

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository: <https://orca.cardiff.ac.uk/id/eprint/165907/>

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

Citation for final published version:

Hong, Eun Pyo, Ramos, Eliana Marisa, Aziz, N. Ahmad, Massey, Thomas H. , McAllister, Branduff, Lobanov, Sergey , Jones, Lesley , Holmans, Peter , Kwak, Seung, Orth, Michael, Ciosi, Marc, Lomeikaite, Vilija, Monckton, Darren G., Long, Jeffrey D., Lucente, Diane, Wheeler, Vanessa C., Gillis, Tammy, MacDonald, Marcy E., Sequeiros, Jorge, Gusella, James F. and Lee, Jong-Min 2024. Modification of Huntington's disease by short tandem repeats. *Brain Communications* 10.1093/braincomms/fcae016

Publishers page: <http://dx.doi.org/10.1093/braincomms/fcae016>

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.



Modification of Huntington's disease by short tandem repeats

Eun Pyo Hong,^{1,2,3,*} Eliana Marisa Ramos,^{1,2,*} N. Ahmad Aziz,^{4,5,*} Thomas H. Massey,⁶
Branduff McAllister,⁶ Sergey Lobanov,⁶ Lesley Jones,⁶ Peter Holmans,⁶ Seung Kwak,⁷
Michael Orth,⁸ Marc Ciosi,⁹ Viliija Lomeikaite,⁹ Darren G. Monckton,⁹ Jeffrey D. Long,¹⁰
Diane Lucente,¹ Vanessa C. Wheeler,^{1,2} Tammy Gillis,¹ Marcy E. MacDonald,^{1,2,3} Jorge
Sequeiros,^{11,12} James F. Gusella,^{1,3,13} and Jong-Min Lee^{1,2,3}

* Eun Pyo Hong, Eliana Marisa Ramos, and N. Ahmad Aziz contributed equally to this work.

Abstract

Expansions of glutamine-coding CAG trinucleotide repeats cause a number of neurodegenerative diseases, including Huntington's disease (HD) and several of the spinocerebellar ataxias (SCAs). In general, age-at-onset of the polyglutamine diseases is inversely correlated with the size of the respective inherited expanded CAG repeat. Expanded CAG repeats are also somatically unstable in certain tissues, and age-at-onset of HD corrected for individual *HTT* CAG repeat length (i.e., residual age-at-onset), is modified by repeat instability-related DNA maintenance/repair genes as demonstrated by recent genome-wide association studies (GWAS). Modification of one polyglutamine disease (e.g., HD) by the repeat length of another (e.g., *ATXN3*, CAG expansions in which cause SCA3) has also been hypothesized. Consequently, we determined whether age-at-onset in HD is modified by the CAG repeats of other polyglutamine disease genes. We found that the CAG measured repeat sizes of other polyglutamine disease genes were polymorphic in HD participants but did not influence HD age-at-onset. Additional analysis focusing specifically on *ATXN3* in a larger sample set ($n=1,388$) confirmed the lack of association between HD residual age-at-

<https://mc.manuscriptcentral.com/braincom>

1
2
3 onset and *ATXN3* CAG repeat length. Additionally, neither our HD onset modifier GWAS
4 single nucleotide polymorphism (SNP) data nor imputed short tandem repeat (STR) data
5 supported involvement of other polyglutamine disease genes in modifying HD. By contrast,
6 our GWAS based on imputed STRs revealed significant modification signals for other
7 genomic regions. Together, our STR GWAS show that modification of HD is associated with
8 STRs that do not involve other polyglutamine disease-causing genes, refining the landscape
9 of HD modification and highlighting the importance of rigorous data analysis, especially in
10 genetic studies testing candidate modifiers.
11
12
13
14
15
16
17
18
19
20
21
22
23
24

25 **Keywords: Huntington's disease; genetic modification; polyglutamine disease; *ATXN3*;**
26 **short tandem repeat**
27
28
29
30
31
32

33 **Author affiliations:**

34
35
36 ¹ Molecular Neurogenetics Unit, Center for Genomic Medicine, Massachusetts General
37 Hospital, Boston, MA 02114, USA
38

39
40
41 ² Department of Neurology, Harvard Medical School, Boston, MA 02115, USA
42

43 ³ Medical and Population Genetics Program, the Broad Institute of M.I.T. and Harvard,
44 Cambridge, MA 02142, USA
45

46
47
48 ⁴ German Center for Neurodegenerative Diseases, Venusberg-Campus 1/99, 53127 Bonn,
49 Germany
50

51
52
53 ⁵ Department of Neurology, Faculty of Medicine, University of Bonn, Bonn, Germany
54

55
56
57 ⁶ Centre for Neuropsychiatric Genetics and Genomics, Division of Psychological Medicine
58 and Clinical Neurosciences, School of Medicine, Cardiff University, Cardiff CF24 4HQ, UK
59
60

1
2
3 ⁷ CHDI Foundation, Princeton, NJ 08540, USA
4

5 ⁸ University Hospital of Old Age Psychiatry and Psychotherapy, Bern University, Bern,
6
7
8 Switzerland
9

10 ⁹ School of Molecular Biosciences, College of Medical, Veterinary and Life Sciences,
11
12
13 University of Glasgow, Glasgow G12 8QQ, UK
14

15 ¹⁰ Department of Psychiatry, Carver College of Medicine and Department of Biostatistics,
16
17
18 College of Public Health, University of Iowa, Iowa City, Iowa 52242, USA
19

20 ¹¹ UnIGENE, IBMC - Institute for Molecular and Cell Biology, i3S - Instituto de Investigação
21
22
23 e Inovação em Saúde, Universidade do Porto, Porto, Portugal
24

25 ¹² ICBAS School of Medicine and Biomedical Sciences, Univ. Porto, Portugal
26

27 ¹³ Department of Genetics, Blavatnik Institute, Harvard Medical School, Boston, MA 02115,
28
29
30 USA
31

32
33
34
35 Correspondence to: Jong-Min Lee, Ph.D.
36

37
38 Center for Genomic Medicine, Massachusetts General Hospital
39

40
41 185 Cambridge Street, Boston, MA 02114, USA
42

43
44 E-mail: jlee51@mgh.harvard.edu
45
46
47
48
49

50 **Short title:** Modification of HD by short tandem repeat
51
52
53
54
55
56
57
58
59
60

Introduction

Expansions of glutamine-encoding CAG repeats cause at least nine neurodegenerative diseases, including Huntington's disease (HD; MIM #143100), several spinocerebellar ataxias (SCAs), and dentatorubral-pallidoluysian atrophy (DRPLA).¹⁻³ The polyglutamine expansion diseases exhibit differences in pathogenesis, susceptible brain regions, and disease symptoms.^{4,5} However, they share a common feature of inverse correlation between age-at-onset and the length of the causative expanded CAG repeat,⁵⁻¹² indicating that increases in CAG repeat size result in accelerated pathogenesis. However, variance in age-at-onset is not fully explained by the glutamine-encoding CAG repeat length alone. For example, the residual variance in HD age-at-onset (i.e., not due to CAG repeat size) showed heritability,¹³ prompting genome-wide association studies (GWAS) to identify genetic modifiers of HD. Several genetic loci discovered to influence HD age-at-onset harbor DNA repair genes, such as *MLH1*, *MSH3*, and *FANL*.^{14,15} These genes have been associated with somatic instability of *HTT* CAG repeats in humans and model systems.¹⁶⁻²⁵ There is also evidence that these and other DNA repair genes may influence somatic CAG repeat expansions and impact other repeat expansion disorders.^{16-19,21,26-39} The striatum, which is severely affected in HD, shows the highest levels of somatic *HTT* CAG repeat expansion; however, expanded repeats in other polyglutamine diseases can also undergo CAG expansion in this brain region.⁴⁰⁻⁴⁴ Interestingly, candidate gene studies have reported modification of HD by normal CAG repeats in *ATXN3* (expansions in which are responsible for SCA3),⁴⁵ and conversely modification of SCA3 by the normal *HTT* CAG repeat.⁴⁶ Although *HTT* and *ATXN3* have potential roles in the DNA damage response,^{37,47} neither *HTT* nor *ATXN3* are known to be directly involved in DNA repair. Therefore, this mutual modification has suggested the possibility of a novel mechanism underlying polyglutamine diseases. Consequently, using a

