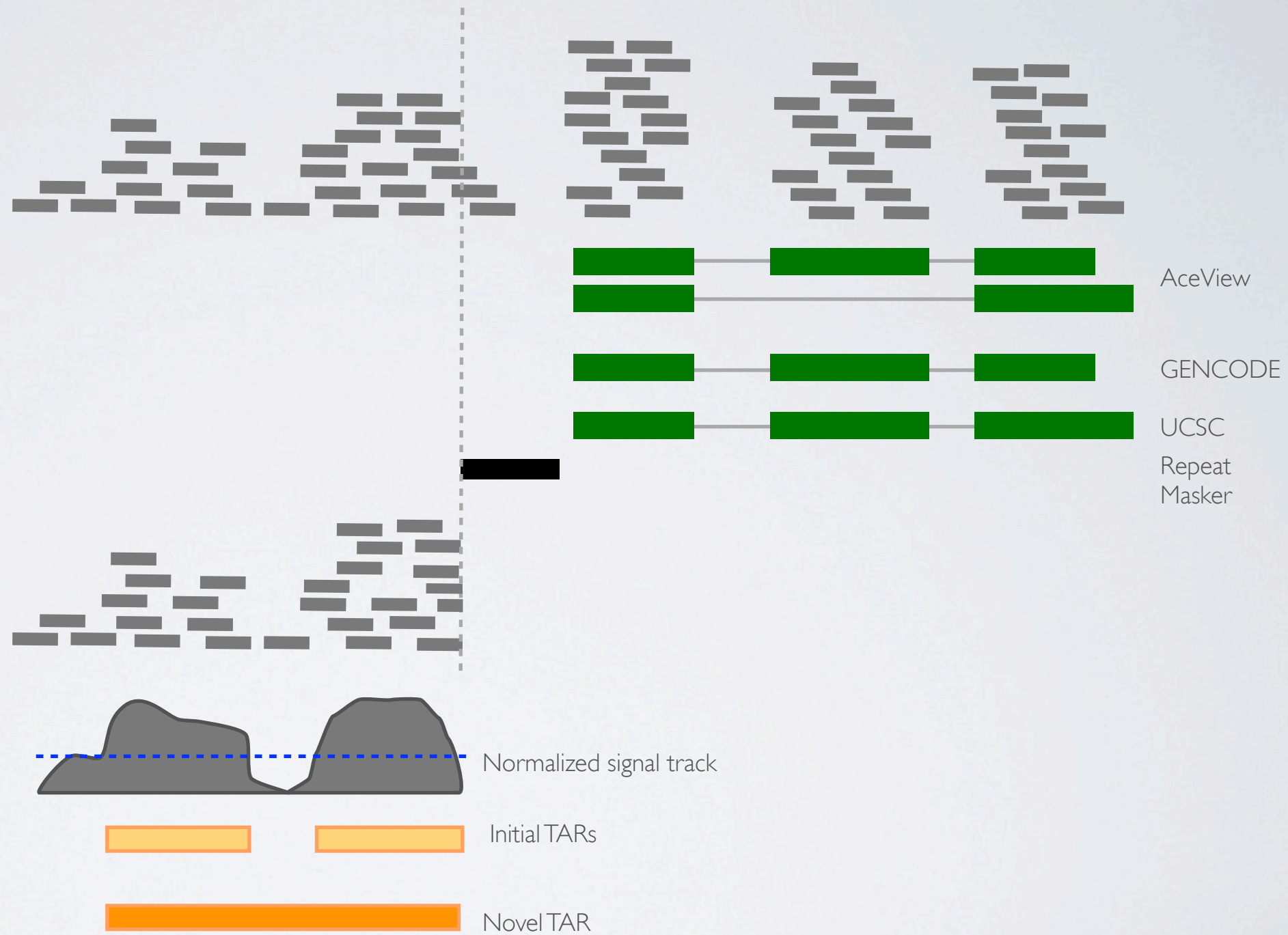


NOVEL TRANSCRIPTIONALLY ACTIVE REGIONS (TARs)

