# Systematic analysis of dark and camouflaged genes: disease-relevant genes hiding in plain sight

Mark T. W. Ebbert[1,2,†,*], Tanner D. Jensen[1,†], Karen Jansen-West[1], Jonathon P. Sens[1], Joseph S. Reddy[1], Perry G. Ridge[3], John S. K. Kauwe[3], Veronique Belzil[1], Luc Pregent[1], Minerva M. Carrasquillo[1], Dirk Keene[4], Eric Larson[5], Paul Crane[5], Yan W. Asmann[6], Nilufer Ertekin-Taner[1,7], Steven G. Younkin[1], Owen A. Ross[1], Rosa Rademakers[1], Leonard Petrucelli[1,2,*], John D. Fryer[1,2,*]


[1]Department of Neuroscience, Mayo Clinic, Jacksonville, FL 32224, USA; [2]Mayo Clinic Graduate School of Biomedical Sciences, Jacksonville, FL 32224, USA; [3]Department of Biology, Brigham Young University, Provo, UT 84602, USA; [4]Department of Pathology, University of Washington, Seattle, Washington 98195, USA; [5]Department of Medicine, University of Washington, Seattle, Washington 98195, USA; [6]Department of Health Sciences Research, Mayo Clinic, Jacksonville, FL 32224, USA; [7]Department of Neurology, Mayo Clinic, Jacksonville, Florida 32224, USA


*Corresponding authors
†Contributed equally

Author emails:
Mark Ebbert: ebbert.mark@mayo.edu
Tanner Jensen: jensen.tanner@mayo.edu
Karen Jansen-West: jansen.karen@mayo.edu
Jonathon Sens: sens.jonathon@mayo.edu
Joseph Reddy: reddy.joseph@mayo.edu
Perry Ridge: perry.ridge@byu.edu
John Kauwe: kauwe@byu.edu
Veronique Belzil: belzil.veronique@mayo.edu
Luc Pregent: pregent.luc@mayo.edu
Minerva Carrasquillo: carrasquillo.minerva@mayo.edu
Dirk Keene: cdkeene@uw.edu
Eric Larson: larson.e@ghc.org
Paul Crane: pcrane@uw.edu
Yan Asmann: asmann.yan@mayo.edu
Nilufer Ertekin-Taner: taner.nilufer@mayo.edu
Steven Younkin: younkin.steven@mayo.edu
Owen Ross: ross.owen@mayo.edu
Rosa Rademakers: rademakers.rosa@mayo.edu
Leonard Petrucelli: petrucelli.leonard@mayo.edu
John Fryer: fryer.john@mayo.edu

# 1 Abstract

2 **Background:** The human genome contains 'dark' gene regions that cannot be adequately assembled or

3 aligned using standard short-read sequencing technologies, preventing researchers from identifying mutations

4 within these gene regions that may be relevant to human disease. Here, we identify regions that are 'dark by

5 depth' (few mappable reads) and others that are 'camouflaged' (ambiguous alignment), and we assess how

6 well long-read technologies resolve these regions. We further present an algorithm to resolve most

7 camouflaged regions (including in short-read data) and apply it to the Alzheimer's Disease Sequencing Project

8 (ADSP; 13 142 samples), as a proof of principle.

9

10 **Results:** Based on standard whole-genome Illumina sequencing data, we identified 37873 dark regions in 5857

11 gene bodies (3635 protein-coding) from pathways important to human health, development, and

12 reproduction. Of the 5857 gene bodies, 494 (8.4%) were 100% dark (142 protein-coding) and 2046 (34.9%)

13 were ≥5% dark (628 protein-coding). Exactly 2757 dark regions were in protein-coding exons (CDS) across 744

14 genes. Long-read sequencing technologies from 10x Genomics, PacBio, and Oxford Nanopore Technologies

15 reduced dark CDS regions to approximately 45.1%, 33.3%, and 18.2% respectively. Applying our algorithm to

16 the ADSP, we rescued 4622 exonic variants from 501 camouflaged genes, including a rare, ten-nucleotide

17 frameshift deletion in *CR1*, a top Alzheimer's disease gene, found in only five ADSP cases and zero controls.

18

19 **Conclusions:** While we could not formally assess the *CR1* frameshift mutation in Alzheimer's disease

20 (insufficient sample-size), we believe it merits investigating in a larger cohort. There remain thousands of

21 potentially important genomic regions overlooked by short-read sequencing that are largely resolved by long-

22 read technologies.

23

24

# Keywords

3

# 1 Background

2 Researchers have known for years that large, complex genomes, including the human genome, contain 'dark'

3 regions—regions where standard high-throughput short-read sequencing technologies cannot be adequately

4 assembled or aligned—thus preventing our ability to identify mutations within these regions that may be

5 relevant to human health and disease. Some dark regions are what we term 'dark by depth' (few or no

6 mappable reads), while others are what we term 'dark by mapping quality' (reads aligned to the region, but

7 with a low mapping quality). Regions that are dark by depth may arise because the region is inherently difficult

8 to sequence at the chemistry level (e.g., high GC content [1, 2]), essentially eliminating sequencing reads from

9 that region altogether. Other dark regions arise, not because the sequencing is inherently problematic, but

10 because of bioinformatic challenges. Specifically, many dark regions arise from duplicated genomic regions,

11 where confidently aligning short reads to a unique location is not possible; we term these regions as

12 'camouflaged'. These camouflaged regions are generally either large contiguous tandem repeats (e.g.,

13 centromeres, telomeres, and other short tandem repeats), or a larger specific DNA region that has been

14 duplicated (e.g., a gene duplication) either in tandem or in a more distal genome region. In fact, many genes in

15 the human genome were duplicated over evolutionary time and are still transcriptionally and translationally

16 active (e.g., heat-shock proteins) [3–9], while others have been duplicated, but are considered inactive (i.e.,

17 pseudogenes). Regardless of whether the duplication is active, however, any genomic region that has been

18 nearly-identically duplicated, and is large enough to prevent sequencing reads from aligning unambiguously

19 will be 'dark', because the aligner cannot determine which genomic region the read originated from.

20

21 When confronted with a read that aligns equally well to two or more camouflaged regions (commonly known

22 as multi-mapping reads [2, 10]), standard next-generation sequence aligners, such as the Burrows-Wheeler

23 Aligner (BWA) [11–13], randomly assign the read to one of the regions and assign a low mapping quality. For

24 BWA, specifically, reads that cannot be uniquely mapped are generally assigned a mapping quality (MAPQ) of

1    0; though, in certain paired-end sequencing scenarios, BWA will assign a high mapping quality if the read mate

2    is confidently mapped nearby (i.e., within the estimated insert-size length).

3

4    Recent work has characterized camouflaged regions, in part, including a study that demonstrates how this

5    issue affects all standard RNA-Seq analyses [10], and another that quantifies the number of nucleotides in

6    human reference GRCh38 that are dark for mapping quality of 0 (camouflaged regions), based on 1000

7    Genome Project data [2]. Robert and Watson demonstrated that expression for 958 genes were either over-

8    or under-represented because of multi-mapping reads across 12 different RNA-Seq processing methods, and

9    no method was immune to the problem [10]. They also demonstrated that many of these genes are directly

10    implicated in human disease. Zheng-Bradley et al. recently re-aligned genomes from the 1000 Genomes

11    Project to GRCh38, and, among other findings, generally demonstrated the breadth of multi-mapping reads

12    across the genome [2]. These data characterize the general problem, and report specific genes affected by this

13    issue.

14

15    Here, we systematically analyze dark and camouflaged genes to more fully characterize the problem, and we

16    highlight many disease-relevant genes that are directly implicated in Alzheimer's disease, autism spectrum

17    disorder, amyotrophic lateral sclerosis (ALS), spinal muscular atrophy (SMA), and others. We also show that

18    long-read sequencing technologies substantially reduce the number of dark and camouflaged regions, and we

19    present a method to address camouflaged regions, even in standard short-read sequencing data. As a proof of

20    concept, we apply our method to the Alzheimer's Disease Sequencing Project (ADSP) data, and identify a rare,

21    ten-nucleotide frameshift deletion in the C3b and C4b binding domain of *CR1*, a top Alzheimer's disease gene

22    [14–22], that is only present in five ADSP cases and zero controls. The ADSP is not large enough to statistically

23    assess association between the *CR1* frameshift mutation and Alzheimer's disease.

24

# 1    Results

2    To quantify the number of dark and camouflaged regions in standard short-read whole-genome sequencing

3    data, we obtained whole-genome sequencing data for ten unrelated males from the Alzheimer's Disease

4    Sequencing Project (ADSP) and scanned each sample for dark and camouflaged regions, averaging across all

5    ten samples; we only used data from males for this study so we could also assess dark and camouflaged

6    regions on the Y chromosome because large portions of the Y chromosome are dark. We ignored incomplete

7    genomic regions (e.g., centromeres). We then limited the dark and camouflaged regions to known gene

8    bodies, based on annotations from build 87 of the Ensembl GRCh37 human reference genome [23]. All ten

9    samples were sequenced using standard Illumina whole-genome sequencing with 100-nucleotide read

10    lengths, where median genome-wide read depths ranged from 35.4x to 42.9x coverage, with an overall

11    median of 39.4x. We performed the same analyses on ten unrelated males from the 1000 Genomes Project

12    [24] that were sequenced using Illumina whole-genome sequencing with 250-nucleotide read lengths, where

13    median genome-wide read depths ranged from 39.3x to 52.6x coverage, with an overall median of 48.9x.

14    Similarly, we assessed how well long-read sequencing technologies, including 10x Genomics (52x median

15    coverage), PacBio (50x median coverage), and ONT (46x median coverage) resolve dark and camouflaged

16    regions. Although we were only able to obtain a single high-depth male genome for each long-read

17    technology, we believe our results are a reasonable estimate for how well each technology addresses dark and

18    camouflaged regions. Larger sequencing studies will further clarify our results.

19

20    We consider a region 'dark' for one of two reasons: (1) insufficient number of reads aligned to the genomic

21    region (dark by depth); and (2) reads aligned to the region, but with insufficient mapping quality for a variant

22    caller to identify mutations in the region (dark by mapping quality). Specifically, we define regions that are

23    dark by depth as those with fewer than five aligned reads (Figure 1a), and regions that are dark by mapping

24    quality as those where $\geq$90% of aligned reads have a mapping quality (MAPQ) <10 (Figure 1b). Defining dark-

6

1    by-depth regions as those with fewer than five reads is a relatively strict cutoff, and likely underestimates the

2    number of dark regions because 20 to 30 reads is often considered a reasonable minimum to confidently

3    identify heterozygous mutations; overall median read depth is an important factor, however, and we believe a

4    strict cutoff provides a more conservative estimate. We used a mapping quality threshold <10 to define

5    regions that are dark by mapping quality because that is the standard cutoff used in the Genome Analysis

6    ToolKit (GATK) [25]. Camouflaged regions are those that are dark by mapping quality because the region has

7    been duplicated in the genome (Figure 1c). We identified sets of camouflaged regions (regions camouflaged by

8    each other) using BLAT [26], where we required at least 98% sequence identity for two regions to be included

9    in the same set.

10

**Standard short-read sequencing leaves 37873 dark regions across 5857 gene bodies, including protein-**

**coding exons from 744 genes**

13    Using whole-genome Illumina sequencing data (100-nucleotide read lengths) from ten unrelated males, we

14    identified 37873 dark regions (>16 million nucleotides) in 5857 gene bodies (based on Ensemble GRCh37 build

15    87 gene annotations) that were either dark by depth or dark by mapping quality (Supplemental Figure 1a;

16    Supplemental Tables 1-2). Stratifying the gene bodies by GENCODE biotype [27], 3635 gene bodies were

17    protein coding, 1102 were pseudogenes, and 720 were long intergenic non-coding RNAs (lincRNA; Figure 2a).

18    Of all 37873 dark gene-body regions, 28598 were intronic, 4113 were in non-coding RNA exons (e.g., lincRNAs

19    and pseudogenes), 2657 were in protein-coding exons (CDS), 1134 were in 3'UTR regions, and 1103 were in

20    5'UTR regions (Figure 2b; Supplemental Table 1). Any dark region that spanned a gene element boundary (e.g.,

21    intron to exon) was split into separate dark regions. Of the 5857 gene bodies, 494 (8.4%) were 100% dark,

22    1560 (26.6%) were at least 25% dark, and 2046 (34.9%) were at least 5% dark (Supplemental Figure 1b;

23    Supplemental Table 1).

24

1    Focusing only on CDS regions, we identified 2757 dark CDS regions (>460000 nucleotides) across 744 protein-

2    coding genes that were dark by either depth or mapping quality (Figure 3a; Supplemental Tables 1-2). Exactly

3    142 (19.1%) of the 744 protein-coding genes were 100% dark in CDS regions, 441 (59.3%) were at least 25%

4    dark in CDS regions, and 628 (84.4%) were at least 5% dark in CDS regions (Figure 3b; Supplemental Table 1).

5    Exactly 474 of the 628 genes that were 5% dark in CDS regions were dark because they were camouflaged.

6

7    **Most dark regions are specifically camouflaged**

8    Regions may be dark because of either low depth or low mapping quality, but the majority of regions are dark

9    because of mapping quality, and specifically because they are camouflaged (low mapping quality because of a

10    duplication). Exactly 3953 of the 5857 dark gene bodies are dark because of mapping quality, where 3252 are,

11    in fact, camouflaged. We also measured the number of times each gene region was duplicated and found that

12    70% of gene regions were replicated three or fewer times in the genome, but 84 regions were duplicated ≥100

13    times (Supplemental Figure 2a), with the most repeated regions (ten separate intronic regions totaling 2235

14    nucleotides from *C5orf48*) being replicated 941 times in aggregate. Limiting to only CDS regions, we estimate

15    that 74.1% are replicated three or fewer times, with 38 replicated ≥10 times (Supplemental Figure 2b) and the

16    most repeated region was from *NBPF12*, in which 173 nucleotides were replicated 37 times.

17

18    **Long-read sequencing technologies resolve substantial portions of the dark regions**

19    Data from the samples sequenced using 250-nucleotide Illumina read lengths reduced the percentage of dark

20    nucleotides by 30.1% and 24.4% for all gene bodies, and for only CDS regions, respectively, leaving 69.9% and

21    75.6% of the nucleotides dark, respectively (Supplemental Figure 1b; Figure 3b; Supplemental Tables 3-4).

22    Comparing long-read sequencing technologies to the standard Illumina 100-nucleotide read lengths, the ONT

23    platform performed best, both when assessing entire gene bodies, and when considering only CDS regions.