1
2
3 variety of genomic data, we set out to determine whether HD is modified by CAG repeat
4
5 length in other polyglutamine disease genes or other short tandem repeats (STR).
6
7
8
9
10
11
12

13 **Materials and methods**

14 15 16 17 **Study subjects**

18
19 To identify genetic modifiers of HD motor onset, a total of 9,058 HD subjects (carrying
20
21 inherited CAG 40 to 55) of European ancestry were previously analyzed in our GWAS.¹⁵
22
23 Among those HD subjects, we analyzed participants of the COHORT study
24
25 (<https://clinicaltrials.gov/ct2/show/NCT00313495>) to test association between HD age-at-
26
27 onset and experimentally determined (i.e., genotyped) CAG repeat lengths of other
28
29 polyglutamine disease-causing genes ($n=606$). For *ATXN3*-focused analysis, we also
30
31 analyzed participants of the REGISTRY study
32
33 (<https://clinicaltrials.gov/ct2/show/NCT01590589>) ($n=885$). Details of study approval,
34
35 genotyping, determination of CAG repeat size, and calculation of residual age-at-onset are
36
37 described elsewhere.¹⁵
38
39
40
41
42
43
44
45
46

47 **Determination of the CAG repeats in the COHORT and** 48 49 **REGISTRY samples**

50
51 We determined the sizes of CAG repeats of *ATNI* [dentatorubral-pallidoluysian atrophy;
52
53 DRPLA MIM #125370], *ATXN1* [spinocerebellar ataxia type 1; SCA1 MIM #164400],
54
55 *ATXN2* [spinocerebellar ataxia type 2; SCA2 MIM #183090], *ATXN3* [Machado-Joseph
56
57 disease aka spinocerebellar ataxia type 3; SCA3 MIM #109150], *CACNA1A* [spinocerebellar
58
59
60

1
2
3 ataxia type 6; SCA6 MIM #183086], and *TBP* [spinocerebellar ataxia type 17; SCA17 MIM
4 #607136] in the participants of the COHORT study ($n=606$). After quality control, 551, 502,
5
6
7 604, 503, 497 and 483 COHORT samples were analyzed for *ATNI*, *ATXN1*, *ATXN2*, *ATXN3*,
8
9
10 *CACNA1A*, and *TBP*, respectively. In addition, REGISTRY samples were analyzed to
11
12 determine the CAG repeat sizes of *ATXN3* ($n=885$). CAG repeat lengths for each
13
14 polyglutamine disease-causing genes were determined by polymerase chain reaction (PCR)
15
16 assays, using fluorescently labelled primers with minor modifications.⁴⁸ PCR products were
17
18 resolved by an ABI PRISM 3730XL automated DNA Sequencer (Applied Biosystems) and
19
20 analyzed using GeneMapper version 3.7 software. A set of genomic DNA standard samples
21
22 were also sequenced for each polyglutamine disease-causing repeat and used as references of
23
24 CAG allele sizes. Expansions above the non-disease associated repeat range were sequenced
25
26 after gel separation to further confirm the number of CAGs and the presence of CAA/CAT
27
28 interruptions.
29
30
31
32
33
34
35

36 **Analysis to determine the modification of HD by the CAG repeats** 37 38 39 **of other polyglutamine disease-causing genes** 40 41

42 Residual age-at-onset of HD, representing age-at-onset corrected for individual *HTT* CAG
43
44 repeat length, was based on the rater's estimation of onset age of motor symptoms and the
45
46 uninterrupted *HTT* CAG repeat size.^{12,15} For example, an HD subject with a positive residual
47
48 age-at-onset of 5 means that the individual developed motor onset 5 years later than expected
49
50 based on his or her uninterrupted CAG repeat size. To determine whether CAG repeat sizes
51
52 of other polyglutamine disease-causing genes modify HD age-at-onset, we used the same
53
54 residual age-at-onset phenotype used in our GWAS to identify onset modifiers of HD.^{15,49}
55
56 Briefly, we modeled residual age-at-onset of HD as a function of the CAG repeat of another
57
58
59
60

1
2
3 polyglutamine disease gene with a set of covariates including 4 genotype-based principal
4 components, sex, and study group in a linear regression analysis. To analyze typed CAG
5 repeats of each polyglutamine disease gene, 3 separate linear regression models were
6 constructed to test the longer, the shorter, and the sum of both repeat lengths. To validate
7 previous dichotomous *ATXN3* CAG repeat association analysis,⁴⁵ we also performed a Mann-
8 Whitney *U* test to compare age-at-onset and residual age-at-onset between HD subjects
9 carrying below vs. above the median of longer *ATXN3* CAG repeat (i.e., 23 CAG).
10 Specifically, based on the longer of two *ATXN3* CAG repeats in a given individual, 1) HD
11 subjects carrying 22 or shorter repeats were assigned as below the median group, and 2) HD
12 subjects with 24 or higher repeats were assigned as above the median group for the
13 dichotomous analysis.
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31

32 **SNP association analysis for other polyglutamine disease-causing** 33 **genes in HD modifier GWAS data** 34 35 36 37 38

39 We evaluated the levels of association between residual age-at-onset of HD and other
40 polyglutamine disease-causing genes by checking SNP data available at the GeM Euro 9K
41 website (cegeme.partners.org/gem.euro.9k).⁴⁹ For each gene, we took a region for the RefSeq
42 select transcript, and identified the SNP with the highest significance in our HD subject-based
43 GWA analysis to assess the levels of modification of HD by other polyglutamine disease
44 genes.
45
46
47
48
49
50
51
52
53
54
55

56 **Imputation of STR lengths from GWAS SNP data and association** 57 **analysis** 58 59 60

1
2
3 For the imputation of STR lengths in our GWAS data, we performed quality control analysis
4 of typed genotype data by taking SNPs with call rate > 95% and minor allele frequency >
5 1%. Allele frequencies of SNPs and reference alleles in the typed data set were compared to
6 those of 1000 Genomes Project data to confirm data quality using the conform-gt
7 (<https://faculty.washington.edu/browning/conform-gt.html>). Then, imputation of autosomal
8 STRs was performed by the Beagle program
9 (<https://faculty.washington.edu/browning/beagle/beagle.html>; v4.1) using the 1000 Genomes
10 Project reference panel consisting of SNPs and STRs.^{50,51} Imputed STR data were further
11 filtered by taking tandem repeats located by the 'Tandem Repeats Finder' algorithm⁵² and
12 annotated as "SimpleRepeat" in the UCSC genome browser
13 (<https://genome.ucsc.edu/index.html>). These procedures generated repeat length genotypes of
14 66,154 tandem repeats for the 9,058 HD subjects. We finally selected 58,894 tandem repeats
15 that were polymorphic in our data for the subsequent association analysis. The proportions of
16 repeats of 1, 2, 3, 4, 5, and 6 nucleotide motif were 0.30, 48.89, 9.47, 29.37, 7.84, and 2.38%,
17 accounting for 98.6% of all analyzed tandem repeats. For association analysis, the sum of two
18 repeat sizes was used as the independent variable (which was similar to the additive model of
19 single SNP association analysis) with the same covariates that were used in our SNP GWAS
20 to explain residual age-at-onset.¹⁵

21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60

Data availability

STR GWAS summary data that support the findings of this study are available from the corresponding author, upon reasonable request.

