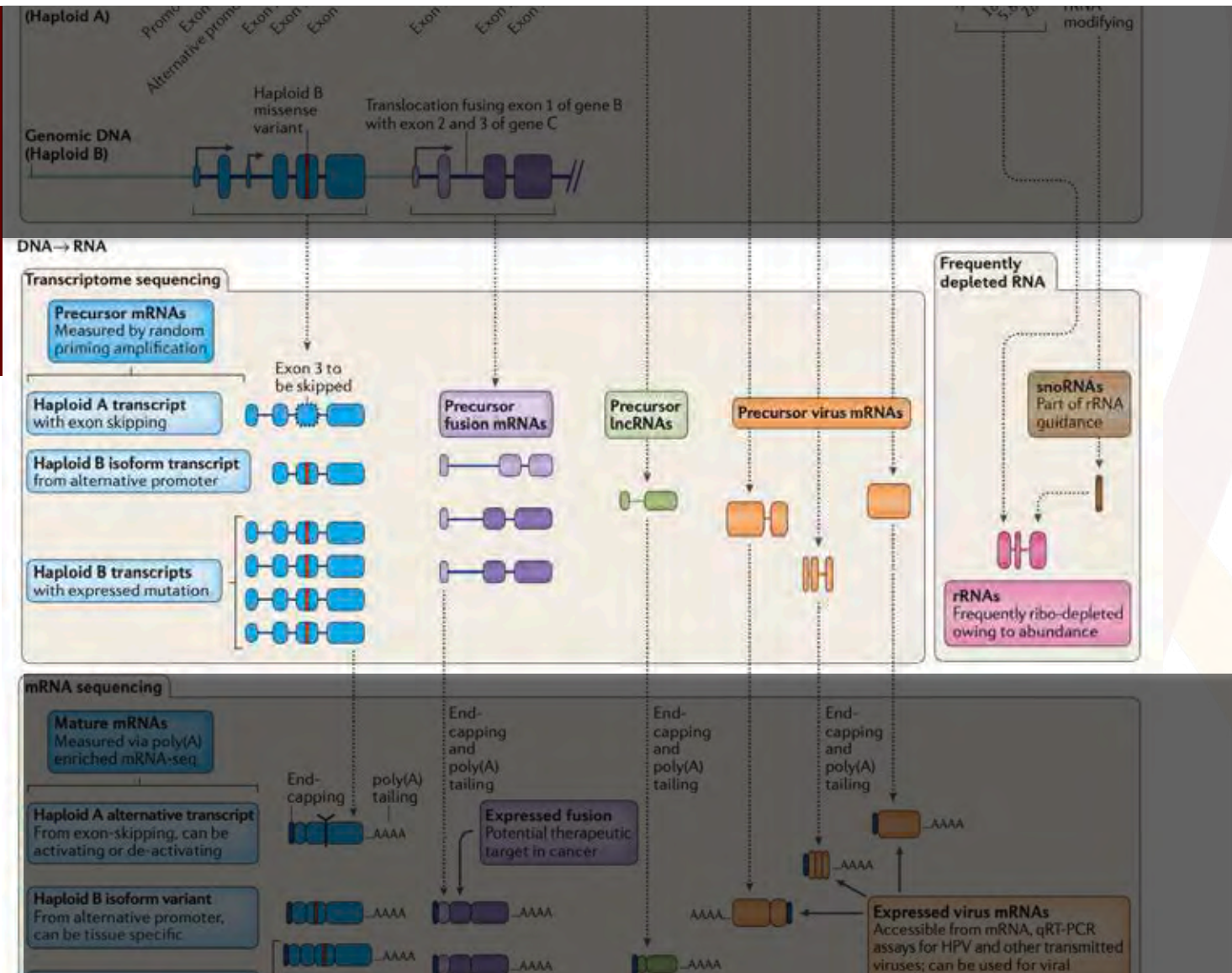
Translating RNA sequencing into clinical diagnostics: opportunities and challenges

Sara A. Blythe¹, Kenneth K. Van Nieuwenhove¹, David An. Engelmann², John P. Gujral¹ and David W. Craig¹

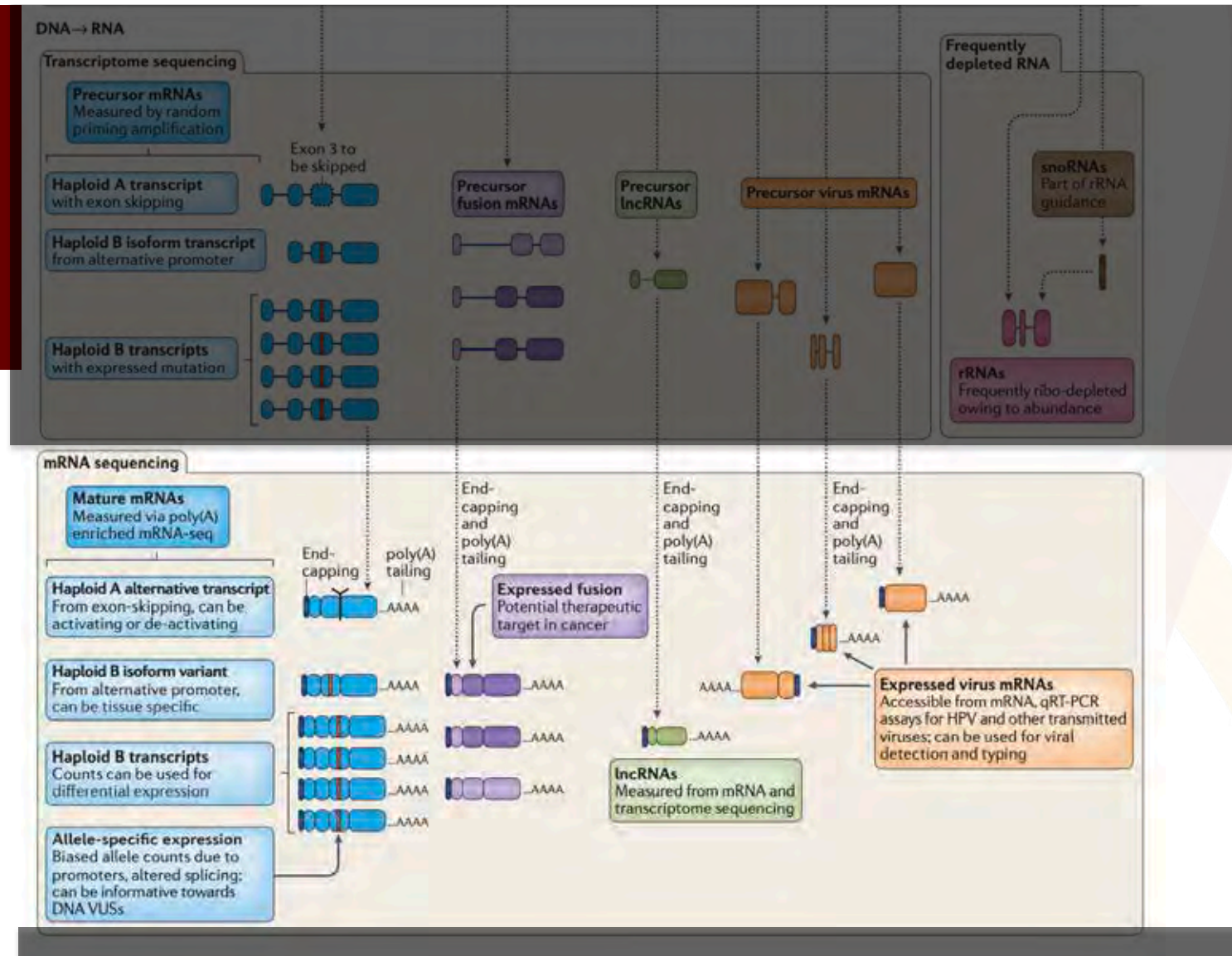
HUMAN VARIATION AT THE RNA LEVEL



HUMAN VARIATION AT THE RNA LEVEL



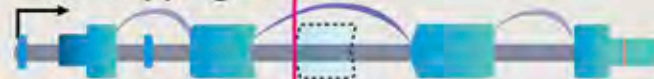
HUMAN VARIATION AT THE RNA LEVEL



TRANSCRIPTOME: SNPs

Ex. RNA Variant Types

Exon Skipping Variants



Alternative 5' Splicing Variants



Alternative 3' Splicing Variants



Alternative Promoter Variants



Intron Retention Variants



Alternative Terminator Variants



RNA Editing Events



INTRON INCLUSION

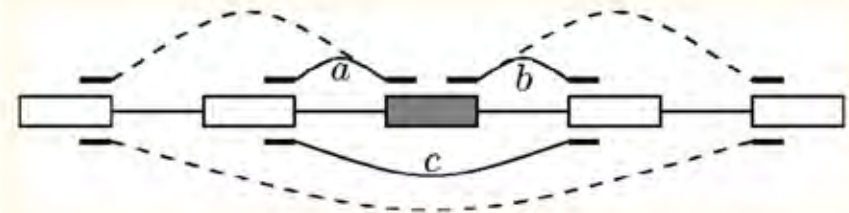
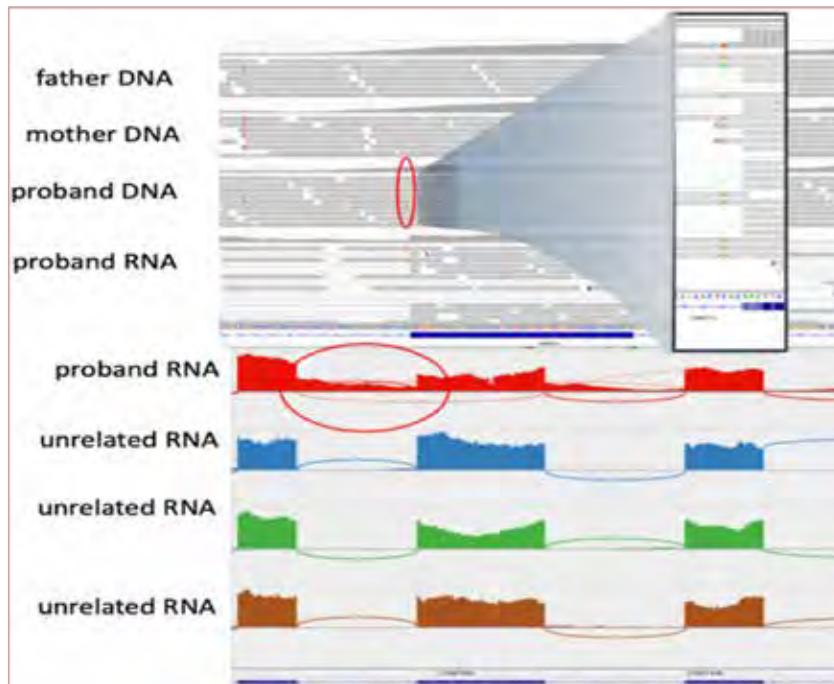


Fig 1.

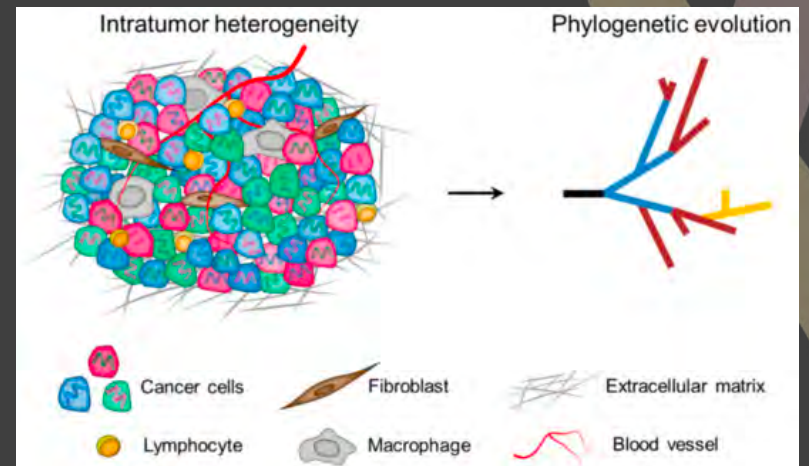
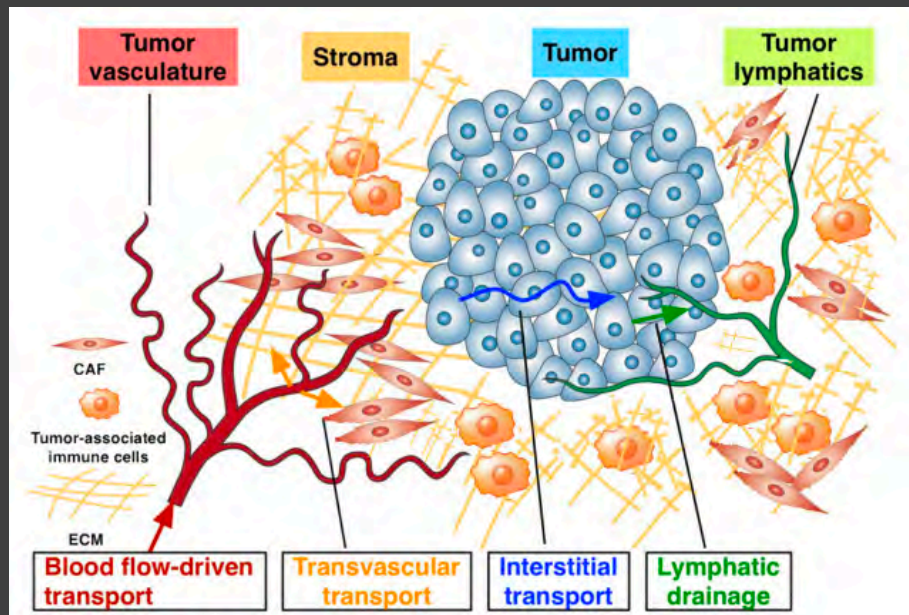
The percent-spliced-in (PSI, Ψ) metric is defined as the number of reads supporting exon inclusion ($a + b$) as the fraction of the combined number of reads supporting inclusion and exclusion (c). The exon of interest is shown in gray. Only reads that span to the adjacent exons (solid arcs) account for Equation (1)

