

Targeted sequencing of highly homologous genes using Cas9 digestion

Chris Boles, Christine Vasselin, Ezra Abrams

Sage Science, Inc., Beverly, MA 01915



Abstract

Next generation sequencing (NGS) is now widely used for genetic testing in clinical settings. Genes with high homology to other regions of the genome present major challenges for current short-read targeted sequencing methods. We have developed a new method for addressing this problem, which relies on the programmable, high-specificity DNA cleavage of the Cas9 enzyme. The process involves four steps: a) isolation of extremely high molecular weight (HMW) DNA, b) digestion of the HMW DNA with customized Cas9 enzymes that cleave unique sites at the boundaries of the DNA sequencing target, c) size selection electrophoresis to isolate the sequencing target, and d) sequencing of the target using an NGS instrument. To enable this method in clinical settings, we have developed an instrument that can accept isolated white blood cells as the input sample and perform DNA extraction, Cas9 digestion, and preparative target purification in an integrated, semi-automated fashion. We report NGS sequencing results on a medically important gene, PKD1, that has six closely related pseudogenes with sequence homology to about ¾ of the functional gene copy.

Medically relevant genes in regions of high homology

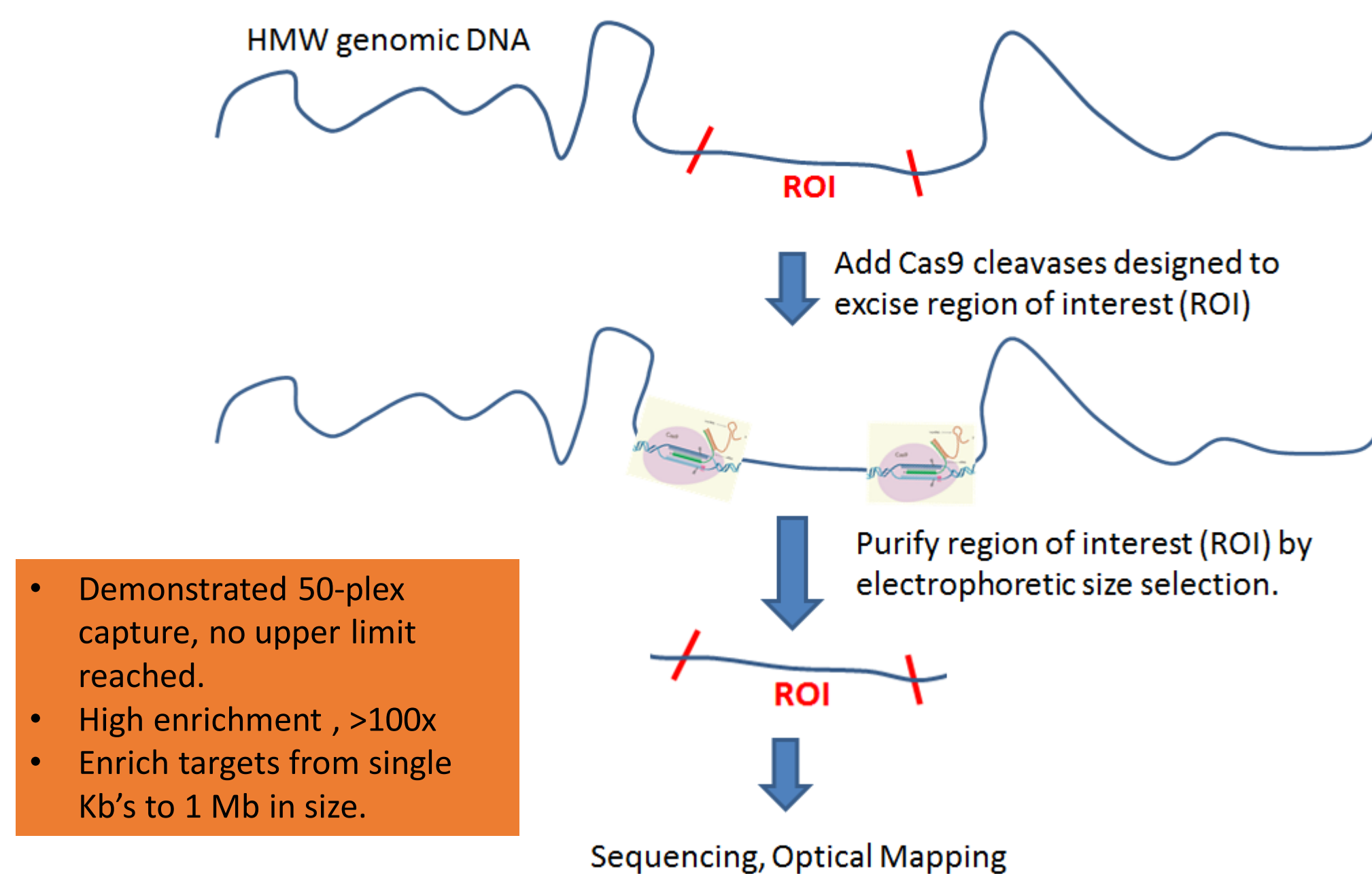
Gene	Affected exons (%)	Affected positions (%)	% Observed low MD	Homology type	Disease(s)	Category
SMN1	14/18 (77.8)	4,826/5,166 (93.4)	92.7	Different gene non-CDS pseudogene	Spinal muscular atrophy	M
RPS17	8/10 (80)	1,850/2,116 (87.4)	76.4	Same gene non-CDS pseudogene	Diamond-blackfan anemia	M
SMN2	14/18 (77.8)	4,826/5,166 (93.4)	92.7	Different gene non-CDS pseudogene	Spinal muscular atrophy	M
KBKG	7/10 (70)	1,921/2,764 (69.5)	63.6	Pseudogene	Incontinentia pigmenti	M
CFP1	5/6 (83.3)	637/1,471 (43.3)	76.7	Different gene	Congenital heart defects	M