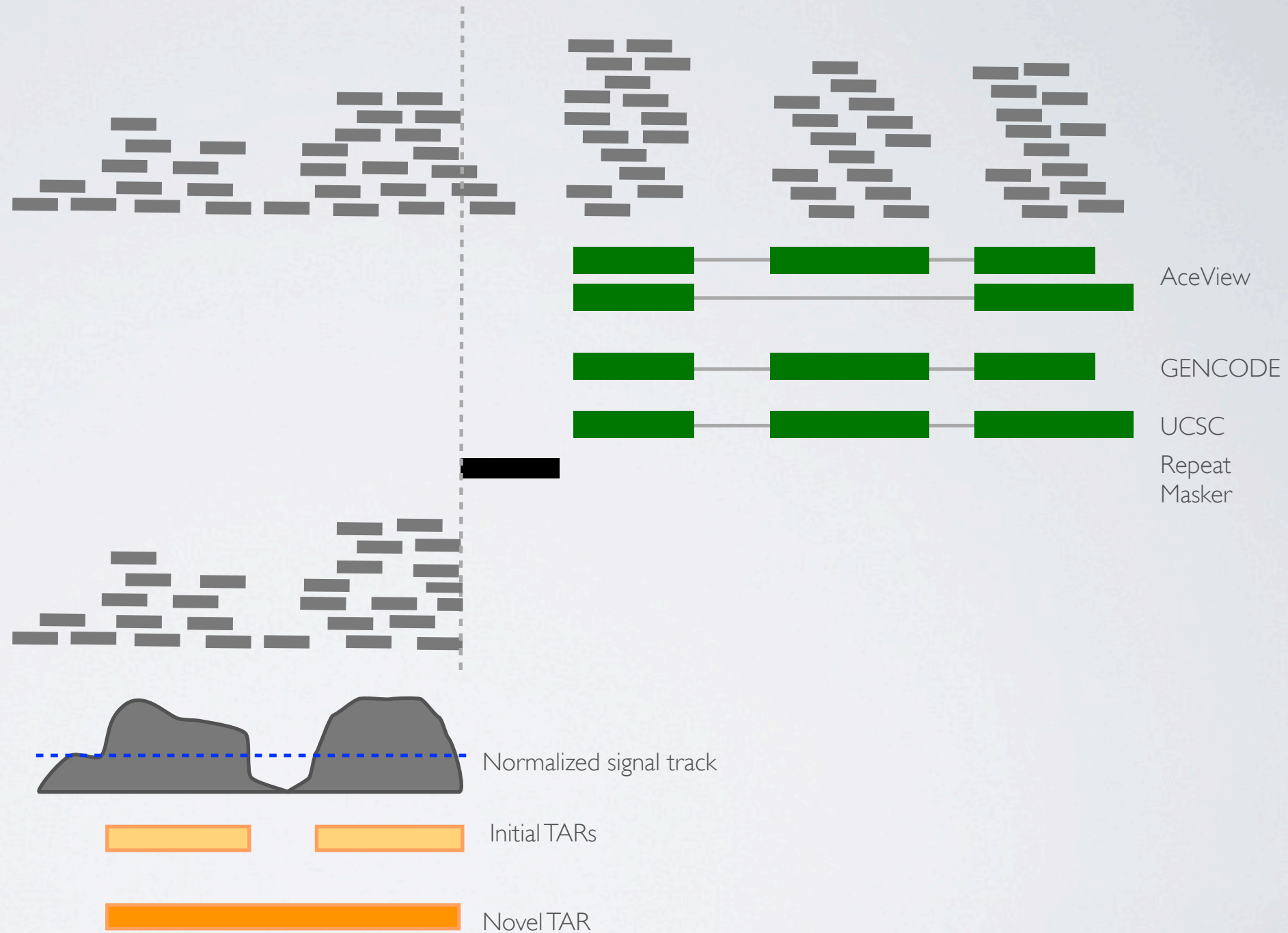
Andrea Sboner - 2011.04.08

DISCOVERY OF NOVEL TARs



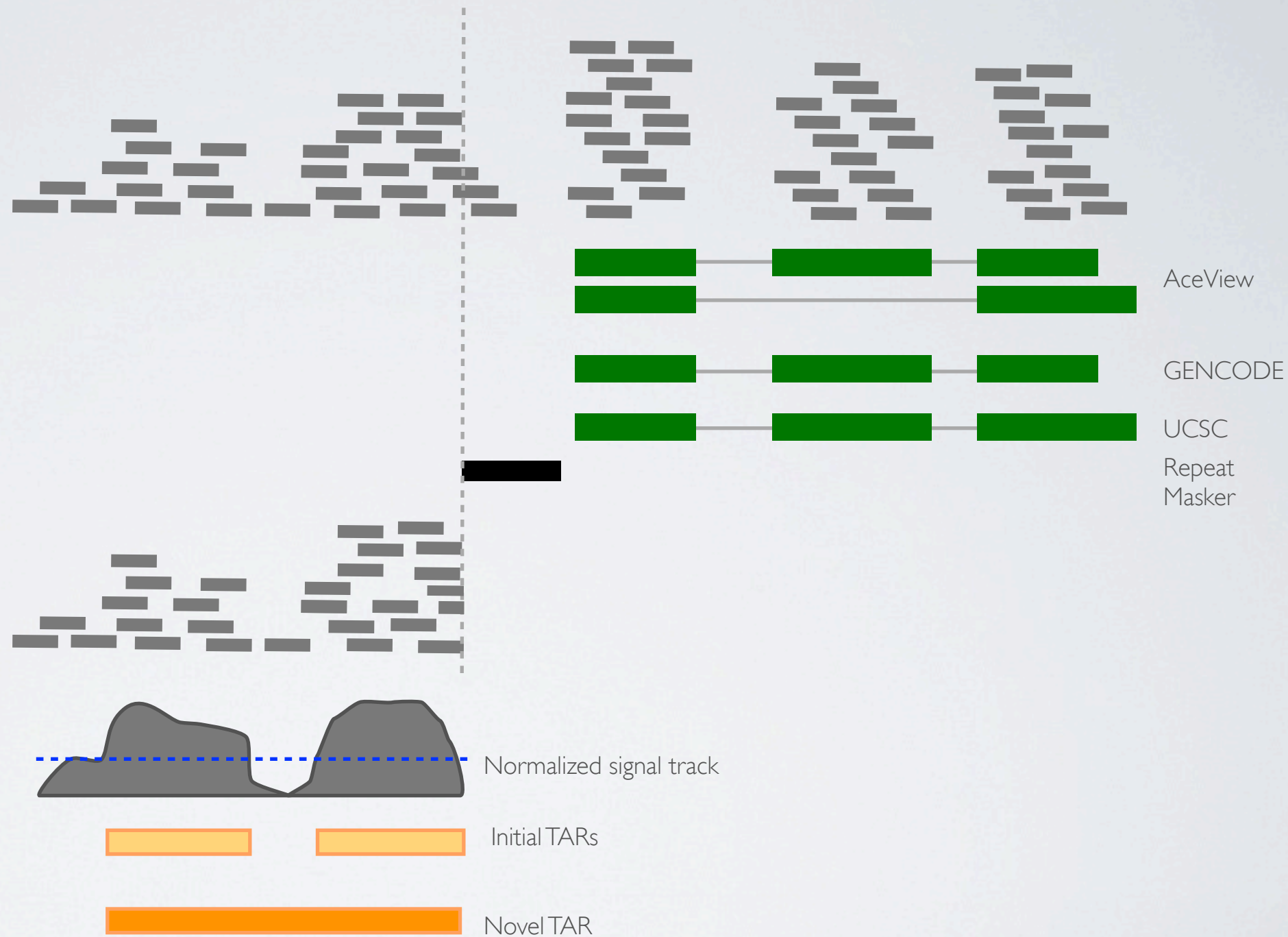
DISCOVERY OF NOVEL TARs

- Reads overlapping gene annotation set -- AceView, GENCODE, UCSC, and repetitive regions are excluded



