# Statistical Genomics
# Total Request Live

## Kasper Daniel Hansen

< khansen@jhsph.edu | www.hansenlab.org >

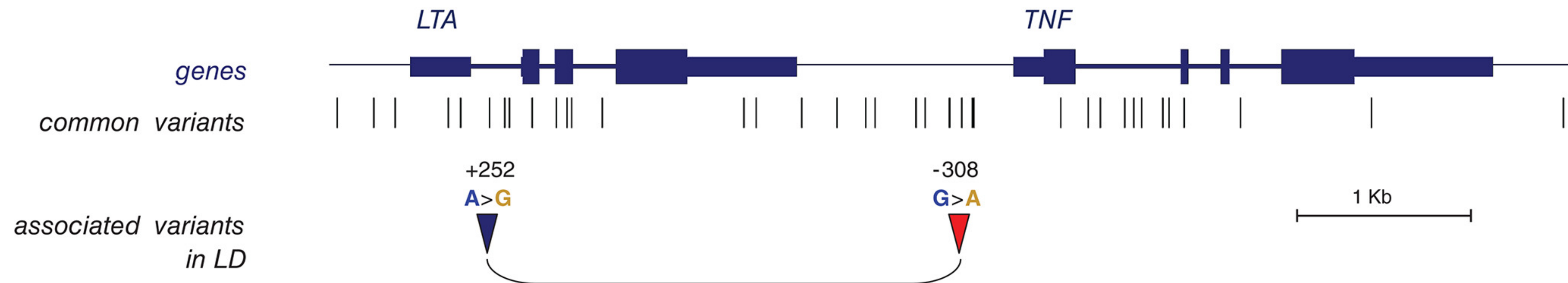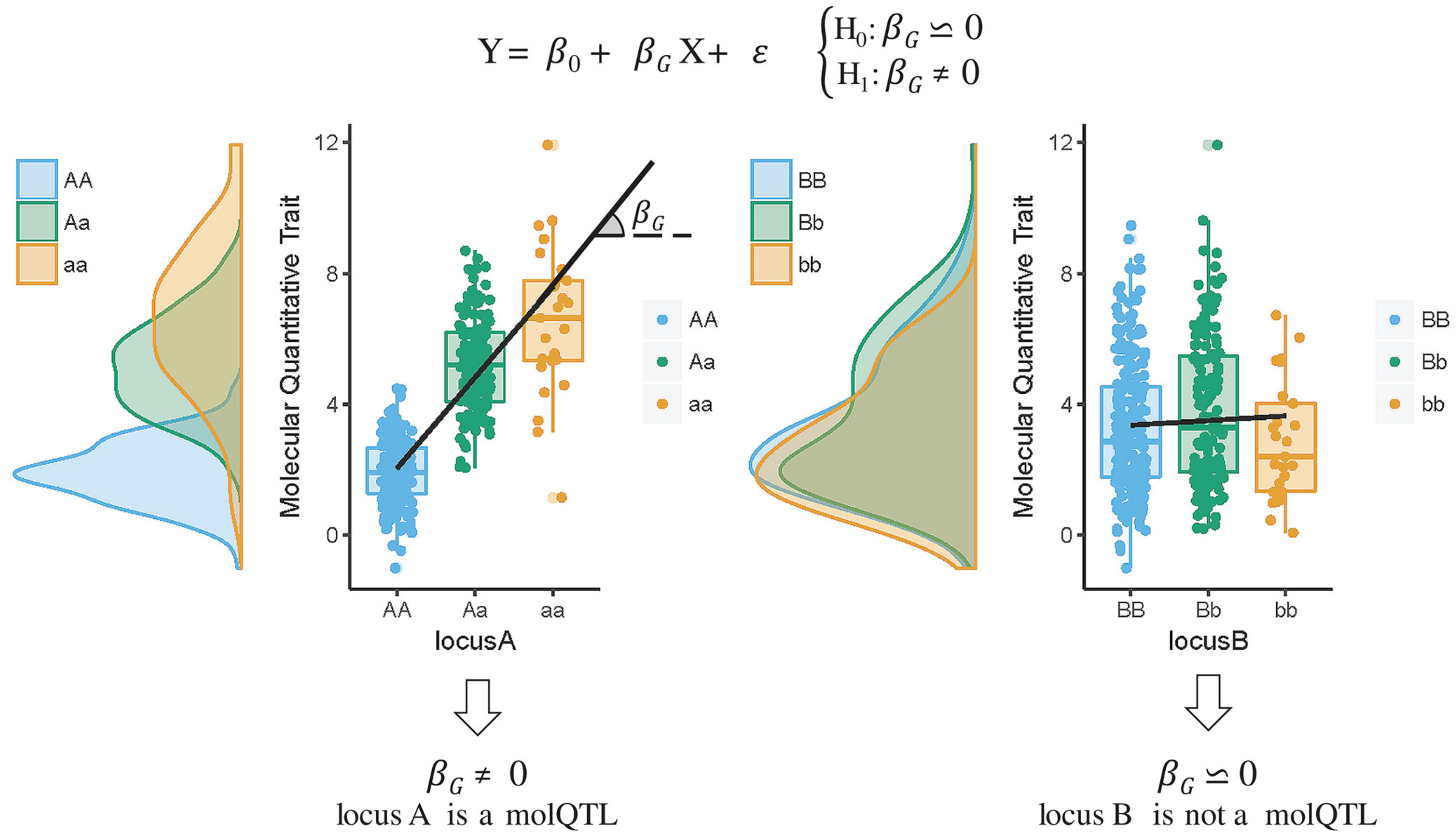Department of Biostatistics