ADAMTSL2	9/18 (50)	2,736/5,196 (52.7)	60.2	Non-CDS	Colorblindness, deutan; blue cone monochromacy	M
OPN1MW	7/12 (58.3)	1,915/3,750 (51.1)	67.2	Same gene different gene	Colorblindness, deutan; blue cone monochromacy	M
STRC	10/20 (50)	1,887/3,948 (47.8)	62.2	Different gene non-CDS pseudogene	Sensorineural hearing loss	M
KRT36	3/9 (33.3)	904/2,831 (31.9)	37.9	Different gene non-CDS pseudogene	Monilethrix	M
TURB2B	1/4 (25)	53/213 (24.9)	22.2	Different gene non-CDS pseudogene	Polymyositis	M
LPA	10/39 (25.6)	3,003/11,193 (26.8)	39.5	Same gene different gene non-CDS pseudogene	Coronary artery disease	A
CHRNA7	2/10 (20)	660/2,996 (22.1)	26.9	Different gene	Mac13.3 microdeletion syndrome	M
KRT81	2/9 (22.2)	342/2,688 (12.7)	38.0	Different gene non-CDS pseudogene	Mac13.3 microdeletion syndrome	M
NCF1	2/11 (18.2)	454/2,603 (17.4)	22.3	Pseudogene	Chronic granulomatous disease	M
OTOD	4/20 (20)	1,247/2,398 (52.1)	28.3	Non-CDS	Sensorineural hearing loss	M
KRTD11	1/9 (11.1)	349/2,505 (13.9)	41.9	Different gene non-CDS pseudogene	HIV disease progression	A
TNFR	10/26 (38.5)	2,892/7,942 (36.3)	25.5	Same gene	Ehlers-danlos syndrome	M
OPN1LW	1/6 (16.7)	2,411/8,751 (27.5)	19.8	Different gene	Blue cone monochromacy	M
NEB	16/51 (31.4)	4,768/15,143 (31.5)	15.3	Same gene	Nemaline myopathy	M
CORO1A	1/10 (10)	235/2,886 (8.2)	7.9	Non-CDS	Immunodeficiency	M
OCLN	1/8 (12.5)	172/2,809 (6.1)	36.0	Pseudogene	Basal-like carcinoma with simplified gyration and polymyositis	M
FLG	1/2 (50)	802/2,446 (32.8)	20.0	Same gene	Ichthyosis vulgaris	M
HYDN	6/86 (7)	1,701/26,643 (6.4)	67.4	Non-CDS pseudogene	Primary ciliary dyskinesia	M
RHCE	1/10 (10)	157/2,554 (6.1)	17.4	Different gene	Rh blood group antigens	M
PMS2	1/15 (6.7)	274/4,539 (6)	20.4	Non-CDS pseudogene	HNPPC	M, S
STAT5B	1/18 (5.6)	266/4,704 (5.7)	15.0	Different gene	Growth hormone insensitivity with immunodeficiency	M
TTN	7/83 (8.4)	1,308/16,621 (7.8)	2.2	Same gene	Dilated cardiomyopathy	M