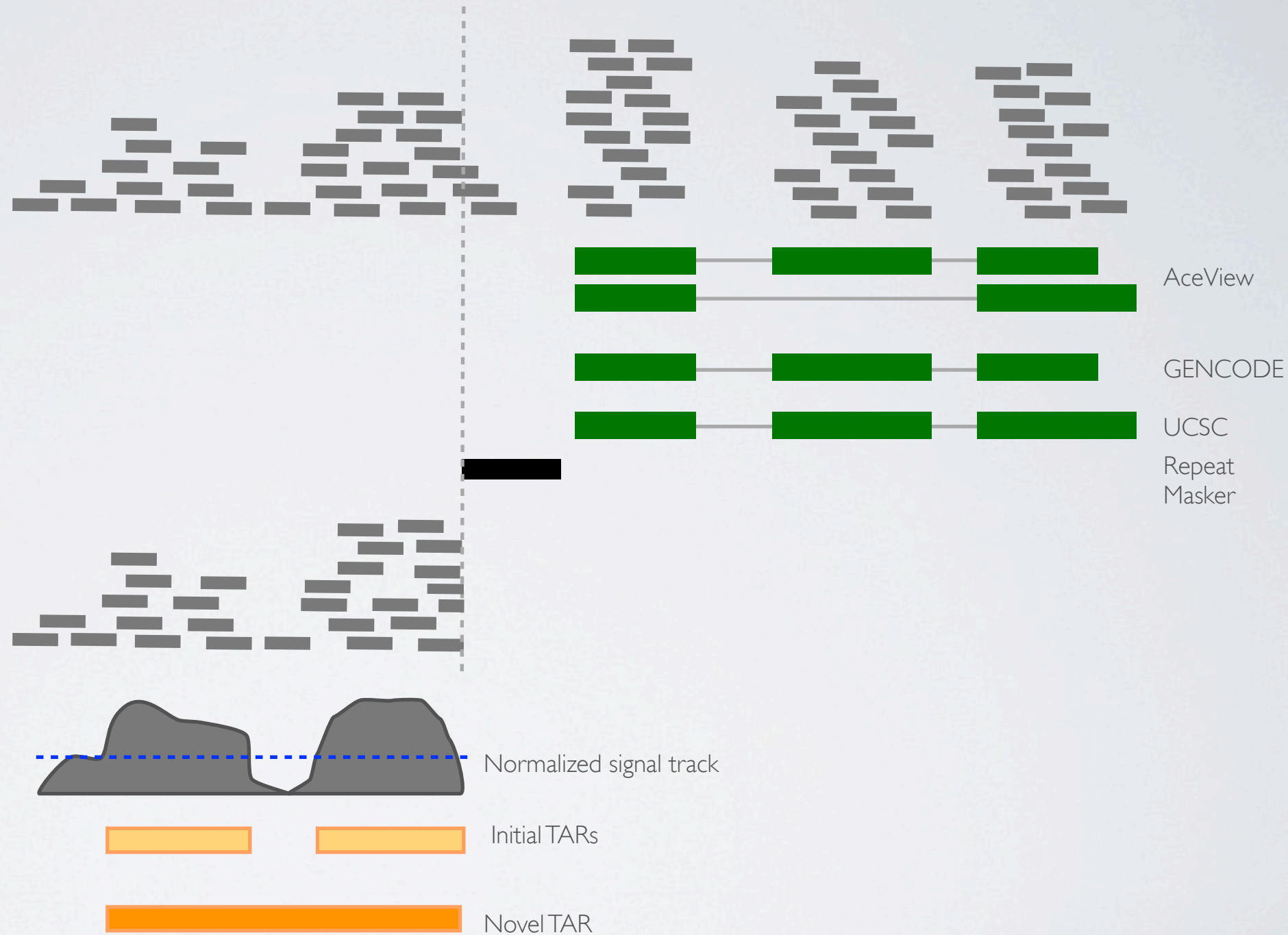
DISCOVERY OF NOVEL TARs

- Reads overlapping gene annotation set -- AceView, GENCODE, UCSC, and repetitive regions are excluded
- Reads from all samples pooled together



