

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository:<https://orca.cardiff.ac.uk/id/eprint/102184/>

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

Chao, Michael J., Gillis, Tammy, Atwal, Ranjit S., Srinidhi Mysore, Jayalakshmi, Arjomand, Jamshid, Harold, Denise, Holmans, Peter Alan , Jones, Lesley , Orth, Michael, Myers, Richard H., Kwak, Seung, Wheeler, Vanessa C., MacDonald, Marcy E., Gusella, James F. and Lee, Jong-Min 2017. Haplotype-based stratification of Huntington's disease. *European Journal of Human Genetics* 25 , pp. 1202-1209. 10.1038/ejhg.2017.125

Publishers page: <https://doi.org/10.1038/ejhg.2017.125>

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies. See <http://orca.cf.ac.uk/policies.html> for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



Haplotype-based stratification of Huntington's disease

Michael J. Chao^{1,2}, Tammy Gillis¹, Ranjit S. Atwal^{1,2}, Jayalakshmi Srinidhi Mysore¹,
Jamshid Arjomand³, Denise Harold^{4,§,&}, Peter Holmans^{4,§}, Lesley Jones^{4,§}, Michael
Orth^{5,§}, Richard H. Myers^{6,§}, Seung Kwak^{7,§}, Vanessa C. Wheeler^{1,2,§}, Marcy E.
MacDonald^{1,2,8,§}, James F. Gusella^{1,8,9,§}, and Jong-Min Lee^{1,2,8,§,*}

¹ Molecular Neurogenetics Unit, Center for Human Genetic Research, Massachusetts
General Hospital, Boston, MA 02114, USA

² Department of Neurology, Harvard Medical School, Boston, MA 02115, USA

³ Genea Biocells, San Diego, CA 92037, USA

⁴ Medical Research Council Centre for Neuropsychiatric Genetics and Genomics,
Department of Psychological Medicine and Neurology, School of Medicine, Cardiff
University, Cardiff, United Kingdom

⁵ Department of Neurology, University of Ulm, Germany

⁶ Department of Neurology and Genome Science Institute, Boston University School of
Medicine, Boston, MA 02118, USA

⁷ CHDI Foundation, Princeton, NJ 08540, USA

⁸ Medical and Population Genetics Program, the Broad Institute of M.I.T. and Harvard,
Cambridge, MA 02142, USA

⁹ Department of Genetics, Harvard Medical School, Boston, MA 02115, USA

[§] Founding GeM-HD Consortium investigators

[&] Present address: School of Biotechnology, Dublin City University, Dublin 9, Ireland

1 *** Correspondence:** Dr. JM Lee, Center for Human Genetic Research, Massachusetts
2 General Hospital, 185 Cambridge Street, Boston, MA 02114, USA. Tel: 617-643-9714;
3 Fax: 617-726-5735; e-mail: jlee51@mgh.harvard.edu
4
5 **Running title:** Haplotype of Huntington's disease
6
7 **Keywords:** Huntington's disease, haplotype, stratification, SNP

1 **ABSTRACT**

2

3 Huntington's disease (HD) is an autosomal dominant neurodegenerative disease caused
4 by expansion of a CAG trinucleotide repeat in *HTT*, resulting in an extended
5 polyglutamine tract in huntingtin. We and others have previously determined that the
6 HD-causing expansion occurs on multiple different haplotype backbones, reflecting
7 more than one ancestral origin of the same type of mutation. In view of the therapeutic
8 potential of mutant allele-specific gene silencing, we have compared and integrated two
9 major systems of *HTT* haplotype definition, combining data from 74 sequence variants
10 to identify the most frequent disease-associated and control chromosome backbones
11 and revealing that there is potential for additional resolution of HD haplotypes. We
12 have used the large collection of 4,078 heterozygous HD subjects analyzed in our
13 recent genome-wide association study of HD age at onset to estimate the frequency of
14 these haplotypes in European subjects, finding that common genetic variation at *HTT*
15 can distinguish the normal and CAG-expanded chromosomes for more than 95% of
16 European HD individuals. As a resource for the HD research community, we have also
17 determined the haplotypes present in a series of publicly available HD subject-derived
18 fibroblasts, induced pluripotent cells, and embryonic stem cells in order to facilitate
19 efforts to develop inclusive methods of allele-specific *HTT* silencing applicable to most
20 HD patients. Our data providing genetic guidance for therapeutic gene-based targeting
21 will significantly contribute to the developments of rational treatments and
22 implementation of precision medicine in HD.

23

1 INTRODUCTION

2

3 Huntington's disease (HD) [MIM 143100] is a progressive neurodegenerative disorder
4 caused by expansion of a CAG repeat in huntingtin (*HTT*) exon 1 that lengthens a
5 normally polymorphic polyglutamine tract in *HTT*¹ and produces characteristic motor
6 disturbances, along with cognitive and psychiatric manifestations.² Both the age at
7 onset and the age at death of HD subjects are inversely correlated with the length of
8 their CAG repeat, while the duration from onset to death, typically 15-20 years is
9 largely independent of the mutation size.^{3,4} Currently, there is no treatment to either
10 delay the onset or slow the progression of HD, but the recent discovery of genetic
11 modifiers of age at onset establishes that the rate of HD pathogenesis can be altered
12 before symptoms appear.⁵ Genetic analysis of large HD cohorts has demonstrated that
13 HD is inherited as a complete dominant where a single mutant *HTT* allele determines
14 the timing of disease onset, with no discernible impact of either the normal *HTT* allele
15 or, when present, a second mutant *HTT* allele.⁴ Consequently, suppression of the
16 expression of mutant *HTT* is an appealing therapeutic strategy which, if achieved in an
17 allele-specific manner,⁶ could avoid any potential negative consequences attributable to
18 deficiency of normal huntingtin activity.

19 *HTT* allele-specific gene silencing strategies can either directly target the
20 expanded CAG repeat or aim at other genetic variants in the surrounding haplotype.⁷⁻⁹
21 While the former is an attractive target that would be applicable in all HD subjects,
22 establishing allele specificity in individuals where the second *HTT* CAG repeat is high
23 in the normal range, and limiting the effect to *HTT* when there are other expressed
24 CAG repeats in the human genome may be technically challenging. However, targeting

1 genetic variants on the mutant *HTT* haplotype can achieve allele-specificity only in
2 those HD individuals who are heterozygous for those variants.¹⁰⁻¹² A multiplicity of
3 *HTT* haplotypes, with normal and expanded repeats, have been observed in HD
4 individuals of European ancestry.¹³⁻¹⁷ We have previously delineated the 8 most
5 common *HTT* haplotypes bearing expanded alleles based upon 21 genetic variants^{14,18}
6 while others have described 3 major haplogroups,¹⁶ which they recently resolved into
7 subtypes using 63 variants.¹¹ The two marker sets, which are only partially overlapping,
8 have each been used to define sites that are most frequently heterozygous in HD
9 subjects as potential targets for allele-specific *HTT* silencing.^{6,11} In order to facilitate
10 research and development towards this goal, we have compared and integrated the two
11 haplotype systems, better estimated *HTT* haplotype frequencies on normal and disease
12 chromosomes in Europeans, and delineated the *HTT* haplotypes present in publicly
13 available cell line resources available to the HD research community.

14

15 **MATERIALS AND METHODS**

16

17 **Definitions of *HTT* haplotypes**

18 Selection of variants, mainly single-nucleotide polymorphisms (SNPs), and samples
19 initially used to characterize *HTT* haplotypes on HD expanded chromosomes and
20 normal chromosomes were described elsewhere.¹⁴ Briefly, twenty SNPs and one 3 bp
21 insertion-deletion that showed significant association with HD in either 1) comparison
22 of all HD vs. controls, or 2) comparison of those HD individuals lacking the major
23 disease haplotype vs. controls, were used for haplotype phasing.¹⁴ The *HTT* CAG
24 repeat sizes in HD individuals were coded as bi-allelic genotypes (expanded and

1 normal), each person being a heterozygote. In contrast, each control individual was
2 coded as homozygous normal for the *HTT* CAG repeat. Haplotype phasing of SNP
3 genotypes was performed by the MaCH program,¹⁹ and the ten most frequent
4 haplotypes on each of expanded chromosomes and normal chromosomes were
5 identified. As four haplotypes overlapped between both disease and normal, the union
6 set comprised 16 distinct haplotypes. Definitions of haplotypes described previously
7 (hap.01 ~ hap.07)¹⁴ are same as those in this study. The phylogeny tree of haplotypes
8 was obtained by the MEGA5 program (neighbor-joining method, P-distance model;
9 <http://www.megasoftware.net/>).