Gene	Affected exons (%)	Affected positions (%)	% Observed low MD	Homology type	Disease(s)	Category
RPS17	8/10 (80)	1,850/2,116 (87.4)	76.4	Same gene non-CDS pseudogene	Diamond-blackfan anemia	M
SMN1	12/18 (66.7)	2,310/3,850 (60)	92.7	Different gene non-CDS pseudogene	Spinal muscular atrophy	M
SMN2	12/18 (66.7)	2,310/3,850 (60)	92.7	Different gene non-CDS pseudogene	Spinal muscular atrophy	M
ADAMTSL2	9/18 (50)	2,736/5,196 (52.7)	60.0	Non-CDS	Geleophysic dysplasia	M
KBKG	6/10 (60)	1,447/2,764 (52.4)	63.6	Pseudogene	Incontinentia pigmenti	M
OPN1MW	5/12 (41.7)	1,431/3,750 (38.2)	67.2	Same gene different gene	Colorblindness, deutan; blue cone monochromacy	M
CFP1	3/6 (50)	367/1,471 (25.0)	76.7	Different gene	Congenital heart defects	M
OPN1LW	2/6 (33.3)	3,961/8,751 (45.3)	19.8	Different gene	Blue cone monochromacy	M
STRC	4/20 (20)	1,888/3,948 (47.8)	62.2	Different gene non-CDS pseudogene	Sensorineural hearing loss	M
CORO1A	2/10 (20)	344/2,886 (11.9)	7.9	Non-CDS	Immunodeficiency	M
GRK1	2/7 (28.6)	300/2,602 (11.5)	0.0	Different gene	Oguchi disease	M
LPA	6/39 (15.4)	1,220/11,193 (10.9)	39.5	Same gene different gene non-CDS pseudogene	Coronary artery disease	A
OCLN	2/8 (25)	113/2,809 (4.0)	36.0	Pseudogene	Basal-like carcinoma with simplified gyration and polymyositis	M
NCF1	1/11 (9.1)	202/2,603 (7.8)	22.3	Pseudogene	Chronic granulomatous disease	M
RHCE	1/10 (10)	157/2,554 (6.1)	17.4	Different gene	Rh blood group antigens	M
NEB	11/51 (21.6)	2,560/15,143 (16.9)	15.3	Same gene	Nemaline myopathy	M
OTOD	2/20 (10)	375/2,398 (15.7)	28.3	Non-CDS	Sensorineural hearing loss	M
TNFR	4/26 (15.4)	4,692/15,143 (31.0)	25.5	Same gene	Ehlers-danlos syndrome	M
PMS2	1/15 (6.7)	274/4,539 (6)	20.4	Non-CDS pseudogene	HNPPC	M, S

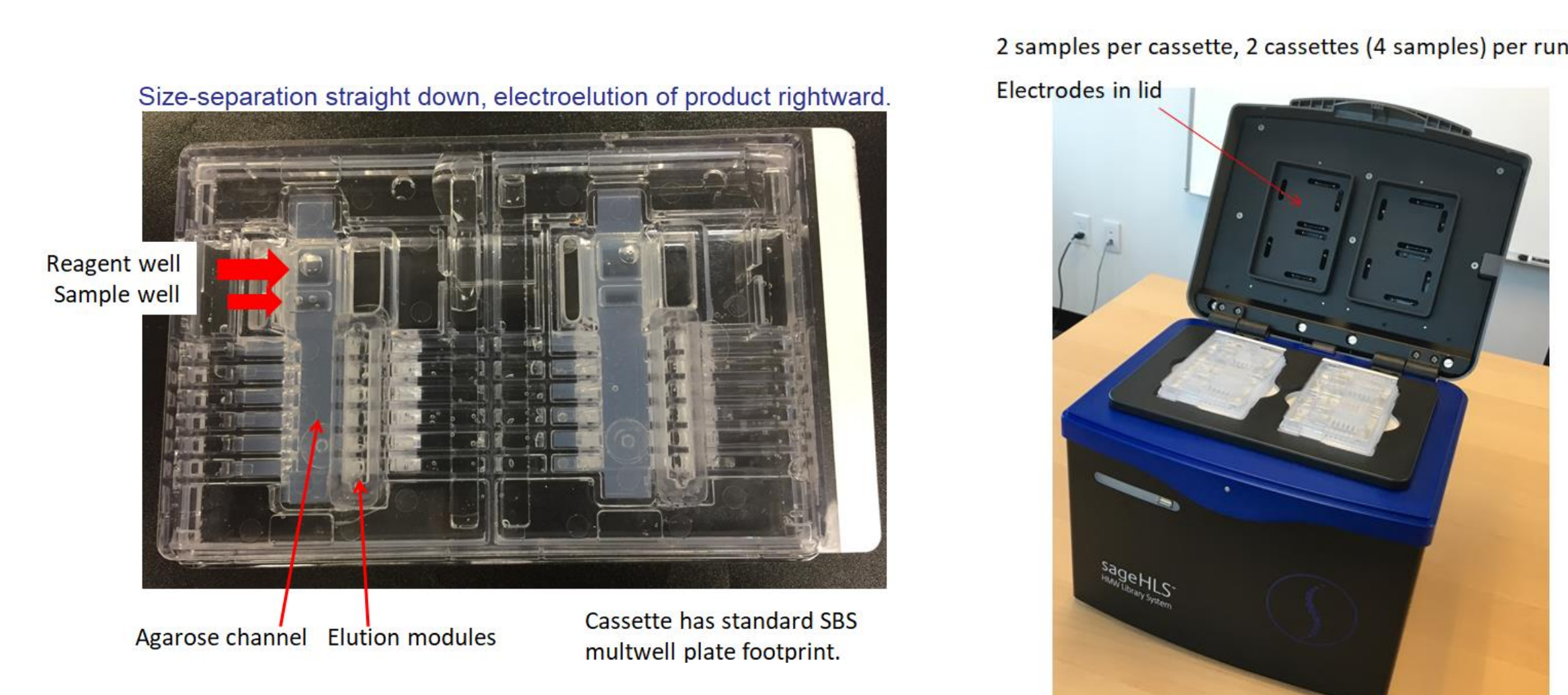
Genes for which design of hybridization baits (NGS dead zone) or PCR amplicons (Sanger dead zone) are complicated by segmental duplication or pseudogenes.