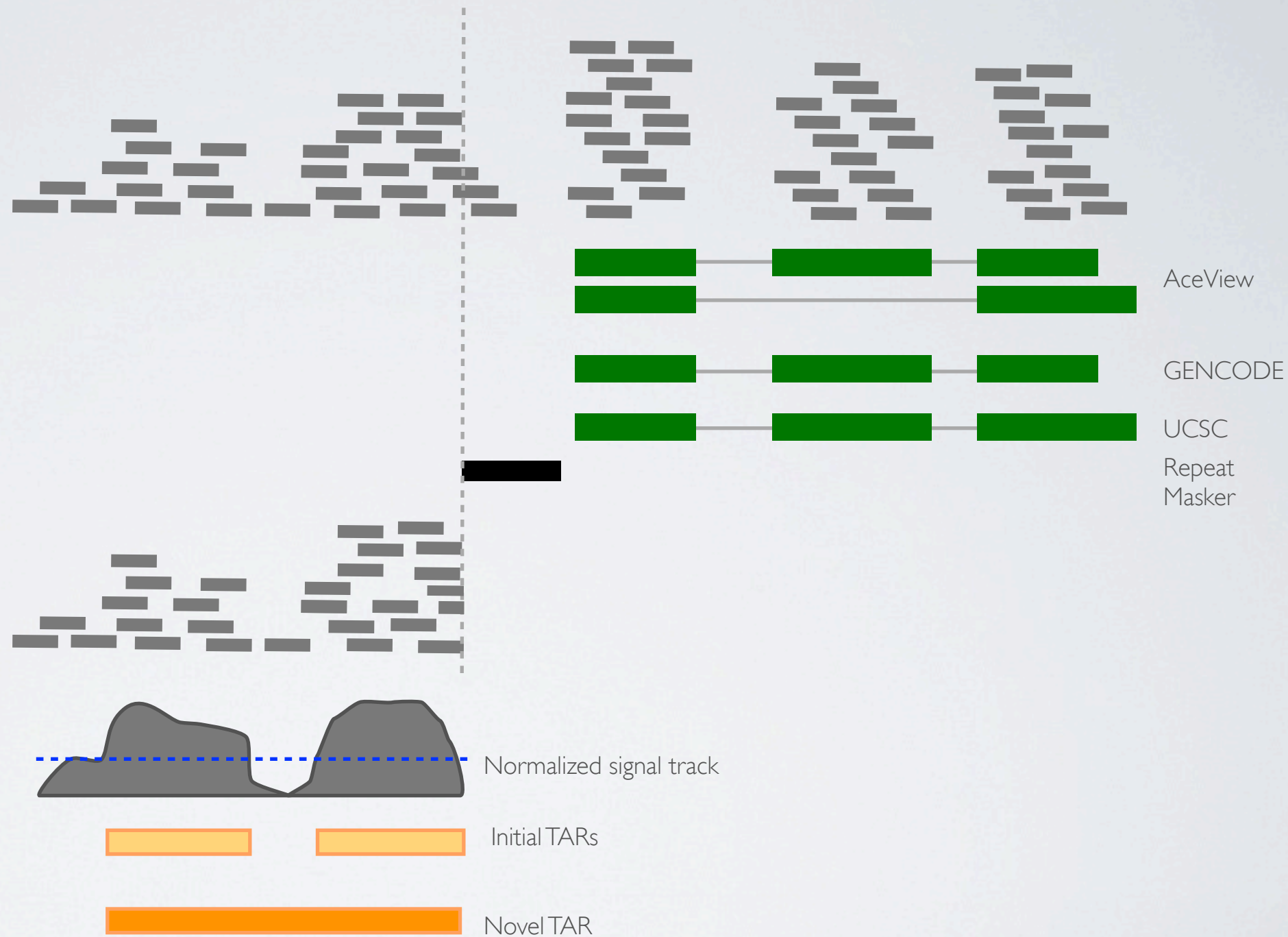
DISCOVERY OF NOVEL TARs

- Reads overlapping gene annotation set -- AceView, GENCODE, UCSC, and repetitive regions are excluded
- Reads from all samples pooled together