24    Specifically, approximately 41.2%, 25.8%, and 24.9% of the nucleotides remained dark for all gene bodies for

8

1    PacBio, 10x Genomics, and ONT, respectively (Supplemental Figure 1b; Supplemental Tables 5-10). Similarly,

2    approximately 42.2%, 31.4%, and 18.5% of CDS nucleotides remained dark for 10x Genomics, PacBio, and

3    ONT, respectively (Figure 3b; Supplemental Tables 5-10). In contrast to overall gene-body results, PacBio

4    outperformed 10x Genomics when looking only at CDS regions (Supplemental Figure 1b; Figure 3b). The long-

5    read technologies improved over Illumina mostly by reducing the percentage of nucleotides that are dark by

6    mapping quality (Supplemental Figure 1c). Surprisingly, the percentage of gene-body regions that are dark

7    because of low depth is higher for all long-read technologies than it is for Illumina (Supplemental Figure 1c).

8

9    We generated a density plot for the length of all dark-by-mapping quality regions to approximate the

10   proportion of regions each sequencing technology should be able to resolve (Supplemental Figure 3), which

11   resulted in a bimodal distribution. The two modes are located at 95 and 538 nucleotides. As expected, median

12   read lengths for the Illumina whole-genome sequencing based on 100-nucleotide and 250-nucleotide read

13   lengths were 100 and 250 nucleotides, respectively. The first mode for the camouflaged region lengths is at

14   95, explaining why 100-nucleotide read lengths are insufficient to unambiguously span most dark-by-mapping

15   quality regions. The 250-nucleotide read lengths fall between the two modes, explaining why 250-nucleotide

16   read lengths resolve a high percentage of camouflaged regions. In other words, 100-nucleotide read lengths

17   are too short to bridge most camouflaged regions, but 250-nucleotide read lengths appear to be sufficient for

18   many. Median read lengths for both the ONT and PacBio genomes we used in this study were 6276 (N50 =

19   33973) and 8511 (N50 = 17467) nucleotides, respectively, which is shorter than expected, but substantially

20   longer than necessary to resolve most camouflaged regions. We believe comparing median read lengths,

21   rather than N50, is more useful in this scenario, because we are interested to know what percentage of reads

22   are likely to bridge a given dark or camouflaged region. Our results suggest that our estimates for the

23   percentage of camouflaged regions ONT and PacBio are able to resolve may be conservative because a longer

24   DNA library should resolve even more camouflaged regions.

9

1

**Important pathways and gene families are affected by dark and camouflaged regions**

Because such a large number of genes are dark, we characterized the pathways for genes that are not fully

represented in standard Illumina short-read sequencing (100-nucleotide reads) datasets. We included all

genes where at least 5% of the CDS regions were dark (670 unique gene symbols) and identified several

pathways that are important in human health, development, and reproductive function (Figure 4a;

Supplemental Table 11). Specific pathways included defensins (R-HSA-1461973; logP = -7.04), gonadal

mesoderm development (GO:0007506; logP = -6.18), base-excision repair (GO:0006284; logP = -5.93),

chromatin silencing (GO:0006342; logP = -5.86), Deubiquitination (R-HSA-5688426; logP = -5.32), NLS-bearing

protein import into nucleus (GO:0006607; logP = -5.31), spindle assembly (GO:0051225; logP = -5.19),

spermatogenesis (GO:0007283; logP = -4.93), and forebrain neuron differentiation (GO:0021879; logP = -4.09).

Some specific gene families involved in these pathways include eleven defensin genes (e.g., DEFA1 and

DEFB4A), five testis specific proteins (e.g., TSPY2), eleven ubiquitin-specific 17-like family members, and

twelve golgin genes (e.g., GOLGA6B; Supplemental Table 11).

Looking specifically at known protein-protein interactions, we found 138 proteins with 212 known interactions

(Supplemental Figure 4), and within those, identified three groups enriched for protein-protein interactions

using the MCODE algorithm [28] (Figure 4b). All three MCODE groups combined are primarily associated with

RNA transport (hsa030313; logP = -17.3; Supplemental Figure 5; accessed December 2018). Individually, the

first group (MCODE1) is enriched for proteins involved in systemic lupus erythematosus (hsa05322; logP = -

6.7), cellular response to stress (R-HSA-2262752; logP = -6.6), and RNA transport (hsa03013; logP = -4.39;

Supplemental Figure 6). The second group (MCODE2) is enriched with proteins involved in NLS-bearing protein

import into nucleus (GO:0006607; logP = -17.1) and protein import into nucleus (GO:0006606; logP = -15.4;

Supplemental Figure 7). The third group does not have significant enrichment associations, likely because little

10

1  is known about them; all four genes (*PRR20B*, *PRR20C*, *PRR20D*, and *PRR20E*) are 100% camouflaged and do

2  not even have known expression measurements in GTEx [29] (Supplemental Figures 8-11).

3

4  **There are 75 genes with known mutations associated with 305 human phenotypes**

5  To assess the potential impact missing mutations in dark genes may have on human disease genetics, we

6  measured the number of dark genes with at least 5% dark CDS that have mutations known to be involved in

7  human disease; we calculated the number of genes that are ≥5% dark CDS with a mutation in the Human Gene

8  Mutation Database (HGMD) [30]. We found 75 genes associated with 305 unique human phenotypes,

9  including 277 diseases (Figure 5a). Some of the diseases with the most known associated genes include autism

10  spectrum disorder, hemophilia A, schizophrenia, hearing loss, spinal muscular atrophy, and inflammatory

11  bowel disease. Some of the diseases most represented in our data are not surprising, given the number of

12  genes involved in the disease, but these data demonstrate the number of diseases impacted by genes that are

13  at least 5% dark CDS. We also performed an enrichment analysis, where the diseases most enriched for dark

14  genes included Hemophilia A, color blindness (protan colour vision defect), and X-linked cone-rod dystrophy

15  (Supplemental Figure 12).

16

17  Similarly, we quantified the number of diseases each gene was associated with (Figure 5b). We identified

18  many disease-relevant genes with large portions of dark CDS regions that may harbor critical disease-

19  modifying mutations that currently go undetected. Some of the genes with the most known disease

20  associations include *ARX* (14.0% dark CDS), *NEB* (9.5% dark CDS), *TBX1* (10.5% dark CDS), *RPGR* (12.9% dark

21  CDS), *HBA2* (12.8% dark CDS), and *CR1* (26.5% dark CDS). The *CR1* gene is particularly notable given that *CR1* is

22  a top-ten Alzheimer's disease gene. Other notable genes include *SMN1* (89.9% dark CDS) and *SMN2* (88.2%

23  dark CDS), which are known to be involved in spinal muscular atrophy (SMA) and ALS. *HSPA1A* (52.8% dark

11

1    CDS) and *HSPA1B* (51.1% dark CDS) also encode two primary 70-kilodalton (kDa) heat-shock proteins, a family

2    of proteins that have been implicated in ALS [31, 32].

3

4    **Camouflaged genes are consistently dark in gnomAD, but dark-by-depth genes may be sample or dataset**

5    **specific**

6    Although most dark genes are specifically camouflaged (Supplemental Tables 12-13), many are dark by depth

7    in the ADSP data; upon manual comparison between whole-genome sequencing data from the ten ADSP

8    males and coverage plots from the gnomAD consortium dataset (http://gnomad.broadinstitute.org/) [33], we

9    found that camouflaged regions in the ADSP males are consistently dark in the gnomAD data, demonstrating

10   that these camouflaged regions are consistent across datasets. The dark-by-depth regions are more variable

11   between samples and datasets, however, suggesting these regions may be sensitive to specific aspects of

12   whole-genome sequencing (e.g., library preparation) or downstream analyses. Specific camouflaged genes

13   include *SMN1* and *SMN2* (89.9% and 88.2% dark CDS, respectively; Figure 6a), *HSPA1A* and *HSPA1B* (52.8%

14   and 51.1% dark CDS, respectively; Figure 6b), *NEB* (9.5% dark CDS; Figure 6c), and *CR1* (26.5% dark CDS; Figure

15   6d). Specific dark-by-depth genes include *HLA-DRB5* (50.2% dark CDS; Figure 6e), *RPGR* (12.9% dark CDS;

16   Figure 6f), *ARX* (14.0% dark CDS; Figure 6g), and *TBX1* (10.5% dark CDS; Figure 6h). All four camouflaged genes

17   are also dark in the gnomAD data. A manual inspection of our dark-by-depth gene list, however, suggests most

18   are not completely dark in gnomAD, but vary by sample or dataset. Specifically, *HLA-DRB5* and *RPGR* in

19   gnomAD appear to be consistent with the ADSP data; *ARX* and *TBX1*, however, only appear to be dark in a

20   portion of the gnomAD samples, where about 30% of samples have ≤5 reads in their respectively defined dark

21   regions (Note: our threshold for dark regions is <5 reads, but the gnomAD plots for *ARX* and *TBX1* are based on

22   ≤5 reads). Dark regions (Figures 6a-h) are either similar or more pronounced in the gnomAD whole-exome

23   data than what we observed in the whole-genome data, highlighting that dark and camouflaged regions are

24   generally magnified in whole-exome data; this is likely because of differences in library preparation and

1    shorter read lengths in exome data. For interest, we also found that *APOE*—the top genetic risk for

2    Alzheimer's disease [34–36]—is approximately 6% dark CDS (by depth) for certain ADSP samples with whole-

3    genome sequencing, and the same region is dark in gnomAD whole-exome data (Supplemental Figure 13). It is

4    possible some of the dark regions we identified in standard short-read whole-genome data are specific to the

5    ADSP samples, but additional work can clarify this issue. In either case, *dark-by-depth regions* (Supplemental

6    Tables 14-15) *should be interrogated within individual datasets, and perhaps for individual samples as a*

7    *quality control measure.*

8

9    *SMN1* and *SMN2* are camouflaged by each other, where both genes are known to contribute to spinal

10    muscular atrophy, and have been implicated in ALS. *HSPA1A* and *HSPA1B* are also camouflaged by each other,

11    and the heat-shock protein family has been implicated in ALS [37, 38]. *NEB* is a special case that is

12    camouflaged by itself (rather than another gene), and is associated with 24 diseases in the HGMD, including

13    nemaline myopathy, a hereditary neuromuscular disorder. *NEB* is a large gene (249151 nucleotides; 25577

14    CDS nucleotides), thus, ~9.5% dark CDS translates to 2424 dark protein-coding bases. *CR1* is a top Alzheimer's

15    disease gene that plays a critical role in the complement cascade as a receptor for the C3b and C4b

16    complement components, and potentially helps clear amyloid-beta (Aß) [39–41]. Like *NEB*, *CR1* is also

17    camouflaged by itself, where the repeated region actually includes the extracellular C3b and C4b binding

18    domain. The number of repeats and density of certain isoforms have been associated with Alzheimer's disease

19    [21, 42–45].

20

21    We found *HLA-DRB5* is dark by depth in the ADSP and gnomAD data, and has been implicated in several

22    diseases, including Alzheimer's disease. *RPGR* is likewise dark in ADSP and gnomAD, and is associated with

23    several eye diseases, including retinitis pigmentosa and cone-rod dystrophy. We identified *ARX* as a dark-by-

1    depth gene, but this gene appears to vary by sample or cohort, as only approximately 30% of gnomAD samples

2    are strictly dark by depth, using our cutoff of <5 reads. *ARX* is associated with diseases including early infantile

3    epileptic encephalopathy 1 (EIEE1) [46] and Partington syndrome [47]. Similarly, *TBX1*, which harbors

4    mutations that cause the same phenotype as 22q11.2 deletion syndrome [48], is dark by depth in only

5    approximately 30% of gnomAD samples.

6

7    **Long-read technologies resolve many camouflaged regions, with variable success**

8    We selected three camouflaged gene regions to highlight common strengths and differences for how well

9    each long-read sequencing technology addresses the camouflaged region, including *SMN1* and *SMN2* (Figure

10    7a), *HSPA1A* and *HSPA1B* (Figure 7b), and *CR1* (Figure 7c). The *SMN1* and *SMN2* genes are camouflaged by

11    each other (gene duplication), as are *HSPA1A* and *HSPA1B*. *CR1*, however, is a special case, where it is

12    camouflaged by a repeated region within itself. Only ONT appeared to be capable of fully addressing the

13    camouflaged region for all three genes. 10x Genomics also performed well under certain circumstances, such

14    as *SMN1* and *SMN2* (regions where the duplication is >50kb away), but did not perform well for *HSPA1A* and

15    *HSPA1B*. PacBio performed well for *CR1* and *HSPA1A/HSPA1B*, but did not perform as well as ONT in the

16    *SMN1/SMN2* region.

17

18    *SMN1* and *SMN2* were 89.9% and 88.2% dark CDS, respectively (Figure 7a), based on standard Illumina

19    sequencing with 100-nucleotide read lengths, and were 84.0% and 83.1% dark CDS based on Illumina 250-

20    nucleotide read lengths (not shown). Both genes were technically 0% dark CDS based on 10x Genomics,

21    PacBio, and ONT data (Figure 7a). PacBio coverage does drop significantly throughout both genes, however.

22

23    *HSPA1A* and *HSPA1B* were 52.8% and 51.1% dark CDS (Figure 7b), respectively, based on standard Illumina

24    100-nucleotide read lengths, and were 50.2% and 49.5% dark CDS based on Illumina 250-nucleotide read

1    lengths (not shown). Both genes were 0% dark CDS based on ONT and PacBio data, and were 45.8% and 51.8%

2    dark CDS based on 10x Genomics data (Figure 7b). In contrast to the results for *SMN1* and *SMN2*, both ONT

3    and PacBio had consistent coverage throughout the camouflaged regions, whereas the camouflaged regions

4    remained dark for 10x Genomics (Figure 7b).

5

6    *CR1* was 26.5% dark CDS based on Illumina 100-nucleotide read lengths (Figure 7c), and was 24.5% dark based

7    on Illumina 250-nucleotide read lengths (not shown). *CR1* was 26.2% dark CDS for 10x Genomics, and 0% for

8    both ONT and PacBio (Figure 7c). While both PacBio and ONT were able fill the camouflaged region, coverage

9    drops dramatically throughout the region, despite both genomes being sequenced at 50x and 46x median

10   depth, which does not presently represent average use case for these technologies. It is likely that the

11   performance of these long-read platforms will be better with longer average sequencing libraries (e.g. >50kb

12   fragment sizes).

13

14   **Many camouflaged regions can be rescued, including in standard short-read sequencing data**

15   There are many large-scale whole-genome or whole-exome sequencing projects across tens of thousands of

16   individuals that are either completed or underway for a variety of diseases, including cancer (e.g., The Cancer

17   Genome Atlas; TCGA), autism spectrum disorder (e.g., The Autism Sequencing Consortium; ASC), Alzheimer's

18   disease (e.g., The Alzheimer's Disease Sequencing Project; ADSP), Parkinson's disease (e.g., The Parkinson's

19   Progression Markers Initiative; PPMI), and ALS (e.g., Target ALS and CReATe). All of these datasets are affected

20   by dark and camouflaged regions that may harbor mutations that are either driving or modify disease in

21   patients. Ideally, all samples would be re-sequenced using the latest technologies over time, but financial

22   resources and biological samples are limited, making it essential to maximize the utility of existing data.