Johns Hopkins University

# eQTL

In genetics, associating genetic variants with a continuous phenotype (like height) is call a quantitative trait loci (QTL) analysis.

We can do this, where our "phenotype" is the expression of all genes. That is called an expression-QTL analysis.

Some numbers: say we have 1M genotypes and 10,000 SNPs.
That yields 10e11 comparisons. So we only do "local" comparisons.

$$Y = \beta_0 + \beta_G X + \varepsilon \qquad \begin{cases} H_0: \beta_G \simeq 0 \\ H_1: \beta_G \neq 0 \end{cases}$$

$\beta_G \neq 0$
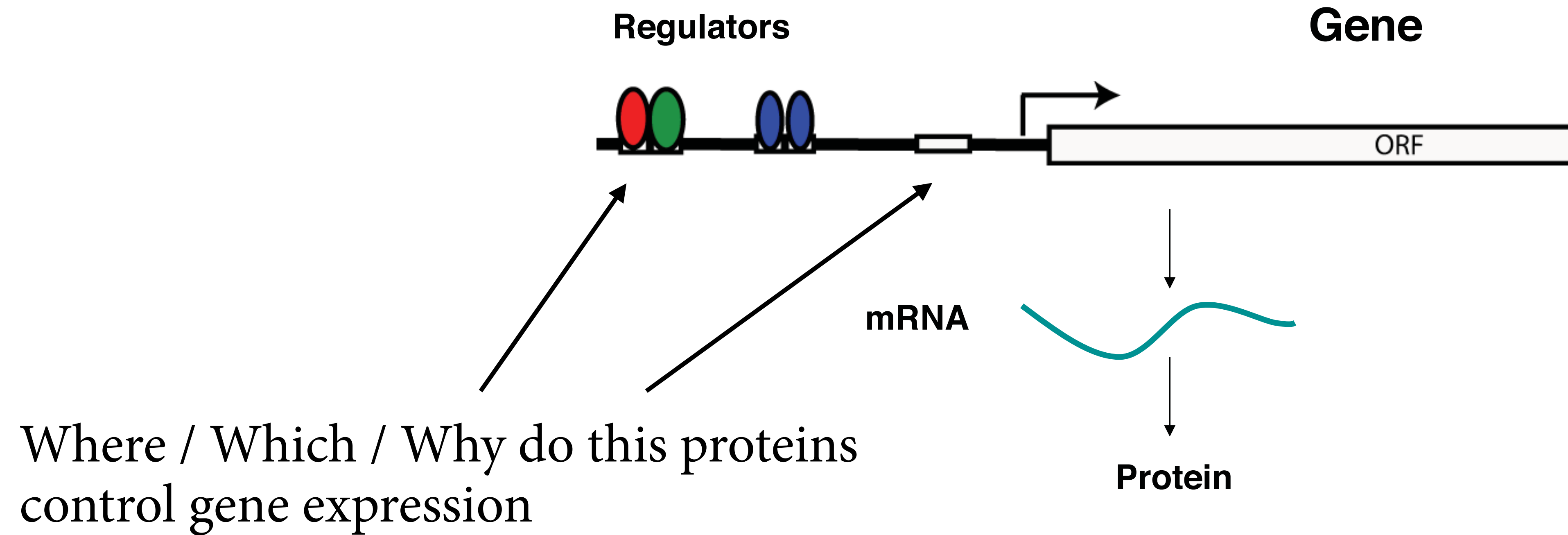locus A is a molQTL

$\beta_G \simeq 0$
locus B is not a molQTL

Lots of book keeping. Lots of tests were both the outcome (expression) AND the covariates (the SNPs) vary.

Standard approach: control for "everything" and see if anything "survives".

Otherwise, there is nothing special.
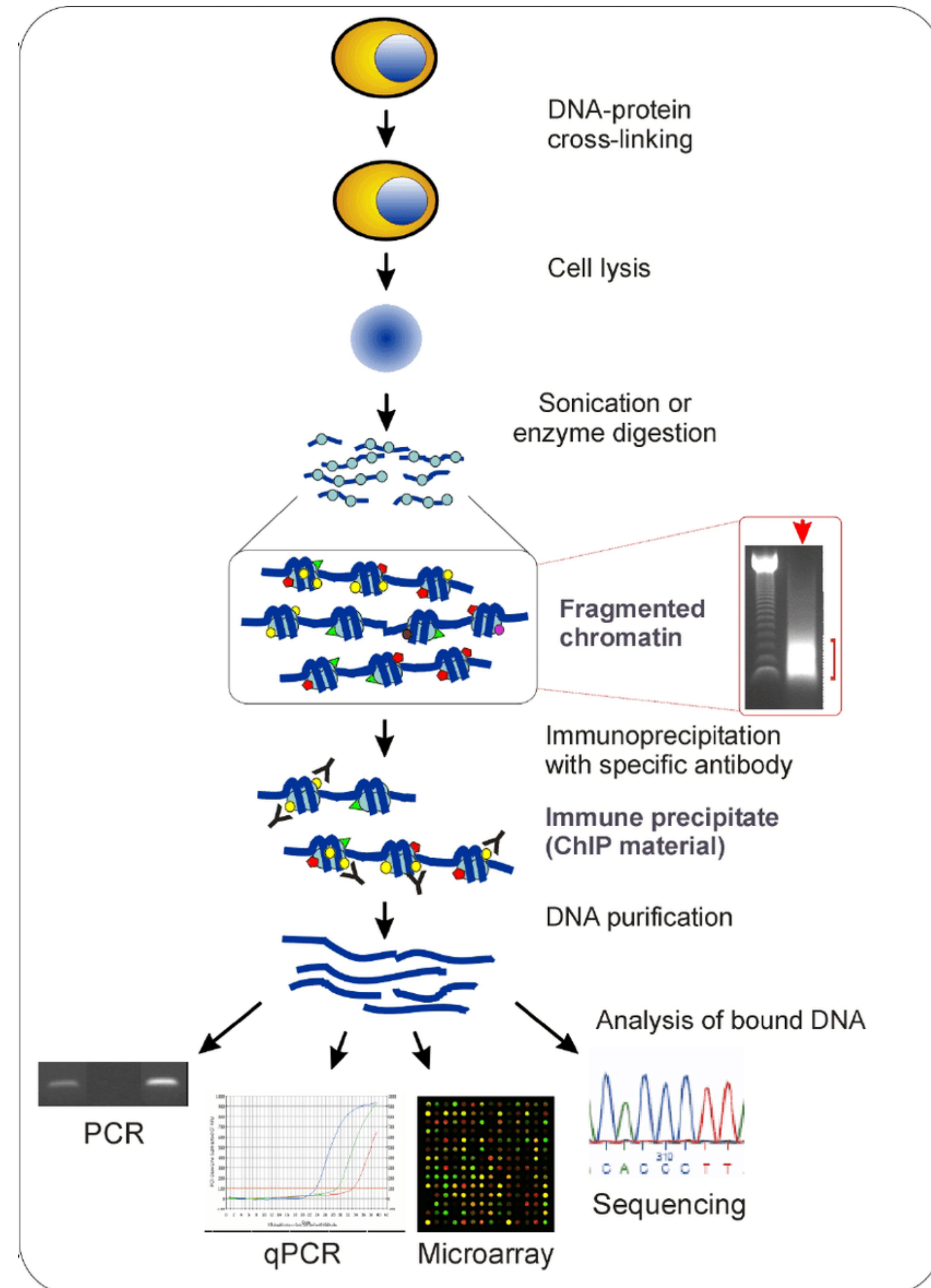
# ATAC (and ChIP) sequencing

