

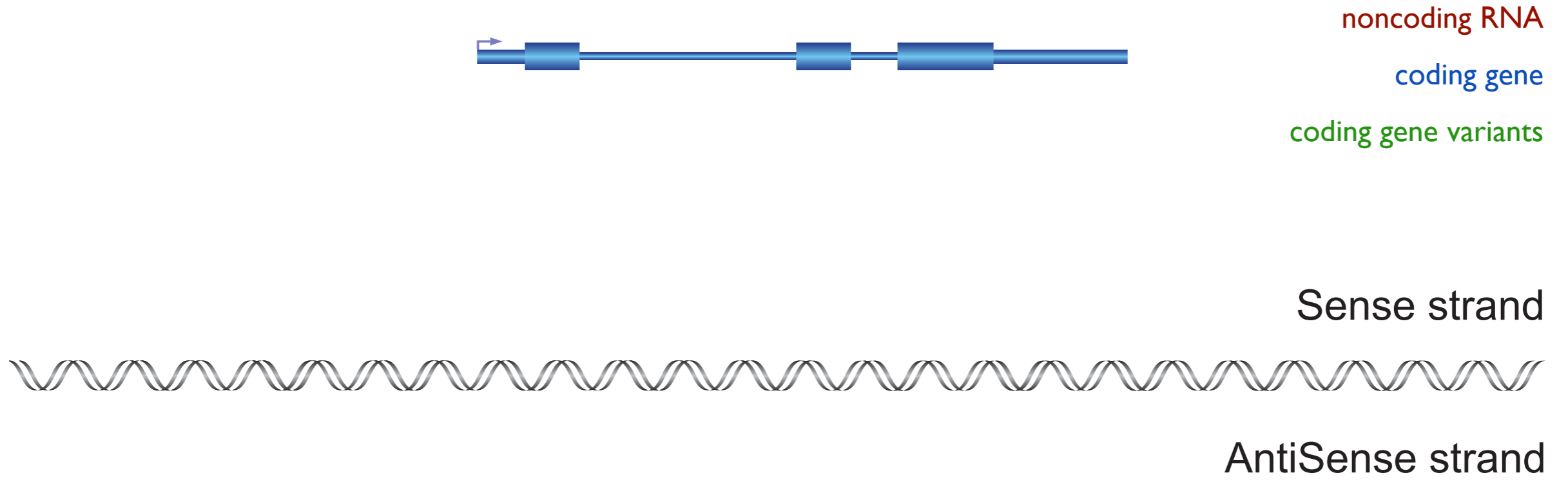
Uncovering hidden genes in intergenic GWAS regions

Mike Clark
Cancer Genomics Research Laboratory
Garvan Institute for Medical Research
m.clark@garvan.org.au

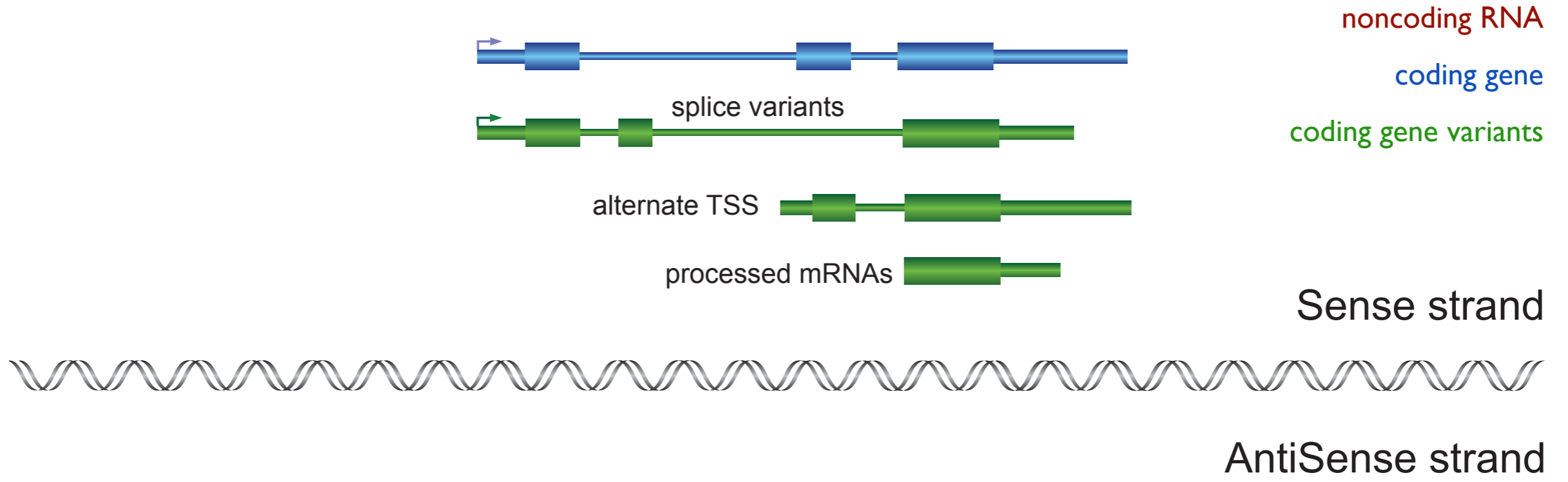
Uncovering hidden genes in intergenic GWAS regions with RNA capture sequencing

- The pervasive transcriptome and “intergenic” genome-wide association study (GWAS) regions
- Detecting transcription with RNA Capture sequencing
- Using capture sequencing for novel gene discovery in human “intergenic” GWAS regions

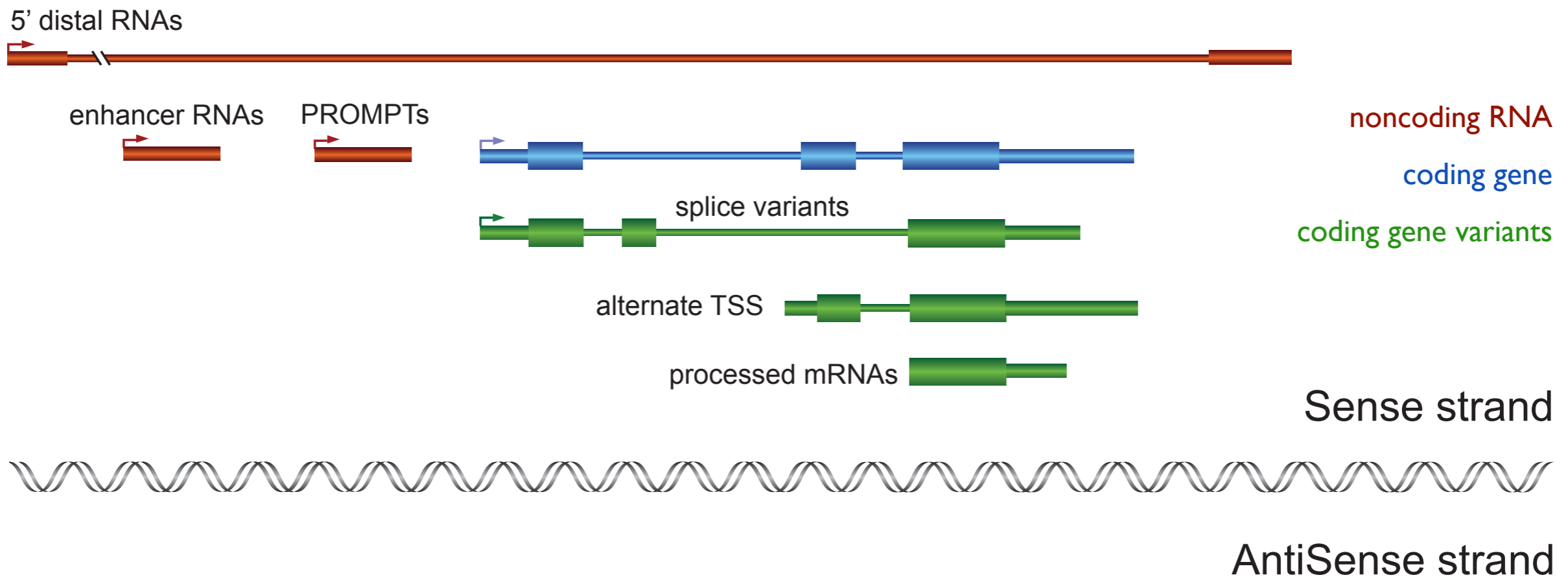
Transcription is complex, and pervasive across the genome



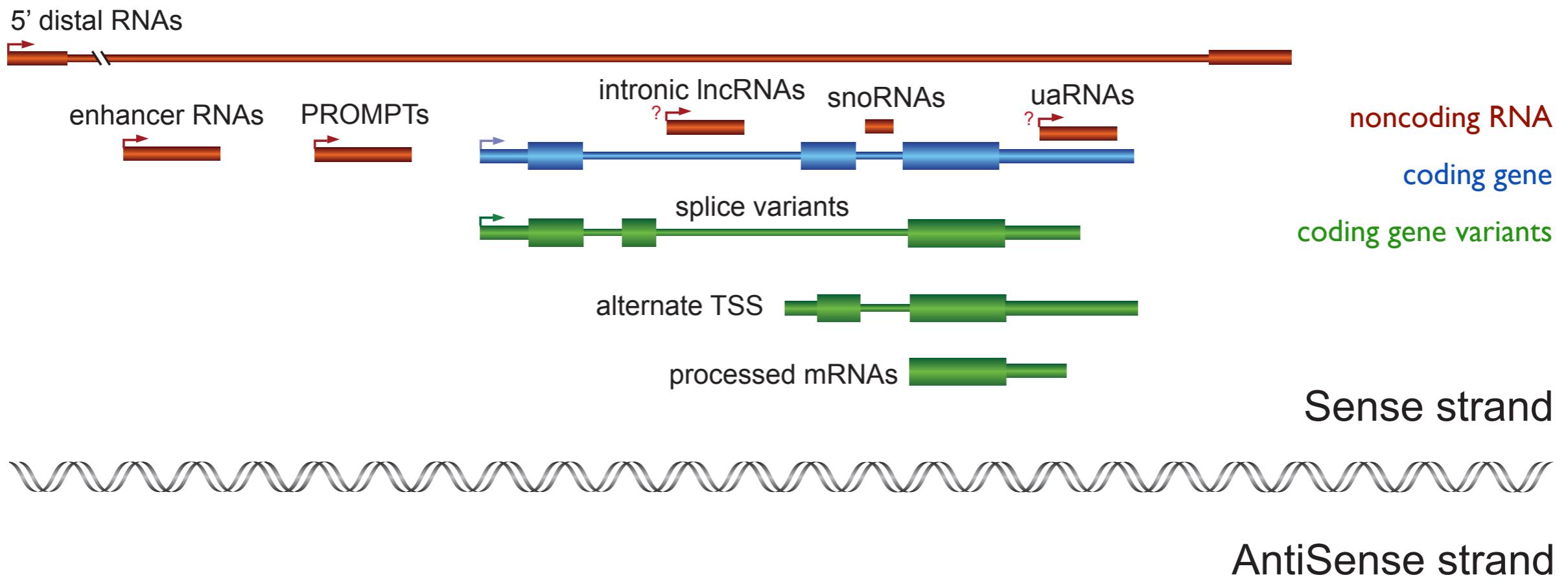
Transcription is complex, and pervasive across the genome



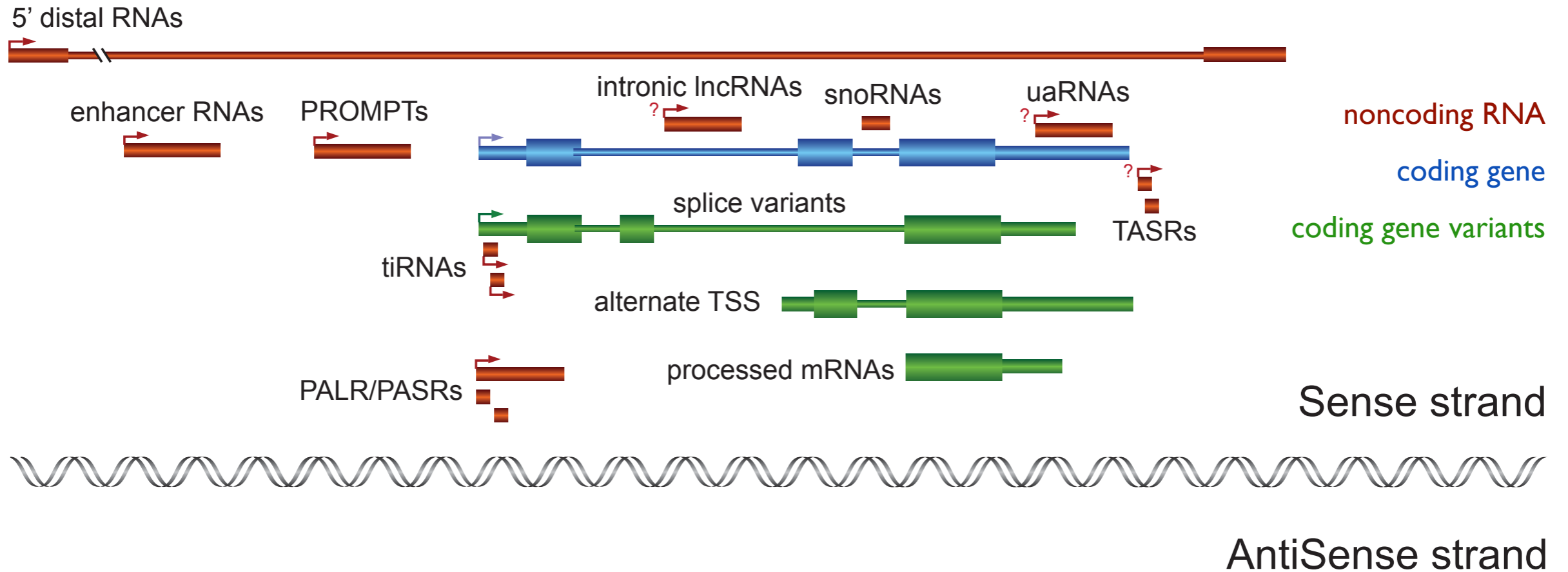
Transcription is complex, and pervasive across the genome



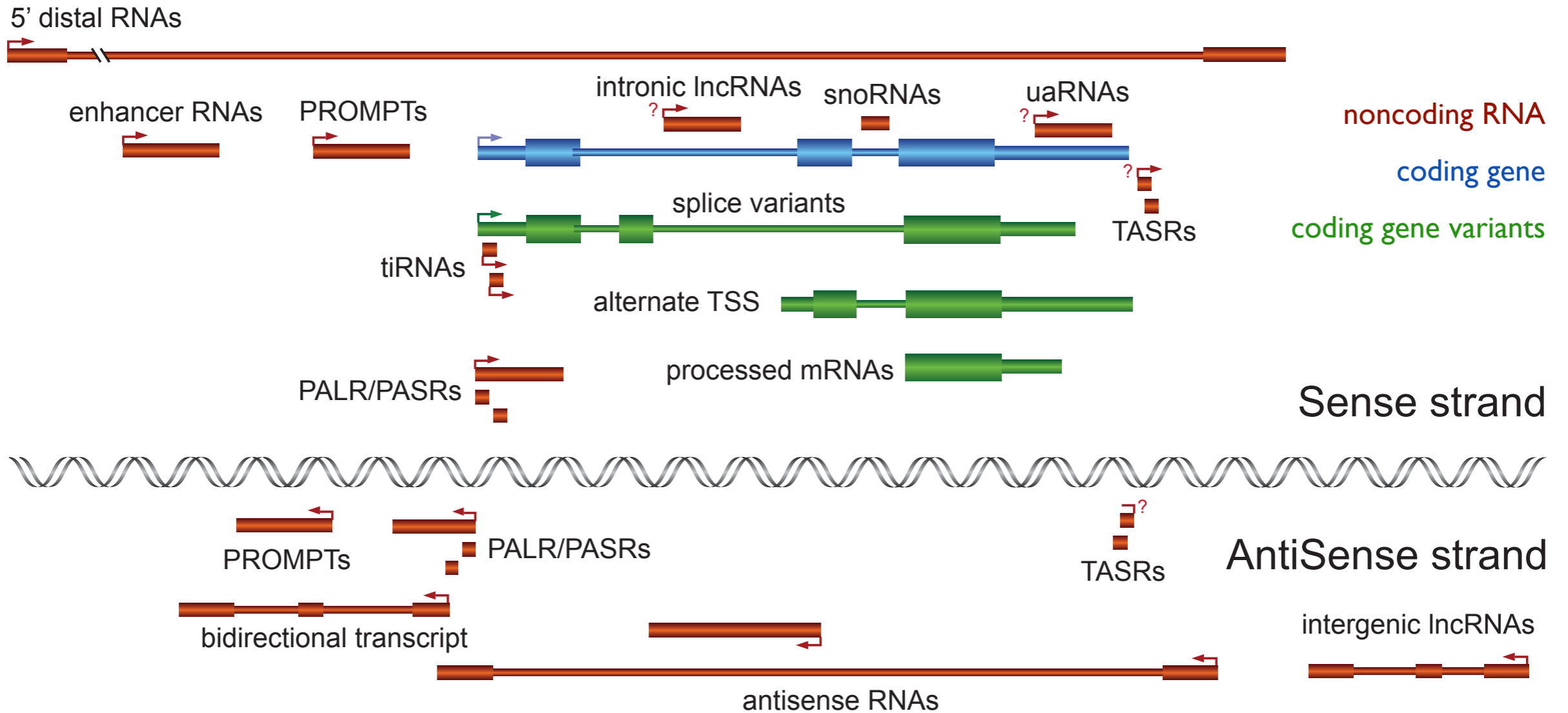
Transcription is complex, and pervasive across the genome



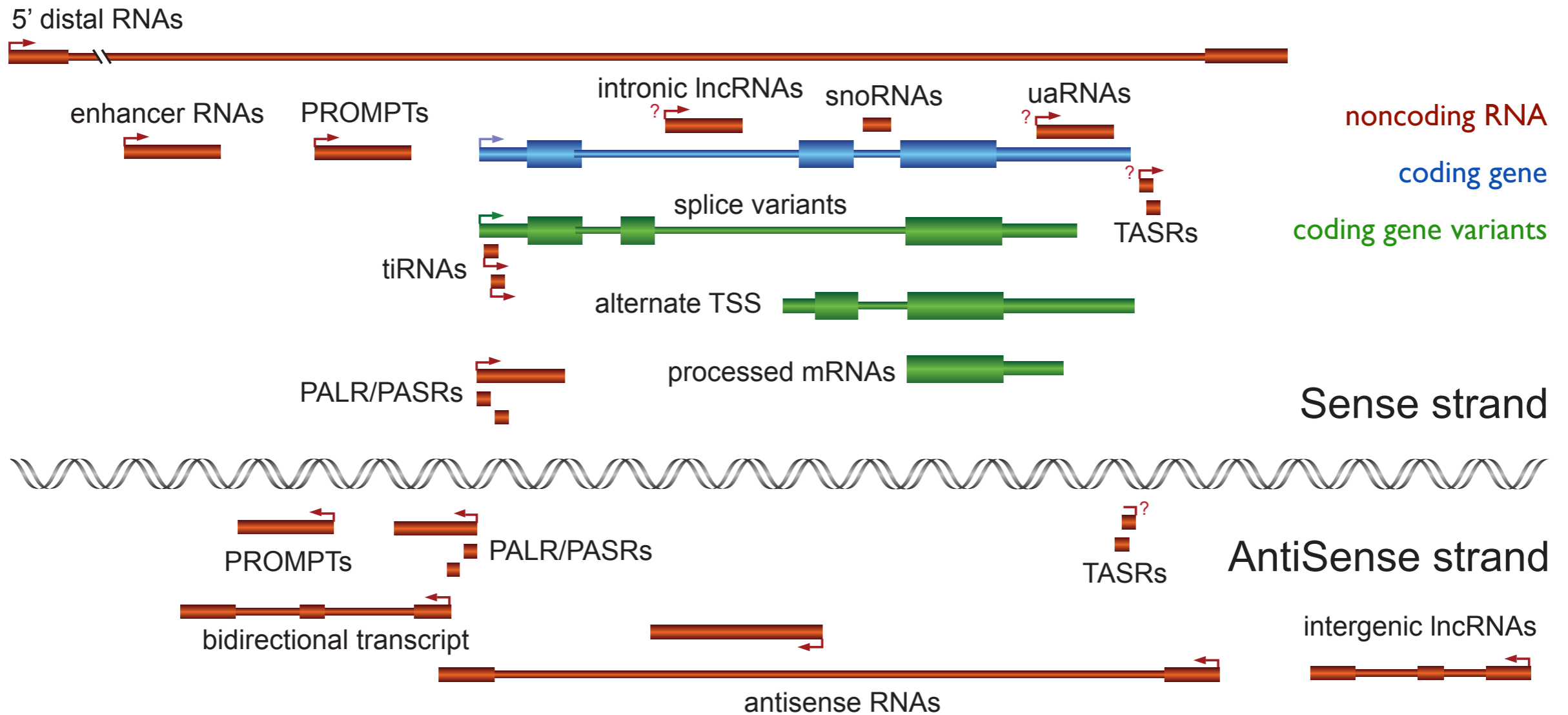
Transcription is complex, and pervasive across the genome