23

1    Using a strategy similar to that proposed by Robert and Watson [10], we have developed a method to rescue

2    mutations in most camouflaged regions, including for standard Illumina short-read sequencing data. When

3    confronted with a sequencing read that aligns to two or more regions equally well (with high confidence),

4    most aligners (e.g., BWA [11–13]) will randomly assign the read to one of the regions and assign a low

5    mapping quality (MAPQ = 0 for BWA, or MAPQ = 1 for novoalign). Because the reads are already aligned to

6    one of the regions, we can use the following steps to rescue mutations in most camouflaged regions (Figure

7    8): (1) extract reads from camouflaged regions; (2) mask all highly similar regions in the reference genome,

8    except one, and re-align the extracted reads; (3) call mutations using standard methods. Without competing

9    camouflaged regions to confuse the aligner, the aligner will assign a high mapping quality, allowing variant

10    callers to behave normally. This will enable researchers to identify mutations that exist in one of the

11    camouflaged regions, but not which specific region (Figure 8). After rescuing these mutations, researchers can

12    then perform association studies to determine whether any of the mutations may be implicated in disease,

13    and follow up with targeted sequencing methods to determine the exact camouflage region a mutation lies in.

14

15    **Re-alignment rescues approximately 4622 exonic variants, including a rare ten-nucleotide frameshift**

16    **deletion in *CR1***

17    As a proof of principle, we applied our method to the Alzheimer's Disease Sequencing Project (ADSP) case-

18    control data [49] to approximate the number of potential mutations our approach could rescue. The ADSP is a

19    large sequencing project organized, in part, to identify functional mutations that influence Alzheimer's disease

20    development. Across 13142 samples from the ADSP, excluding all variants with a quality by depth (QD) <2.5,

21    we were able to rescue approximately 4622 exonic variants with a transition-transversion ration (Ti/Tv) of 1.97

22    from 147 camouflaged region sets, that are spread across 501 camouflaged genes (Supplemental Figure 14;

23    VCF will be provided to the ADSP). Using a more stringent QD (excluding variants with QD <5), we rescued

16

1    3152 variants with a Ti/Tv ratio of 2.17. We only included camouflaged regions from CDS exons for all genes,

2    including those that are <5% dark CDS.

3

4    Because *CR1* is a top-10 Alzheimer's disease gene, we then specifically interrogated it using our method

5    (Figure 8) for any functional mutations that could be involved in Alzheimer's disease, and identified a rare ten-

6    nucleotide frameshift deletion that is only found in five cases and zero controls, all of which are heterozygous

7    (Figure 8d). Thus, the estimated minor allele frequency for this mutation is 5 / (13142 * 2) = 0.00019, making it

8    more rare than the *TREM2* R47H allele [50–52]. For interest, only one of the individuals carried a single

9    *APOEε4* allele (ε3/ε4). The other four individuals were homozygous for *APOEε3* (ε3/ε3). We were able to

10    determine that the frameshift deletion is in one of exons 10, 18, or 26. Briefly, our method extracts all reads

11    with a low mapping quality (MAPQ < 10) from all three exons, masks all but one of the camouflaged regions

12    within each set of camouflaged regions, and aligns all reads from each set to only one of the regions (Figure 8).

13    Without identical competing regions to confuse the aligner, the mapping qualities are high enough for a

14    variant caller (e.g., GATK HaplotypeCaller) to identify whether a mutation exists. For example, reads harboring

15    the ten-nucleotide frameshift mutation were originally randomly scattered across exons 10, 18, and 26 from

16    the original alignment (Figure 8). We masked exons 18 and 26, leaving exon 10 unmasked; this allowed reads

17    from each of the three exons to align to only exon 10, so we could perform variant calling. We estimate a

18    cohort of approximately 70000 cases and controls would have approximately 80% statistical power to formally

19    assess this mutation's involvement in Alzheimer's disease, assuming a Relative Risk (RR) of 3.3, at an alpha of

20    0.0001. We provide the .bed files in GRCh37 and GRCh38, along with scripts that will enable researchers to

21    perform similar analyses in any sequencing dataset at

22    https://github.com/mebbert/Dark_and_Camouflaged_genes.

23

24    **Discussion**

1    While researchers have known for years that dark regions exist in standard short-read sequencing data, little

2    work has been done to characterize the breadth of the issue, and to develop possible solutions until more

3    financially-feasible long-read sequencing options are available. Short-read sequencing is unable to adequately

4    address camouflaged regions because the reads cannot fully span camouflaged regions to properly align

5    homologous nucleotides. Long-read sequencing technologies, such as those from 10x Genomics (synthetic

6    long reads), Oxford Nanopore Technologies (ONT), and Pacific Biosciences (PacBio) have the potential to

7    address many camouflaged regions because these technologies have median read lengths measured in

8    thousands of nucleotides, rather than only 100-300 nucleotides from standard short-read sequencing

9    technologies (e.g., Illumina). Recent work has even demonstrated that mappable ONT reads can exceed two

10   million nucleotides (e.g, 2272580) [53, 54], showing future potential for addressing large camouflaged regions.

11

12   In this study, we systematically characterized dark and camouflaged gene regions and proposed a method to

13   address most camouflaged regions in long- or short-read sequencing data. Our solution is specifically

14   applicable to camouflaged regions, not regions that are dark by depth, simply because there are no reads

15   available in regions that are dark by depth. While our solution is conceptually simple, implementing the

16   solution systematically was challenging because of many intricate details, including increased zygosity, and

17   would ideally be integrated into the original alignment and variant-calling process. While the original

18   implementation was challenging, we provide the resulting .bed files for both GRCh37 and GRCh38 that are

19   necessary to rescue mutations from camouflaged regions in any human re-sequencing dataset

20   (https://github.com/mebbert/Dark_and_Camouflaged_genes). We also provide all of our data and source

21   code. The .bed files and source code should make implementing our method relatively straightforward for

22   other groups. As a proof of concept, we were able to rescue approximately 4622 variants in the ADSP dataset

23   from 147 sets of camouflaged gene regions, which are spread across 501 camouflaged genes. Included in

1    these rescued mutations is a ten-nucleotide frameshift deletion in *CR1* found in five ADSP cases and zero

2    controls.

3

4    The number of genes affected by dark and camouflaged regions was surprisingly high. We identified 37873

5    total dark regions across 5857 gene bodies, nearly 4000 of which were protein coding genes. Exactly 28751 of

6    the dark regions were intronic and 2657 were in protein-coding exons (CDS). Others were in pseudogenes

7    (1134) and lincRNAs (732). While most of the dark regions were non-coding (e.g., intronic), these regions may

8    still harbor important mutations that drive or modify human diseases. For example, there are many examples

9    of mutations in non-coding regions driving disease, including repeat expansions [1, 55–62], splice-site

10   mutations (these may be intronic or exonic) [63–77], and regulatory mutations (e.g., UTR regions) [78–87].

11   There are also many lincRNAs associated with disease [88–97].

12

13   There are many patients with diseases known to be genetically inherited that remain genetically unexplained

14   because the patients do not have any of the known mutations. Many of the genes we identified as being at

15   least partially dark are known to be involved in numerous diseases, including Alzheimer's disease, ALS, SMA,

16   hemophilia A, autism spectrum disorder, schizophrenia, and others; functional mutations that modify disease

17   likely lie in some of these dark and camouflaged regions. For example, *SMN1* and *SMN2* are mostly dark

18   (camouflaged) and are known to harbor mutations that cause disease [63, 65–67]. *CR1* is another dark gene

19   that is 26.5% dark CDS, being camouflaged to itself, and is strongly implicated in Alzheimer's disease. In fact,

20   the *CR1* camouflaged region includes the C3b and C4b protein binding sites, repeated several times.

21   Interestingly, the *C4B* gene (encodes the C4b protein) is also 72.8% dark CDS (camouflaged) and may be

22   involved in disease [98, 99]. We are confident that rescuing mutations from camouflaged regions will have a

23   meaningful impact on disease research, and may explain some of the missing heritability of Alzheimer's

24   disease [18, 100–102] and other diseases.

19

1

2    A large number of gene bodies (494) were 100% dark, which means they are entirely overlooked in standard

3    whole-exome, whole-genome, and RNA sequencing studies [10]. Additionally, more than 1500 gene bodies, or

4    nearly 27%, were at least 25% dark and more than 2000 (34.9%) were at least 5% dark; of these, 628 protein-

5    coding genes were at least 5% dark within CDS regions. Understanding what role these genes play in human

6    health and disease will require being able to resolve them in DNA and RNA sequencing experiments.

7

8    A critical decision for future large-scale sequencing projects will be regarding which long-read technology is

9    ideal to maximize the probability of identifying functional mutations driving disease. Unfortunately, the

10    answer is not clear, as each technology has its pros and cons. Based on our results, the ONT platform

11    performed best, overall, resolving 71.4% of dark gene-body regions. Current costs will likely be prohibitive for

12    large studies, however. The 10x Genomics platform resolved 66.3% of dark gene-body regions, when

13    compared to standard Illumina sequencing. PacBio resolved 49.0% of dark gene-body regions. Even increasing

14    Illumina read lengths from 100 to 250 made a sizeable difference, overall, resolving 21.1% of dark gene-body

15    regions. Both the PacBio and ONT data used in this study had shorter median read lengths than expected,

16    suggesting both technologies can likely perform better than our estimates.

17

18    Focusing only on CDS regions, there were 2757 dark CDS regions across 744 protein-coding genes, based on

19    Illumina 100-nucleotide read lengths. ONT outperformed other long-read technologies, resolving 81.8% of

20    dark CDS regions. PacBio and 10x Genomics resolved 66.6% and 54.9%, respectively. We found that 10x

21    Genomics performed well in the *SMN1* and *SMN2* genes (Figure 7), attaining consistently deep, high-quality

22    coverage throughout. Both ONT and PacBio coverage declined in the interior regions of the genes. In other

23    cases, such as *CR1* and *NEB*, 10x Genomics was unable to improve on standard Illumina sequencing, but

24    PacBio and ONT were able to largely resolve the region—albeit requiring higher than normal sequencing

1    depth. We believe that 10x Genomics can correct the issues we observed in *CR1* and *NEB*, by implementing a

2    more sophisticated version of our method that also incorporates evidence from their synthetic long-read

3    technology.

4

5    Whether each technology is able to reliably resolve dark and camouflaged regions is an important

6    consideration for choosing the best long-read technology, but we should also consider how reliably each

7    technology is able to resolve structural mutations. In a previous study, we tested how well ONT and PacBio are

8    able to traverse challenging repeat expansions, and whether they are amenable to genetic discovery [1]. We

9    found that both technologies are well-suited, but we have not assessed performance of the 10x Genomics

10    platform across long repeat expansions.

11

12    The primary challenge with ONT and PacBio long-read sequencing is, of course, the high error rate, which can

13    be overcome through deeper sequencing because errors in ONT and PacBio sequencing are mostly random

14    [103, 104]. Ultimately, we are confident that, as long-read error rates improve, and costs continue to decline,

15    long-read technologies will be the preferred sequencing choice for large-scale sequencing projects, especially

16    when considering structural mutations.

17

18    We identified dark and camouflaged regions in this study by averaging data across ten males with deep

19    Illumina whole-genome sequencing, using 100-nucleotide read lengths. We assessed how well long-read

20    sequencing technologies (PacBio, ONT, and 10X genomics) resolve these regions, but our measurements

21    should only be considered estimates. While long-read sequencing technologies are becoming more common,

22    we were unable to find more than one male individual for each long-read technology; we needed male

23    samples to assess all chromosomes, including the Y chromosome. Additionally, the samples we used for each

24    long-read technology were sequenced at a much higher depth than is currently typical for a re-sequencing

1    effort, which is likely over estimating the number of dark regions they resolve for the average use case. Our

2    measurements should be a reasonable estimate of reality, however, and future analyses will be able to refine

3    our estimates.

4

5    We used whole-genome sequencing to assess dark and camouflaged regions, but this problem is magnified in

6    whole-exome data, which many large-scale sequencing studies are based on, either completely, or in part.

7    Whole-exome data are typically generated using even shorter read lengths. They are also generally based on

8    capture, which means certain exons are not fully represented. *APOE* is a prime example, where it is typically

9    well-covered in whole-genome data, but a portion is dark in whole-exome data (Supplemental Figure 13).

10    With *APOE* harboring the largest genetic risk factors for Alzheimer's disease, it is important to properly

11    characterize the entire gene.

12

13    In this study, we characterized dark and camouflaged gene bodies, and demonstrated several disease-relevant

14    genes where a significant portion is dark in standard short-read sequencing data, including *SMN1* and *SMN2*,

15    *CR1*, and sometimes even *APOE*. We also identified a rare ten-nucleotide frameshift deletion in *CR1* that is

16    found only in five ADSP cases and zero controls, as a proof of principle (Figure 8d). Using our method (Figure

17    8), we were able to determine that the frameshift deletion is in one of exons 10, 18, or 26. With *CR1* being a

18    top Alzheimer's disease gene without any known functional mutations, we believe it will be important to

19    assess this mutation in a large cohort, to determine whether it plays a role in disease development and

20    progression. We have also proposed a solution to address most camouflaged genes in sequencing data, and

21    believe that our approach has the potential to identify functional mutations that are influencing development

22    across a range of diseases, but are currently entirely overlooked by standard short-read sequencing

23    approaches.

24

# Conclusion

There remain thousands of potentially important genomic regions that are overlooked with short-read sequencing, but are largely resolved by long-read technologies. While these regions represent only a small portion of the entire genome or exome, many of these regions are known to be important in human health and disease. Equally important, however, is that the impact of many other genes is entirely unknown because they are 100% dark. We presented a method that can resolve most camouflaged regions that we believe will help researchers identify mutations that are involved in disease. As a proof of principle, we rescued approximately 4622 variants in the ADSP dataset, including a ten-nucleotide frameshift mutation in *CR1*. While we cannot formally assess the *CR1* frameshift mutation in Alzheimer's disease (insufficient sample-size), we believe it is worth investigating in a larger cohort. In the long-term, we believe long-read sequencing technologies will be the best solution for resolving dark and camouflaged regions.

# Methods

### Sample selection and preparation

To identify dark and camouflaged regions, and to assess how well other technologies address them, we selected samples from each technology and read length. All samples were aligned to hg19/GRCh37. To assess dark and camouflaged regions in standard Illumina sequencing with 100-nucleotide read lengths, we selected ten unrelated male control samples from the Alzheimer's Disease Sequencing Project (ADSP) where deep whole-genome sequencing had been performed by randomly selecting one male from ten random families. All ten males were from either the "Health/Medical/Biomedical" (HMB-IRB) or "Health/Medical/Biomedical" for non-profit organizations (HMB-IRB-NPU) consent groups, indicated as groups C1 and C2 in the ADSP pedigree files (available through dbGAP). We selected samples from the ADSP because we required samples that met the following criteria: (1) had been sequenced using standard paired-end Illumina sequencing with 100-

23

1   nucleotide read lengths, (2) had been sequenced with a median depth >30x, and (3) were publicly available.