(From Mandelker, et al., 2016, Genetics In Medicine 18(2): 1282-1289.)

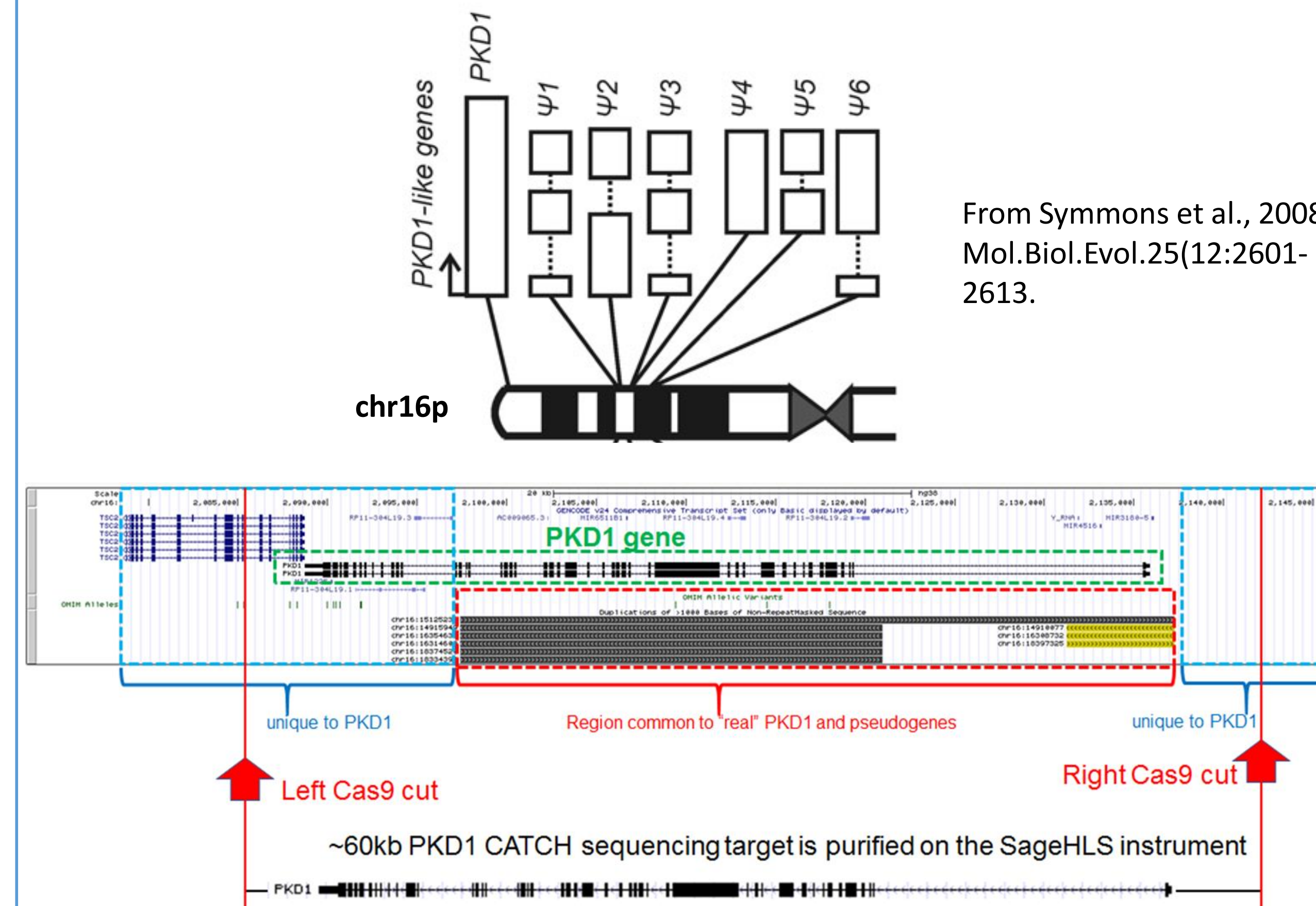
CATCH concept for isolation of genomic DNA targets



