

# DCEG Statistical Genetics Workshop 2023

## Session 9: Functional genomics

**Laboratory of Translational Genomics:**

Ludmila Prokunina-Olsson, PhD

Kevin Brown, PhD

Jiyeon Choi, PhD

Laufey Amundadottir, PhD

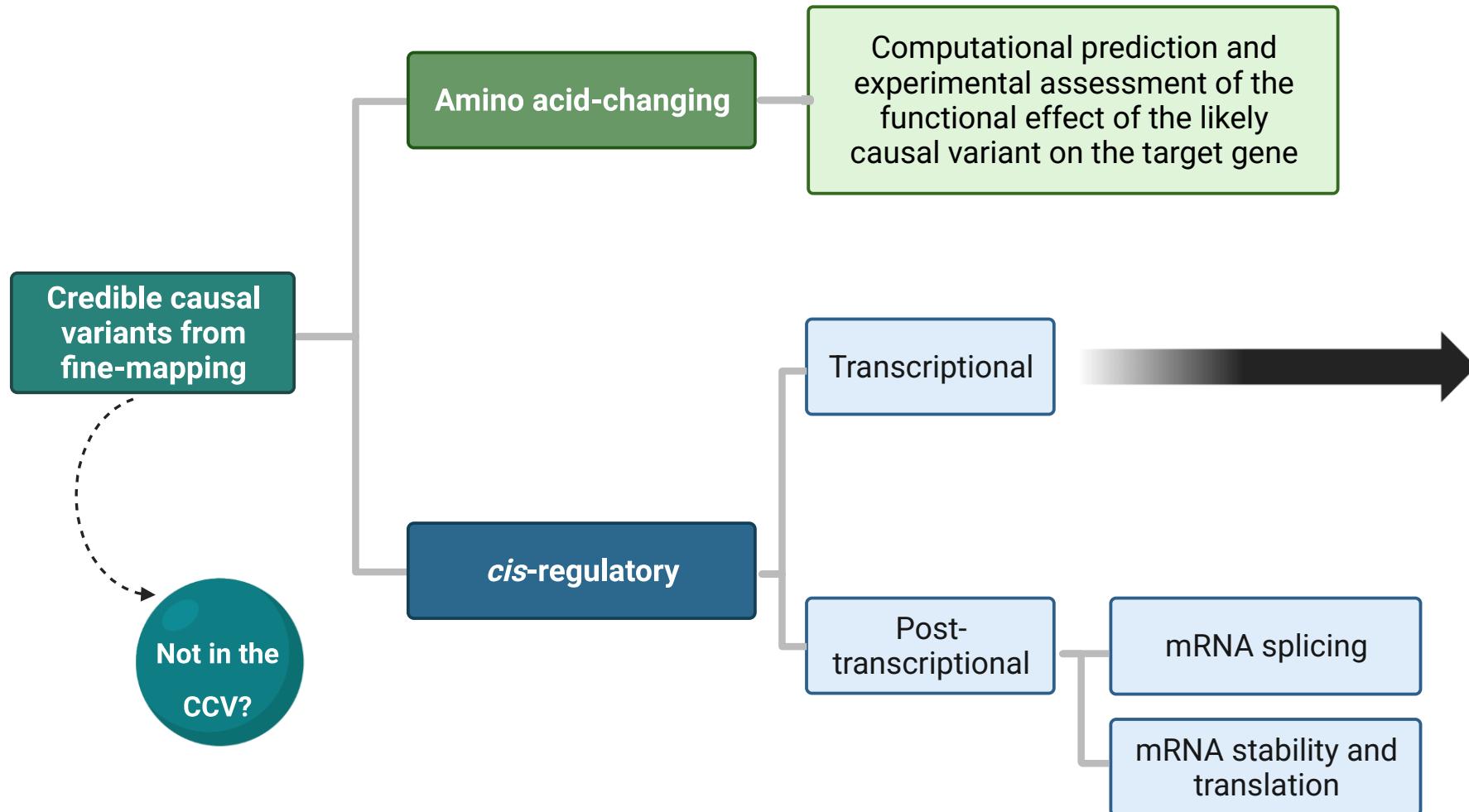
Oscar Florez-Vargas, PhD

December 13, 2023

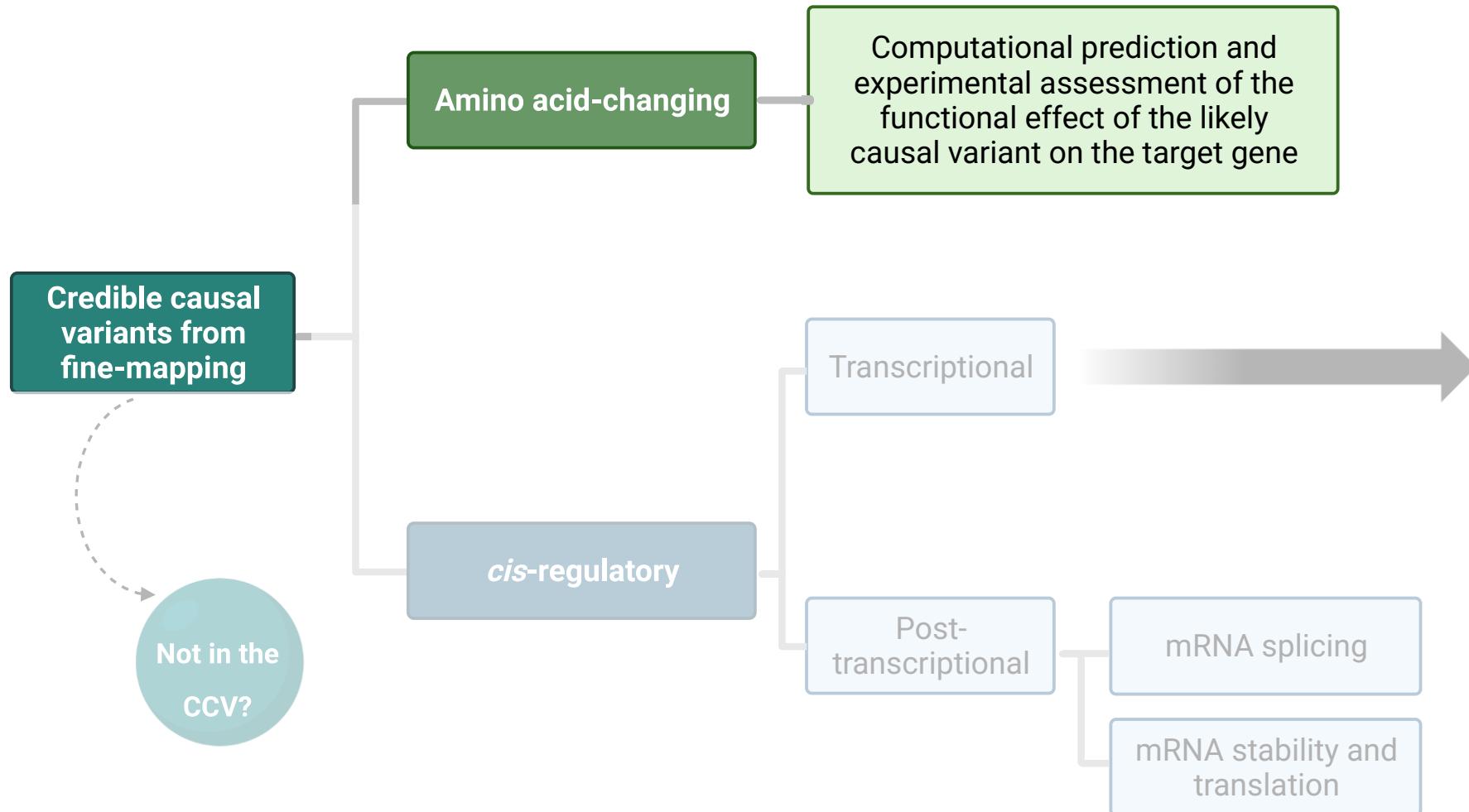
# Post-GWAS questions:

- What are the “functional” variant(s) lead GWAS SNPs are tagging?
- What are the gene(s) through which underlying genetic variation influence risk?
- How does the functional variant influence gene function?
- What role does the gene play in cancer associated phenotypes?

# Mechanisms of variant function



# Mechanisms of variant function



# Protein-coding variants: pathogenicity prediction

There are MANY tools to predict potential impact of non-synonymous protein-coding variants on protein function, based on:

- DNA and or protein sequence conservation
- Protein structure
- Other features (including SNPs and pathogenic variants)

Common variants from GWAS -> small effect sizes

- Various tools often give conflicting interpretations

Experimental validation required, many potential mechanisms by which to influence protein function

Note: protein-coding variants can also function via most of the other mechanisms by which non-coding variants can, including:

- be *cis*-regulatory!
- Regulate splicing

Nice review and list of tools (see Table 1):

<https://www.frontiersin.org/articles/10.3389/fgene.2022.1010327/full>

# Relatively few fine-mapped variants are non-synonymous

**Table 1** The genomic context in which a variant is found can be used as preliminary functional analysis

Classification	Approximate percentages <sup>a</sup>	Approximate numbers <sup>a</sup>
Intronic	40	1,047
Intergenic	32	838
Within non-coding sequence of a gene	10	262
Upstream	8	210
Downstream	4	105
<b>Non-synonymous coding</b>	<b>3</b>	<b>79</b>
3' untranslated region	~1	26
Synonymous coding	~1	26
5' untranslated region		
Regulatory region		
Nonsense-mediated decay transcript		
Unknown	~1	26
Splice site		
Gained stop codon		
Frameshift in a coding sequence		

The table broadly summarizes the genomic context of disease- and trait-associated SNPs annotated in the Catalog of Genome-Wide Association Studies (<http://www.genome.gov/gwastudies/>) as of December 9th, 2010: 1,212 published genome-wide associations with  $P < 5 \times 10^{-8}$  for 210 traits totaling 2,619 SNPs. Most of the SNPs are located in intergenic and intronic positions, but a small percentage are located upstream and downstream of genes, as well as in regulatory regions and splice sites. SNPs in these locations can be analyzed in more detail using more specific bioinformatics tools.

<sup>a</sup>Values are indicative and dependent on genomic boundaries used.

# Non-protein coding mechanisms

## **cis-regulation:**

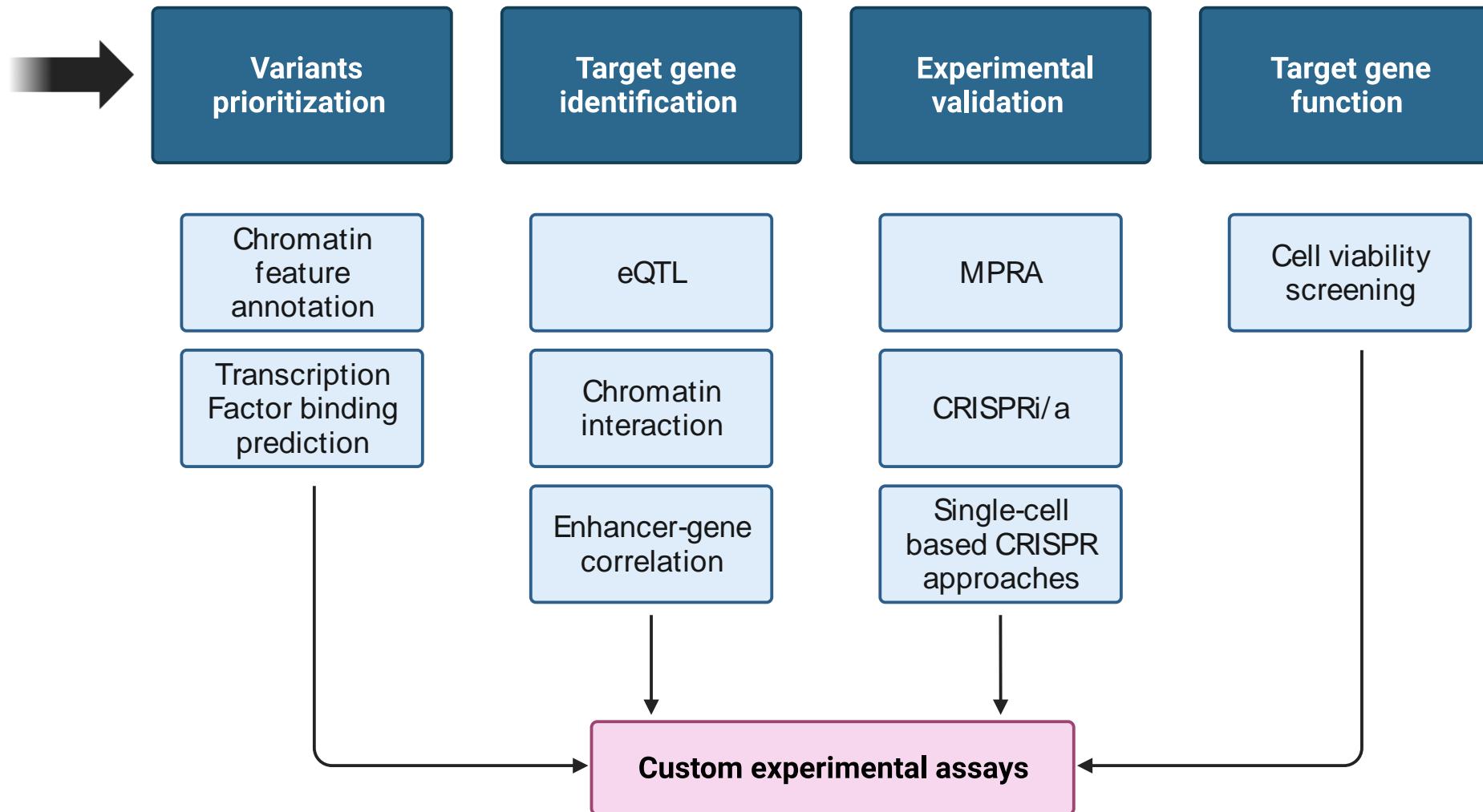
- Altered TF binding/regulatory activity in enhancer/promoter
  - Direct regulation of transcription
  - Regulation of chromatin accessibility
- Alteration of 3-Dimensional chromatin structure and localization
  - ->Changes to enhancer-promoter associations
- Altered CpG methylation

## **Post-transcriptional:**

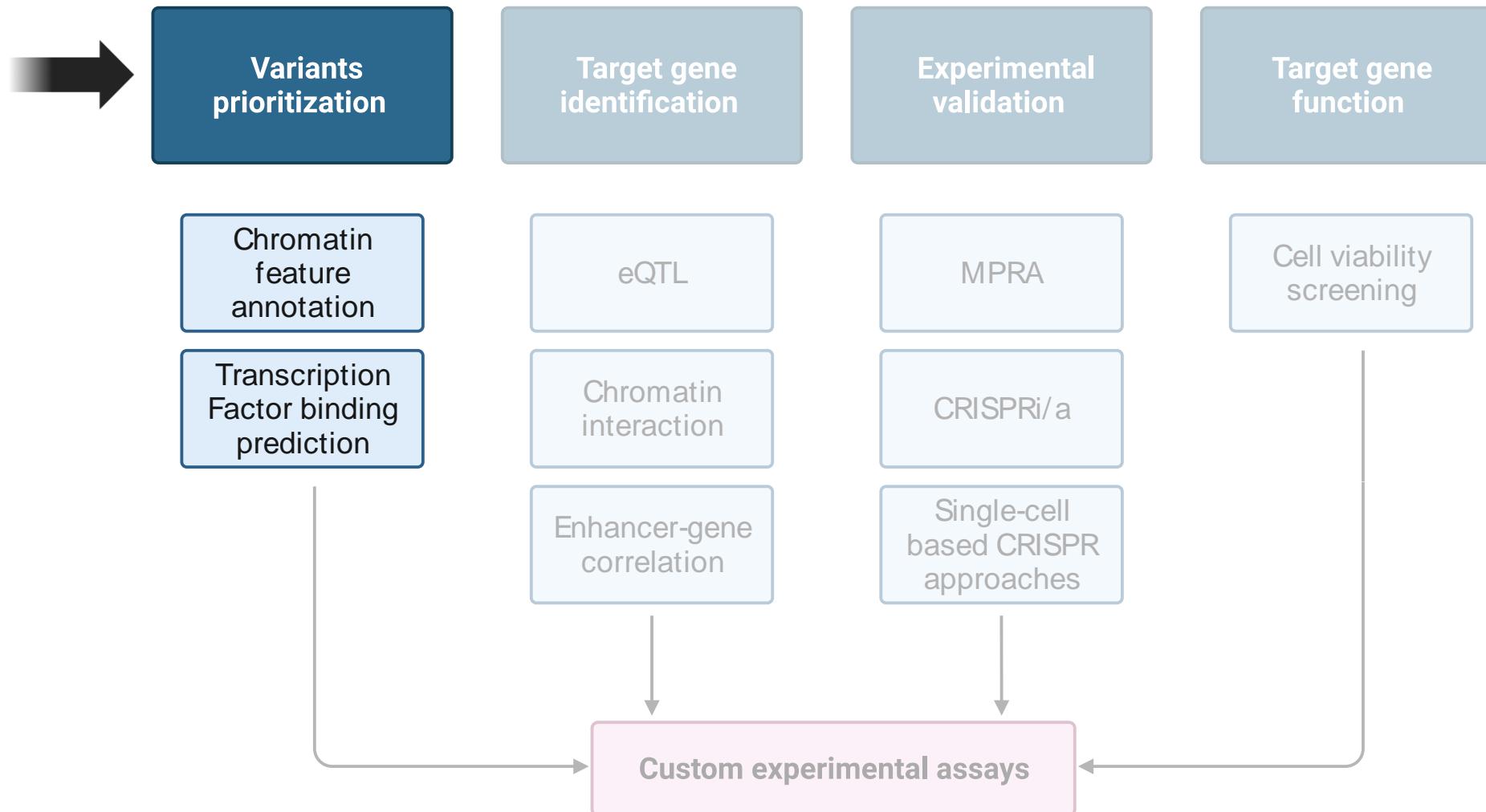
- Splicing (splice variants can be in both introns or exons)
- Post-transcriptional modifications, polyadenylation, etc.
- Transcript stability
- Translation efficiency
- Alternative codon usage, translational efficiency, and protein folding

## **Altered function of non protein-coding transcripts**

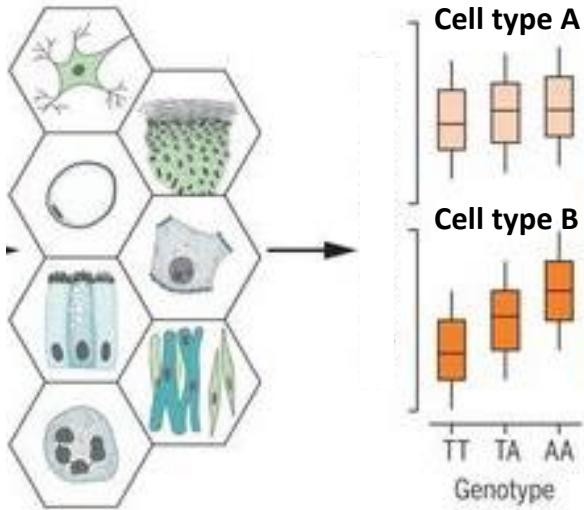
# Identification of "causal" variants and target genes



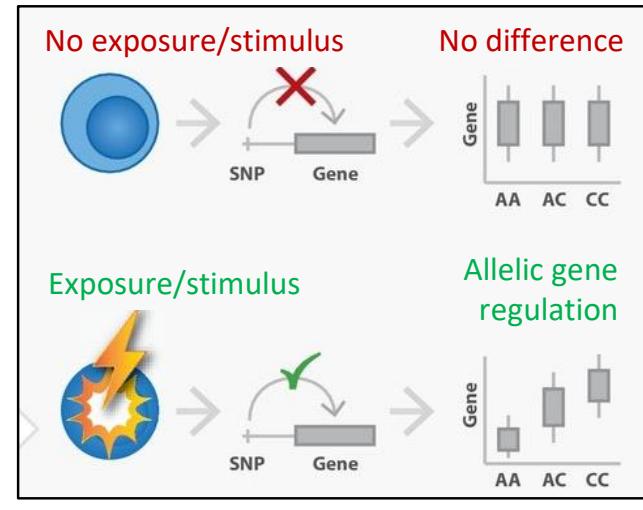
# Identification of "causal" variants and target genes



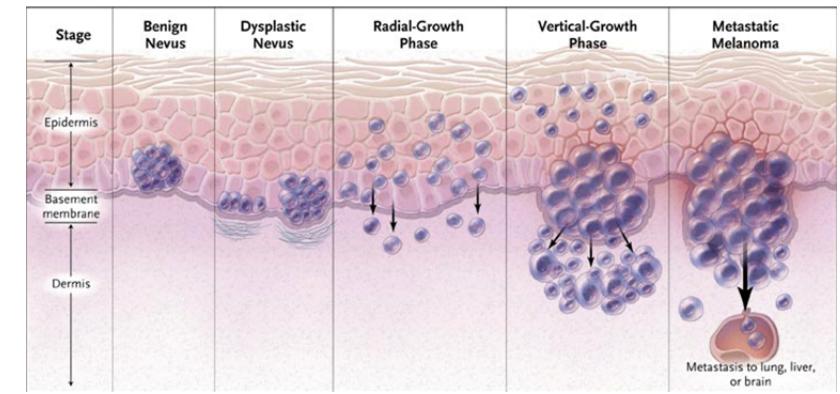
# Complex function of cancer-associated variants



Modified from Kim-Hellmuth *et al.* *Science* 2020



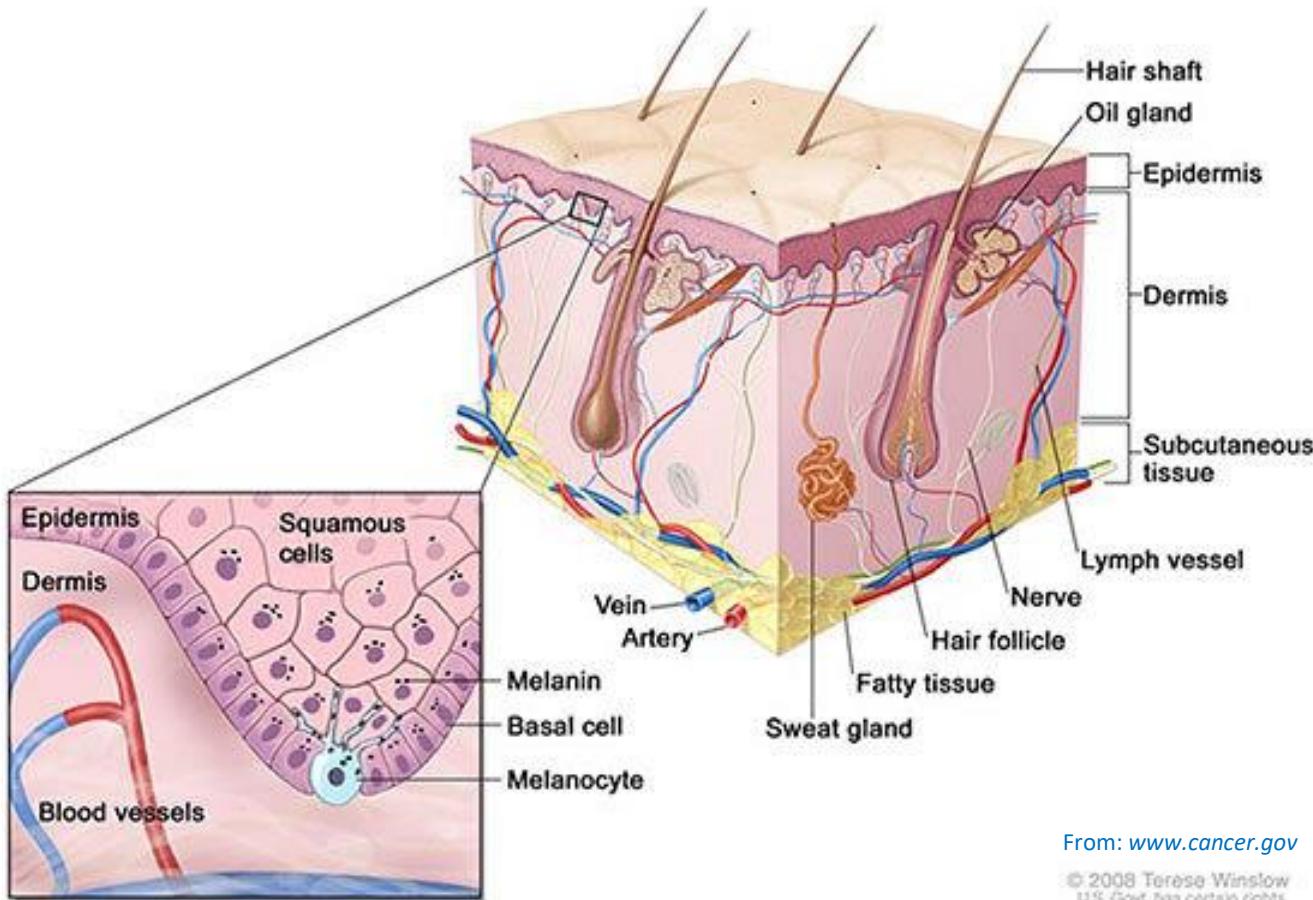
Modified from Wijst *et al.* *eLife* 2020