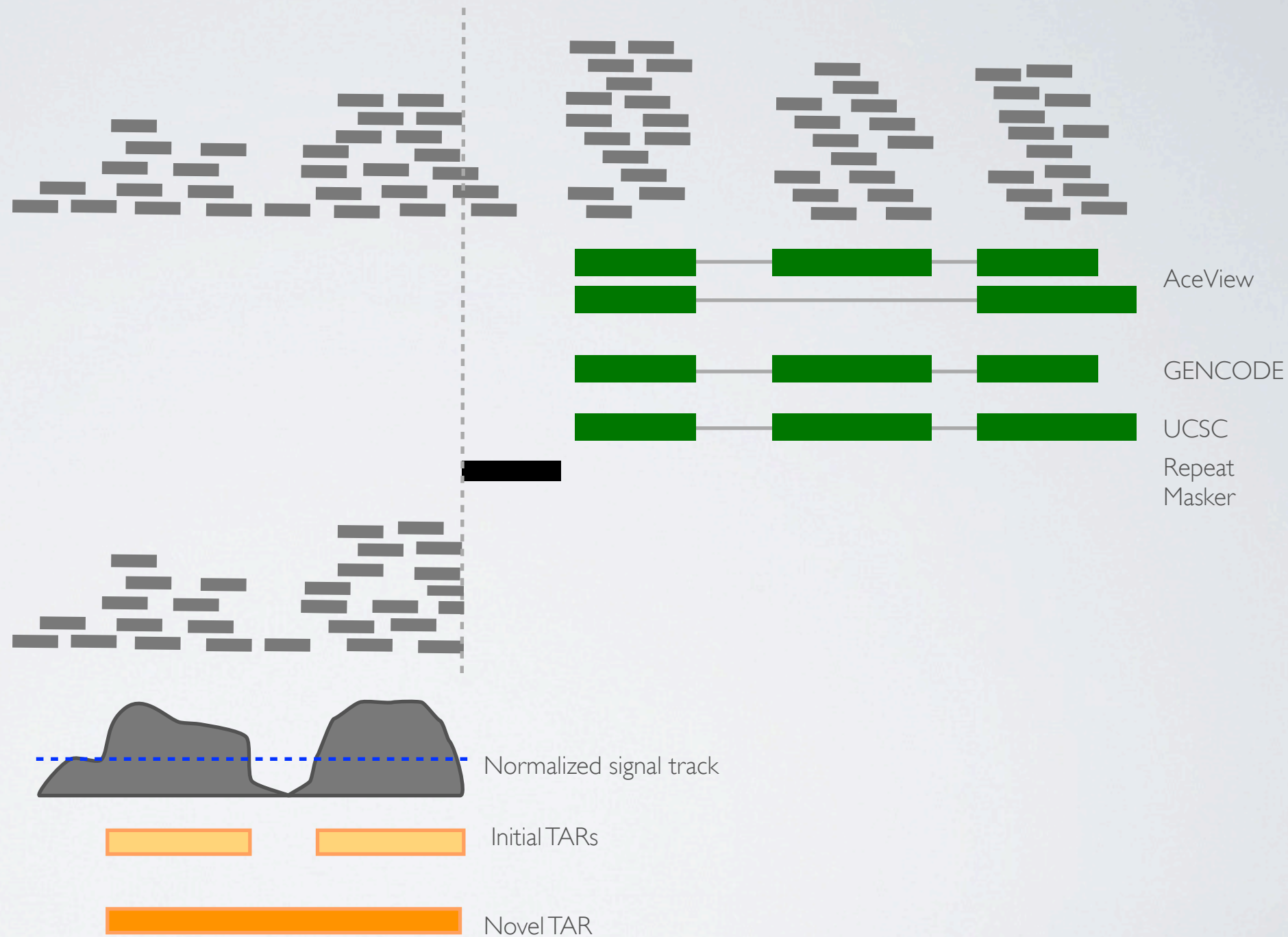
DISCOVERY OF NOVEL TARs

- Reads overlapping gene annotation set -- AceView, GENCODE, UCSC, and repetitive regions are excluded
- Reads from all samples pooled together



DISCOVERY OF NOVEL TARs

- Reads overlapping gene annotation set -- AceView, GENCODE, UCSC, and repetitive regions are excluded
- Reads from all samples pooled together
- Max-gap, min-run algorithm to identify TARs