2   Median genome-wide read depths ranged from 35.4x to 42.9x, with a median of 39.4x. Samples were

3   prepared and sequenced as part of the ADSP [49]. These samples were aligned using BWA (v0.5.9). We could

4   not find samples from the 1000 Genomes Project [24] that met these criteria; sequencing depths were either

5   too shallow, or read lengths were too long or short. The ADSP sample IDs we used were: A-CUHS-CU000406,

6   A-CUHS-CU002997, A-CUHS-CU000779, A-CUHS-CU000208, A-CUHS-CU001010, A-CUHS-CU002031, A-CUHS-

7   CU002707, A-CUHS-CU003023, A-CUHS-CU003090, A-CUHS-CU003128.

8

9   To assess dark and camouflaged regions in samples sequenced using Illumina 250-nucleotide read lengths, we

10  selected ten samples from the 1000 Genomes Project that had been sequenced with 250-nucleotide read

11  lengths, and had a median depth >30x. All ten samples were aligned using BWA (v 0.7.5a-r428) [2, 11–13].

12  Median genome-wide read depths ranged from 39.3 to 52.6, with a median of 48.9x. Sample IDs for the

13  Illumina 250-nucleotide read lengths were: NA20845

14  (ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/data/NA20845/high_coverage_alignment/), HG01112

15  (ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/data/HG01112/high_coverage_alignment/), HG01583

16  (ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/data/HG01583/high_coverage_alignment/), HG01051

17  (ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/data/HG01051/high_coverage_alignment/), HG03742

18  (ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/data/HG03742/high_coverage_alignment/), HG00096

19  (ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/data/HG00096/high_coverage_alignment/), HG01565

20  (ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/data/HG01565/high_coverage_alignment/), HG01879

21  (ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/data/HG01879/high_coverage_alignment/), HG01500

22  (ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/data/HG01500/high_coverage_alignment/), and HG03006

23  (ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/data/HG03006/high_coverage_alignment/).

24

1    We also selected samples generated using the 10x Genomics synthetic long-read sequencing platform, and

2    ONT and PacBio long-read sequencing platforms that were either prepared by, and publicly available from the

3    respective company, or prepared using standard practice. Specifically, we downloaded HG00512 raw FASTQ

4    data from 10x Genomics (https://support.10xgenomics.com/de-novo-assembly/datasets/1.1.0/msHG00512;

5    http://s3-us-west-2.amazonaws.com/10x.files/samples/assembly/2.1.0/chi/chi_fastqs.tar) and aligned it

6    according to 10x Genomics' standard practices. We used longranger (v2.2.2) and aligned to GRCh37

7    (longranger wgs --id HG00512 --description="Han Chinese" --sex="male" --

8    fastqs=chi/HNKHFCCXX/,chi/HWHFTCCXX/ --reference="10x-b37-2.1.0/" --jobmode=sge --mempercore=125 –

9    downsample=385). Median depth for HG00512 was 52x, after downsampling. For ONT, we downloaded the

10   final Cliveome v2 from ONT's official GitHub page (http://cliveo.me/; https://github.com/nanoporetech/ONT-

11   HG1/blob/master/CONTENTS.md), which was prepared by ONT. Cliveome v2 was sequenced to a median

12   depth of 36x. To increase the median read depth to more closely match those of other technologies, we

13   merged reads from HG002 (https://ftp-

14   trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/HG002_NA24385_son/Ultralong_OxfordNanopore/com

15   bined_2018-08-10/HG002_ONTrel2_16x_RG_HP10xtrioRTG.cram) [105, 106] and aligned using minimap2

16   [107] (ALIGN_OPTS="x map-pb -a --eqx -L -O 5,56 -E 4,1 -B 5 --secondary=no -z 400,50 -r 2k -Y";

17   REF=g1kv37/g1kv37.fa; minimap2 -d ${REF}.mmi ${ALIGN_OPTS} ${REF}; minimap2 ${ALIGN_OPTS} -a

18   ${REF}.mmi <reads.fq> | samtools view -T {REF} -F 2308 > output_file). The merged sample had 46x median

19   depth. We used the same alignment options recommended for PacBio because we found the recommended

20   'map-ont' option in minimap2 performed substantially worse. We used PacBio data generated from HG005

21   (ftp://ftp-

22   trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/HG005_NA24631_son/MtSinai_PacBio/PacBio_minimap2_b

23   am/) [105], which was sequenced to a median depth of 50x and aligned using minimap2 [107] (pbsv fasta

24   [movie].subreads.bam | minimap2 -t 8 -x map-pb -a --eqx -L -O 5,56 -E 4,1 -B 5 --secondary=no -z 400,50 -r 2k -

1    Y

2    ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/reference/phase2_reference_assembly_sequence/hs37d

3    5.fa.gz - | samtools sort > HG005_PacBio_GRCh37.bam). Neither the ONT nor the PacBio alignments include

4    secondary alignments.

5

6    **Identifying dark and camouflaged gene body regions**

7    To identify dark and camouflaged gene body regions in standard Illumina 100-nucleotide read length data, we

8    first scanned all ten ADSP whole-genome sequence samples for genomic positions that met either of the

9    following criteria: (1) had <5 reads, and (2) had ≥90% of reads with a mapping quality (MAPQ) <10. We then

10   averaged the depth and count of low MAPQ reads across all samples for each position. We used strict cutoffs

11   to identify regions that are clearly dark, but there are many additional regions that fall just beyond our

12   thresholds. This analysis was performed using the Dark Region Finder (DRF;

13   https://github.com/mebbert/DarkRegionFinder; mapq=9; dark_mass=90; camo_mass=50; dark_depth=5; java

14   -jar -Xmx20g CamoGeneFinder.jar -i <sample>.bam --human-ref genome.fa --min-region-size 1 --camo-mapq-

15   threshold $mapq --min-dark-mapq-mass $dark_mass --min-camo-mapq-mass $camo_mass --dark-depth

16   $dark_depth --camo-bed-output <sample>-camo-dark_depth_${dark_depth}-dark_mass_${dark_mass}-

17   camo_mass_${camo_mass}-mapq_${mapq}.b37.bed --dark-bed-output <sample>-dark-

18   dark_depth_${dark_depth}-dark_mass_${dark_mass}.b37.bed --incomplete-bed-output <sample>-

19   incomplete.b37.bed). Any position that met either criteria was considered dark and categorized as either dark

20   by depth or dark by mapping quality. We then limited the dark regions to gene bodies by intersecting dark

21   regions identified by Dark Region Finder with Ensembl's GRCh37 build 87 gene annotations. We converted the

22   transcript-level annotations to gene-level annotations using bedtools [108] and custom scripts that are

23   available. Any dark region that spanned a gene body element region (e.g., intron-exon boundary) was split into

24   two separate dark regions so we could estimate the number of dark bases in each type of gene body region

1    (e.g., introns, exons, UTRs, etc.). For most analyses, we only included dark regions with ≥20 contiguous bases.

2    The only exception is for Supplemental Tables 1, 3, 5, 7, 9, 12, and 14, where we calculate total percentage of

3    each gene body that is dark, in which we include all dark positions. To identify camouflaged regions,

4    specifically, we used BLAT [26] to identify all genomic regions that were highly similar to any given gene body

5    region that was dark by mapping quality. Any region that was ≥98% identical (-minIdentity = 98), and that was

6    considered dark (≥90% of reads with MAPQ <10), was considered a match. We generated .bed files for

7    GRCh37 using this method. We also converted the GRCh37 .bed file to GRCh38 using a custom script, based

8    off the Ensembl build 87 GRCh38 gene annotations. All code and .bed files can be found at

9    https://github.com/mebbert/Dark_and_Camouflaged_genes.

10

11   **Statistics**

12   We quantified the percentage of each gene body that was dark by summing the total number of dark bases in

13   the gene (i.e., between the 5'UTR to the 3'UTR start and end, respectively) and dividing by the total number of

14   bases in the gene. We similarly calculated the percentage of intronic, exonic (including CDS and UTR), and only

15   CDS exons by dividing the total number of dark bases in each category within the gene by the total number of

16   bases within that category. We performed these calculations for data based on Illumina 100-nucleotide reads

17   for all dark regions combined (Supplemental Tables 1-2), dark by depth only (Supplemental Tables 14-15), dark

18   by mapping quality (Supplemental Tables 16-17), and only camouflaged regions (Supplemental Tables 12-13).

19   We performed identical calculations for the samples from Illumina 250-nucleotide read length data, 10x

20   Genomics, ONT, and PacBio (Supplemental Tables 3-10, 18-41). We identified diseases that were known to be

21   associated with genes that are at least 5% dark CDS by searching for mutations in the Human Gene Mutation

22   Database (HGMD) [30].

23

1    Coverage plots from gnomAD data were obtained from gnomAD-old.broadinstitute.org [33]. We used the old

2    version because the current version of gnomAD (accessed December 2018) does not allow the user to view

3    median read depths, nor the percentage of samples with greater than a given coverage depth. Sequence

4    pileups in representative samples were generated using the Integrative Genomics Viewer (IGV) [109], where

5    reads with a MAPQ < 10 were filtered, and insertions, deletions, and mismatches were not shown. Karyotype

6    plots showing genomic locations for dark and camouflaged regions were generated using KaryotypeR (v1.6.2)

7    [110] in R (v3.5.1). Bar plots were made using ggplot2 (v3.0.0). Pathway analyses and resulting plots were

8    generated using Metascape (accessed December 2018) [111]. Word clouds were generated at

9    wordclouds.com. Gene schematics were generated using the Gene Structure Display Server (GSDS; v2) [112].

10

11    We performed an enrichment analysis to assess whether genes that are ≥5% dark CDS are enriched for specific

12    diseases. Because we identified 75 genes that have a known mutation associated with disease, and that are

13    ≥5% dark CDS, we randomly selected 75 genes from the with known HGMD mutations and measured the

14    number of genes with known mutation associated with each disease. We repeated this process 10000 times

15    and used the following metric as our enrichment score: -10*log10(empirical_pvalue), rounded to the nearest

16    whole number.

17

18    **Screening ADSP for functional *CR1* mutations in camouflaged region**

19    After discovering that more than 25% of the *CR1* gene's CDS is camouflaged, we screened all ADSP samples for

20    rare functional mutations that could play a role in Alzheimer's disease development and progression by

21    applying our proposed method (Figure 8). To apply our method, we extracted all reads with a mapping quality

22    (MAPQ) <10 from each camouflaged region within *CR1*, and from each of the respective camouflage mate

23    regions, using samtools and the GRCh37 .bed file we generated that identifies all camouflaged regions. An

24    example of camouflaged mate regions in *CR1* includes exons 10, 18, and 26, which are identical in the

28

1    reference genome (Figure 8). As previously mentioned, *CR1* is a special case that is camouflaged by regions

2    duplicated within itself, rather than being camouflaged by a different gene; thus, we knew that any mutations

3    we discovered would be from *CR1*. Our approach works the same regardless of whether a gene is camouflaged

4    by itself or another gene, but we mention that *CR1* is camouflaged by itself, for interest. After extracting reads

5    from each camouflaged region, using the .bed file we provide, we then masked all camouflaged regions within

6    *CR1* in the reference genome, except for one from each set of camouflaged mates. For example, between

7    exons 10, 18, and 26, we masked exons 18 and 26 in the reference genome, allowing reads from all three

8    exons to align only to exon 10; without competing camouflaged regions to confuse the aligner, all reads from

9    exons 10, 18, and 26 mapped to exon 10 with high quality. Masking regions of the reference genome simply

10    means to change nucleotides to an unmappable character (usually 'N'), to prevent any reads from aligning to

11    that region.

12

13    After aligning all reads to a single region within each set of camouflaged regions, we were able to perform

14    standard variant calling using the GATK HaplotypeCaller [25], with one exception: instead of treating each

15    camouflaged region as diploid, we increased the ploidy setting in HaplotypeCaller according to the number of

16    copies within a given set of camouflaged regions. Referring again to our *CR1* example, because there are three

17    regions (exons 10, 18, and 26), we set the HaplotypeCaller ploidy to hexaploid. Increasing the ploidy is

18    essential for increased sensitivity, since the number of reads harboring a given variant—which only originate

19    from one of the camouflaged regions—will be overwhelmed by reads from the others, thus preventing the

20    variant caller from identifying the mutation under the assumption that the data are from a diploid region. In

21    other words, if a mutation exists in exon 26, we would expect only approximately 1/6$^{th}$ of reads from exons

22    10, 18, and 26 to harbor that mutation, rather than approximately 1/2. Because the ADSP is mostly exome

23    data, we limited HaplotypeCaller to CDS exons only. According to the current ADSP phenotype data, one of the

1    samples harboring the *CR1* frameshift mutation is a control. The individual has since been officially diagnosed

2    with Alzheimer's disease, however.

3

4    # Abbreviations

5    MAPQ: mapping quality; CDS: coding sequence; ALS: amyotrophic lateral sclerosis; FTD: frontotemporal

6    dementia; PacBio: Pacific Biosciences; ONT: Oxford Nanopore Technologies; ADSP: Alzheimer's Disease

7    Sequencing Project;

8

9    # Declarations

10    **Ethics approval and consent to participate**

11    The Mayo Clinic Institutional Review Board (IRB) approved all procedures for this study and we followed all

12    appropriate protocols.

13

14    **Consent for publication**

15    All participants were properly consented for this study.

16

17    **Availability of data and materials**

18    The ADSP dataset supporting the conclusions of this article (including the whole-genome and whole-exome

19    data) are available in the National Institute on Aging Genetics of Alzheimer's Disease Storage (NIAGADS) site,

20    and may be requested therein: https://www.niagads.org/adsp/. Public links to the high-coverage whole-

21    genome data from the 1000 Genomes Project (illumina 250bp read lengths) are listed in the Methods. Raw

22    data from the 10x Genomics sample (HG00512) used within this article was downloaded directly from the 10x

23    Genomics website at: https://support.10xgenomics.com/de-novo-assembly/datasets/2.1.0/chi. The Cliveome2

1  data was downloaded from the official Oxford Nanopore Technologies GitHub page:

2  https://github.com/nanoporetech/ONT-HG1/blob/master/. The PacBio data used in this publication (HG005)

3  was downloaded from ftp://ftp-

4  trace.ncbi.nlm.nih.gov/giab/ftp/data/ChineseTrio/HG005_NA24631_son/MtSinai_PacBio/PacBio_minimap2_b

5  am/HG005_PacBio_GRCh37.bam.