10

11 **Haplotype-specific SNP sites for mutant allele-selective silencing**

12 Previously, based on cumulative heterozygosity analyses of HD subjects with European
13 ancestry, we revealed 20 SNP sites that can be targeted for mutant allele-specific *HTT*
14 silencing / lowering.¹⁸ In order to relate alleles of target SNPs to haplotypes, we
15 determined consensus alleles of those 20 SNPs (10 exon SNPs and 10 intron SNPs) for
16 each haplotype. Briefly, for a given haplotype, we extracted chromosomes from 1000
17 Genomes Project data (phase 1; <http://www.internationalgenome.org/data/>) to
18 determine consensus alleles by taking the most frequent allele of each of 20 target SNP
19 sites. Some of SNP sites are not variable among 16 haplotypes, and consensus alleles of
20 variable SNPs are indicated in Figure 1B together with 2 exon SNPs used to define
21 haplotypes. In 1000 Genomes data, hap.10 is not present, and therefore excluded in this
22 analysis.

23

24 **Haplotypes of publicly available cell lines**

1 We assembled genotypes of 21 tagging SNPs, either from genome-wide association
2 data⁵ or from specific TaqMan assays applied to DNA from blood, lymphoblasts,
3 fibroblasts, induced pluripotent stem cells (iPSC) or derived neural progenitor cells.
4 Those cell line data described in this study represent 59 individuals whose fibroblast
5 cell lines are available in public repositories and 7 human embryonic stem cell (hESC)
6 lines from Genea Biocells Inc. (<http://geneabiocells.com/>). *HTT* CAG repeat length was
7 also determined as described previously.²⁰ Cell line genotype data and the *HTT* CAG
8 repeat genotype coded as a bi-allelic system (expanded or normal) were combined for
9 haplotype phasing in order to identify haplotype carrying expanded CAG or normal
10 repeat. Genotype data for HD and control subjects that were used to define haplotypes¹⁴
11 were also included to increase the accuracy of computational population phasing by the
12 MaCH program.¹⁹ Familial relationships (S. Table 1) were further considered when the
13 relationships between CAG repeats and haplotypes were ambiguous to determine the
14 phase of CAG repeats and haplotypes (e.g., control subjects).

15

16 **Frequencies of haplotypes and haplogroups in control samples**

17 Fully phased 1000 Genomes Project data (Phase 1;
18 <http://www.internationalgenome.org/>) were used to estimate population frequencies of
19 *HTT* haplotypes defined in this study and haplogroups described by Kay et. al.¹¹ Each
20 chromosome was classified into haplotypes based on 21 SNPs, and further summarized
21 for each population group (i.e., Europeans, Asians, Africans, and Ad Mixed
22 Americans). The haplogroup of each chromosome was determined similarly based on
23 63 SNPs, permitting direct delineation of correspondence between haplotype systems
24 on normal chromosomes.

1

2 **Genotype imputation of HD samples**

3 Genotypes on chromosome 4 were imputed for HD samples with European ancestry
4 used in a recent onset-modifier genome-wide association (GWA) study (4082
5 Europeans)⁵ and control samples (1676 Europeans)²¹ using the Michigan Imputation
6 Server.²² Pre-phasing was performed by Eagle2²³ and imputation was performed by
7 Minimac3 using 1000 Genomes Phase 1 as a reference panel (all populations).²⁴ A set
8 of SNPs used for haplotype and/or haplogroup analysis were then extracted from
9 imputed data to determine relationships between haplotypes and haplogroups.

10

11 **Determination of the relationship between haplotypes and haplogroups**

12 Twenty-one genetic variations from our study¹⁴ and 63 tagging SNPs from Kay and
13 colleagues¹¹ were used to classify haplotypes of samples used in the HD modifier
14 GWA study.⁵ There were 10 shared SNPs between the two haplotype systems
15 (rs2798296, GRCh37 chr4:g.3062165A>G; rs3856973, GRCh37 chr4:g.3080173G>A;
16 rs2285086, GRCh37 chr4:g.3089259A>G; rs10015979, GRCh37 chr4:g.3109442A>G;
17 rs11731237, GRCh37 chr4:g.3151813C>T; rs363096, GRCh37 chr4:g.3180021T>C;
18 rs2298969, GRCh37 chr4:g.3186244A>G; rs363092, GRCh37 chr4:g.3196029A>C;
19 rs916171, GRCh37 chr4:g.3216815C>G; and rs362272, GRCh37 chr4:g.3234980G>A)
20 so genotypes for a total of 74 variants were extracted from the imputed data. Then the
21 recoded bi-allelic *HTT* CAG repeat length genotype (expanded or normal) was added to
22 the imputed genotype data, and haplotype phasing was performed for the 75 variant
23 sites in the HD samples (4082 Europeans),⁵ control samples (1676 Europeans),²¹ and
24 1000 Genomes Phase 1 samples (379 Europeans, 181 Ad Mixed Americans, 246

1 Africans, 286 Asians)²⁴ by the Beagle program.²⁵ Subsequently, the CAG-expanded
2 and normal chromosomes from each HD heterozygous subject (4078 Europeans) were
3 named based on 1) our haplotype definitions, and 2) haplotype definitions used by Kay
4 and colleagues¹¹ in order to delineate the relationships between the two haplotype
5 systems.

6

7 **Description of SNPs, Website, and public access**

8 Detailed description of SNPs used in this study can be found in S. Table 2. In addition,
9 description of SNPs, definition of haplotypes, and genotype data are available at
10 chgr.partners.org/htt.haplotype.html. The genotype data set is also available at the
11 European Variation Archive (<http://www.ebi.ac.uk/eva/>) (accession number:
12 PRJEB20817).

13

14

1 RESULTS

2

3 Common SNP-based haplotypes

4 We previously defined the 8 most frequent haplotypes (hap.01 to hap.08) on HD
5 disease-causing chromosomes using 21 common genetic variants, including 20 SNPs
6 and one 3 bp indel, genotyped in 699 unrelated HD subjects and 1,676 population
7 controls of European ancestry.^{14,18} Approximate locations and alleles of DNA
8 variations that were used for haplotype analysis are summarized in Figure 1A. Here, we
9 extend the definitions in that dataset to the most frequent 10 HD and the 10 most
10 frequent normal European *HTT* haplotypes. Four haplotypes were shared between the
11 two groups, so the union created a single set of 16 different haplotypes. These were
12 named based first upon decreasing frequency on CAG-expanded disease chromosomes
13 in this initial HD dataset (hap.01 through hap.10) and then, after excluding the four
14 shared haplotypes (hap.08, hap.02, hap.03 and hap.01, in order of normal frequency),
15 based upon decreasing frequency on normal chromosomes (hap.11 through hap.16); all
16 other rare haplotypes were grouped as “hap.other”. A comparison of the potential
17 relationships between these 16 haplotypes was achieved by phylogeny analysis using a
18 neighbor-joining algorithm. The result is a dendrogram with two main branches
19 containing different-sized sub-clusters (Figure 1A). For example, hap.01, the most
20 common haplotype on the HD disease chromosomes forms a cluster with hap.05 and
21 hap.10, whereas hap.08, the most common haplotype on normal chromosomes is a part
22 of a cluster of haplotypes involving hap.04, hap.16, and hap.14. Divergence of related
23 haplotypes could potentially be explained by a single marker allele change in some
24 cases (e.g., hap.01, hap.05, and hap.10; hap.11 and hap.12; hap.02 and hap.07; hap.04

1 and hap.16), by insertion/deletion of a simple repeat (e.g., hap.01 and hap.12), and by
2 combinations of various genetic events including local recombination or gene
3 conversion. The two main branches of the dendrogram suggest at least two different
4 ancestral origins of *de novo* CAG expansion mutation. However, it is likely that the
5 haplotype diversity within the subclusters reflects the occurrence of many more *de*
6 *novo* expansions rather than being the result of haplotype decay since we have
7 previously demonstrated *de novo* CAG expansion on both hap.01 and hap.05.¹⁸

8

9 **Haplotype-specific target SNP sites for allele-specific silencing**

10 Initial selection of SNPs for haplotyping was based on comparisons between HD
11 subjects and normal control individuals, and therefore did not represent the
12 combination of disease chromosomes and normal chromosomes in HD subjects.¹⁴
13 Subsequently, we performed iterative heterozygosity analyses aiming at revealing a
14 minimal number of SNPs covering the maximum proportion of HD patients in allele-
15 specific gene targeting therapies.¹⁸ Cumulatively, 10 exon SNPs and 10 intron SNPs
16 covered 93.8% and 97.% of HD subjects, respectively, indicating that the vast majority
17 of HD subjects with European ancestry carry at least one heterozygous SNP site among
18 20 nominated targetable locations.¹⁸ However, the heterozygosity analysis did not
19 immediately show mutant alleles to target. Here, we determined consensus alleles of 20
20 targetable SNP sites on each of haplotypes based on 1000 Genomes Project data, and
21 mapped variable alleles on each haplotype. As summarized in Figure 1B, target SNP
22 sites for each diplotype can be selected immediately by comparing two haplotypes
23 (assuming one is mutant chromosome and the other is normal chromosome). For
24 example, if a HD individual carries mutant hap.01 and normal hap.08 chromosomes,

there are 12 SNP sites that can be used to distinguish mutant allele from normal allele (Figure 1B).