Results

No significant association between HD and CAG repeat lengths of other polyglutamine disease-causing genes

To determine whether age-at-onset in HD is modified by the CAG repeats of other polyglutamine disease-causing genes, we determined directly the length of CAG repeats in *ATN1* (DRPLA), *ATXN1* (SCA1), *ATXN2* (SCA2), *ATXN3* (SCA3), *CACNA1A* (SCA6), and *TBP* (SCA17) in the HD individuals who participated in both the COHORT study (<https://clinicaltrials.gov/ct2/show/NCT00313495>) and our recent HD modifier GWA analysis.¹⁵ The CAG repeat sizes of the polyglutamine disease-causing genes showed distinct distribution patterns. For example, *CACNA1A* and *TBP* showed the smallest and the largest median repeat sizes, while *ATXN2* and *ATXN3*, respectively showed different ranges of repeat lengths despite similar median repeat sizes (Supplementary Figure 1). Next, we performed statistical analyses of HD subjects with European ancestry to determine whether 1) the longer repeat, 2) the shorter repeat, or 3) the sum of the two repeat alleles of other polyglutamine disease-causing genes were associated with residual age-at-onset of HD. As the primary phenotype of the analysis, we used residual age-at-onset of HD motor symptoms representing age-at-onset that was corrected for individual pathogenic *HTT* CAG repeat size. In linear regression models corrected for genetic ancestry and other potential confounding factors, residual age-at-onset of HD was not significantly associated with the longer, the shorter, or the sum of the two repeat alleles of any of the tested genes (Table 1). In contrast, HD age-at-onset was significantly associated the size of expanded *HTT* CAG repeat (p-value, 2.1E-111), consistent with our previous report.¹² As expected, residual age-at-onset, representing onset

1
2
3 age corrected for the length of expanded *HTT* CAG repeat, was not significantly associated
4 with *HTT* CAG repeat size (expanded repeat p-value, 0.6948; normal repeat p-value, 0.7171).
5
6
7
8
9

10 11 **Association analysis of *ATXN3* CAG repeat**

12
13
14
15 It has been proposed that the length of CAG repeat in *ATXN3* is associated with age-at-onset
16 of HD.⁴⁵ In contrast, our initial analysis of the COHORT participants ($n = 503$) did not show
17 statistically significant associations between HD residual age-at-onset and *ATXN3* CAG
18 repeat lengths. To confirm this lack of association in a larger sample set, we also analyzed
19 REGISTRY participants who were also part of our recent GWA study. A total of 1388 (706
20 males and 682 females) unique HD individuals with European ancestry (503 COHORT and
21 885 REGISTRY) were analyzed for *ATXN3* CAG repeats. Consistent with our initial
22 observations, linear regression analyses to test the longer, the shorter, or the sum of the two
23 repeats of *ATXN3* showed no statistically significant associations with residual age-at-onset
24 of motor signs (Fig. 1). We further performed a dichotomous analysis to test whether age-at-
25 onset or residual age-at-onset was significantly different between HD individuals carrying
26 *ATXN3* repeats above and below the median length. As shown in Supplementary Figure 2,
27 these two groups of HD study participants were not significantly different for age-at-onset
28 (Supplementary Figure 2A) or residual age-at-onset (Supplementary Figure 2B), arguing
29 against modification of HD by the *ATXN3* CAG repeats.
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53

54 **No SNP association signals at other polyglutamine disease-causing** 55 **genes in HD onset modifier GWAS data** 56 57 58 59 60

1
2
3 Previously, we discovered genetic modifiers of HD through genome-wide association
4 analysis of SNPs,¹⁵ and subsequently generated a website to make these GWAS results
5 publicly accessible.⁴⁹ We used this resource to check HD modification signals at the
6 polyglutamine disease-causing gene assessed above, plus spinocerebellar ataxia type 7
7 (SCA7, MIM #607640) and spinal and bulbar muscular atrophy (*AR*, MIM #313200). We
8 also included *PPP2R2B* and *DMPK*, expansions of CAG•CTG repeats within which underlie
9 spinocerebellar ataxia type 12 (SCA12, MIM #604326) and myotonic dystrophy type 1
10 (DM1, MIM #160900), respectively, but are not translated into polyglutamine. For each
11 locus, we evaluated the RefSeq select ('ncbiRefSeqSelect' in the UCSC genome browser) as
12 the representative transcript and identified the SNP with the smallest association *P*-value in
13 the transcript region. The top SNPs at other CAG repeat expansion disease-causing genes
14 were relatively infrequent except *AR*. We observed nominally significant *P*-values for
15 association of some loci with HD residual age-at-onset (Supplementary Table 1), but when
16 these were corrected for the gene size and number of SNPs in the region, none remained
17 statistically significant. Together, our HD modifier SNP GWAS data did not support
18 modification of HD by variants in these other CAG repeat expansion disease-causing genes.
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

44 **Genome-wide STR association analysis**

45
46
47 Though our GWAS, representing the largest dataset of HD individuals with genome-wide
48 genotype and phenotype,¹⁵ successfully identified associated SNPs, those analyses could have
49 limited power to assess effects of polymorphic STRs on modification of HD. Recently,
50 methods for imputation of STR lengths from genome-wide SNP data have been developed
51 and further optimized.⁵¹ Therefore, in order to test the association between HD onset and
52 repeat size of polyglutamine disease-causing genes and other STRs, we imputed ~60,000
53
54
55
56
57
58
59
60

1
2
3 STRs for the HD subjects who participated in our HD modifier GWAS ($n=9058$). We then
4 performed association analysis using the same residual age-at-onset phenotype as in our
5 GWAS¹⁵ and STR genotypes as the independent variable with a set of covariates.
6
7 Specifically, we used the sum of the two repeat lengths, which is similar to the additive
8 model in standard SNP analysis. Notably, the genomic regions containing other
9 polyglutamine disease-causing genes showed no significant STR association signals (Fig. 2,
10 white triangles). The imputed CAG repeat sizes of other polyglutamine disease-causing genes
11 were also not significantly associated with HD residual age-at-onset (Table 2) based on the
12 Bonferroni multiple correction method (P -value, $8.5E-7$). We further evaluated the levels of
13 association between age-at-onset of HD and imputed CAG repeat of *ATXN3*. To confirm the
14 quality of STR imputation, we compared the genotyped and the imputed *ATXN3* CAG repeat
15 in 1388 HD individuals where both estimates were available. The longer, the shorter, and the
16 sum of the two alleles of the *ATXN3* repeat showed 74.6%, 87.2%, and 69.2% concordance
17 between experimentally determined and imputed repeat lengths (Supplementary Figure 3).
18 Moreover, more than 90% of the observed differences were fewer than 5 repeats
19 (Supplementary Figure 3D), suggesting relatively high levels of accuracy in STR imputation.
20 Like the analysis using typed data, association analyses to test the longer (Supplementary
21 Figure 4A), the shorter (Supplementary Figure 4B), or the sum of the two STR alleles (data,
22 not shown) in the imputed data did not reveal statistically significant association with HD
23 age-at-onset. Furthermore, HD participants who carry longer *ATXN3* repeats smaller than the
24 median of the longer repeat (i.e., 23 CAG) showed the same age-at-onset and residual age-at-
25 onset compared to those with longer *ATXN3* repeats larger than median (Supplementary
26 Figure 4C and 4D). Together, our genetic analyses using both genotyped and imputed STR
27 data strongly indicated that *ATXN3* CAG repeat length does not modify HD age-at-onset, in
28 contrast to the previous report.
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 In contrast, as shown in Fig. 2, we identified 3 Bonferroni significant modification
4 signals (P -value $< 8.5E-7$) tagged by STRs on chromosome 2, chromosome 8 and
5 chromosome 15, that potentially captured effects of modification by *PMS1*, *RRM2B*, and
6 *FANI*,¹⁵ which were also implicated by SNP association. Similarly, the previously implicated
7 *MLH1* region of chromosome 3 and the *MSH3* region at chromosome 5 showed STR
8 association signals at suggestive significance (P -value, $E-5$). Interestingly, a near significant
9 new signal (uncorrected P -value, $8.7E-7$) was evident on distal chromosome 15q with a GA
10 repeat (chr15:91711532-91711561) in *SV2B* gene, which encodes synaptic vesicle
11 glycoprotein 2B. Still, confirmation analysis by direct repeat genotyping will be required to
12 establish their roles in modifying HD, considering the inherent uncertainty in imputing
13 tandem repeats from SNP data.
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33