6

7  All scripts are available at: https://github.com/mebbert/Dark_and_Camouflaged_genes.

8

9  **Competing interests**

10  All authors declare they have no conflicts of interest.

11

12  **Funding**

22

23  **Authors' contributions**

1    ME, LP, and JF developed and designed the study, and wrote the manuscript. ME and TJ performed all

2    analyses. JR, SY, NT, YA, VB, EL, DK, PC, LP, PR, JK, MC, OA, and RR contributed important intellectual ideas and

3    feedback. SY, YA, NT, OA, and RR helped obtain data. KW and JS performed experiments. EL, DK, and PC

4    provided samples. All authors read and approved the final manuscript.

5

34

3

4

5

## References

1. Ebbert MTW, Farrugia SL, Sens JP, Jansen-West K, Gendron TF, Prudencio M, et al. Long-read sequencing across the C9orf72 "GGGGCC" repeat expansion: implications for clinical use and genetic discovery efforts in human disease. Mol Neurodegener. 2018;13:46. doi:10.1186/s13024-018-0274-4.

2. Zheng-Bradley X, Streeter I, Fairley S, Richardson D, Clarke L, Flicek P, et al. Alignment of 1000 Genomes Project reads to reference assembly GRCh38. Gigascience. 2017;6:1–8. doi:10.1093/gigascience/gix038.

3. Callaway E. Human brain shaped by duplicate genes. Nature. 2012. doi:10.1038/nature.2012.10584.

4. Charrier C, Joshi K, Coutinho-Budd J, Kim J-E, Lambert N, de Marchena J, et al. Inhibition of SRGAP2 function by its human-specific paralogs induces neoteny during spine maturation. Cell. 2012;149:923–35. doi:10.1016/j.cell.2012.03.034.

5. Dennis MY, Nuttle X, Sudmant PH, Antonacci F, Graves TA, Nefedov M, et al. Evolution of human-specific neural SRGAP2 genes by incomplete segmental duplication. Cell. 2012;149:912–22. doi:10.1016/j.cell.2012.03.033.

6. Karlin S, Brocchieri L. Heat shock protein 60 sequence comparisons: duplications, lateral transfer, and mitochondrial evolution. Proc Natl Acad Sci USA. 2000;97:11348–53. doi:10.1073/pnas.97.21.11348.

7. Lin Y, Cheng Y, Jin J, Jin X, Jiang H, Yan H, et al. Genome duplication and gene loss affect the evolution of heat shock transcription factor genes in legumes. PLoS One. 2014;9:e102825. doi:10.1371/journal.pone.0102825.

8. Nguyen AD, Gotelli NJ, Cahan SH. The evolution of heat shock protein sequences, cis-regulatory elements, and expression profiles in the eusocial Hymenoptera. BMC Evol Biol. 2016;16:15. doi:10.1186/s12862-015-0573-0.

9. Sørensen JG, Kristensen TN, Loeschcke V. The evolutionary and ecological role of heat shock proteins. Ecol Lett. 2003;6:1025–37. doi:10.1046/j.1461-0248.2003.00528.x.

10. Robert C, Watson M. Errors in RNA-Seq quantification affect genes of relevance to human disease. Genome Biol. 2015;16:177. doi:10.1186/s13059-015-0734-x.

11. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. Bioinformatics. 2009;25:1754–60. doi:10.1093/bioinformatics/btp324.

12. Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler transform. Bioinformatics. 2010;26:589–95. doi:10.1093/bioinformatics/btp698.

13. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. arXiv. 2013.

14. Lambert J-C, Heath S, Even G, Campion D, Sleegers K, Hiltunen M, et al. Genome-wide association study identifies variants at CLU and CR1 associated with Alzheimer's disease. Nat Genet. 2009;41:1094–9. doi:10.1038/ng.439.

15. Hollingworth P, Harold D, Sims R, Gerrish A, Lambert J-C, Carrasquillo MM, et al. Common variants at ABCA7, MS4A6A/MS4A4E, EPHA1, CD33 and CD2AP are associated with Alzheimer's disease. Nat Genet. 2011;43:429–35. doi:10.1038/ng.803.

16. Lambert JC, Ibrahim-Verbaas CA, Harold D, Naj AC, Sims R, Bellenguez C, et al. Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for Alzheimer's disease. Nat Genet. 2013;45:1452–8. doi:10.1038/ng.2802.

17. Kauwe JSK, Cruchaga C, Karch CM, Sadler B, Lee M, Mayo K, et al. Fine mapping of genetic variants in BIN1, CLU, CR1 and PICALM for association with cerebrospinal fluid biomarkers for Alzheimer's disease.

PLoS One. 2011;6:e15918. doi:10.1371/journal.pone.0015918.

18. Ridge PG, Hoyt KB, Boehme K, Mukherjee S, Crane PK, Haines JL, et al. Assessment of the genetic variance of late-onset Alzheimer's disease. Neurobiol Aging. 2016;41:200.e13-200.e20. doi:10.1016/j.neurobiolaging.2016.02.024.

19. Ebbert MTW, Ridge PG, Wilson AR, Sharp AR, Bailey M, Norton MC, et al. Population-based analysis of Alzheimer's disease risk alleles implicates genetic interactions. Biol Psychiatry. 2014;75:732–7. doi:10.1016/j.biopsych.2013.07.008.

20. Ridge PG, Ebbert MTW, Kauwe JSK. Genetics of Alzheimer's disease. Biomed Res Int. 2013;2013:254954. doi:10.1155/2013/254954.

21. Mahmoudi R, Feldman S, Kisserli A, Duret V, Tabary T, Bertholon L-A, et al. Inherited and Acquired Decrease in Complement Receptor 1 (CR1) Density on Red Blood Cells Associated with High Levels of Soluble CR1 in Alzheimer's Disease. Int J Mol Sci. 2018;19. doi:10.3390/ijms19082175.

22. Naj AC, Jun G, Beecham GW, Wang L-S, Vardarajan BN, Buros J, et al. Common variants at MS4A4/MS4A6E, CD2AP, CD33 and EPHA1 are associated with late-onset Alzheimer's disease. Nat Genet. 2011;43:436–41. doi:10.1038/ng.801.

23. Zerbino DR, Achuthan P, Akanni W, Amode MR, Barrell D, Bhai J, et al. Ensembl 2018. Nucleic Acids Res. 2018;46:D754–61. doi:10.1093/nar/gkx1098.

24. 1000 Genomes Project Consortium, Abecasis GR, Auton A, Brooks LD, DePristo MA, Durbin RM, et al. An integrated map of genetic variation from 1,092 human genomes. Nature. 2012;491:56–65. doi:10.1038/nature11632.

25. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. Genome Res. 2010;20:1297–303. doi:10.1101/gr.107524.110.

26. Kent WJ. BLAT--the BLAST-like alignment tool. Genome Res. 2002;12:656–64. doi:10.1101/gr.229202.

27. Frankish A, Diekhans M, Ferreira A-M, Johnson R, Jungreis I, Loveland J, et al. GENCODE reference annotation for the human and mouse genomes. Nucleic Acids Res. 2018. doi:10.1093/nar/gky955.

28. Bader GD, Hogue CWV. An automated method for finding molecular complexes in large protein interaction networks. BMC Bioinformatics. 2003;4:2.

29. Carithers LJ, Ardlie K, Barcus M, Branton PA, Britton A, Buia SA, et al. A Novel Approach to High-Quality Postmortem Tissue Procurement: The GTEx Project. Biopreserv Biobank. 2015;13:311–9. doi:10.1089/bio.2015.0032.

30. Stenson PD, Ball EV, Mort M, Phillips AD, Shiel JA, Thomas NST, et al. Human Gene Mutation Database (HGMD): 2003 update. Hum Mutat. 2003;21:577–81. doi:10.1002/humu.10212.

31. Seminary ER, Sison SL, Ebert AD. Modeling Protein Aggregation and the Heat Shock Response in ALS iPSC-Derived Motor Neurons. Front Neurosci. 2018;12:86. doi:10.3389/fnins.2018.00086.

32. Kalmar B, Lu C-H, Greensmith L. The role of heat shock proteins in Amyotrophic Lateral Sclerosis: The therapeutic potential of Arimoclomol. Pharmacol Ther. 2014;141:40–54. doi:10.1016/j.pharmthera.2013.08.003.

33. Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, et al. Analysis of protein-coding genetic variation in 60,706 humans. Nature. 2016;536:285–91. doi:10.1038/nature19057.

34. Roses AD. Apolipoprotein E alleles as risk factors in Alzheimer's disease. Annu Rev Med. 1996;47:387–400. doi:10.1146/annurev.med.47.1.387.

35. Roses AD, Saunders AM. APOE is a major susceptibility gene for Alzheimer's disease. Curr Opin Biotechnol. 1994;5:663–7. doi:10.1016/0958-1669(94)90091-4.

36. Strittmatter WJ, Saunders AM, Schmechel D, Pericak-Vance M, Enghild J, Salvesen GS, et al. Apolipoprotein E: high-avidity binding to beta-amyloid and increased frequency of type 4 allele in late-onset familial Alzheimer disease. Proc Natl Acad Sci USA. 1993;90:1977–81.

37. Blauw HM, Barnes CP, van Vught PWJ, van Rheenen W, Verheul M, Cuppen E, et al. SMN1 gene duplications are associated with sporadic ALS. Neurology. 2012;78:776–80. doi:10.1212/WNL.0b013e318249f697.

38. Corcia P, Camu W, Halimi JM, Vourc'h P, Antar C, Vedrine S, et al. SMN1 gene, but not SMN2, is a risk

factor for sporadic ALS. Neurology. 2006;67:1147–50. doi:10.1212/01.wnl.0000233830.85206.1e.

39. Rogers J, Cooper NR, Webster S, Schultz J, McGeer PL, Styren SD, et al. Complement activation by beta-amyloid in Alzheimer disease. Proc Natl Acad Sci USA. 1992;89:10016–20.

40. Rogers J, Li R, Mastroeni D, Grover A, Leonard B, Ahern G, et al. Peripheral clearance of amyloid beta peptide by complement C3-dependent adherence to erythrocytes. Neurobiol Aging. 2006;27:1733–9. doi:10.1016/j.neurobiolaging.2005.09.043.

41. Kisserli A, Tabary T, Cohen JHM, Duret V, Mahmoudi R. High-resolution Melting PCR for Complement Receptor 1 Length Polymorphism Genotyping: An Innovative Tool for Alzheimer's Disease Gene Susceptibility Assessment. J Vis Exp. 2017. doi:10.3791/56012.

42. Fonseca MI, Chu S, Pierce AL, Brubaker WD, Hauhart RE, Mastroeni D, et al. Analysis of the putative role of CR1 in alzheimer's disease: genetic association, expression and function. PLoS One. 2016;11:e0149792. doi:10.1371/journal.pone.0149792.

43. Brouwers N, Van Cauwenberghe C, Engelborghs S, Lambert JC, Bettens K, Le Bastard N, et al. Alzheimer risk associated with a copy number variation in the complement receptor 1 increasing C3b/C4b binding sites. Mol Psychiatry. 2012;17:223–33. doi:10.1038/mp.2011.24.

44. Kucukkilic E, Brookes K, Barber I, Guetta-Baranes T, ARUK Consortium, Morgan K, et al. Complement receptor 1 gene (CR1) intragenic duplication and risk of Alzheimer's disease. Hum Genet. 2018;137:305–14. doi:10.1007/s00439-018-1883-2.

45. Crane A, Brubaker WD, Johansson JU, Trigunaite A, Ceballos J, Bradt B, et al. Peripheral complement interactions with amyloid β peptide in Alzheimer's disease: 2. Relationship to amyloid β immunotherapy. Alzheimers Dement. 2018;14:243–52. doi:10.1016/j.jalz.2017.04.015.

46. Kato M, Saitoh S, Kamei A, Shiraishi H, Ueda Y, Akasaka M, et al. A longer polyalanine expansion mutation in the ARX gene causes early infantile epileptic encephalopathy with suppression-burst pattern (Ohtahara syndrome). Am J Hum Genet. 2007;81:361–6. doi:10.1086/518903.

47. Partington MW, Turner G, Boyle J, Gécz J. Three new families with X-linked mental retardation caused by the 428-451dup(24bp) mutation in ARX. Clin Genet. 2004;66:39–45. doi:10.1111/j.0009-9163.2004.00268.x.

48. Zweier C, Sticht H, Aydin-Yaylagül I, Campbell CE, Rauch A. Human TBX1 missense mutations cause gain of function resulting in the same phenotype as 22q11.2 deletions. Am J Hum Genet. 2007;80:510–7. doi:10.1086/511993.

49. Naj AC, Lin H, Vardarajan BN, White S, Lancour D, Ma Y, et al. Quality control and integration of genotypes from two calling pipelines for whole genome sequence data in the Alzheimer's disease sequencing project. Genomics. 2018. doi:10.1016/j.ygeno.2018.05.004.

50. Guerreiro R, Wojtas A, Bras J, Carrasquillo M, Rogaeva E, Majounie E, et al. TREM2 variants in Alzheimer's disease. N Engl J Med. 2013;368:117–27. doi:10.1056/NEJMoa1211851.

51. Jonsson T, Stefansson H, Steinberg S, Jonsdottir I, Jonsson PV, Snaedal J, et al. Variant of TREM2 associated with the risk of Alzheimer's disease. N Engl J Med. 2013;368:107–16. doi:10.1056/NEJMoa1211103.

52. Gonzalez Murcia JD, Schmutz C, Munger C, Perkes A, Gustin A, Peterson M, et al. Assessment of TREM2 rs75932628 association with Alzheimer's disease in a population-based sample: the Cache County Study. Neurobiol Aging. 2013;34:2889.e11-3. doi:10.1016/j.neurobiolaging.2013.06.004.

53. Payne A, Holmes N, Rakyan V, Loose M. Whale watching with BulkVis: A graphical viewer for Oxford Nanopore bulk fast5 files. BioRxiv. 2018. doi:10.1101/312256.

54. Payne A, Holmes N, Rakyan V, Loose M. BulkVis: a graphical viewer for Oxford nanopore bulk FAST5 files. Bioinformatics. 2018. doi:10.1093/bioinformatics/bty841.

55. La Spada AR, Taylor JP. Repeat expansion disease: progress and puzzles in disease pathogenesis. Nat Rev Genet. 2010;11:247–58. doi:10.1038/nrg2748.

56. Orr HT, Chung MY, Banfi S, Kwiatkowski TJ, Servadio A, Beaudet AL, et al. Expansion of an unstable trinucleotide CAG repeat in spinocerebellar ataxia type 1. Nat Genet. 1993;4:221–6. doi:10.1038/ng0793-221.

57. Lindblad K, Savontaus ML, Stevanin G, Holmberg M, Digre K, Zander C, et al. An expanded CAG repeat sequence in spinocerebellar ataxia type 7. Genome Res. 1996;6:965–71.

58. Squitieri F, Andrew SE, Goldberg YP, Kremer B, Spence N, Zeisler J, et al. DNA haplotype analysis of

Huntington disease reveals clues to the origins and mechanisms of CAG expansion and reasons for geographic variations of prevalence. Hum Mol Genet. 1994;3:2103–14.