$$\Psi = \frac{a + b}{a + b + 2c}$$

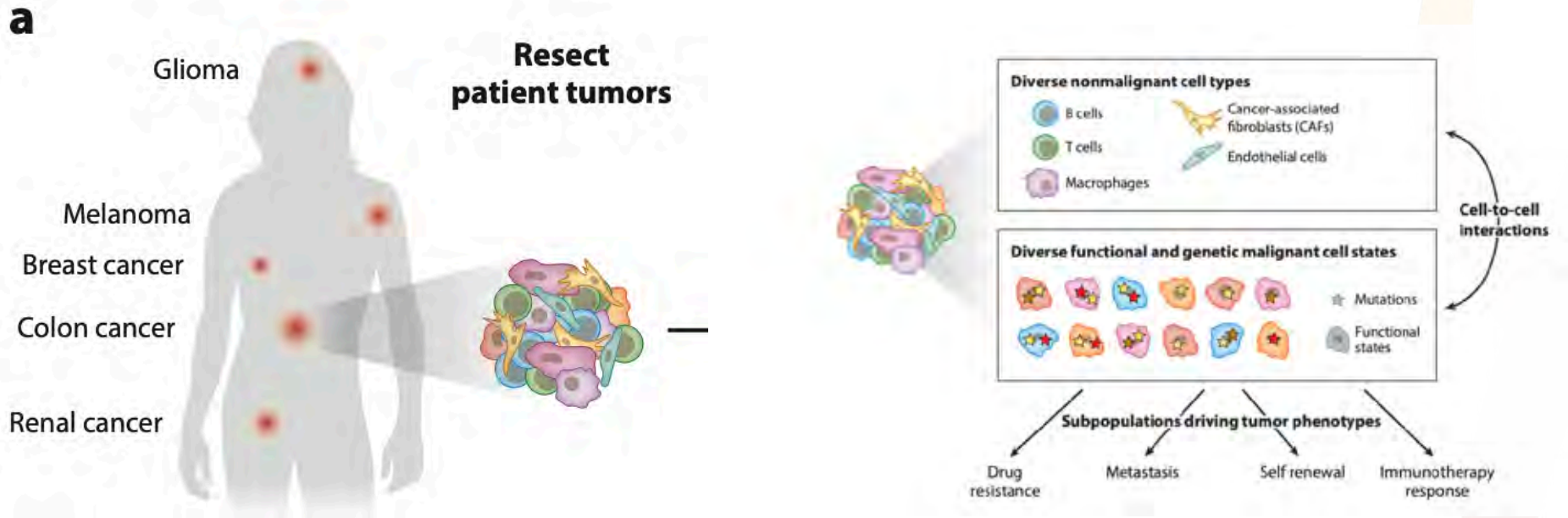
Bioinformatics. 2013 Jan 15; 29(2): 273–274



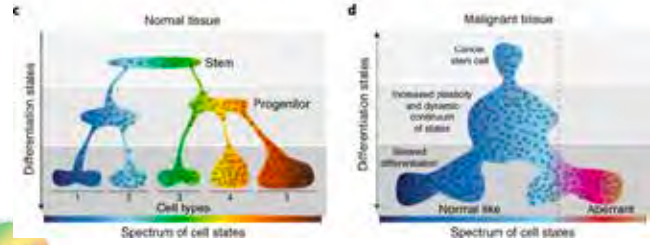
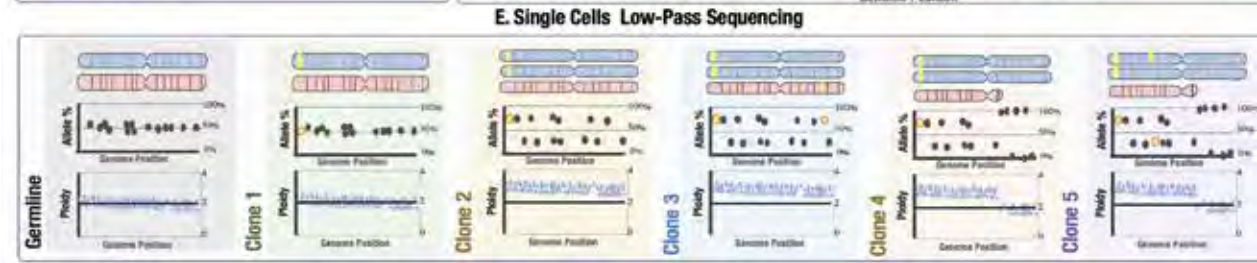
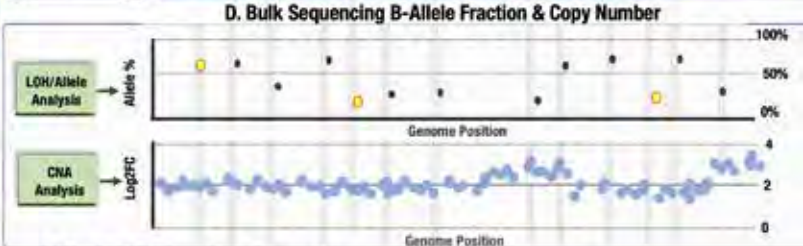
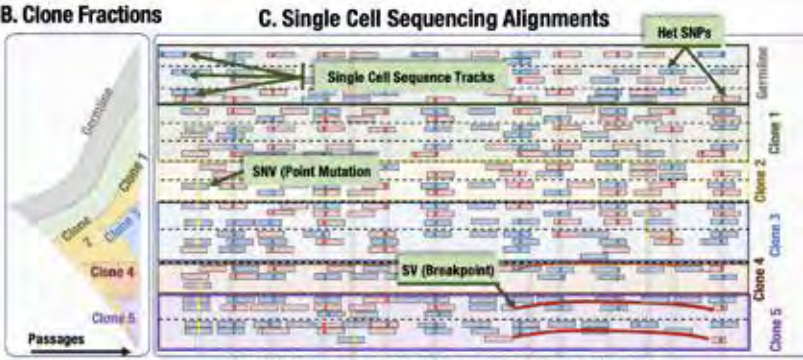
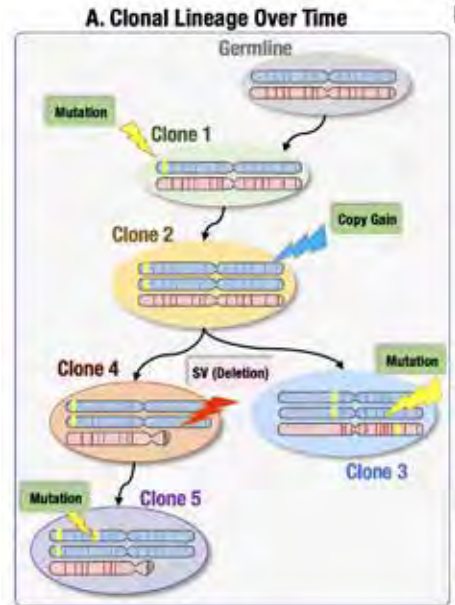
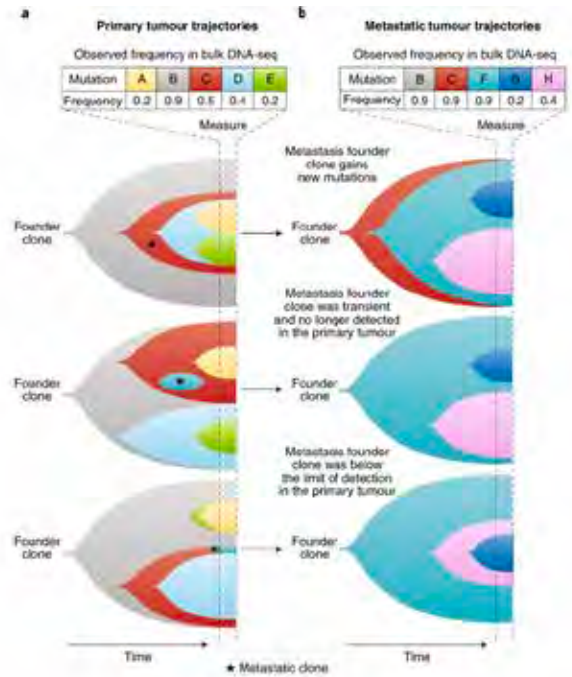
CANCER TRANSCRIPTOMICS



TUMOR BIOLOGY IS DRIVEN BY SUB-POPULATIONS AND HETEROGENEITY



CANCER IS A MIXTURE OF TUMOR AND HEALTHY CELLS



David Wesley Craig (davidwcr@usc.edu)

USC Translational Genomics

EASY EXAMPLE: FUSIONS

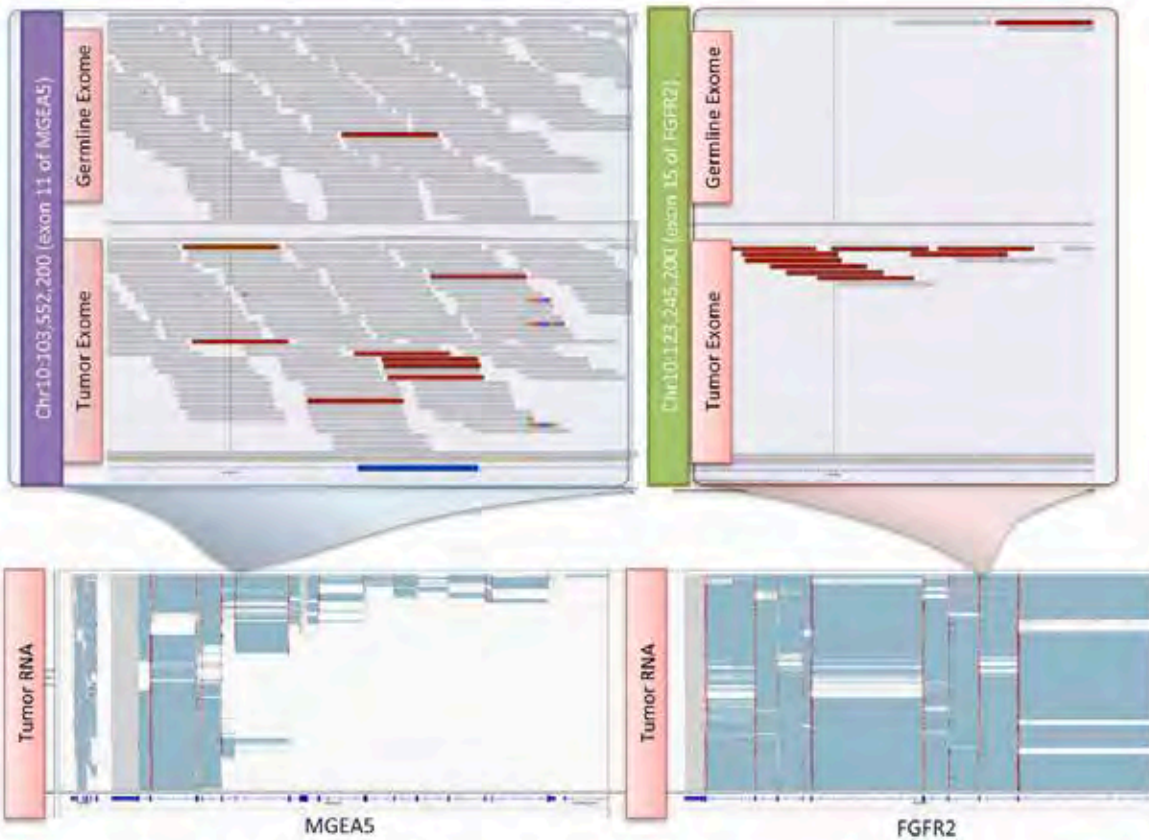


Figure 4. Visualisation of FGFR2-MGEA5 fusion in the Integrative Genomics Viewer (IGV).