MIN-GAP, MAX-RUN ALGORITHM: THRESHOLD DEFINITION

$$|\vec{L}| = 0.129Kb$$

MIN-GAP, MAX-RUN ALGORITHM: THRESHOLD DEFINITION

$$r(g) = \frac{n(g)}{l(g)*M}$$

$$|\vec{L}| = 0.129Kb$$

MIN-GAP, MAX-RUN ALGORITHM: THRESHOLD DEFINITION

$$r(g) = \frac{n(g)}{l(g)*M} \quad x(g) = \log_2(r(g) + 1)$$

$$|\vec{L}| = 0.129Kb$$

MIN-GAP, MAX-RUN ALGORITHM: THRESHOLD DEFINITION

$$r(g) = \frac{n(g)}{l(g) * M} \quad x(g) = \log_2(r(g) + 1)$$

- From exon expression values:

Determine the median value of each gene:

$$\tilde{x}(g)$$

$$x_c = \text{quantile}(\vec{\tilde{x}}, 0.05)$$

$$\forall \tilde{x}(g) > 0$$

$$r_c = 2^{x_c} - 1$$

$$|\vec{L}| = 0.129Kb$$

$$t_{bgr} = |\vec{L}| \cdot r_c$$

MIN-GAP, MAX-RUN ALGORITHM: THRESHOLD DEFINITION

$$r(g) = \frac{n(g)}{l(g) * M} \quad x(g) = \log_2(r(g) + 1)$$

- From exon expression values:

Determine the median value of each gene:

$$\tilde{x}(g)$$

- Computed the 5th percentile of expressed genes:

$$x_c = \text{quantile}(\vec{\tilde{x}}, 0.05)$$

$$\forall \tilde{x}(g) > 0$$

$$r_c = 2^{x_c} - 1$$

$$|\vec{L}| = 0.129Kb$$

$$t_{bgr} = |\vec{L}| \cdot r_c$$

MIN-GAP, MAX-RUN ALGORITHM: THRESHOLD DEFINITION

$$r(g) = \frac{n(g)}{l(g) * M} \quad x(g) = \log_2(r(g) + 1)$$

- From exon expression values:

Determine the median value of each gene:

$$\tilde{x}(g)$$

- Computed the 5th percentile of expressed genes: $x_c = \text{quantile}(\vec{\tilde{x}}, 0.05)$
 $\forall \tilde{x}(g) > 0$
- Compute the corresponding RPKM value: $r_c = 2^{x_c} - 1$

$$|\vec{L}| = 0.129Kb$$

$$t_{bgr} = |\vec{L}| \cdot r_c$$

MIN-GAP, MAX-RUN ALGORITHM: THRESHOLD DEFINITION

$$r(g) = \frac{n(g)}{l(g)*M} \quad x(g) = \log_2(r(g) + 1)$$

- From exon expression values:

Determine the median value of each gene:

$$\tilde{x}(g)$$

- Computed the 5th percentile of expressed genes: $x_c = \text{quantile}(\vec{\tilde{x}}, 0.05)$
 $\forall \tilde{x}(g) > 0$
- Compute the corresponding RPKM value: $r_c = 2^{x_c} - 1$
- Consider the median length of GENCODE exons: $|\vec{L}| = 0.129Kb$