Miller *et al.* *NEJM* 2006

- Cell-of-tumor-origin
- Primary vs. cancer cells
- Cells in TME
- Exposure/stimulus
- Age, sex,  
race/ethnicity/ancestry
- Initiation, progression, outcome
- Germline x somatic
- Specific driver events

# Example: melanoma originates from melanocytes



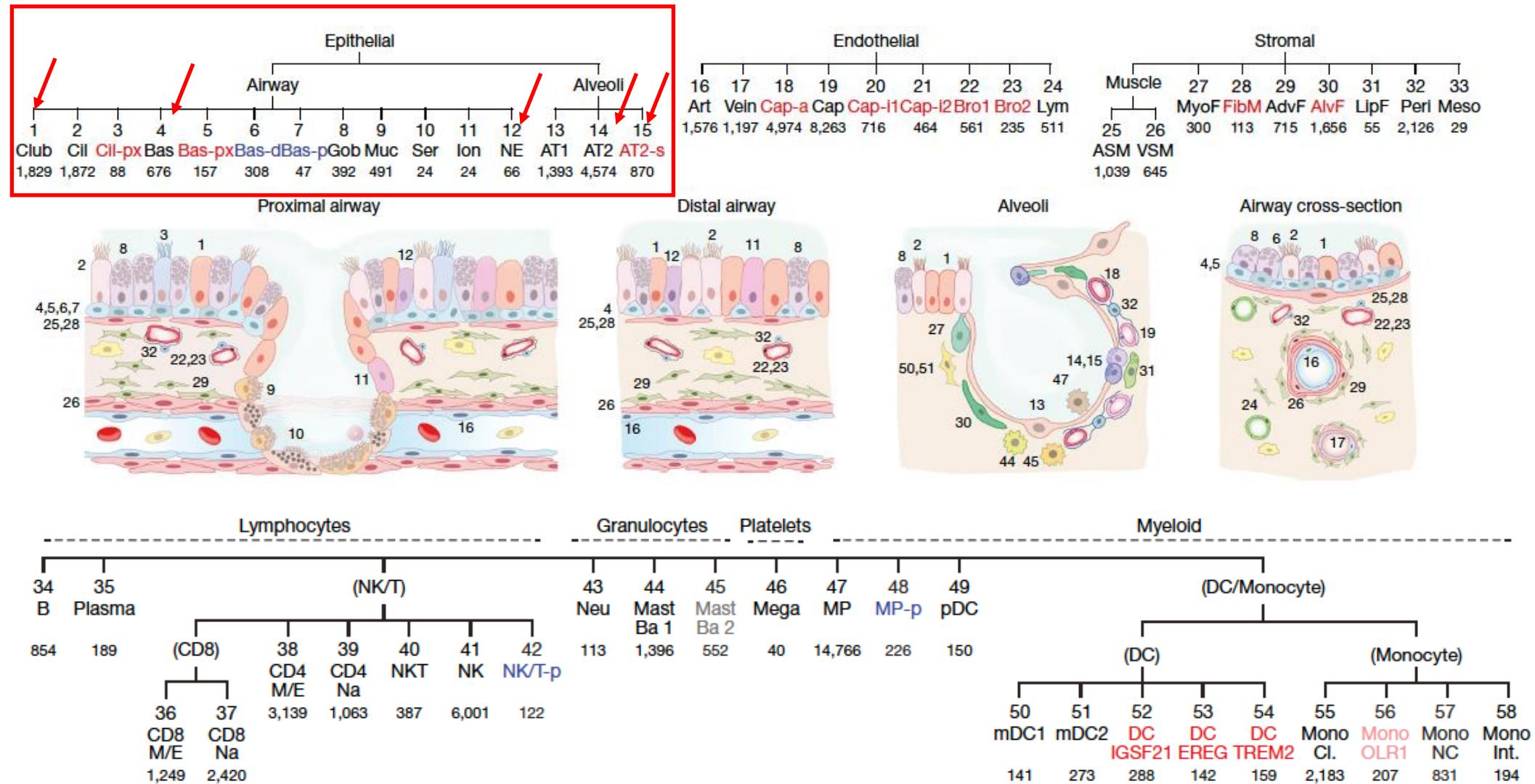
Melanocytes are a small fraction of the cells in a skin biopsy

Melanocyte-specific contributions to expression and allelic expression may not be well-reflected in bulk skin tissue for many genes

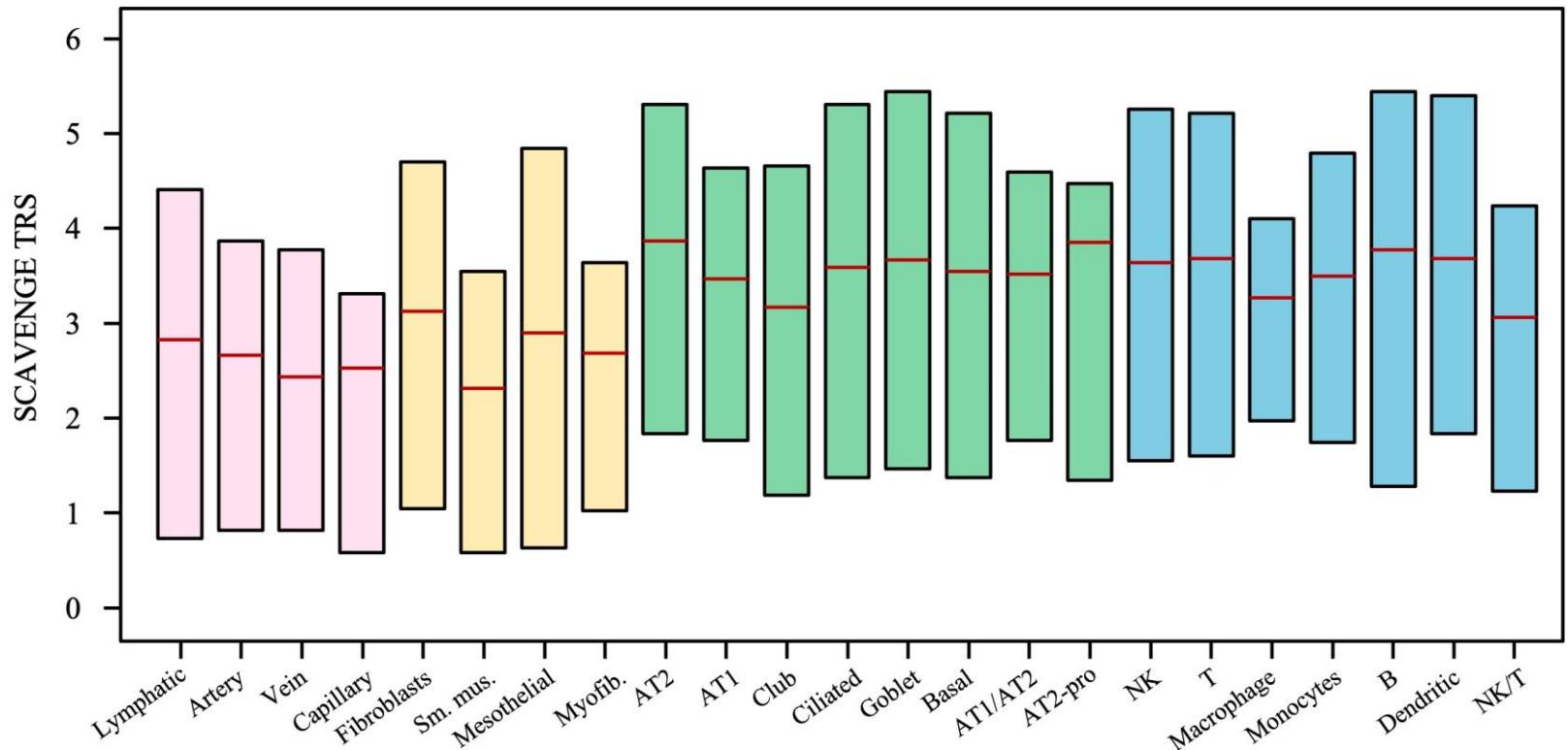
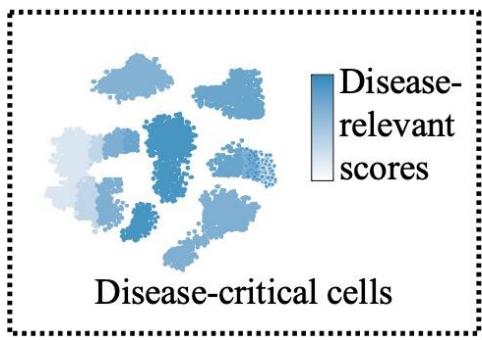
We can culture melanocytes and measure expression for QTLs!... But:

- Cultured cells -> artificial system
- Not possible for all cell types
- Cell of origin may not be known

# Cells of lung cancer origin are diverse



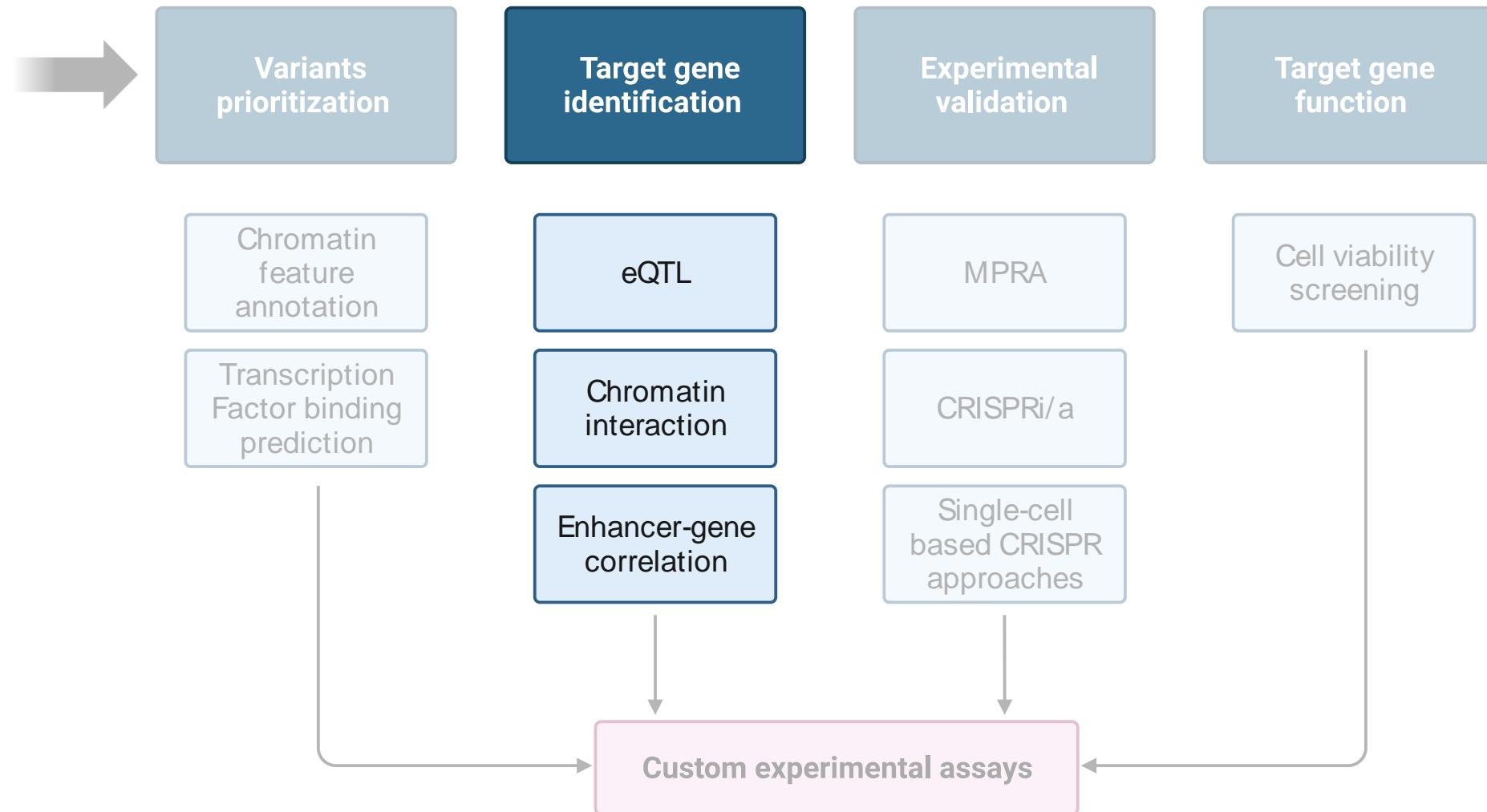
# Single-cell level trait relevance scores using single-cell ATACseq



# Variant prioritization tools

- SNP look-up tools with LD information
  - Ldlink (<https://ldlink.nih.gov/>)
  - Haploreg (<https://pubs.broadinstitute.org/mammals/haploreg/haploreg.php>)
  - Regulome DB (functional scores) (<https://regulomedb.org/regulome-search/>)
  - FORGEdb (functional scores) (<https://forgedb.cancer.gov/>)
- Public epigenomic resources
  - ENCODE (<https://www.encodeproject.org/>)
    - <http://encodec.encodeproject.org/>
    - <https://screen.encodeproject.org/>
  - RoadMap Epigenomics Project
    - <http://www.roadmapepigenomics.org/>
  - Data from both are also housed in the UCSC Genome Browser
    - <https://genome.ucsc.edu/>
- Single cell-level resources
  - <http://catlas.org/humanenhancer/#/>

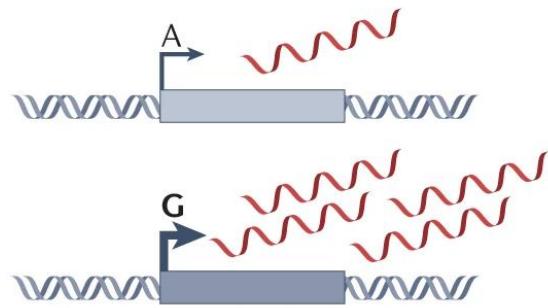
# Identification of "causal" variants and target genes



# Types of molecular QTLs:

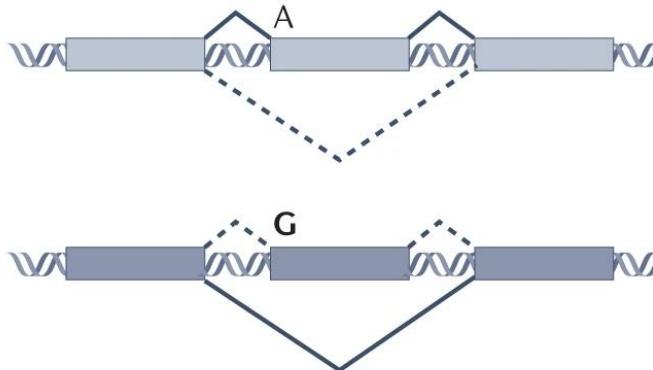
## Expression QTL (eQTL)

RNA expression level of a gene or a transcript



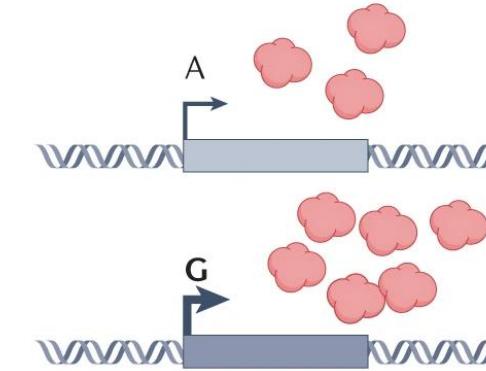
## Splicing QTL (sQTL)

Inclusion ratio of an exon, ratio of transcript levels or intron length



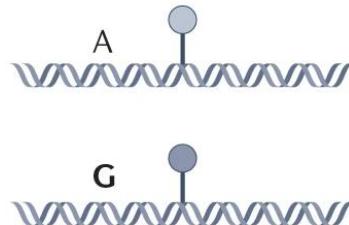
## Protein QTL (pQTL)

Protein expression level of a gene



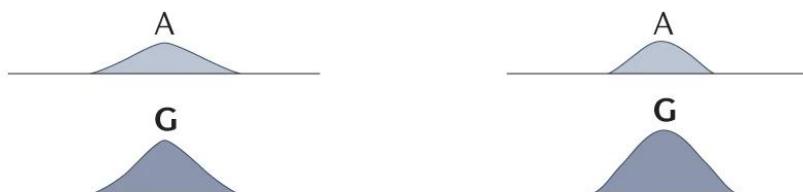
## Methylation QTL (meQTL)

Methylation ratio of a CpG site



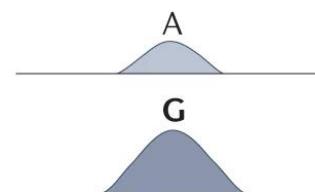
## Chromatin accessibility QTL (caQTL or chQTL)

Chromatin accessibility measured by ATAC-seq, DNase I sensitivity, etc.



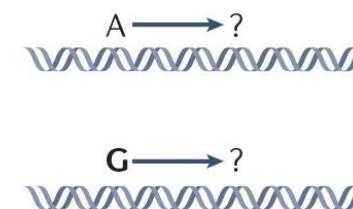
## Histone modification QTL (hQTL or cQTL)

Histone mark ChIP-seq peak height



## Molecular QTL (molQTL)

Any molecular trait with a locus in the genome



# Target gene identification tools

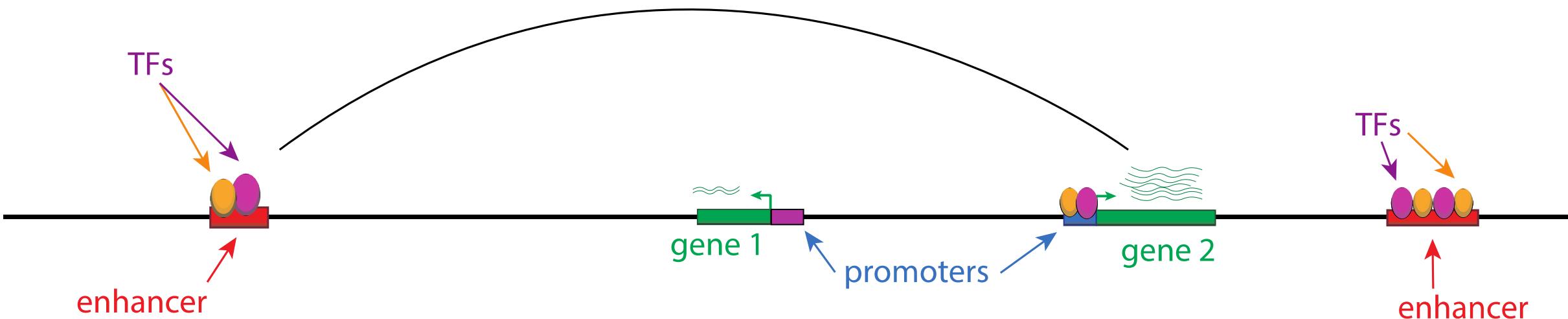
Some eQTL resources:

- GTEx portal (<https://gtexportal.org/home/>)
- eQTLgen (<https://www.eqtldgen.org/>)
- psychENCODE (<http://resource.psychencode.org/>)
- TCGA eQTL ([http://gong\\_lab.hzau.edu.cn/PancanQTL/](http://gong_lab.hzau.edu.cn/PancanQTL/))
- Single cell eQTLs (<https://eqtlgen.org/sc/>)

QTLS + Colocalization:

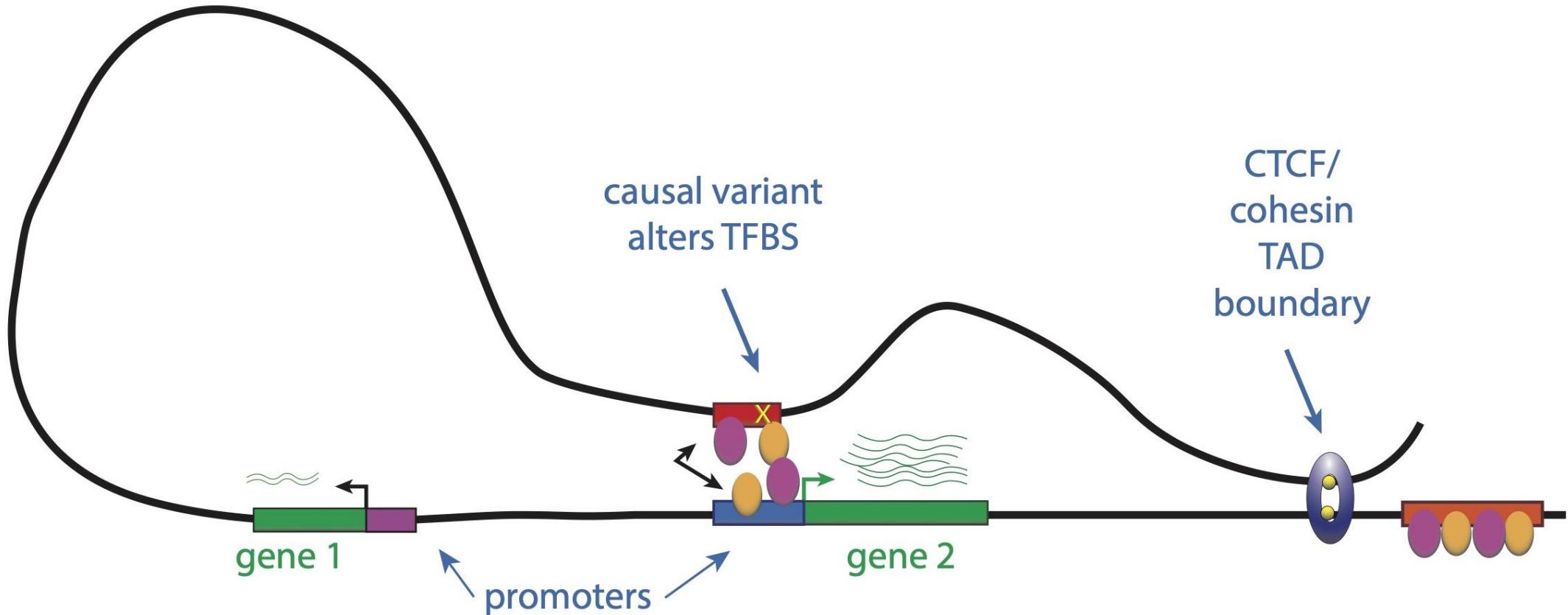
- ezQTL (<https://analysistools.cancer.gov/ezqtl/>)
- LocusFocus (<https://locusfocus.research.sickkids.ca/>)
- LocusCompare (<http://locuscompare.com/>)

# Enhancer-gene chromatin interactions



The genome is not linear!