34 Discussion

35
36
37 The existence of multiple neurodegenerative disorders associated with lengthened
38 polyglutamine segments in different proteins has implicated common unifying mechanisms
39 of pathogenesis^{3,4,53} and also raised the related possibility that CAG repeat disease genes
40 show functional interaction with one another in modulating disease manifestations. For
41 example, in early studies, SCA3 fasciculations were reported to be associated with normal
42 *ATXN2* CAG repeat length,⁵⁴ while SCA2 onset was reported to be influenced by the normal
43 *CACNA1A* CAG repeat which, when expanded, causes SCA6.⁵⁵ Subsequently, age-at-onset in
44 several SCAs was reported to be influenced by CAG repeat length variation in various
45 polyglutamine disease-causing genes.^{46,56,57} There have also been interactions suggested
46 between coding and non-coding CAG repeats, as the *ATXN1* 31 CAG repeat allele was
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 reported to be enriched in myotonic dystrophy, which is caused by an expanded CTG repeat
4 in the *DMI* 3'- untranslated region of the *DMPK* gene,⁵⁸ although these findings were
5 challenged later.⁵⁹ Specifically in HD, age-at-onset appeared to be modified by the normal
6 CAG repeat in *ATXN3*, which is responsible for causing SCA3 when expanded.⁴⁵ Given the
7 variability observed across these studies and the possibility that such genetic interactions
8 could provide important insights into both underlying disease mechanisms and potential
9 therapeutic directions, we reasoned that the possibility of modification of HD by other CAG
10 repeat disease-causing genes was an important subject for rigorous investigation.

11
12
13
14
15
16
17
18
19
20
21
22
23 Using both a candidate approach based on CAG length and an unbiased SNP-based
24 GWAS, we did not detect any significant influence of polyglutamine disease-causing genes
25 and other CAG repeat expansion disease-causing genes on the age-at-onset of HD. The lack
26 of replication of candidate modifiers of HD has been reported before. For example, candidate
27 studies suggested modification of HD by *ADORA2A*,^{60,61} *ATG7*,⁶² *BDNF*,⁶³ *GRIK2*,⁶⁴
28 *GRIN2A*,⁶⁴⁻⁶⁶ *GRIN2B*,^{64,65} *HAPI*,⁶⁷ *HIP1*,⁶⁴ *LINC01559*,⁶⁴ *NPY2R*,⁶⁸ *PPARGCIA*,⁶⁹⁻⁷² and
29 *UCHL1*.⁷³ However, none of these genes generated significant onset modification signals in
30 our large scale unbiased genetic analysis.^{14,15,49,74} Interestingly, one candidate modifier that
31 showed a trend of association^{66,75} and was replicated by GWAS is *TCERG1*, which harbors a
32 complex coding hexamer repeat that appears to be the source of the influence on HD age-at-
33 onset.⁷⁶ Unfortunately, the hexamer repeat with potential association was not imputed in our
34 data because this repeat was not present in the imputation reference panel that we used.
35 Nevertheless, the lack of replication for most candidates could be due to spurious signals
36 from underpowered studies, confounded by ancestry differences, lack of multiple test
37 correction, and/or outlier effects.^{74,77} Outlier effects are particularly significant when using
38 continuous variables, as we observed that a single data point could change insignificant signal
39 into significant association.¹² This may explain the lack of replication of modifying effects of
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 the *ATXN3* repeat on HD age-at-onset. The overall high rate of failure to replicate reinforces
4 the importance of rigorous data quality control and stringent statistical analysis for the
5 association analysis of human data.
6
7
8
9

10
11 Although our data did not validate the modification of HD by other polyglutamine
12 disease-causing genes, imputed STR data did show Bonferroni significant signals on
13 chromosomes 2, 8, and 15. These significant STRs appear to tag previously identified
14 modifier haplotypes of *PMS1*, *RRM2B*, and *FANI*, which were detected in our SNP-based
15 GWAS. Considering that this original association analysis tested more than 10 million SNPs,
16 the detection of these 3 significant association signals from testing ~60,000 genetic
17 polymorphisms supports the levels of power and efficiency of the STR approach and argues
18 for its use in modifier studies of other disorders. More than sixty diseases are known to be
19 caused by expansions of tandem repeats, and additional disease-causing repeats are being
20 discovered with the advance of genomic technologies.^{78,79} In addition to tandem repeats that
21 cause Mendelian disorders, repeat polymorphism may contribute to the missing heritability of
22 common polygenic disorders.⁸⁰⁻⁸² Importantly, changes in tandem repeats represent one of the
23 major sources of *de novo* mutation with clinical significances.^{83,84} For example, somatically
24 expanded tandem repeats influence disease age of onset and tissue specificity of pathogenic
25 features^{78,79} Furthermore, significant genome-wide excess of tandem repeat mutations has
26 been reported in the autism spectrum disorder,⁸⁵⁻⁸⁷ implying that tandem repeats may have
27 profound effects on human health beyond the well characterized repeat expansion disorders.
28 Many GWAS signals are due to effects on gene expression levels, and recent findings of a
29 role for length and motif composition of the tandem repeats including variable number
30 tandem repeats (VNTRs) in regulating gene expression⁸⁸⁻⁹² suggest that some HD onset
31 modification signals might be caused by altered expression of modifier genes due to
32 polymorphic tandem repeat lengths. Therefore, investigating a potential role for tandem
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 repeats in regulating the expression levels of the *PMS1*, *RRM2B*, and *FANI* modifier genes
4
5 implicated at the chromosome 2, 8 and 15 loci may reveal an important underlying source of
6
7 HD modification. Finally, the STR association signal on distal chromosome 15q suggests the
8
9 possibility of an HD modifier effect due to synaptic vesicle glycoprotein 2B (*SV2B*). The
10
11 *SV2B* protein localizes to synaptic vesicles, where it is believed to function in the regulation
12
13 of vesicle trafficking and exocytosis. While a role at synapses makes it an attractive candidate
14
15 for involvement in the HD damage mechanism(s) precipitated by the expanded CAG repeat,
16
17 this locus will require replication or other confirmation to ensure that it harbors a *bona fide*
18
19 genetic modifier of HD.
20
21
22
23
24

25 In summary, we expanded our approaches for identifying genetic modifiers of HD
26
27 using typed and imputed repeats. We focused on STRs in this study primarily due to their
28
29 clinical significance and available resources for genome-wide imputation.^{51,78,79,84} However,
30
31 other genetic variations (i.e., VNTR and structural variations) also generate biological
32
33 consequences in humans.^{93,94} Therefore, investigation of other types DNA polymorphisms
34
35 may yield a more complete map of genetic modifiers of HD. Nevertheless, although
36
37 our data clearly show the lack of modification of HD onset by CAG repeat size
38
39 polymorphisms in other polyglutamine disease genes, they do point, along with the complex
40
41 coding hexamer repeat in *TCERG1*⁷⁶ and a complex nonamer coding repeat in *MSH3*,²¹ to
42
43 the potential of finer delineation of other tandem repeats across the genome as a potential
44
45 source of modifiers that could further refine the HD landscape and inform the development of
46
47 treatments for HD.
48
49
50
51
52
53
54
55
56
57
58
59
60

Acknowledgements

This study would not be possible without the vital contribution of the research participants and their families. This manuscript is dedicated to the late Dr. Lesley Jones, who had made invaluable contributions to this study and countless discoveries in Huntington's disease.