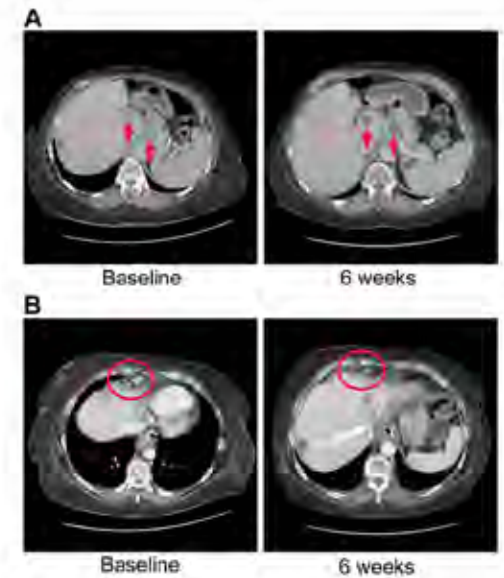


Figure 7. Anti-tumor activity in Patient 4 harboring an *FGFR2-MGEA5* fusion, to FGFR inhibitors. **A)** CT images of patient 4, whose tumor possessed an *FGFR2-MGEA5* fusion, at baseline and 6 weeks demonstrate central necrosis of a caudate liver lobe mass (left arrow), 2.6 cm at baseline and 6 weeks, and shrinkage of a metastatic supraceliac axis lymph node (right arrow), 3.1 cm and 2.9 cm at baseline and 6 weeks respectively. **B)** CT images of patient 4 showing shrinkage of metastatic lymph nodes involving the right cardiophrenic angle (red circles), 1.3 cm and 0.5 cm at baseline and 6 weeks respectively.
doi:10.1371/journal.pgen.1004135.g007

OPEN Clinical Implementation of Integrated Genomic Profiling in Patients with Advanced Cancers

Mitch J. Borad^{1,2}, Jan B. Egan¹, Rachel M. Condeelis¹, Winnie S. Liang¹, Rafael Fonseca^{1,3}, Nicole R. Rizucci⁴, Ann L. McCallough¹, Michael T. Barrett^{1,5}, Katherine S. Hunt⁶, Mia D. Chempman^{1,7}, Malray D. Patel¹, Scott W. Young¹, Alvin C. Silva¹, Thai H. Ho^{1,8,9}, Thorvardur R. Halldorsson^{10,11}, Robert R. McWilliams¹², Konstantinos N. Lazaridis¹³, Ramesh K. Ramasathan¹⁴, Angela Baker¹⁵, Jessica Albrich¹⁶, Ahmet Kurdoglu¹⁷, Tyler Izatt¹⁸, Alexis Christofides¹⁹, Irma Cherni²⁰, Sara Nasser²¹, Rebecca Remar²², Lori Conroy²³, Gregory McDonald²⁴, Jonathan Adams²⁵, Stephen D. Manley²⁶, Marcario Valdez²⁷, Dawn E. Aniszewski²⁸, Daniel D. Von Hoff²⁹, David W. Coig³⁰, A. Keith Stewart^{31,32}, John D. Carpenter³³ & Alan H. Bryce³⁴

EXAMPLE CHOLANGIOCARCINOMA

35 somatic coding mutations

- ↳ Two COSMIC genes
- ↳ None in known commercial cancer panels
- ↳ One flagged inferred therapeutic context

One focal copy number event on chr3

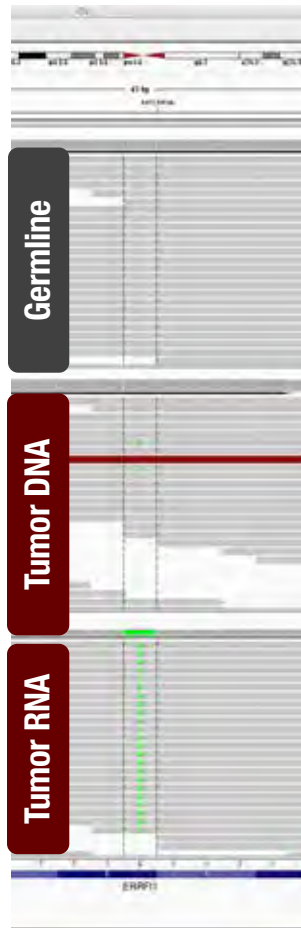
- ↳ None flagged inferred therapeutic context

RNA-seq data Differential Expression

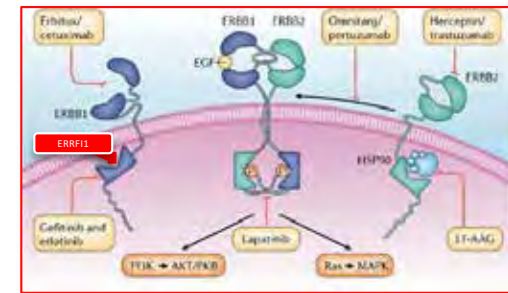
- ↳ None flagged inferred therapeutic context

Structural Variants and Fusions

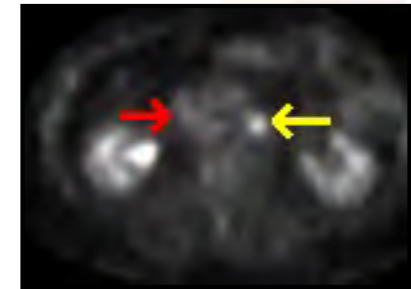
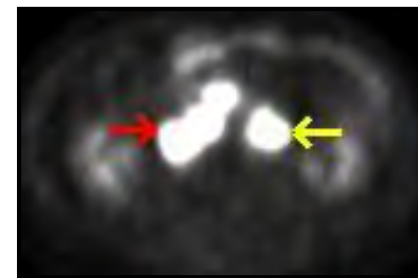
- ↳ None flagged inferred therapeutic context



ERFF1 – Context for EGFR Activation



↳ EGFR inhibitor was recommended By Tumor Board



↳ Previously bulky much smaller and much less metabolically active at 10/25/12 on PET/CT.

↳ >60% reduction of retroperitoneal left periaortic node from 08/06/12 to 10/25/12 on PET scan.

SAMPLE PREP TYPES

Library Preparation

TABLE 3.1 RNA-seq Library Protocols

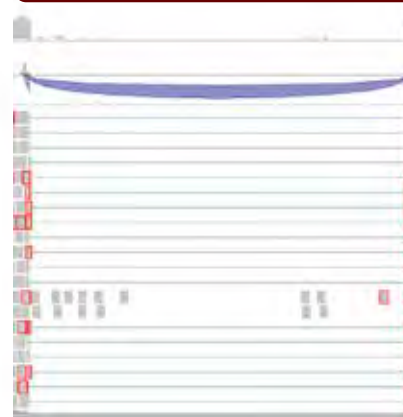
Library Design	Pros	Cons
Ribo-depletion	Removes ribosomal RNA from sample and provides access to mRNA and non-coding RNA	Pre-mRNA retained providing reads for junctions
Poly-A selection	Removes all non-polyadenylated RNA. Enriches for mRNA	Limited quantification of non-coding RNA
Treatment with double-stranded nuclease	Reduces expression level of highly abundant transcripts	Not a well-utilized protocol; biases/impact on global expression not characterized
MicroRNA capture	Provides access to microRNA	MicroRNA recovered depends on capture technique
mRNA capture	Provides selective representation of specific transcripts and potentially unknown fusion partners	Potentially biases quantification, reduces representation of untargeted genes, which may be of interest
Multiplexed using genetic barcoding	Many samples processed together	Some read pairs might be erroneously swapped if run cluster density is too high

FFPE Considerations

Fresh/Frozen/Cryo

A fraction of types of assays

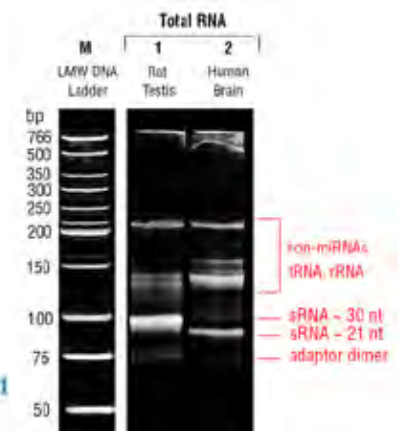
Poly-A



Random Priming



Size Selection – smallRNA



Cancer Genomics. DOI: <http://dx.doi.org/10.101>
© 2014 Elsevier Inc. All rights reserved.

EXAMPLE: CERVICAL CARCINOMA

Integration of HPV (Detailed)

HPV18 DNA is evident as is expression of RNA



Genomic basis for RNA alterations in cancer

<https://doi.org/10.1038/s41586-020-1970-0>

Received: 29 March 2018

Accepted: 11 December 2019

Published online: 5 February 2020

Open access

PCAWG Transcriptome Core Group^{1,2,3}, Claudia Calabrese^{2,3,9}, Natalie R. Davidson^{2,4,5,6,7,35}, Deniz Demircioğlu^{8,9,35}, Nuno A. Fonseca^{2,3,9}, Yao He^{10,35}, André Kahles^{4,6,8,35}, Kiong-Yan Lehmann^{1,4,6,7,35}, Fenglin Liu^{10,35}, Yuichi Shiraiishi^{11,35}, Cameron M. Soulette^{12,35}, Lara Urban^{2,3}, Liliana Greger⁷, Silang Li^{13,14}, Dongbing Liu^{13,14}, Marc D. Perry^{15,36}, Qian Xiang¹⁵, Fan Zhang¹⁶, Junjun Zhang¹⁶, Peter Bailey¹⁷, Serap Erkek¹⁸, Katherine A. Hoadley¹⁹, Yong Hou^{13,14}, Matthew R. Huska²⁰, Helena Kilpinen²¹, Jan O. Korbel¹⁹, Maximilian G. Marin²¹, Julia Markowski²², Tannistha Nandi⁹, Qiang Pan-Hammarström^{13,22}, Chandra Sekhar Pedamallu^{23,28,29}, Reiner Siebert²⁴, Stefan G. Stark^{2,4,6,7}, Hong Su^{30,34}, Patrick Tan²⁵, Sebastian M. Waszak²⁶, Christina Yung²⁶, Shida Zhu^{31,34}, Philip Awadalla^{15,26}, Chad J. Creighton²⁷, Matthew Meyerson^{23,28,29}, B. F. Francis Ouellette³², Kui Wu^{33,34}, Huanming Yang³¹, PCAWG Transcriptome Working Group¹, Alvis Brazma^{2,3,6*}, Angela N. Brooks^{12,23,28,30*}, Jonathan Göke^{8,21,36}, Gunnar Rätsch^{4,6,8,12,36*}, Roland F. Schwarz^{2,30,32,33,36}, Oliver Stegle^{2,33,33,36}, Zemin Zhang^{10,36} & PCAWG Consortium³⁴

Transcript alterations often result from somatic changes in cancer genomes¹. Various forms of RNA alterations have been described in cancer, including overexpression², altered splicing³ and gene fusions⁴; however, it is difficult to attribute these to underlying genomic changes owing to heterogeneity among patients and tumour types, and the relatively small cohorts of patients for whom samples have been analysed by both transcriptome and whole-genome sequencing. Here we present, to our knowledge, the most comprehensive catalogue of cancer-associated gene alterations to date, obtained by characterizing tumour transcriptomes from 1,188 donors of the Pan-Cancer Analysis of Whole Genomes (PCAWG) Consortium of the International Cancer Genome Consortium (ICGC) and The Cancer Genome Atlas (TCGA)⁵. Using matched whole-genome sequencing data, we associated several categories of RNA alterations with germline and somatic DNA alterations, and identified probable genetic mechanisms. Somatic copy-number alterations were the major drivers of variations in total gene and allele-specific expression. We identified 649 associations of somatic single-nucleotide variants with gene expression in *cis*, of which 68.4% involved associations with flanking non-coding regions of the gene. We found 1,900 splicing alterations associated with somatic mutations, including the formation of exons within introns in proximity to Alu elements. In addition, 82% of gene fusions were associated with structural variants, including 75 of a new class, termed ‘bridged’ fusions, in which a third genomic location bridges two genes. We observed transcriptomic alteration signatures that differ between cancer types and have associations with variations in DNA mutational signatures. This compendium of RNA alterations in the genomic context provides a rich resource for identifying genes and mechanisms that are functionally implicated in cancer.

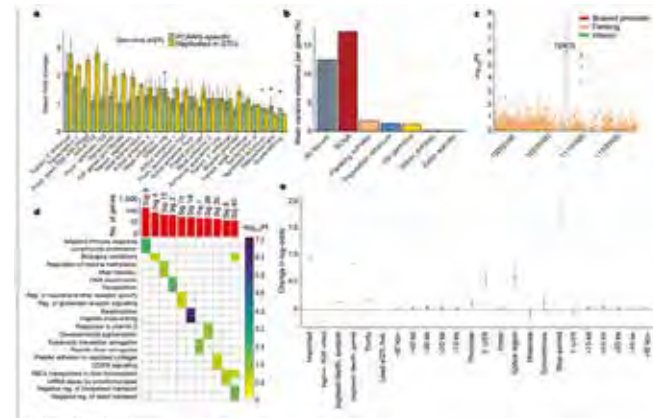


Fig. 1 | Position-specific effect of somatic mutations on alternative splicing. **a**, Top, proportion of mutations near exon–intron junctions and at branch sites that are associated with exon-skipping events. Mutations with associated splicing changes are those in which the percentages (plotted in derived (z-score) ≥ 3 (dark blue). Asterisks denote intron positions significantly enriched for splicing changes relative to background based on a permutation test. * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$. Bottom, sequence motifs of regions. **b**, Example of an

exonization event in the tumour-suppressor gene *5TK11*. The RNA-seq read coverage for a part of the gene is shown in red for a donor carrying the alternative allele, and in grey for a random donor with reference allele. The cassette exon event is shown as a schematic below. **c**, Enrichment of SINE elements in SAVs compared to sequence background (BG). Shown for SINE elements overlapping in sense (middle) and antisense (right) directions.