SageHLS Cassette and Instrument



Cas9 guide RNA design for PKD1, a gene with six pseudogenes

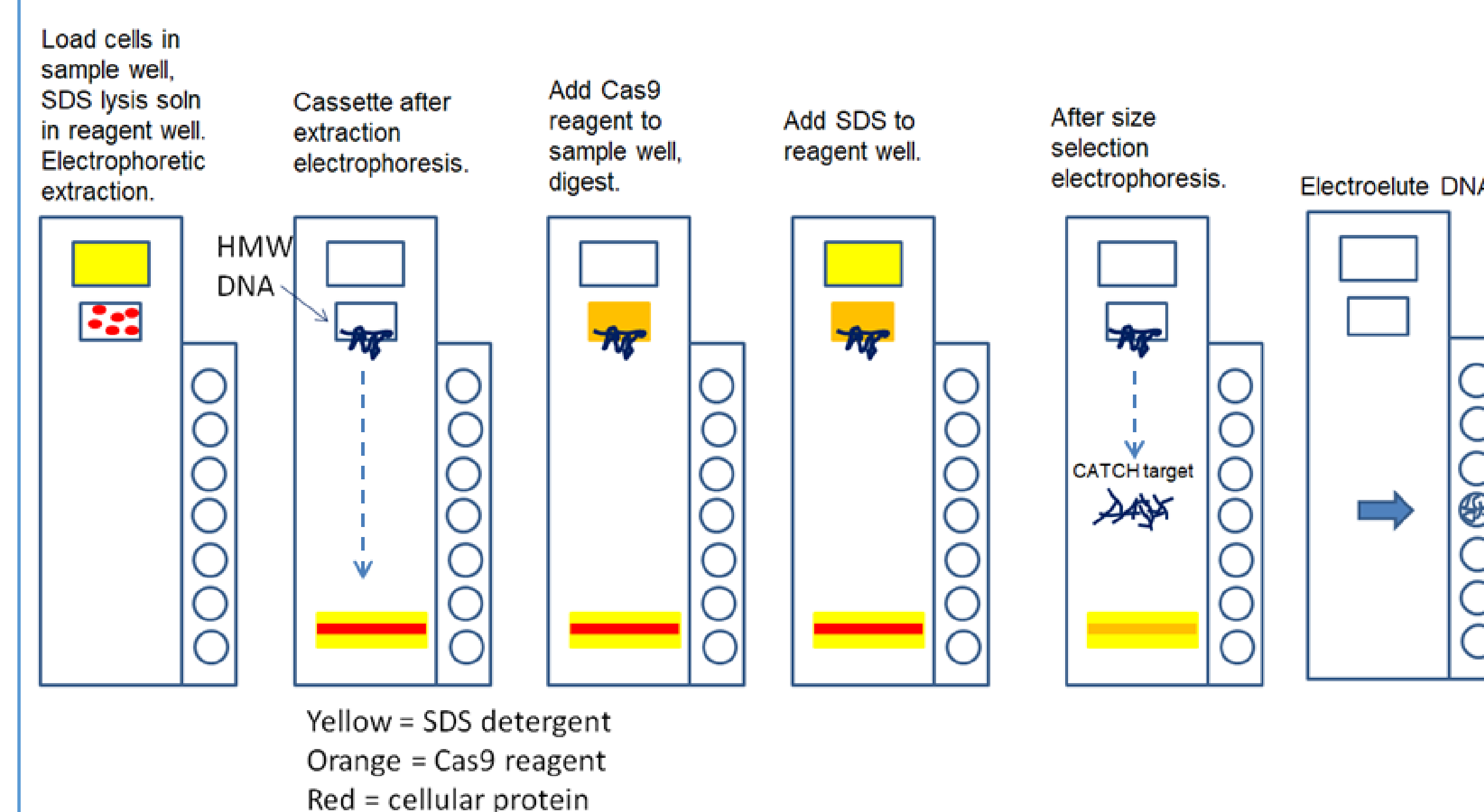


About ¾ of the active PKD1 gene sequence is shared with six pseudogenes, all located on chr16. To excise the active gene, we select gRNA sequences from unique sequence regions that flank the regions that are duplicated in the various PKD-pseudogenes. Using this strategy, a ~60kb CATCH target can be isolated and sequenced without complications from the six other repeated regions.