Transcription is complex, and pervasive across the genome

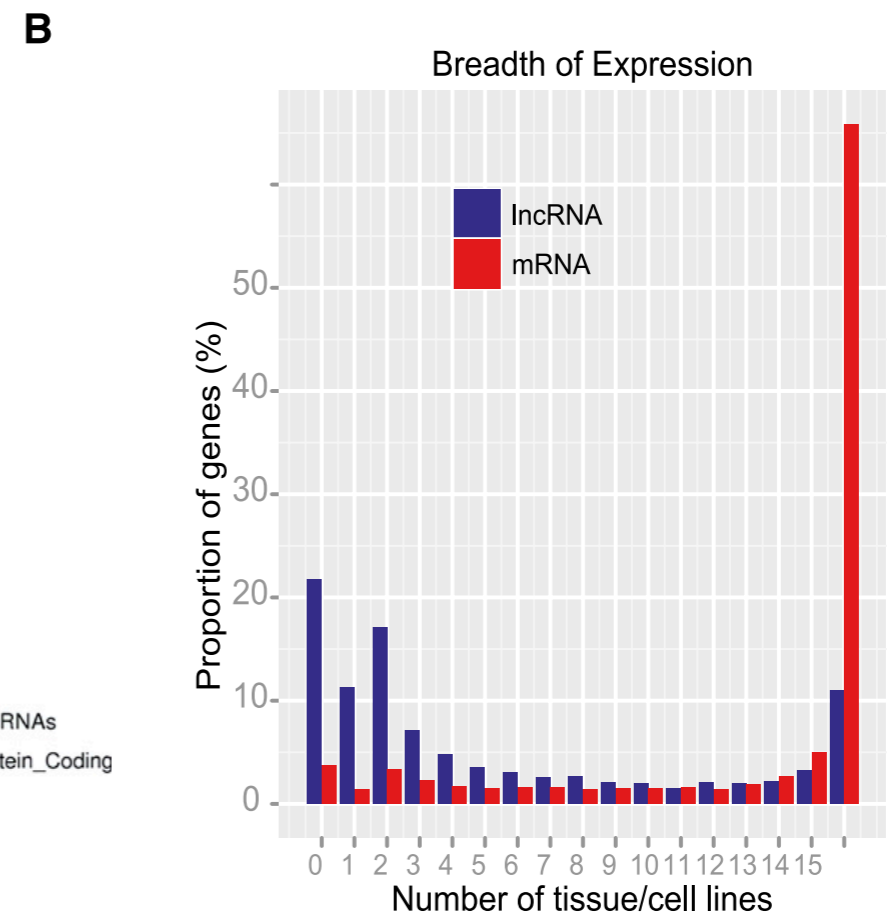
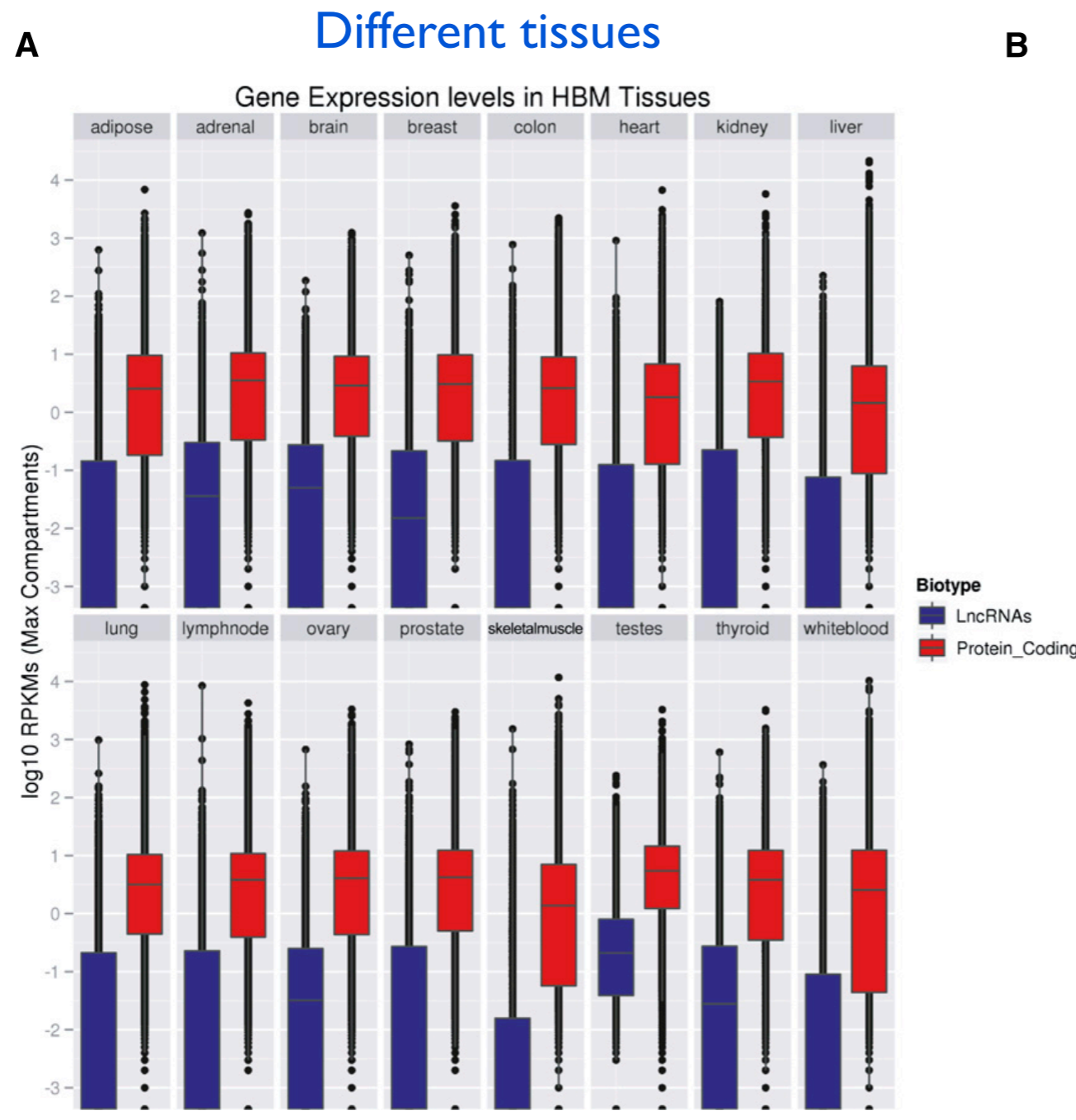


Transcription is complex, and pervasive across the genome



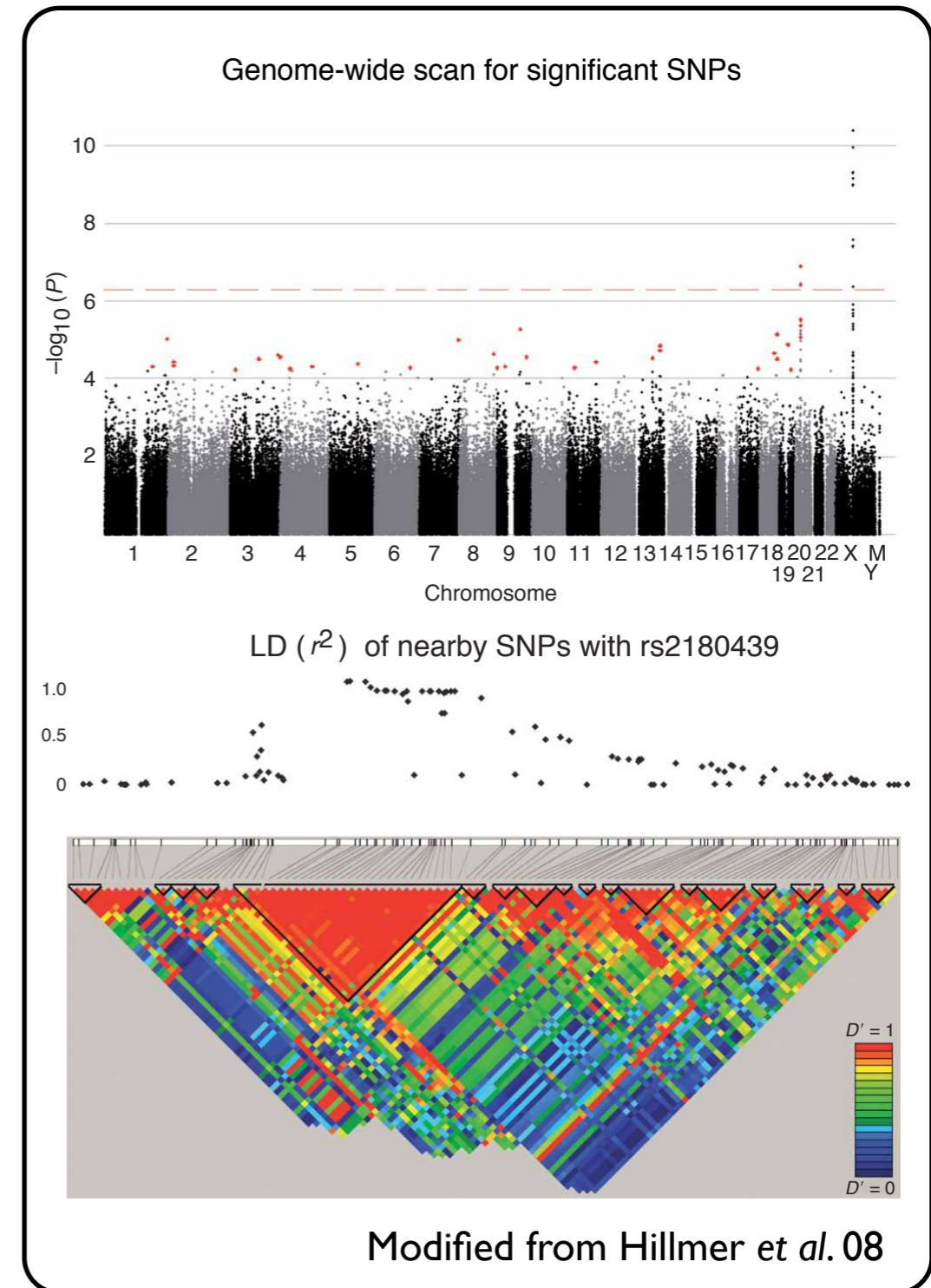
- At least 75% of the human genome is transcribed, though a lot of this is transcribed at low levels (ENCODE 2012)

Transcripts can be lowly and specifically expressed



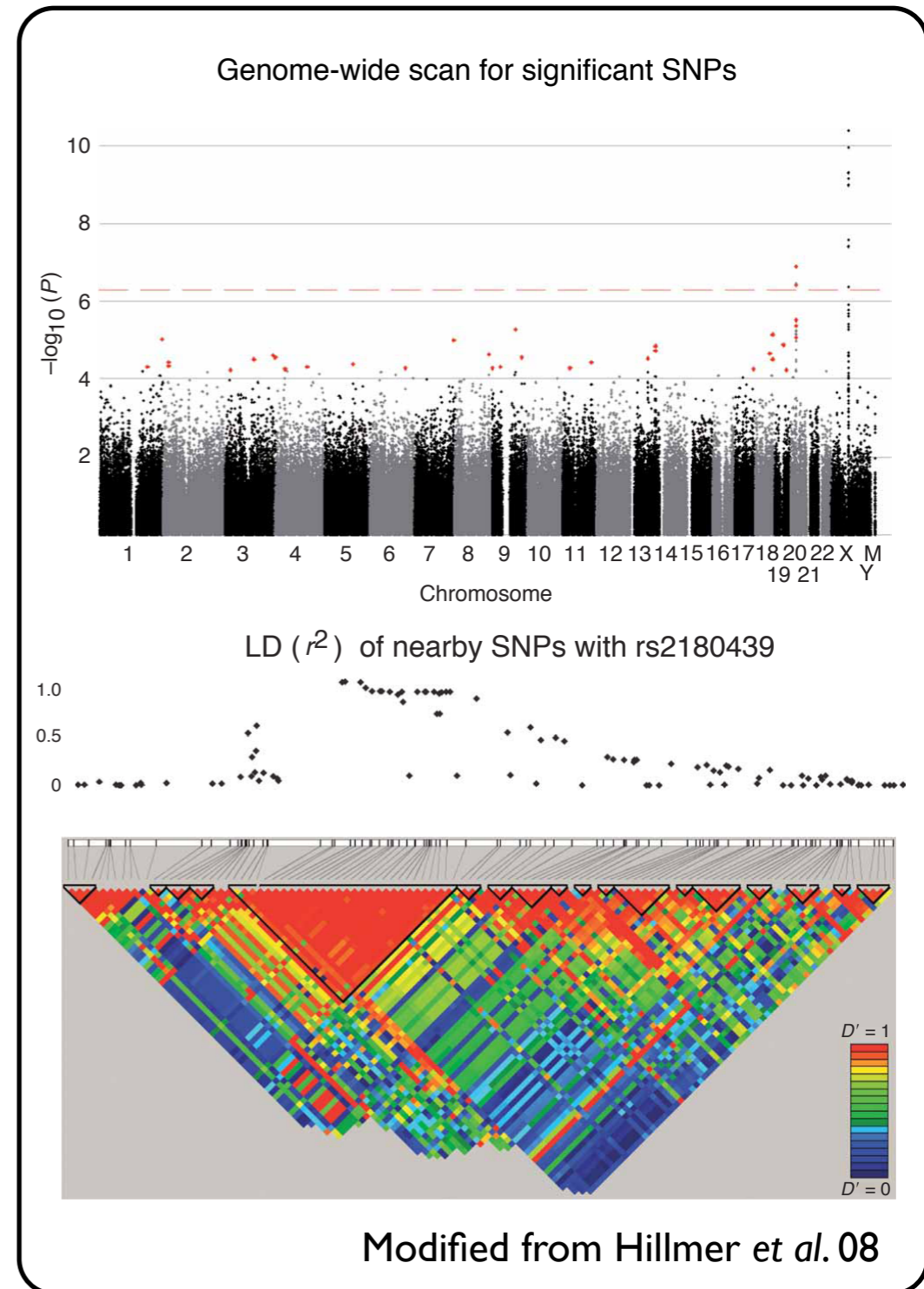
Red: coding genes. Blue: LncRNAs

“Intergenic” GWAS regions



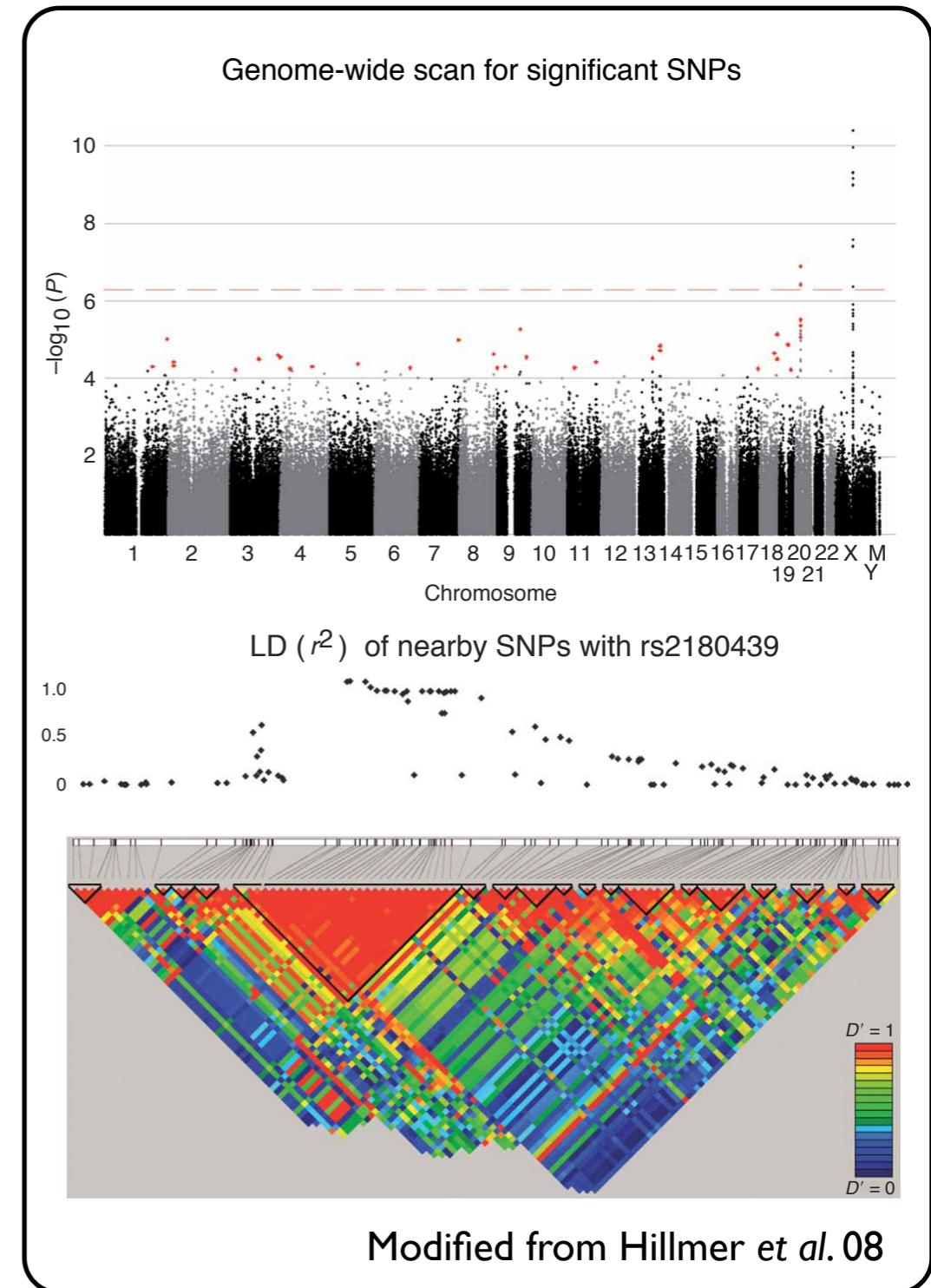
“Intergenic” GWAS regions

- Genome-wide association studies identify SNPs with a significant association to a disease or trait



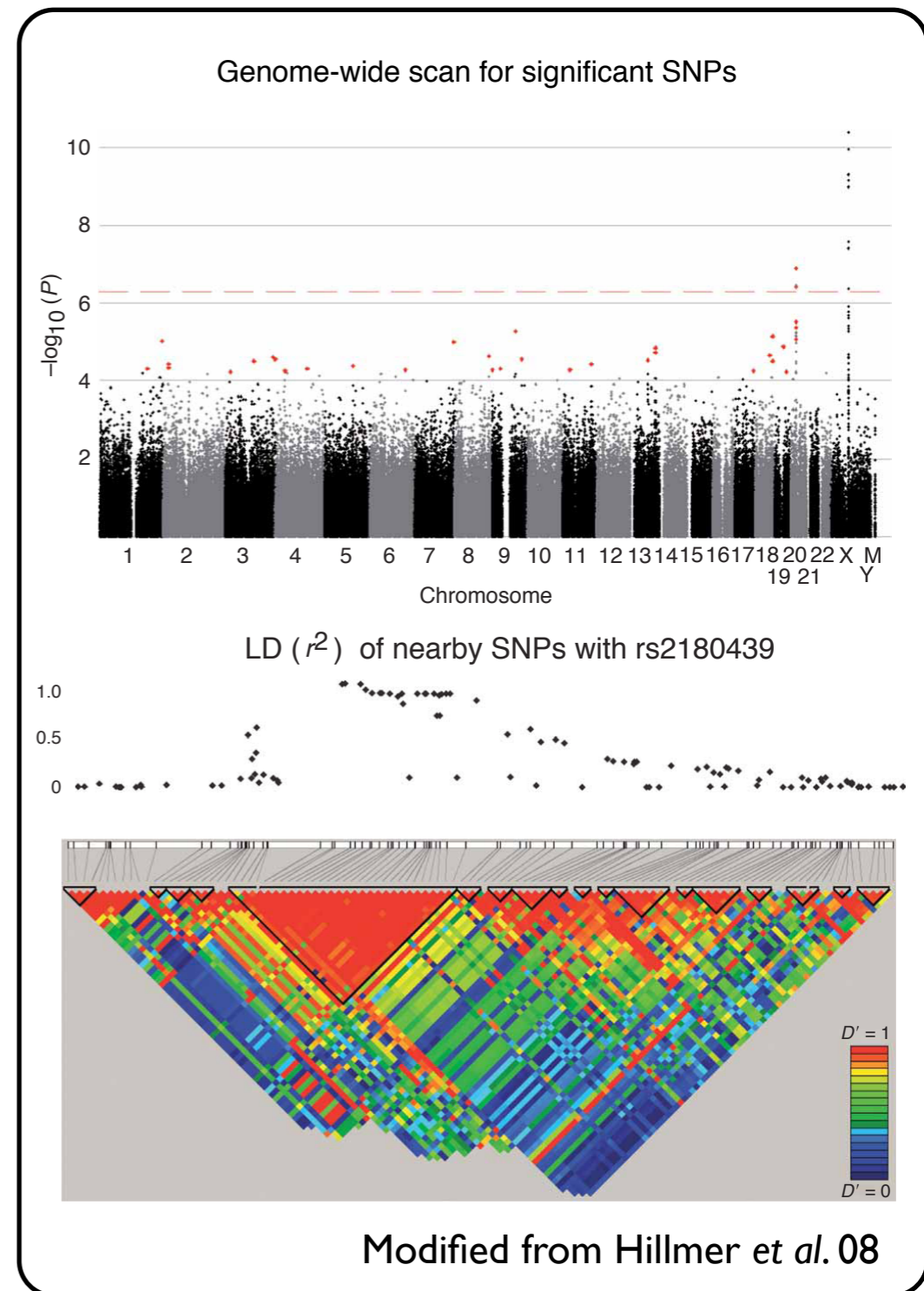
“Intergenic” GWAS regions

- Genome-wide association studies identify SNPs with a significant association to a disease or trait
- 88% of identified SNPs are intergenic or intronic (Hindorff *et al.* 09)