PCAWG

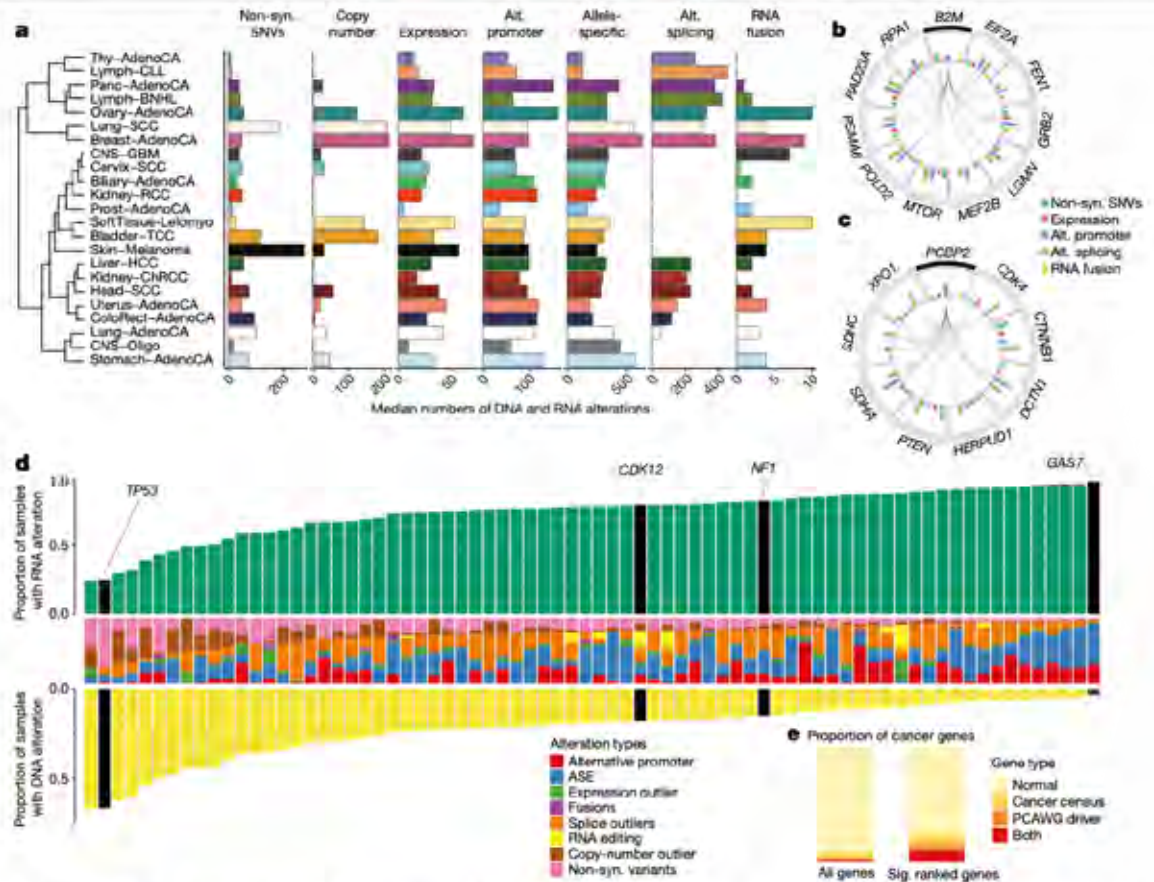
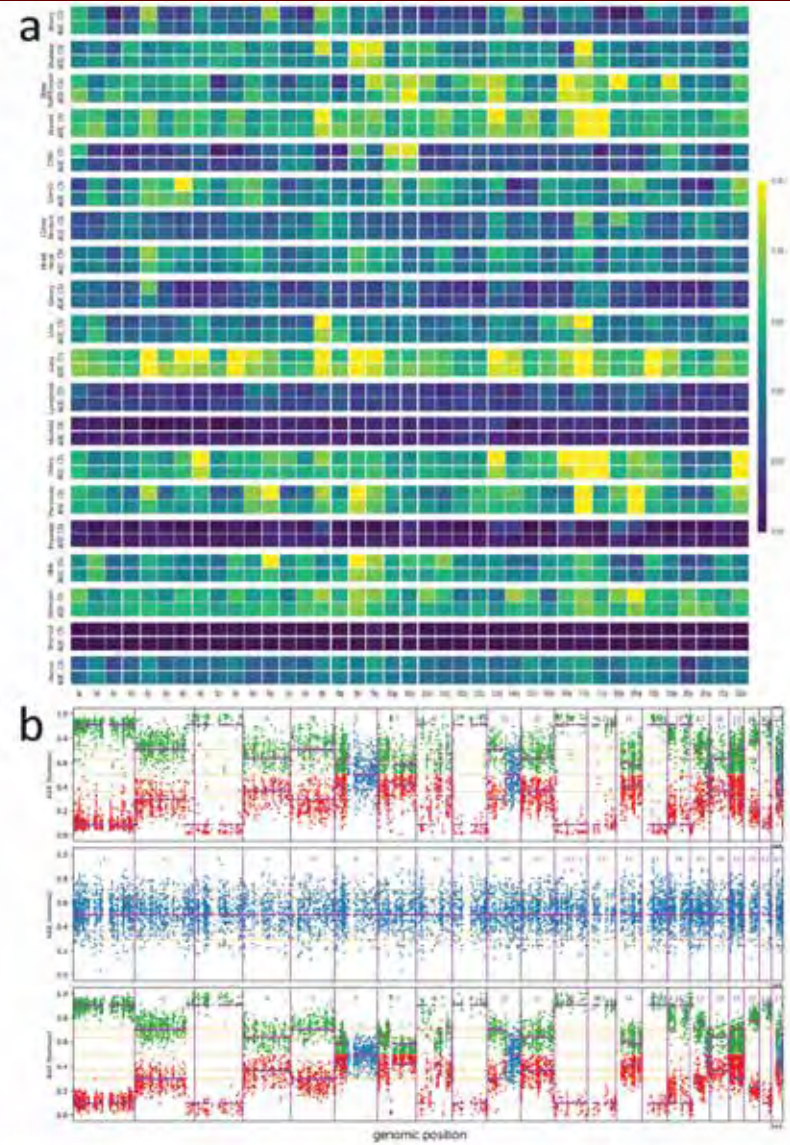
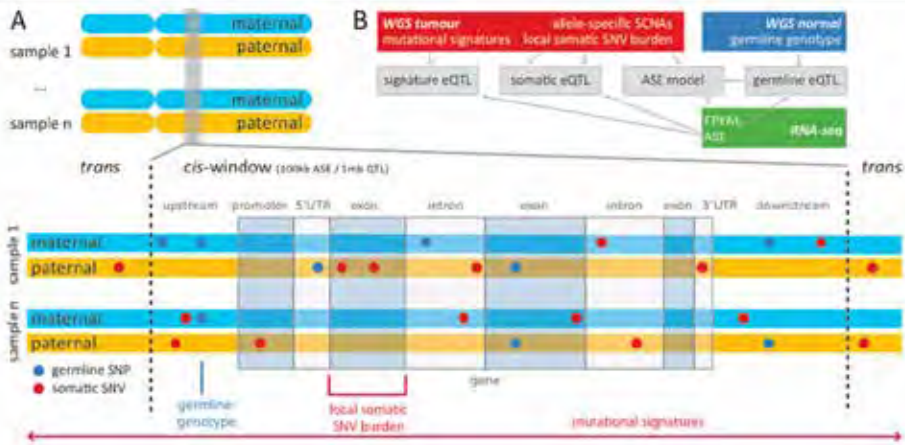


Fig. 4 | Global view of DNA and RNA alterations that affect tumours. a, The median numbers of different alterations across histotypes. Histotypes are ordered by hierarchical clustering based on the pattern of different types of alteration. Only histotypes with more than 10 donors are shown. Alt., alternative; non-syn, non-synonymous. Cancer-type abbreviations are listed in Supplementary Table 23. **b, c,** Circular representations of the selected genes significantly co-occurred with *B2M* (**b**) and *PCBP2* (**c**). Connecting lines indicate the specific types of co-occurrence of alteration pairs. The inner histograms

indicate the frequencies of incidences of different alteration types shown in different colours. **d,** All 74 Catalogue of Somatic Mutations in Cancer (COSMIC) cancer census genes or PCAWG driver genes that are both frequently and heterogeneously altered across both RNA- and DNA-level alterations. Yellow bars indicate the proportion of samples that had DNA-level alterations, and green bars indicate the proportion of samples with RNA-level alterations. Middle column is the proportion of each alteration type observed for that gene. **e,** The enrichment of cancer genes within our list of significantly recurrent genes.

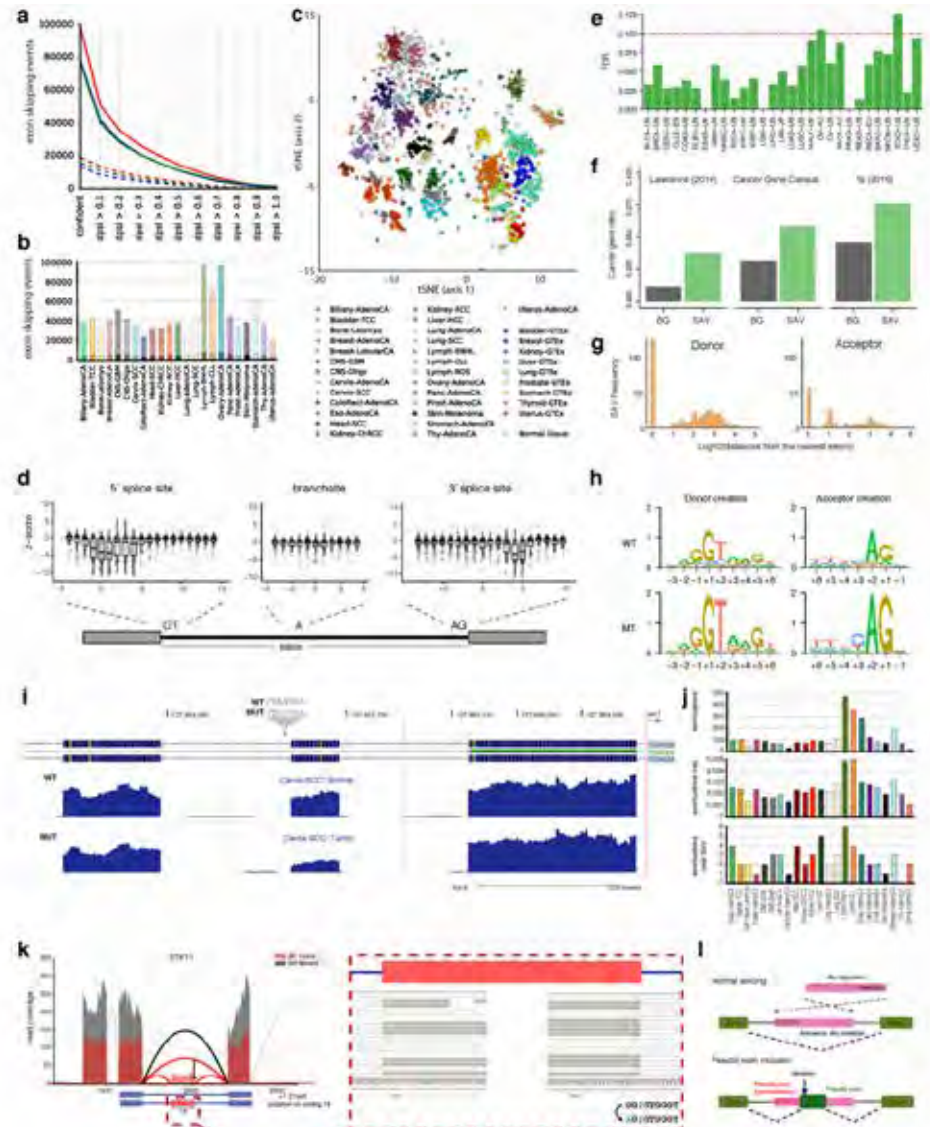




David Wesley Craig (c

Genomics

ALT SPLICING IN PCAWG



David Wesley Craig (dwcraig@usc.edu)

USC Translational Genomics



BIOMARKERS

Table 1 | Selected examples of current RNA-based clinical tests

RNA biomolecule	Method	Examples	Use
Viral RNA	qRT-PCR	<ul style="list-style-type: none"> Influenza virus¹⁰ Dengue virus¹¹ HIV¹² Ebola virus¹³ 	Viral detection and typing
miRNA	qRT-PCR	<ul style="list-style-type: none"> AlloMap (CareDx; heart transplant)^{14,15} Cancer Type ID (BioTherapeutics)¹⁶ 	Diagnosis
	Microarray	Affirma Thyroid Nodule Assessment (Veracyte) ¹⁷	Diagnosis
	qRT-PCR	<ul style="list-style-type: none"> OncotypeDx (Genome Health; breast, prostate and colon cancer)^{18,19} Breast Cancer Index (BioTherapeutics)²⁰ Prolaris (Myriad; prostate cancer)²¹ 	Prognosis
	Digital barcoded mRNA analysis	Prosigna Breast Cancer Prognostic Gene Signature (Nanostring) ²²	Prognosis
miRNA	Microarray	<ul style="list-style-type: none"> MammaPrint (Agendia; breast cancer)²³ ColoPrint (Agendia; colon cancer)²⁴ Decipher (Genome Dx; prostate cancer)²⁵ 	Prognosis
	Fusion transcript	<ul style="list-style-type: none"> qRT-PCR: AML (RUNX1-RUNX1T1)²⁶ qRT-PCR: BCR-ABL1 (REF-21) 	Diagnosis Monitoring molecular response during therapy
miRNA	qRT-PCR (exosomal RNA)	ExoDx Lung (AUG; Exosome Dx) ²⁷	Fusion detection
	RNA-seq	FoundationOne Home ²⁸	Fusion detection

Analytical validity

Accuracy and reliability of a test to measure a specific biomarker

Clinical validity

The accuracy of how well a test detects or predicts clinical diagnosis or outcome

Clinical utility

The likelihood the test is to inform clinical decisions and improve outcome

Analytical sensitivity

How often is the test positive when the biomarker is present?