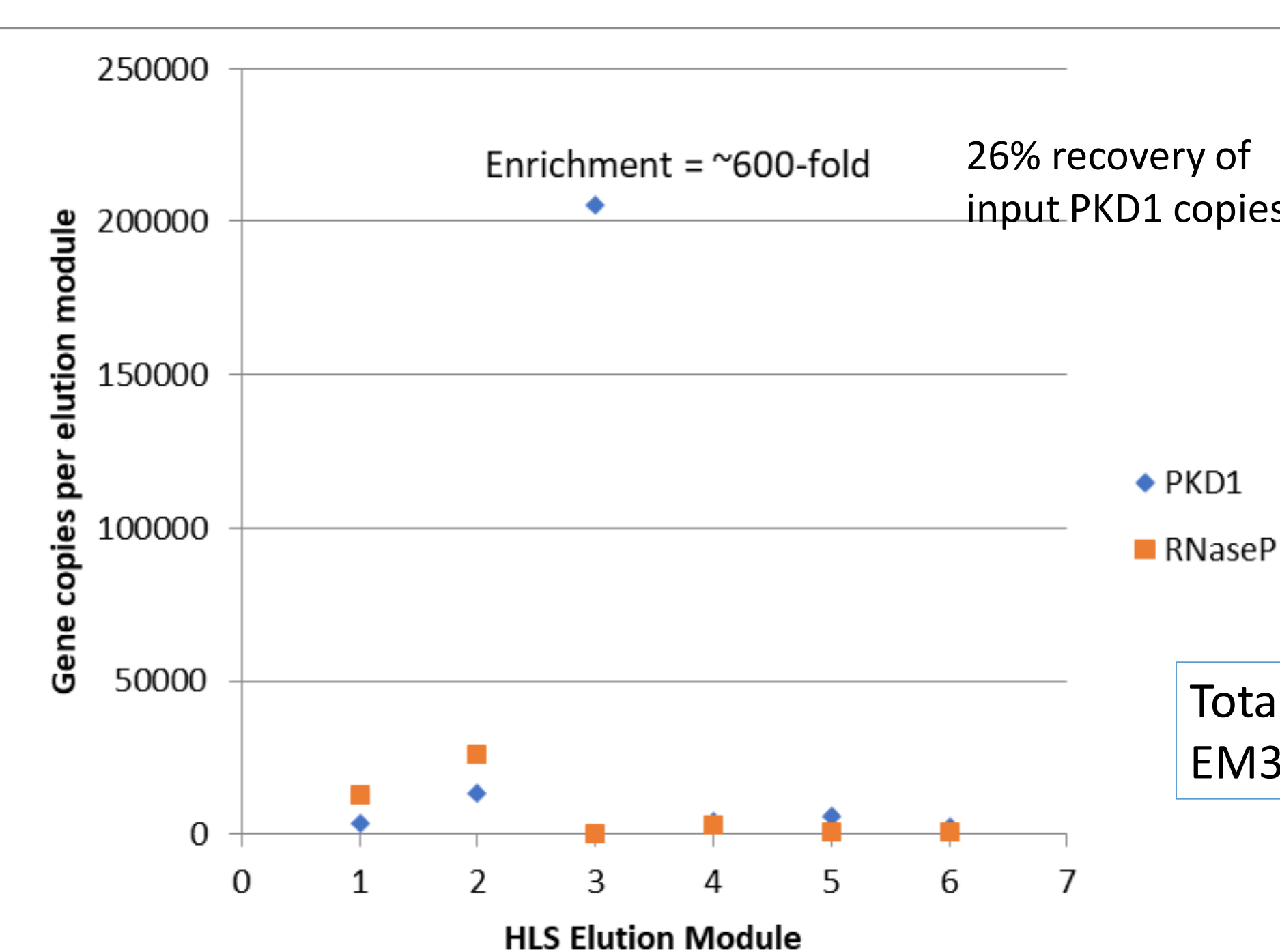
Integrated DNA extraction, digestion, and size-selection using

HLS-CATCH

~375,000 NA12878 human lymphoblastoid cells were loaded into the sample well of an HLS cassette. A lysis buffer containing SDS was loaded into the upstream reagent well. Electrophoretic extraction was carried out for 3 hours, so that the SDS was driven through the sample well. After extraction, the HMW DNA, which remained trapped in the sample wall, was digested with Cas9-gRNA complexes specific for the PKD1 gene. After digestion, the excised PKD1 gene was recovered using an automated HLS size-selection + electroelution program.