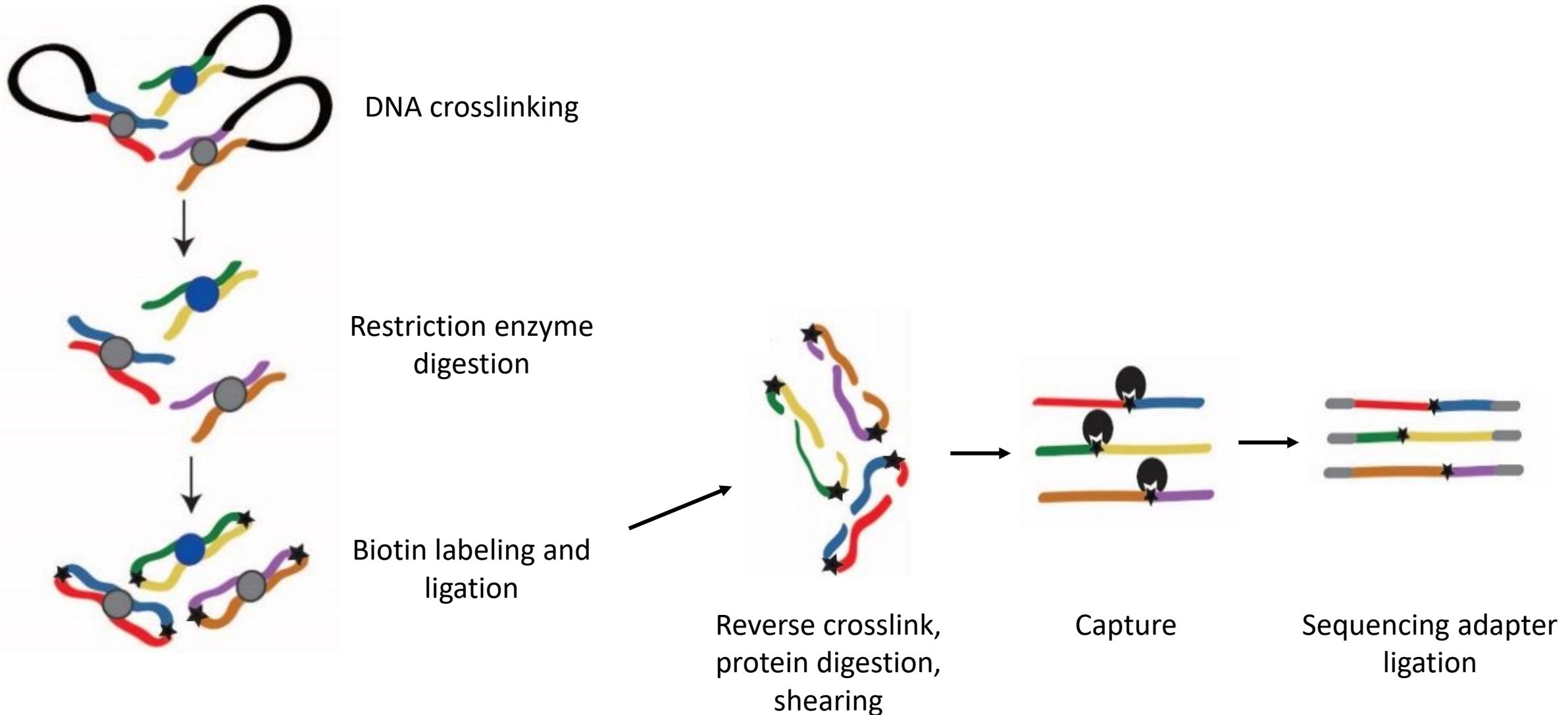
# The non-linear genome



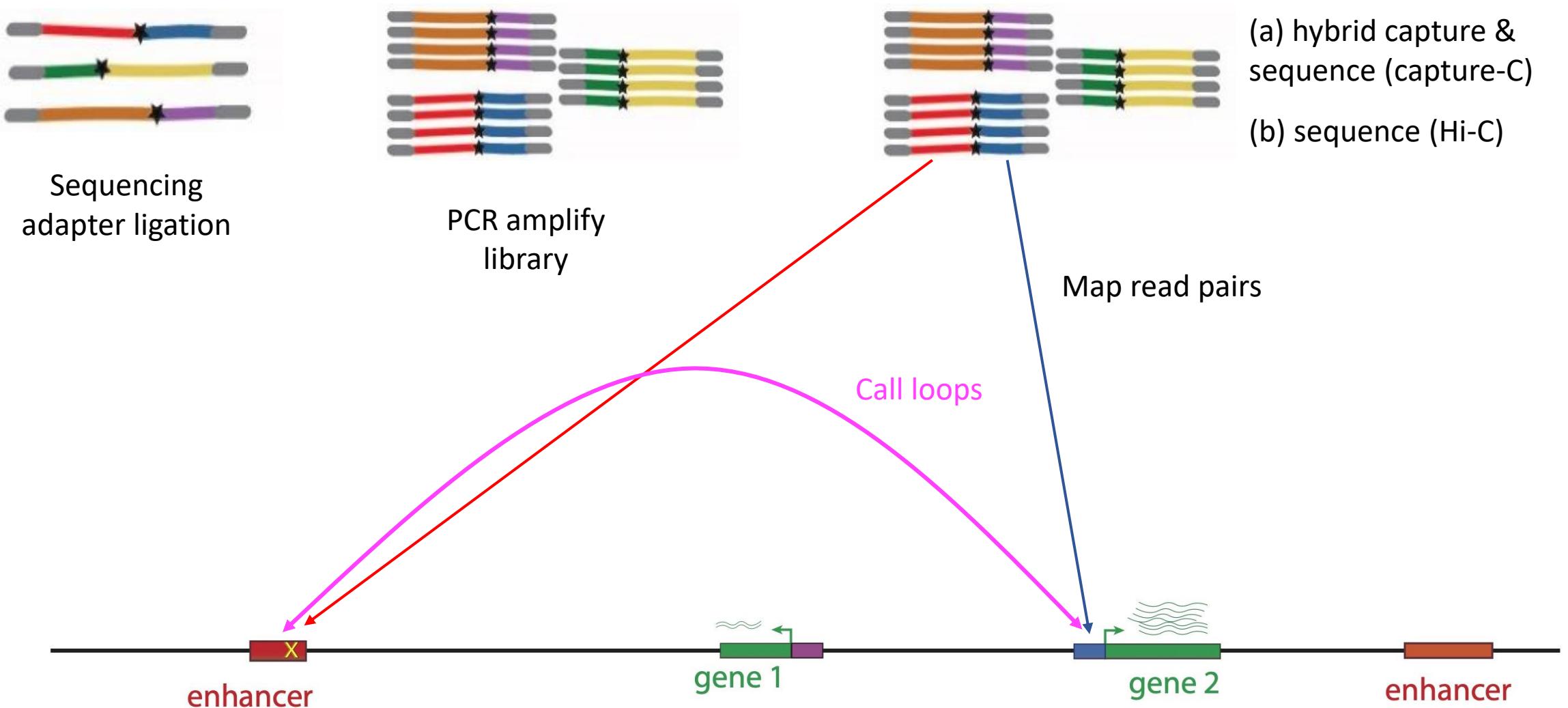
Features of regulatory enhancer-gene interactions:

- enhancer and gene in close physical proximity
- enhancer activity (chromatin openness, histone marks) correlated with gene expression

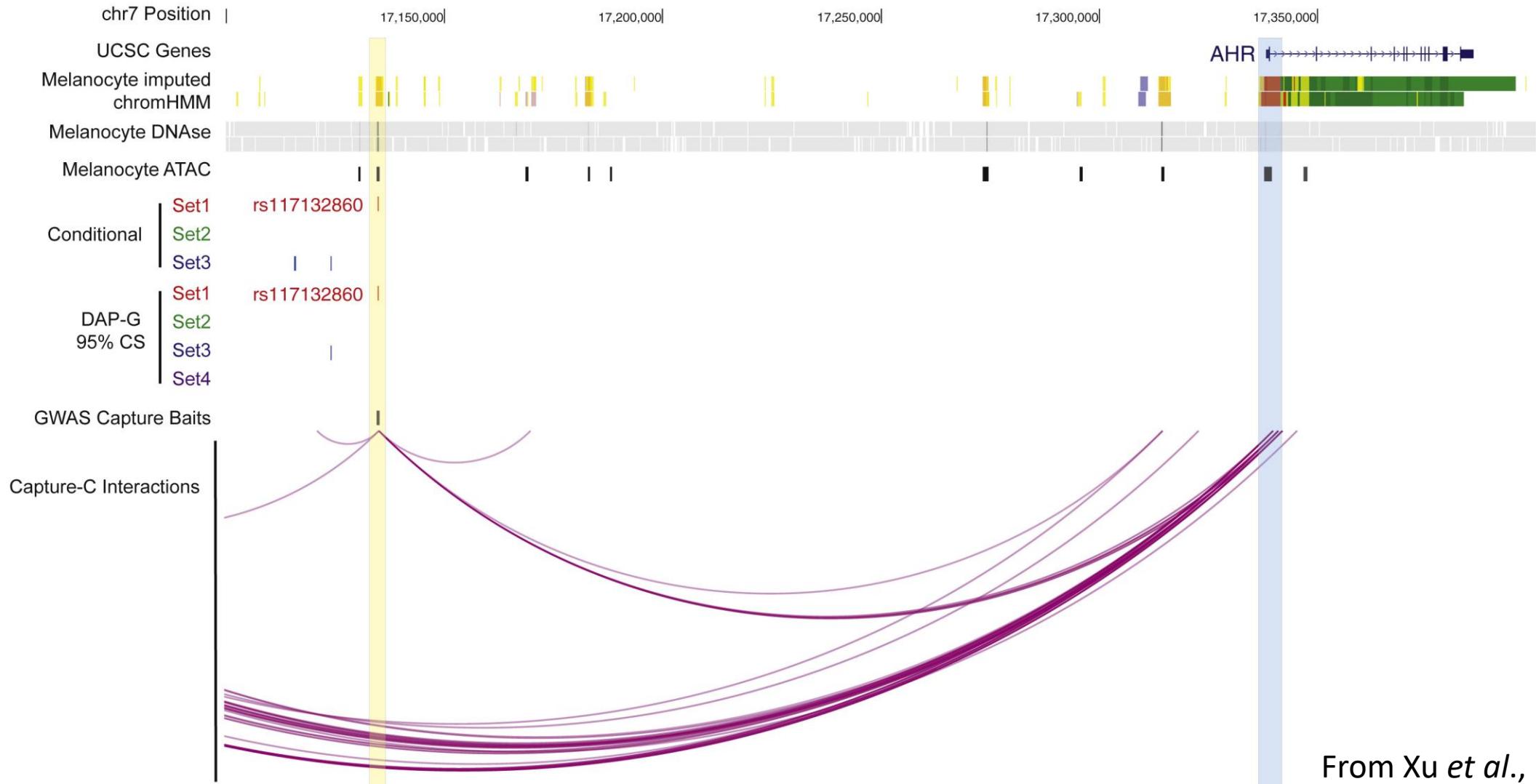
# Chromatin conformation-based methods to identify enhancer-gene interactions



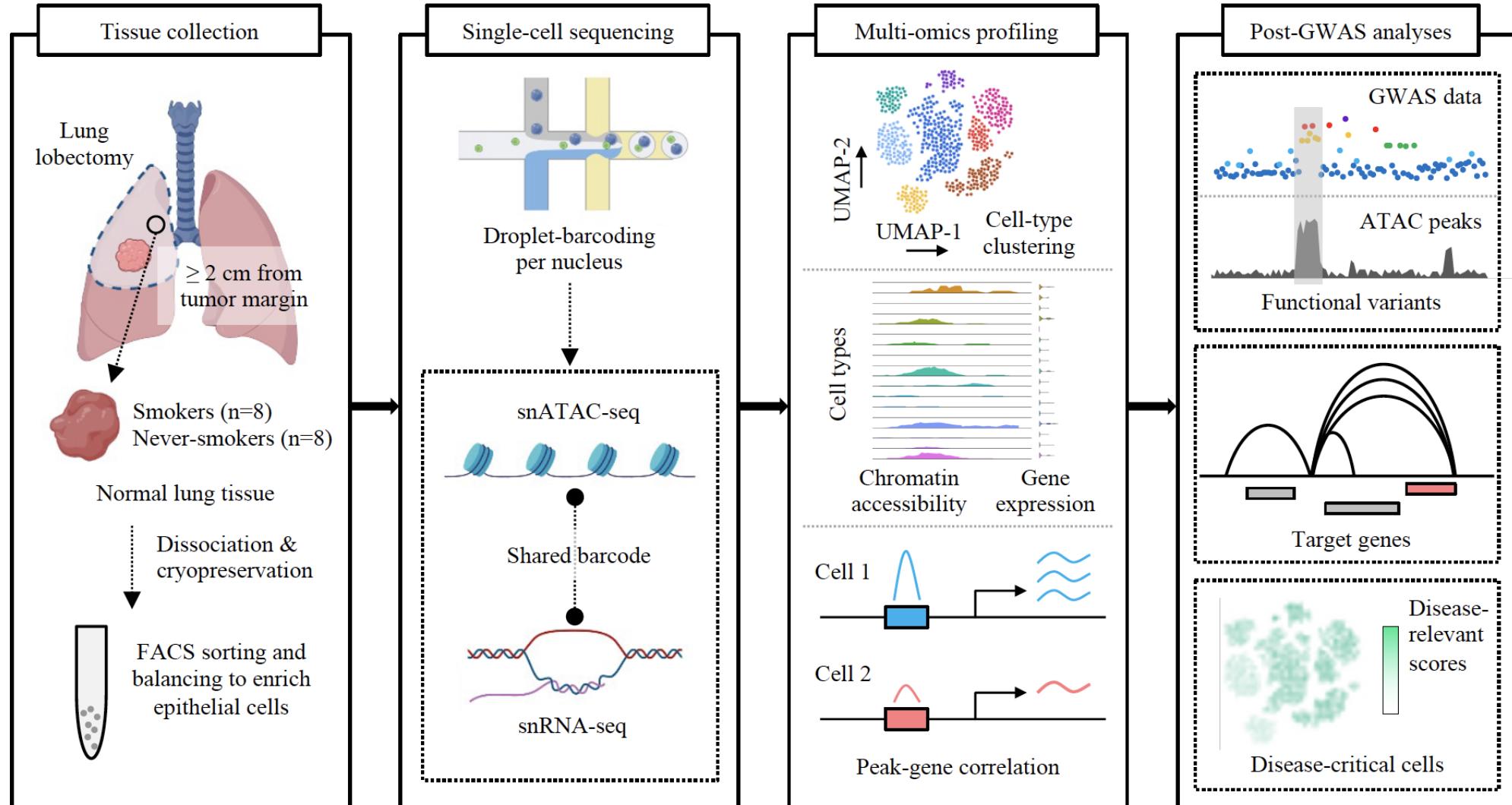
# Sequencing-based chromatin conformation assay (HiC/Capture-C)



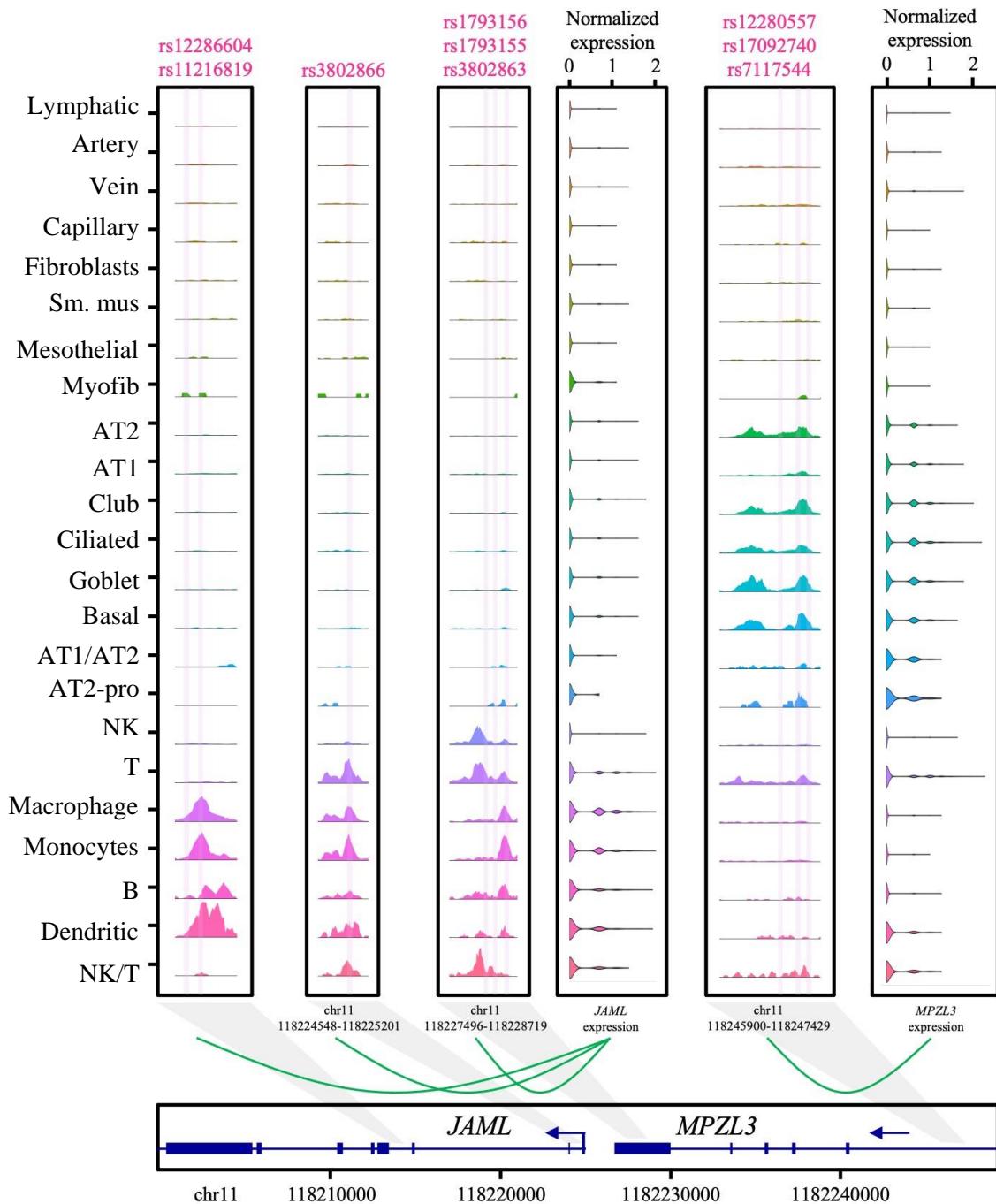
# Example: capture-C in melanocytes identifies *AHR* as a target for a functional melanoma risk variant



# Single-cell ATACseq/RNaseq for variant-gene connection



# Cell type-specific target genes in a single locus



# Public resources for chromatin interaction and enhancer-gene correlation

## Bulk cell or tissue data:

3D Genome Browser (<http://3dgenome.fsm.northwestern.edu/>)

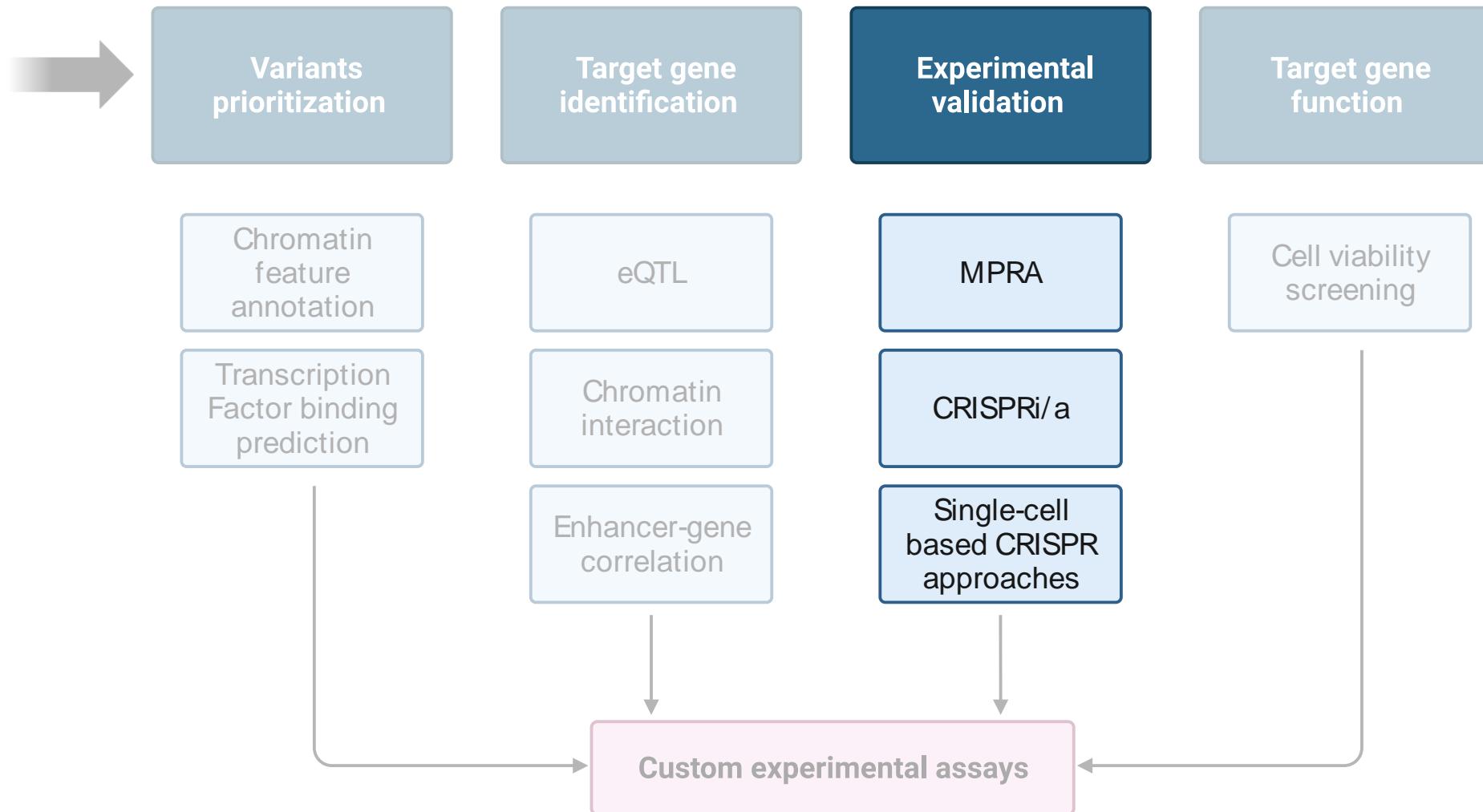
- Incorporates chromatin conformation data from many cell types
- Can assess interactions for specific genetic variants