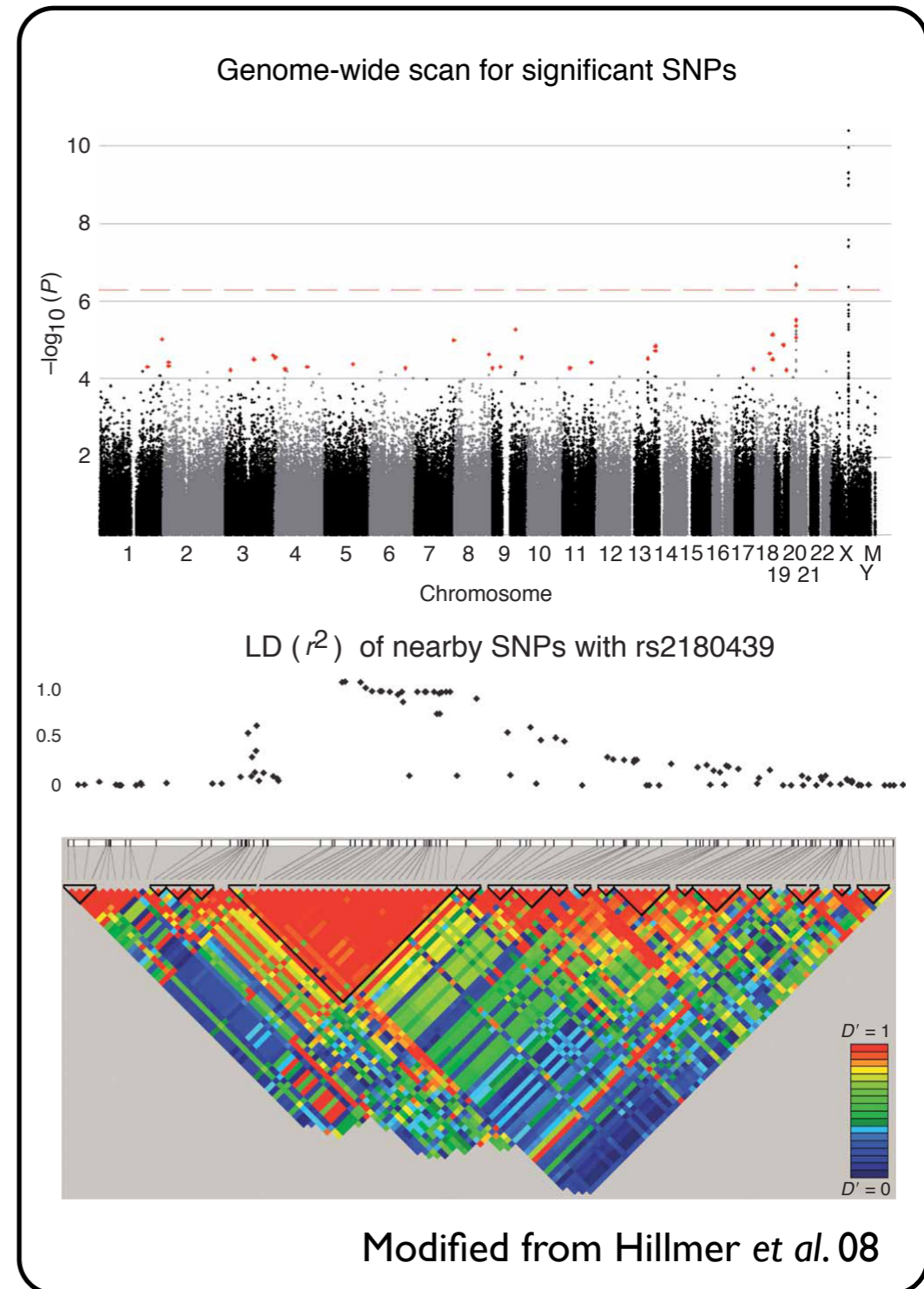
“Intergenic” GWAS regions

- Genome-wide association studies identify SNPs with a significant association to a disease or trait
- 88% of identified SNPs are intergenic or intronic (Hindorff *et al.* 09)
- SNP may not be causative allele, but this should be in a region of linkage disequilibrium (LD) around SNP



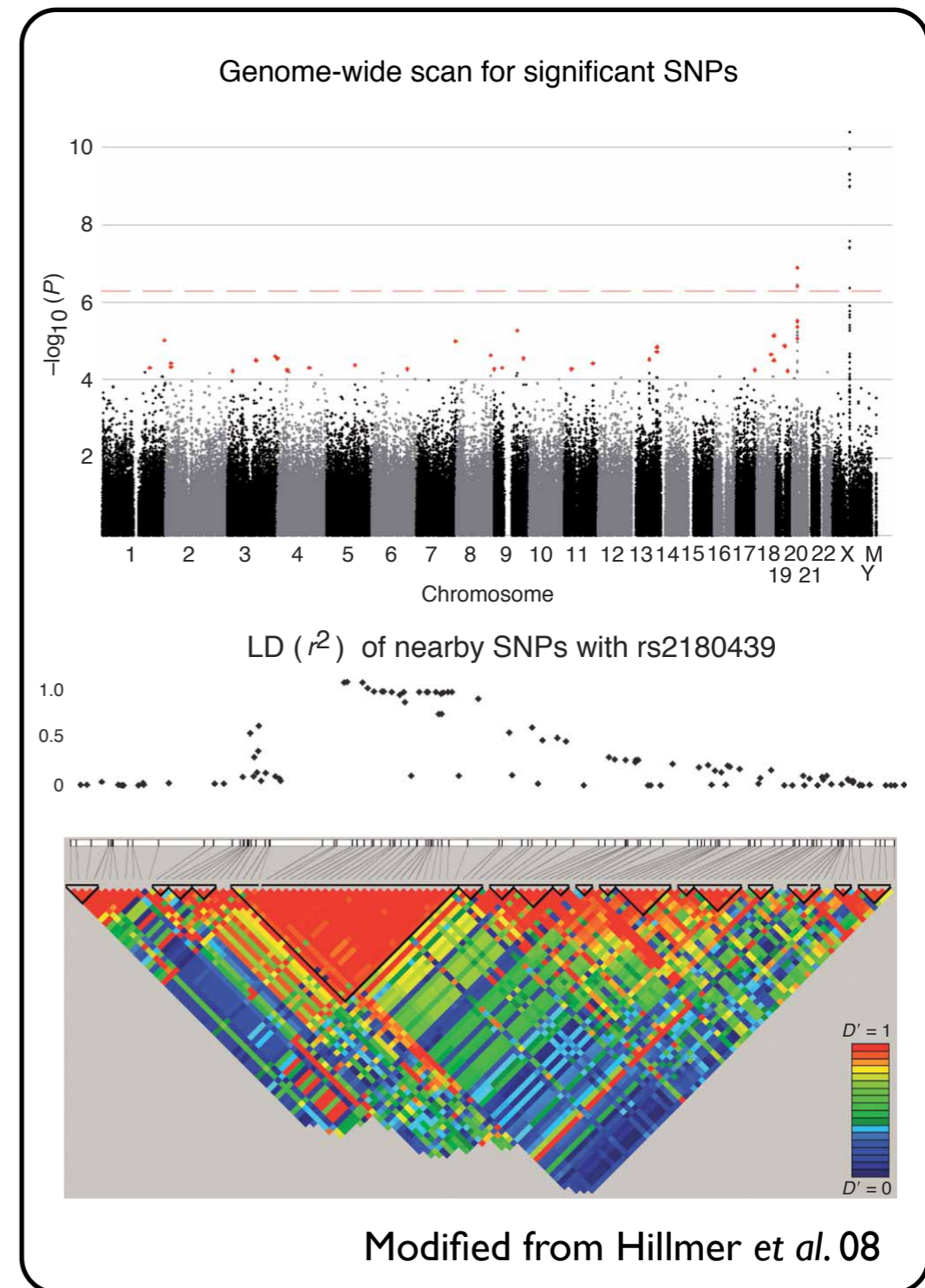
“Intergenic” GWAS regions

- Genome-wide association studies identify SNPs with a significant association to a disease or trait
- 88% of identified SNPs are intergenic or intronic (Hindorff *et al.* 09)
- SNP may not be causative allele, but this should be in a region of linkage disequilibrium (LD) around SNP
- Hundreds of GWAS LD regions are intergenic.



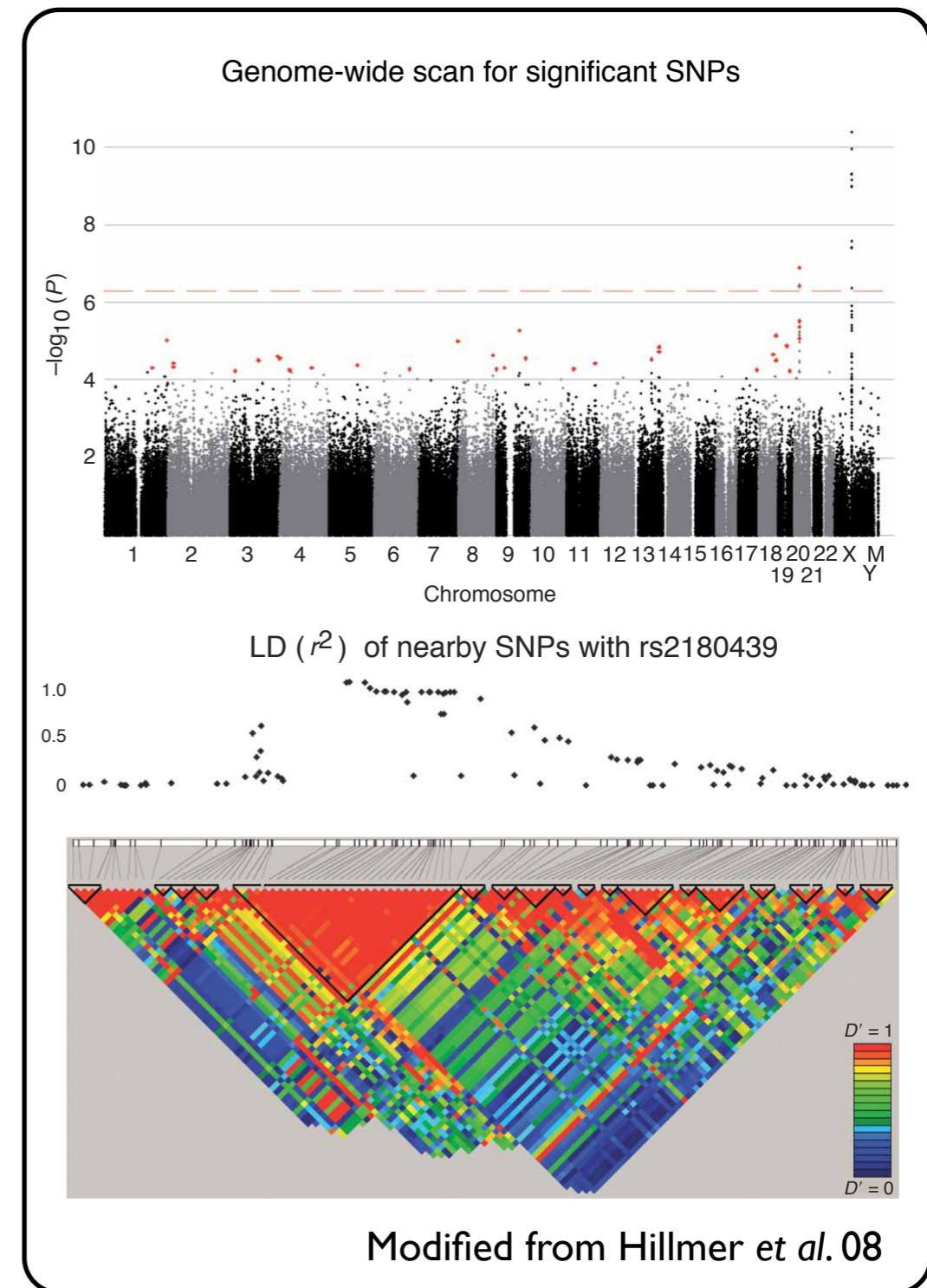
“Intergenic” GWAS regions

- Genome-wide association studies identify SNPs with a significant association to a disease or trait
- 88% of identified SNPs are intergenic or intronic (Hindorff *et al.* 09)
- SNP may not be causative allele, but this should be in a region of linkage disequilibrium (LD) around SNP
- Hundreds of GWAS LD regions are intergenic.
- Are these really empty of genes, or have we just not found them yet?



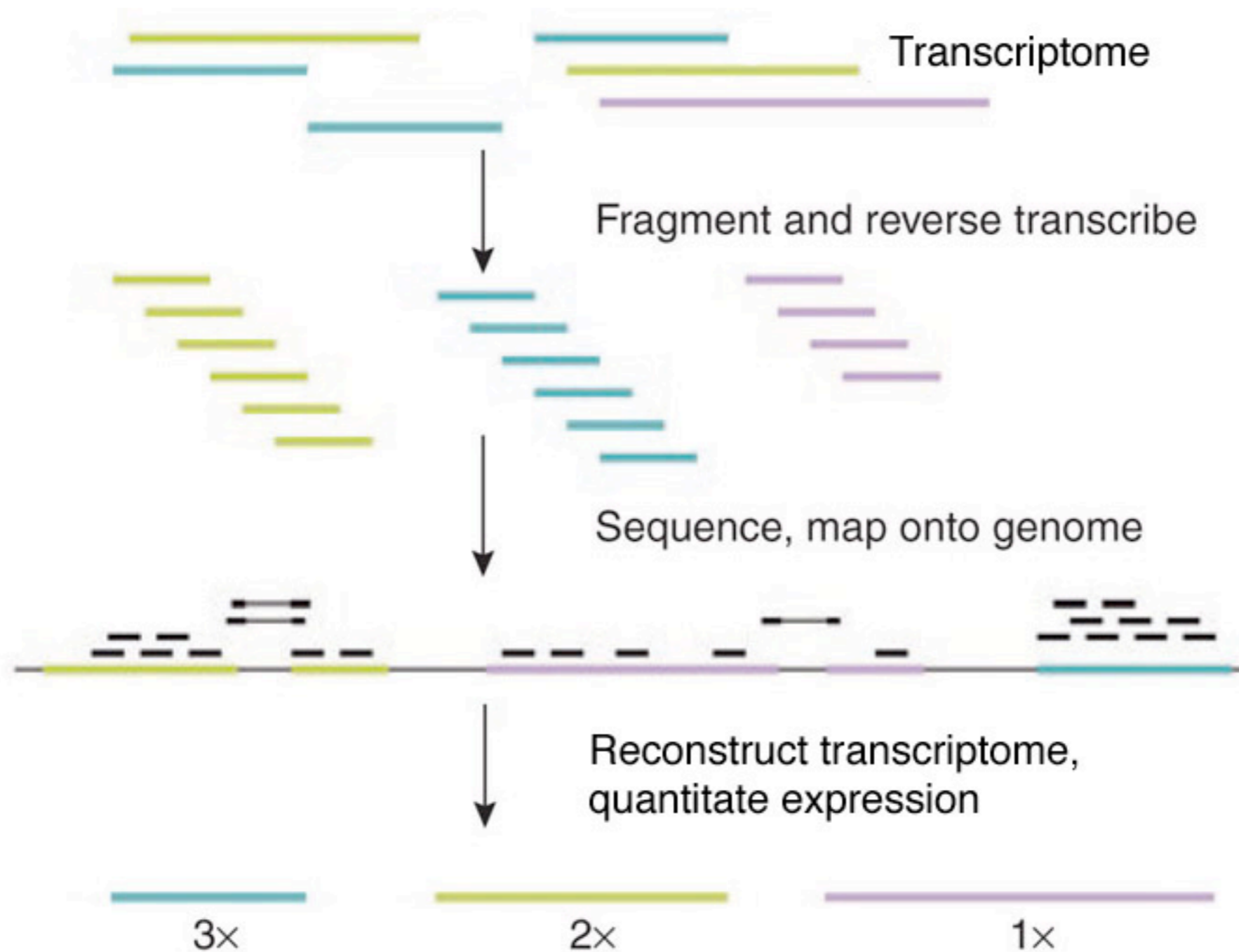
“Intergenic” GWAS regions

- Genome-wide association studies identify SNPs with a significant association to a disease or trait
- 88% of identified SNPs are intergenic or intronic (Hindorff *et al.* 09)
- SNP may not be causative allele, but this should be in a region of linkage disequilibrium (LD) around SNP
- Hundreds of GWAS LD regions are intergenic.
- Are these really empty of genes, or have we just not found them yet?



- Aim to discover novel genes within these intergenic GWAS regions

Standard RNA sequencing outline



RNAseq is mostly re-sequencing of highly expressed genes



The human transcriptome is complex. RNA Sequencing is like eating some of these jelly beans at random. You have a good chance of picking common colours but a low chance of finding rare ones.

RNAseq is mostly re-sequencing of highly expressed genes



The human transcriptome is complex. RNA Sequencing is like eating some of these jelly beans at random. You have a good chance of picking common colours but a low chance of finding rare ones.

RNAseq is mostly re-sequencing of highly expressed genes



The human transcriptome is complex. RNA Sequencing is like eating some of these jelly beans at random. You have a good chance of picking common colours but a low chance of finding rare ones.

Limitations of RNAseq: Rare or restricted transcripts

Limitations of RNAseq: Rare or restricted transcripts

- Rare transcripts:

The most abundant 1.5% of transcripts take up almost half the RNAseq reads, while only 1% of reads are measuring the least abundant 44% of detected transcripts (Jiang *et al.* 11).

Limitations of RNAseq: Rare or restricted transcripts

- Rare transcripts:

The most abundant 1.5% of transcripts take up almost half the RNAseq reads, while only 1% of reads are measuring the least abundant 44% of detected transcripts (Jiang *et al.* 11).

- Transcripts with restricted expression:

A transcript expressed at an average level (10 copies per cell) in 0.1% of brain cells (~170 million cells), requires ~400 million reads for 1x coverage.

8x coverage is required for robust assembly (Jiang *et al.* 11)..

Limitations of RNAseq: Rare or restricted transcripts

- Rare transcripts:

The most abundant 1.5% of transcripts take up almost half the RNAseq reads, while only 1% of reads are measuring the least abundant 44% of detected transcripts (Jiang *et al.* 11).

- Transcripts with restricted expression:

A transcript expressed at an average level (10 copies per cell) in 0.1% of brain cells (~170 million cells), requires ~400 million reads for 1x coverage.

8x coverage is required for robust assembly (Jiang *et al.* 11)..

- So how do we find the rare or restricted transcripts likely present in intergenic GWAS regions?

Limitations of RNAseq: Rare or restricted transcripts

- Rare transcripts:

The most abundant 1.5% of transcripts take up almost half the RNAseq reads, while only 1% of reads are measuring the least abundant 44% of detected transcripts (Jiang *et al.* 11).

- Transcripts with restricted expression:

A transcript expressed at an average level (10 copies per cell) in 0.1% of brain cells (~170 million cells), requires ~400 million reads for 1x coverage.

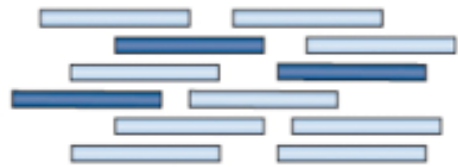
8x coverage is required for robust assembly (Jiang *et al.* 11)..

- RNA Capture Sequencing provides a potential method to detect transcripts with rare or restricted expression.

RNA CaptureSeq

CaptureSeq selects a portion of the transcriptome for focused sequencing, thereby achieving a huge increase in sequencing depth and coverage.

RNA transcripts
(cDNA)

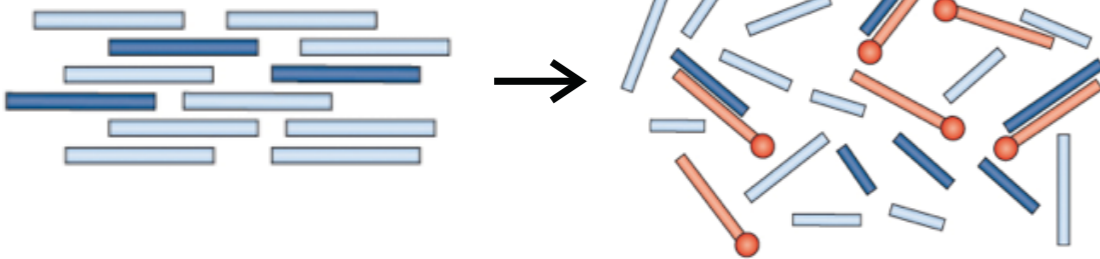


RNA CaptureSeq

CaptureSeq selects a portion of the transcriptome for focused sequencing, thereby achieving a huge increase in sequencing depth and coverage.