Analytical specificity

How often is the test negative when the biomarker is not present?

Robustness

Repeatability and reproducibility of the assay within and across laboratories.

Limits of detection

Lowest level of reliable detection of transcripts.

Stability

Collection, handling, transport of sample and impact on robustness.

Gold standards

Reference sets for assessing sensitivity and specificity.

Clinical sensitivity

How often is the test positive in patients with the disease or clinical outcome?

Clinical specificity

How often is the test negative in patients without the disease or clinical outcome?

Prevalence

The proportion of individuals that will have a disease or outcome.

Positive predictive value

Given prevalence, the probability that subjects with a positive test result for a disorder or outcome will have the disease or outcome.

Negative predictive value

For negative tests, the probability that subjects truly will not have the disease or outcome.

Penetrance

The proportion of subjects with the biomarker that have the predicted outcome or diagnosis.

Appropriate intervention

Assessment of test impact on patient care, publishing of clinical trials.

Quality assurance

Quality control measures for tests, reagents and/or facilities.

Monitoring

Long-term monitoring of patients and establishment of guidelines for performance.

Economics

Financial costs and economic benefits associated with test.

Education

Educational materials and informed consent requirements.

ELSI

Assessment of ethical, legal and societal implications that arise in the context of the test.

ANALYSIS COMMON FEATURES

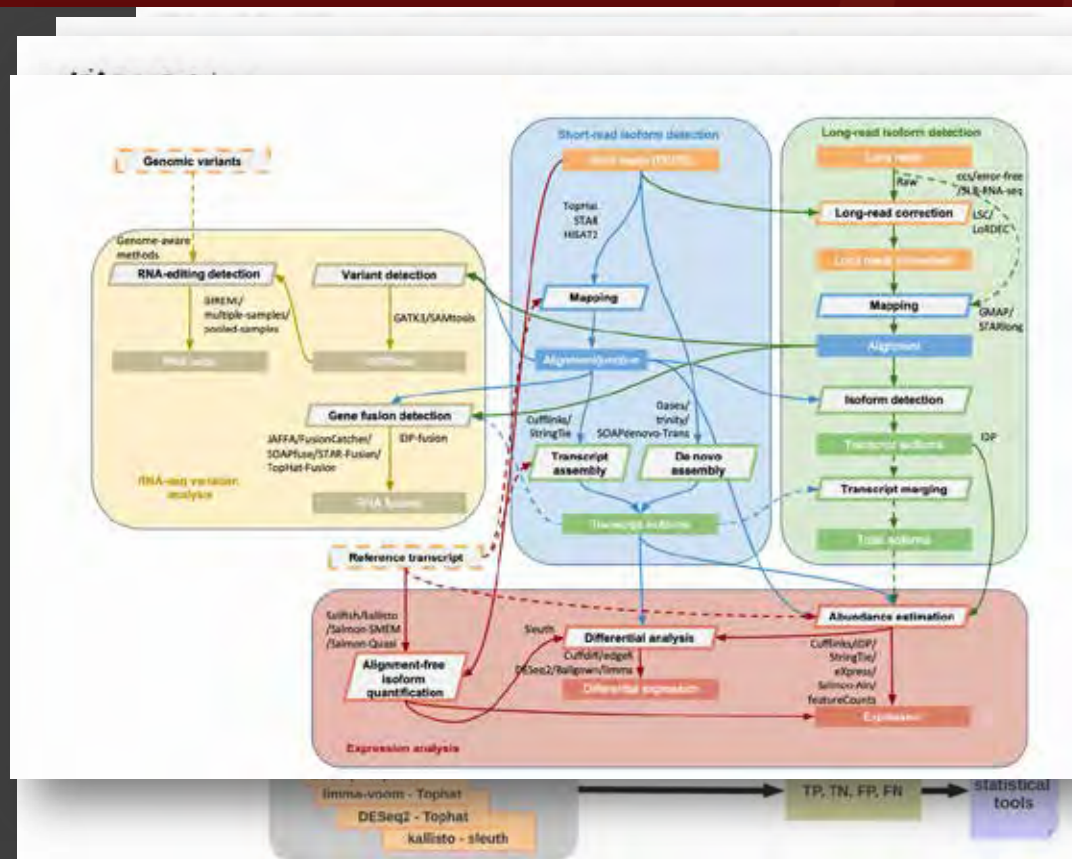
Quality Control

Alignment/Assembly

Detection & Abundance

Normalization

Significance Of Model/Hypothesis



RNA-SEQ 2020

- RNA-Seq
- CaptureSeq
- RASL-Seq
- ClickSeq
- 3Seq
- cP-RNA
- 3P-Seq
- 2P-Seq
- 3'-Seq
- TIF-Seq
- PEAT
- SMORE-Seq
- TL-Seq
- TATL-Seq
- RARseq
- TAIL-Seq
- PAL-Seq
- ChIRP
- CHART
- RAP
- GRO-seq
- Bru-Seq
- 3'NT Method
- NET-Seq
- mNET-Seq
- PARE-Seq
- GMUCT
- Ribo-Seq or ARTSeq
- RIP-Seq
- CLIP-Seq or HITS-CLIP
- Pol II CLIP
- miR-CLIP
- PIP-Seq
- hiCLIP
- RBNS
- TRIBE
- HiTS-RAP
- TRAP-Seq
- DLAF
- miTRAP
- CLASH
- 71 RNA Modifications
- MaRIP-Seq
- icSHAPE
- CIRS-Seq
- SHAPE-MaP
- Structure-Seq/DMS-Seq
- SPARE
- PARS-Seq
- Cap-Seq
- CIP-TAP
- scRNA-Seq
- SUPeR-Seq
- CirSeq
- TIVA
- PAIR
- CLaP
- CytoSeq
- Drop-Seq
- Hi-SCL
- InDrop
- snRNA-Seq
- Nuc-Seq
- Div-Seq
- SCRBS-Seq
- G&T-Seq
- scM&T-Seq
- scTrio-seq
- TCR Chain Pairing
- TCR-LA-MC PCR

Analysis Involves:

- Domain Specific Knowledge (Biology/Data)
- Often home-brew software, layered over standard pipelines
- R/Bioconductor Key Tools

UTILIZING IDEP FRAMEWORK FOR LEARNING

iDEP: an integrated web application for differential expression and pathway analysis of RNA-Seq data

[Steven Xijin Ge](#)  [Eun Wo Son](#) & [Ruihan Yao](#)

BMC Bioinformatics **19**, Article number: 534 (2018) | [Cite this article](#)

12k Accesses | **26** Citations | **23** Altmetric | [Metrics](#)

LEARNING BY EXAMPLE

iDEP:90 Load Data Pre-Process Heatmap k-Means PCA DEG1 DEG2 Pathway Genome Bicluster Network R

[Click here to load demo data](#)
and just click the tabs for some magic!

1. Select or search for your species.
Best matching species

2. Choose data type
 Read counts data (recommended)
 Normalized expression values (RNA-seq FPKM, microarray, etc.)
 Fold-changes and corrected P values from CuffDiff or any other program

3. Upload expression data (CSV or text)
Browse... No file selected

Analyze public RNA-seq datasets for 9 species

Optional: Upload an experiment design file(CSV or text)
Browse... No file selected

Loading R packages
Done. Ready to load data files.
New! Massively upgraded annotation database! V0.90 includes 315 organisms in Ensembl release 96, plus all species from STRINGdb (v10):115 archaeal, 1678 bacterial, and 238 eukaryotic species
Now published on BMC Bioinformatics!
Due to lack of funding, iDEP has not been thoroughly tested.
We are happy to help prepare your data for iDEP. Please email us.

Read counts or FPKM
Remove lowly expressed genes
Convert gene IDs to Ensembl IDs
DESeq2 /limma
Transcript clustering
Fold changes
Gene lists
K-means clustering
PCA
clustering
Pathway analysis
Enrichment analysis
Biclustering & Co-expression networks
Visualize on genomes
GO terms, KEGG pathways, TF target genes, miRNA target genes
iDEP workflow

iDEP: an integrated web application for differential expression and pathway analysis of RNA-Seq data

Steven Xijin Ge, Eun Wo Son & Bunsan Yoo

BMC Bioinformatics 19, Article number: 534 (2018) | [Cite this article](#)

12k Accesses | 26 Citations | 23 Altmetric | [Metrics](#)

iDEP:90 Load Data Pre-Process Heatmap k-Means PCA DEG1 DEG2 Pathway Genome Bicluster Network R

Email us for questions, suggestions, or data contributions. Stay connected via [user group](#) or [Twitter](#). Visit our [GitHub](#) page to see source code, install a local version, or report bugs and request features. iDEP is being developed by a very small team: Dr. Xijin Ge and a graduate student ([Homepage](#)).

R as in Reproducibility

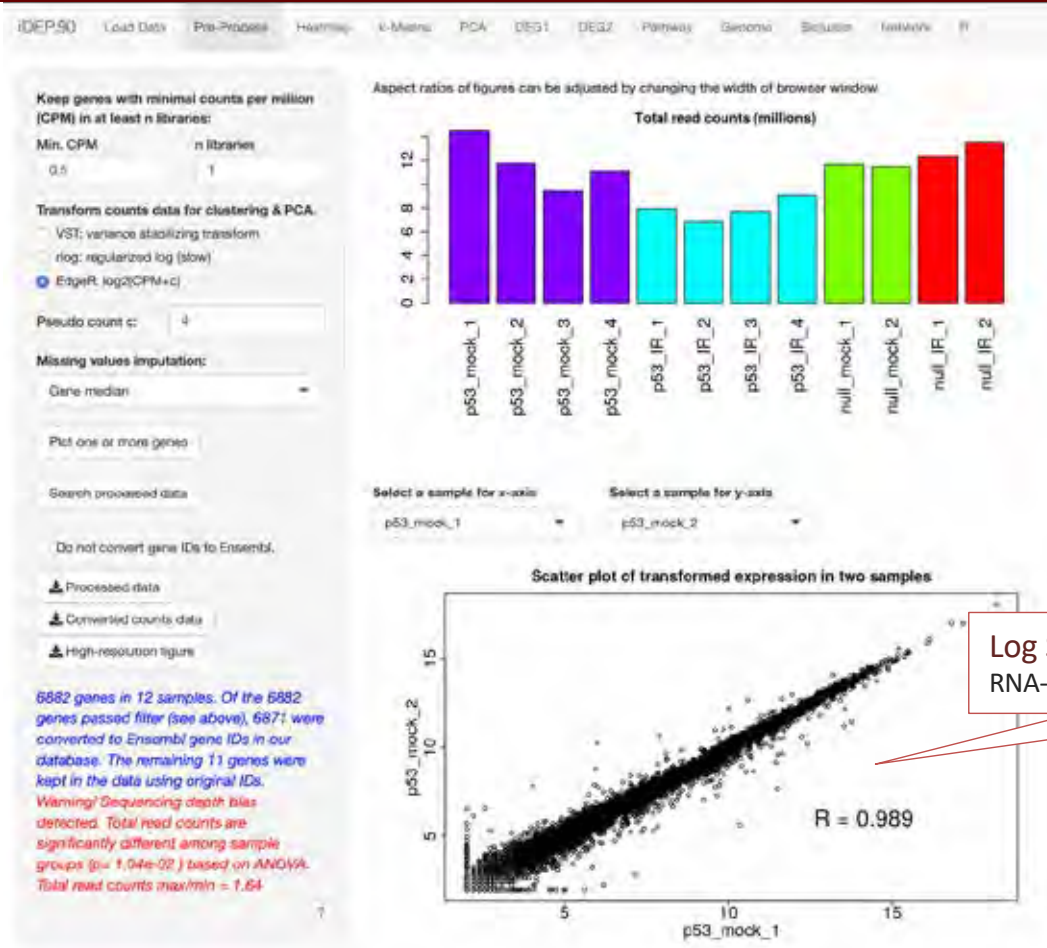
Documentation site. Source code on GitHub, where users can also report bugs or request features.