## Single-cell data:

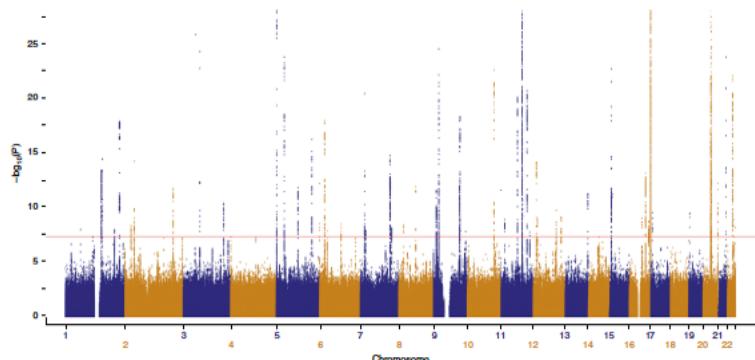
Activity-by-Contact (ABC) Browser (<https://www.engreitzlab.org/resources>)

- Incorporates both single-cell level enhancer-gene correlation and chromatin contact data

# Identification of "causal" variants and target genes

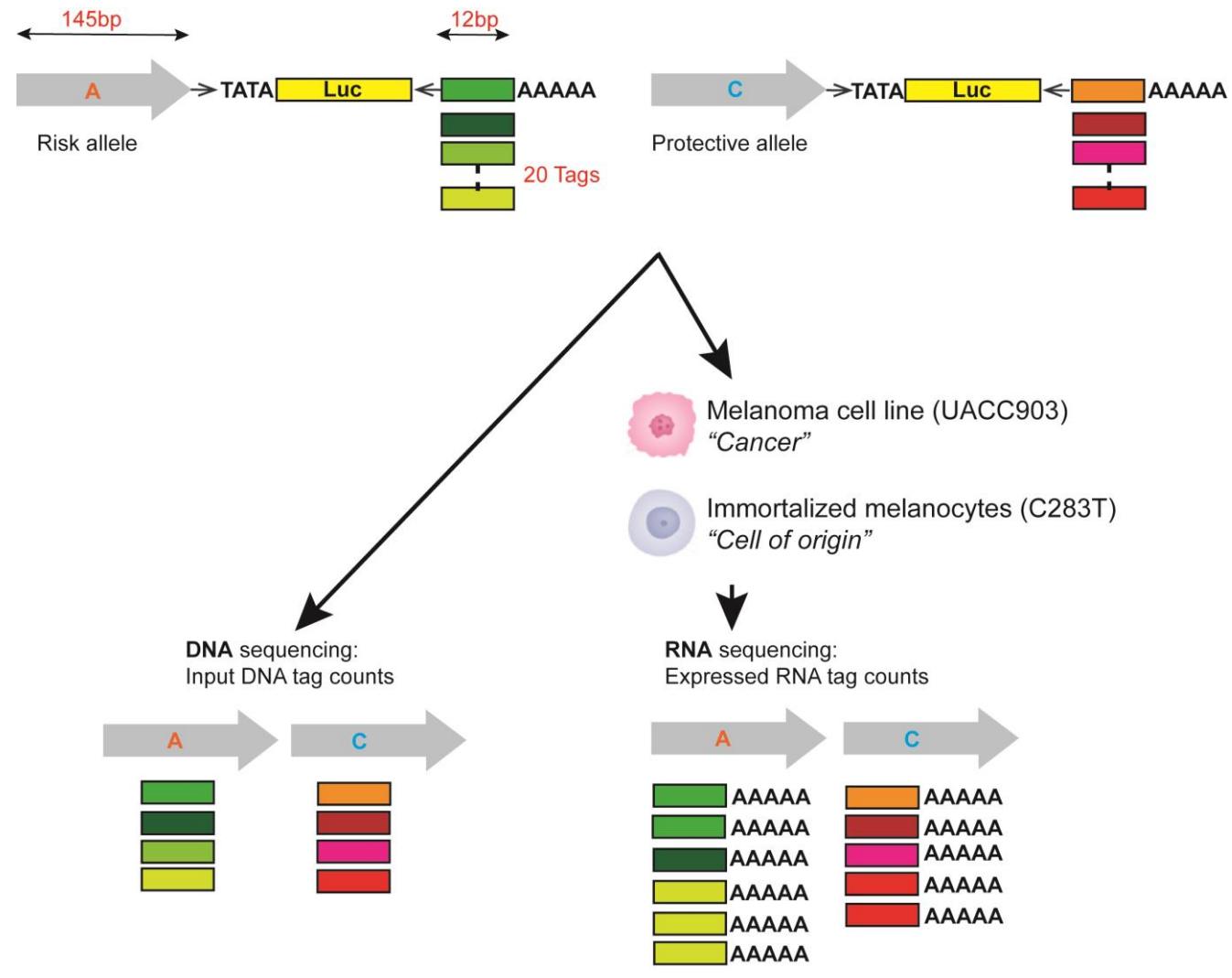
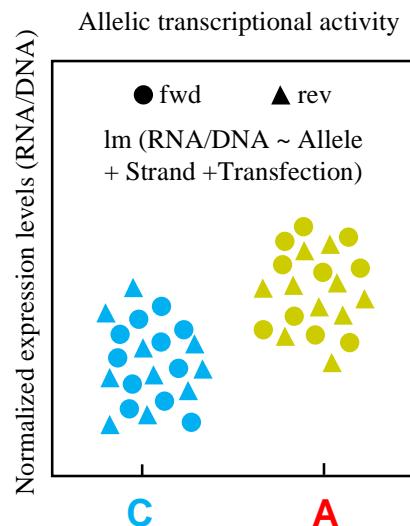


# Massively Parallel Reporter Assays (MPRA) can test allelic transcriptional activity GWAS variants



Melanoma GWAS  
(Landi et al., 2020 *Nature Genetics*)  
36,760 cases/ 375,188 controls

54 loci  
 $R^2 > 0.8$   
LLR < 1:1000  
1,992 variants  
~244,000 oligos



Long and Yin...Choi, AJHG, 2022

# MPRA-based resources

- MPRA
  - <https://www.varianteffect.org/resources> (region-specific, limited)

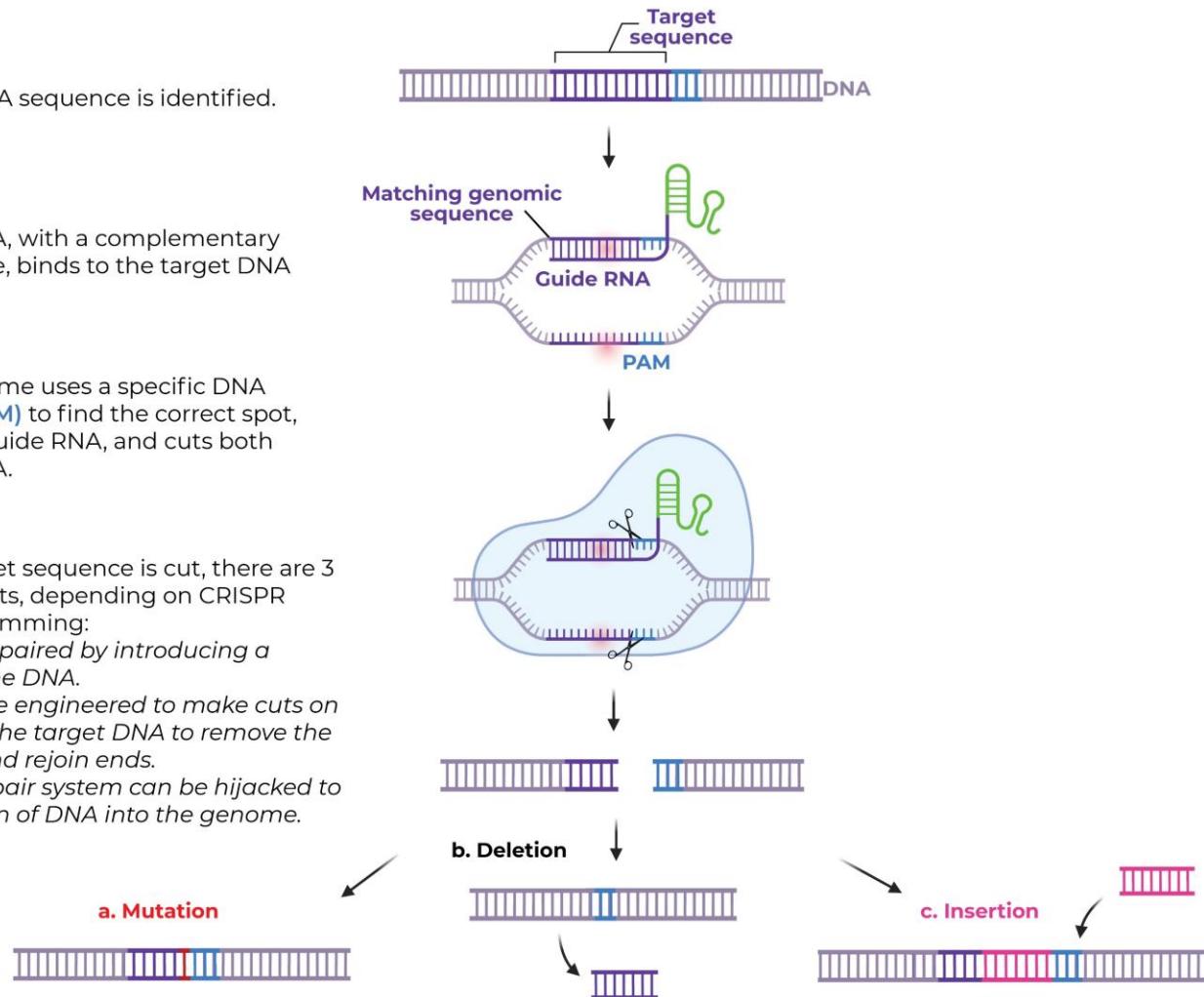
# CRISPR as a tool for discovery and validation

① The target DNA sequence is identified.

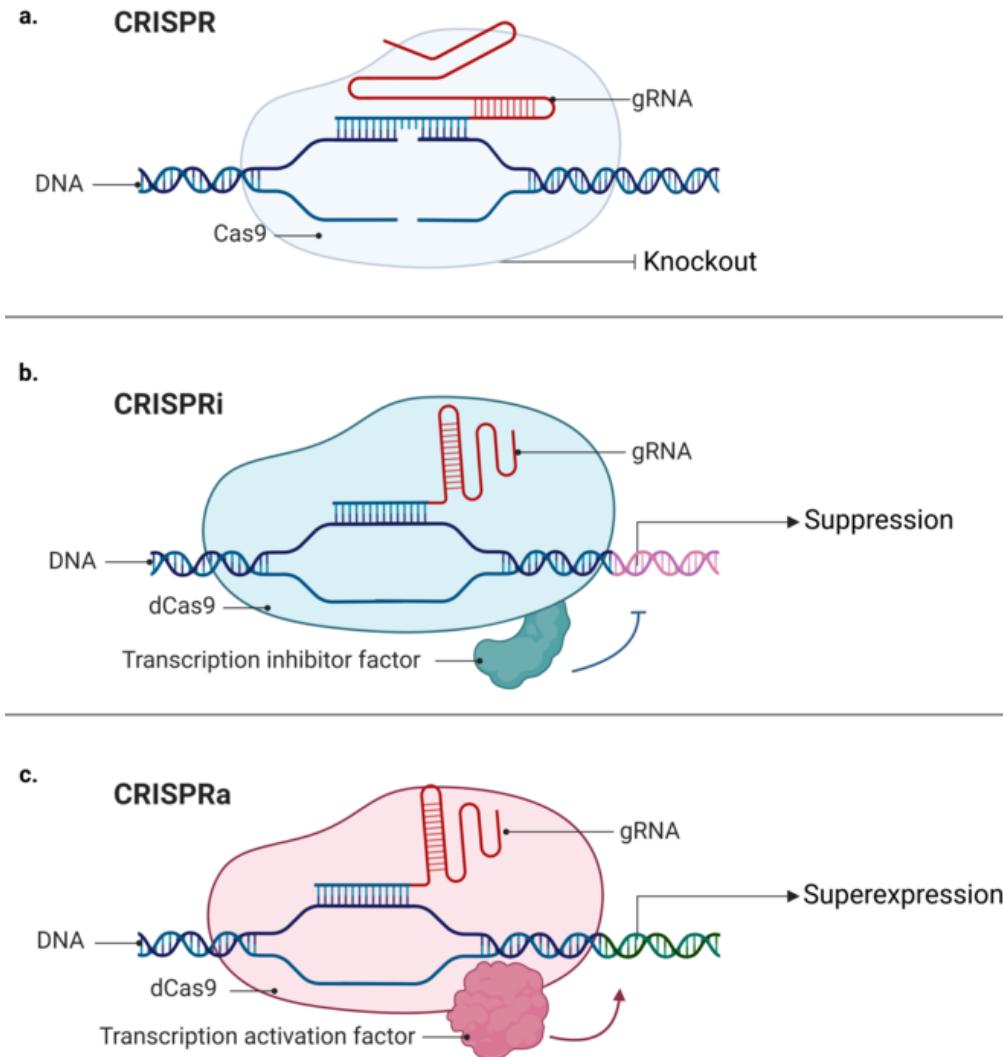
② The guide RNA, with a complementary DNA sequence, binds to the target DNA sequence.

③ The Cas9 enzyme uses a specific DNA sequence (PAM) to find the correct spot, binds to the guide RNA, and cuts both strands of DNA.

④ Once the target sequence is cut, there are 3 potential results, depending on CRISPR system programming:  
a. The cut is repaired by introducing a mutation in the DNA.  
b. Enzymes are engineered to make cuts on either side of the target DNA to remove the target DNA and rejoin ends.  
c. The DNA repair system can be hijacked to insert a section of DNA into the genome.



# CRISPR, CRISPR-inhibition, CRISPR-activation



**Knock-out:**

- SNP/regulatory region
- Gene (introduce truncating mutations)

**Inhibits:**

- SNP/regulatory activity
- Gene (targeting gene promoters)

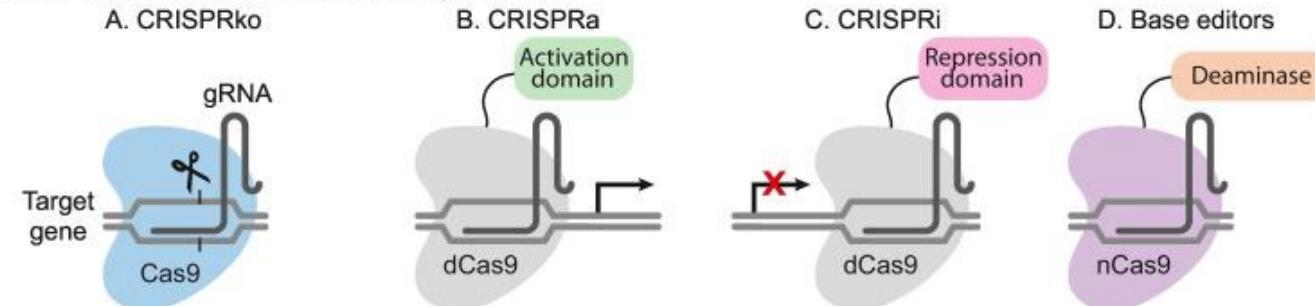
**Activates:**

- SNP/regulatory activity
- Gene (targeting gene promoters)

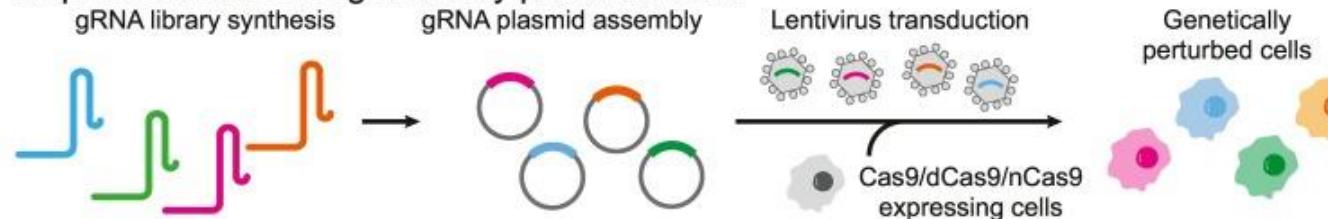
# Discovery: CRISPR screens

## (A) CRISPR screens

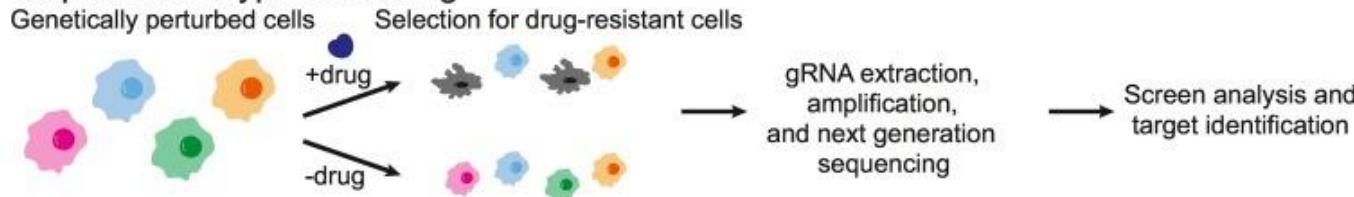
### Step 1: Selection of screening system



### Step 2: Generation of genetically perturbed cells



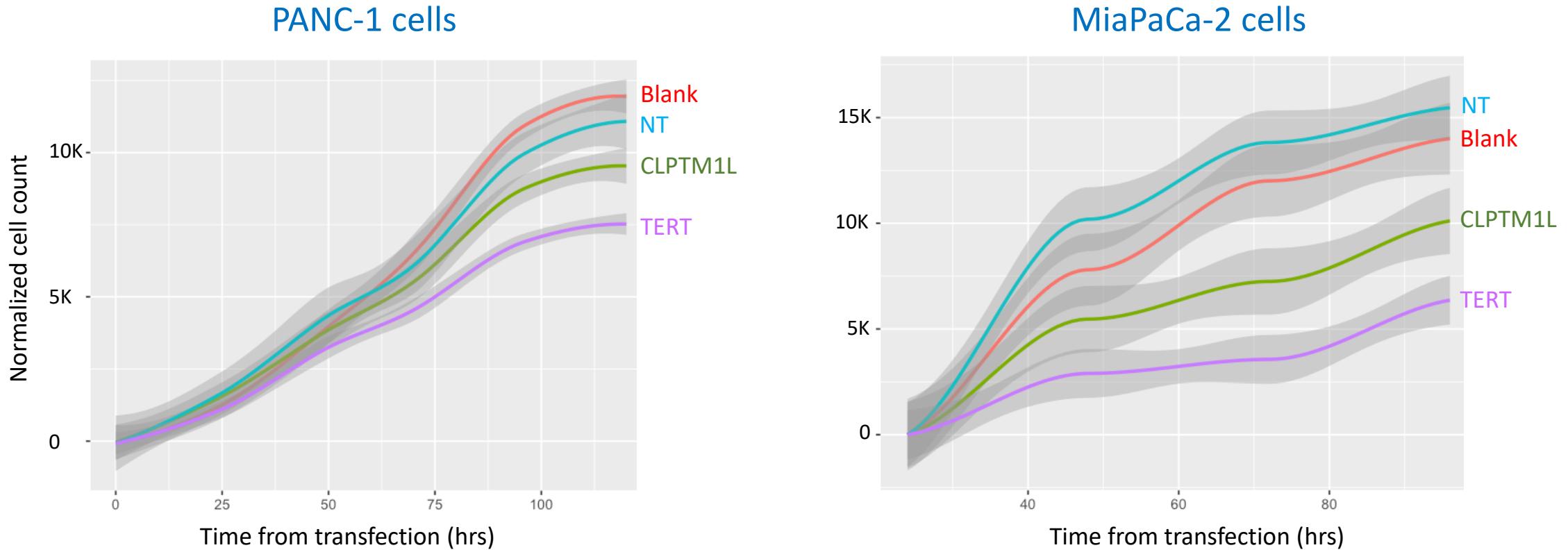
### Step 3: Phenotypic screening



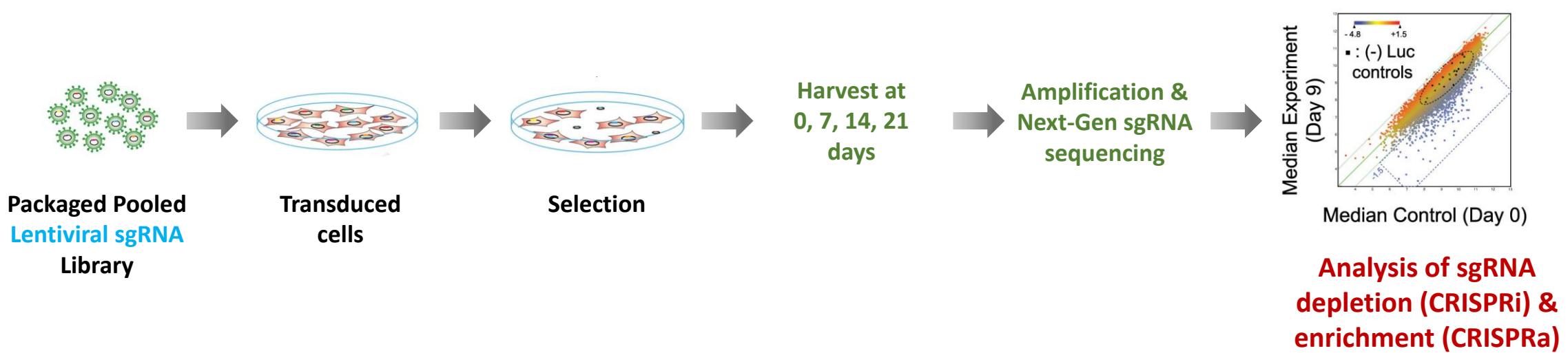
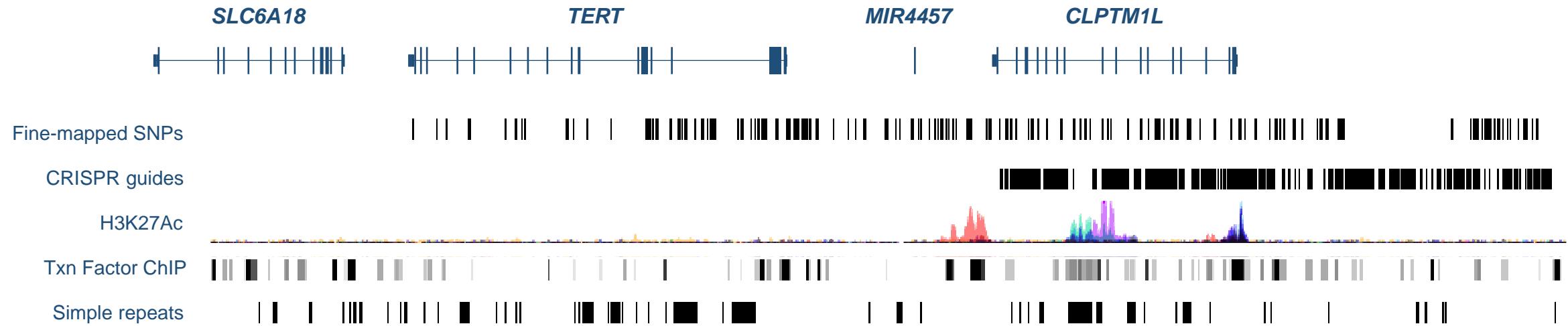
# CRISPRi Screen (variant based) example – chr5p15.33