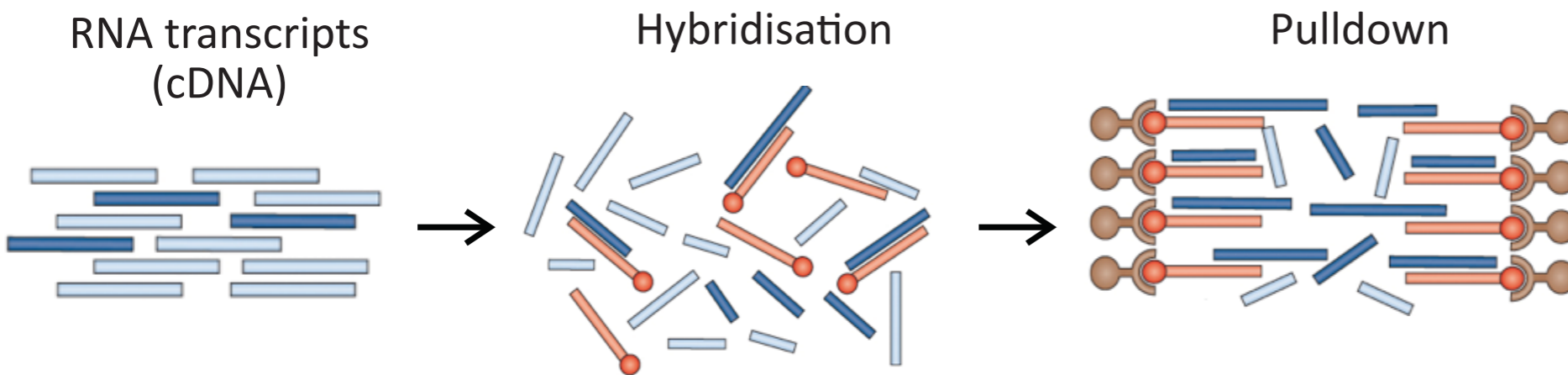
RNA transcripts
(cDNA)

Hybridisation



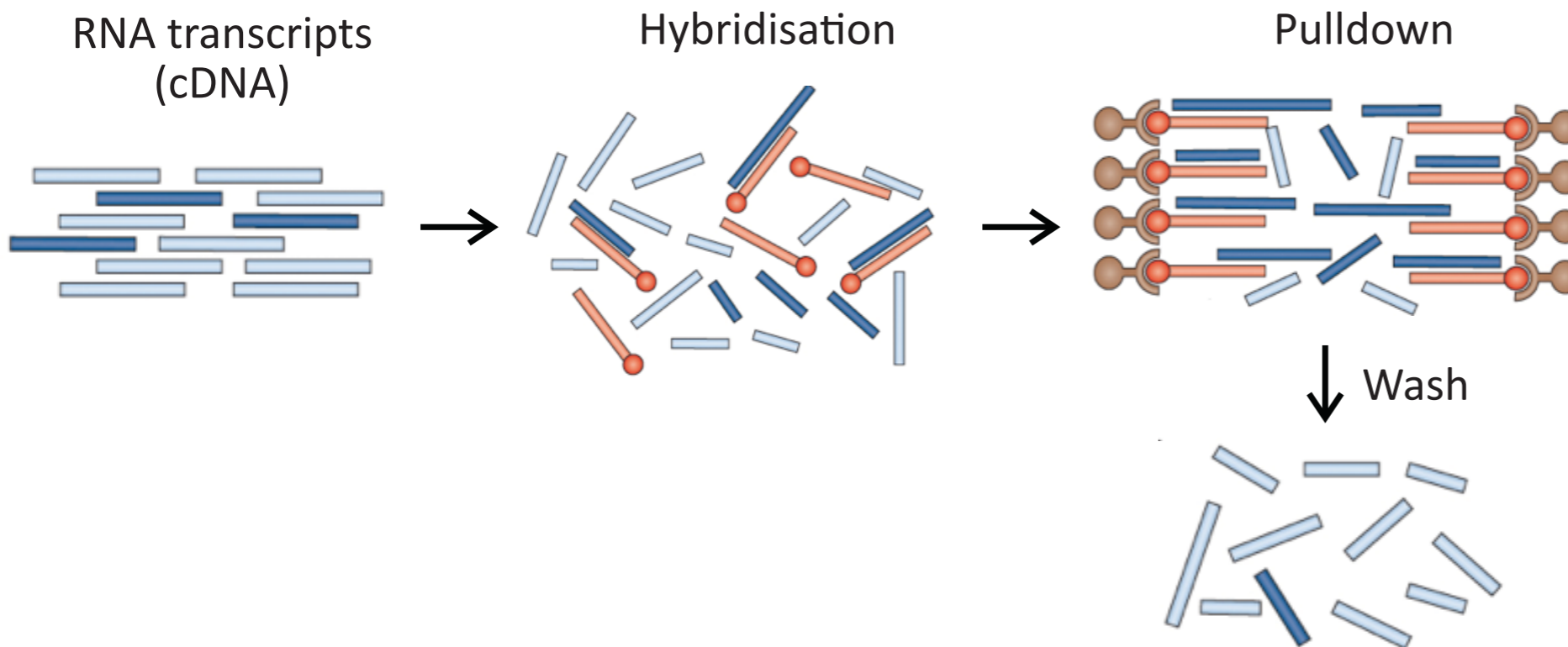
RNA CaptureSeq

CaptureSeq selects a portion of the transcriptome for focused sequencing, thereby achieving a huge increase in sequencing depth and coverage.



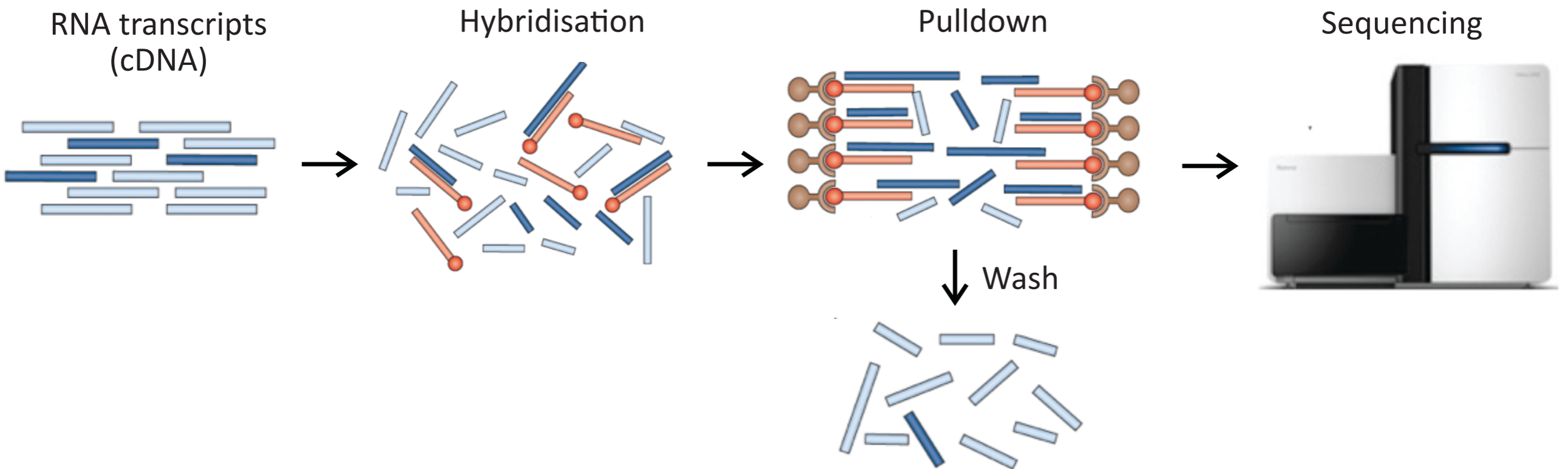
RNA CaptureSeq

CaptureSeq selects a portion of the transcriptome for focused sequencing, thereby achieving a huge increase in sequencing depth and coverage.



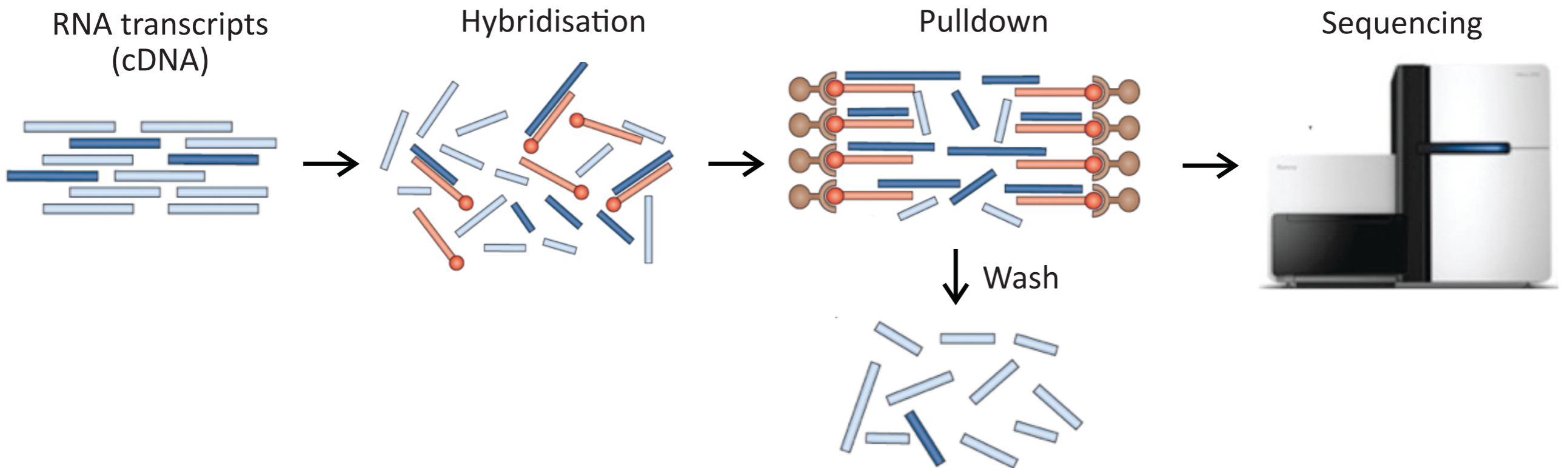
RNA CaptureSeq

CaptureSeq selects a portion of the transcriptome for focused sequencing, thereby achieving a huge increase in sequencing depth and coverage.

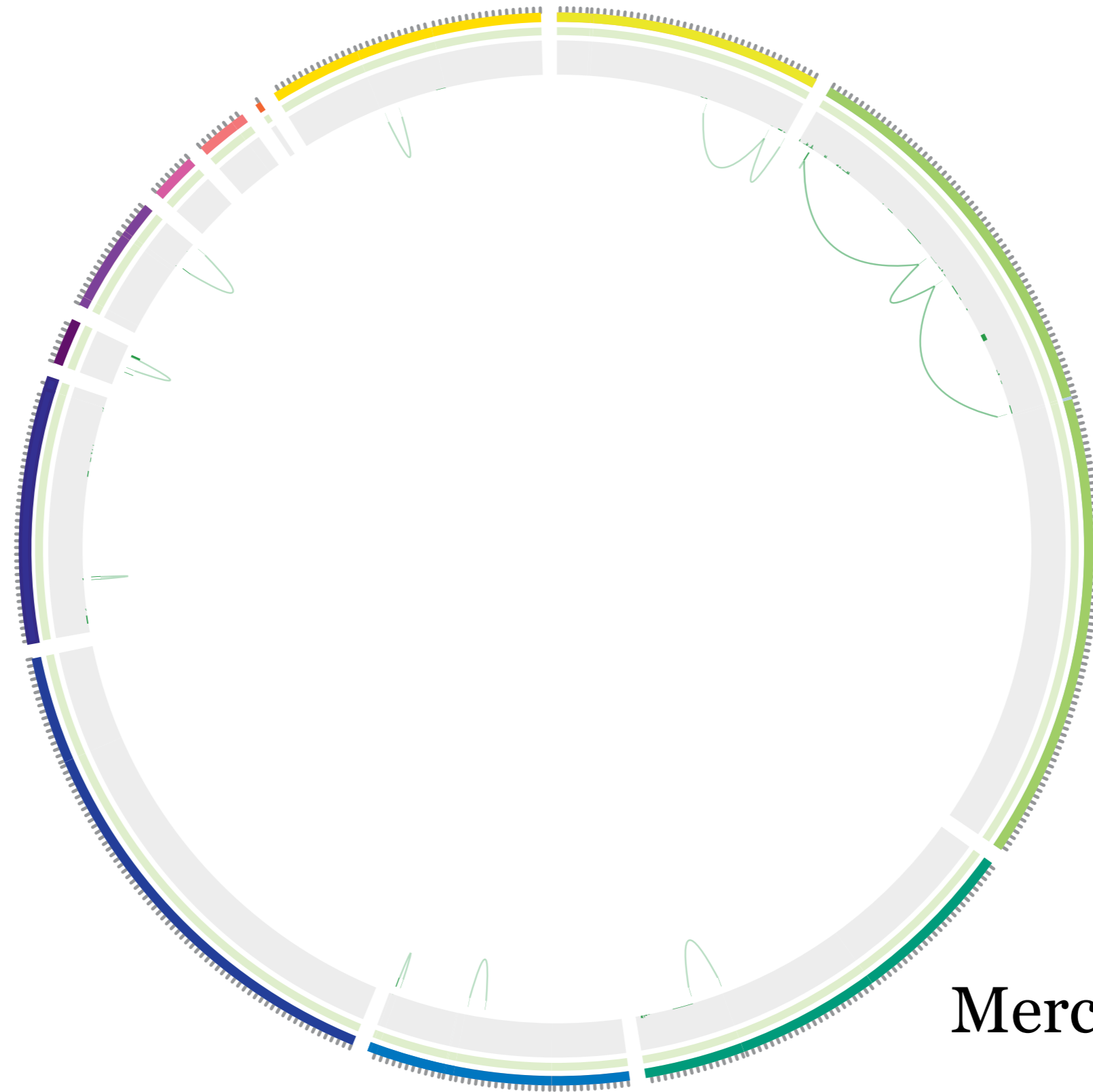


RNA CaptureSeq

This approach can characterise transcripts with rare or transient expression that is below the detection limits of conventional sequencing approaches.



RNA capture sequencing can discover new transcripts



Mercer et al. 2011

Tracks
(from outer edge):

■ Annotated gene

■ RNaseq

■ Assembled exon

— Splice junction

■ Probed region

■ RNaseq

■ Novel assembled exon

— Novel splice junction

RNA capture sequencing can discover new transcripts



Mercer et al. 2011

Tracks
(from outer edge):

■ Annotated gene

■ RNAseq

■ Assembled exon

— Splice junction

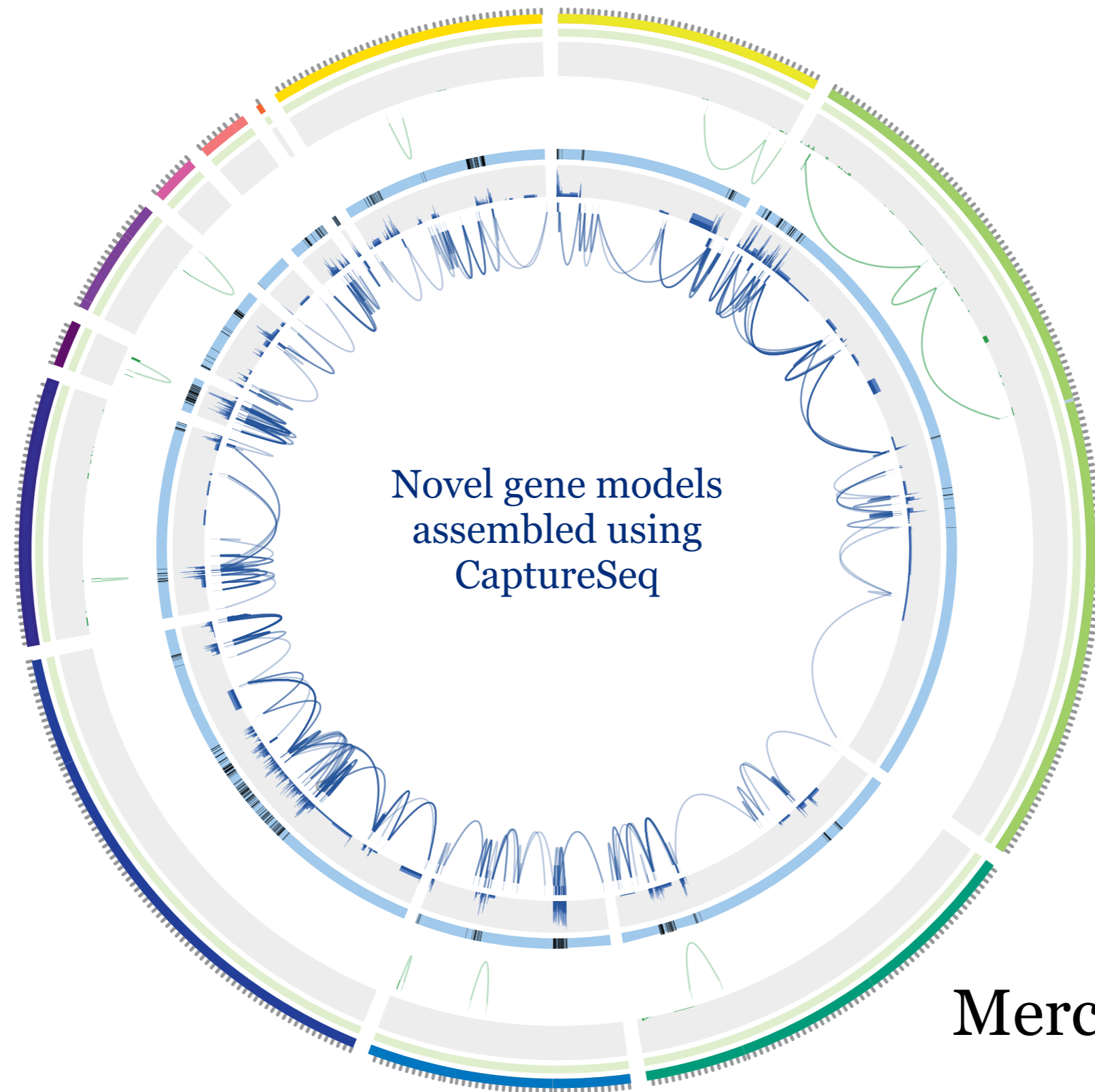
■ Probed region

■ RNAseq

■ Novel assembled exon

— Novel splice junction

RNA capture sequencing can discover new transcripts



Mercer et al. 2011

Tracks
(from outer edge):

■ Annotated gene

■ RNAseq

■ Assembled exon

— Splice junction

■ Probed region

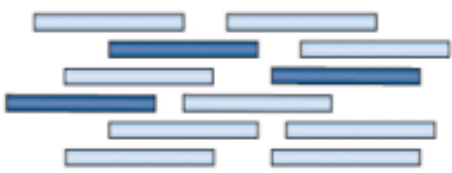
■ RNAseq

■ Novel assembled exon

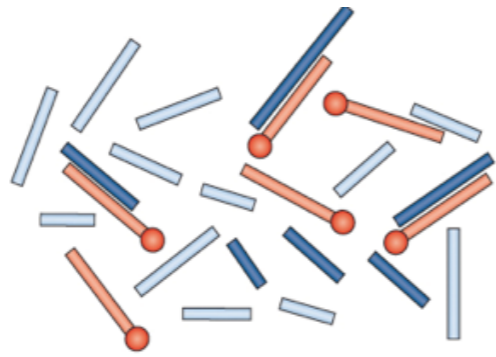
— Novel splice junction

Identifying novel genes in “intergenic” GWAS regions

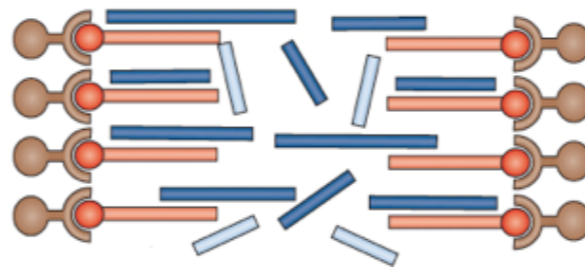
RNA transcripts
(cDNA)



Hybridisation



Pulldown



Sequencing



Identifying novel genes in “intergenic” GWAS regions

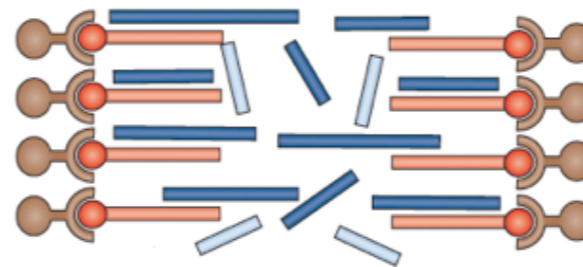
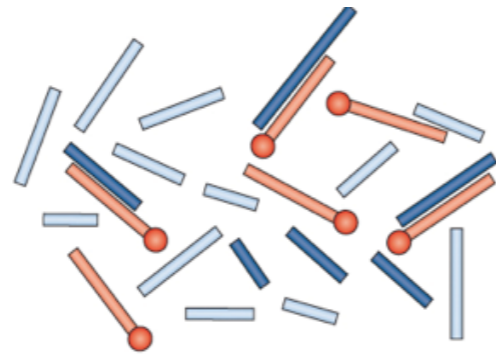
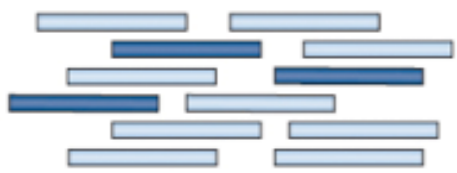
Ribo-depleted
RNA from 20
normal human
tissues.