59. DeJesus-Hernandez M, Mackenzie IR, Boeve BF, Boxer AL, Baker M, Rutherford NJ, et al. Expanded GGGGCC hexanucleotide repeat in noncoding region of C9ORF72 causes chromosome 9p-linked FTD and ALS. Neuron. 2011;72:245–56. doi:10.1016/j.neuron.2011.09.011.

60. Renton AE, Majounie E, Waite A, Simón-Sánchez J, Rollinson S, Gibbs JR, et al. A hexanucleotide repeat expansion in C9ORF72 is the cause of chromosome 9p21-linked ALS-FTD. Neuron. 2011;72:257–68. doi:10.1016/j.neuron.2011.09.010.

61. Campuzano V, Montermini L, Moltò MD, Pianese L, Cossée M, Cavalcanti F, et al. Friedreich's ataxia: autosomal recessive disease caused by an intronic GAA triplet repeat expansion. Science. 1996;271:1423–7. doi:10.1126/science.271.5254.1423.

62. Mahadevan M, Tsilfidis C, Sabourin L, Shutler G, Amemiya C, Jansen G, et al. Myotonic dystrophy mutation: an unstable CTG repeat in the 3' untranslated region of the gene. Science. 1992;255:1253–5.

63. Kashima T, Rao N, David CJ, Manley JL. hnRNP A1 functions with specificity in repression of SMN2 exon 7 splicing. Hum Mol Genet. 2007;16:3149–59. doi:10.1093/hmg/ddm276.

64. Ward AJ, Cooper TA. The pathobiology of splicing. J Pathol. 2010;220:152–63. doi:10.1002/path.2649.

65. Cartegni L, Hastings ML, Calarco JA, de Stanchina E, Krainer AR. Determinants of exon 7 splicing in the spinal muscular atrophy genes, SMN1 and SMN2. Am J Hum Genet. 2006;78:63–77. doi:10.1086/498853.

66. Kashima T, Manley JL. A negative element in SMN2 exon 7 inhibits splicing in spinal muscular atrophy. Nat Genet. 2003;34:460–3. doi:10.1038/ng1207.

67. Cartegni L, Krainer AR. Disruption of an SF2/ASF-dependent exonic splicing enhancer in SMN2 causes spinal muscular atrophy in the absence of SMN1. Nat Genet. 2002;30:377–84. doi:10.1038/ng854.

68. Takahara K, Schwarze U, Imamura Y, Hoffman GG, Toriello H, Smith LT, et al. Order of intron removal influences multiple splice outcomes, including a two-exon skip, in a COL5A1 acceptor-site mutation that results in abnormal pro-alpha1(V) N-propeptides and Ehlers-Danlos syndrome type I. Am J Hum Genet. 2002;71:451–65. doi:10.1086/342099.

69. Habara Y, Takeshima Y, Awano H, Okizuka Y, Zhang Z, Saiki K, et al. In vitro splicing analysis showed that availability of a cryptic splice site is not a determinant for alternative splicing patterns caused by +1G-->A mutations in introns of the dystrophin gene. J Med Genet. 2009;46:542–7. doi:10.1136/jmg.2008.061259.

70. Anna A, Monika G. Splicing mutations in human genetic disorders: examples, detection, and confirmation. J Appl Genet. 2018;59:253–68. doi:10.1007/s13353-018-0444-7.

71. Zeng L, Liu W, Feng W, Wang X, Dang H, Gao L, et al. A novel donor splice-site mutation of major intrinsic protein gene associated with congenital cataract in a Chinese family. Mol Vis. 2013;19:2244–9.

72. Hori T, Fukao T, Murase K, Sakaguchi N, Harding CO, Kondo N. Molecular basis of two-exon skipping (exons 12 and 13) by c.1248+5g>a in OXCT1 gene: study on intermediates of OXCT1 transcripts in fibroblasts. Hum Mutat. 2013;34:473–80. doi:10.1002/humu.22258.

73. Känsäkoski J, Jääskeläinen J, Jääskeläinen T, Tommiska J, Saarinen L, Lehtonen R, et al. Complete androgen insensitivity syndrome caused by a deep intronic pseudoexon-activating mutation in the androgen receptor gene. Sci Rep. 2016;6:32819. doi:10.1038/srep32819.

74. Fang LJ, Simard MJ, Vidaud D, Assouline B, Lemieux B, Vidaud M, et al. A novel mutation in the neurofibromatosis type 1 (NF1) gene promotes skipping of two exons by preventing exon definition. J Mol Biol. 2001;307:1261–70. doi:10.1006/jmbi.2001.4561.

75. Symoens S, Malfait F, Vlummens P, Hermanns-Lê T, Syx D, De Paepe A. A novel splice variant in the N-propeptide of COL5A1 causes an EDS phenotype with severe kyphoscoliosis and eye involvement. PLoS One. 2011;6:e20121. doi:10.1371/journal.pone.0020121.

76. Sanz DJ, Hollywood JA, Scallan MF, Harrison PT. Cas9/gRNA targeted excision of cystic fibrosis-causing deep-intronic splicing mutations restores normal splicing of CFTR mRNA. PLoS One. 2017;12:e0184009. doi:10.1371/journal.pone.0184009.

77. Ramalho AS, Beck S, Penque D, Gonska T, Seydewitz HH, Mall M, et al. Transcript analysis of the cystic fibrosis splicing mutation 1525-1G>A shows use of multiple alternative splicing sites and suggests a putative role of exonic splicing enhancers. J Med Genet. 2003;40:e88.

78. Ridge PG, Karch CM, Hsu S, Arano I, Teerlink CC, Ebbert MTW, et al. Linkage, whole genome sequence, and biological data implicate variants in RAB10 in Alzheimer's disease resilience. Genome Med. 2017;9:100. doi:10.1186/s13073-017-0486-1.

79. Lettice LA, Heaney SJH, Purdie LA, Li L, de Beer P, Oostra BA, et al. A long-range Shh enhancer regulates expression in the developing limb and fin and is associated with preaxial polydactyly. Hum Mol Genet. 2003;12:1725–35. doi:10.1093/hmg/ddg180.

80. Emison ES, McCallion AS, Kashuk CS, Bush RT, Grice E, Lin S, et al. A common sex-dependent mutation in a RET enhancer underlies Hirschsprung disease risk. Nature. 2005;434:857–63. doi:10.1038/nature03467.

81. de Vooght KMK, van Wijk R, van Solinge WW. Management of gene promoter mutations in molecular diagnostics. Clin Chem. 2009;55:698–708. doi:10.1373/clinchem.2008.120931.

82. Short PJ, McRae JF, Gallone G, Sifrim A, Won H, Geschwind DH, et al. De novo mutations in regulatory elements in neurodevelopmental disorders. Nature. 2018;555:611–6. doi:10.1038/nature25983.

83. De Gobbi M, Viprakasit V, Hughes JR, Fisher C, Buckle VJ, Ayyub H, et al. A regulatory SNP causes a human genetic disease by creating a new transcriptional promoter. Science. 2006;312:1215–7. doi:10.1126/science.1126431.

84. Grant SF, Reid DM, Blake G, Herd R, Fogelman I, Ralston SH. Reduced bone density and osteoporosis associated with a polymorphic Sp1 binding site in the collagen type I alpha 1 gene. Nat Genet. 1996;14:203–5. doi:10.1038/ng1096-203.

85. Benko S, Fantes JA, Amiel J, Kleinjan D-J, Thomas S, Ramsay J, et al. Highly conserved non-coding elements on either side of SOX9 associated with Pierre Robin sequence. Nat Genet. 2009;41:359–64. doi:10.1038/ng.329.

86. Jeong Y, Leskow FC, El-Jaick K, Roessler E, Muenke M, Yocum A, et al. Regulation of a remote Shh forebrain enhancer by the Six3 homeoprotein. Nat Genet. 2008;40:1348–53. doi:10.1038/ng.230.

87. Rahimov F, Marazita ML, Visel A, Cooper ME, Hitchler MJ, Rubini M, et al. Disruption of an AP-2alpha binding site in an IRF6 enhancer is associated with cleft lip. Nat Genet. 2008;40:1341–7. doi:10.1038/ng.242.

88. Faghihi MA, Modarresi F, Khalil AM, Wood DE, Sahagan BG, Morgan TE, et al. Expression of a noncoding RNA is elevated in Alzheimer's disease and drives rapid feed-forward regulation of beta-secretase. Nat Med. 2008;14:723–30. doi:10.1038/nm1784.

89. Chen W-L, Lin J-W, Huang H-J, Wang S-M, Su M-T, Lee-Chen G-J, et al. SCA8 mRNA expression suggests an antisense regulation of KLHL1 and correlates to SCA8 pathology. Brain Res. 2008;1233:176–84. doi:10.1016/j.brainres.2008.07.096.

90. Kogo R, Shimamura T, Mimori K, Kawahara K, Imoto S, Sudo T, et al. Long noncoding RNA HOTAIR regulates polycomb-dependent chromatin modification and is associated with poor prognosis in colorectal cancers. Cancer Res. 2011;71:6320–6. doi:10.1158/0008-5472.CAN-11-1021.

91. Ishii N, Ozaki K, Sato H, Mizuno H, Saito S, Takahashi A, et al. Identification of a novel non-coding RNA, MIAT, that confers risk of myocardial infarction. J Hum Genet. 2006;51:1087–99. doi:10.1007/s10038-006-0070-9.

92. Khalil AM, Faghihi MA, Modarresi F, Brothers SP, Wahlestedt C. A novel RNA transcript with antiapoptotic function is silenced in fragile X syndrome. PLoS One. 2008;3:e1486. doi:10.1371/journal.pone.0001486.

93. Chubb JE, Bradshaw NJ, Soares DC, Porteous DJ, Millar JK. The DISC locus in psychiatric illness. Mol Psychiatry. 2008;13:36–64. doi:10.1038/sj.mp.4002106.

94. Tripathi V, Ellis JD, Shen Z, Song DY, Pan Q, Watt AT, et al. The nuclear-retained noncoding RNA MALAT1 regulates alternative splicing by modulating SR splicing factor phosphorylation. Mol Cell. 2010;39:925–38. doi:10.1016/j.molcel.2010.08.011.

95. Matouk IJ, DeGroot N, Mezan S, Ayesh S, Abu-lail R, Hochberg A, et al. The H19 non-coding RNA is essential for human tumor growth. PLoS One. 2007;2:e845. doi:10.1371/journal.pone.0000845.

96. Lin R, Maeda S, Liu C, Karin M, Edgington TS. A large noncoding RNA is a marker for murine hepatocellular carcinomas and a spectrum of human carcinomas. Oncogene. 2007;26:851–8. doi:10.1038/sj.onc.1209846.

97. Yang Z, Zhou L, Wu L-M, Lai M-C, Xie H-Y, Zhang F, et al. Overexpression of long non-coding RNA

HOTAIR predicts tumor recurrence in hepatocellular carcinoma patients following liver transplantation. Ann Surg Oncol. 2011;18:1243–50. doi:10.1245/s10434-011-1581-y.

98. Zorzetto M, Datturi F, Divizia L, Pistono C, Campo I, De Silvestri A, et al. Complement C4A and C4B gene copy number study in alzheimer's disease patients. Curr Alzheimer Res. 2017;14:303–8. doi:10.2174/1567205013666161013091934.

99. Trouw LA, Nielsen HM, Minthon L, Londos E, Landberg G, Veerhuis R, et al. C4b-binding protein in Alzheimer's disease: binding to Abeta1-42 and to dead cells. Mol Immunol. 2008;45:3649–60. doi:10.1016/j.molimm.2008.04.025.

100. Ridge PG, Mukherjee S, Crane PK, Kauwe JSK, Alzheimer's Disease Genetics Consortium. Alzheimer's disease: analyzing the missing heritability. PLoS One. 2013;8:e79771. doi:10.1371/journal.pone.0079771.

101. Ebbert MTW, Ridge PG, Kauwe JSK. Bridging the gap between statistical and biological epistasis in Alzheimer's disease. Biomed Res Int. 2015;2015:870123. doi:10.1155/2015/870123.

102. Ebbert MTW, Boehme KL, Wadsworth ME, Staley LA, Alzheimer's Disease Neuroimaging Initiative, Alzheimer's Disease Genetics Consortium, et al. Interaction between variants in CLU and MS4A4E modulates Alzheimer's disease risk. Alzheimers Dement. 2016;12:121–9. doi:10.1016/j.jalz.2015.08.163.

103. Weirather JL, de Cesare M, Wang Y, Piazza P, Sebastiano V, Wang X-J, et al. Comprehensive comparison of Pacific Biosciences and Oxford Nanopore Technologies and their applications to transcriptome analysis. [version 2; referees: 2 approved]. F1000Res. 2017;6:100. doi:10.12688/f1000research.10571.2.

104. Ardui S, Ameur A, Vermeesch JR, Hestand MS. Single molecule real-time (SMRT) sequencing comes of age: applications and utilities for medical diagnostics. Nucleic Acids Res. 2018;46:2159–68. doi:10.1093/nar/gky066.

105. Zook JM, Catoe D, McDaniel J, Vang L, Spies N, Sidow A, et al. Extensive sequencing of seven human genomes to characterize benchmark reference materials. Sci Data. 2016;3:160025. doi:10.1038/sdata.2016.25.

106. Zook J, McDaniel J, Parikh H, Heaton H, Irvine SA, Trigg L, et al. Reproducible integration of multiple sequencing datasets to form high-confidence SNP, indel, and reference calls for five human genome reference materials. BioRxiv. 2018. doi:10.1101/281006.

107. Li H. Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics. 2018;34:3094–100. doi:10.1093/bioinformatics/bty191.

108. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. Bioinformatics. 2010;26:841–2. doi:10.1093/bioinformatics/btq033.

109. Thorvaldsdóttir H, Robinson JT, Mesirov JP. Integrative Genomics Viewer (IGV): high-performance genomics data visualization and exploration. Brief Bioinformatics. 2013;14:178–92. doi:10.1093/bib/bbs017.

110. Gel B, Serra E. karyoploteR: an R/Bioconductor package to plot customizable genomes displaying arbitrary data. Bioinformatics. 2017;33:3088–90. doi:10.1093/bioinformatics/btx346.

111. Tripathi S, Pohl MO, Zhou Y, Rodriguez-Frandsen A, Wang G, Stein DA, et al. Meta- and Orthogonal Integration of Influenza "OMICs" Data Defines a Role for UBR4 in Virus Budding. Cell Host Microbe. 2015;18:723–35. doi:10.1016/j.chom.2015.11.002.

112. Hu B, Jin J, Guo A-Y, Zhang H, Luo J, Gao G. GSDS 2.0: an upgraded gene feature visualization server. Bioinformatics. 2015;31:1296–7. doi:10.1093/bioinformatics/btu817.