**Hypothesis:** functional variants underlying chr5p15.33 signals mediate risk via noncoding regulatory effects on *TERT* and/or *CLPTM1L* expression leading to changes in cellular growth

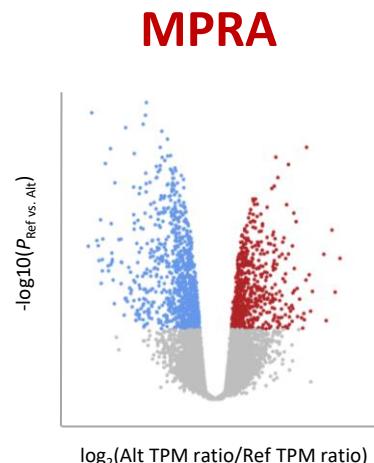
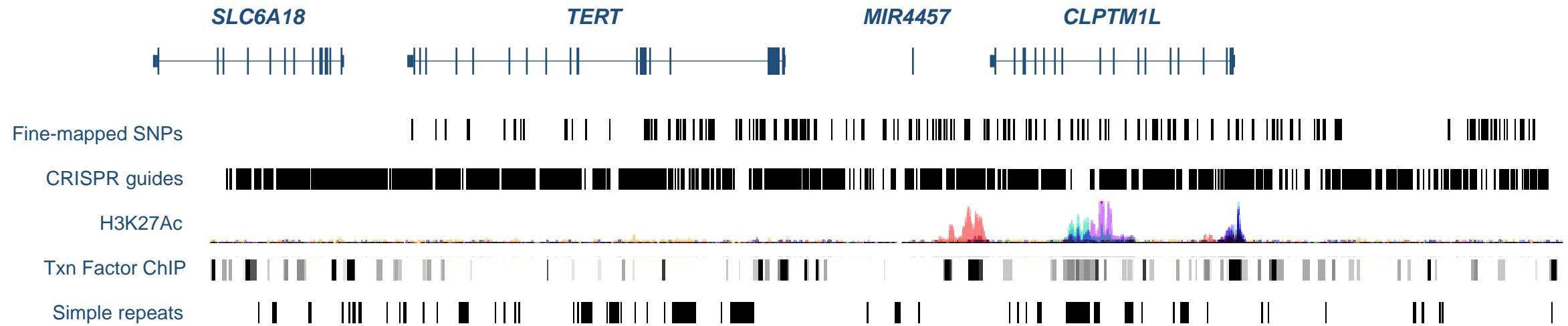
**Approach:** CRISPR and MPRA screens to target all known chr5p15.33 signals simultaneously



# Tiled CRISPRi screen across 10 GWAS signals at chr5p15.33



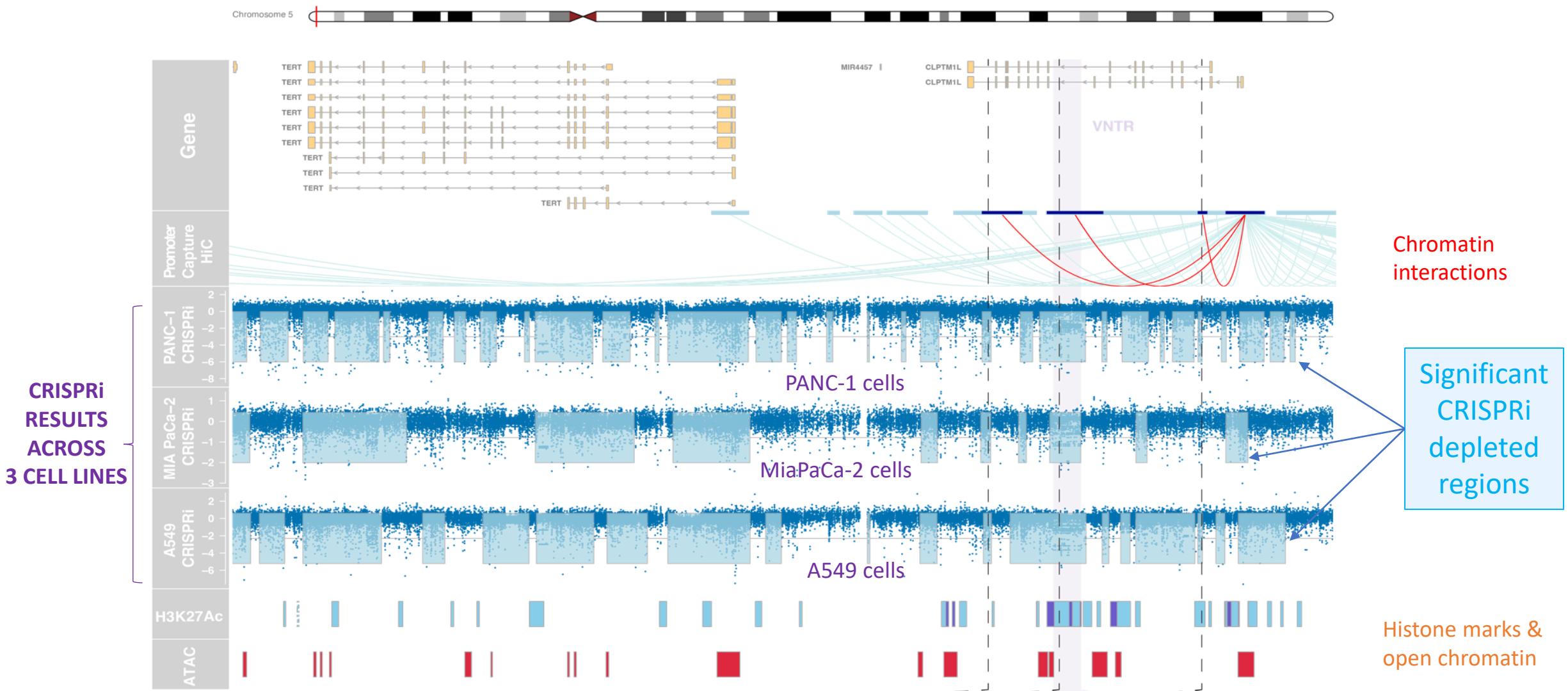
# Massively Parallel Reporter Assay (MPRA) screen across 5p15.33



**Assess allele specific gene regulatory effects across all signals**

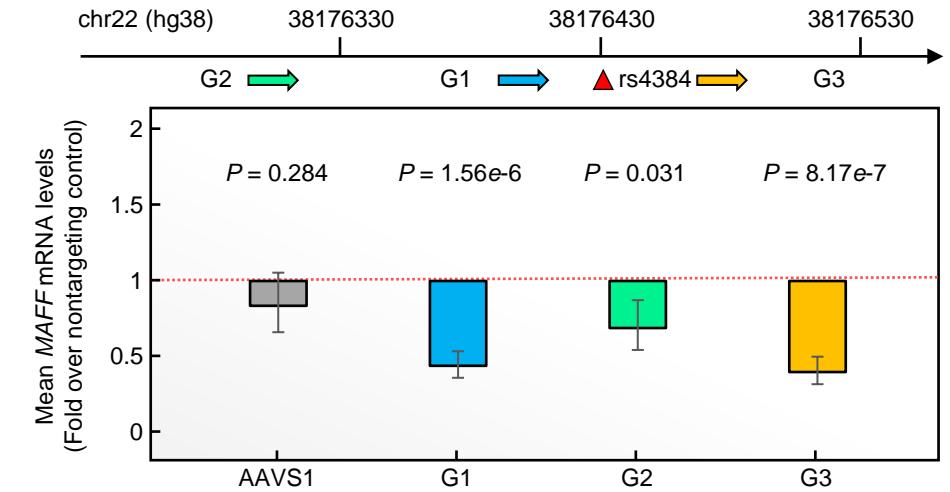
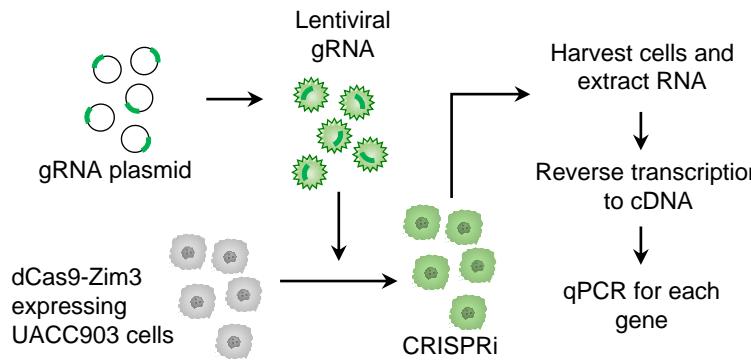
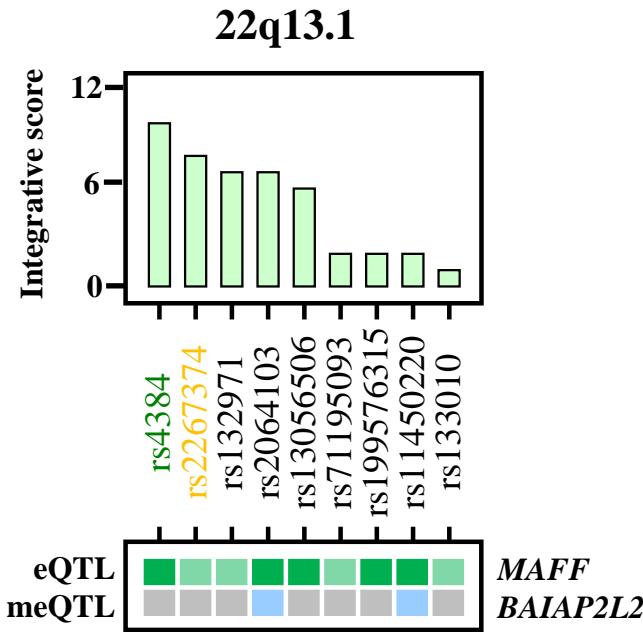
- 175 fine-mapped variants
- ~60,000 reporter constructs
- 2-3 cell lines each used for pancreatic cancer, melanoma, bladder and lung cancer

# Chr5p15.33 - CRISPRi screen results in pancreatic and lung cancer cells

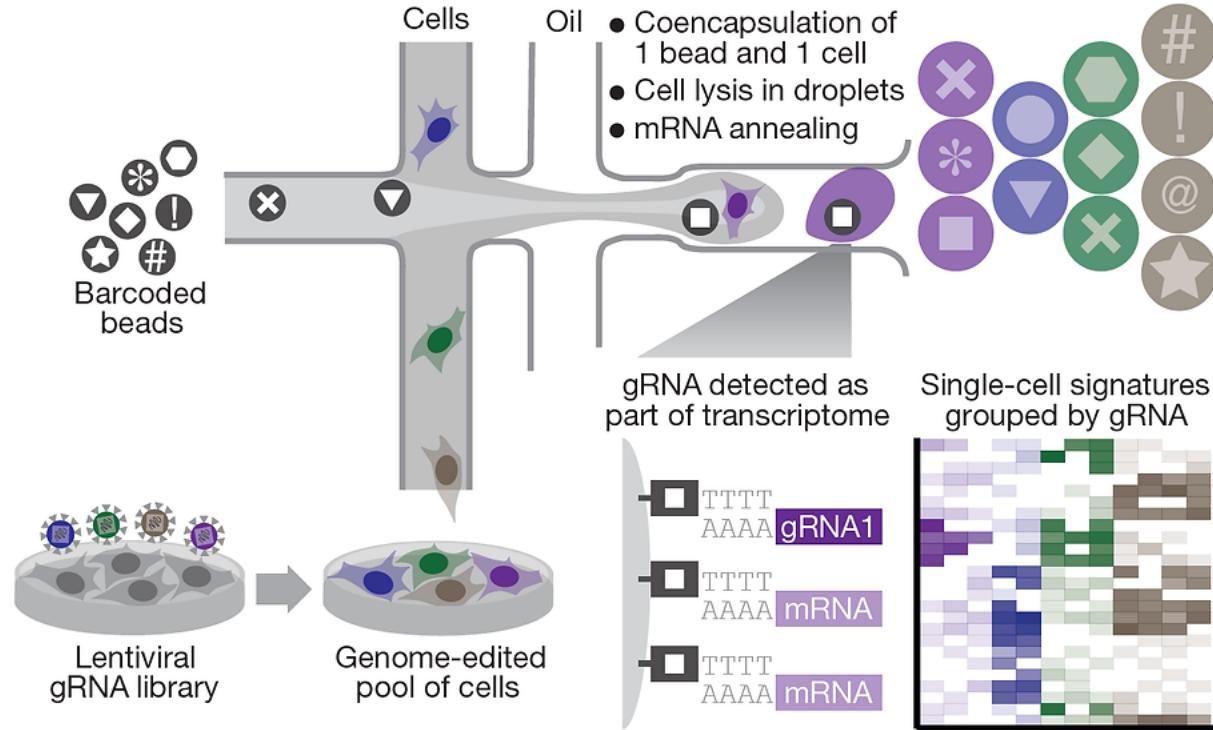


→ Next steps: individual CRISPRi validation for 10-30 SNPs to assess effects on *TERT* and *CLPTM1L* expression

# Validation: CRISPRi can validate variant gene connection

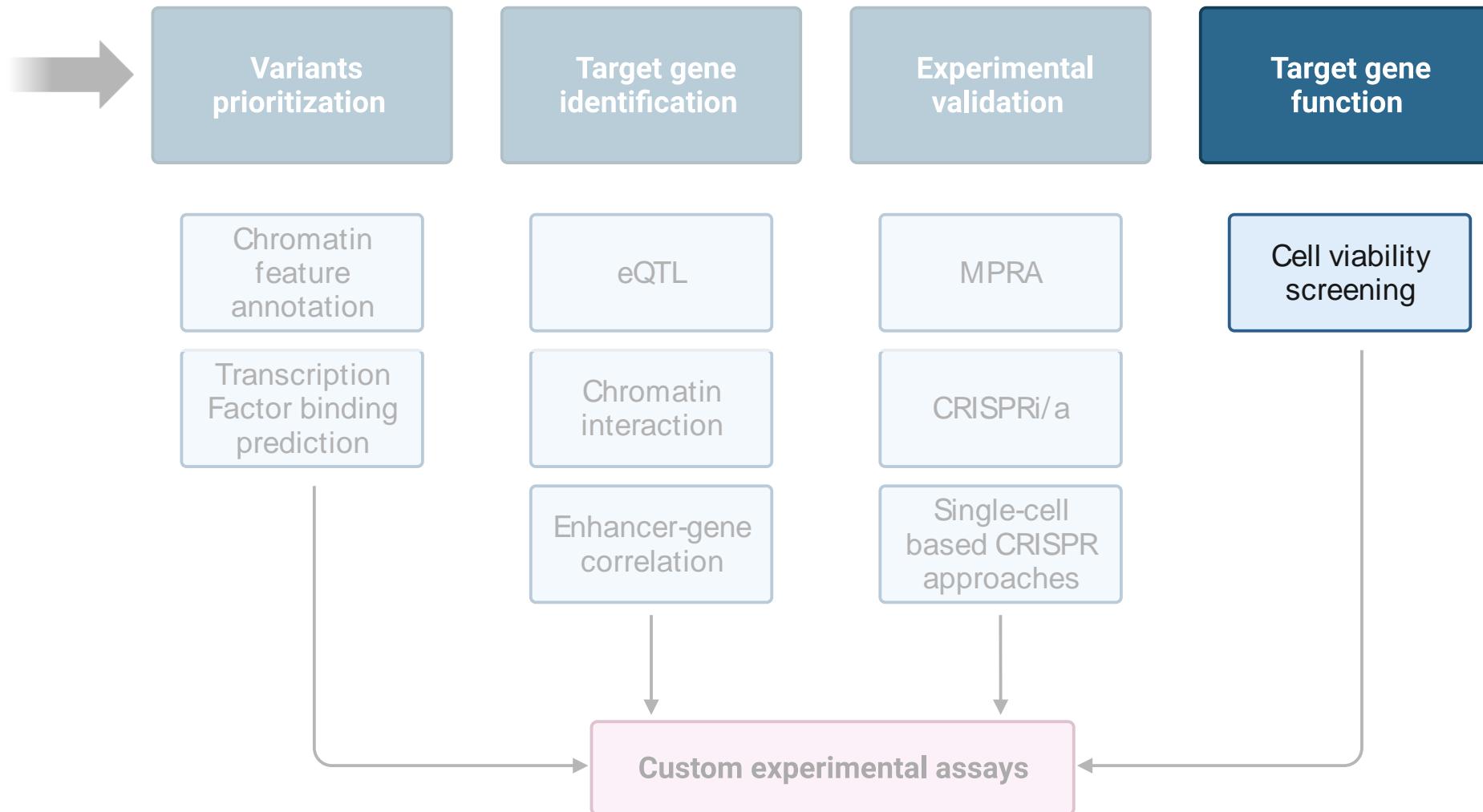


# CROP-seq (CRISPR screening with scRNAseq readout)



<https://www.bocklab.org/resources/crop-seq>

# Identification of "causal" variants and target genes



# Target Gene Function

## High-throughput (typically CRISPR) methods:

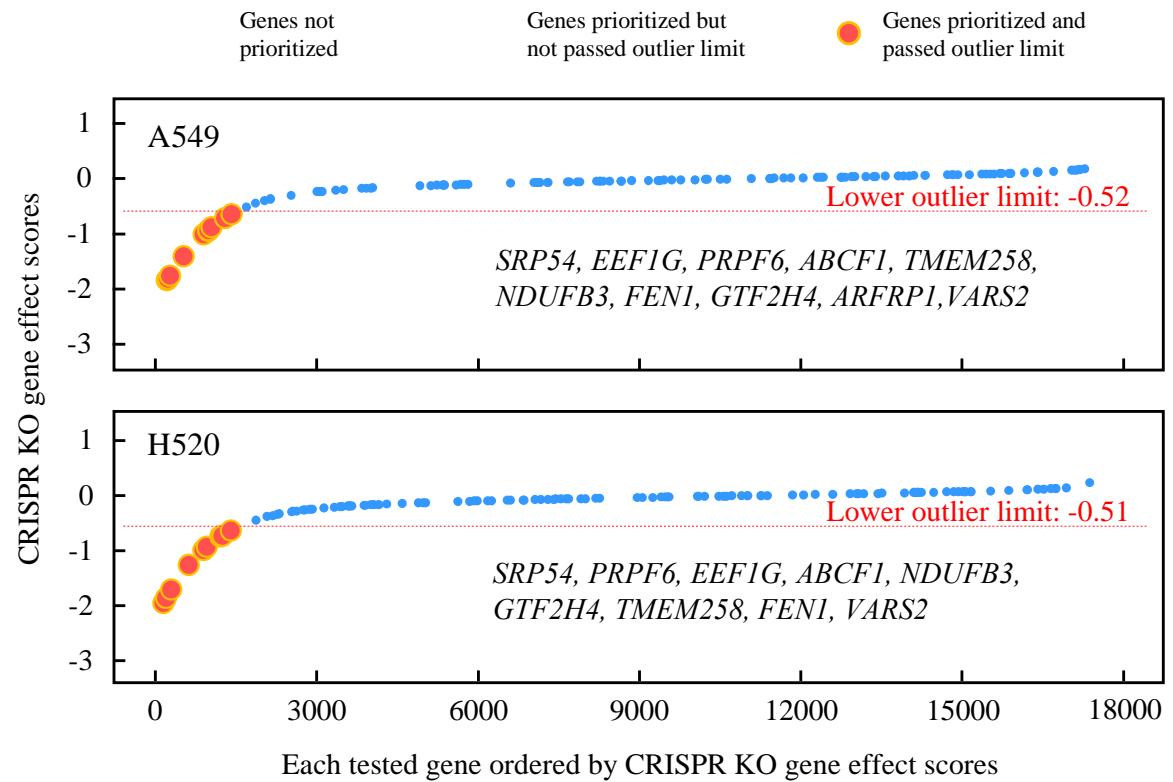
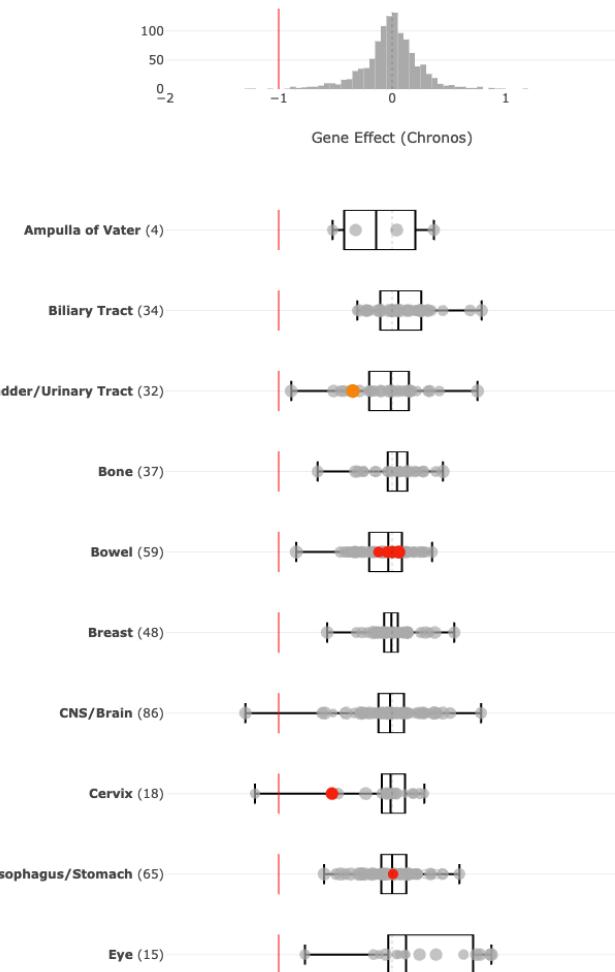
- Publicly-available screen data for effect on cell growth/viability
  - DepMap: Cancer Dependency Map Project (<https://depmap.org/portal/>)
- Custom pooled CRISPR-based screens -> need selectable phenotype

## Individual gene assessment:

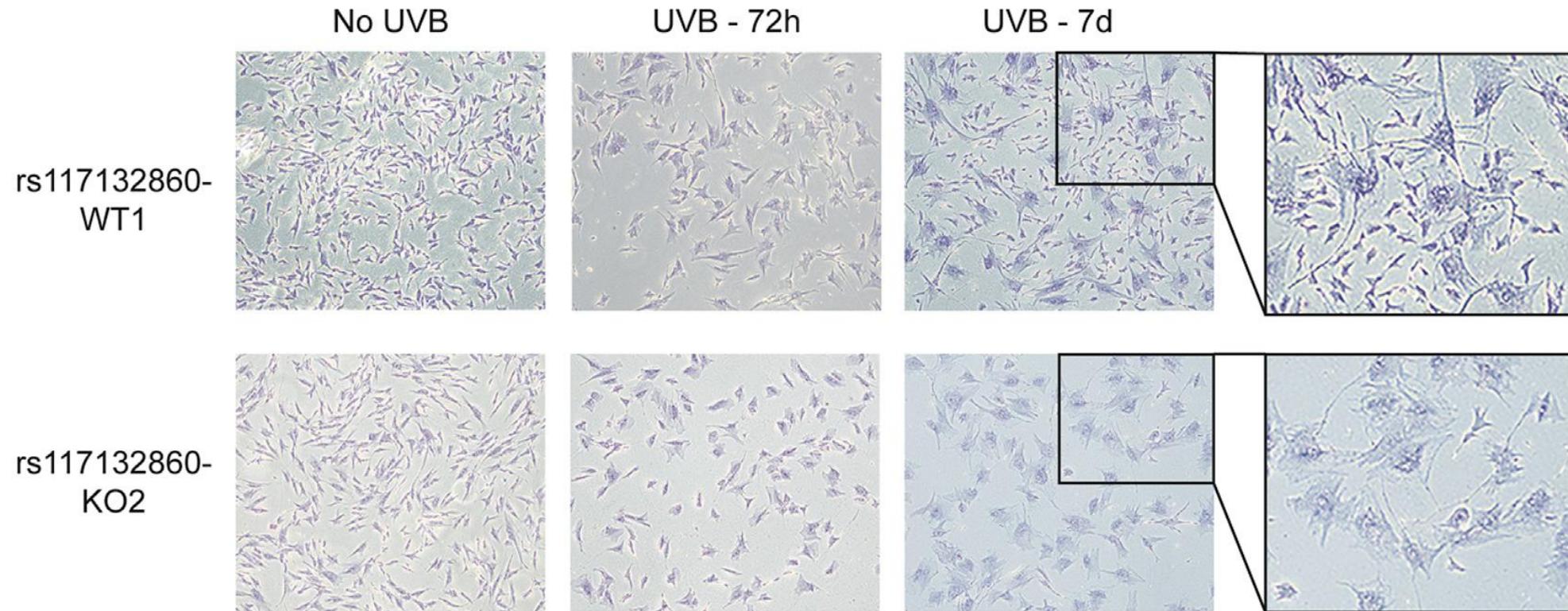
- Typically over-expression/knockdown and assess cellular phenotypes
- There is no one-size-fits-all assay
- Cell type and cell context matters!