Funding

This research was supported by the CHDI Foundation Inc., the U.S. National Institutes of Health (NS082079, NS091161, NS016367, NS049206, NS105709, NS119471), the Medical Research Council (UK MR/L010305/1 and fellowship MR/P001629/1), and a Cardiff University School of Medicine studentship. EMR was the recipient of a scholarship from Fundação para a Ciência e a Tecnologia (SFRH/BD/44335/2008). NAA is partly supported by an Alzheimer's Association Research Grant (Award Number: AARG-19-616534) and a European Research Council Starting Grant (Number: 101041677).

Competing interests

J.F.G. was a Scientific Advisory Board member and had a financial interest in Triplet Therapeutics, Inc. His NIH-funded project is using genetic and genomic approaches to uncover other genes that significantly influence when diagnosable symptoms emerge and how rapidly they worsen in Huntington disease. The company is developing new therapeutic approaches to address triplet repeat disorders such Huntington's disease, myotonic dystrophy and spinocerebellar ataxias. His interests were reviewed and are managed by Massachusetts General Hospital and Mass General Brigham in accordance with their conflict of interest policies. J.F.G. has also been a consultant for Wave Life Sciences USA, Inc., Biogen, Inc.

1
2
3 and Pfizer, Inc. Within the last five years D.G.M. has been a scientific consultant and/or
4 received an honoraria/stock options from AMO Pharma, Dyne, F. Hoffman-La Roche,
5 LoQus23, Novartis, Ono Pharmaceuticals, Rgenta Therapeutics, Sanofi, Sarepta Therapeutics
6 Inc, Script Biosciences, Triplet Therapeutics, and Vertex Pharmaceuticals and held research
7 contracts with AMO Pharma and Vertex Pharmaceuticals. J.D.L. is a paid Advisory Board
8 member for F. Hoffmann-La Roche Ltd and uniQure biopharma B.V., and he is a paid
9 consultant for Vaccinex Inc, Wave Life Sciences USA Inc, Genentech Inc, Triplet Inc, and
10 PTC Therapeutics Inc. T.H.M. is an associate member of the scientific advisory board of
11 LoQus23 Therapeutics. L.J. was a member of the scientific advisory boards of LoQus23
12 Therapeutics and Triplet Therapeutics. V.C.W. was a Scientific Advisory Board member of
13 Triplet Therapeutics, Inc., a company developing new therapeutic approaches to address
14 triplet repeat disorders such Huntington's disease and myotonic dystrophy. Her financial
15 interests in Triplet Therapeutics were reviewed and are managed by Massachusetts General
16 Hospital and Mass General Brigham in accordance with their conflict of interest policies. She
17 is a scientific advisory board member of LoQus23 Therapeutics and has provided paid
18 consulting services to Alnylam, Acadia Pharmaceuticals Inc., Alnylam Inc., Biogen Inc. and
19 Passage Bio. J.M.L. consults for Life Edit Therapeutics and serves in the scientific advisory
20 board of GenEdit, Inc.
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

References

1. Depienne C, Mandel JL. 30 years of repeat expansion disorders: What have we learned and what are the remaining challenges? *American journal of human genetics*. May 6 2021;108(5):764-785. doi:10.1016/j.ajhg.2021.03.011
2. Di Prospero NA, Fischbeck KH. Therapeutics development for triplet repeat expansion diseases. *Nature reviews Genetics*. Oct 2005;6(10):756-65. doi:10.1038/nrg1690
3. Ross CA. Polyglutamine pathogenesis: emergence of unifying mechanisms for Huntington's disease and related disorders. *Neuron*. Aug 29 2002;35(5):819-22. doi:10.1016/s0896-6273(02)00872-3
4. Paulson HL, Bonini NM, Roth KA. Polyglutamine disease and neuronal cell death. *Proceedings of the National Academy of Sciences of the United States of America*. Nov 21 2000;97(24):12957-8. doi:10.1073/pnas.210395797
5. Orr HT, Zoghbi HY. Trinucleotide repeat disorders. *Annual review of neuroscience*. 2007;30:575-621. doi:10.1146/annurev.neuro.29.051605.113042
6. Andrew SE, Goldberg YP, Kremer B, *et al*. The relationship between trinucleotide (CAG) repeat length and clinical features of Huntington's disease. *Nature genetics*. Aug 1993;4(4):398-403. doi:10.1038/ng0893-398
7. Duyao M, Ambrose C, Myers R, *et al*. Trinucleotide repeat length instability and age of onset in Huntington's disease. *Nature genetics*. Aug 1993;4(4):387-92. doi:10.1038/ng0893-387
8. Snell RG, MacMillan JC, Cheadle JP, *et al*. Relationship between trinucleotide repeat expansion and phenotypic variation in Huntington's disease. *Nature genetics*. Aug 1993;4(4):393-7. doi:10.1038/ng0893-393
9. Persichetti F, Srinidhi J, Kanaley L, *et al*. Huntington's disease CAG trinucleotide repeats in pathologically confirmed post-mortem brains. *Neurobiology of disease*. Dec 1994;1(3):159-66. doi:10.1006/nbdi.1994.0019
10. Gusella JF, MacDonald ME. Molecular genetics: unmasking polyglutamine triggers in neurodegenerative disease. *Nature reviews Neuroscience*. Nov 2000;1(2):109-15. doi:10.1038/35039051
11. Stevanin G, Durr A, Brice A. Clinical and molecular advances in autosomal dominant cerebellar ataxias: from genotype to phenotype and physiopathology. *European journal of human genetics : EJHG*. Jan 2000;8(1):4-18. doi:10.1038/sj.ejhg.5200403
12. Lee JM, Ramos EM, Lee JH, *et al*. CAG repeat expansion in Huntington disease determines age at onset in a fully dominant fashion. *Neurology*. Mar 6 2012;78(10):690-5. doi:10.1212/WNL.0b013e318249f683
13. Li JL, Hayden MR, Almqvist EW, *et al*. A genome scan for modifiers of age at onset in Huntington disease: The HD MAPS study. *American journal of human genetics*. Sep 2003;73(3):682-7. doi:10.1086/378133
14. Genetic Modifiers of Huntington's Disease C. Identification of Genetic Factors that Modify Clinical Onset of Huntington's Disease. *Cell*. Jul 30 2015;162(3):516-26. doi:10.1016/j.cell.2015.07.003
15. Genetic Modifiers of Huntington's Disease Consortium. Electronic address ghmhe, Genetic Modifiers of Huntington's Disease C. CAG Repeat Not Polyglutamine Length Determines Timing of Huntington's Disease Onset. *Cell*. Aug 8 2019;178(4):887-900 e14. doi:10.1016/j.cell.2019.06.036