113. Lefebvre S, Bürglen L, Reboullet S, Clermont O, Burlet P, Viollet L, et al. Identification and characterization of a spinal muscular atrophy-determining gene. Cell. 1995;80:155–65. doi:10.1016/0092-8674(95)90460-3.

**Figure 1. Genomic regions may be 'dark' by depth or mapping quality, many of which are 'camouflaged'.**

Large, complex genomes are known to contain 'dark' regions where standard high-throughput short-read

1    sequencing technologies cannot be adequately assembled or aligned. We split these dark regions into two

2    types: (1) dark because of low depth; and (2) dark because of low mapping quality (MAPQ), which are mostly

3    'camouflaged'. **(a)** *HLA-DRB5* encodes a Major Histocompatibility Complex protein that plays an important role

4    in immune-response and has been associated with several diseases, including Alzheimer's disease. It is well

5    known to be dark (low depth); specifically, when performing whole-genome sequencing using standard short-

6    read sequencing technologies, an insufficient number of reads align, preventing variant callers from assessing

7    mutations. We calculated sequencing depth across *HLA-DRB5* for ten male samples from the Alzheimer's

8    Disease Sequencing Project (ADSP) that were sequenced using standard Illumina whole-genome sequencing

9    with 100-nucleotide read lengths. Approximately 62.0% (50.2% of coding sequence) of *HLA-DRB5* is dark by

10    depth (<5 aligned reads; indicated by red lines). **(b)** *HSPA1A* is a heat-shock protein from the 70-kilodalton

11    (kDa) heat-shock protein family, and plays an important role in stabilizing proteins against aggregation.

12    *HSPA1A* is dark because of low mapping quality (MAPQ <10 for ≥90% of reads at a given position).

13    Approximately 41.8% (52.8% coding sequence) of *HSPA1A* is dark by mapping quality (indicated by red line).

14    Dark gray bars indicate sequencing reads with a relatively high mapping quality, whereas white bars indicate

15    reads with a low mapping quality (MAPQ = 0). **(c)** Many genomic regions that are dark because of mapping

16    quality arise because they have been duplicated in the genome, which we term 'camouflaged' (or 'camo

17    genes'). When confronted with a read that aligns equally well to more than one location, standard sequence

18    aligners randomly assign the read to one location and give it a low mapping quality. Thus, it is unclear from

19    which gene any of the reads indicated by white bars originated from. *HSPA1A* and *HSPA1B* are clear examples

20    of camouflaged genes arising from a tandem duplication. The two genes are approximately 14kb apart and

21    approximately 50% of the genes are identical.


22

23    **Figure 2. Many dark regions involve protein-coding gene regions.** We identified 37873 dark regions (>16

24    million nucleotides) in 5857 gene bodies that were either dark by depth or dark by mapping quality. **(a)**

1 Stratifying the gene bodies by GENCODE biotype, 3635 gene bodies were protein coding, 1102 were

2 pseudogenes, and 720 were long intergenic non-coding RNAs (lincRNA). **(b)** Of all 37873 dark regions, 28598

3 were intronic, 4114 were in lincRNA exons, 2657 were in protein-coding exons (CDS), 1134 were in 3'UTR

4 regions, and 1103 were in 5'UTR regions. Any dark region that spanned a gene element boundary (e.g., intron

5 to exon) was split into separate dark regions.

6

7 **Figure 3. Dark coding regions occur throughout the genome, and are largely resolved with long-read**

8 **sequencing technologies.** We identified 2757 dark coding (CDS) regions (>460000 nucleotides) in 744 protein-

9 coding genes that were dark by either depth or mapping quality (Supplemental Tables 1-2). Exactly 142

10 (19.1%) of the 744 protein-coding genes were 100% dark in CDS regions, 441 (59.3%) were at least 25% dark in

11 CDS regions, and 628 (84.4%) were at least 5% dark in CDS regions (Supplemental Table 1). **(a)** We mapped all

12 protein-coding gene bodies with a dark coding exon to the genome to visualize their genomic location, and are

13 generally spread throughout. There are several tight clusters of dark CDS regions on chromosomes 1, 9, 10,

14 and Y, however. **(b)** We assessed how well increasing read lengths would resolve dark regions by assessing

15 samples sequenced with Illumina whole-genome sequencing using 250-nucleotided read lengths, as well as

16 long-read technologies 10x Genomics, Oxford Nanopore Technologies (ONT), and Pacific Biosciences (PacBio).

17 Data from the samples sequenced using 250-nucleotide Illumina read lengths reduced the area under the

18 curve by 23.2% in CDS regions; this translates to a 24.4% reduction in dark CDS nucleotides. Comparing long-

19 read sequencing technologies to the standard Illumina 100-nucleotide read lengths, 10x Genomics, PacBio,

20 and ONT reduced the area under the curve for CDS regions by approximately 54.9%, 66.7%, and 81.8%,

21 respectively; this translates to a 57.8%, 68.6%, and 81.4% reduction in dark CDS nucleotides, respectively. The

22 area under the curve (AUC) for each technology is scaled in reference to Illumina sequencing based on 100-

23 nucleotide read lengths (i.e., AUC for Illumina 100-nucleotide read lengths = 1). In contrast to overall results,

1 PacBio outperformed 10x Genomics when looking only at CDS regions (see text). Most analyses focused on

2 genes where at least 5% of the CDS nucleotides are dark, indicated by the dashed line.

3

4 **Figure 4. Pathways relevant to human health, development, and reproductive function are affected by dark**

5 **and camouflaged genes.** We characterized the pathways for dark and camouflaged genes using

6 Metascape.org, including only genes where at least 5% of the CDS regions were dark (670 unique gene

7 symbols; based on standard Illumina 100 nucleotide read lengths). **(a)** We identified several pathways that are

8 important in human health, development, and reproductive function (Supplemental Table 11). Specific

9 pathways included defensins (R-HSA-1461973; logP = -7.04), gonadal mesoderm development (GO:0007506;

10 logP = -6.18), base-excision repair (GO:0006284; logP = -5.93), chromatin silencing (GO:0006342; logP = -5.86),

11 Deubiquitination (R-HSA-5688426; logP = -5.32), NLS-bearing protein import into nucleus (GO:0006607; logP =

12 -5.31), spindle assembly (GO:0051225; logP = -5.19), spermatogenesis (GO:0007283; logP = -4.93), and

13 forebrain neuron differentiation (GO:0021879; logP = -4.09). **(b)** Looking specifically at known protein-protein

14 interactions, Metascape identified 138 proteins with 212 known interactions (Supplemental Figure 4), and

15 within those, identified three groups enriched for protein-protein interactions using the MCODE algorithm. All

16 three MCODE groups combined are primarily associated with RNA transport (hsa030313; logP = -17.3;

17 Supplemental Figure 5). Individually, the first group (MCODE1) is enriched for proteins involved in systemic

18 lupus erythematosus (hsa05322; logP = -6.7), cellular response to stress (R-HSA-2262752; logP = -6.6), and

19 RNA transport (hsa03013; logP = -4.39; Supplemental Figure 6). The second group (MCODE2) is enriched with

20 proteins involved in NLS-bearing protein import into nucleus (GO:0006607; logP = -17.1) and protein import

21 into nucleus (GO:0006606; logP = -15.4; Supplemental Figure 7). The third group does not have significant

22 enrichment associations, likely because little is known about them; all four (*PRR20B, PRR20C, PRR20D*, and

23 *PRR20E*) are 100% camouflaged and do not even have known expression measurements in GTEx [29]

1    (Supplemental Figures 8-11).

2

3    **Figure 5. Seventy-five dark genes (≥5% CDS) are associated with 305 human phenotypes, including autism,**

4    **inflammatory bowel disease, and others.** We found 75 genes ≥5% dark CDS that harbor mutations associated

5    with 305 unique human phenotypes, including 277 diseases, according to the Human Gene Mutation Database

6    (HGMD). **(a)** Some of the diseases with the most known associated genes include autism spectrum disorder,

7    hemophilia A, schizophrenia, hearing loss, spinal muscular atrophy, and inflammatory bowel disease. Word

8    size represents the number of genes associated with each disease. Some of the diseases most represented in

9    our data are not surprising, given the number of genes involved in the disease, but these data demonstrate

10   the number of diseases impacted by genes that are at least 5% dark CDS, and how important it is to

11   completely resolve dark regions. We also performed an enrichment analysis, where the diseases most

12   enriched for dark genes included Hemophilia A, color blindness (protan colour vision defect), and X-linked

13   cone-rod dystrophy (Supplemental Figure 12). **(b)** Similarly, we quantified the number of diseases each gene

14   was associated with, and identified many disease-relevant genes with large portions of dark CDS regions that

15   may harbor critical disease-modifying mutations that currently go undetected. Some of the genes with the

16   most known disease associations include *ARX* (14.0% dark CDS), *NEB* (9.5% dark CDS), *TBX1* (10.5% dark CDS),

17   *RPGR* (12.9% dark CDS), *HBA2* (12.8% dark CDS), and *CR1* (26.5% dark CDS). *CR1* is particularly notable for

18   neuroscientists and Alzheimer's disease geneticists, patients, and their caregivers, given that CR1 is a top-ten

19   Alzheimer's disease gene. Other notable genes include SMN1 (89.9% dark CDS) and SMN2 (88.2% dark CDS),

20   which are known to harbor mutations (in camouflaged regions) that are involved in spinal muscular atrophy

21   (SMA) [65, 66, 113]. HSPA1A (52.8% dark CDS) and HSPA1B (51.1% dark CDS) also encode two primary 70-

22   kilodalton (kDa) heat-shock proteins. Heat-shock proteins have been implicated in ALS [31, 32].

23

1 **Figure 6. Camouflaged genes are consistently dark in gnomAD, but dark-by-depth genes may be sample or**

2 **dataset specific**. Most dark genes are specifically camouflaged (Supplemental Tables 12-13), but many are

3 dark by depth; we found that camouflaged regions in the ADSP are consistently dark in the gnomAD

4 consortium data (http://gnomad.broadinstitute.org/) [33]. Dark-by-depth regions may be more variable

5 between samples and datasets, however, suggesting these regions may be sensitive to specific aspects of

6 whole-genome sequencing (e.g., library preparation) or downstream analyses. **(a)** *SMN1* and *SMN2* are

7 camouflaged by each other (89.9% and 88.2% dark CDS, respectively; only *SMN1* shown). Both genes

8 contribute to spinal muscular atrophy, and have been implicated in ALS. **(b)** *HSPA1A* and *HSPA1B* are also

9 camouflaged by each other (52.8% and 51.1% dark CDS, respectively; only *HSPA1A* shown). The heat-shock

10 protein family has been implicated in ALS. **(c)** *NEB* (9.5% dark CDS) is a special case that is camouflaged by

11 itself. *NEB* is associated with 24 diseases in the HGMD, including nemaline myopathy, a hereditary

12 neuromuscular disorder. *NEB* is a large gene, thus, 9.5% dark CDS translates to 2424 protein-coding bases. **(d)**

13 *CR1* is a top Alzheimer's disease gene that plays a critical role in the complement cascade as a receptor for the

14 C3b and C4b complement components, and potentially helps clear amyloid-beta (Aß) [39–41]. *CR1* is also

15 camouflaged by itself (26.5% dark CDS), where the repeated region includes the extracellular C3b and C4b

16 binding domain. The number of repeats and density of certain isoforms have been associated with Alzheimer's

17 disease [21, 42–45]. **(e)** *HLA-DRB5* is dark by depth in the ADSP and gnomAD data (50.2% dark CDS). *HLA-DRB5*

18 has been implicated in several diseases, including Alzheimer's disease. **(f)** *RPGR* is likewise dark in ADSP and

19 gnomAD (12.9% dark CDS), and is associated with several eye diseases, including retinitis pigmentosa and

20 cone-rod dystrophy. **(g)** *ARX* is dark-by-depth (14.0% dark CDS), but varies by sample or cohort, as

21 approximately 70% of gnomAD samples are not strictly dark by depth. *ARX* is associated with diseases

22 including early infantile epileptic encephalopathy 1 (EIEE1) and Partington syndrome. **(h)** Similarly, *TBX1* is not

23 strictly dark by depth in approximately 70% of gnomAD samples (10.5% dark CDS). The Y axes for figures **a-f**

24 indicate median coverage in gnomAD (blue = exomes; green = genomes), whereas the Y axes in **g-h** represent
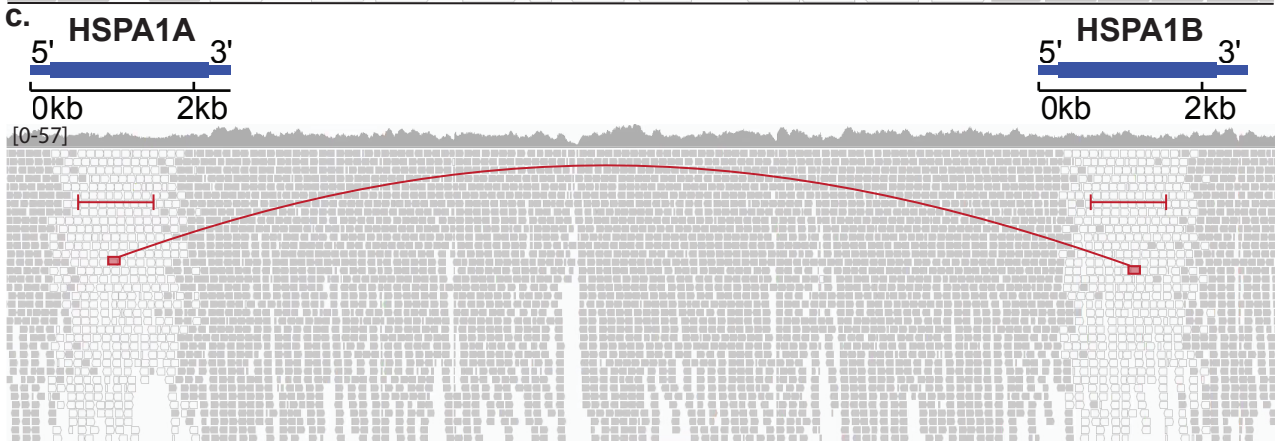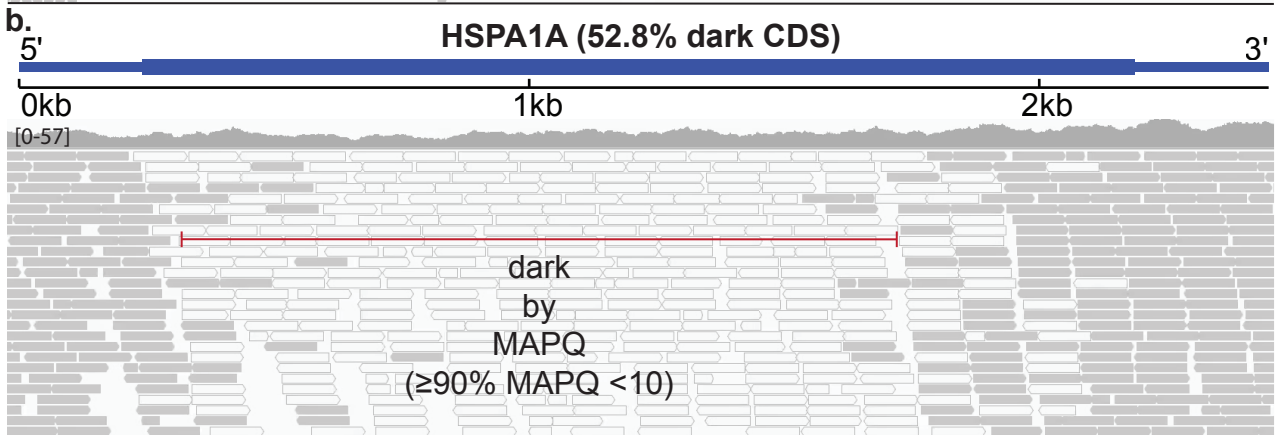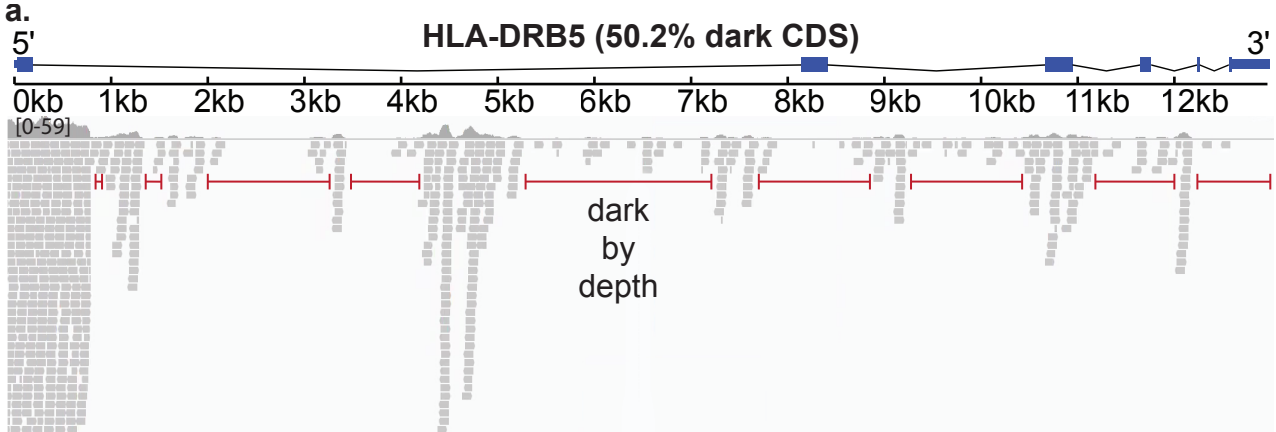
46