# Publicly available CRISPR viability screen data

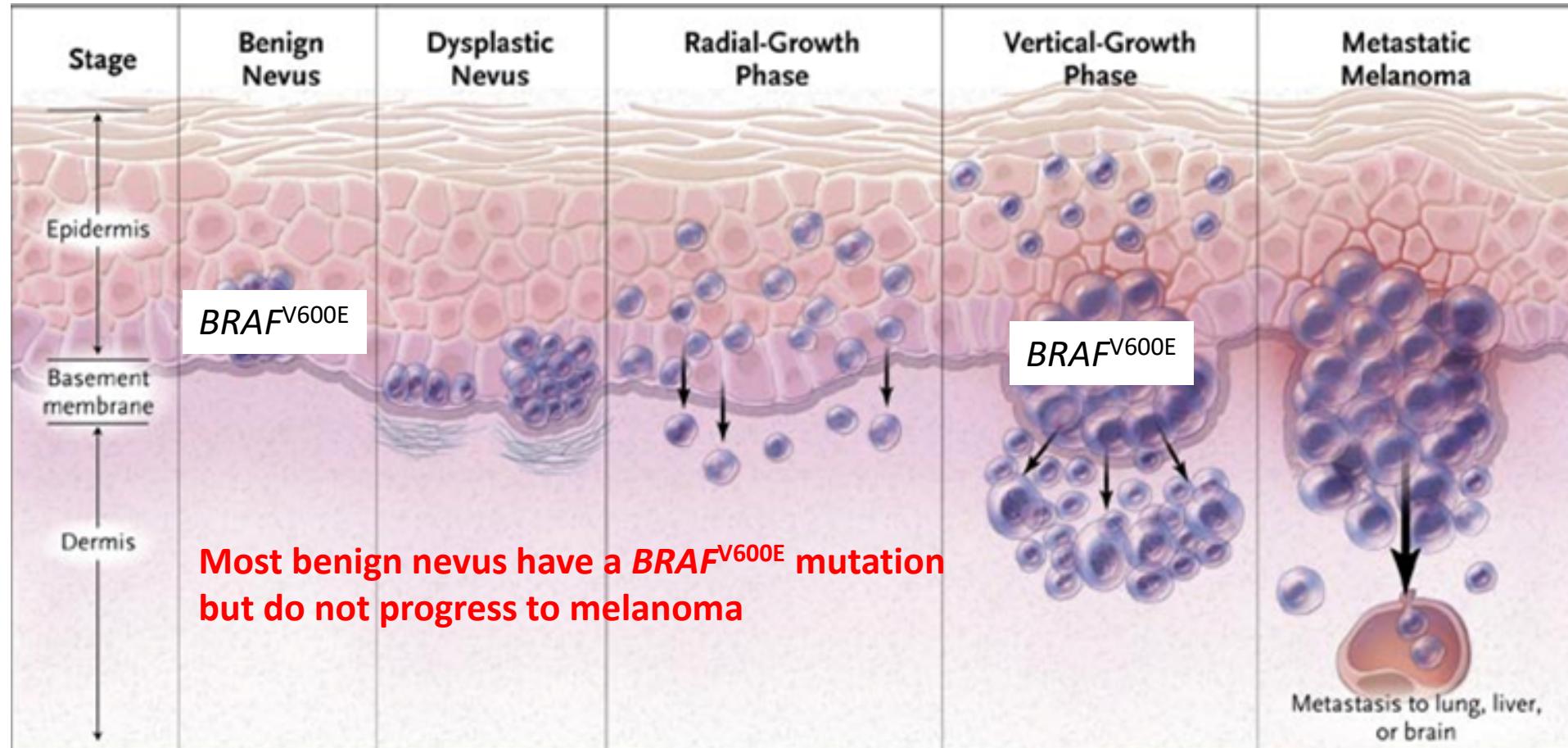
- DepMap cell viability screening (<https://depmap.org/portal/>)



# UV exposure-specific function of *AHR* melanoma-risk variant



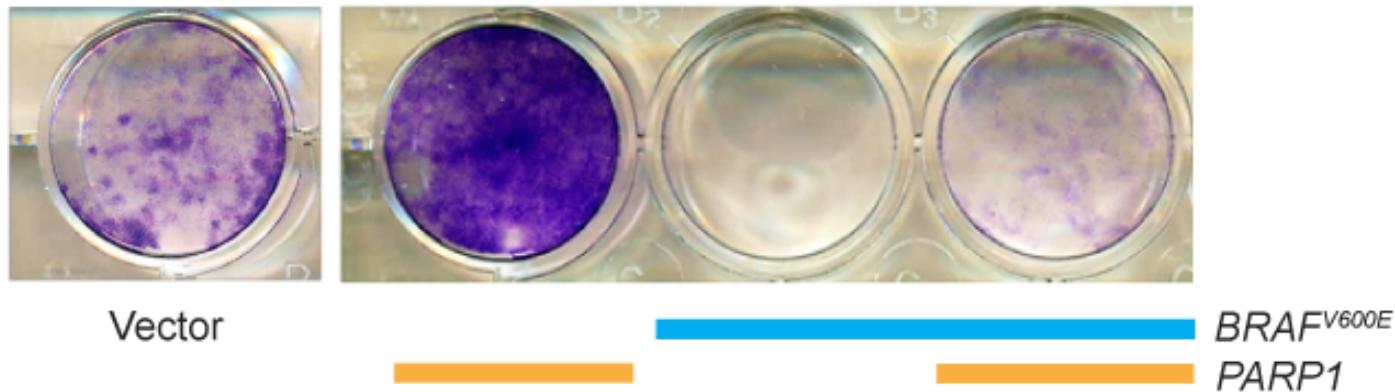
# $BRAF^{V600E}$ -induced senescence is a barrier to melanoma progression



# PARP1 function: promotes melanocyte growth, bypass senescence, and transforms immortalized melanocytes

✓ Does PARP1 help bypass  $BRAF^{V600E}$ -induced senescence in primary melanocytes?

Primary melanocytes  
growing attached to  
culture dish  
**More purple: more  
cell growth**



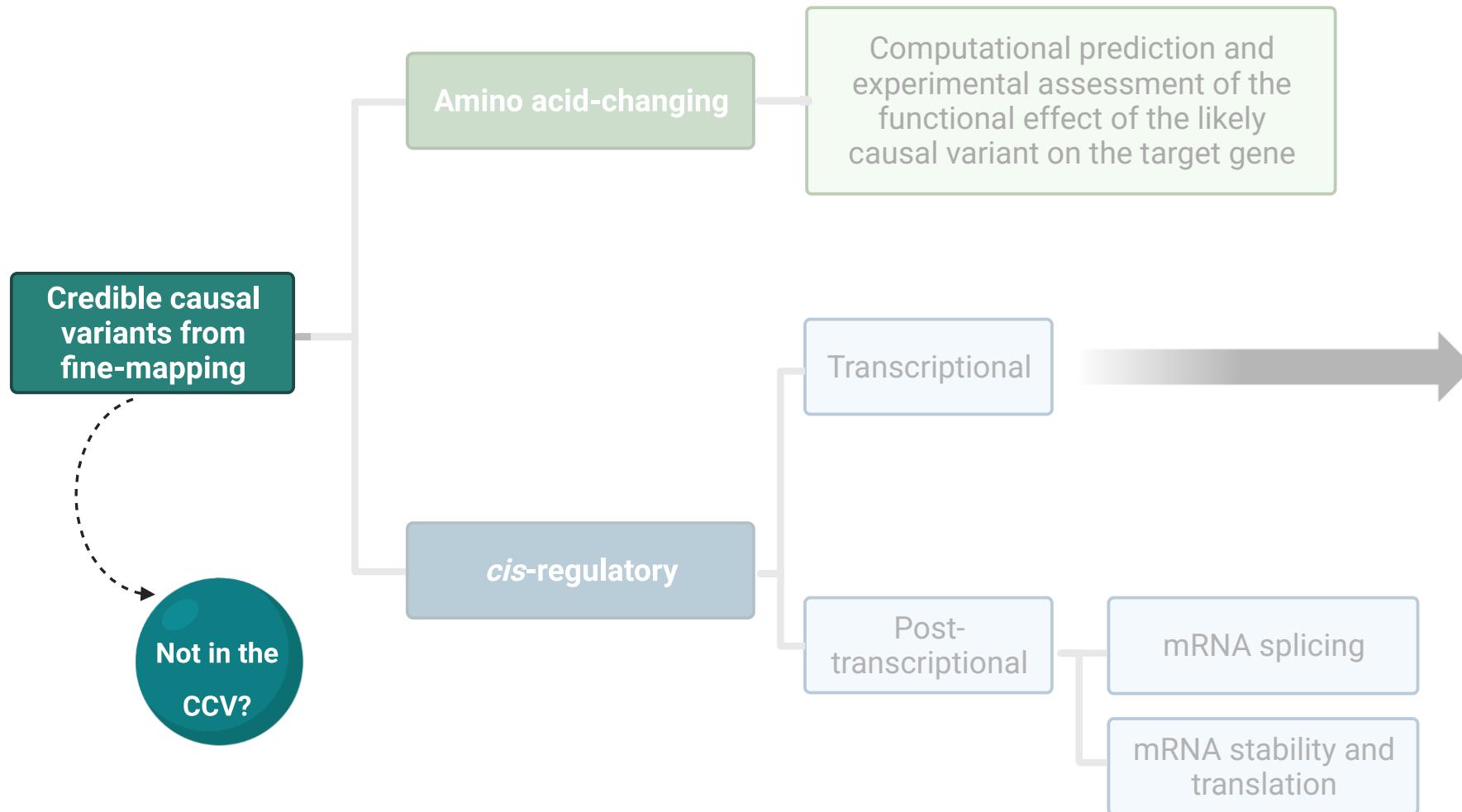
✓ Does PARP1 promote malignant transformation of immortalized melanocytes?

Soft agar assay  
(tumorous potential by growth  
in suspension, unattached to  
any matrix  
**More dots of colonies: more  
transformation)**



Mai Xu (LTG)

# Mechanisms of variant function



# From statistics to function

Genetic variants directly genotyped

Genetic variants imputed based on various panels and tools

Which of the available variant(s) can statistically account for the signal?



Predicted functionality of the candidate variants

Assumption:  
We know all the variants in the region, they can be scored or well imputed in relevant sets

Examples

Assumption:  
statistical association  
=functionality=causality;  
Single vs. multiple variants (haplotypes) that might need to work together

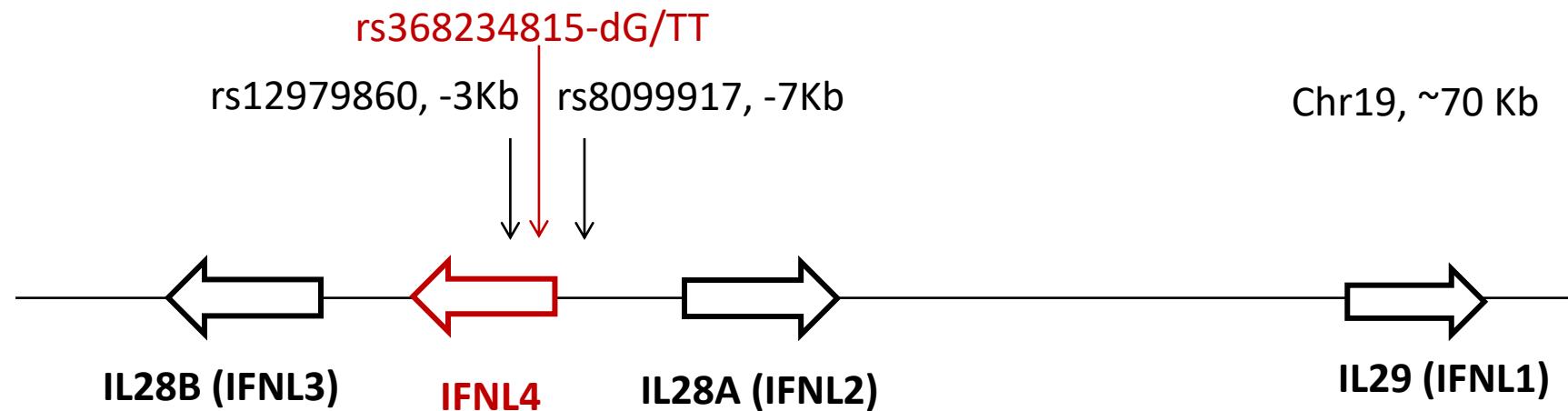
Examples

# What is functional? Almost anything at least in some conditions. But is some functionality somewhere = causality for a specific phenotype?

- Level:
- mRNA – eQTLs, stability, modifications, splicing, polyadenylation
- DNA – binding sites for various regulatory factors, methylation
- Chromatin – open/closed regions, long-range interactions
- Protein – pQTLs, stability, localization, modifications, etc
- Derived phenotypes – anything that can be quantified
- Somatic mutagenesis and mutation signatures; immune cell ratios, cell viability, proliferation, senescence or stemness scores, relative telomere length, hormone ratios, etc

# Not all variants are SNPs and we still don't know all of them.

## Example 1: A small polymorphic indel



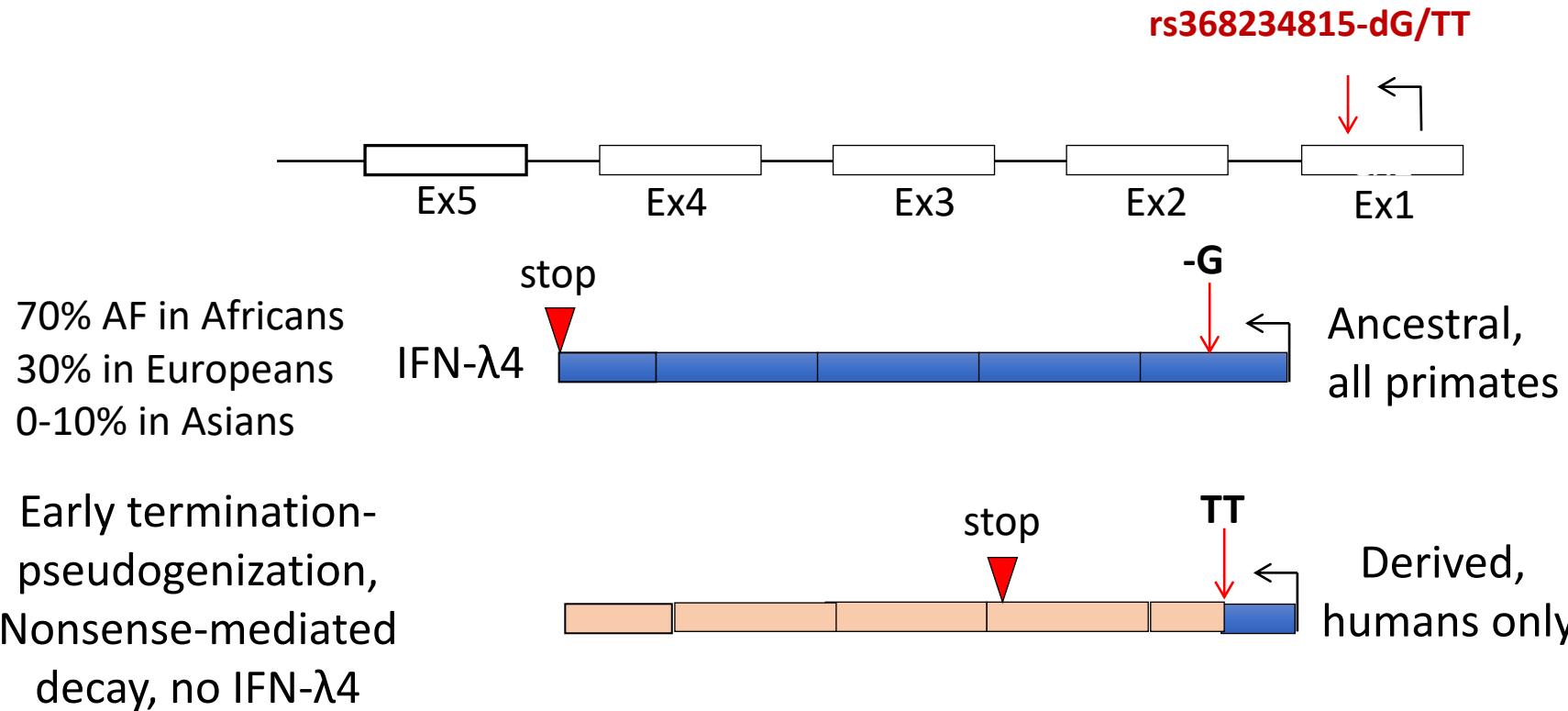
**2009:** GWAS for spontaneous and treatment-induced clearance of hepatitis C virus (HCV). One of the strongest GWAS signals ever reported, ORs ~3-10  
70 Kb region encoding type-III interferons IL28A, IL28B, IL29, ~95% similarity

**A number of efforts to functionally characterize this GWAS signal**

**2013:** We discovered a novel gene, Interferon Lambda 4 (IFNL4)  
GWAS SNP - within first intron of IFNL4, only 367 bp from the novel variant  
rs368234815-dG/TT in first exon, explains and improves GWAS signal, function - frame-shift variant

Ge et al, Nature, 2009; Thomas et al, Nature, 2009  
Kotenko et al, Sheppard et al, Nat Immun, 2002  
Prokunina-Olsson et al, Nat Gen, 2013

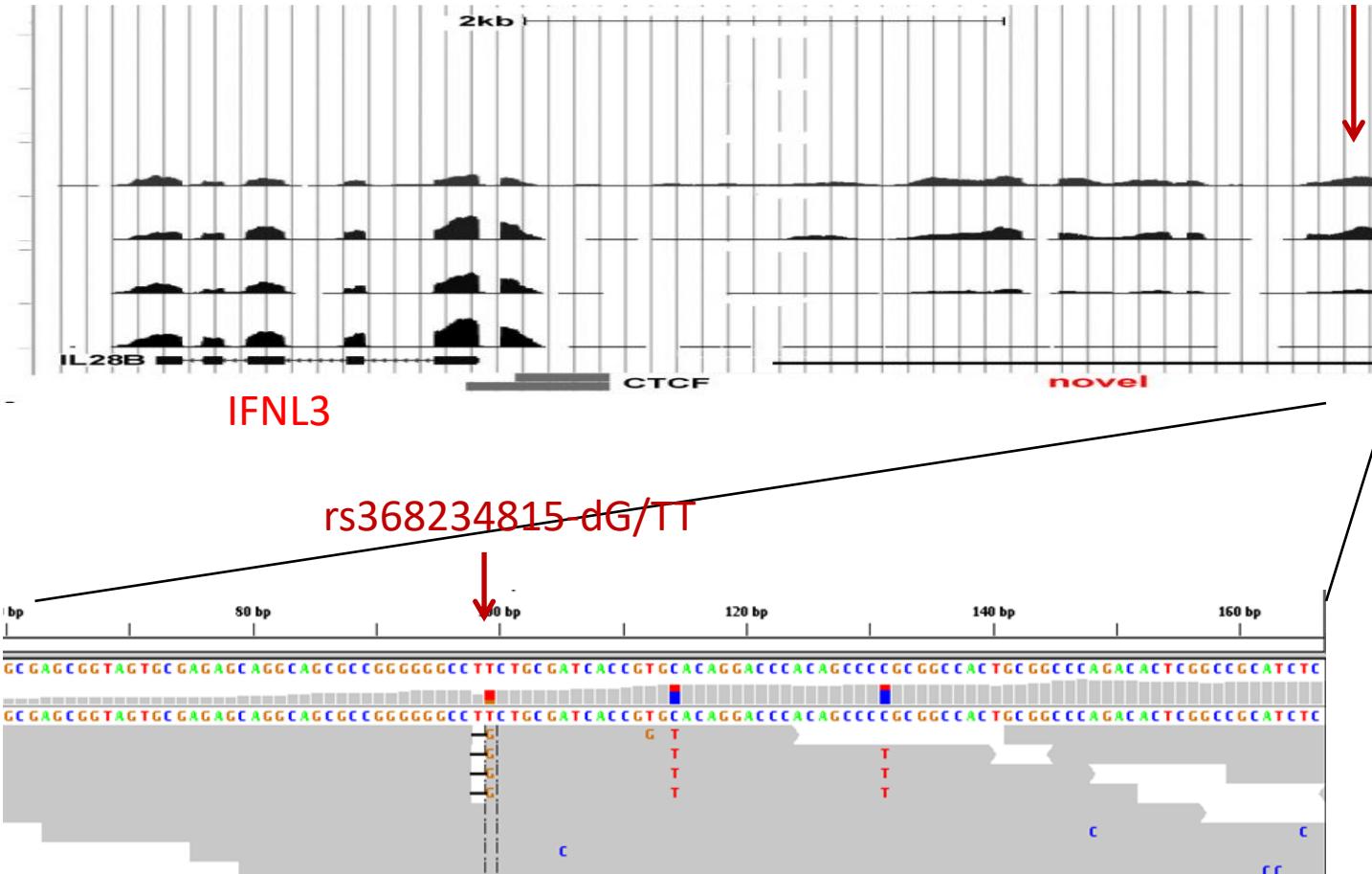
# Molecular phenotype of the GWAS signal: creating or eliminating IFNL4 protein



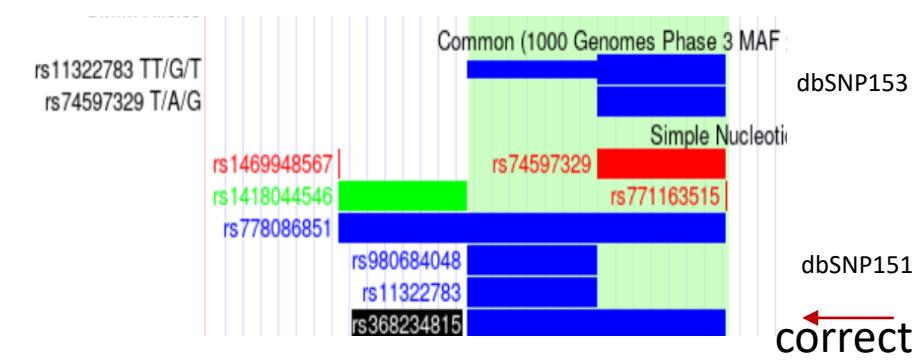
From UCSC: Homo sapiens interferon, lambda 4 (gene/pseudogene) (IFNL4)

**RefSeq Summary (NR\_074079):** This gene is a polymorphic pseudogene which, in some humans, encodes the interferon (IFN) lambda 4 protein. Humans are polymorphic for the dinucleotide TT/deltaG allele. Compared to the ancestral state in non-human primates, the TT allele produces a frameshift in the coding region of this gene which is predicted to induce nonsense-mediated mRNA decay.