Haplotypes of publicly available cell line resources

Results of mutant allele-specific gene silencing studies have produced promising results in animal models,¹⁰ encouraging the application of this approach to human HD. In this context, cell lines derived from HD subjects provide valuable tools to test the specificity and efficacy of allele-specific silencing reagents in pre-clinical experiments. Thus, we performed haplotype analysis using our haplotype system for HD cell lines readily available from various public repositories. Table 1 gives the *HTT* haplotypes for 59 fibroblast lines available from the NIGMS Repository at the Coriell Institute (<https://catalog.coriell.org/1/NIGMS>) or the NINDS Human Cell and Data Repository at RUCDR Infinite Biologics (<https://nindsgenetics.org/>). These include 43 lines representing individuals (from 26 families) with an expanded *HTT* repeat, whose allele lengths range from 38 to 180 CAGs. The remaining 16 lines from 10 families represent control individuals with CAG repeat lengths 33 or shorter. Where possible the phase of the CAG repeat with respect to the *HTT* haplotype was confirmed from family relationships (S. Table 1). In the remaining instances (unrelated subjects; noted by * on the sample ID in Table 1), the phase of the expanded repeat was assigned probabilistically using MaCH program (see methods) or the phase of distinguishable normal alleles was assigned arbitrarily for control individuals. As expected from HD population data, the most frequent haplotype on the disease and normal chromosomes in these families are hap.01 and hap.08, respectively, and this most common HD diplotype, hap.01/hap.08, is present in multiple lines from independent families.

1 However, many other HD haplotypes and diplotypes are also represented. Only 5 of the
2 HD individuals are homozygous for the same haplotype, and 4 of these, two of which
3 are also homozygous for an expanded CAG repeat, derive from the large Venezuela
4 HD kindreds in which the disease segregates with hap.03 haplotype.

5 Induced pluripotent stem cell (iPSC) lines are already available to the research
6 community from the above repositories or from the Cedars-Sinai iPSC Core
7 ([https://www.cedars-sinai.edu/Research/Research-Cores/Induced-Pluripotent-Stem-](https://www.cedars-sinai.edu/Research/Research-Cores/Induced-Pluripotent-Stem-Cell-Core/)
8 [Cell-Core-/](https://www.cedars-sinai.edu/Research/Research-Cores/Induced-Pluripotent-Stem-Cell-Core/)) for 11 of the subjects with expanded repeat fibroblast lines and 5 of the
9 normals, as noted in Table 1. In addition, we have performed haplotyping for 7 human
10 embryonic stem cell (hESC) lines, with expanded CAG alleles ranging from 40 to 48
11 repeats, available from Genea Biocells, as shown in Table 2. The HD mutation in these
12 lines resides either on hap.01 (4 independent lines) or hap.02 (3 lines from the same
13 family).

14

15 **Haplogroup definition of HD chromosomes**

16 A different set of genetic markers (S. Table 2) has been used by others to define
17 haplogroups A, B and C, each of which represents a cluster of similar haplotypes.^{13,16,17}
18 Recently, Kay *et al.* performed a more detailed analysis of the haplogroup system in
19 738 European reference haplotypes from the 1000 Genomes Project and 2,364
20 haplotypes from HD patients and relatives in Canada and Europe to define individual
21 subtypes within each haplogroup based upon 63 genetic variants across *HTT*.¹¹ Across
22 the Canadian and European HD subjects, selected subtypes from the A haplogroup
23 accounted for 86% of all CAG-expanded chromosomes, but the remaining HD

1 chromosomes fell into haplogroup B or C subtypes or rarely, into none of the three
2 major haplogroups (“Other”).

3

4 **Comparison and integration of the two *HTT* haplotype systems**

5 Between the 21 markers used in our haplotype system and the 63 markers used in the
6 recent subdividing of A, B, and C haplogroups, only 10 markers are overlapping (S.
7 Table 2). In order to maximize the utility of both haplotype systems, we have directly
8 compared them by examining the fully phased 1000 Genomes Project haplotype data
9 (Phase 1, Release v3; <http://www.internationalgenome.org/>). To extend the analysis
10 across all available populations rather than only Europeans, we analyzed a total of
11 1,092 control individuals (2,184 normal chromosomes) consisting of Africans (ASW,
12 LWK and YRI), Ad Mixed Americans (CLM, MXL, and PUR), East Asians (CHB,
13 CHS, and JPT) and Europeans (CEU, FIN, GBR, IBS, and TSI). Each 1000 Genomes
14 chromosome was independently classified into our haplotypes using 21 variants sites
15 and haplogroup subtypes using 63 variant sites. The hap.01-hap.16 designations
16 encompassed almost 76% of European chromosomes and more than 60% of Ad Mixed
17 American and Asian chromosomes, but only 23% of African chromosomes, which
18 display far greater genetic complexity (S. Table 3). A similar pattern was evident using
19 haplogroup subtypes which accounted for almost 67% of European chromosomes,
20 about half of Ad Mixed American and Asian chromosomes and only about 7% of
21 African chromosomes (S. Table 3). Subsequently, we delineated the relationships
22 between the two haplotype systems by calculating the percentage of chromosomes with
23 each haplotype defined in our system that distributed to each haplotype defined in the
24 haplogroup system (S. Table 4; Figure 2) and, vice versa (S. Table 5). For example,

1 93.6% and 100% of chromosomes defined as bearing the related hap.01 or hap.05
2 haplotypes are classified as haplogroup subtype A1a (S. Table 4; Figure 2). Among
3 only Europeans, the same correspondence is 100% for both haplotypes. The third
4 related member of this haplotype subcluster from Figure 1, hap.10, was originally
5 defined from HD chromosomes but was not seen on any 1000 Genomes Project
6 chromosomes and so is not reflected in the Tables. The haplotype most common on
7 European normal chromosomes, hap.08, corresponds 95.1% of the time with the C1
8 subtype designation in the haplogroup system (S. Table 4). However for some other
9 haplotypes, the correspondence is not so direct, as some hap.02 chromosomes (25.6%)
10 are classified as haplogroup subtype A2a while others (61.0%) are classified as A2b.
11 Similarly, hap.06 also divides between these two related A2 subtypes, but is primarily
12 assigned to the “Other” class, not being classified as haplogroup A, B or C.
13 Interestingly, haplotypes hap.04, hap.07 and hap.09, which were named by decreasing
14 order of their frequency on HD disease chromosomes in our original study all
15 correspond to the “Other” class of haplogroups except 12.5% of the hap.04 group,
16 which are designated as C4b.

17 Considering the reverse comparison of chromosomes named by the haplogroup
18 system to our haplotypes (S. Table 5), the correspondence is similar to the above, with
19 the A1, A2 and A3 designations, which are the most common on European HD
20 chromosomes, corresponding largely to hap.01+hap.05, hap.02+hap.06 and hap.03,
21 respectively, encompassing most of the haplotypes seen frequently on European HD
22 chromosomes. Those haplogroup subtypes rarely seen on HD chromosomes, such as
23 A4a, A4b, A5a, A5b, B1a, C2, C4 and C6 correspond largely with haplotypes seen on
24 the normal chromosomes in our HD dataset (hap.12, hap.11, hap.12, hap.15, hap.13,

hap.14, hap.16, and hap.14, respectively), or in the cases of B1b, B2, C3, C5, C7 and C8, among the mixed hap.other group of less frequent normal haplotypes.

Overall, these comparisons indicate that the haplotypes most frequently associated with HD disease chromosomes (i.e., hap.01, hap.02, and hap.03) correspond in general with haplogroup subtypes A1, A2, and A3. However, there is the potential for additional resolution in both systems, as illustrated by the fact that the subtypes of A2 (A2a and A2b) subdivide the hap.02 chromosomes, but each subtype (A2a and A2b) is also classified into either hap.02 or hap.06. Similarly, the lack of strong correspondence with haplogroup subtypes of some of the rarer haplotypes identified on HD chromosomes in our studies suggests yet greater diversity among disease chromosomes, and predicts that additional genetic variants can further subdivide the defined haplotypes and haplogroups, particularly in non-European populations.