To improve reproducibility, iDEP generates custom R code based on your data and choices of parameters. Users with some R coding experience should be able to re-run most analyses by downloading all of the files below. If Ensembl IDs is not used in users' original file, we should use the converted data file. Click through all the tabs and then download all these file to a folder. Run the Customized R code or the Markdown file. [R Markdown example](#)

- Customized R code
- Customized R code(Markdown)
- iDEP core functions
- Gene info file
- Pathway file (large)
- Converted counts
- Experiment design file

<http://bioinformatics.sdstate.edu/idep/>

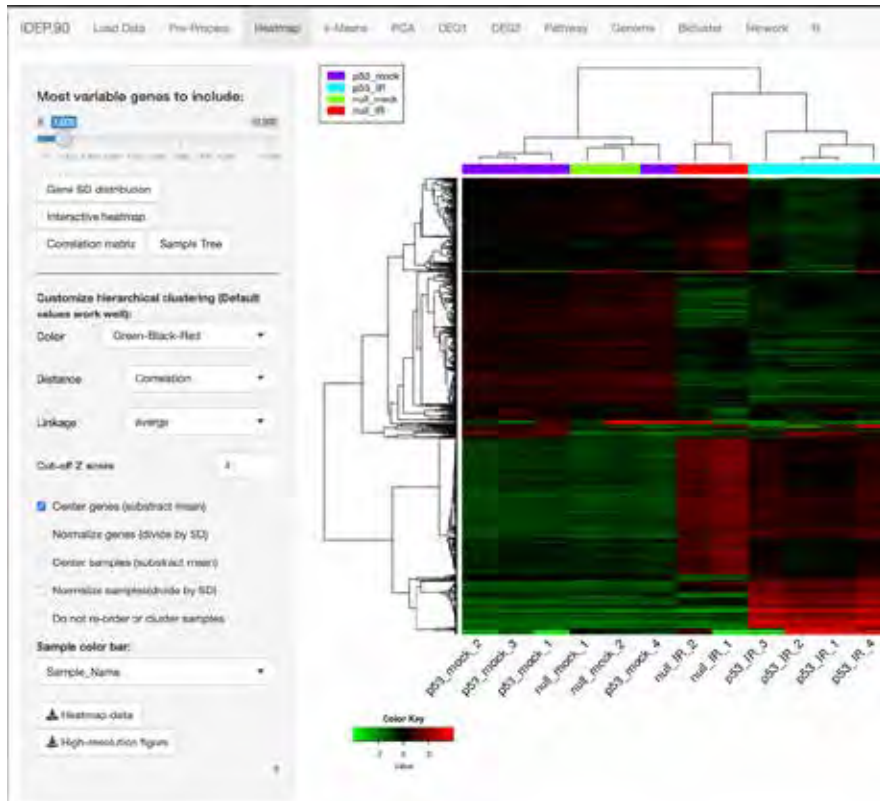
PRE-PROCESSING



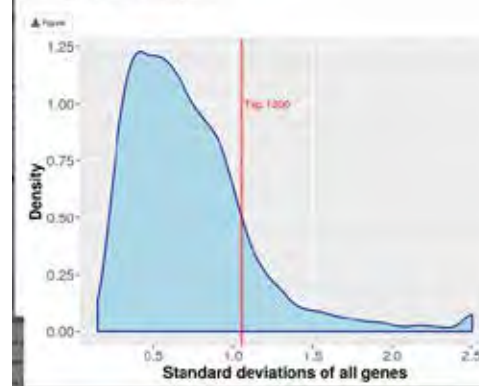
Log Scales
RNA-Seq spans many orders of magnitude

UNSUPERVISED: HEATMAPS, PCA, HIERARCHICAL CLUSTERING

Heatmaps



Reproducibility/Filtering



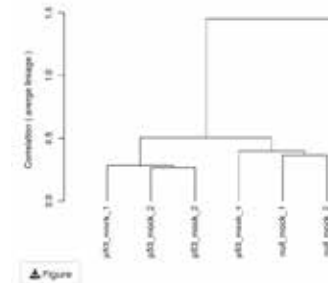
Correlation



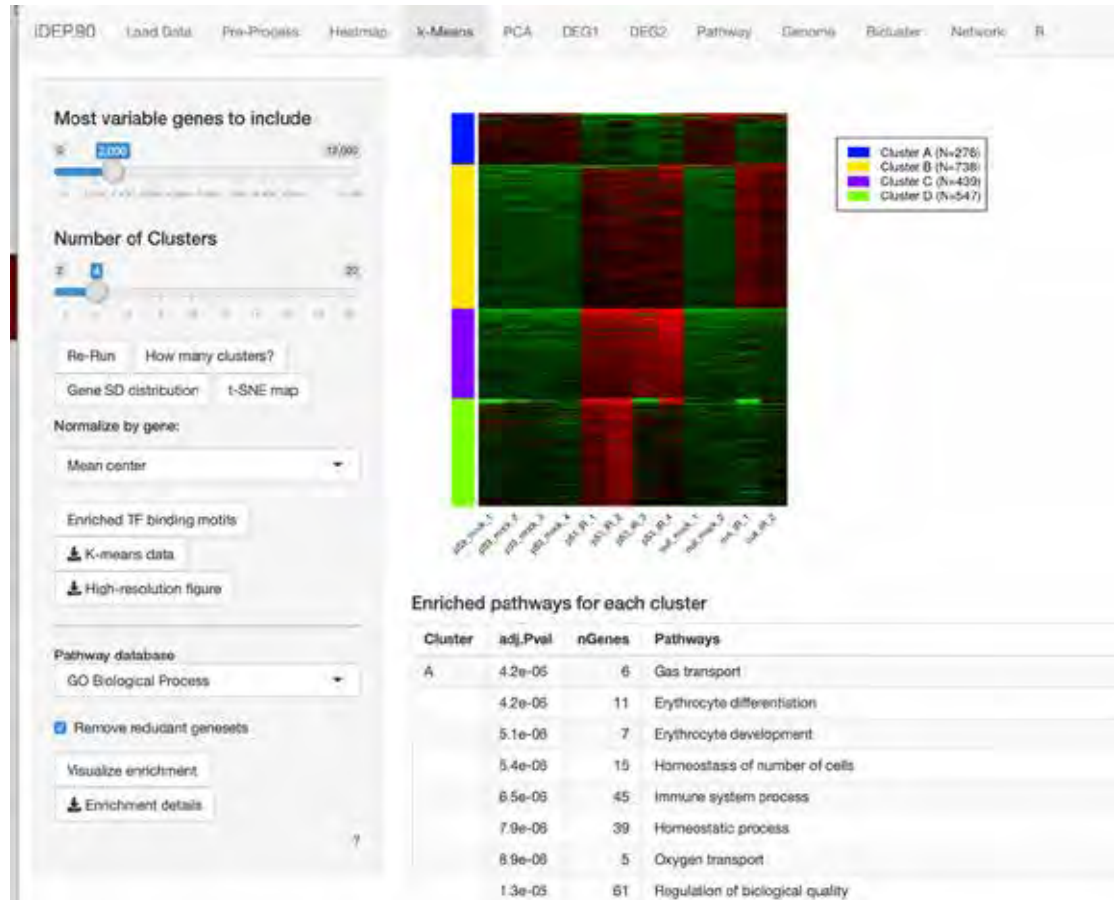
Hierarchical Clustering

Hierarchical clustering tree.

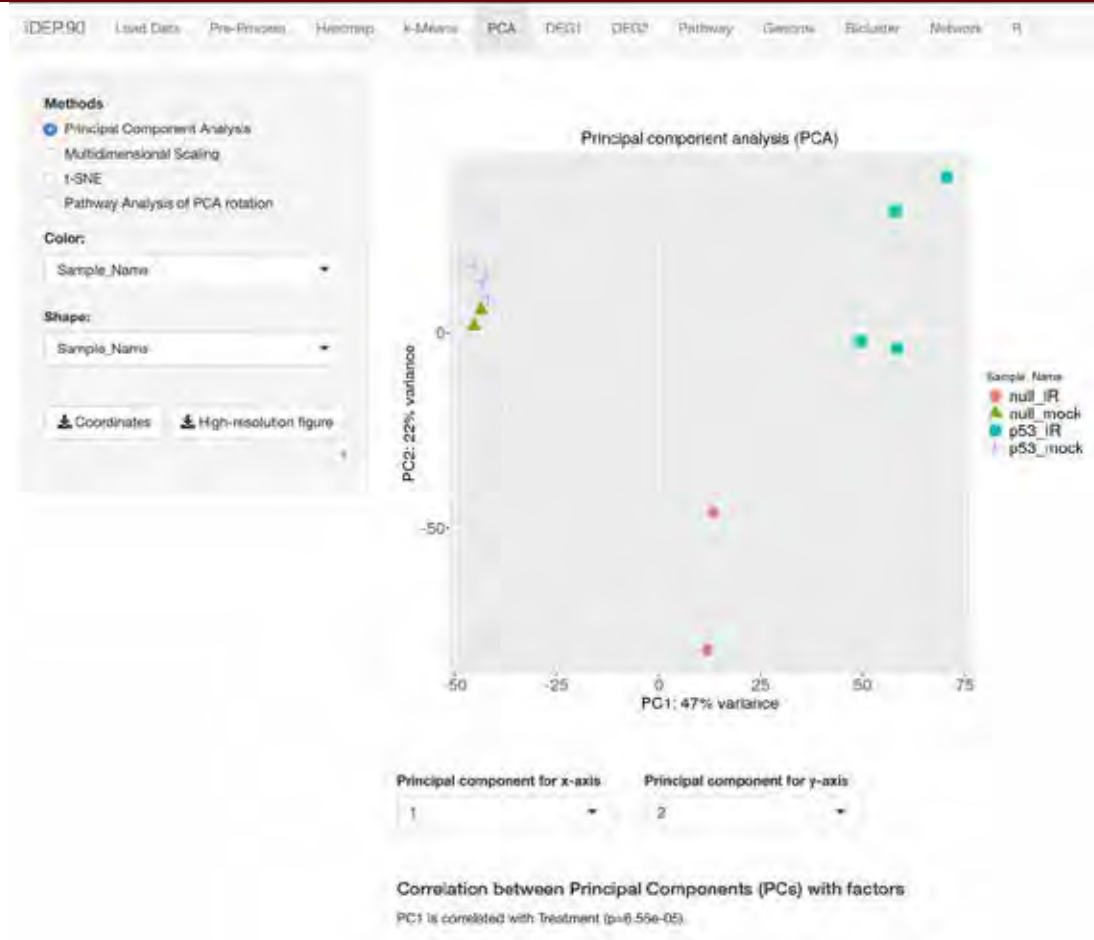
Using genes with maximum expression level at the top 75%. Data is transformed and clustered as specified in the main page.



CLUSTERING



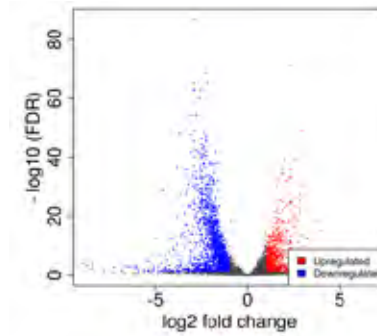
DIMENSION REDUCTION: PCA, T-SNE, MDS



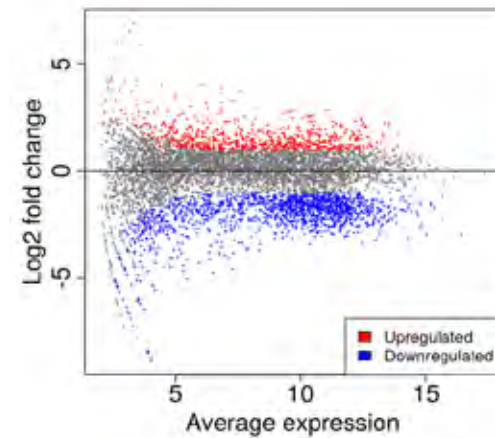
DIFFERENTIAL EXPRESSION



Volcano Plots



MA-Plots



PATHWAY ANALYSIS

KEGG Pathway Analysis Interface

Select a combination to analyze: null, null-null, R1

Select method: GAGE

Select genes (Choose KEGG to show pathway diagram): GO Biological Process

Connect size: Min. 10, Max. 2000

Pathway algorithm used: ZORA 3.5

Number of top pathways to show: 20

Use absolute values of fold changes for GSEA and GAGE

Remove genes with big FDR before pathway analysis

Pathway list: Pathway network

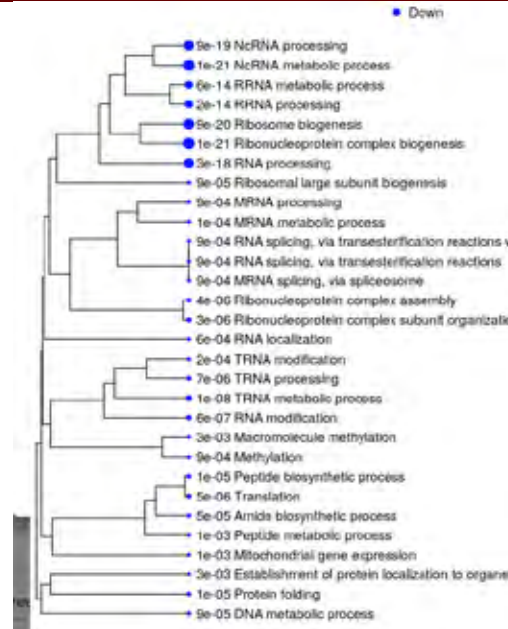
Warning: The many combinations can lead to false positives in pathway analysis.