# Variant and gene discovery through functional studies



RNA-sequencing in time series,  
In primary human  
hepatocytes from a single donor.  
In vitro treated with polyIC to  
mimic infection with RNA-viruses



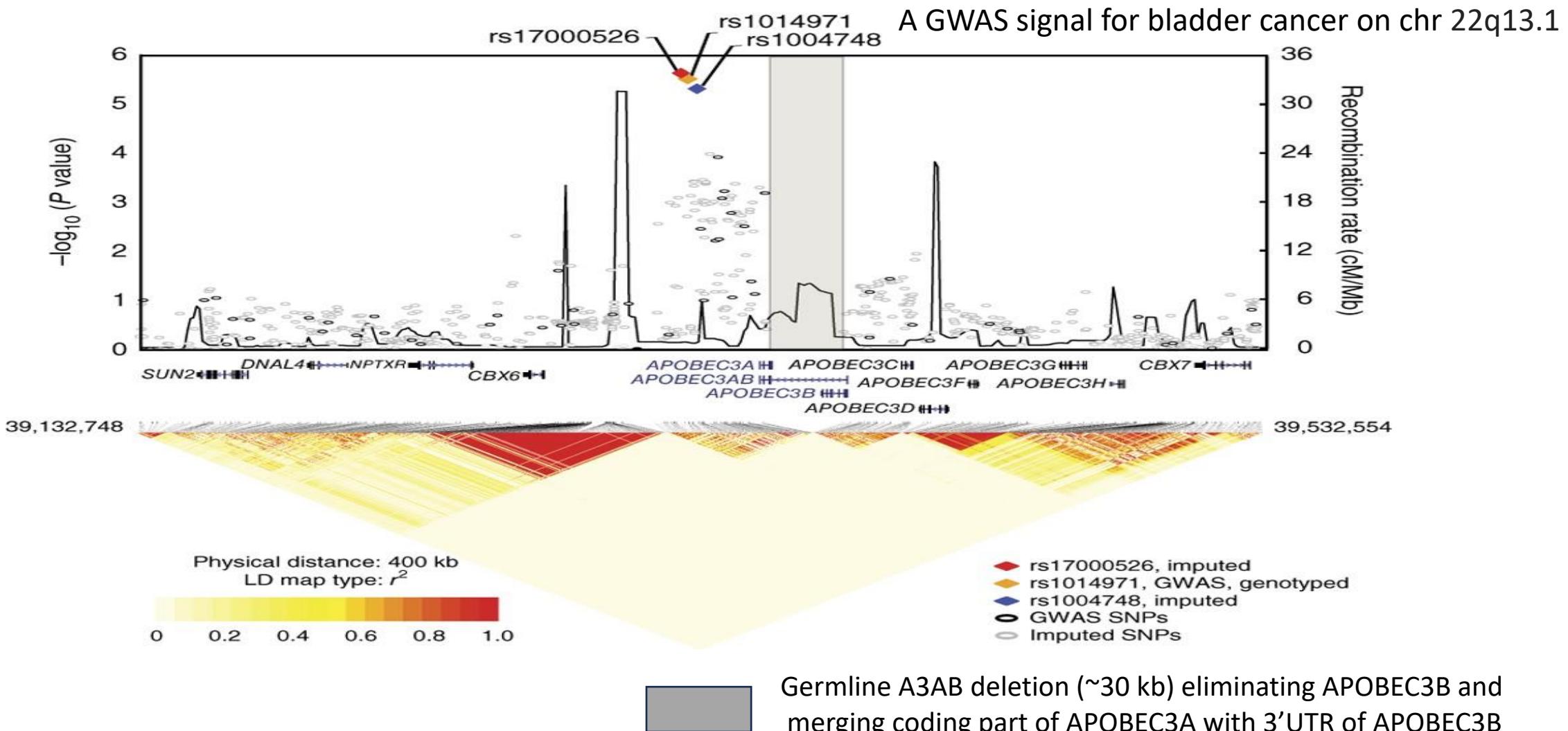
Currently: in 1KG **rs368234815-TT/-G** is represented by  
[rs74597329: T/G \(and A\)](#) – no Forgedb annotation, HaploReg 4  
[rs11322783 – TT/T \(and G\)](#) - no Forgedb annotation

## Summary on the Example 1

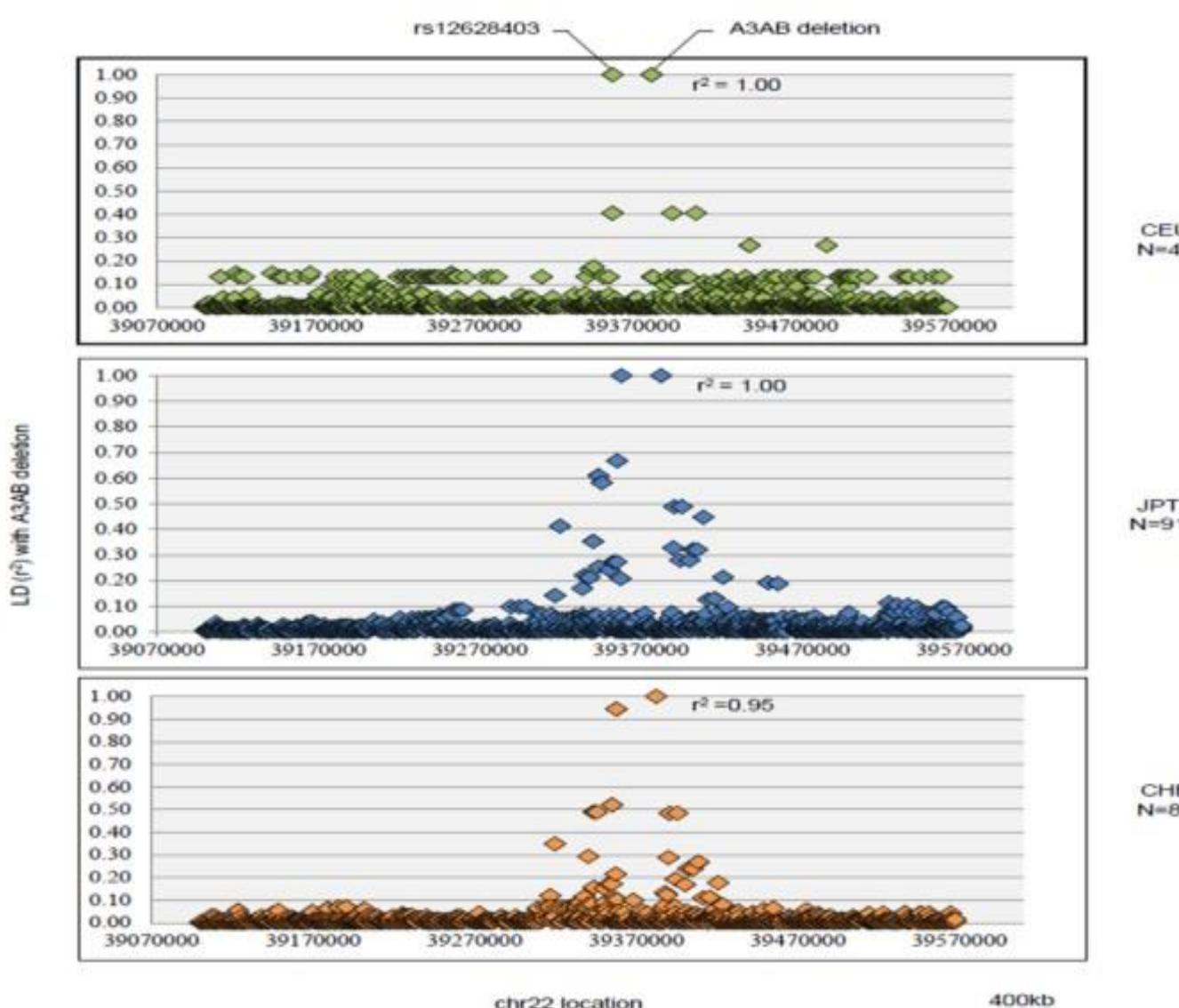
- Even with extremely strong GWAS signals (ORs=3-10) functional interpretations are not obvious
- The genome data and annotations maybe incomplete/misleading or confusing
- Functional analysis = discovery of the causal variant, novel gene and mechanism
- Without knowing this variant, the functional studies would be very different and not informative
- Even after 10 years after discovery of IFNL4 and the causal variant, genome annotation (dbSNP153 is still confusing and unhelpful) and negatively affecting downstream analyses
- The initial GWAS SNP, rs12979860 (now intron of IFNL4) is not “the causal” variant but a good proxy for it (TT/dG). Rs12979860 = “IL28B marker” is a patented marker for predicting outcomes of HCV clearance on several therapies
- Prediction markers are not necessarily the causal ones – matter of timing (priority), convenience of genotyping, commerce. Hard to change once established.

# Not all variants are SNPs and we still don't know all of them.

## Example 2: Polymorphic indels



# Intronic SNP rs12628403 is the only proxy for the A3AB deletion



LD between rs12628403 and A3AB deletion:

In Europeans,  $r^2 \sim 0.92-1.00$ , 7%

East Asians,  $r^2 \sim 0.95-1.00$ , 34%

Africans, no proxy, CNV has a frequency of 4.2% while rs12628403 is monomorphic

If both unavailable - cannot be imputed and have to be directly genotyped

Genotype landscape of this region is affected by the A3AB deletion

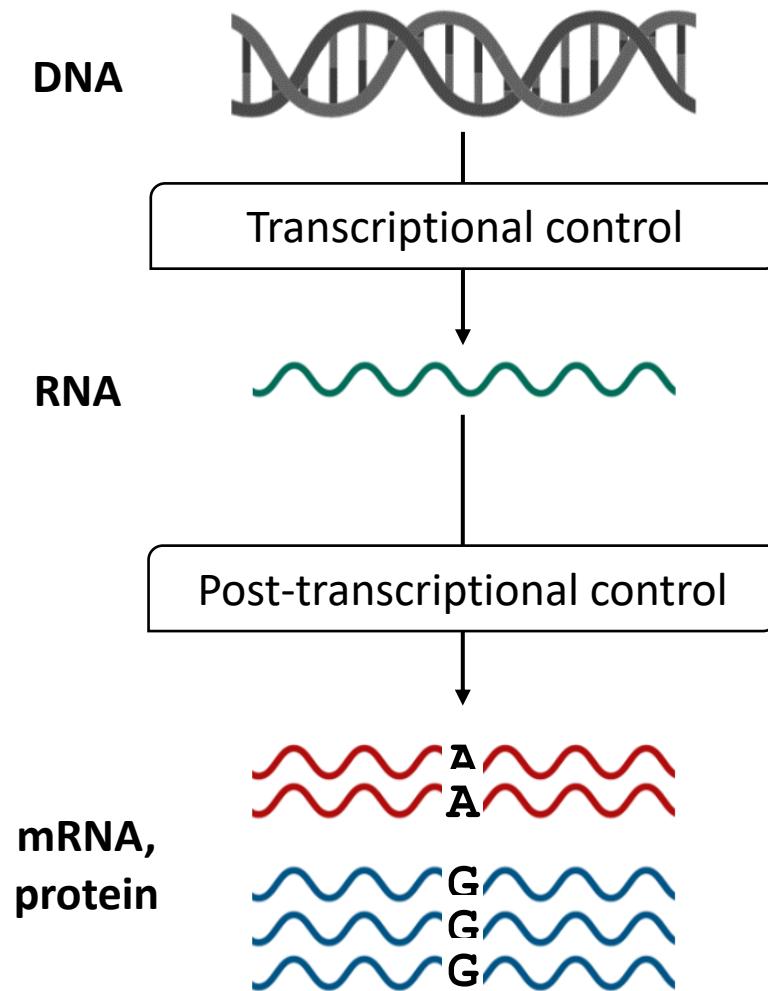
Very different functional consequences that need to be considered

Molecular phenotype: quantitative trait for genome-wide somatic mutations caused by APOBEC3B enzymes

## Example 3. It make take more than one variant ...

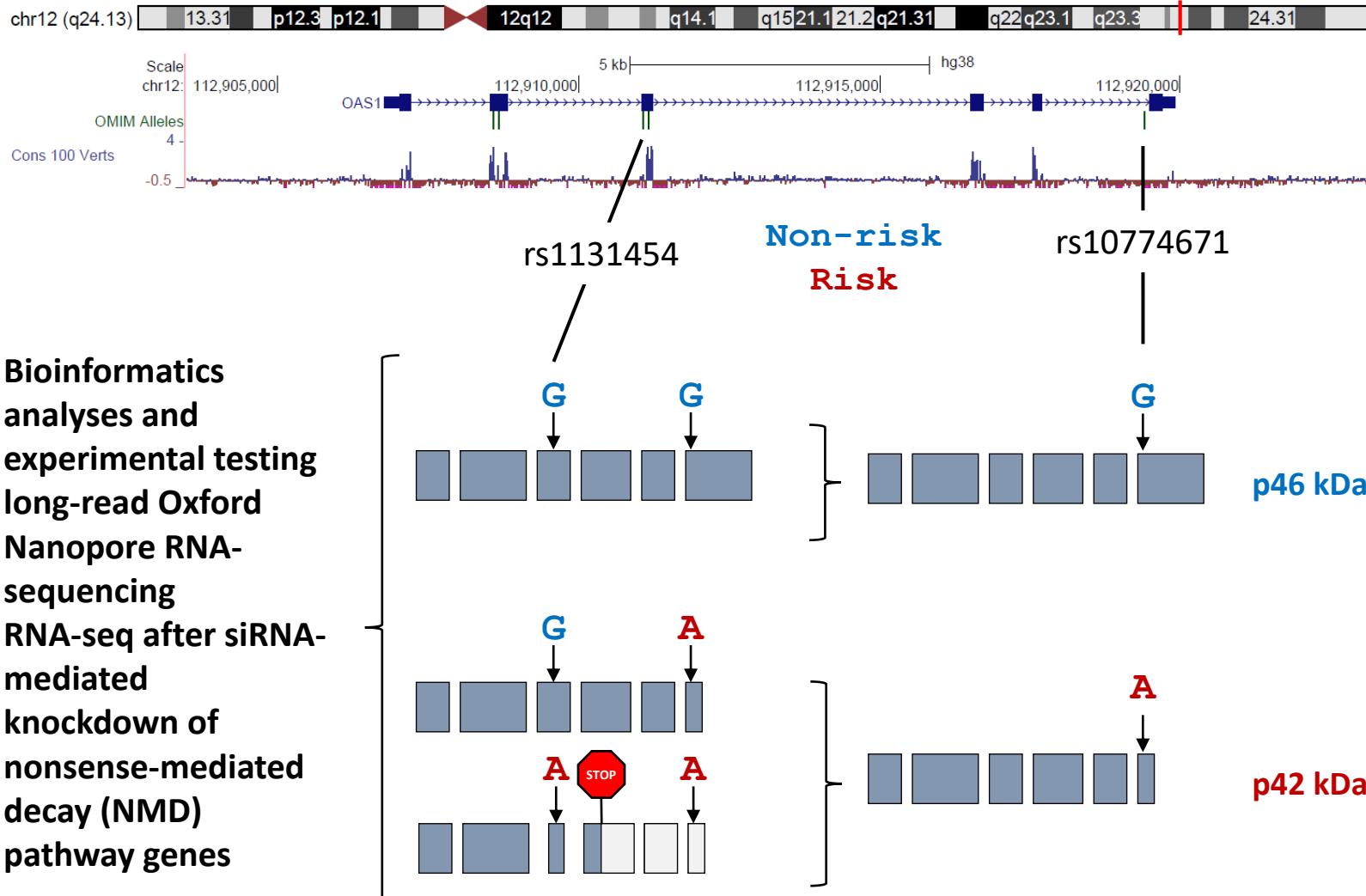
- We tend to assign all the blame on a single variant that statistically explains the association signal, other variants are ignored
- Functionally – the effects of several variants might be required to achieve a nuanced regulation through different but complementary mechanisms
- Risk of severe COVID-19 and OAS1 region

# Multiple approaches to explore the molecular contribution of rs10774671-A and rs1131454-A to the risk of severe COVID-19

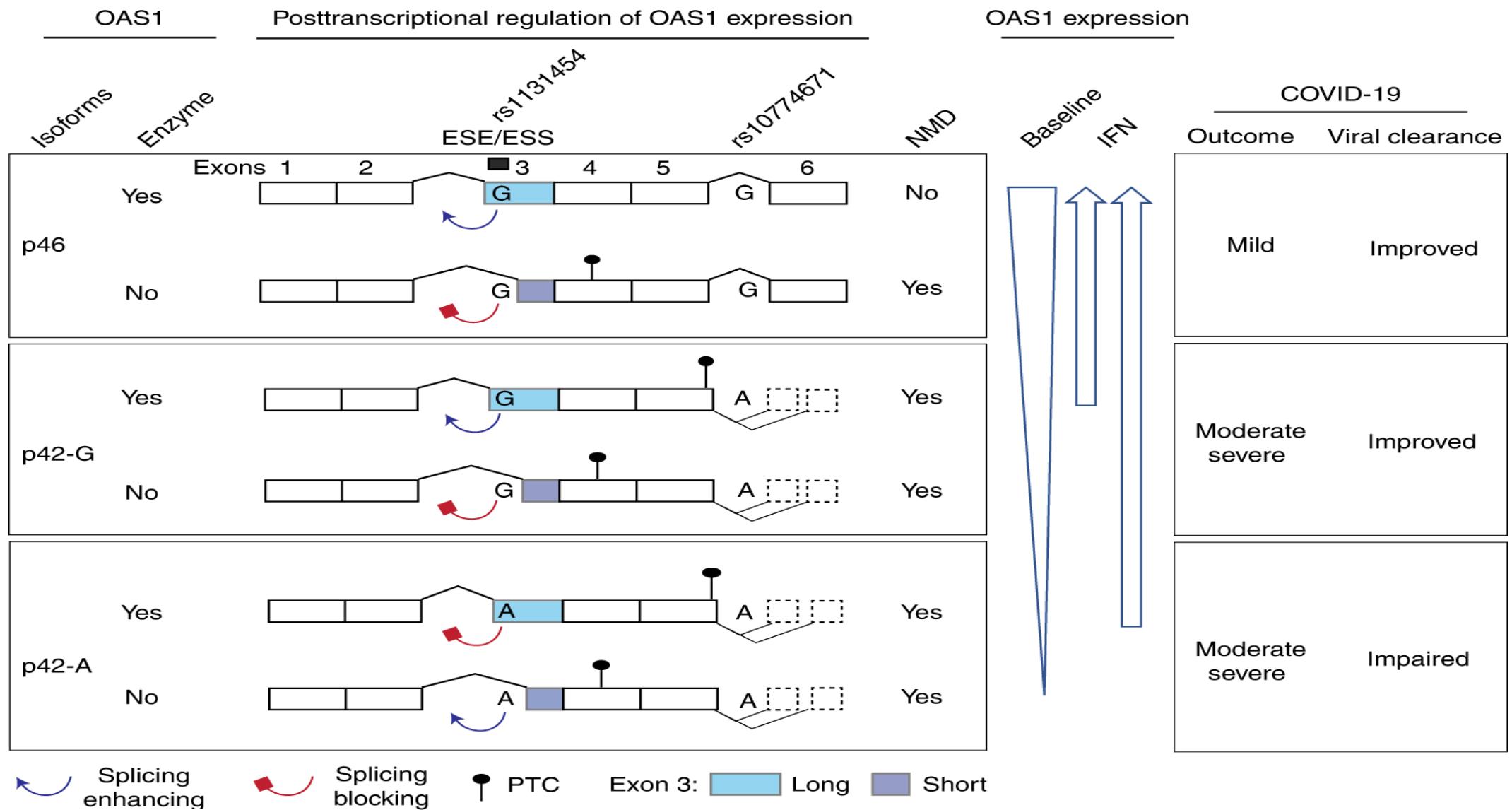


- Bioinformatic analyses of multiple datasets:
  - ATAC-seq: chromatin accessibility analysis
  - ChIP-seq: DNA-protein interactions
  - Hi-C: chromatin interactions
  - Conclusion: these mechanisms are not important here
- 
- *In silico* and *in vitro*:
  - rs10774671 – affects splicing of OAS1
  - rs1131454 creates a putative exonic splicing enhancer/silencer of exon 3
  - Bioinformatic analyses, ExonTrap and RT-PCR:
    - ✓ rs1131454-G allele creates a long exon 3 isoform
    - ✓ rs1131454-A allele creates short exon 3 isoform
  - Conclusion: these mechanisms are important for nonsense-mediated decay of mRNA carrying the risk alleles

# Through different mechanisms risk alleles of both SNPs contribute to nonsense-mediated decay (NMD) resulting in mRNA degradation, decreased protein expression and impaired sensing and clearance of SARS-CoV-19



# Two SNPs on the same haplotype contribute to COVID-19 outcomes



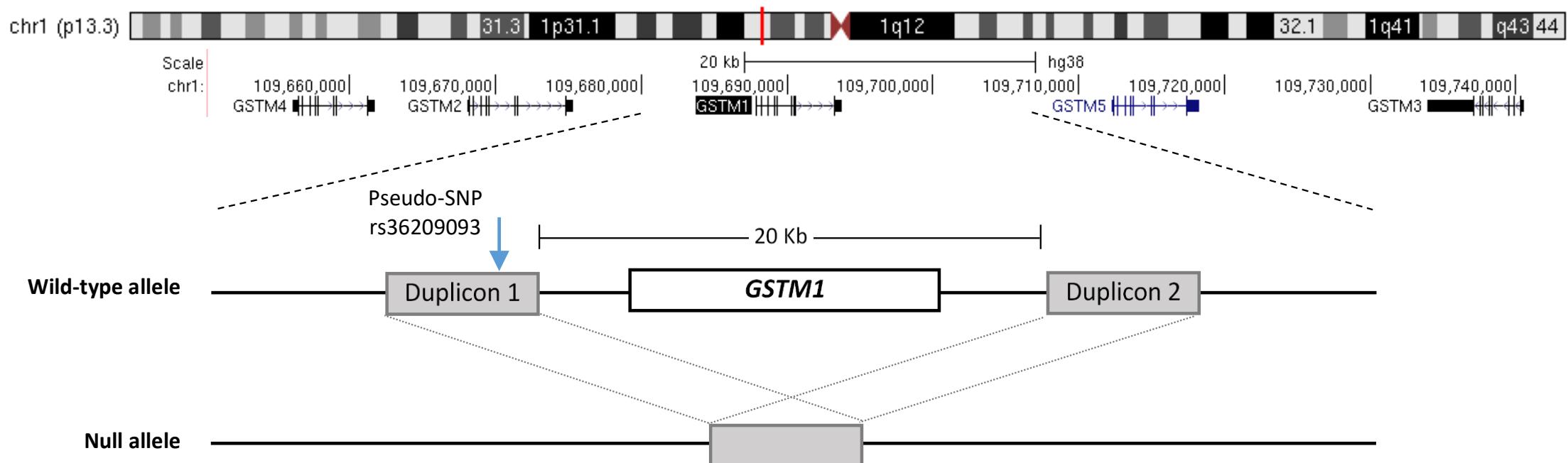
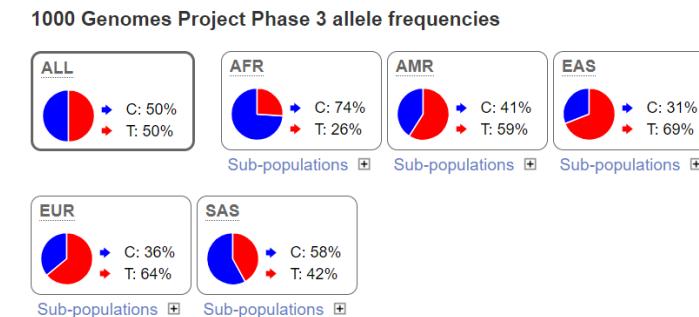
# Example 4. Fake SNPs and synthetic associations

rs36209093

Bladder cancer –  $p = 3.21E-18$ , in our meta-analysis, Koutros et al, Eur Urology, 2023