16. Pearson CE, Nichol Edamura K, Cleary JD. Repeat instability: mechanisms of dynamic mutations. *Nature reviews Genetics*. Oct 2005;6(10):729-42. doi:10.1038/nrg1689
17. Kovtun IV, McMurray CT. Features of trinucleotide repeat instability in vivo. *Cell research*. Jan 2008;18(1):198-213. doi:10.1038/cr.2008.5
18. Dragileva E, Hendricks A, Teed A, *et al.* Intergenerational and striatal CAG repeat instability in Huntington's disease knock-in mice involve different DNA repair genes. *Neurobiology of disease*. Jan 2009;33(1):37-47. doi:10.1016/j.nbd.2008.09.014
19. Pinto RM, Dragileva E, Kirby A, *et al.* Mismatch repair genes Mlh1 and Mlh3 modify CAG instability in Huntington's disease mice: genome-wide and candidate approaches. *PLoS genetics*. Oct 2013;9(10):e1003930. doi:10.1371/journal.pgen.1003930
20. Ciosi M, Maxwell A, Cumming SA, *et al.* A genetic association study of glutamine-encoding DNA sequence structures, somatic CAG expansion, and DNA repair gene variants, with Huntington disease clinical outcomes. *EBioMedicine*. Oct 2019;48:568-580. doi:10.1016/j.ebiom.2019.09.020
21. Flower M, Lomeikaite V, Ciosi M, *et al.* MSH3 modifies somatic instability and disease severity in Huntington's and myotonic dystrophy type 1. *Brain : a journal of neurology*. Jun 19 2019;142(7):1876-86. doi:10.1093/brain/awz115
22. Goold R, Flower M, Moss DH, *et al.* FAN1 modifies Huntington's disease progression by stabilizing the expanded HTT CAG repeat. *Human molecular genetics*. Feb 15 2019;28(4):650-661. doi:10.1093/hmg/ddy375
23. Kim KH, Hong EP, Shin JW, *et al.* Genetic and Functional Analyses Point to FAN1 as the Source of Multiple Huntington Disease Modifier Effects. *American journal of human genetics*. Jul 2 2020;107(1):96-110. doi:10.1016/j.ajhg.2020.05.012
24. Goold R, Hamilton J, Menneteau T, *et al.* FAN1 controls mismatch repair complex assembly via MLH1 retention to stabilize CAG repeat expansion in Huntington's disease. *Cell reports*. Aug 31 2021;36(9):109649. doi:10.1016/j.celrep.2021.109649
25. McAllister B, Donaldson J, Binda CS, *et al.* Exome sequencing of individuals with Huntington's disease implicates FAN1 nuclease activity in slowing CAG expansion and disease onset. *Nature neuroscience*. Apr 2022;25(4):446-457. doi:10.1038/s41593-022-01033-5
26. Takano H, Onodera O, Takahashi H, *et al.* Somatic mosaicism of expanded CAG repeats in brains of patients with dentatorubral-pallidoluysian atrophy: cellular population-dependent dynamics of mitotic instability. *American journal of human genetics*. Jun 1996;58(6):1212-22.
27. Maciel P, Lopes-Cendes I, Kish S, Sequeiros J, Rouleau GA. Mosaicism of the CAG repeat in CNS tissue in relation to age at death in spinocerebellar ataxia type 1 and Machado-Joseph disease patients. *American journal of human genetics*. Apr 1997;60(4):993-6.
28. Manley K, Shirley TL, Flaherty L, Messer A. Msh2 deficiency prevents in vivo somatic instability of the CAG repeat in Huntington disease transgenic mice. *Nature genetics*. Dec 1999;23(4):471-3. doi:10.1038/70598
29. Matsuura T, Sasaki H, Yabe I, *et al.* Mosaicism of unstable CAG repeats in the brain of spinocerebellar ataxia type 2. *Journal of neurology*. Sep 1999;246(9):835-9. doi:10.1007/s004150050464
30. Kennedy L, Evans E, Chen CM, *et al.* Dramatic tissue-specific mutation length increases are an early molecular event in Huntington disease pathogenesis. *Human molecular genetics*. Dec 15 2003;12(24):3359-67. doi:10.1093/hmg/ddg352
31. Wheeler VC, Lebel LA, Vrbanac V, Teed A, te Riele H, MacDonald ME. Mismatch repair gene Msh2 modifies the timing of early disease in Hdh(Q111) striatum. *Human molecular genetics*. Feb 1 2003;12(3):273-81. doi:10.1093/hmg/ddg056

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
32. Gomes-Pereira M, Fortune MT, Ingram L, McAbney JP, Monckton DG. Pms2 is a genetic enhancer of trinucleotide CAG/CTG repeat somatic mosaicism: implications for the mechanism of triplet repeat expansion. *Human molecular genetics*. Aug 15 2004;13(16):1815-25. doi:10.1093/hmg/ddh186
33. Owen BA, Yang Z, Lai M, *et al.* (CAG)(n)-hairpin DNA binds to Msh2-Msh3 and changes properties of mismatch recognition. *Nature structural & molecular biology*. Aug 2005;12(8):663-70. doi:10.1038/nsmb965
34. Tome S, Manley K, Simard JP, *et al.* MSH3 polymorphisms and protein levels affect CAG repeat instability in Huntington's disease mice. *PLoS genetics*. 2013;9(2):e1003280. doi:10.1371/journal.pgen.1003280
35. Bettencourt C, Hensman-Moss D, Flower M, *et al.* DNA repair pathways underlie a common genetic mechanism modulating onset in polyglutamine diseases. *Annals of neurology*. Jun 2016;79(6):983-90. doi:10.1002/ana.24656
36. Morales F, Vasquez M, Santamaria C, Cuenca P, Corrales E, Monckton DG. A polymorphism in the MSH3 mismatch repair gene is associated with the levels of somatic instability of the expanded CTG repeat in the blood DNA of myotonic dystrophy type 1 patients. *DNA repair*. Apr 2016;40:57-66. doi:10.1016/j.dnarep.2016.01.001
37. Massey TH, Jones L. The central role of DNA damage and repair in CAG repeat diseases. *Disease models & mechanisms*. Jan 30 2018;11(1)doi:10.1242/dmm.031930
38. Laabs BH, Klein C, Pozojevic J, *et al.* Identifying genetic modifiers of age-associated penetrance in X-linked dystonia-parkinsonism. *Nature communications*. May 28 2021;12(1):3216. doi:10.1038/s41467-021-23491-4
39. Roy JCL, Vitalo A, Andrew MA, *et al.* Somatic CAG expansion in Huntington's disease is dependent on the MLH3 endonuclease domain, which can be excluded via splice redirection. *Nucleic acids research*. Apr 19 2021;49(7):3907-3918. doi:10.1093/nar/gkab152
40. Telenius H, Kremer B, Goldberg YP, *et al.* Somatic and gonadal mosaicism of the Huntington disease gene CAG repeat in brain and sperm. *Nature genetics*. Apr 1994;6(4):409-14. doi:10.1038/ng0494-409
41. Fortune MT, Vassilopoulos C, Coolbaugh MI, Siciliano MJ, Monckton DG. Dramatic, expansion-biased, age-dependent, tissue-specific somatic mosaicism in a transgenic mouse model of triplet repeat instability. *Human molecular genetics*. Feb 12 2000;9(3):439-45. doi:10.1093/hmg/9.3.439
42. Kennedy L, Shelbourne PF. Dramatic mutation instability in HD mouse striatum: does polyglutamine load contribute to cell-specific vulnerability in Huntington's disease? *Human molecular genetics*. Oct 12 2000;9(17):2539-44. doi:10.1093/hmg/9.17.2539
43. Watase K, Venken KJ, Sun Y, Orr HT, Zoghbi HY. Regional differences of somatic CAG repeat instability do not account for selective neuronal vulnerability in a knock-in mouse model of SCA1. *Human molecular genetics*. Nov 1 2003;12(21):2789-95. doi:10.1093/hmg/ddg300
44. Mouro Pinto R, Arning L, Giordano JV, *et al.* Patterns of CAG repeat instability in the central nervous system and periphery in Huntington's disease and in spinocerebellar ataxia type 1. *Human molecular genetics*. Aug 29 2020;29(15):2551-2567. doi:10.1093/hmg/ddaa139
45. Stuitje G, van Belzen MJ, Gardiner SL, *et al.* Age of onset in Huntington's disease is influenced by CAG repeat variations in other polyglutamine disease-associated genes. *Brain : a journal of neurology*. Jul 1 2017;140(7):e42. doi:10.1093/brain/awx122
46. Tezenas du Montcel S, Durr A, Bauer P, *et al.* Modulation of the age at onset in spinocerebellar ataxia by CAG tracts in various genes. *Brain : a journal of neurology*. Sep 2014;137(Pt 9):2444-55. doi:10.1093/brain/awu174