RNA transcripts
(cDNA)

Hybridisation

Pulldown

Sequencing



Identifying novel genes in “intergenic” GWAS regions

Ribo-depleted RNA from 20 normal human tissues.

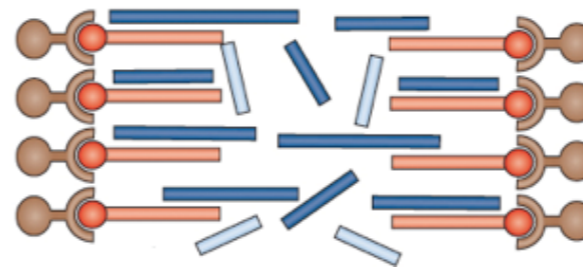
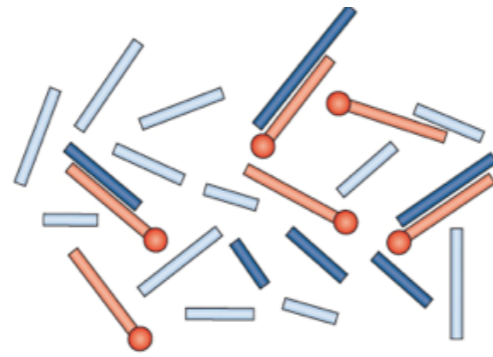
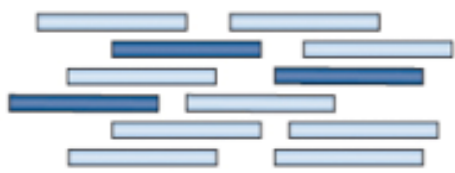
Probes to 339 regions in linkage disequilibrium with GWAS SNP that contained no genes.

RNA transcripts (cDNA)

Hybridisation

Pulldown

Sequencing



Identifying novel genes in “intergenic” GWAS regions

Ribo-depleted RNA from 20 normal human tissues.

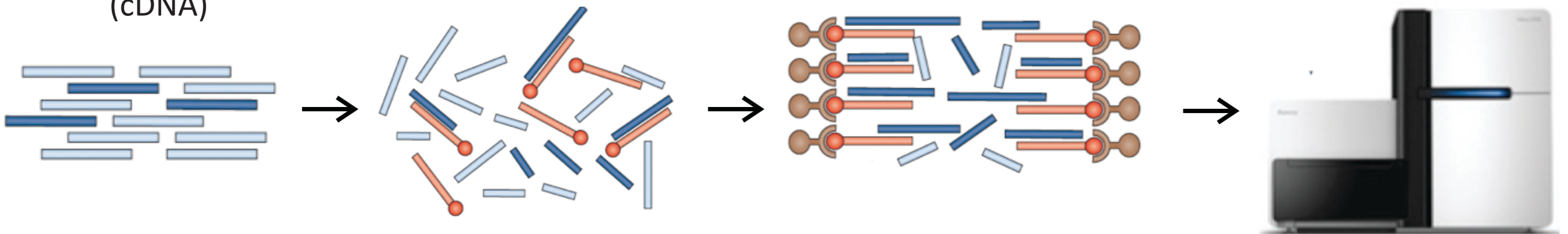
Probes to 339 regions in linkage disequilibrium with GWAS SNP that contained no genes.

RNA transcripts (cDNA)

Hybridisation

Pulldown

Sequencing



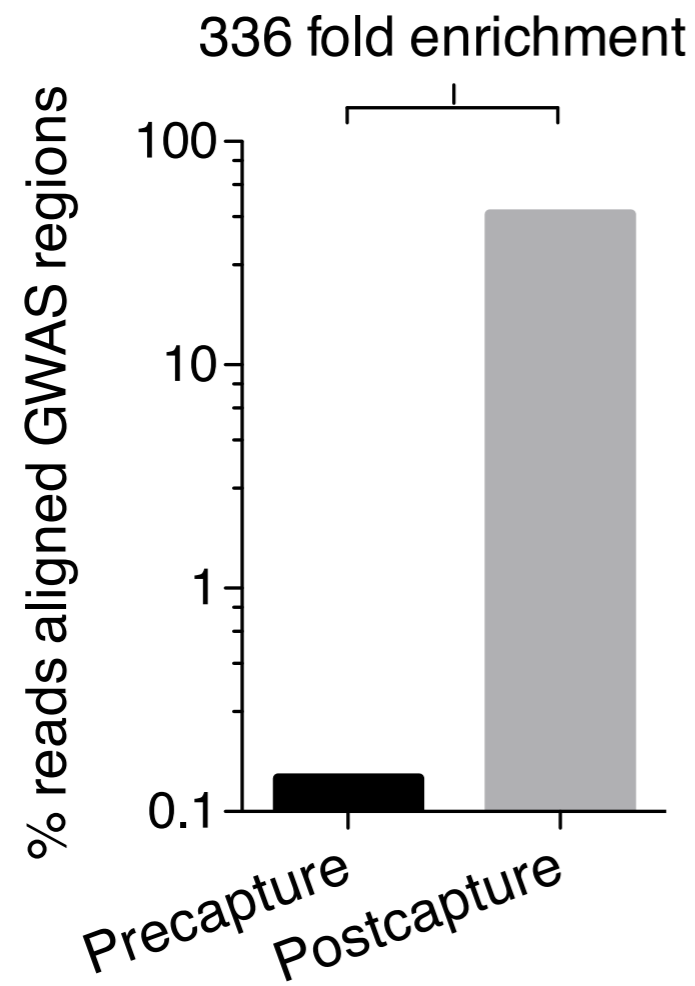
Advantage of GWAS regions over other intergenic regions is we know these regions are functionally relevant.

Targeting “intergenic” GWAS regions for ~150 traits and diseases

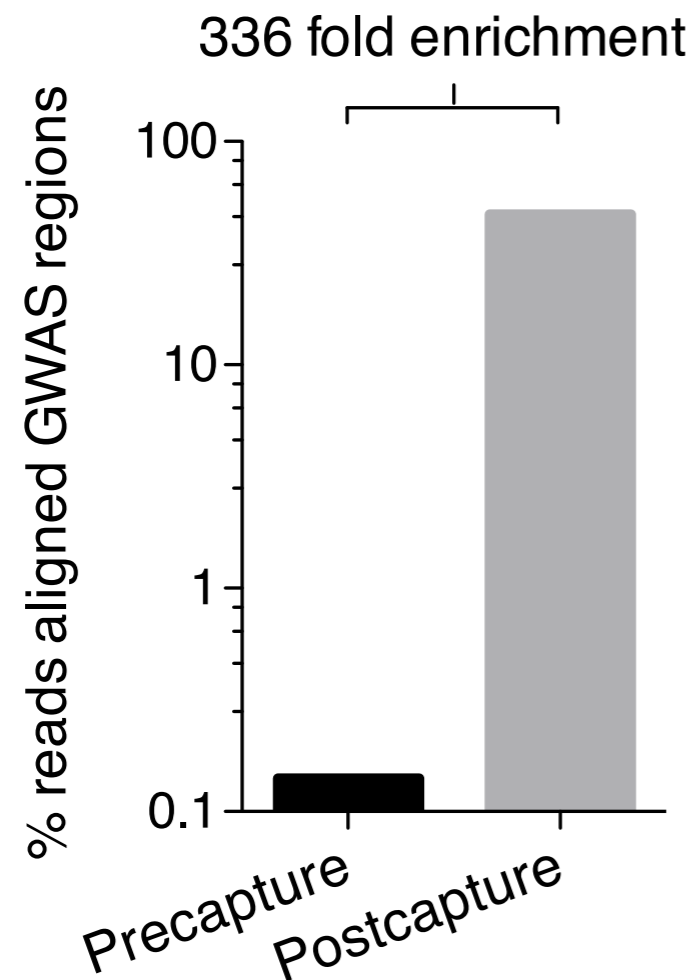
AB1-42	Chemerin levels	Iron status biomarkers	Response to antineoplastic agents
Adiponectin levels	Cholesterol, total	LDL cholesterol	Response to antipsychotic therapy
Age-related macular degeneration	Chronic lymphocytic leukemia	Lipid metabolism phenotypes	Response to antipsychotic treatment
Aging (time to event)	Cleft lip	Liver enzyme levels	Response to citalopram treatment
Aging traits	Cognitive performance	Major depressive disorder	Response to hepatitis C treatment
AIDS	Cognitive test performance	Menarche (age at onset)	Response to statin therapy
Alcoholism	Colorectal cancer	Menarche and menopause (age at onset)	Resting heart rate
Allergic rhinitis	Common traits (Other)	Metabolic syndrome	Restless legs syndrome
Alzheimer's disease	Conduct disorder	Metabolic traits	Rheumatoid arthritis
Amyotrophic lateral sclerosis	Corneal structure	Metabolite levels	RR interval (heart rate)
Ankylosing spondylitis	Coronary heart disease	Morbidity-free survival	Schizophrenia
Anthropometric traits	Crohn's disease	Multiple sclerosis	Self-rated health
Aortic root size	D-dimer levels	Myocardial infarction	Sleep duration
Aortic stiffness	Depression--quantitative trait	Nevirapine-induced rash	Smoking behavior
Asthma	Diabetes related insulin traits	Non-alcoholic fatty liver disease	Sphingolipid levels
Atopic dermatitis	Diabetic retinopathy	Nonsyndromic cleft lip with or without cleft palate	Stroke
Atrial fibrillation	Echocardiographic traits	Optic disc size (rim)	Subclinical atherosclerosis traits (other)
Atrioventricular conduction	Electrocardiographic conduction measures	Ovarian cancer	Sudden cardiac arrest
Attention deficit hyperactivity disorder	Endometriosis	Ovarian reserve	Systemic lupus erythematosus
Autism	Factor VII	P-tau181p	Systemic sclerosis
Bilirubin levels	Fasting plasma glucose	Panic disorder	T-tau
Biochemical measures	HDL cholesterol	Parkinson's disease	Tanning
Bipolar disorder	HDL Cholesterol - Triglycerides (HDLC-TG)	Peripheral artery disease	Temperament-related traits
Blood pressure	Heart failure	Permanent tooth development	Tonometry
Body mass index	Height	Personality dimensions	Triglycerides
Bone mineral density	Hematological and biochemical traits	Phospholipid levels (plasma)	Triglycerides-Blood Pressure (TG-BP)
Brain structure	Hemostatic factors and hematological phenotypes	PR interval	Type 1 diabetes
Breast cancer	Hip geometry	Premature ovarian failure	Type 2 diabetes
C-reactive protein	Hippocampal atrophy	Primary biliary cirrhosis	Ulcerative colitis
Caffeine consumption	HIV-1 control	Primary tooth development	Urate levels
Calcium levels	Hoarding	Prostate cancer	Urinary albumin excretion
Cannabis dependence	Hyperactive-impulsive symptoms	Psoriasis	Vaccine-related adverse events
Cardiac hypertrophy	Hypertension	Pulmonary function	Waist circumference
Cardiac repolarization	IgE grass sensitization	QT interval	Waist-hip ratio
Cardiovascular disease risk factors	Immunoglobulin A	Quantitative traits	Weight
Carotid atherosclerosis in HIV infection	Inattentive symptoms	Reasoning	White blood cell types
Carotid intima media thickness	Inflammatory bowel disease	Renal cell carcinoma	Working memory
Celiac disease	Information processing speed	Response to antidepressants	

339 intergenic GWAS linkage blocks, covering a total of ~50 Mbp of the genome (~1.7%).

High global enrichment for transcription from captured GWAS regions

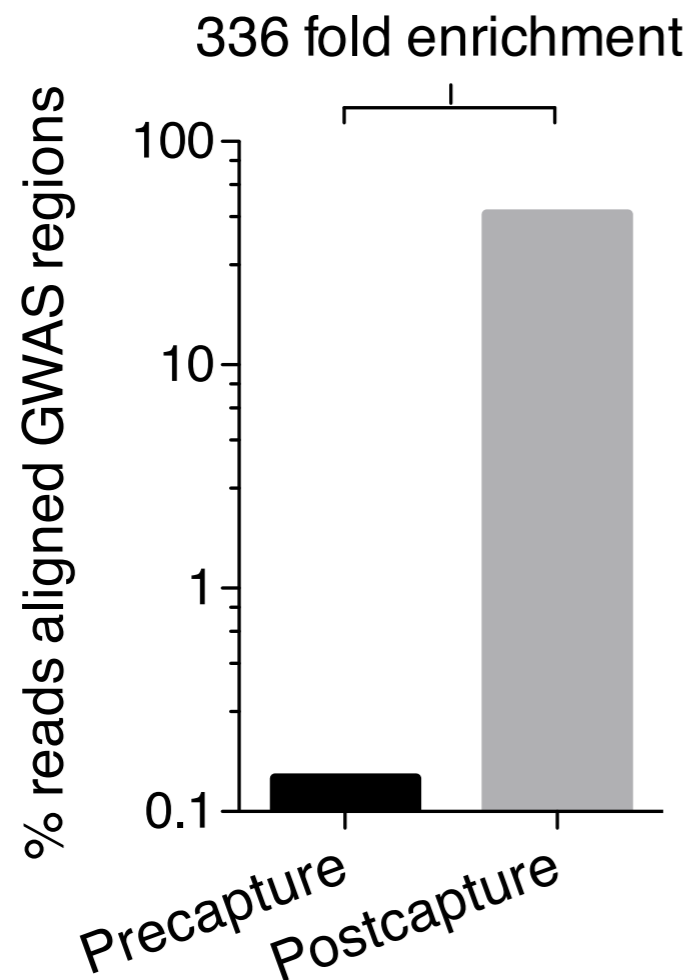


High global enrichment for transcription from captured GWAS regions



Comparing the percentage of mapped reads from targeted regions pre- and post-capture

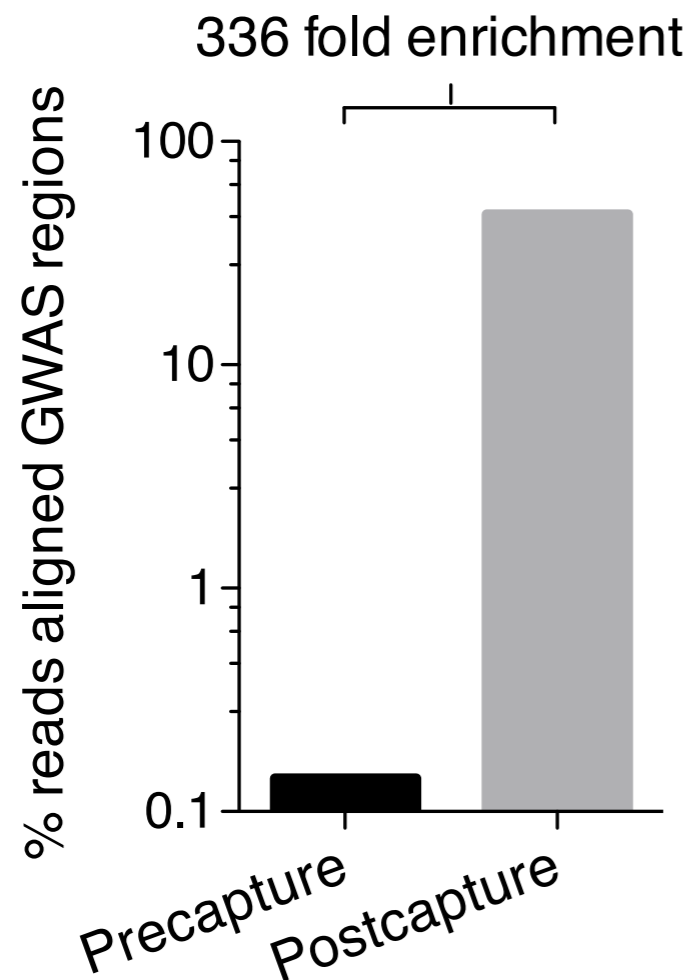
High global enrichment for transcription from captured GWAS regions



Comparing the percentage of mapped reads from targeted regions pre- and post-capture

~0.2% of the RNA in pre-capture library now comprises ~60% in the post-capture library.

High global enrichment for transcription from captured GWAS regions

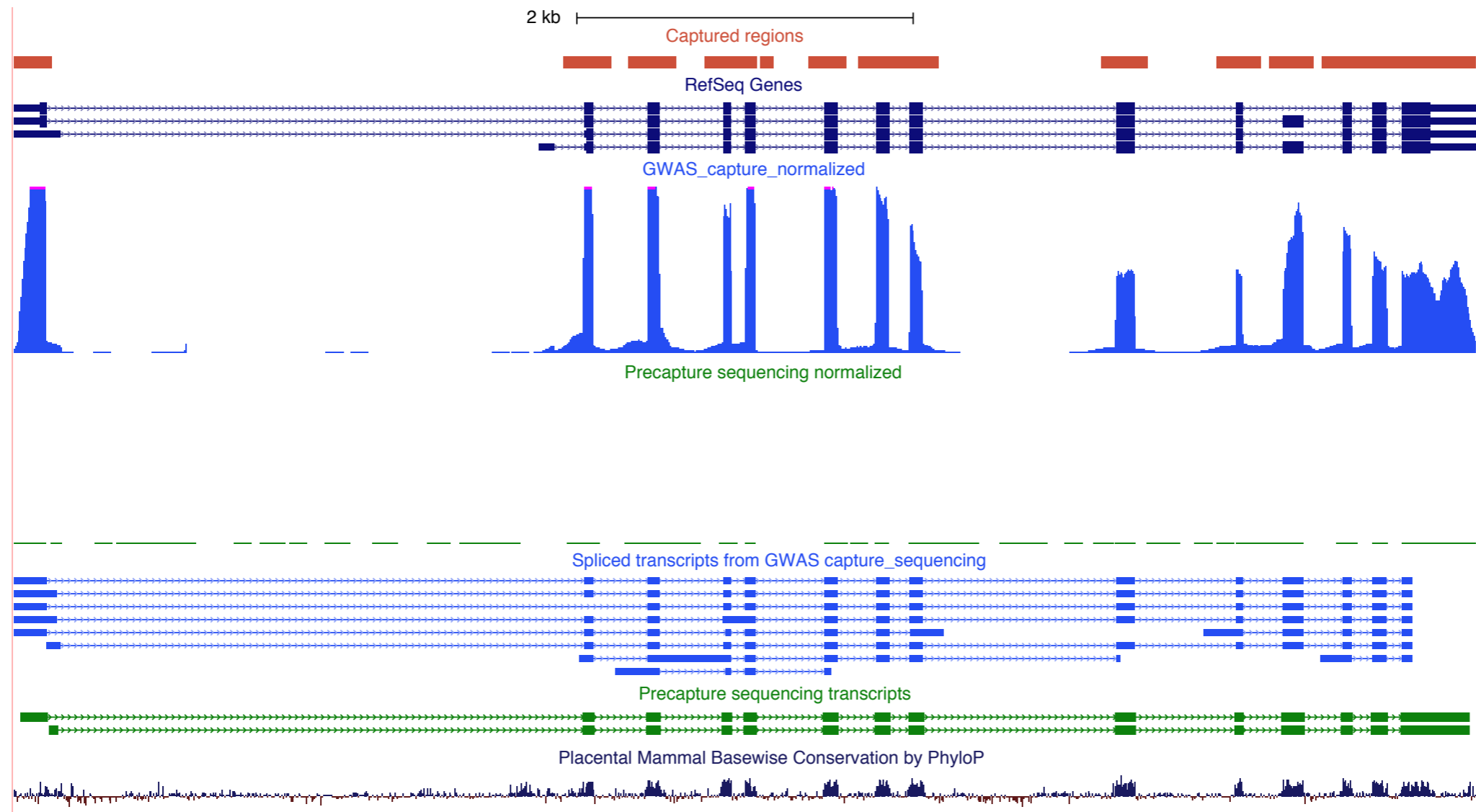


Comparing the percentage of mapped reads from targeted regions pre- and post-capture

~0.2% of the RNA in pre-capture library now comprises ~60% in the post-capture library.