Enrichment of the PKD1 CATCH product as measured by qPCR

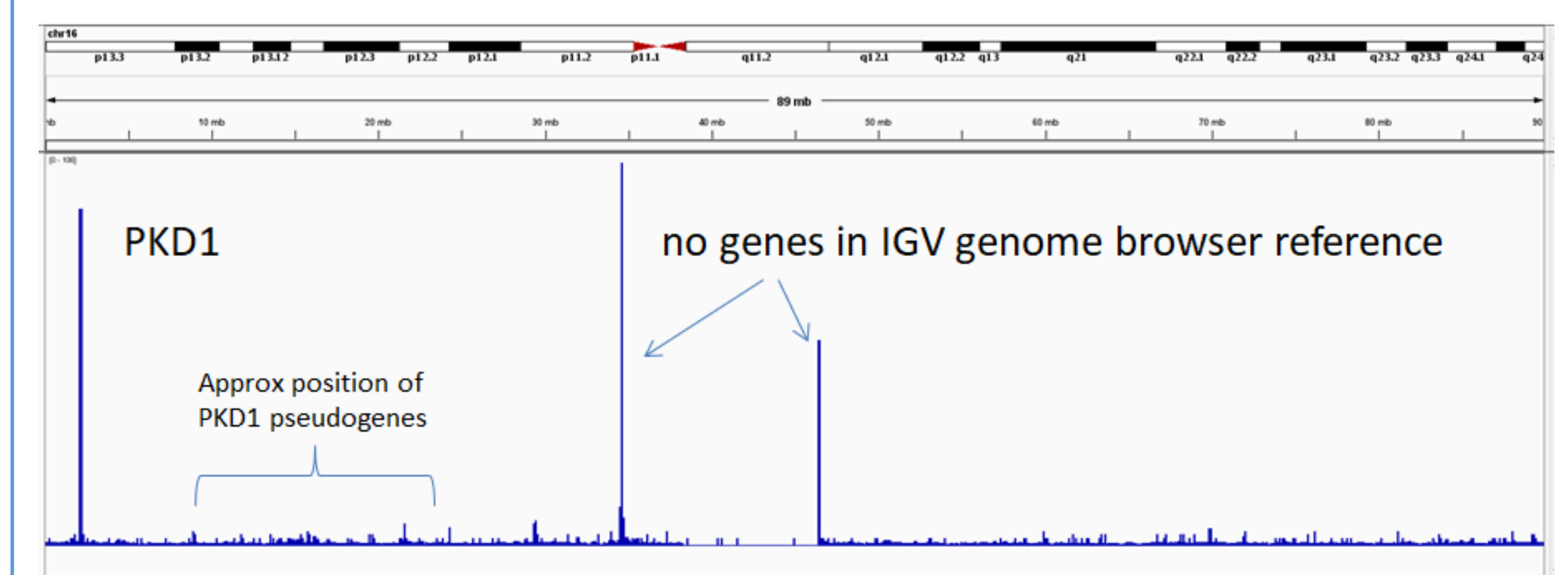


Sequencing

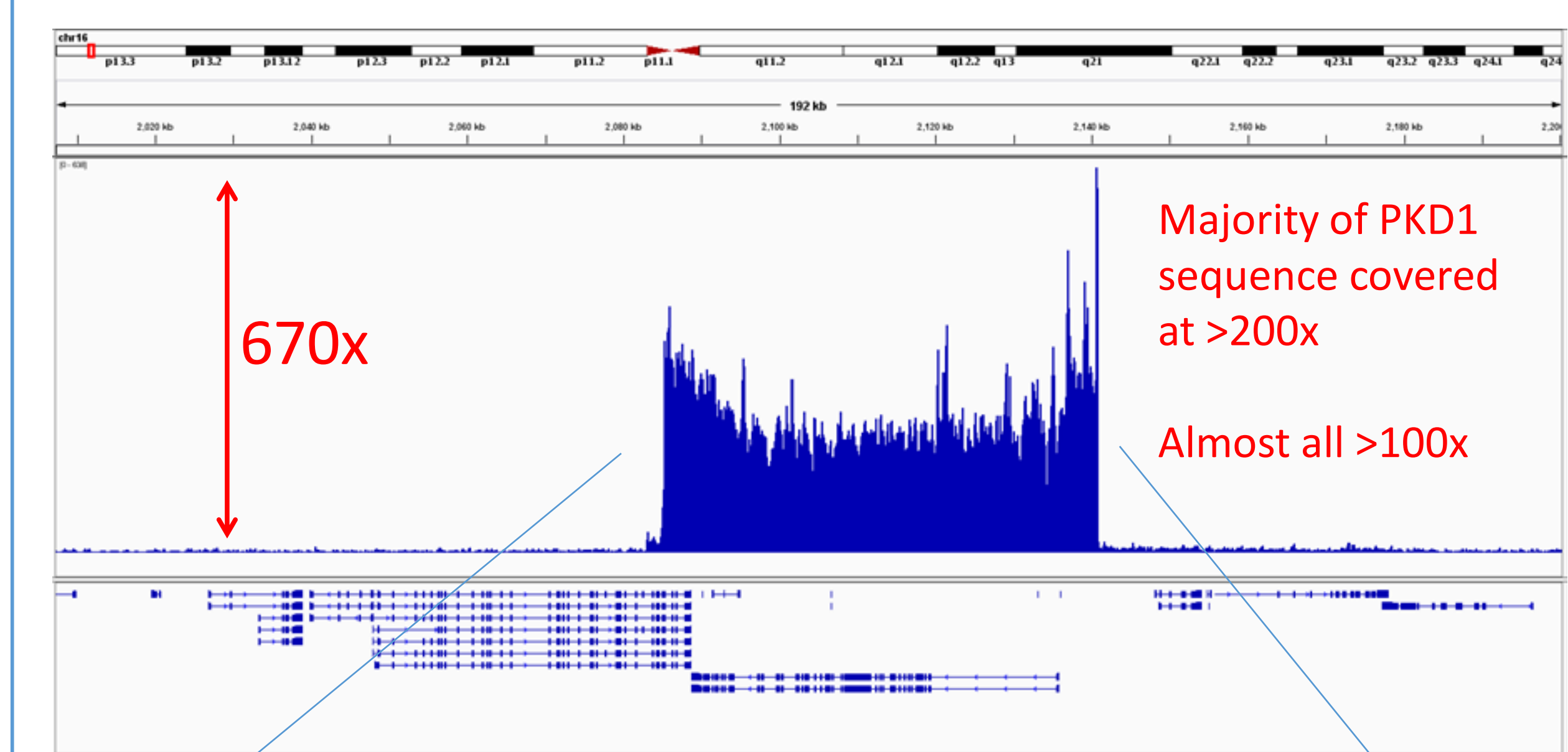
1. Identify target-enriched fractions by qPCR
2. Shear to 300bp with Covaris
3. Ethanol precipitate
4. Library construction with Kapa HyperPrep Library kit
5. QC with Bioanalyzer, Kapa Illumina Library Quant kit
6. Miseq, 2x150 bp PE sequencing
7. Data analysis (BWA MEM, Samtools, Bedtools, IGV)