***HTT* haplotype frequencies on CAG-expanded and normal chromosomes**

The comparisons in S. Tables 4 and 5 relied on fully phased control chromosomes to define haplotype/haplogroup relationships in samples with various ancestries. To estimate the frequency of these groupings on HD chromosomes of European ancestry, we examined the imputed genotypes of 4,078 heterozygous HD subjects recently studied in a GWA study of HD modifiers.⁵ We extracted a unionset of 74 SNPs (representing 21 SNPs used to define our *HTT* haplotypes and the 53 non-overlapping variant sites used by Kay et. al.), and performed probabilistic phasing of the marker alleles using the Beagle program.²⁵ Each of 74-SNP haplotypes (either expanded CAG or normal CAG chromosome) was assigned to both a haplotype and a haplogroup subtype, generating a data set that permitted assessment of the frequency of each

1 haplotype/haplogroup subtype combination. When focusing on our haplotypes defined
2 in this study (Figure 1), frequencies of haplotypes of the expanded and normal
3 chromosomes based on a large collection of HD subjects with European ancestry
4 revealed that HD expansion mutation sits on diverse haplotypes that are also present in
5 normal chromosomes (Figure 3). In addition, comparisons of haplotype frequencies
6 revealed overrepresented and underrepresented haplotypes in HD. For example, hap.01
7 and hap.08 are enriched in disease and normal chromosomes, respectively (Figure 3).
8 Frequency data predicted that the most common diplotype in heterozygous HD subjects
9 would be expanded CAG repeat on hap.01 and normal CAG repeat on hap.08. When
10 comparing our haplotypes to haplogroup, overall, 78% and 71% of European HD and
11 normal chromosomes, respectively, were assignable to discrete ‘super’-haplotype
12 backbones that combined discrete haplotypes and haplogroup subtypes, excluding the
13 uncertain hap.other and haplogroup ‘Other’ catch-all categories (Table 3). As expected,
14 the most frequent HD chromosome backbone was hap.01/A1a and comprised over 38%
15 of European HD chromosomes from the GWA study. Similarly, the most frequent
16 control backbone hap.08/C1 accounted for about 25% of normal chromosomes.
17 Examination of diplotypes of the 4,078 European HD individuals revealed that 56%
18 possessed HD and normal chromosomes that could both be assigned to a fully-defined
19 haplotype/haplogroup backbone, without the uncertainty of the hap.other and
20 haplogroup ‘Other’ categories (Table 4). Notably, less than 5% of these HD subjects
21 had fully-defined chromosomal backbones that were identical on disease and normal
22 chromosomes, being homozygous for all tagging markers. If all 4,078 heterozygous
23 HD subjects were analyzed, 4.9% of them carry identical alleles for 74 SNPs,
24 suggesting that the majority of HD subjects of European descents are eligible for allele-

specific gene targeting strategies. Our previous full sequence analysis of HD hap.01 chromosomes suggests that many of the individuals with the same haplotype backbone on the normal and disease chromosomes could harbour heterozygous variants not considered in the current haplotypes/haplogroups,¹⁸ further implying an additional likelihood of allele discrimination.

DISCUSSION

Huntingtin (*HTT*) shows evolutionarily conserved structural characteristics, and deficiency or hypomorphism of huntingtin are associated with pleiotropic effects involving a number of critical biological processes,²⁶ suggesting that *HTT* silencing approaches to treat HD may need to be specific to the mutant allele. Allele-specific silencing of *HTT* can be achieved either by directly targeting the CAG repeats or, alternatively, by targeting polymorphisms in linkage disequilibrium (LD) with the CAG expansion.^{8,10} Because the HD mutation can occur across a wide range of pathogenic sizes, and CAG repeats are found in many other genes, directly targeting the CAG expansion could result in variable levels of allele selectivity and off-target effects. Previous studies have demonstrated the feasibility of silencing the expression of the expanded allele by targeting a variation on the expanded chromosome.^{10,27-29} Recently, SNP heterozygosity analysis has revealed that the disease chromosome can be distinguished from the normal chromosome in most HD subjects of European ancestry.^{11,16,18,28,30} Therapeutic strategy leveraging a SNP-targeting approach is therefore possible (Figure 1B), but would require knowledge about presence of target SNP site, haplotype phasing, and preferably additional exon SNP sites for outcome measurements (i.e., levels of mutant *HTT*) for a given HD individual. Still, analytical

1 pipelines to identify variant alleles on the CAG-expanded chromosome of an HD
2 individual are yet to be developed because simple genotyping assays do not
3 differentiate allelic phase unless family members are also analyzed. This limitation can
4 be overcome by computational haplotype phasing approaches, because haplotype
5 phasing with a large collection of HD data allows relatively accurate inference of the
6 disease and normal chromosome. Results described here on haplotype phasing of large
7 population of HD individuals, can help populate attributes on HD patient database, and
8 inform where patient groups enriched for targeting SNP can be sought. Subsequent
9 sequence analysis of representative common HD haplotypes and pair-wise comparisons
10 then provide a comprehensive list of targetable sites for each diplotype. In addition,
11 development of allele-specific *HTT* quantification assays to assess the efficacy and
12 allele specificity of silencing reagents require knowledge of variations and their
13 relationships to the expanded chromosomes. Therefore, haplotypes of expanded
14 chromosomes, individual-level diplotype data, and our analytical pipelines provide
15 guidance for identifying targets for mutant allele-specific *HTT* lowering strategies and a
16 route to developing allele-specific readouts to assess specificity of silencing reagents.
17 In addition, genome-wide genotyping assays for HD subjects in a large observational
18 study is on going (i.e., ENROLL-HD), and our pipelines can efficiently identify each
19 individual's expanded and normal chromosomes. Such individual level diplotype data
20 will be critically important in stratifying subjects to identify optimal study populations
21 in clinical trials.

22 In summary, we performed individual level haplotype analyses on a large
23 cohort of HD subjects to evaluate the power of haplotype-based genetics in stratifying
24 HD subjects. Our haplotypes based on a relatively small number of SNPs were able to

1 distinguish mutant chromosomes from their normal counterparts, and confirmed that
2 the majority of HD subjects carry two different haplotypes, further supporting the
3 conclusion from population-based SNP analysis that most HD individuals could be
4 eligible for allele-specific gene silencing¹⁸ and demonstrating the efficiency of
5 haplotype-based approaches. By providing the HD haplotypes of commonly-used
6 publicly available cell lines and a haplotype conversion tables for the comparable
7 haplogroup classification strategy, we hope to promote and facilitate the use of these
8 resources to accelerate pre-clinical allele-specific gene silencing studies and a true
9 precision medicine approach to HD.
10

1 **CONFLICT OF INTEREST**

2 The authors declare no conflict of interest.

3

4 **ACKNOWLEDGEMENTS**

5 We would like to thank all HD patients and their families who generously participated
6 in this study. The full list of clinical investigators contributing samples to the
7 generation of genetic data sets used in this study can be found at PMC4524551. This
8 work was supported by the CHDI Foundation, by grants U01NS082079,
9 R01NS091161, R01HG002449, and P50NS016367 from the National Institutes of
10 Health (USA), and by grants G0801418 and MR/L010305/1 from the Medical
11 Research Council (UK).

12

1 REFERENCES

2

- 3 1. The Huntington's Disease Collaborative Research Group: A novel gene
4 containing a trinucleotide repeat that is expanded and unstable on
5 Huntington's disease chromosomes. *Cell* 1993; **72**: 971-983.
- 6 2. Bates GP, Dorsey R, Gusella JF *et al*: Huntington disease. *Nature Reviews*
7 *Disease Primers* 2015; 15005.
- 8 3. Keum JW, Shin A, Gillis T *et al*: The HTT CAG-Expansion Mutation
9 Determines Age at Death but Not Disease Duration in Huntington Disease.
10 *Am J Hum Genet* 2016; **98**: 287-298.
- 11 4. Lee JM, Ramos EM, Lee JH *et al*: CAG repeat expansion in Huntington
12 disease determines age at onset in a fully dominant fashion. *Neurology*
13 2012; **78**: 690-695.
- 14 5. Genetic Modifiers of Huntington's Disease (GeM-HD) Consortium:
15 Identification of Genetic Factors that Modify Clinical Onset of Huntington's
16 Disease. *Cell* 2015; **162**: 516-526.
- 17 6. Shin JW, Kim KH, Chao MJ *et al*: Permanent inactivation of Huntington's
18 disease mutation by personalized allele-specific CRISPR/Cas9. *Hum Mol*
19 *Genet* 2016.
- 20 7. Hu J, Matsui M, Gagnon KT *et al*: Allele-specific silencing of mutant
21 huntingtin and ataxin-3 genes by targeting expanded CAG repeats in
22 mRNAs. *Nat Biotechnol* 2009; **27**: 478-484.