- 1
2
3 47. Gao R, Chakraborty A, Geater C, *et al.* Mutant huntingtin impairs PNKP and
4 ATXN3, disrupting DNA repair and transcription. *eLife*. Apr 17
5 2019;8doi:10.7554/eLife.42988
6
7 48. Sequeiros J, Seneca S, Martindale J. Consensus and controversies in best practices for
8 molecular genetic testing of spinocerebellar ataxias. *European journal of human genetics* :
9 *EJHG*. Nov 2010;18(11):1188-95. doi:10.1038/ejhg.2010.10
10
11 49. Hong EP, MacDonald ME, Wheeler VC, *et al.* Huntington's Disease Pathogenesis:
12 Two Sequential Components. *Journal of Huntington's disease*. 2021;10(1):35-51.
13 doi:10.3233/JHD-200427
14
15 50. Browning SR, Browning BL. Rapid and accurate haplotype phasing and missing-data
16 inference for whole-genome association studies by use of localized haplotype clustering.
17 *American journal of human genetics*. Nov 2007;81(5):1084-97. doi:10.1086/521987
18
19 51. Saini S, Mitra I, Mousavi N, Fotsing SF, Gymrek M. A reference haplotype panel for
20 genome-wide imputation of short tandem repeats. *Nature communications*. Oct 23
21 2018;9(1):4397. doi:10.1038/s41467-018-06694-0
22
23 52. Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic
24 acids research*. Jan 15 1999;27(2):573-80. doi:10.1093/nar/27.2.573
25
26 53. Gatchel JR, Zoghbi HY. Diseases of unstable repeat expansion: mechanisms and
27 common principles. *Nature reviews Genetics*. Oct 2005;6(10):743-55. doi:10.1038/nrg1691
28
29 54. Jardim L, Silveira I, Pereira ML, *et al.* Searching for modulating effects of SCA2,
30 SCA6 and DRPLA CAG tracts on the Machado-Joseph disease (SCA3) phenotype. *Acta
31 neurologica Scandinavica*. Mar 2003;107(3):211-4. doi:10.1034/j.1600-0404.2003.00046.x
32
33 55. Pulst SM, Santos N, Wang D, *et al.* Spinocerebellar ataxia type 2: polyQ repeat
34 variation in the CACNA1A calcium channel modifies age of onset. *Brain : a journal of
35 neurology*. Oct 2005;128(Pt 10):2297-303. doi:10.1093/brain/awh586
36
37 56. Raposo M, Ramos A, Bettencourt C, Lima M. Replicating studies of genetic
38 modifiers in spinocerebellar ataxia type 3: can homogeneous cohorts aid? *Brain : a journal of
39 neurology*. Dec 2015;138(Pt 12):e398. doi:10.1093/brain/awv206
40
41 57. Chen Z, Zheng C, Long Z, *et al.* (CAG)_n loci as genetic modifiers of age-at-onset in
42 patients with Machado-Joseph disease from mainland China. *Brain : a journal of neurology*.
43 Aug 2016;139(Pt 8):e41. doi:10.1093/brain/aww087
44
45 58. Savic D, Topisirovic I, Keckarevic M, *et al.* Is the 31 CAG repeat allele of the
46 spinocerebellar ataxia 1 (SCA1) gene locus non-specifically associated with trinucleotide
47 expansion diseases? *Psychiatric genetics*. Dec 2001;11(4):201-5. doi:10.1097/00041444-
48 200112000-00004
49
50 59. Hellenbroich Y, Kaulich M, Opitz S, Schwinger E, Zuhlke C. No association of the
51 SCA1 (CAG)₃₁ allele with Huntington's disease, myotonic dystrophy type 1 and
52 spinocerebellar ataxia type 3. *Psychiatric genetics*. Jun 2004;14(2):61-3.
53 doi:10.1097/01.ypg.0000128763.69225.77
54
55 60. Dhaenens CM, Burnouf S, Simonin C, *et al.* A genetic variation in the ADORA2A
56 gene modifies age at onset in Huntington's disease. *Neurobiology of disease*. Sep
57 2009;35(3):474-6. doi:10.1016/j.nbd.2009.06.009
58
59 61. Taherzadeh-Fard E, Saft C, Wiczorek S, Epplen JT, Arning L. Age at onset in
60 Huntington's disease: replication study on the associations of ADORA2A, HAP1 and OGG1.
Neurogenetics. Oct 2010;11(4):435-9. doi:10.1007/s10048-010-0248-3
62. Metzger S, Saukko M, Van Che H, *et al.* Age at onset in Huntington's disease is
modified by the autophagy pathway: implication of the V471A polymorphism in Atg7.
Human genetics. Oct 2010;128(4):453-9. doi:10.1007/s00439-010-0873-9

63. Alberch J, Lopez M, Badenas C, *et al.* Association between BDNF Val66Met polymorphism and age at onset in Huntington disease. *Neurology*. Sep 27 2005;65(6):964-5. doi:10.1212/01.wnl.0000175977.57661.b1
64. Valcarcel-Ocete L, Alkorta-Aranburu G, Iriando M, *et al.* Exploring Genetic Factors Involved in Huntington Disease Age of Onset: E2F2 as a New Potential Modifier Gene. *PLoS one*. 2015;10(7):e0131573. doi:10.1371/journal.pone.0131573
65. Arning L, Kraus PH, Valentin S, Saft C, Andrich J, Epplen JT. NR2A and NR2B receptor gene variations modify age at onset in Huntington disease. *Neurogenetics*. Feb 2005;6(1):25-8. doi:10.1007/s10048-004-0198-8
66. Andresen JM, Gayan J, Cherny SS, *et al.* Replication of twelve association studies for Huntington's disease residual age of onset in large Venezuelan kindreds. *Journal of medical genetics*. Jan 2007;44(1):44-50. doi:10.1136/jmg.2006.045153
67. Metzger S, Rong J, Nguyen HP, *et al.* Huntingtin-associated protein-1 is a modifier of the age-at-onset of Huntington's disease. *Human molecular genetics*. Apr 15 2008;17(8):1137-46. doi:10.1093/hmg/ddn003
68. Kloster E, Saft C, Akkad DA, Epplen JT, Arning L. Association of age at onset in Huntington disease with functional promoter variations in NPY and NPY2R. *Journal of molecular medicine*. Feb 2014;92(2):177-84. doi:10.1007/s00109-013-1092-3
69. Taherzadeh-Fard E, Saft C, Andrich J, Wiczorek S, Arning L. PGC-1alpha as modifier of onset age in Huntington disease. *Molecular neurodegeneration*. Feb 6 2009;4:10. doi:10.1186/1750-1326-4-10
70. Weydt P, Soyal SM, Gellera C, *et al.* The gene coding for PGC-1alpha modifies age at onset in Huntington's Disease. *Molecular neurodegeneration*. Jan 8 2009;4:3. doi:10.1186/1750-1326-4-3
71. Che HV, Metzger S, Portal E, Deyle C, Riess O, Nguyen HP. Localization of sequence variations in PGC-1alpha influence their modifying effect in Huntington disease. *Molecular neurodegeneration*. Jan 6 2011;6(1):1. doi:10.1186/1750-1326-6-1
72. Weydt P, Soyal SM, Landwehrmeyer GB, Patsch W, European Huntington Disease N. A single nucleotide polymorphism in the coding region of PGC-1alpha is a male-specific modifier of Huntington disease age-at-onset in a large European cohort. *BMC neurology*. Jan 2 2014;14:1. doi:10.1186/1471-2377-14-1
73. Metzger S, Bauer P, Tomiuk J, *et al.* The S18Y polymorphism in the UCHL1 gene is a genetic modifier in Huntington's disease. *Neurogenetics*. Mar 2006;7(1):27-30. doi:10.1007/s10048-005-0023-z
74. Correia K, Harold D, Kim KH, *et al.* The Genetic Modifiers of Motor OnsetAge (GeM MOA) Website: Genome-wide Association Analysis for Genetic Modifiers of Huntington's Disease. *Journal of Huntington's disease*. 2015;4(3):279-84. doi:10.3233/JHD-150169
75. Holbert S, D Nghien I, Kiechle T, *et al.* The Gln-Ala repeat transcriptional activator CA150 interacts with huntingtin: neuropathologic and genetic evidence for a role in Huntington's disease pathogenesis. *Proceedings of the National Academy of Sciences of the United States of America*. Feb 13 2001;98(4):1811-6. doi:10.1073/pnas.98.4.1811
76. Lobanov SV, McAllister B, McDade-Kumar M, *et al.* Huntington's disease age at motor onset is modified by the tandem hexamer repeat in TCERG1. *NPJ genomic medicine*. Sep 5 2022;7(1):53. doi:10.1038/s41525-022-00317-w
77. Gusella JF, MacDonald ME, Lee JM. Genetic modifiers of Huntington's disease. *Movement disorders : official journal of the Movement Disorder Society*. Sep 15 2014;29(11):1359-65. doi:10.1002/mds.26001