Direction	Gene	Statistic	Gene	adj.Pval
Down	miRNA metabolic process	-11.2173	185	8.8e-22
	Ribosome biogenesis	-11.1908	177	8.4e-22
	Ribosome biogenesis	-11.0723	122	8.7e-20
	miRNA processing	-10.2228	127	8.2e-19
	mRNA processing	-9.8888	97	6.7e-19
	RNA processing	-9.0123	84	2.5e-14
	RNA metabolic process	-8.9802	88	2.5e-14
	RNA metabolic process	-7.8919	61	1.2e-08
	RNA metabolic process	-6.7253	48	3.9e-07
	RNA processing	-6.1129	37	7.2e-06
	Ribosome biogenesis	-5.3228	66	7.9e-05
	Ribosome biogenesis	-5.1712	82	3.4e-04
	RNA metabolic process	-5.0991	57	1.8e-04
	Protein folding	-5.0880	81	1.8e-04
	Translational	-5.0687	178	6.4e-06
	Ribosome large subunit biogenesis	-5.2915	51	8.4e-05
	Peptide biosynthetic process	-5.0347	182	1.1e-05
	Amide biosynthetic process	-5.0514	200	5.7e-05
	DNA metabolic process	-4.9479	298	8.9e-05
	RNA metabolic process	-4.9989	140	1.9e-04
	Mitochondrial gene expression	-4.7447	24	1.8e-03
	RNA splicing	-4.8281	62	6.2e-04
	RNA splicing, via transesterification reactions	-4.1149	61	8.6e-04
	RNA splicing, via transesterification reactions with bulged adenosine as nucleophile	-4.1149	61	8.6e-04
	mRNA splicing, via transesterification reactions	-4.1149	61	8.6e-04
	mRNA processing	-4.1001	94	8.9e-04
	Metabolism	-4.0229	192	9.3e-04
	Peptide metabolic process	-4.1181	178	1.2e-03
	Macromolecule methylation	-4.0579	80	8.6e-03
	Establishment of protein localization to organelle	-4.0771	126	2.9e-03

Select a pathway to show expression pattern of related genes on a heatmap or a KEGG pathway diagram: miRNA metabolic process

Expression data for genes in selected pathway

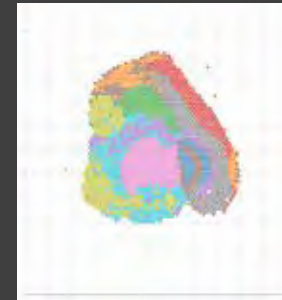
Heatmap showing expression levels for genes in the miRNA metabolic process pathway.



DATABASES, ANALYSIS, AND SOFTWARE TOOLS

Resource	Description	URL	Refs	Analysis			
Annotation				OncoPrint	Web application for user-friendly analysis and exploration of cancer transcriptomes	https://www.oncoprint.org/resource/login.html	180
RefSeq	Curated reference sequence database (transcriptome-centric, that is, defined by transcript sequence)	https://www.ncbi.nlm.nih.gov/refseq/	75	Xena	UCSC Xena: versatile genomic data mining and analysis portal	https://xenabrowser.ucsf.edu/	267
Gencode	Curated reference gene annotation (genome-centric, that is, defined by alignment to reference genome)	http://www.gencodegenes.org/	272	Data warehouse			
MiTranscriptome	Automated reference transcriptome based on sequence assembly, includes long non-coding RNAs	http://mitranscriptome.org/	76	ENCODE	Repository of diverse functional genomics data, including RNA-seq, from the ENCODE project	https://www.encodeproject.org/	59
Reference data				GDC	Genomic Data Commons: provides access to raw and harmonized data for multiple genomic projects, including RNA-seq data processed using a standard pipeline	https://portal.gdc.cancer.gov/	295
MSigDB	Collection of experimental and curated gene sets (signatures)	http://software.broadinstitute.org/gsea/msigdb	179	FANTOM5	Repository of CAGE data from the FANTOM5 project	http://fantom5.gsc.nriks.jp/	271
Human Protein Atlas	Compendium of proteomic and transcriptomic data in diverse normal tissues	http://www.proteinatlas.org/	63	ArrayExpress	Standard repositories of functional genomic and transcriptome profiling data	http://www.ebi.ac.uk/arrayexpress/	76
CCLE	Genomic and transcriptomic data on hundreds of cancer cell lines	https://portals.broadinstitute.org/ccle/home	60	GEO	Standard repositories of functional genomic and transcriptome profiling data	https://www.ncbi.nlm.nih.gov/geo/	77
GTEX	Transcriptomic data (RNA-seq) from normal human tissues from a large number of individuals	https://gtexportal.org/home/	62	CAGE, cap analysis of gene expression; CCLE, Cancer Cell Line Encyclopedia; ENCODE, Encyclopedia of DNA Elements; FANTOM5, Functional Annotation of the Mammalian Genome 5; GENCODE, the genome annotation project of ENCODE; GEO, Gene Expression Omnibus; GTEX, Genotype-Tissue Expression Project; Limma, Linear Models for Microarray Data; MSigDB, Molecular Signatures Database; PARADIGM, Pathway Recognize Algorithm using Data Integration on Genomic Models; QoRTs, Quality of RNA-seq Toolset; RNA-seq, RNA sequencing; STAR, Spliced Transcripts Alignment to a Reference; UCSC, University of California, Santa Cruz.			
Mitelman	Database of gene fusions and chromosomal aberrations	https://cgap.nc.nih.gov/Chromosomes/Mitelman	288				
COSMIC	Catalogue of somatic mutations in cancer patients and cell lines, including gene fusions	http://cancer.sanger.ac.uk/cosmic/classic#fus	289				
Tool							
QoRTs	Comprehensive collection of RNA-seq quality control functions	http://hertleya.github.io/QoRTs/index.html	290				
STAR	Fast and accurate splice-aware sequence aligner	https://github.com/AlexDobin/STAR	262				
featureCounts	Fast read counting for gene-level or exon-level expression estimates	http://bioinf.welhi.edu.au/featureCounts/	291				
Kallisto	Pseudo-alignment-based quantification at the transcript level	https://pachterlab.github.io/kallisto/	292				
EdgeR	Differential expression using the negative binomial distribution (see also DESeq2)	http://bioconductor.org/packages/release/bioc/html/edgeR.html	170				
Limma	Flexible linear modelling and empirical Bayes moderation to assess differential expression by use of precision weights for RNA-seq data (Voom)	http://bioconductor.org/packages/release/bioc/html/limma.html	167, 168				
CIBERSORT	In silico transcriptome deconvolution into relative abundances of different immune cell types	https://cibersort.stanford.edu/	260				
MIXCR	T cell and B cell CDR3 sequences assembler; enables repertoire profiling from RNA-seq data	https://millsboon.com/software/mixcr/	261				
GSEA	Gene set enrichment analysis	http://www.broad.mit.edu/GSEA	273				
PARADIGM	Computational tool for the inference of patient-specific pathway activities	https://sbenz.github.io/Paradigm	186				
FusionCatcher	A sensitive and specific tool for the detection of gene fusions	https://github.com/ndaniel/fusioncatcher	293				
TopHat-Fusion	A very sensitive tool for the detection of gene fusions	http://ccb.jhu.edu/software/tophat/fusion_index.html	294				

CUTTING EDGE: RNA-SEQ EMPOWERING ANALYSIS OF HETEROGENEITY



Single cell RNA-seq
Spatial Transcriptomes

SINGLE-CELL RNA-SEQ

Functional Studies w/ snRNA-seq

Study Types

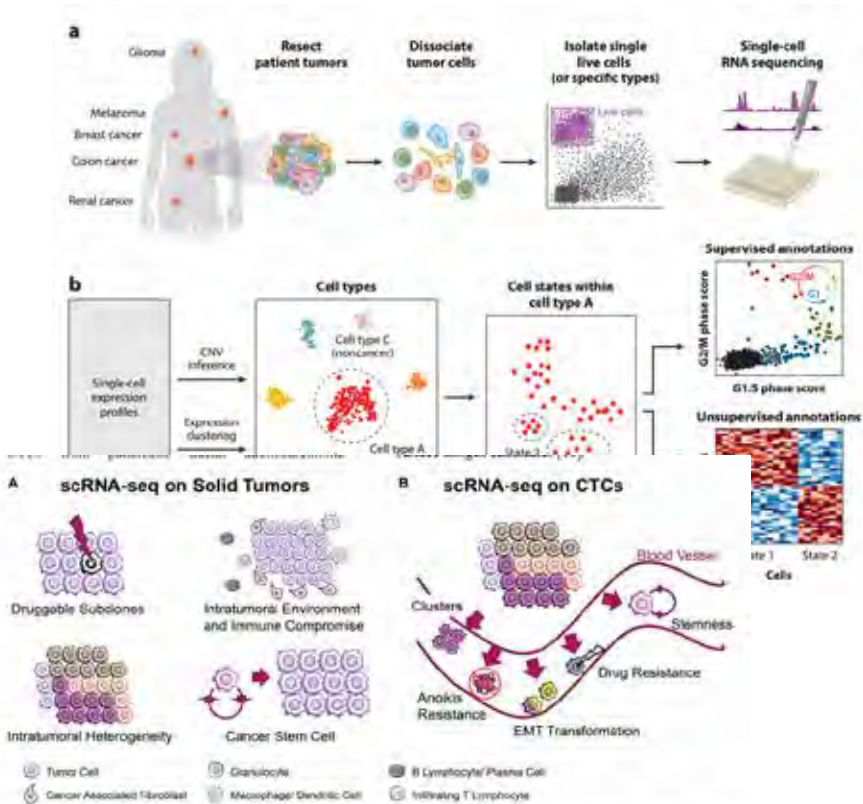
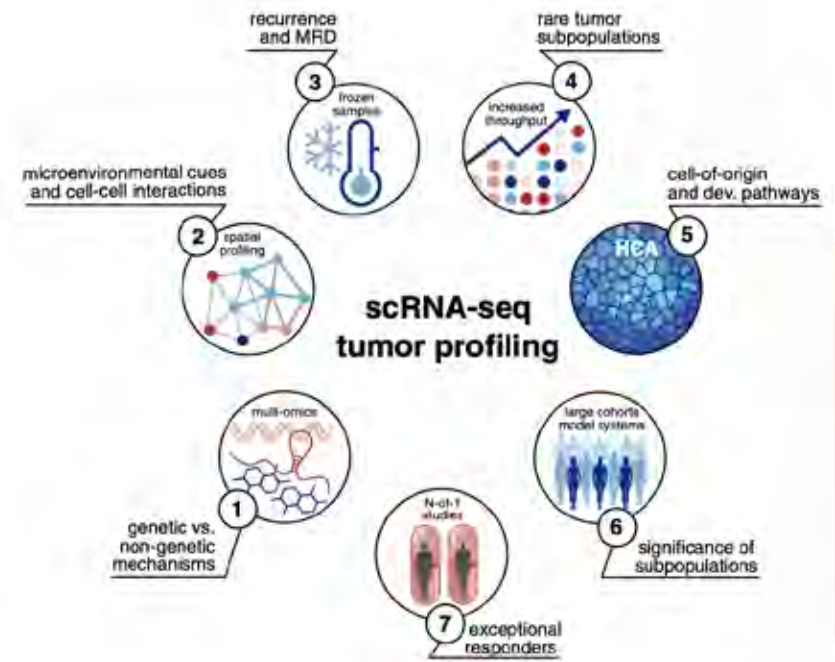


Figure 1: scRNA-seq technology facilitates cancer research when it comes to... www.annualreviews.org • Single-Cell Expression Profiling in Cancer