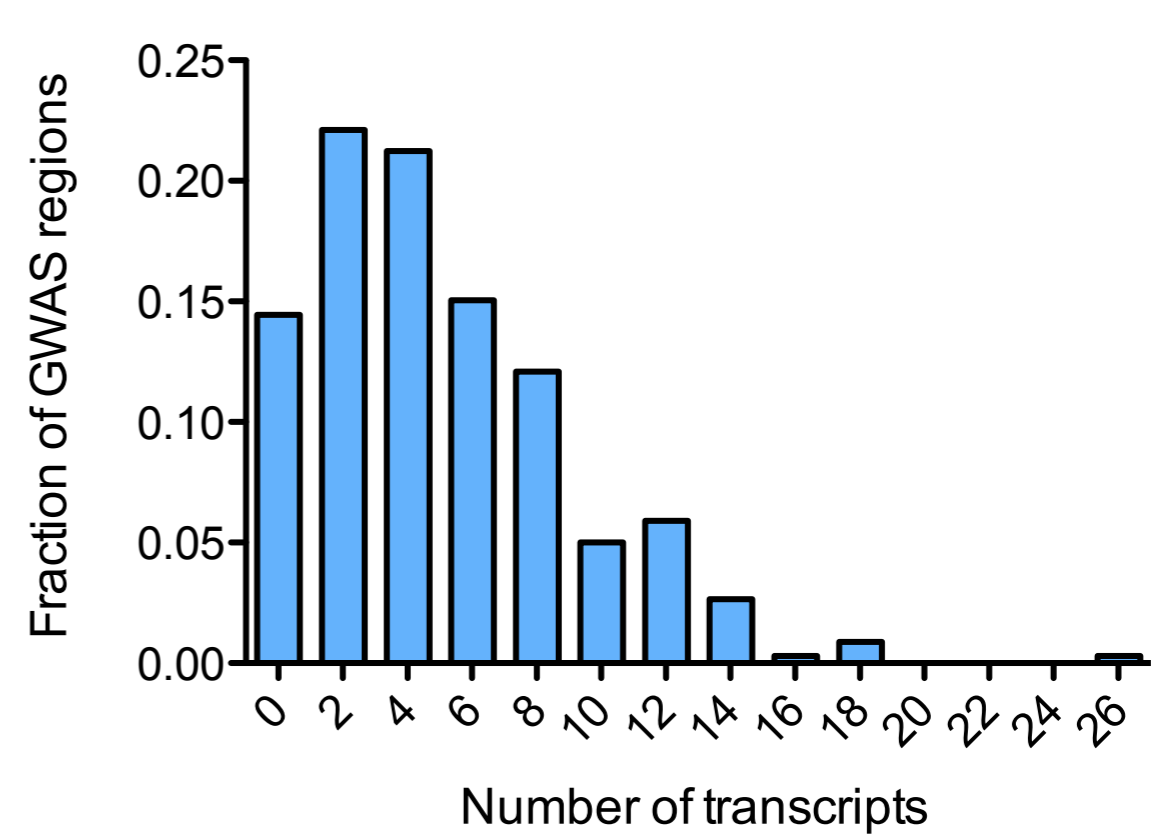
336 fold enrichment from CaptureSeq

Sequencing based enrichment: HMBS control loci



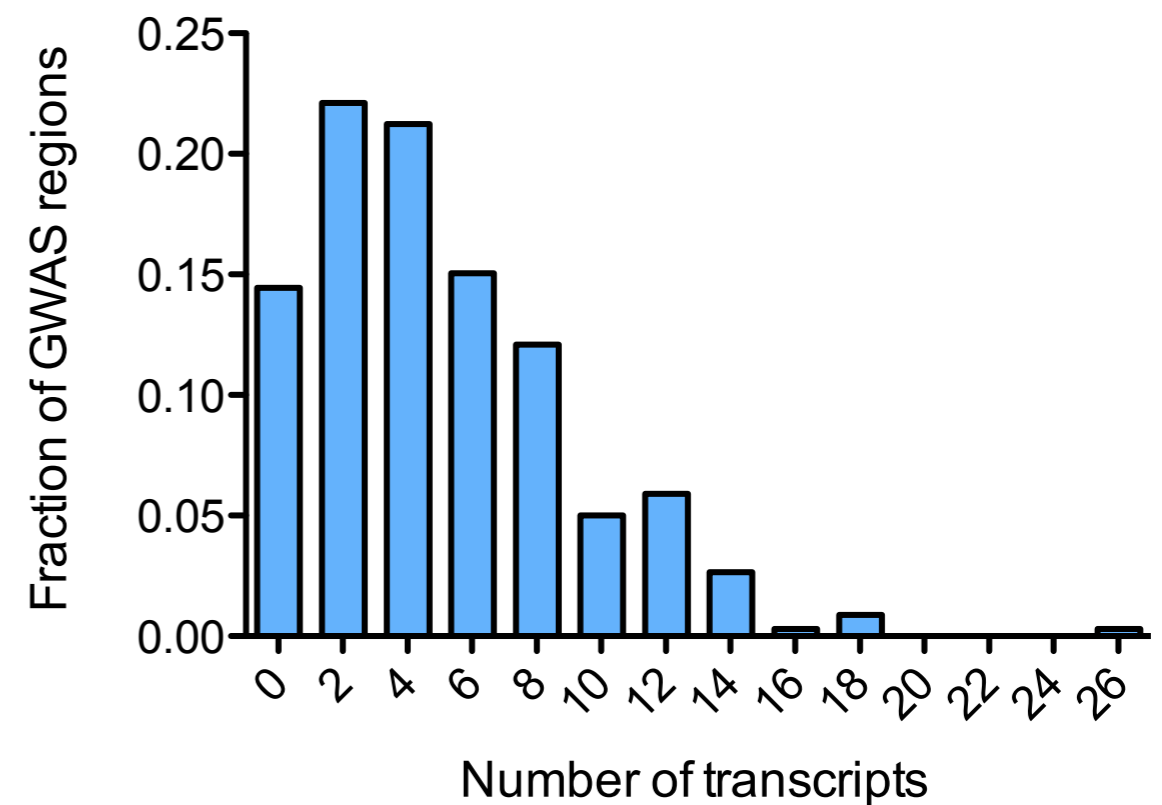
- HMBS is a control loci that we capture
- Sequencing confirms large enrichment measured globally and by qPCR
- Assemble realistic transcripts from Capture sequencing

Extensive transcription within intergenic GWAS regions



Extensive transcription within intergenic GWAS regions

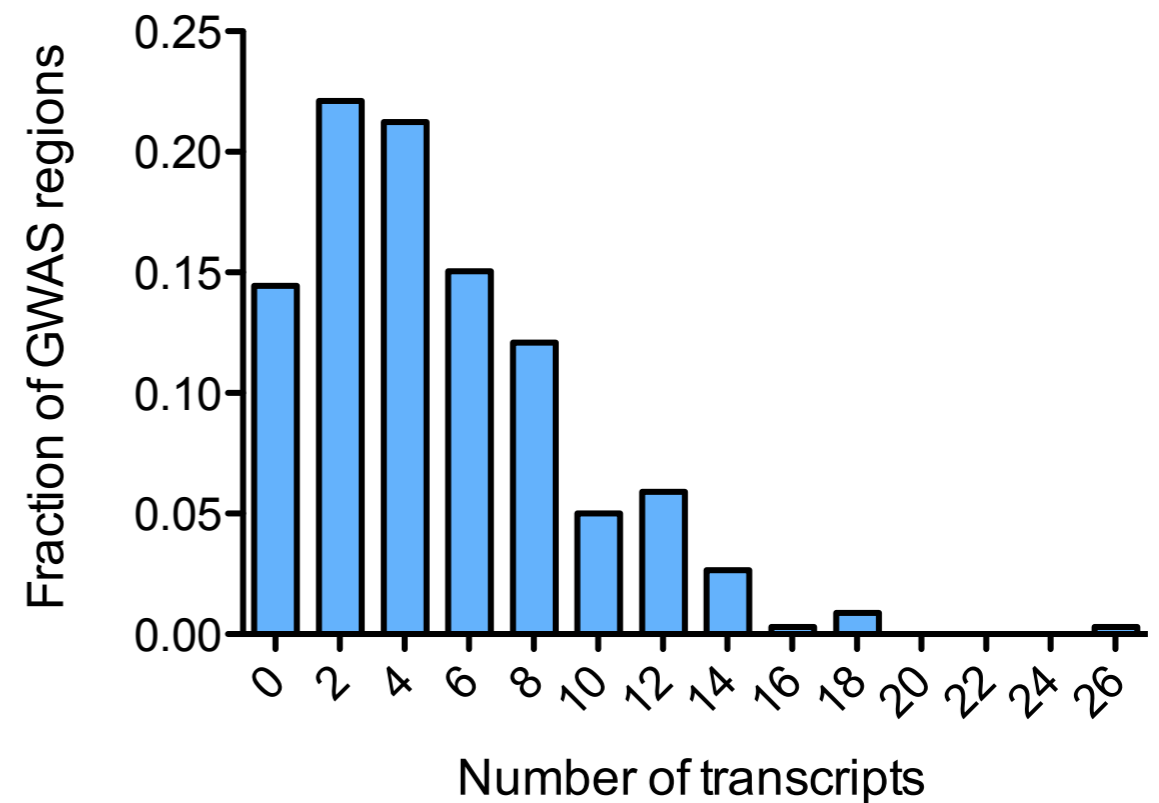
85% of 339 captured GWAS regions contained spliced transcripts.



Extensive transcription within intergenic GWAS regions

85% of 339 captured GWAS regions
contained spliced transcripts.

Median of 4 spliced transcripts per locus

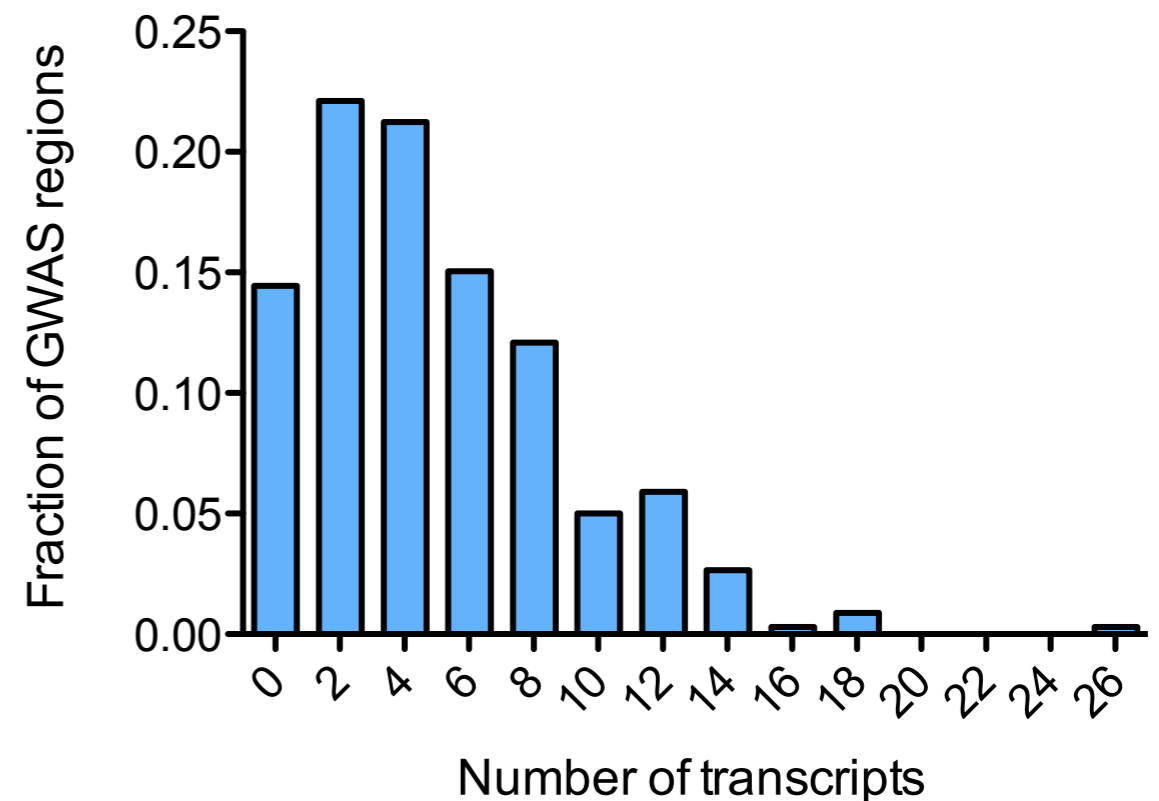


Extensive transcription within intergenic GWAS regions

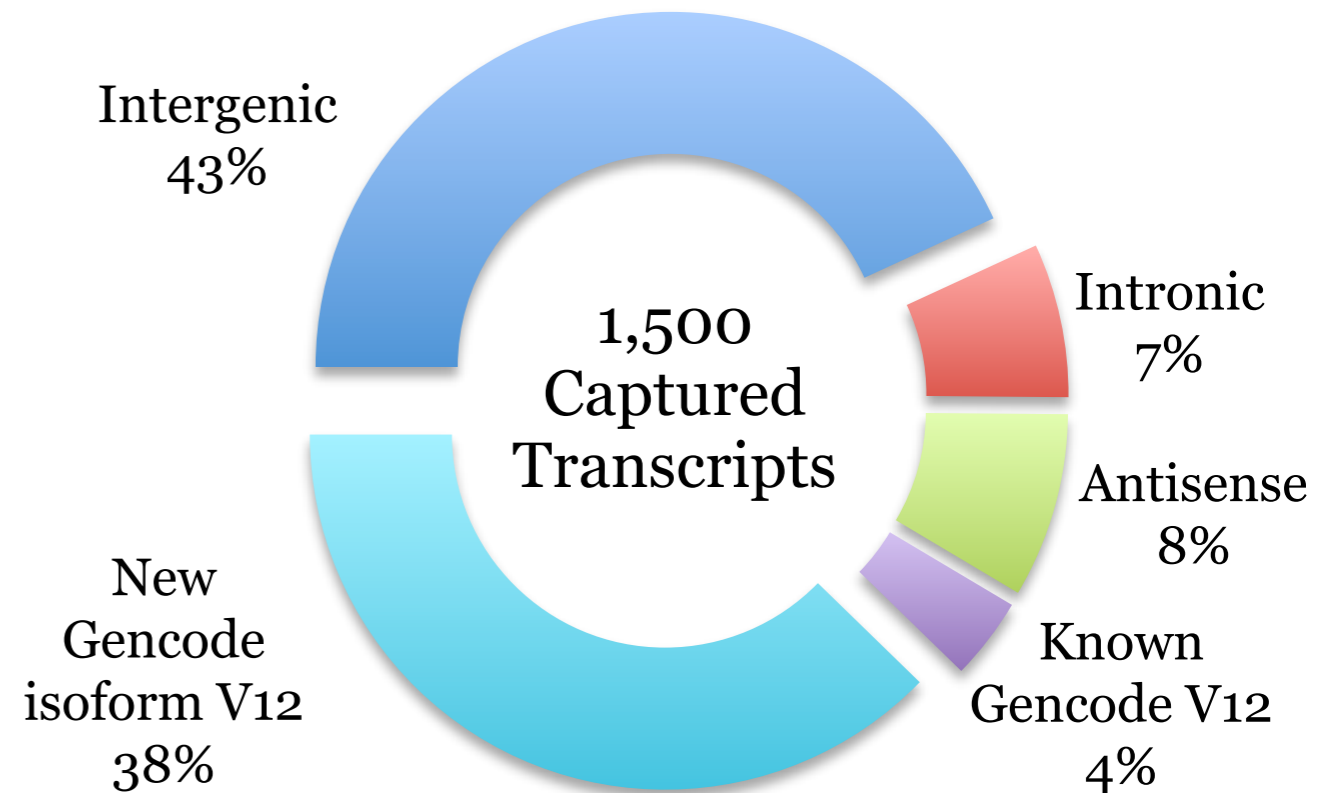
85% of 339 captured GWAS regions
contained spliced transcripts.

Median of 4 spliced transcripts per locus

Spliced transcripts covered a total of 55.6
Mb with many transcripts transcribed
into, or out of, the captured regions.

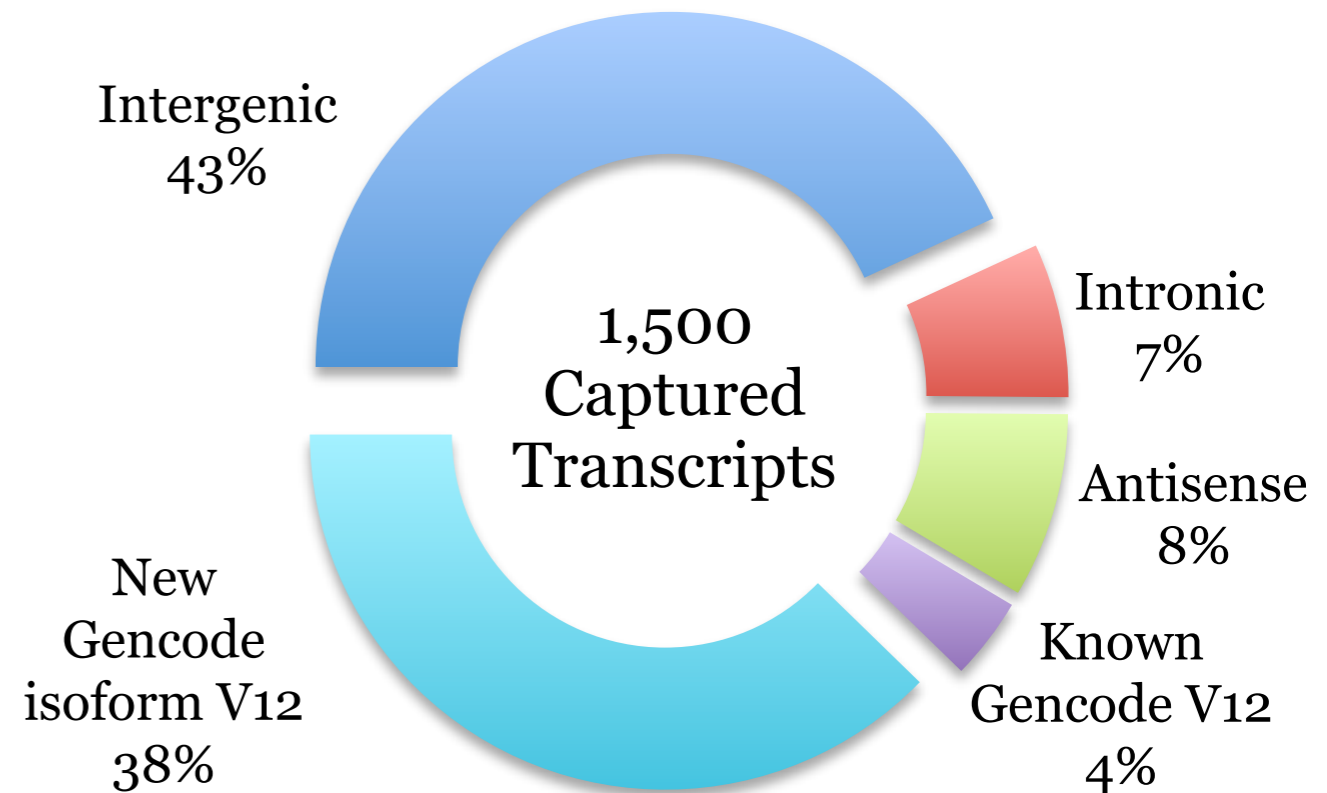


1500 transcripts identified in “intergenic” GWAS regions



1500 transcripts identified in “intergenic” GWAS regions

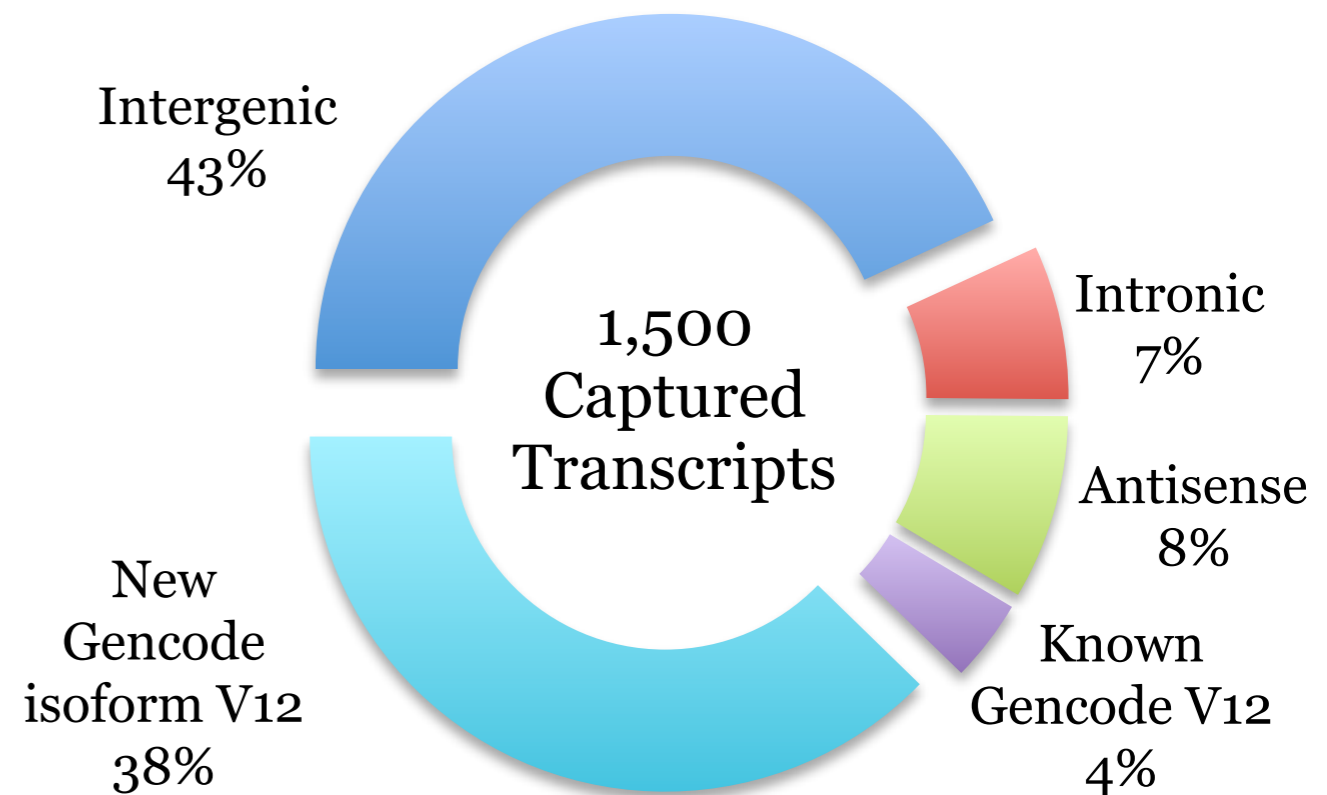
1,500 transcripts were identified (compared to 110 transcript in pre-capture).