- 1 8. Keiser MS, Kordasiewicz HB, McBride JL: Gene suppression strategies for
2 dominantly inherited neurodegenerative diseases: lessons from
3 Huntington's disease and spinocerebellar ataxia. *Hum Mol Genet* 2015.
- 4 9. Yu D, Pendergraft H, Liu J *et al*: Single-stranded RNAs use RNAi to potently
5 and allele-selectively inhibit mutant huntingtin expression. *Cell* 2012; **150**:
6 895-908.
- 7 10. Carroll JB, Warby SC, Southwell AL *et al*: Potent and selective antisense
8 oligonucleotides targeting single-nucleotide polymorphisms in the
9 Huntington disease gene / allele-specific silencing of mutant huntingtin.
10 *Mol Ther* 2011; **19**: 2178-2185.
- 11 11. Kay C, Collins JA, Skotte NH *et al*: Huntingtin Haplotypes Provide
12 Prioritized Target Panels for Allele-specific Silencing in Huntington Disease
13 Patients of European Ancestry. *Mol Ther* 2015; **23**: 1759-1771.
- 14 12. Southwell AL, Skotte NH, Kordasiewicz HB *et al*: In vivo evaluation of
15 candidate allele-specific mutant huntingtin gene silencing antisense
16 oligonucleotides. *Mol Ther* 2014; **22**: 2093-2106.
- 17 13. Baine FK, Kay C, Ketelaar ME *et al*: Huntington disease in the South African
18 population occurs on diverse and ethnically distinct genetic haplotypes.
19 *Eur J Hum Genet* 2013; **21**: 1120-1127.
- 20 14. Lee JM, Gillis T, Mysore JS *et al*: Common SNP-based haplotype analysis of
21 the 4p16.3 Huntington disease gene region. *Am J Hum Genet* 2012; **90**: 434-
22 444.

- 1 15. Ramos EM, Gillis T, Mysore JS *et al*: Prevalence of Huntington's disease
2 gene CAG trinucleotide repeat alleles in patients with bipolar disorder.
3 *Bipolar Disord* 2015; **17**: 403-408.
- 4 16. Warby SC, Montpetit A, Hayden AR *et al*: CAG expansion in the Huntington
5 disease gene is associated with a specific and targetable predisposing
6 haplogroup. *Am J Hum Genet* 2009; **84**: 351-366.
- 7 17. Warby SC, Visscher H, Collins JA *et al*: HTT haplotypes contribute to
8 differences in Huntington disease prevalence between Europe and East
9 Asia. *Eur J Hum Genet* 2011; **19**: 561-566.
- 10 18. Lee JM, Kim KH, Shin A *et al*: Sequence-Level Analysis of the Major
11 European Huntington Disease Haplotype. *Am J Hum Genet* 2015; **97**: 435-
12 444.
- 13 19. Li Y, Willer CJ, Ding J, Scheet P, Abecasis GR: MaCH: using sequence and
14 genotype data to estimate haplotypes and unobserved genotypes. *Genet*
15 *Epidemiol* 2010; **34**: 816-834.
- 16 20. Perlis RH, Smoller JW, Mysore J *et al*: Prevalence of incompletely penetrant
17 Huntington's disease alleles among individuals with major depressive
18 disorder. *Am J Psychiatry* 2010; **167**: 574-579.
- 19 21. Myocardial Infarction Genetics Consortium, Kathiresan S, Voight BF *et al*:
20 Genome-wide association of early-onset myocardial infarction with single
21 nucleotide polymorphisms and copy number variants. *Nat Genet* 2009; **41**:
22 334-341.
- 23 22. Das S, Forer L, Schonherr S *et al*: Next-generation genotype imputation
24 service and methods. *Nat Genet* 2016; **48**: 1284-1287.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24

23. Loh PR, Danecek P, Palamara PF *et al*: Reference-based phasing using the Haplotype Reference Consortium panel. *Nat Genet* 2016; **48**: 1443-1448.
24. The 1000 Genomes Project Consortium: An integrated map of genetic variation from 1,092 human genomes. *Nature* 2012; **491**: 56-65.
25. Browning SR, Browning BL: Rapid and accurate haplotype phasing and missing-data inference for whole-genome association studies by use of localized haplotype clustering. *Am J Hum Genet* 2007; **81**: 1084-1097.
26. Rodan LH, Cohen J, Fatemi A *et al*: A novel neurodevelopmental disorder associated with compound heterozygous variants in the huntingtin gene. *Eur J Hum Genet* 2016; **24**: 1833.
27. Ostergaard ME, Southwell AL, Kordasiewicz H *et al*: Rational design of antisense oligonucleotides targeting single nucleotide polymorphisms for potent and allele selective suppression of mutant Huntingtin in the CNS. *Nucleic Acids Res* 2013; **41**: 9634-9650.
28. Pfister EL, Kennington L, Straubhaar J *et al*: Five siRNAs targeting three SNPs may provide therapy for three-quarters of Huntington's disease patients. *Curr Biol* 2009; **19**: 774-778.
29. Zhang Y, Engelman J, Friedlander RM: Allele-specific silencing of mutant Huntington's disease gene. *J Neurochem* 2009; **108**: 82-90.
30. Lombardi MS, Jaspers L, Spronkmans C *et al*: A majority of Huntington's disease patients may be treatable by individualized allele-specific RNA interference. *Exp Neurol* 2009; **217**: 312-319.

1 **Figure Legends**

2

3 **Figure 1. Definitions and sequence relationships of *HTT* haplotypes.**

4 (A) Twenty-one SNPs, one 3bp indel (rs149109767, alleles R-reference & D-deletion)
5 and the CAG repeat polymorphism are shown at their genomic locations relative to that
6 of the *HTT* RefSeq transcript (NM_002111). Genotype at each marker on each of 16
7 *HTT* haplotypes, defined in the text, is shown above the marker. Haplotypes are
8 ordered based upon a neighbor-joining method (p-distance model) in a dendrogram
9 with two main branches, each with different sizes of sub-clusters. Alleles in red
10 represent differences from hap.01, the most frequent haplotype on CAG-expanded HD
11 chromosomes.

12 (B) Consensus alleles of 10 exon SNPs and 10 intron SNPs that showed the biggest
13 cumulative heterozygosity were determined for each haplotype based on 1000
14 Genomes Project data. A consensus allele for a given SNP site represents the most
15 frequent allele among a collection of chromosomes with same haplotpye. Since hap.10
16 is not present in 1000 Genomes data (Phase 1), hap.10 was excluded in this analysis.
17 Subsequently, alleles of SNPs that show variable alleles in 15 haplotypes and alleles of
18 two exon SNPs that were used to define the haplotypes are indicated. SNPs in orange
19 and black font colors represent SNPs on exons and introns of RefSeq NM_002111,
20 respectively.

21

22 **Figure 2. Correspondences of haplotypes and haplogroups.**

23 Based on S. Table 4, correspondences of haplotypes to haplogroups were summarized.
24 "hap.other" and "Other" were excluded to focus on distinct haplotypes. Thickness of an

1 arrow represents relative proportion of a specific haplotype-haplogroup correspondence
2 for a given haplotype. For example, most of hap.02 is classified as haplogroup A2b,
3 and a small portion of hap.02 is classified as haplogroup A2a. Actual haplotype-
4 haplogroup correspondence data can be found in S. Table 4.

5

6 **Figure 3. Frequencies of haplotypes in HD disease and normal chromosomes.**

7 HD subjects carrying one expanded and one normal chromosome were included in this
8 analysis to estimate overall frequencies of haplotypes. From haplotypes
9 probabilisitically determined based on union set of 74 SNPs, we used our haplotype
10 definitions to calssify each chromosome. Subsequently, frequencies of our haplotypes
11 in HD disease chromosomes (A) and normal chromosomes in HD subjects (B) were
12 calcuated and summarized.