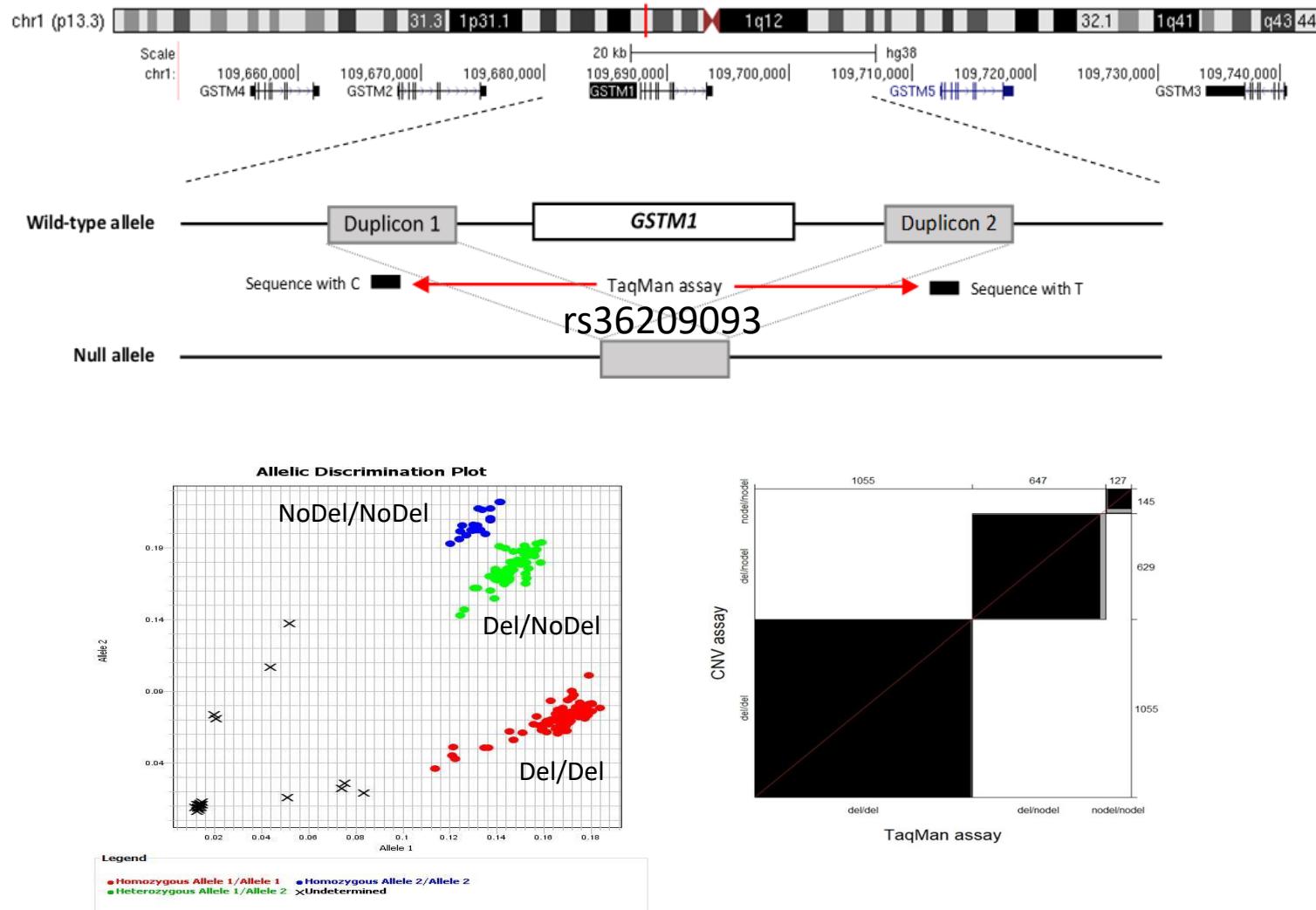
Other GWAS-significant associations in GWAS catalog:

Metabolite measurements; HDL, urinary metabolites, free and total cholesterol, cholesterol ratio, GSTM1 levels



Koutros et al, Eur Urology, 2023  
Florez-Vargas et al, in preparation

# Example 4. rs36209093 is a fake SNP but is still useful as a proxy



Custom TaqMan assay, 98.1% concordance rate in 1102 samples previously genotyped with a CNV assay.

Koutros et al, Eur Uro, 2023 and Florez-Vargas in preparation

Our custom TaqMan genotyping assay for rs36209093

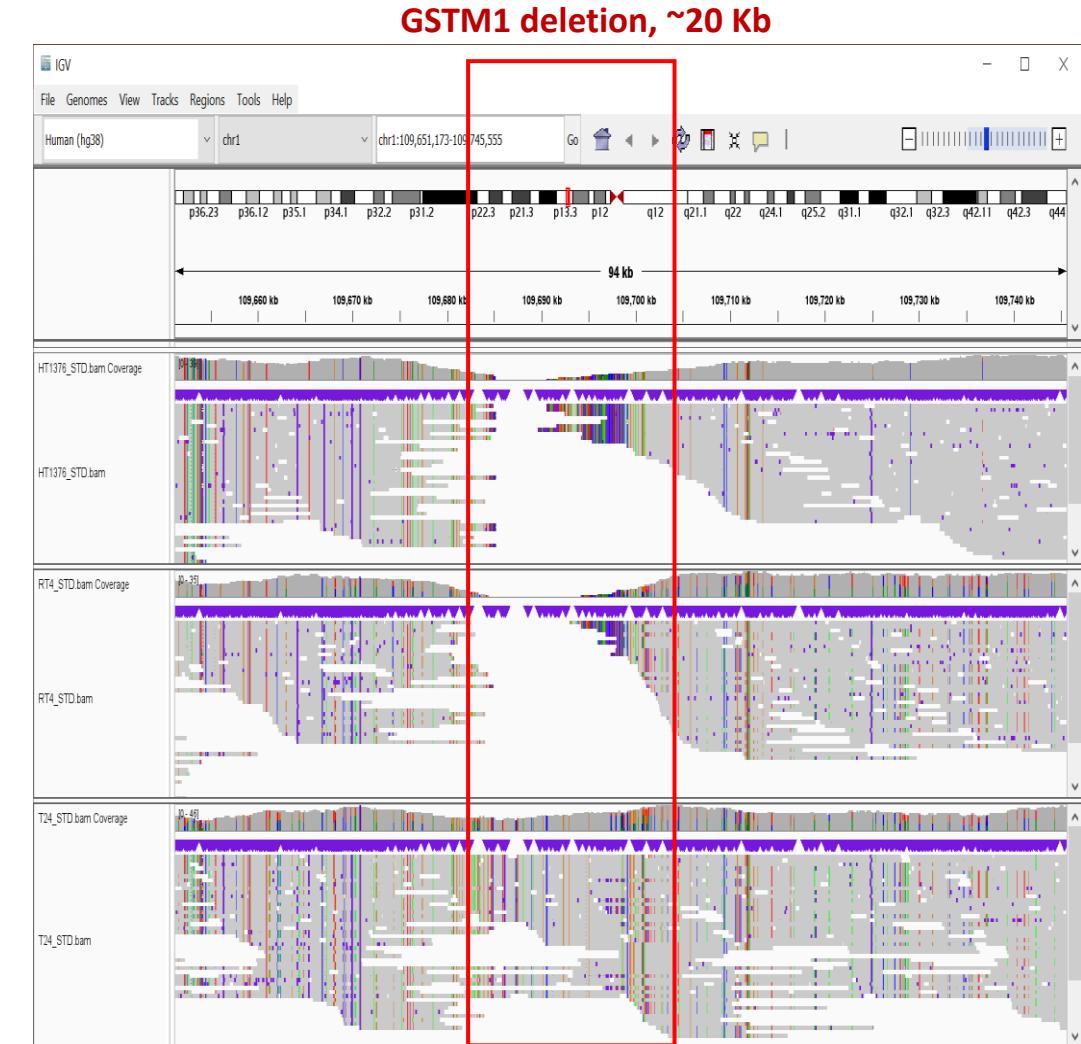
Dupl 1	C	2C:2T
Dupl 1	C	No del/
Dupl 2	T	No del
Dupl 2	T	
Dupl 1	C	No del/
Dupl 1	C	Del
Dupl 2	T	2C:1T
Dupl 1	C	
Dupl 1	C	Del/
Dupl 2	del	Del
Dupl 2	del	2C:0T

# Our experience with exploring the GSTM1 GWAS signal

Targeted PacBio long-read sequencing (~100Kb) in 1KG samples – poor coverage, low confidence about borders



Whole-genome long-read PacBio sequencing in several cell lines – looks better, still hard to define borders



# Human Pangenomes – a new reference resource

Article

## A draft human pangenome reference

<https://doi.org/10.1038/s41586-023-05896-x>

Received: 9 July 2022

Accepted: 28 February 2023

Published online: 10 May 2023

Open access

Check for updates

Wen-Wei Liao<sup>1,2,3,60</sup>, Mobin Asri<sup>4,60</sup>, Jana Ebler<sup>5,6,60</sup>, Daniel Doerr<sup>2,6</sup>, Marina Haukness<sup>4</sup>, Glenn Hickey<sup>4</sup>, Shuanglia Lu<sup>2</sup>, Julian K. Lucas<sup>4</sup>, Jean Monlong<sup>4</sup>, Haley J. Abel<sup>1</sup>, Silvia Buonaiuto<sup>6</sup>, Xian H. Chang<sup>4</sup>, Haoyu Cheng<sup>3,10</sup>, Justin Chu<sup>9</sup>, Vincenza Colonna<sup>8,11</sup>, Jordan M. Eizenga<sup>9</sup>, Xiaowen Feng<sup>10,9</sup>, Christian Fischer<sup>11</sup>, Robert S. Fulton<sup>12,13</sup>, Shilpa Garg<sup>14</sup>, Cristian Groza<sup>10</sup>, Andrea Guaracino<sup>10,8</sup>, William T. Harvey<sup>7</sup>, Simon Heumos<sup>10,13</sup>, Kerstin Howe<sup>10</sup>, Miten Jain<sup>2</sup>, Tsung-Yu Lu<sup>2</sup>, Charles Markello<sup>4</sup>, Fergal J. Martin<sup>23</sup>, Matthew W. Mitchell<sup>24</sup>, Katherine M. Munson<sup>17</sup>, Moses Njagi Mwaniki<sup>25</sup>, Adam M. Novak<sup>4</sup>, Hugh E. Olsen<sup>4</sup>, Trevor Pesout<sup>4</sup>, David Porubsky<sup>9</sup>, Piotr Princ<sup>11</sup>, Jonas A. Stibbeisen<sup>26</sup>, Jouni Sirén<sup>4</sup>, Chad Tomlinson<sup>2</sup>, Flavia Villani<sup>11</sup>, Mitchell R. Vollger<sup>17,27</sup>, Lucinda L. Antonacci-Fulton<sup>12</sup>, Gunjan Bald<sup>28</sup>, Carl A. Baker<sup>17</sup>, Anastasiya Belyaeva<sup>28</sup>, Konstantinos Billis<sup>23</sup>, Andrew Carroll<sup>28</sup>, Pi-Chuan Chang<sup>28</sup>, Sarah Cody<sup>12</sup>, Daniel E. Cook<sup>28</sup>, Robert M. Cook-Deegan<sup>29</sup>, Omar E. Cornejo<sup>29</sup>, Mark Diekhans<sup>30</sup>, Peter Ebert<sup>2,6,31</sup>, Susan Fairley<sup>32</sup>, Olivier Fedrigoto<sup>32</sup>, Carlos Garcia Giron<sup>23</sup>, Richard E. Green<sup>38,39</sup>, Yan Gao<sup>34</sup>, Nanibar A. Garrison<sup>35,36,37</sup>, Carlos Garcia Giron<sup>23</sup>, Richard E. Green<sup>38,39</sup>, Leanne Haggerty<sup>23</sup>, Kendra Hoekzema<sup>27</sup>, Thibaut Hourlier<sup>23</sup>, Hanlee P. Ji<sup>40</sup>, Einear E. Kenny<sup>41</sup>, Barbara A. Koenig<sup>42</sup>, Alexey Kolesnikov<sup>28</sup>, Jan O. Korbel<sup>23,43</sup>, Jennifer Kordosky<sup>17</sup>, Sergey Koren<sup>44</sup>, Ho-Joon Lee<sup>40</sup>, Alexandra P. Lewis<sup>17</sup>, Hugo Magalhães<sup>3,6</sup>, Santiago Marco-Sola<sup>45,46</sup>, Pierre Marjion<sup>33</sup>, Ann McCartney<sup>44</sup>, Jennifer McDaniel<sup>47</sup>, Jacqueline Mountcastle<sup>23</sup>, Maria Nattestad<sup>29</sup>, Sergey Nurk<sup>44</sup>, Nathan D. Olson<sup>9</sup>, Alice B. Popejoy<sup>48</sup>, Daniela Putz<sup>49</sup>, Mikko Rautialainen<sup>44</sup>, Allison A. Regier<sup>49</sup>, Arang Rhie<sup>44</sup>, Samuel Sacco<sup>30</sup>, Ashley D. Sanders<sup>50</sup>, Valerie A. Schneider<sup>31</sup>, Baerent I. Schultz<sup>33</sup>, Kishwar Shafin<sup>48</sup>, Michael W. Smith<sup>33</sup>, Heidi J. Sofia<sup>33</sup>, Ahmad N. Abou Tayoun<sup>32,33</sup>, Françoise Thibaud-Nissen<sup>31</sup>, Francesca Floriana Tricomi<sup>23</sup>, Justin Wagner<sup>47</sup>, Brian Walenz<sup>44</sup>, Jonathan D. Wood<sup>20</sup>, Aleksey V. Zimin<sup>49,54</sup>, Guillaume Bourque<sup>35,36,37</sup>, Mark J. P. Chaisson<sup>22</sup>, Paul Flicek<sup>23</sup>, Adam M. Phillip<sup>44</sup>, Justin M. Zook<sup>47</sup>, Evan E. Eichler<sup>75,58</sup>, David Haussler<sup>4,56</sup>, Ting Wang<sup>72,53</sup>, Erich D. Jarvis<sup>72,53,59</sup>, Karen H. Miga<sup>1</sup>, Erik Garrison<sup>10,53</sup>, Tobias Marschall<sup>1,6,32</sup>, Ira M. Hall<sup>1,2,50</sup>, Heng Li<sup>1,3,50</sup> & Benedict Paten<sup>42</sup>

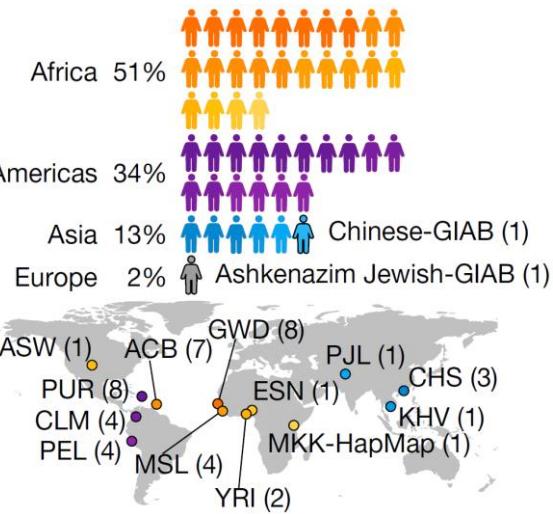
Here the Human Pangenome Reference Consortium presents a first draft of the human pangenome reference. The pangenome contains 47 phased, diploid assemblies from a cohort of genetically diverse individuals<sup>1</sup>. These assemblies cover more than 99% of the expected sequence in each genome and are more than 99% accurate at the structural and base pair levels. Based on alignments of the assemblies, we generate a draft pangenome that captures known variants and haplotypes and reveals new alleles at structurally complex loci. We also add 119 million base pairs of euchromatic polymorphic sequences and 1,115 gene duplications relative to the existing reference GRCh38. Roughly 90 million of the additional base pairs are derived from structural variation. Using our draft pangenome to analyse short-read data reduced small variant discovery errors by 34% and increased the number of structural variants detected per haplotype by 104% compared with GRCh38-based workflows, which enabled the typing of the vast majority of structural variant alleles per sample.

The human reference genome has formed the backbone of human genomics since its initial draft release more than 20 years ago<sup>2</sup>. The primary sequences are a mosaic representation of individual haplotypes containing one representative scaffold sequence for each chromosome. There are 210 Mb of gap or unknown (151 Mb) or computationally simulated sequences (59 Mb) within the current GRCh38 release, constituting 6.7% of the primary chromosome scaffolds. Missing reference sequences create an observational bias, or streetlamp effect, which limits studies to be within the boundaries of the reference. Recently,

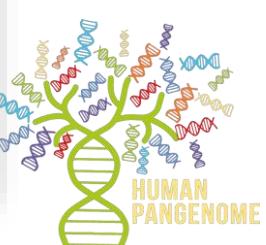
the Telomere-to-Telomere (T2T) consortium finished the first complete sequence of a haploid human genome, T2T-CHM13, which provides a contiguous representation of each autosome and of chromosome X, with the exception of some ribosomal DNA arrays that remain to be fully resolved<sup>3</sup>. Using T2T-CHM13 directly improves genomic analyses; for example, discovering 3.7 million additional single-nucleotide polymorphisms (SNPs) in regions non-syntenic to GRCh38 and better representing the true copy number variants (CNVs) of samples from the 1000 Genomes Project (1KG) compared with GRCh38 (refs. 1, 4).

A list of affiliations appears at the end of the paper.

312 | Nature | Vol 617 | 11 May 2023



- The Human Pangenome Reference Consortium (HPRC) presents a first draft of the human pangenome reference.
- **47 phased, diploid assemblies.**
- The data is based on Pacific Biosciences (PacBio) high-fidelity (HiFi) and Oxford Nanopore Technologies (ONT) long-read sequencing.
- HPRC add **119 million** base pairs of euchromatic polymorphic sequences and 1,115 gene duplications relative to the existing reference GRCh38. Roughly **90 million of the additional base pairs are derived from structural variation**.



<https://humanpangenome.org/>

N = 47 individuals  
98% are non-Europeans

with the goal of increasing to  
350 individuals by mid-2024

**Practical session with Oscar Florez-Vargas:  
Explore the GSTM1 region based on the HPRC assemblies**

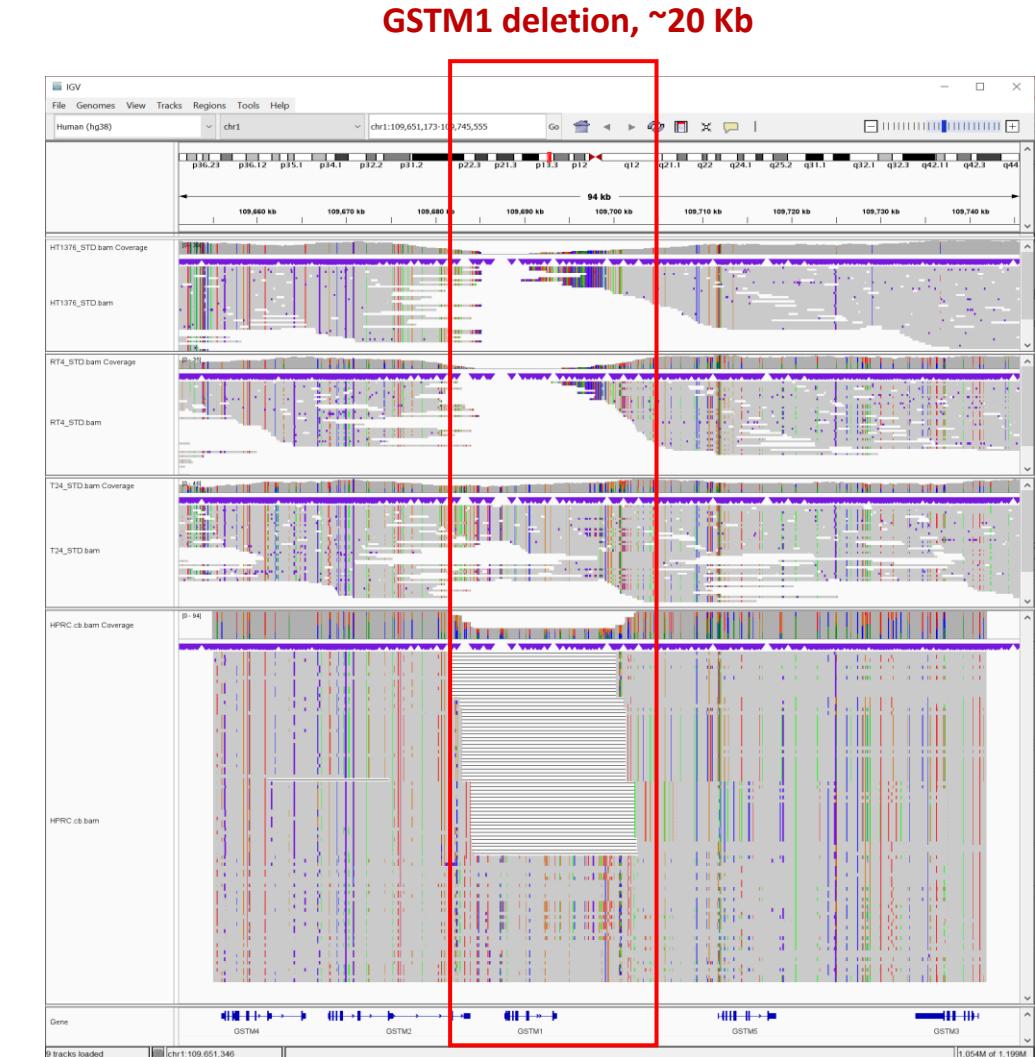
Liao et al., Nature 2023

# Better results based on long-read genome assemblies

Targeted PacBio long-read sequencing (~100Kb) in 1KG samples – poor coverage, low confidence about borders



Whole-genome long-read PacBio sequencing in several cell lines – looks better, still hard to define borders



Assemblies  
1 read = 1 chromosome