Sequencing Results

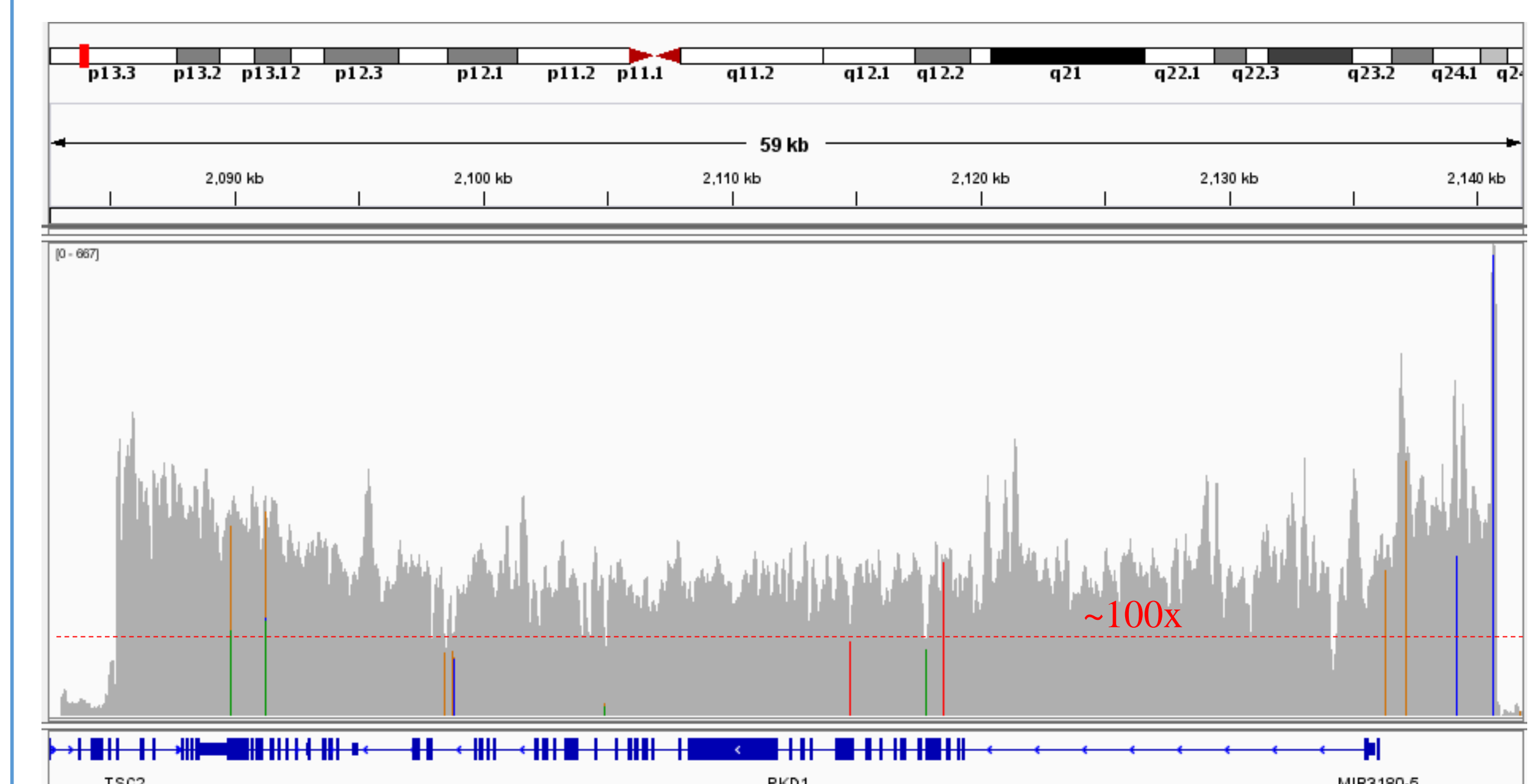
Bedgraph coverage over entire chr16



Bedgraph coverage near PKD1



IGV view of coverage and variants



Acknowledgements

We thank Dr. Marjorie Beggs of Arkana Laboratories (Little Rock, AR) for drawing our attention to the PKD1 gene as a medically important gene that could be addressed with the HLS-CATCH workflow.