13

14

15

1 **Table 1. *HTT* haplotypes of publicly available cell resources.**

Sample ID	Gender	Chromosome 1		Chromosome 2	
		CAG	Haplotype	CAG	Haplotype
GM02077	Female	44	hap.01	17	hap.03
GM02079	Female	44	hap.01	21	hap.02
GM02147	Male	43	hap.06	15	hap.08
GM02149	Female	18	hap.08	18	hap.13
GM02151	Female	45	hap.06	18	hap.08
GM06274	Female	43	hap.06	16	hap.11
GM02153	Female	32	hap.other	16	hap.14
GM02155	Female	17	hap.12	16	hap.14
GM02157	Female	17	hap.12	16	hap.14
GM02159	Female	32	hap.other	17	hap.12
GM02161	Male	17	hap.12	16	hap.14
GM02163	Male	47	hap.other	32	hap.other
GM02165	Male	44	hap.01	33	hap.other
GM02177	Male	44	hap.01	18	hap.08
GM02183	Female	33	hap.other	18	hap.08
GM03621	Female	60	hap.01	18	hap.other
GM02171	Female	17	hap.14	17	hap.08
GM02173	Female	44	hap.01	17	hap.08
GM02175	Male	20	hap.other	17	hap.14
GM00305*	Female	43	hap.other	17	hap.12
GM01061*	Male	44	hap.01	18	hap.other

GM01085*	Male	44	hap.other	21	hap.other
GM03814*	Female	28	hap.01	22	hap.02
GM03864	Female	45	hap.01	15	hap.15
GM03866	Female	43	hap.01	21	hap.02
GM03868	Female	47	hap.01	17	hap.11
GM03872	Male	21	hap.02	15	hap.15
GM04188	Female	16	hap.14	15	hap.08
GM04200	Male	45	hap.07	16	hap.14
GM04281	Female	78	hap.03	17	hap.11
GM04689	Female	45	hap.03	16	hap.other
GM04717	Female	42	hap.03	29	hap.02
GM04773	Female	38	hap.03	15	hap.08
GM04777	Male	44	hap.03	18	hap.03
GM04805	Female	29	hap.other	15	hap.08
GM04807	Male	50	hap.03	38	hap.03
GM04887	Female	45	hap.03	21	hap.02
GM04287	Male	51	hap.03	17	hap.03
GM04687	Female	50	hap.03	15	hap.15
GM04723	Female	69	hap.03	15	hap.08
GM04729*	Female	17	hap.08	17	hap.11
GM04797	Female	17	hap.08	15	hap.15
GM04849	Female	51	hap.03	45	hap.03
GM08330*	Male	17	hap.08	17	hap.other
GM21756*	Female	69	hap.03	15	hap.08

GM21757*	Male	63	hap.03	16	hap.12
GM01168*	Male	48	hap.other	25	hap.11
GM01169*	Male	44	hap.01	17	hap.08
GM01187*	Male	46	hap.01	18	hap.12
GM05539*	Male	96	hap.03	22	hap.02
ND31551*	Male	39	hap.08	18	hap.13
ND33947*	Female	40	hap.01	18	hap.08
ND30013*	Male	43	hap.09	17	hap.other
ND30259*	Female	38	hap.06	21	hap.other
ND30626	Male	41	hap.02	17	hap.11
ND31038	Female	44	hap.02	19	hap.02
ND29970*	Male	40	hap.04	17	hap.other
ND33392*	Female	56	hap.07	17	hap.08
GM09197*	Male	180	hap.01	18	hap.08

1

2 Computational haplotype phasing analysis was performed using 21 SNPs and biallele
3 coding genotype of *HTT* CAG repeats as described in the method section.
4 Subsequently, phased alleles of CAG repeat genotype and family relationships (refer to
5 S. Table 1) were considered to determine and confirm the phase of CAG repeat size
6 and *HTT* haplotype. In some cases, genotypes and haplotypes of relatives were not
7 available, or family relationship was not informative in determining the phase. In such
8 cases, only probabilistic population phasing results are shown (samples marked by *).

9 * population probabilistic phasing for HD; arbitrary phasing for controls.

10

1 **Table 2. Sample information and phased haplotypes of hESC from of Genea**
2 **Biocells.**

Sample ID	Gender	Sibship	Disease		Normal	
			Chromosome		Chromosome	
			CAG	haplotype	CAG	haplotype
GENEA017	Male		40	hap.01	12	hap.other
GENEA018	Female		46	hap.01	17	hap.other
GENEA020	Female		48	hap.01	17	hap.03
GENEA046	Female		45	hap.01	23	hap.02
GENEA089	Female	Sib to 090 & 091	42	hap.02	18	hap.08
GENEA090	Female	Sib to 089 & 091	46	hap.02	19	hap.02
GENEA091	Female	Sib to 089 & 090	41	hap.02	19	hap.02

3
4 Computational haplotype phasing analysis was performed to determine haplotypes and
5 corresponding CAG repeat sizes of embryonic stem cell lines from Genea Biocells. Each
6 phased disease (CAG-expanded) or normal chromosome consists of two components:
7 CAG and *HTT*. This collection include related individuals (siblings) as shown in the
8 Sibship column.

9
10
11
12

1 **Table 3. Frequency of combined haplotype/haplogroup system backbones on**
2 **CAG- expanded and normal chromosomes in European HD subjects.**

3

'Super'-haplotype of CAG-expanded chromosomes			
# HD disease			
Haplotype	Haplogroup	chromosomes	Percent
hap.01	A1a	1556	38.16%
hap.02	A2b	553	13.56%
hap.03	A3a	323	7.92%
hap.02	A2a	291	7.14%
hap.05	A1a	164	4.02%
hap.08	C1	136	3.33%
hap.06	A2b	107	2.62%
hap.06	A2a	60	1.47%
hap.11	A4b	3	0.07%
hap.12	A5a	3	0.07%
hap.15	A5b	2	0.05%
hap.12	A1a	1	0.02%
Sum			78.45%

**Haplotype "hap.other" or haplogroup "Other" categories of CAG-expanded
chromosomes**

# HD disease			
Haplotype	Haplogroup	chromosomes	Percent

hap.other	Other	380	9.32%
hap.04	Other	134	3.29%
hap.07	Other	134	3.29%
hap.12	Other	51	1.25%
hap.02	Other	33	0.81%
hap.06	Other	30	0.74%
hap.other	B2	22	0.54%
hap.09	Other	18	0.44%
hap.11	Other	17	0.42%
hap.14	Other	14	0.34%
hap.01	Other	13	0.32%
hap.16	Other	10	0.25%
hap.other	C5	7	0.17%
hap.05	Other	4	0.10%
hap.03	Other	3	0.07%
hap.other	B1b	3	0.07%
hap.other	A2b	2	0.05%
hap.other	A5a	2	0.05%
hap.08	Other	1	0.02%
hap.other	A1a	1	0.02%
Sum			21.55%

'Super'-haplotype of normal chromosomes

Haplotype	Haplogroup	# Normal chromosomes	Percent
-----------	------------	----------------------	---------

hap.08	C1	1034	25.36%
hap.03	A3a	479	11.75%
hap.02	A2b	242	5.93%
hap.11	A4b	197	4.83%
hap.02	A2a	153	3.75%
hap.13	B1a	141	3.46%
hap.01	A1a	82	2.01%
hap.12	A5a	81	1.99%
hap.12	A4a	76	1.86%
hap.15	A5b	71	1.74%
hap.06	A2a	63	1.54%
hap.14	C6	61	1.50%
hap.06	A2b	53	1.30%
hap.16	C4b	51	1.25%
hap.16	C4a	47	1.15%
hap.05	A1a	42	1.03%
hap.14	C2	26	0.64%
hap.10	A1a	1	0.02%
hap.12	A1a	1	0.02%
Sum			71.14%

**Haplotype "hap.other" or haplogroup "Other" categories of normal
chromosomes**

Haplotype	Haplogroup	# Normal chromosomes	Percent
------------------	-------------------	-----------------------------	----------------

hap.other	Other	559	13.71%
hap.12	Other	187	4.59%
hap.14	Other	96	2.35%
hap.other	B1b	73	1.79%
hap.07	Other	54	1.32%
hap.11	Other	47	1.15%
hap.other	C8	36	0.88%
hap.03	Other	21	0.51%
hap.06	Other	16	0.39%
hap.other	C5	14	0.34%
hap.16	Other	11	0.27%
hap.04	Other	10	0.25%
hap.other	C7	10	0.25%
hap.02	Other	9	0.22%
hap.01	Other	8	0.20%
hap.other	C1	6	0.15%
hap.08	Other	4	0.10%
hap.13	Other	3	0.07%
hap.other	B1a	3	0.07%
hap.15	Other	2	0.05%
hap.other	A5a	2	0.05%
hap.other	B2	2	0.05%
hap.other	A2a	1	0.02%
hap.other	A4b	1	0.02%

hap.other	A5b	1	0.02%
Sum			28.84%

1

2 Phased haplotypes of subjects (4078 heterozygous HD) were grouped into HD disease

3 chromosomes and normal chromosomes. Subsequently, the frequency of each

4 combined haplotype/haplogroup (i.e., 'super'-haplotype) was calculated for HD disease

5 and normal chromosomes. Frequency and corresponding percentage value of each

6 'super'-haplotype were based on 1) haplotypes not involving "hap.other" or "Other" and

7 2) haplotypes involving "hap.other" or "Other".