1    the proportion of gnomAD samples that have >5x coverage. Dark and camouflaged regions, as well as the

2    percentage of each gene's CDS region that is dark, are indicated by red lines. Dark regions in exome data are

3    either similar or more pronounced than what we observed in whole-genome data, highlighting that dark and

4    camouflaged regions are generally magnified in whole-exome data. For interest, we also discovered that

5    *APOE*—the top genetic risk for Alzheimer's disease [34–36]—is approximately 6% dark CDS (by depth) for

6    certain ADSP samples with whole-genome sequencing, and the same region is dark in gnomAD whole-exome

7    data (Supplemental Figure 13).

8

9    **Figure 7. Long-read technologies resolve many camouflaged regions, with variable success.** We found that

10   ONT's long-read technology appeared to resolve all camouflaged regions well with the high sequencing depth.

11   PacBio performed similarly well, and 10x Genomics performs well under certain circumstances. **(a)** *SMN1* and

12   *SMN2* were 89.9% and 88.2% dark CDS, respectively, based on standard Illumina sequencing with 100-

13   nucleotide read lengths (illuminaRL100), and were 84.0% and 83.1% dark CDS based on Illumina 250-

14   nucleotide read lengths (illuminaRL250; not shown). Both genes were technically 0% dark CDS for 10x

15   Genomics, PacBio, and ONT data. **(b)** *HSPA1A* and *HSPA1B* were 52.8% and 51.1% dark CDS, respectively,

16   based on illuminaRL100 data, and were 50.2% and 49.5% dark CDS based on illuminaRL250 (not shown). Both

17   genes were 0% dark CDS based on ONT and PacBio data, and were 45.8% and 51.8% dark CDS based on 10x

18   Genomics data. In contrast to the results for *SMN1* and *SMN2*, both ONT and PacBio had consistent coverage

19   throughout the camouflaged regions, whereas the camouflaged regions remain dark for 10x Genomics. **(c)** *CR1*

20   was 26.5% dark CDS based on illuminaRL100, and was 24.5% dark based on illuminaRL250 (not shown). 10x

21   Genomics did not improve coverage for *CR1*; the region remained 26.2% dark CDS. Both ONT and PacBio were

22   0% dark CDS. While both PacBio and ONT were able fill the camouflaged region, coverage dropped

23   dramatically throughout the region, despite both genomes being sequenced at 50x and 46x median depth,

1 which does not presently represent average use case for these technologies. The duplicated region is indicated

2 by blue bars, where white lines indicate regions that have diverged sufficiently for reads to align uniquely. It is

3 likely that the performance for ONT and PacBio long-read platforms will be better with longer sequencing

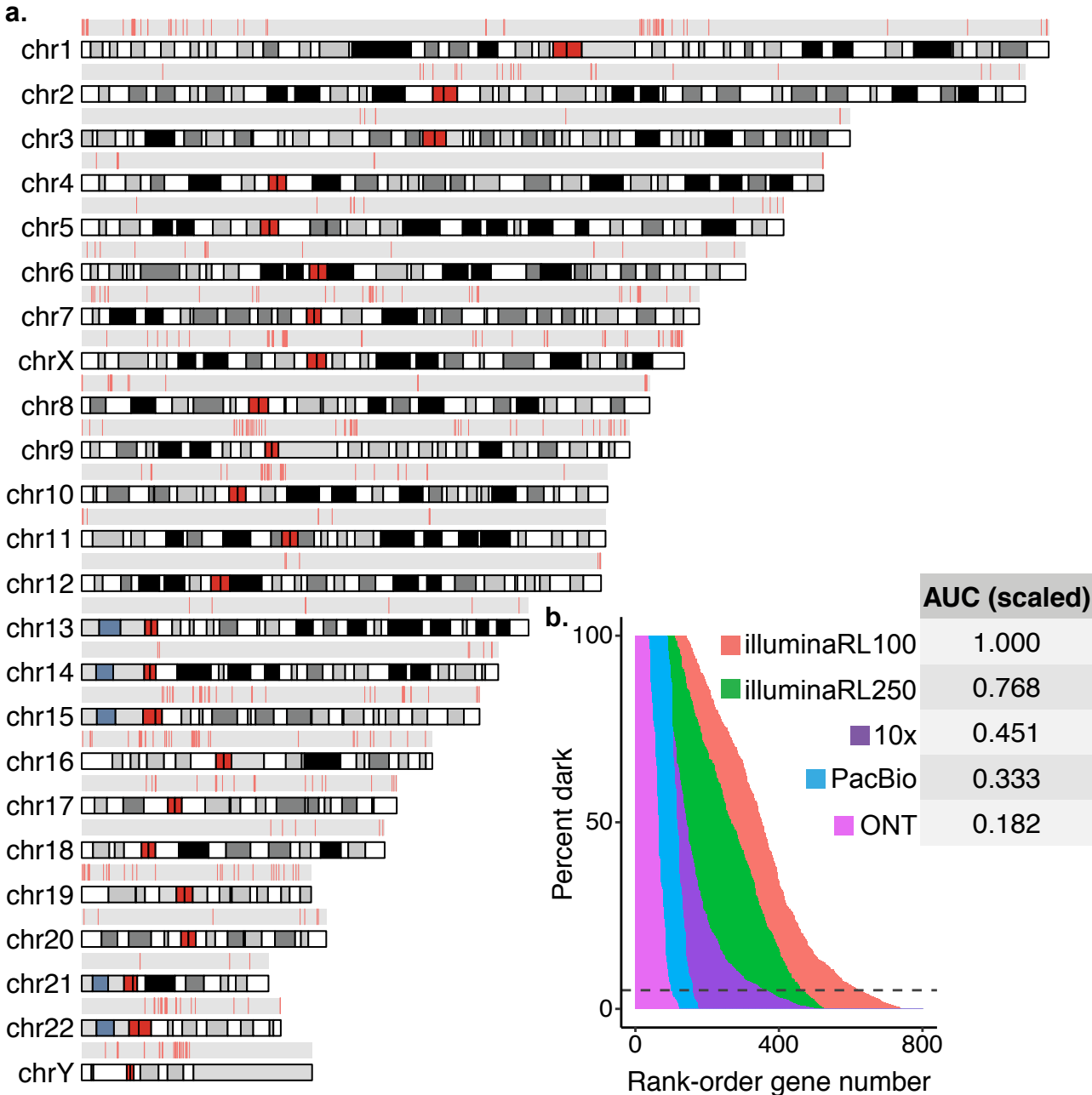4 libraries (e.g. >50kb fragment sizes). Regions were visualized with IGV. Reads with a MAPQ < 10 were filtered,

5 and insertions, deletions, and mismatches are not shown.

6

7 **Figure 8. Many camouflaged regions can be rescued, including *CR1*, even in standard short-read sequencing**

8 **data**. Many large-scale whole-genome or whole-exome sequencing projects exist, covering tens of thousands

9 of individuals. All of these datasets are affected by dark and camouflaged regions that may harbor mutations

10 that either drive or modify disease in patients. Ideally, all samples would be re-sequenced using the latest

11 technologies over time, but financial and biological samples are limited, making it essential to maximize the

12 utility of existing data. We developed a method to rescue mutations in most camouflaged regions, including

13 for standard short-read sequencing data. When confronted with a sequencing read that aligns to two or more

14 regions equally well (with high confidence), most aligners (e.g., BWA [11–13]) will randomly assign the read to

15 one of the regions with a low mapping quality (e.g., MAPQ = 0 for BWA). **(a)** Because the reads are already

16 aligned to one of the regions, we can use the following steps to rescue mutations in most camouflaged

17 regions: (1) extract reads from camouflaged regions; (2) mask all highly similar regions in the reference

18 genome, except one, and re-align the extracted reads; (3) call mutations using standard methods (adjusting

19 for ploidy); and (4) determine precise location using targeted sequencing (e.g., long-range PCR combined with

20 Sanger, or targeted long-read sequencing [1]). Without competing camouflaged regions to confuse the aligner,

21 the aligner will assign a high mapping quality, allowing variant callers to behave normally. **(b)** Exons 10, 18,

22 and 26 in *CR1* are identical, according to the reference genome. Standard aligners will randomly scatter reads

23 matching that sequence across these exons and assign a low mapping quality (e.g., MAPQ = 0 for BWA;

24 indicated as hollow reads). Red lines indicate an individual's mutation that exists in one of these exons, but

48

1    reads containing this mutation also get scattered and assigned a low mapping quality. **(c)** By masking exons 18

2    and 26, we can align all of these reads to exon 10 with high mapping qualities to determine whether a

3    mutation exists. We cannot determine at this stage which of the three exons the mutation is actually located

4    in, but researchers can test association with a given disease to determine whether the mutation is worth

5    further investigation. **(d)** As a proof of principle, we rescued approximately 4622 exonic variants in the ADSP

6    (TiTv = 1.97) using our method, including a frameshift mutation in *CR1* (MAF = 0.00019) that is only found in

7    five cases and zero controls (three representative samples shown). The frameshift results in a stop codon

8    shortly downstream. The ADSP is not large enough to formally assess association between the *CR1* frameshift

9    and Alzheimer's disease, but we believe the mutation merits follow-up studies given its location (*CR1* binding

10   domain) and *CR1*'s strong association with disease.

**a.** HLA-DRB5 (50.2% dark CDS)

dark
by
depth

**b.** HSPA1A (52.8% dark CDS)

dark
by
MAPQ
(≥90% MAPQ <10)

**c.** HSPA1A    HSPA1B

a.

b.

| | AUC (scaled) |
|---|---|
| illuminaRL100 | 1.000 |
| illuminaRL250 | 0.768 |
| 10x | 0.451 |
| PacBio | 0.333 |
| ONT | 0.182 |

**a.**

| Term | -log10(P) |
|---|---|
| R-HSA-1461973: Defensins | |
| GO:0007506: gonadal mesoderm development | |
| GO:0006284: base-excision repair | |
| GO:0006607: NLS-bearing protein import into nucleus | |
| GO:0051225: spindle assembly | |
| GO:0007283: spermatogenesis | |
| GO:0021879: forebrain neuron differentiation | |
| GO:0061408: pos. reg. of transc. from RNA pol II promoter in resp. to heat stress | |
| GO:0006346: methylation-dependent chromatin silencing | |
| R-HSA-189085: Digestion of dietary carbohydrate | |
| GO:0097113: AMPA glutamate receptor clustering | |
| GO:0034472: snRNA 3'-end processing | |
| GO:1903311: regulation of mRNA metabolic process | |
| R-HSA-390696: Adrenoceptors | |
| GO:0006600: creatine metabolic process | |
| GO:0090092: reg. of transmem. receptor protein Ser/Thr kinase signaling pathway | |
| CORUM:626: LSD1 complex | |
| GO:0007548: sex differentiation | |
| GO:0018298: protein-chromophore linkage | |
| R-HSA-5683826: Surfactant metabolism | |

**b.**

MCODE1
MCODE2
MCODE3

**a.**

**b.**

**a.** SMN1 — camouflaged (89.9%)

**b.** HSPA1A — camouflaged (52.8%)

**c.** NEB — camouflaged (9.5%)

**d.** CR1 — camouflaged (26.5%)

**e.** HLA-DRB5 — dark (50.2%)

**f.** RPGR — dark (12.9%)

**g.** ARX — dark (14.0%)

**h.** TBX1 — dark (10.5%)

Legend: ■ exomes ■ genomes

**a.**

SMN1 / SMN2 gene tracks with coverage from Illumina, 10x, ONT, and PacBio sequencing platforms.

SMN1: Illumina [0-71], 10x [0-48], ONT [0-25], PacBio [0-40]
SMN2: Illumina [0-89], 10x [0-83], ONT [0-65], PacBio [0-61]

**b.**

HSPA1A / HSPA1B gene tracks with coverage from Illumina, 10x, ONT, and PacBio sequencing platforms.

Illumina [0-61], 10x [0-52], ONT [0-58], PacBio [0-45]

**c.**

CR1 gene track with coverage from Illumina, 10x, ONT, and PacBio sequencing platforms.

Illumina [0-72], 10x [0-74], ONT [0-79], PacBio [0-81]

**Partial CR1 sequence
from exons 10, 18, and 26**