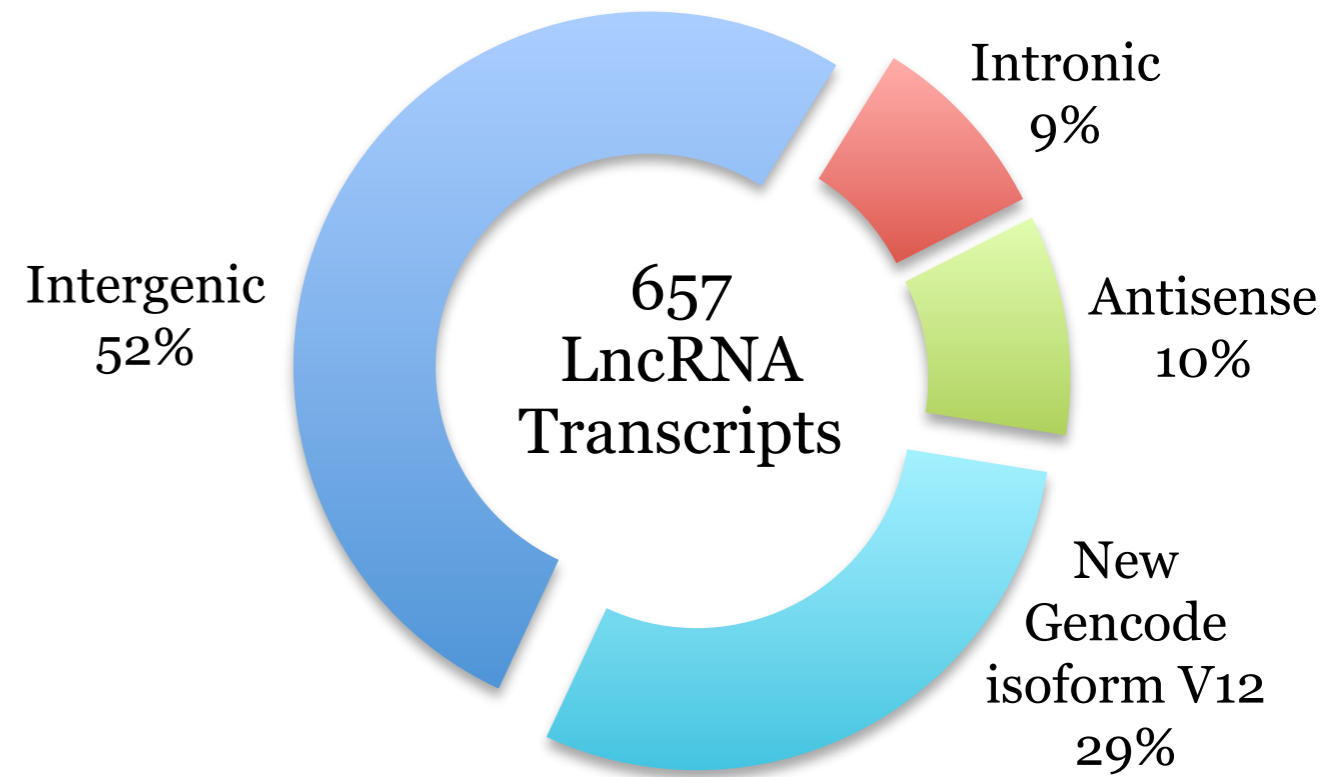
1500 transcripts identified in “intergenic” GWAS regions

1,500 transcripts were identified (compared to 110 transcript in pre-capture).

Majority of captured transcripts are entirely novel or contain novel exons of distal genes.

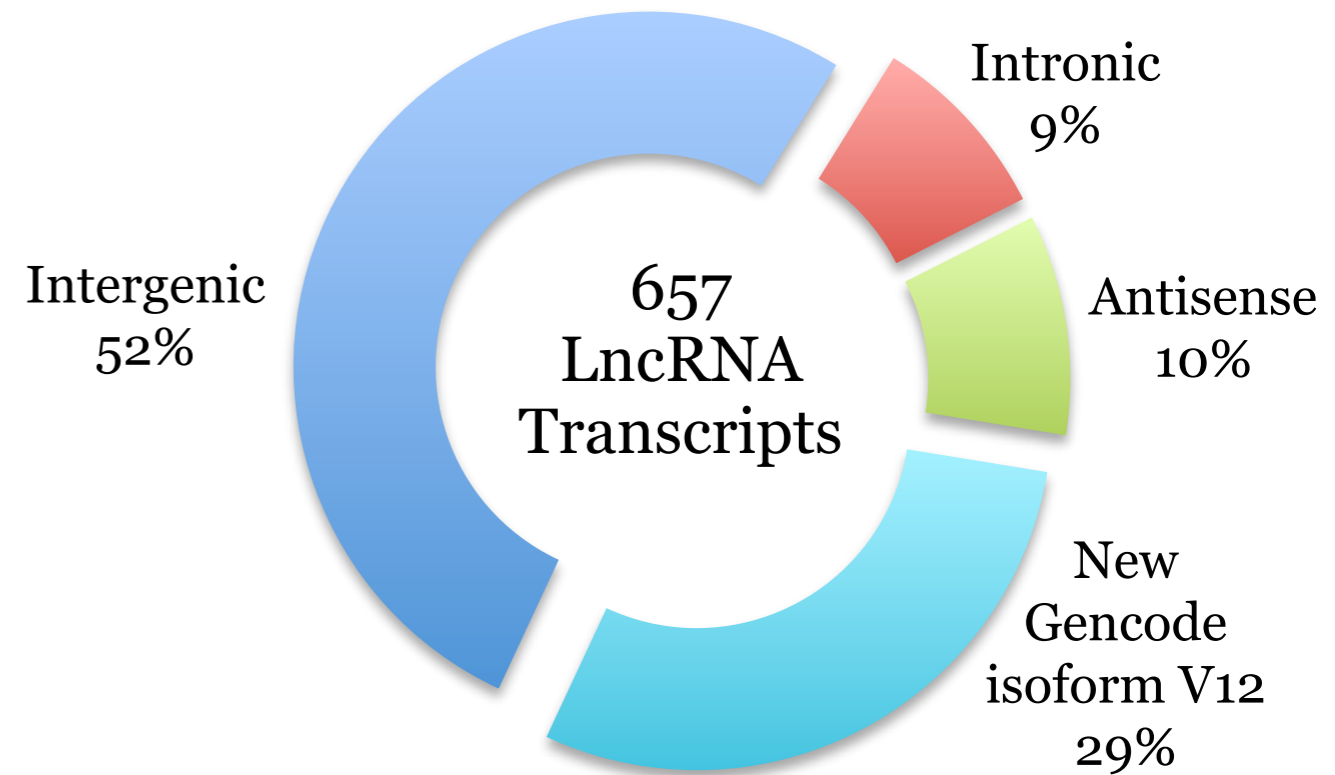


650 novel lncRNAs identified in “intergenic” GWAS regions



650 novel lncRNAs identified in “intergenic” GWAS regions

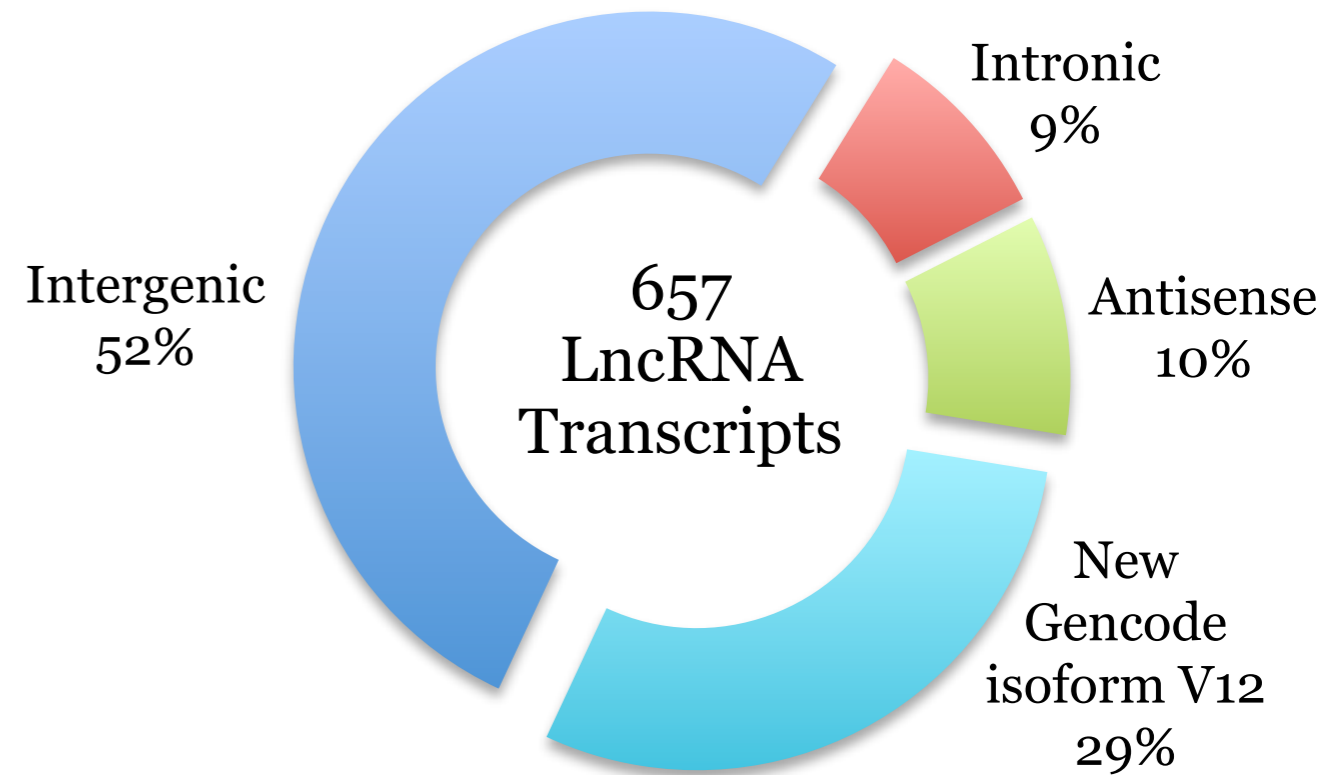
Stringent filtering identifies 657 new lncRNAs in 369 loci



650 novel lncRNAs identified in “intergenic” GWAS regions

Stringent filtering identifies 657 new lncRNAs in 369 loci

Vast majority are completely novel

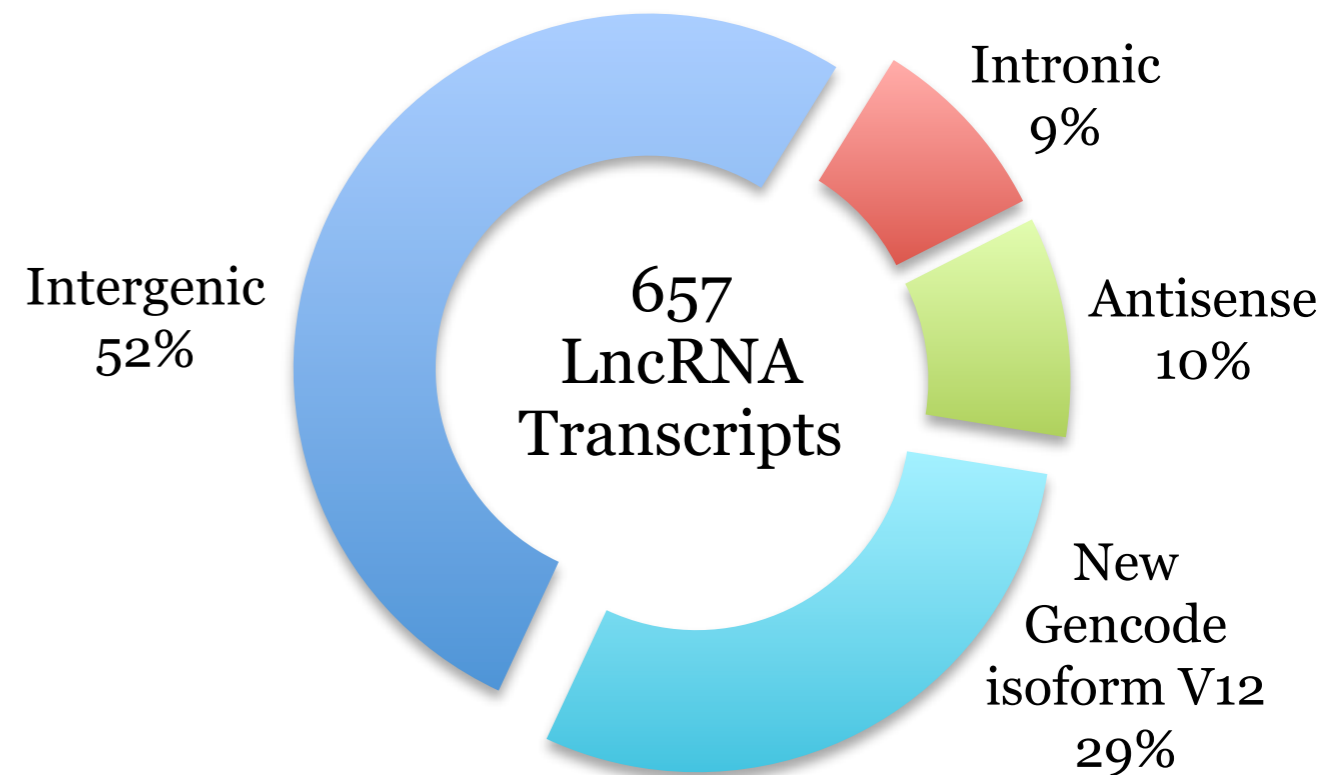


650 novel lncRNAs identified in “intergenic” GWAS regions

Stringent filtering identifies 657 new lncRNAs in 369 loci

Vast majority are completely novel

Remainder of ~1500 transcripts are novel coding gene isoforms and transcripts of uncertain coding potential.



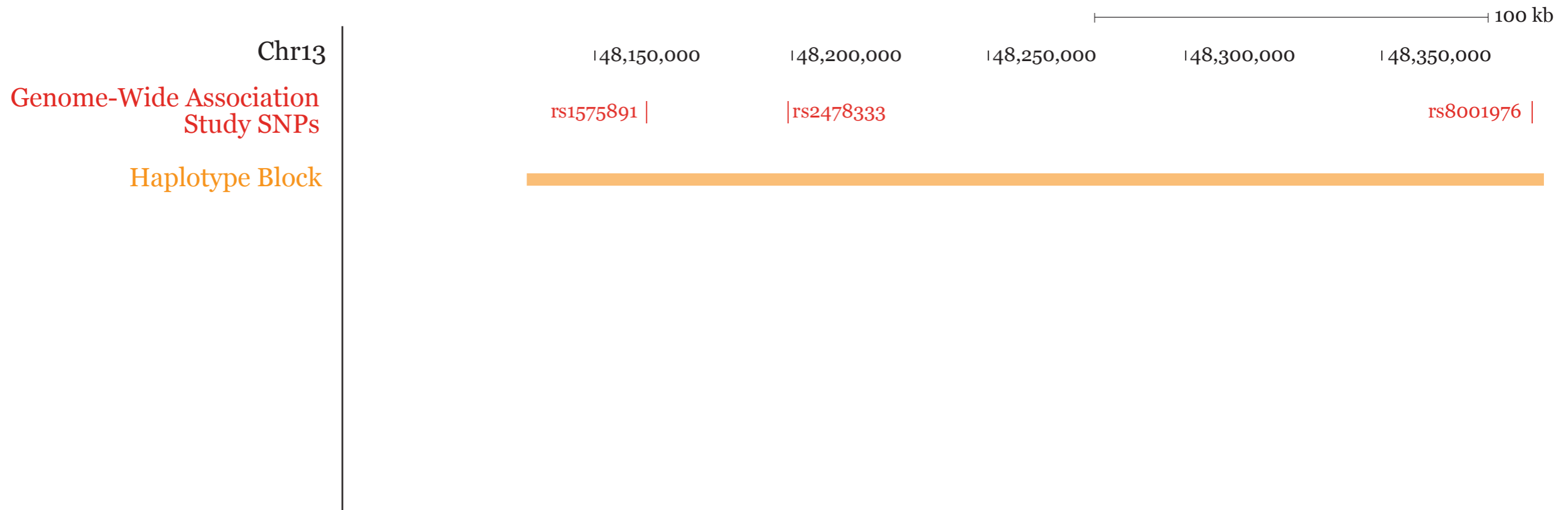
Identifying novel human genes with capture seq

Human loci associated with heart QT length, but doesn't contain any known genes.



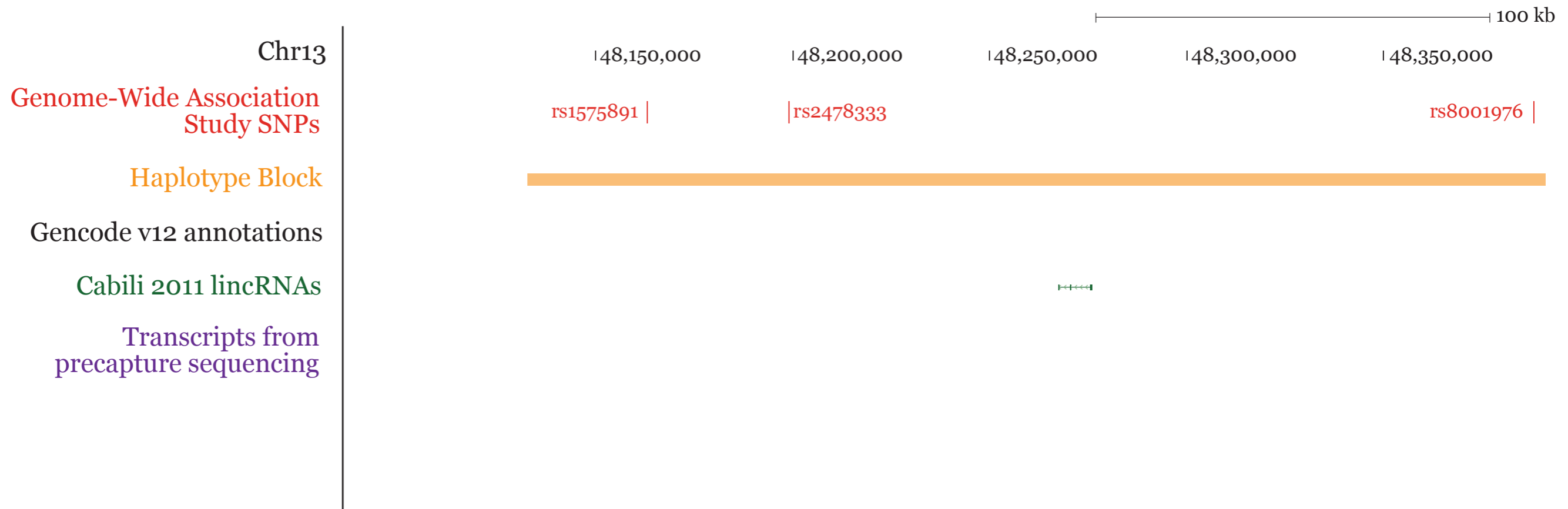
Identifying novel human genes with capture seq

Human loci associated with heart QT length, but doesn't contain any known genes.