8

1 **Table 4. Fully-defined diplotypes in HD subjects with European ancestry.**

2

HD		Normal		# Subjects	% of 4078 subjects
hap.01	A1a	hap.08	C1	350	8.58%
hap.01	A1a	hap.03	A3a	179	4.39%
hap.02	A2b	hap.08	C1	133	3.26%
hap.01	A1a	hap.02	A2b	121	2.97%
hap.03	A3a	hap.08	C1	114	2.80%
hap.02	A2b	hap.03	A3a	79	1.94%
hap.02	A2a	hap.08	C1	75	1.84%
hap.01	A1a	hap.11	A4b	70	1.72%
hap.01	A1a	hap.01	A1a	69	1.69%
hap.01	A1a	hap.02	A2a	69	1.69%
hap.01	A1a	hap.13	B1a	49	1.20%
hap.08	C1	hap.08	C1	47	1.15%
hap.02	A2a	hap.03	A3a	40	0.98%
hap.03	A3a	hap.03	A3a	40	0.98%
hap.05	A1a	hap.08	C1	38	0.93%
hap.01	A1a	hap.05	A1a	34	0.83%
hap.06	A2b	hap.08	C1	34	0.83%
hap.01	A1a	hap.12	A4a	33	0.81%
hap.02	A2b	hap.02	A2b	29	0.71%
hap.01	A1a	hap.15	A5b	28	0.69%
hap.02	A2b	hap.11	A4b	28	0.69%

hap.01	A1a	hap.12	A5a	27	0.66%
hap.01	A1a	hap.06	A2b	24	0.59%
hap.02	A2b	hap.02	A2a	24	0.59%
hap.01	A1a	hap.06	A2a	23	0.56%
hap.01	A1a	hap.14	C6	23	0.56%
hap.02	A2b	hap.13	B1a	19	0.47%
hap.03	A3a	hap.11	A4b	18	0.44%
hap.03	A3a	hap.13	B1a	17	0.42%
hap.06	A2b	hap.03	A3a	17	0.42%
hap.01	A1a	hap.16	C4a	16	0.39%
hap.05	A1a	hap.03	A3a	16	0.39%
hap.06	A2a	hap.03	A3a	16	0.39%
hap.01	A1a	hap.16	C4b	15	0.37%
hap.02	A2a	hap.02	A2b	15	0.37%
hap.02	A2b	hap.16	C4a	14	0.34%
hap.02	A2b	hap.14	C6	13	0.32%
hap.06	A2a	hap.08	C1	13	0.32%
hap.02	A2a	hap.11	A4b	12	0.29%
hap.02	A2b	hap.12	A5a	12	0.29%
hap.03	A3a	hap.12	A4a	11	0.27%
hap.05	A1a	hap.02	A2a	11	0.27%
hap.02	A2a	hap.02	A2a	10	0.25%
hap.02	A2b	hap.06	A2a	10	0.25%
hap.02	A2b	hap.06	A2b	10	0.25%

hap.05	A1a	hap.02	A2b	10	0.25%
hap.05	A1a	hap.11	A4b	10	0.25%
hap.06	A2b	hap.02	A2b	10	0.25%
hap.02	A2b	hap.12	A4a	9	0.22%
hap.08	C1	hap.11	A4b	9	0.22%
hap.08	C1	hap.13	B1a	9	0.22%
hap.02	A2b	hap.15	A5b	8	0.20%
hap.03	A3a	hap.16	C4b	8	0.20%
hap.02	A2a	hap.06	A2a	7	0.17%
hap.02	A2a	hap.13	B1a	7	0.17%
hap.03	A3a	hap.12	A5a	7	0.17%
hap.05	A1a	hap.15	A5b	7	0.17%
hap.08	C1	hap.12	A5a	7	0.17%
hap.05	A1a	hap.12	A4a	6	0.15%
hap.05	A1a	hap.13	B1a	6	0.15%
hap.06	A2b	hap.02	A2a	6	0.15%
hap.01	A1a	hap.14	C2	5	0.12%
hap.02	A2a	hap.16	C4a	5	0.12%
hap.05	A1a	hap.12	A5a	5	0.12%
hap.08	C1	hap.16	C4b	5	0.12%
hap.02	A2a	hap.12	A4a	4	0.10%
hap.02	A2a	hap.15	A5b	4	0.10%
hap.02	A2b	hap.16	C4b	4	0.10%
hap.03	A3a	hap.14	C6	4	0.10%

hap.03	A3a	hap.15	A5b	4	0.10%
hap.06	A2a	hap.11	A4b	4	0.10%
hap.08	C1	hap.15	A5b	4	0.10%
hap.02	A2a	hap.14	C2	3	0.07%
hap.02	A2a	hap.16	C4b	3	0.07%
hap.02	A2b	hap.05	A1a	3	0.07%
hap.03	A3a	hap.06	A2a	3	0.07%
hap.03	A3a	hap.14	C2	3	0.07%
hap.05	A1a	hap.14	C6	3	0.07%
hap.05	A1a	hap.16	C4b	3	0.07%
hap.06	A2a	hap.02	A2b	3	0.07%
hap.06	A2b	hap.11	A4b	3	0.07%
hap.08	C1	hap.14	C2	3	0.07%
hap.08	C1	hap.14	C6	3	0.07%
hap.02	A2a	hap.06	A2b	2	0.05%
hap.02	A2a	hap.12	A5a	2	0.05%
hap.02	A2a	hap.14	C6	2	0.05%
hap.02	A2b	hap.14	C2	2	0.05%
hap.03	A3a	hap.16	C4a	2	0.05%
hap.05	A1a	hap.06	A2b	2	0.05%
hap.05	A1a	hap.14	C2	2	0.05%
hap.06	A2a	hap.13	B1a	2	0.05%
hap.06	A2a	hap.16	C4b	2	0.05%
hap.06	A2b	hap.12	A4a	2	0.05%

hap.06	A2b	hap.16	C4b	2	0.05%
hap.08	C1	hap.16	C4a	2	0.05%
hap.01	A1a	hap.12	A1a	1	0.02%
hap.03	A3a	hap.02	A2a	1	0.02%
hap.05	A1a	hap.16	C4a	1	0.02%
hap.06	A2a	hap.06	A2a	1	0.02%
hap.06	A2a	hap.12	A5a	1	0.02%
hap.06	A2b	hap.06	A2a	1	0.02%
hap.06	A2b	hap.06	A2b	1	0.02%
hap.06	A2b	hap.12	A5a	1	0.02%
hap.06	A2b	hap.13	B1a	1	0.02%
hap.06	A2b	hap.15	A5b	1	0.02%
hap.06	A2b	hap.16	C4a	1	0.02%
hap.08	C1	hap.12	A4a	1	0.02%
hap.12	A5a	hap.11	A4b	1	0.02%
hap.12	A5a	hap.13	B1a	1	0.02%
hap.15	A5b	hap.13	B1a	1	0.02%
Total				2291	56.18%
Total homozygous diplotype				197	4.83%

1

2 To determine the proportion of HD subjects eligible for allele-specific gene targeting

3 approaches, diplotype of each HD subject (i.e., HD disease chromosome and normal

4 chromosome in a given HD subject) in our data was constructed based on 'super'-

5 haplotype system. Subsequently, the frequency of each unique diplotype was calculated.

- 1 Diplotypes in bold and italic represent HD subjects who carry the same haplotypes for
- 2 disease and normal chrommosomes.
- 3