$$t_{bgr} = |\vec{L}| \cdot r_c$$

MIN-GAP, MAX-RUN ALGORITHM: THRESHOLD DEFINITION

$$r(g) = \frac{n(g)}{l(g)*M} \quad x(g) = \log_2(r(g) + 1)$$

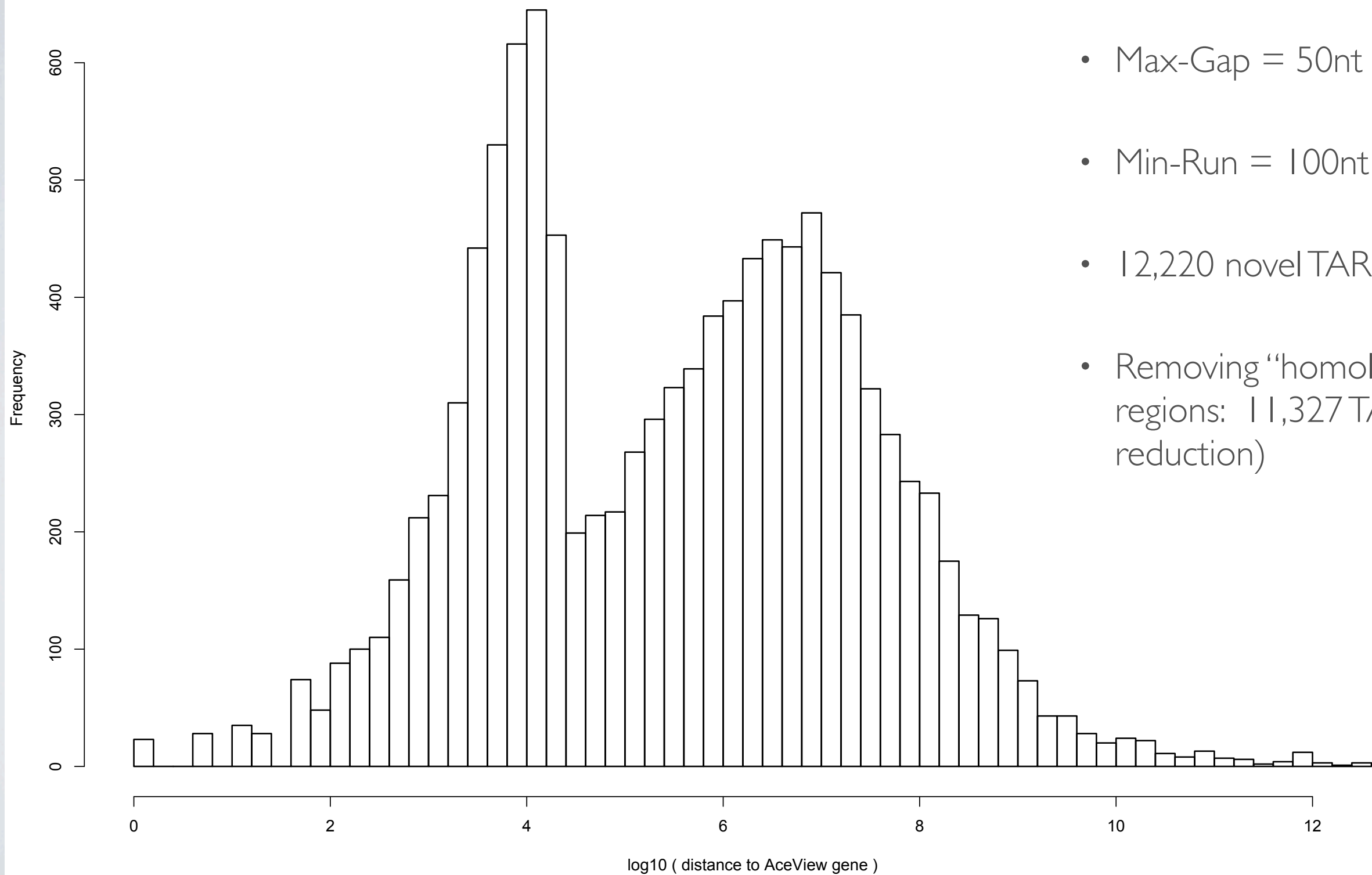
- From exon expression values:

Determine the median value of each gene:

$$\tilde{x}(g)$$

- Computed the 5th percentile of expressed genes: $x_c = \text{quantile}(\vec{\tilde{x}}, 0.05)$
 $\forall \tilde{x}(g) > 0$
- Compute the corresponding RPKM value: $r_c = 2^{x_c} - 1$
- Consider the median length of GENCODE exons: $|\vec{L}| = 0.129Kb$
- Define *threshold* (normalized per million mapped reads): $t_{bgr} = |\vec{L}| \cdot r_c$

RESULTS



- Threshold = 0.019617
- Max-Gap = 50nt
- Min-Run = 100nt
- 12,220 novel TARs
- Removing “homologous” regions: 11,327 TARs (~7% reduction)

EXAMPLE