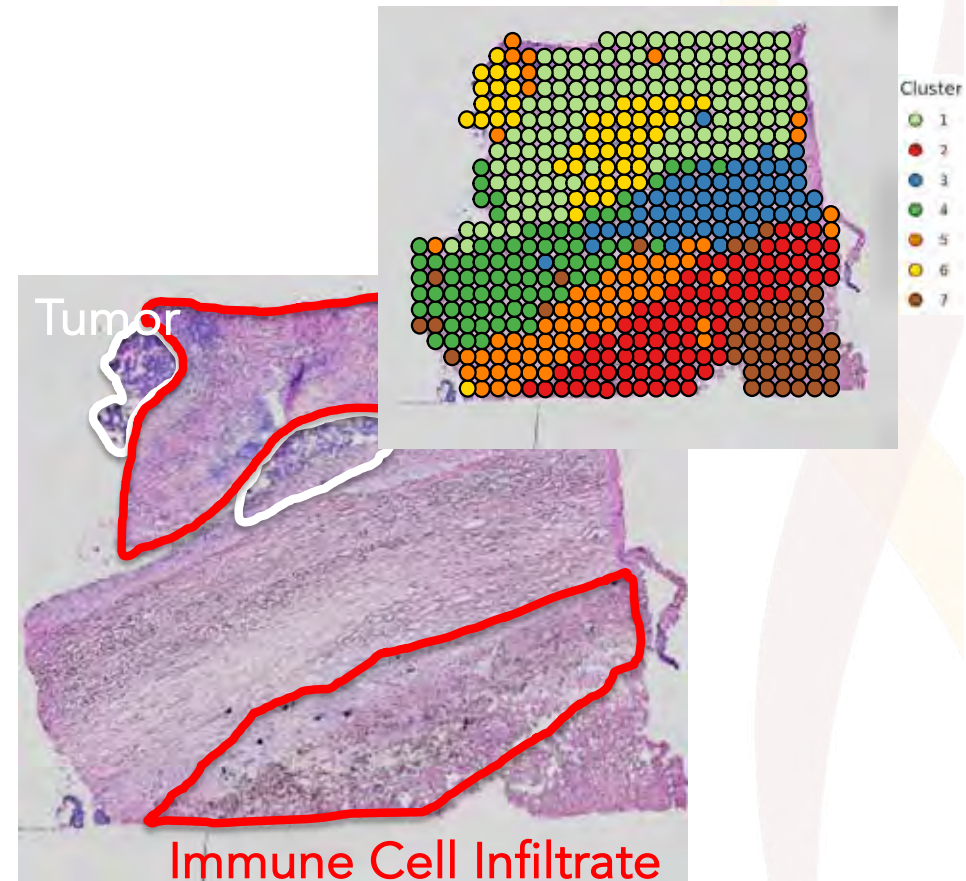
ASSESSING TUMOR HETEROGENEITY

Molecular Annotation of Cluster Data

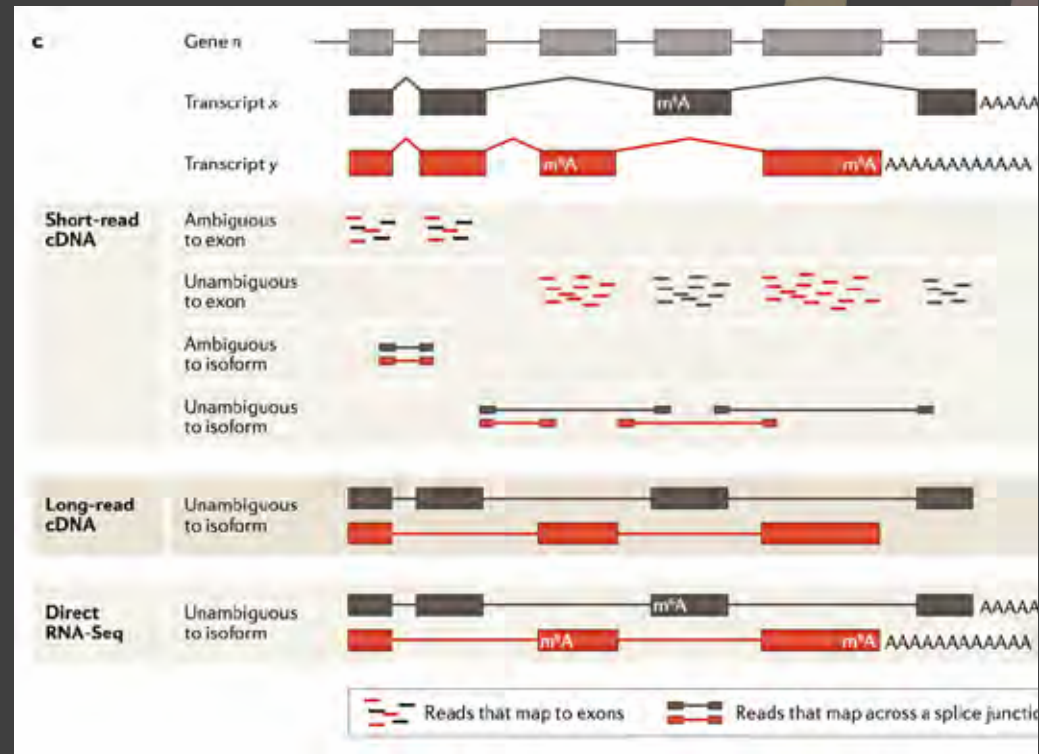
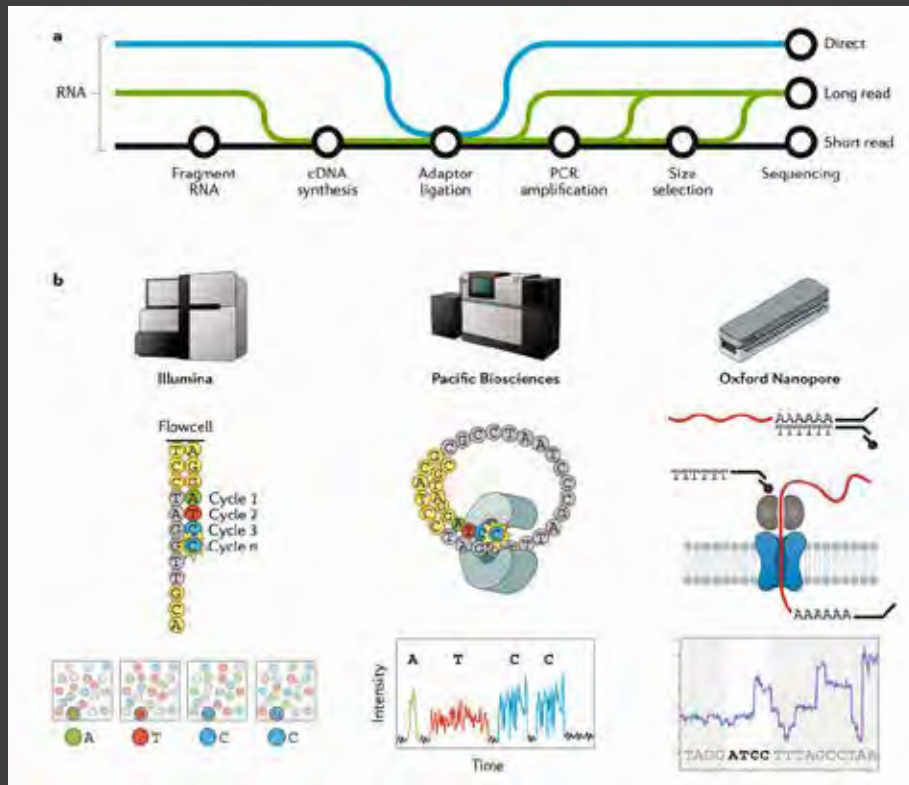
- Immune cluster profiling
- Spatial Gene Expression Maps
 - ESTIMATE Yoshihara et al., Nature Communications; 4 (2013): 2612
- Cibersort
 - Newman et al., Nature Methods. 2015; 12:453–457 (2015)
- xCell
 - Aran et al., Genome Biol. 2017;18(1):220.
- Inflammation and Immune Scoring
 - Ayers et al., J Clinical Investigation. 2017; 127(8):2930-2940.

Tumor cluster profiling

- GSEA
 - Subramanian, Tamayo, et al. PNAS. 2005; 102, 15545-15550.



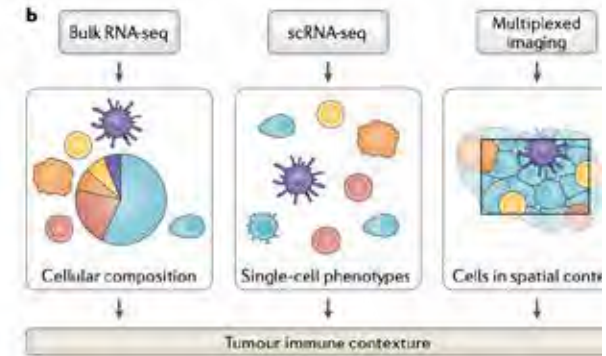
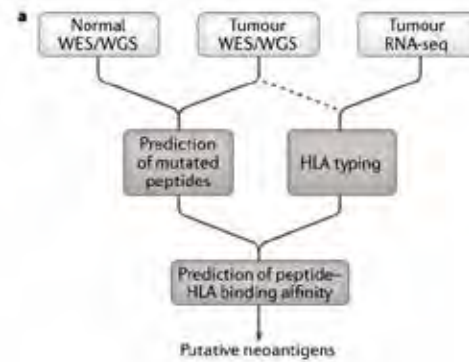
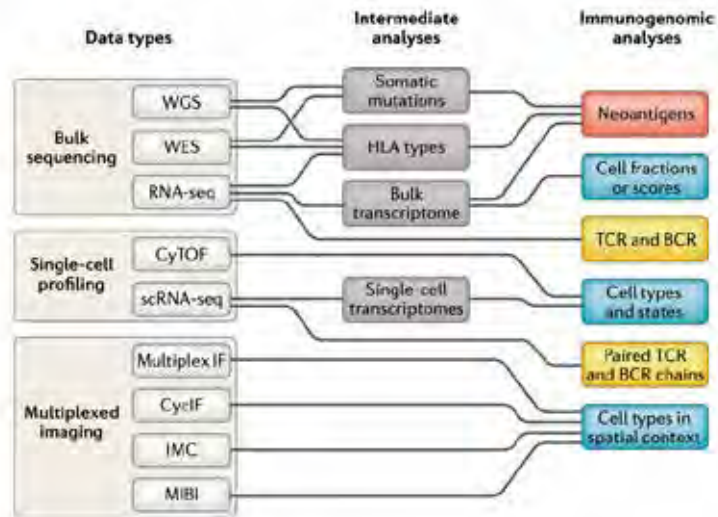
BLEEDING EDGE LONG-READS REALTIME

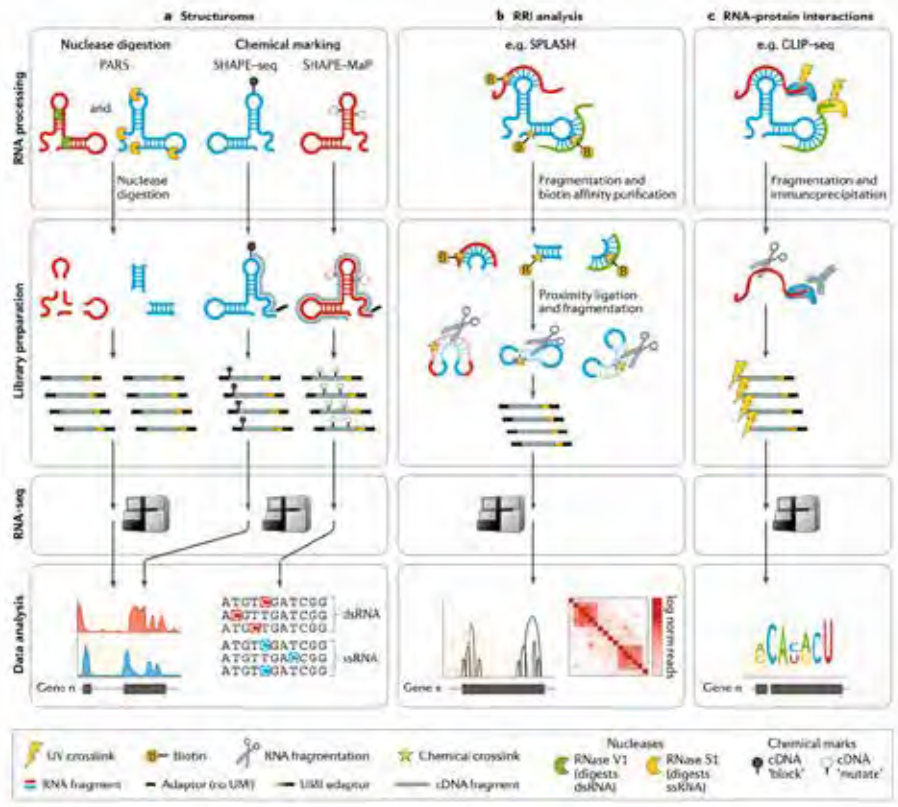
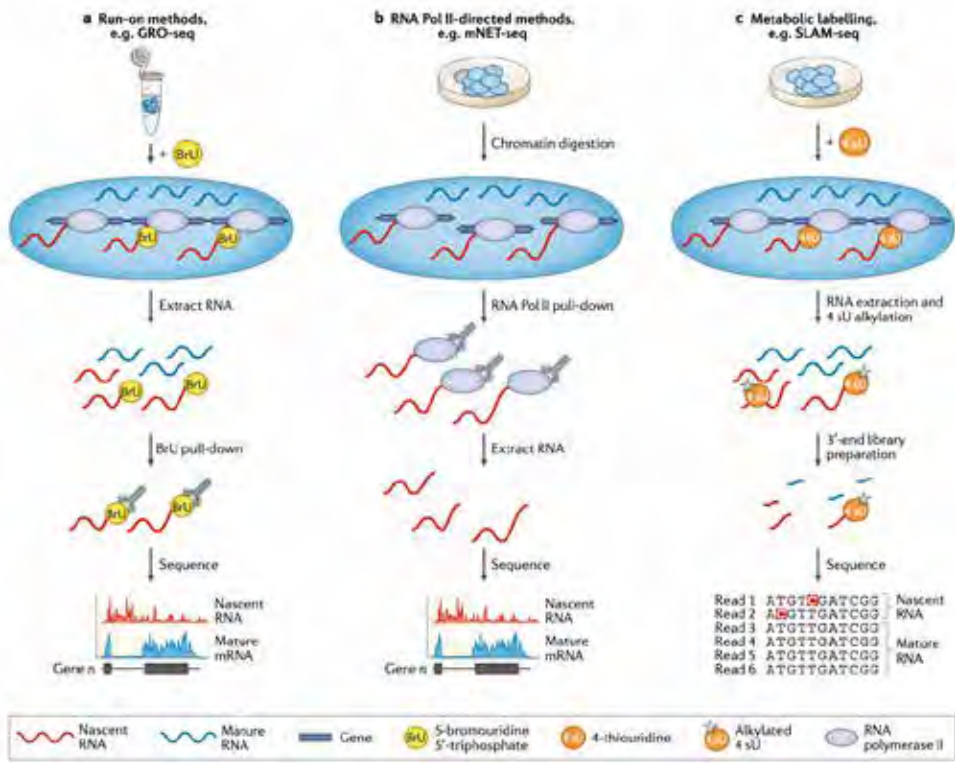


RNA sequencing: the teenage

Amy Stacy, Mirna Grzelak, and James Hadfield

BLEEDING EDGE IMMUNE SINGLE CELL





IN CLASS ACTIVITY