**Regulators**

**Gene**

ORF

**mRNA**

**Protein**

Where / Which / Why do this proteins
control gene expression

3

ChIP: measures binding of a protein to DNA. Used to study
   Transcription factors
   Histone modifications (narrow and broad)

DNAse and ATAC: measures open chromatin
(where can proteins potentially bind)

Cross-linking

Fragmentation

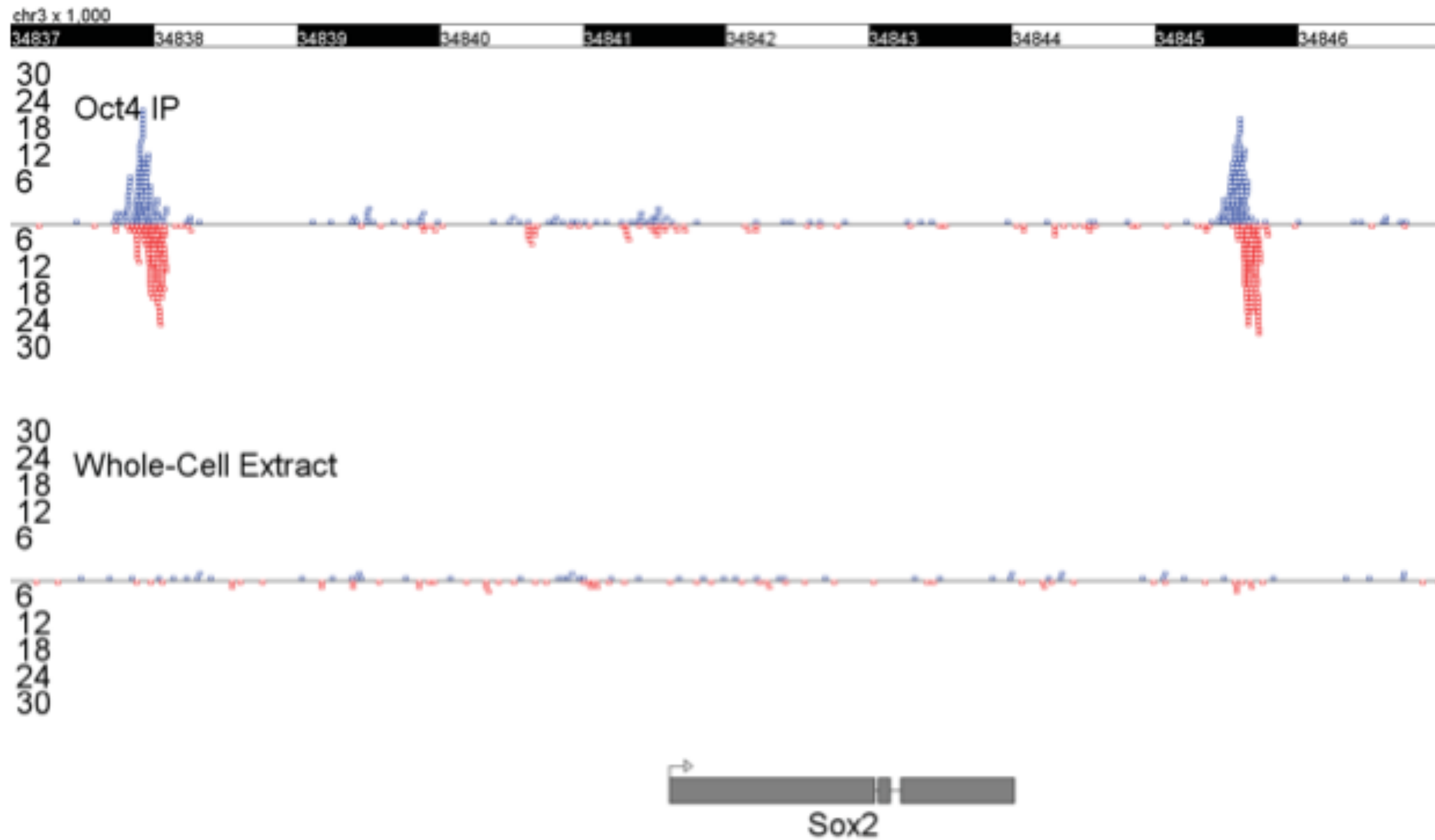IP enrichment (requires antibody)

Sequencing

Typical controls
  Input (don't do IP)
  Igf (non-specific antibody)

We see peaks in control samples

Usually we think of the control sample as measuring background peaks, and most peak callers try to separate signal peaks from background peaks.

Some regions of the genome are highly susceptible to background peaks. ENCODE maintains a "black list regions" which is usually excluded.

ChIP requires a good, specific, antibody. That can be hard to obtain

ChIP requires large amount of material. For this reason, we tend to see ChIP in model organisms or cell lines.

Because of this, ChIP analysis is usually done as "single-condition" (where is the protein bound in this sample) and with few (usually 2) technical replicates.

The first step here is "peak calling".

On technical replicates: these are usually "growth" replicates of cell lines. In the ChIP literature, this is called "biological replicates".

Tons of software.