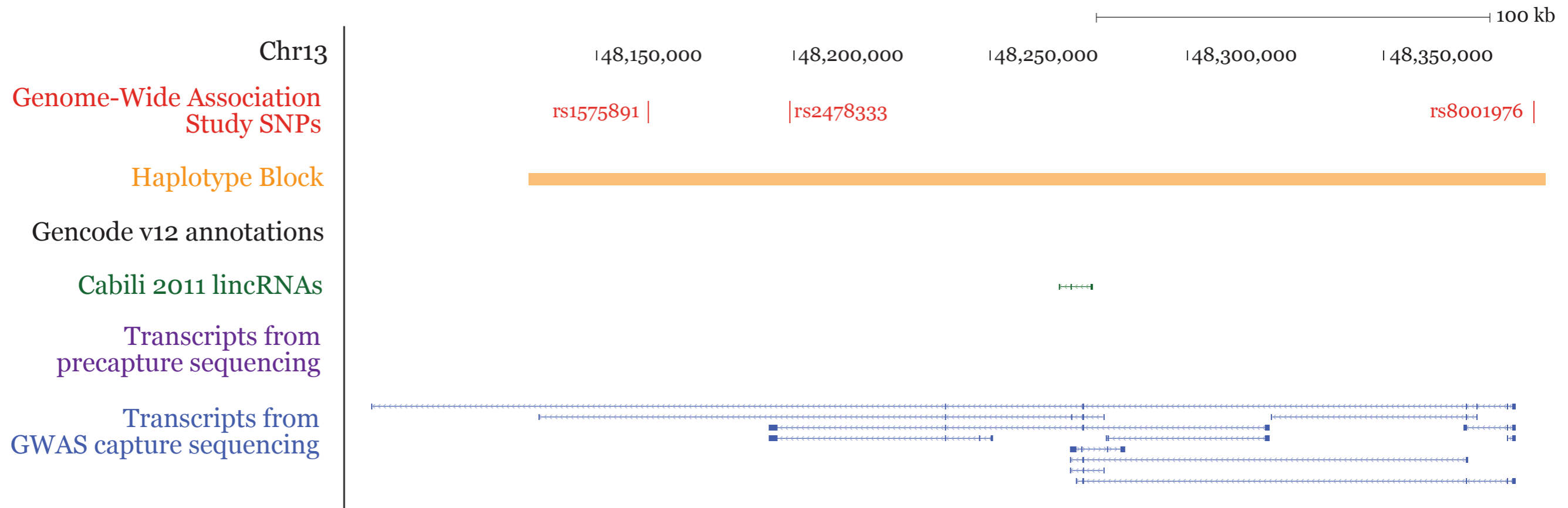
Identifying novel human genes with capture seq

Human loci associated with heart QT length, but doesn't contain any known genes.



Identifying novel human genes with capture seq

Human loci associated with heart QT length, but doesn't contain any known genes.



Capture sequencing finds many new transcripts, these can now be functionally tested.

Identifying novel human genes with capture seq

GWAS regions contains novel 5' or 3' extensions of known genes.



200 kb

Identifying novel human genes with capture seq

GWAS regions contains novel 5' or 3' extensions of known genes.

Genome-Wide Association
Study SNPs

Haplotype Block

rs16920624 |

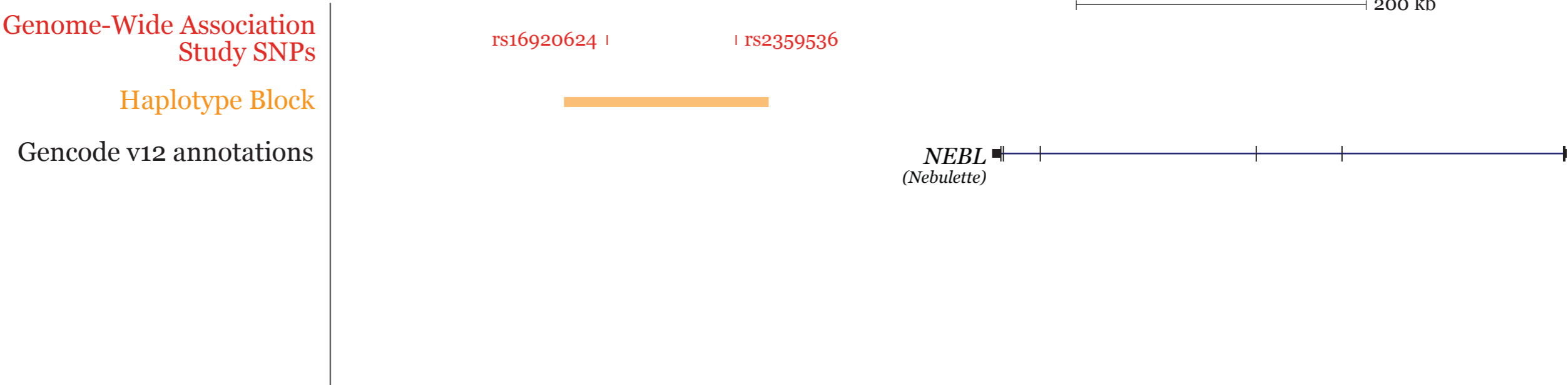
| rs2359536



200 kb

Identifying novel human genes with capture seq

GWAS regions contains novel 5' or 3' extensions of known genes.



Identifying novel human genes with capture seq

GWAS regions contains novel 5' or 3' extensions of known genes.

Genome-Wide Association
Study SNPs

Haplotype Block

Gencode v12 annotations

Transcripts from
precapture sequencing

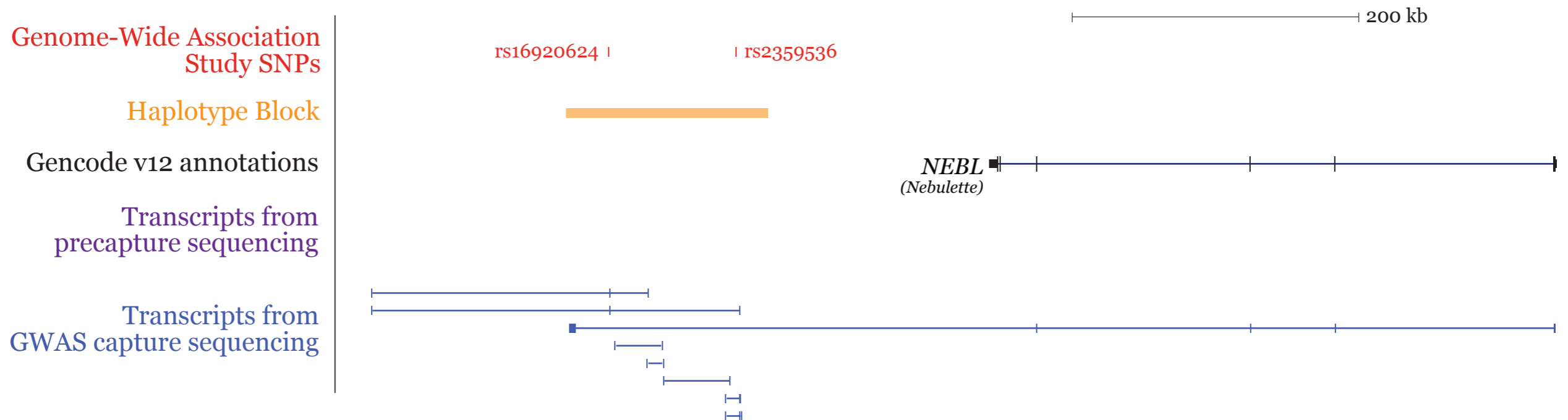
rs16920624 | | rs2359536



200 kb

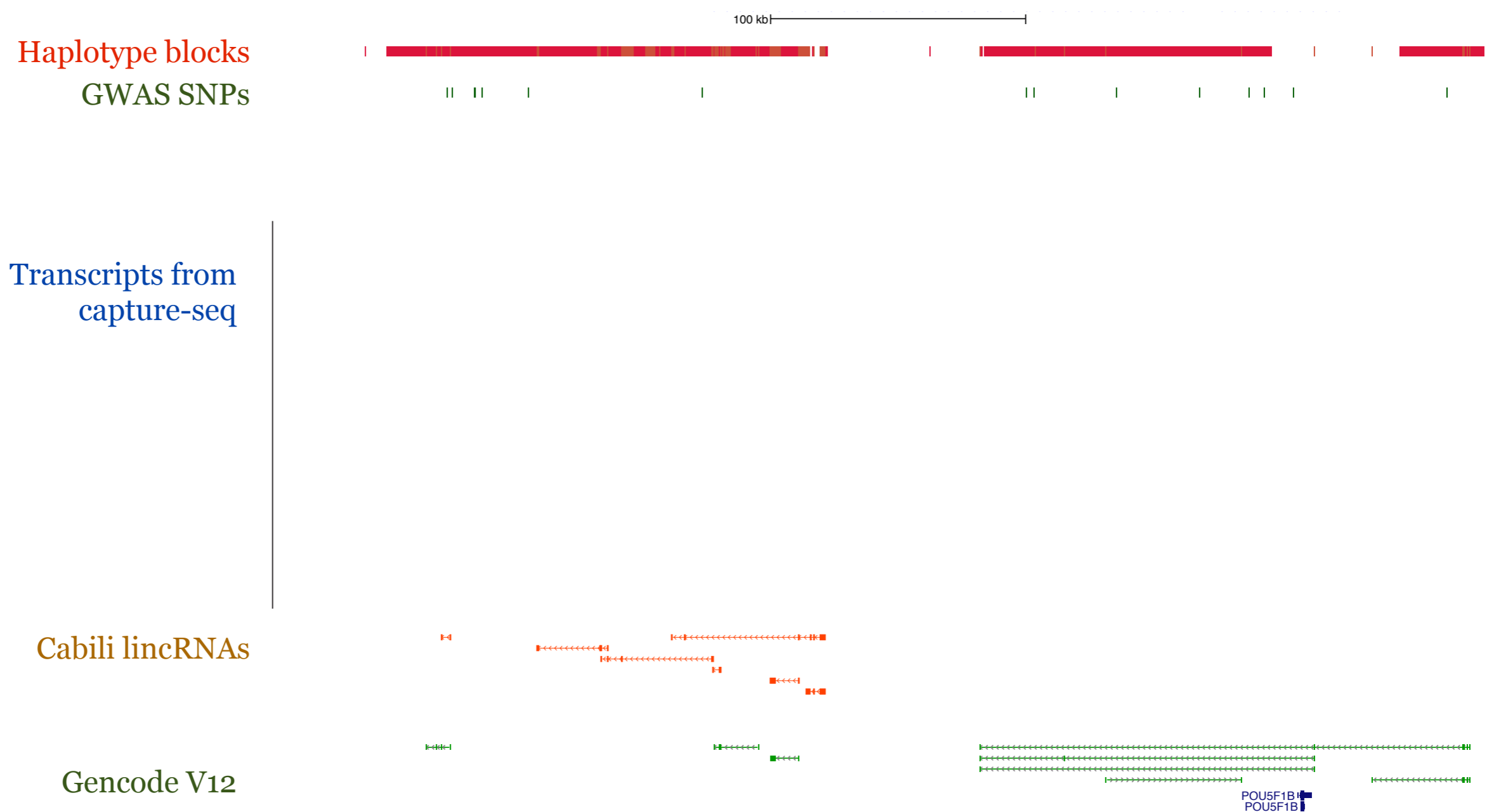
Identifying novel human genes with capture seq

GWAS regions contains novel 5' or 3' extensions of known genes.



- Captured GWAS region contains a novel isoform of the coding gene Nebulette (NEBL).
- A very recently annotated putative miRNA is present in the GWAS region on +ve strand. Likely we have also identified its host transcript(s).

Capture seq resolves fragmented annotations



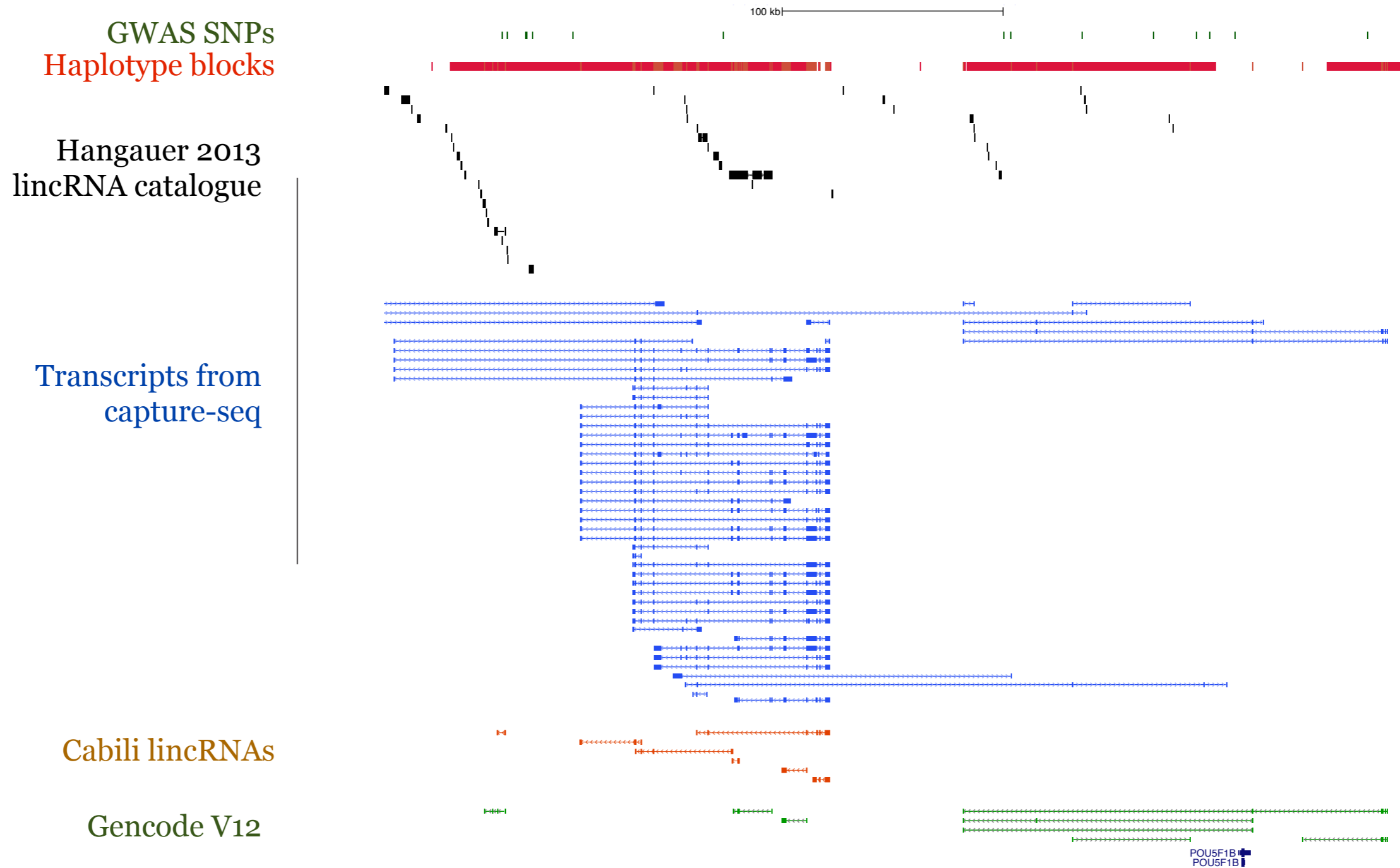
- Region of chr8 associated with prostate cancer by a number of studies.

Capture seq resolves fragmented annotations

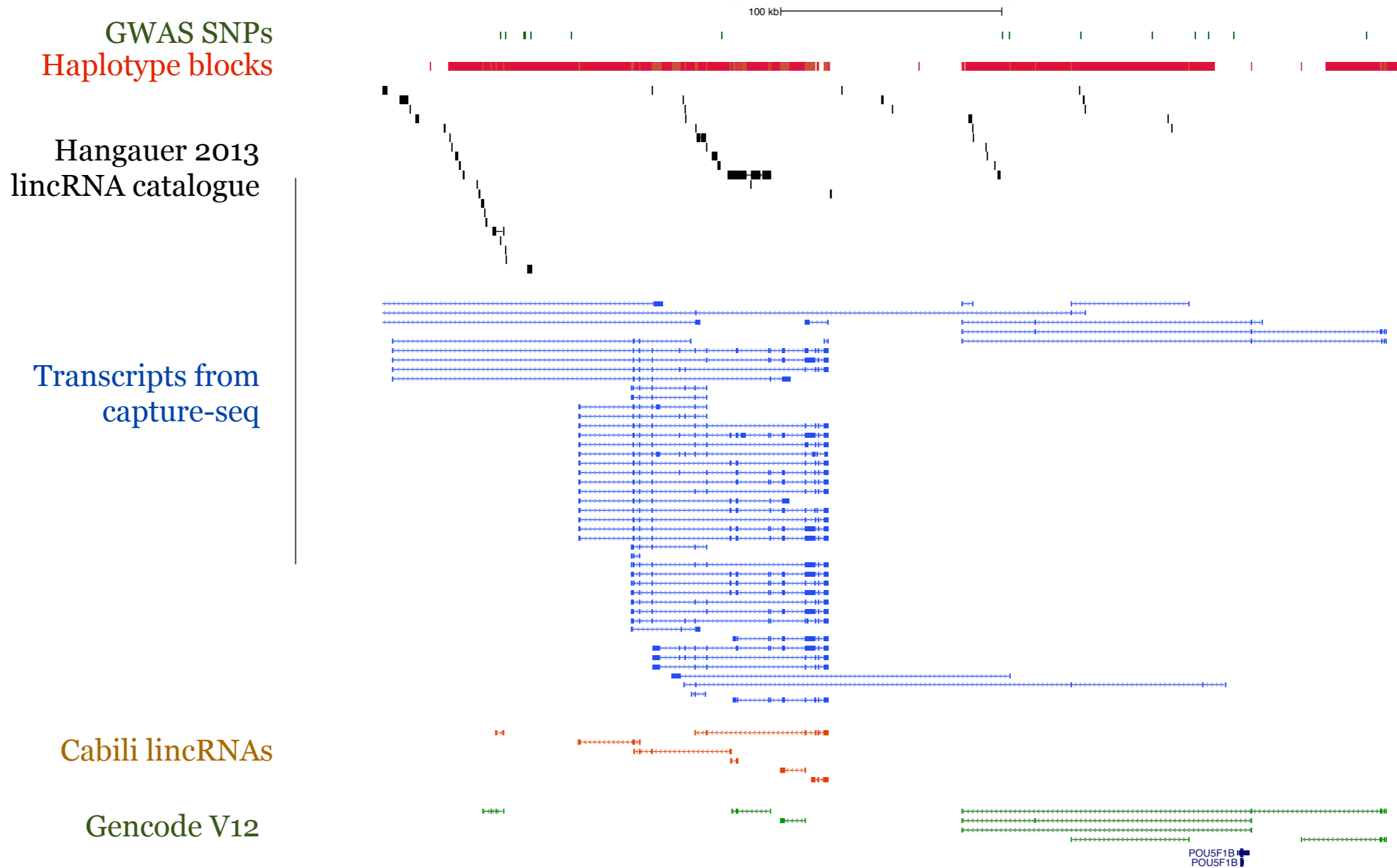


- Region of chr8 associated with prostate cancer by a number of studies.
- Capture sequencing combines previously annotated lincRNA fragments in complex alternatively spliced loci

Capture-seq vs Hangauer 2013 lincRNA catalogue



Capture-seq vs Hangauer 2013 lincRNA catalogue



Capture-seq outperforms other novel gene catalogues

Summary: Intergenic GWAS capture sequencing

Summary: Intergenic GWAS capture sequencing

- Capture-seq allows unprecedented depth of coverage and enriches transcripts standard RNA-seq cannot identify.

Summary: Intergenic GWAS capture sequencing

- Capture-seq allows unprecedented depth of coverage and enriches transcripts standard RNA-seq cannot identify.
- We find extensive transcription within and across “intergenic” GWAS regions, identifying ~1500 mostly novel transcripts.

Summary: Intergenic GWAS capture sequencing

- Capture-seq allows unprecedented depth of coverage and enriches transcripts standard RNA-seq cannot identify.
- We find extensive transcription within and across “intergenic” GWAS regions, identifying ~1500 mostly novel transcripts.
- These novel transcripts provide candidates to explain the traits associated with some of these regions.

Summary: Intergenic GWAS capture sequencing

- Capture-seq allows unprecedented depth of coverage and enriches transcripts standard RNA-seq cannot identify.
- We find extensive transcription within and across “intergenic” GWAS regions, identifying ~1500 mostly novel transcripts.
- These novel transcripts provide candidates to explain the traits associated with some of these regions.
- Investigating intergenic GWAS regions with capture seq is a good method for identifying potentially functional lncRNAs.

Acknowledgements

Garvan Institute

John Mattick

Marcel Dinger

Tim Mercer

Institute for Molecular Bioscience

Ryan Taft

Jo Crawford

Kelin Ru

Roche

Jeffrey Jeddloh

Mater Medical Research Institute

Daniel Gerhardt

Acknowledgements



Garvan Institute
John Mattick
Marcel Dinger
Tim Mercer

Institute for Molecular Bioscience
Ryan Taft
Jo Crawford
Kelin Ru

Roche
Jeffrey Jeddelloh

Mater Medical Research Institute
Daniel Gerhardt