1 **Supplementary Tables**

2 S. Table 1. Familial relationships and other IDs of publicly available cell resources

3 S. Table 2. Description of SNPs used in this study.

4 S. Table 3. *HTT* haplotypes and haplogroup subtypes in the 1000 Genomes Project
5 Data.

6 S. Table 4. Distribution of *HTT* haplogroup subtypes relative to *HTT* haplotypes.

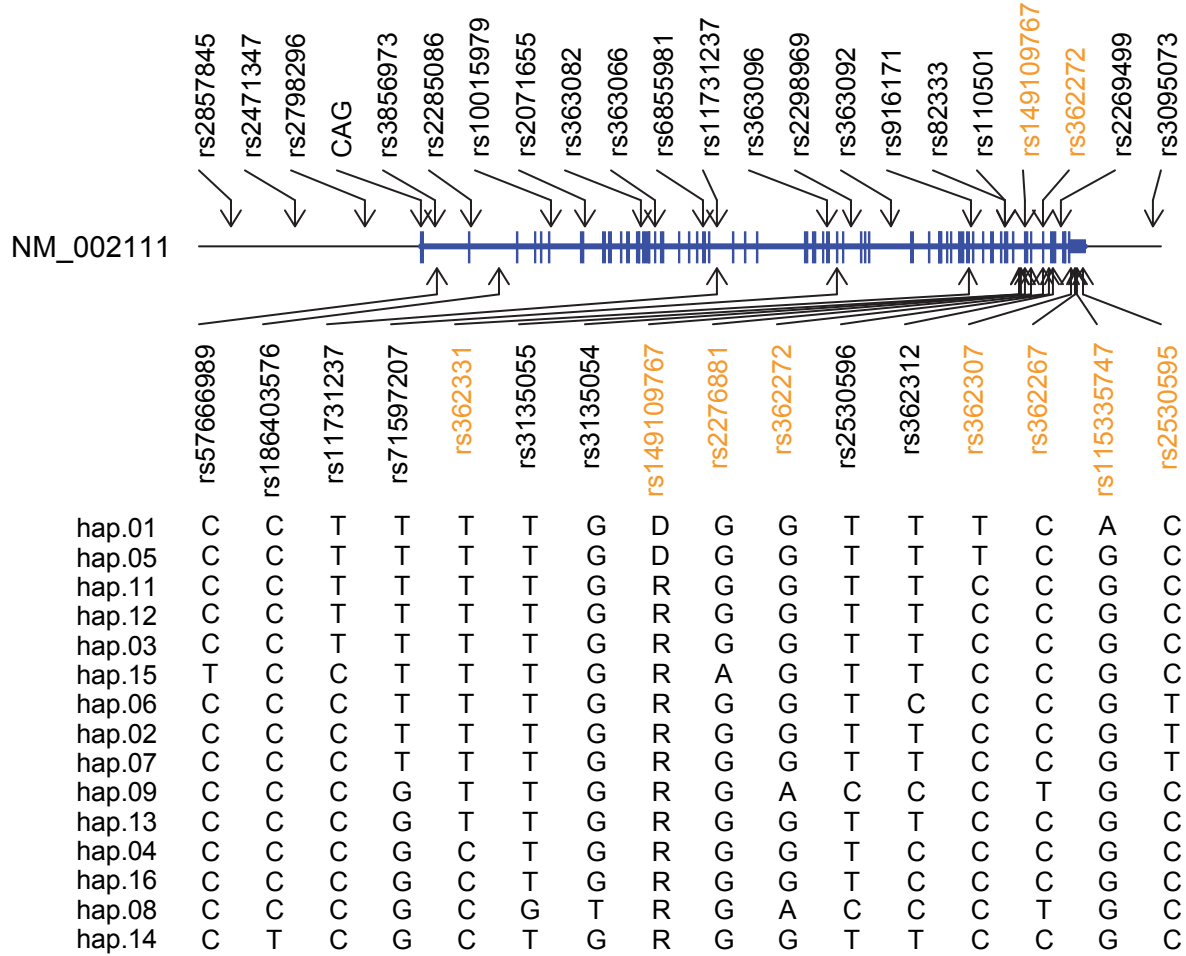
7 S. Table 5. Distribution of *HTT* haplotypes relative to *HTT* haplogroup subtypes.

8

A

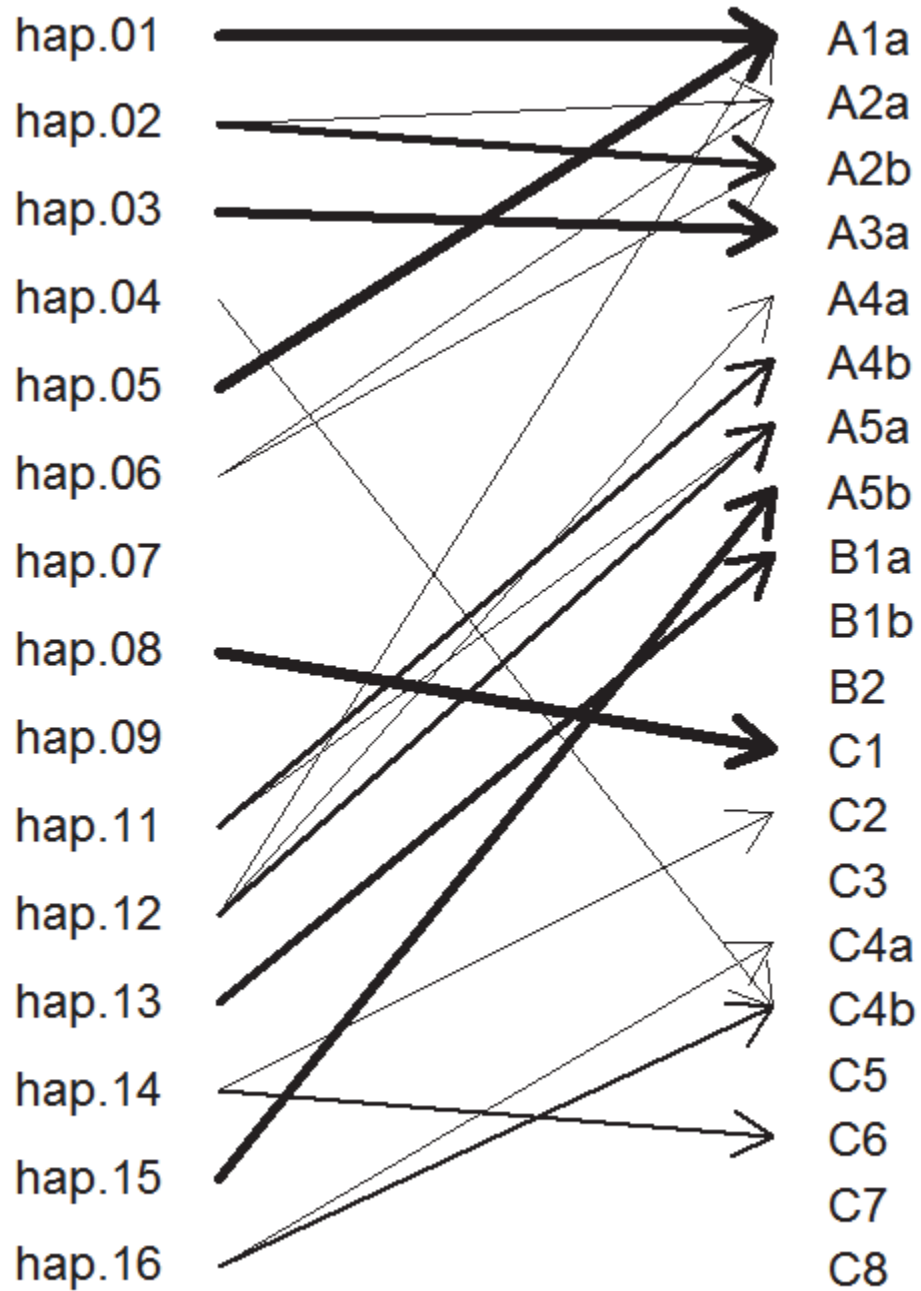


B

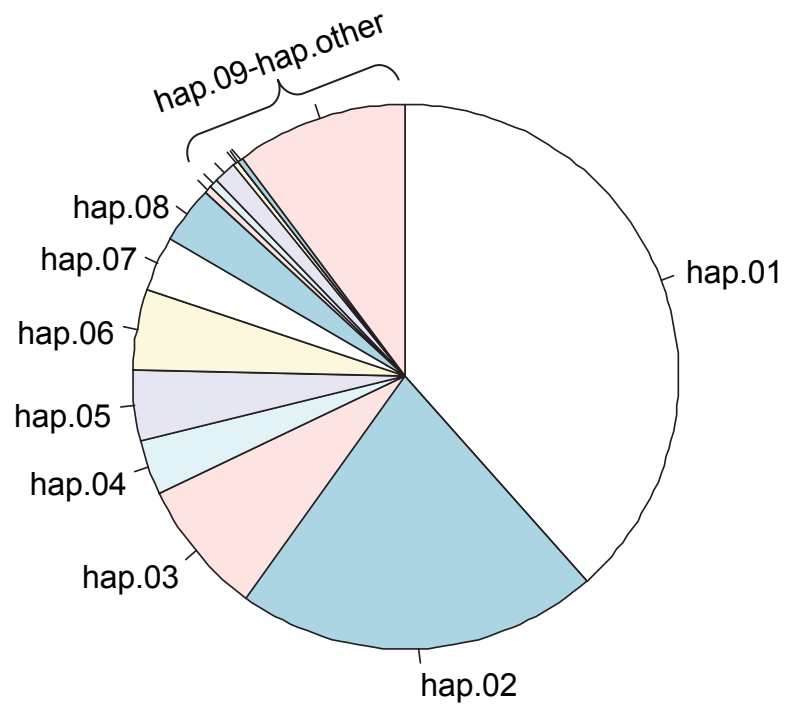


Haplotype

Haplogroup



A. Disease chromosomes



B. Normal chromosomes