Despite all the work, I don't consider this a "approximately solved" problem.
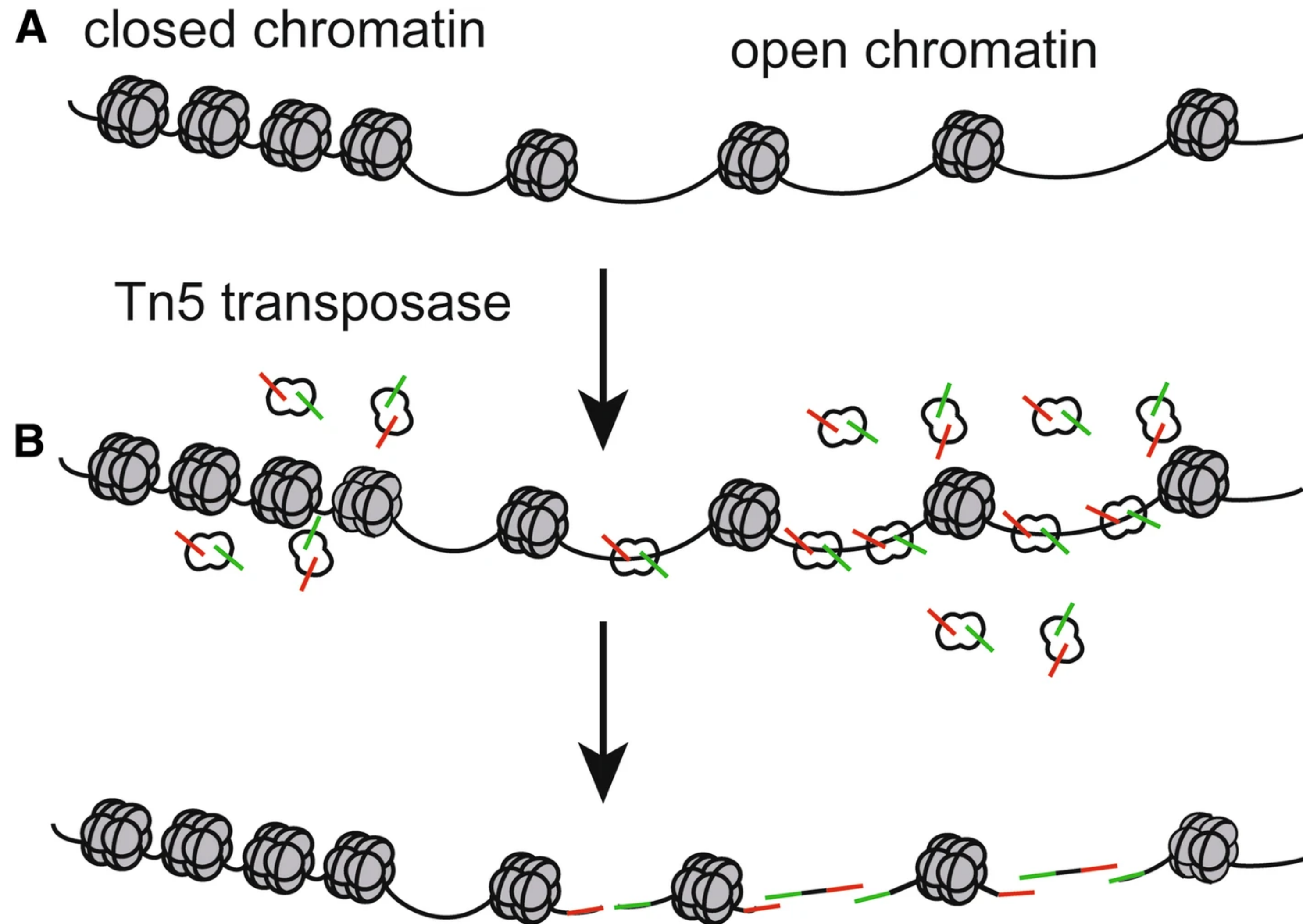
But it seems low reward to work on…

Irreducible discovery rate: Qunhua Li et al (2011) Annals of applied statistics
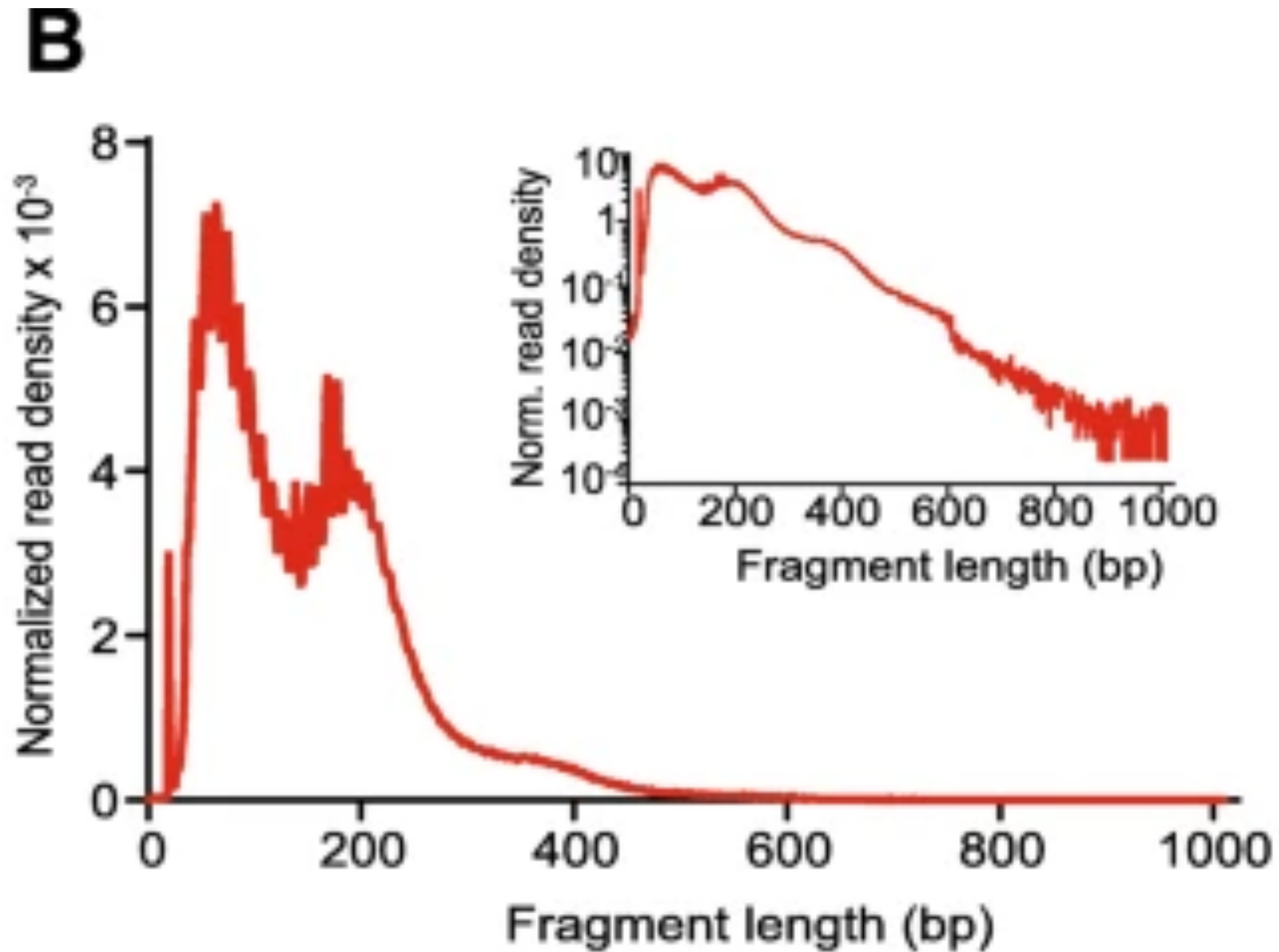
Question: how do we identify peaks which are "reproducible" between 2 replicates?

Both assays measure open chromatin.

DNAse: classic assay. Requires large amounts of material. REALLY HARD to do well.

ATAC: a modern alternative to DNAse. Easy, low input requirement. Massive game-changer

A  closed chromatin   open chromatin

Tn5 transposase

B

Nucleosome free fragmenets, mono-nucleosome fragments, etc

Two approaches:
  1) peak calling, followed by differential analysis of peaks
  2) side-step peak calling by binning the genome

Both creates a features by samples matrix for use with differential expression tools.

Peak calling for ATAC tends to use tools for ChIP without controls.