- 1
2
3 78. Gall-Duncan T, Sato N, Yuen RKC, Pearson CE. Advancing genomic technologies
4 and clinical awareness accelerates discovery of disease-associated tandem repeat sequences.
5 *Genome research*. Jan 2022;32(1):1-27. doi:10.1101/gr.269530.120
- 6 79. Malik I, Kelley CP, Wang ET, Todd PK. Molecular mechanisms underlying
7 nucleotide repeat expansion disorders. *Nature reviews Molecular cell biology*. Sep
8 2021;22(9):589-607. doi:10.1038/s41580-021-00382-6
- 9 80. Hannan AJ. Tandem repeat polymorphisms: modulators of disease susceptibility and
10 candidates for 'missing heritability'. *Trends in genetics : TIG*. Feb 2010;26(2):59-65.
11 doi:10.1016/j.tig.2009.11.008
- 12 81. Hannan AJ. Tandem repeats mediating genetic plasticity in health and disease. *Nature*
13 *reviews Genetics*. May 2018;19(5):286-298. doi:10.1038/nrg.2017.115
- 14 82. Mukamel RE, Handsaker RE, Sherman MA, *et al*. Protein-coding repeat
15 polymorphisms strongly shape diverse human phenotypes. *Science*. Sep 24
16 2021;373(6562):1499-1505. doi:10.1126/science.abg8289
- 17 83. Verbiest M, Maksimov M, Jin Y, Anisimova M, Gymrek M, Bilgin Sonay T.
18 Mutation and selection processes regulating short tandem repeats give rise to genetic and
19 phenotypic diversity across species. *Journal of evolutionary biology*. Feb 2023;36(2):321-
20 336. doi:10.1111/jeb.14106
- 21 84. Gymrek M. A genomic view of short tandem repeats. *Current opinion in genetics &*
22 *development*. Jun 2017;44:9-16. doi:10.1016/j.gde.2017.01.012
- 23 85. Mitra I, Huang B, Mousavi N, *et al*. Patterns of de novo tandem repeat mutations and
24 their role in autism. *Nature*. Jan 2021;589(7841):246-250. doi:10.1038/s41586-020-03078-7
- 25 86. Hannan AJ. Repeat DNA expands our understanding of autism spectrum disorder.
26 *Nature*. Jan 2021;589(7841):200-202. doi:10.1038/d41586-020-03658-7
- 27 87. Trost B, Engchuan W, Nguyen CM, *et al*. Genome-wide detection of tandem DNA
28 repeats that are expanded in autism. *Nature*. Oct 2020;586(7827):80-86. doi:10.1038/s41586-
29 020-2579-z
- 30 88. Fotsing SF, Margoliash J, Wang C, *et al*. The impact of short tandem repeat variation
31 on gene expression. *Nature genetics*. Nov 2019;51(11):1652-1659. doi:10.1038/s41588-019-
32 0521-9
- 33 89. Bakhtiari M, Park J, Ding YC, *et al*. Variable number tandem repeats mediate the
34 expression of proximal genes. *Nature communications*. Apr 6 2021;12(1):2075.
35 doi:10.1038/s41467-021-22206-z
- 36 90. Eslami Rasekh M, Hernandez Y, Drinan SD, Fuxman Bass JI, Benson G. Genome-
37 wide characterization of human minisatellite VNTRs: population-specific alleles and gene
38 expression differences. *Nucleic acids research*. May 7 2021;49(8):4308-4324.
39 doi:10.1093/nar/gkab224
- 40 91. Garg P, Martin-Trujillo A, Rodriguez OL, *et al*. Pervasive cis effects of variation in
41 copy number of large tandem repeats on local DNA methylation and gene expression.
42 *American journal of human genetics*. May 6 2021;108(5):809-824.
43 doi:10.1016/j.ajhg.2021.03.016
- 44 92. Lu TY, Smaruj PN, Fudenberg G, Mancuso N, Chaisson MJP. The motif composition
45 of variable number tandem repeats impacts gene expression. *Genome research*. Apr
46 2023;33(4):511-524. doi:10.1101/gr.276768.122
- 47 93. Hurles ME, Dermitzakis ET, Tyler-Smith C. The functional impact of structural
48 variation in humans. *Trends in genetics : TIG*. May 2008;24(5):238-45.
49 doi:10.1016/j.tig.2008.03.001
- 50 94. Marshall JN, Lopez AI, Pfaff AL, Koks S, Quinn JP, Bubb VJ. Variable number
51 tandem repeats - Their emerging role in sickness and health. *Experimental biology and*
52 *medicine*. Jun 2021;246(12):1368-1376. doi:10.1177/15353702211003511
- 53
54
55
56
57
58
59
60

Figure legends

Fig. 1. The lack of modification of HD onset by *ATXN3* CAG repeats.

To validate the COHORT data analysis of *ATXN3* results, we also analyzed the *ATXN3* repeat in the REGISTRY samples. Subsequently, we performed linear regression analysis using the combined data to determine whether the longer, shorter or sum of the two repeat alleles could explain residual age-at-onset of HD (n=1,388). Twenty-six repeat alleles were observed; the most and second most frequent repeat alleles were 23 and 14 CAGs, accounting for 53% of all repeat alleles. The Y-axis shows residual age-at-onset of HD subjects, representing age-at-onset corrected for individual *HTT* CAG repeat size (i.e., years). The X-axis represents the length of CAG repeat of *ATXN3* (i.e., [CAG]n).

Fig. 2. Genome-wide STR association analysis of residual age-at-onset of HD.

STRs were imputed based on the typed SNP data and subsequently used as the predictor variable with other covariates to explain residual age-at-onset of HD. For this analysis, we used the sum of two alleles for a given STR (i.e., additive model). Y-axis represent significance levels of association, expressed as $-\log_{10}(P\text{-value})$. A dotted red line represents Bonferroni significance ($P\text{-value}$, $8.48E-7$ based on 58,894 tests). White triangles mark the polyglutamine disease-causing repeats.

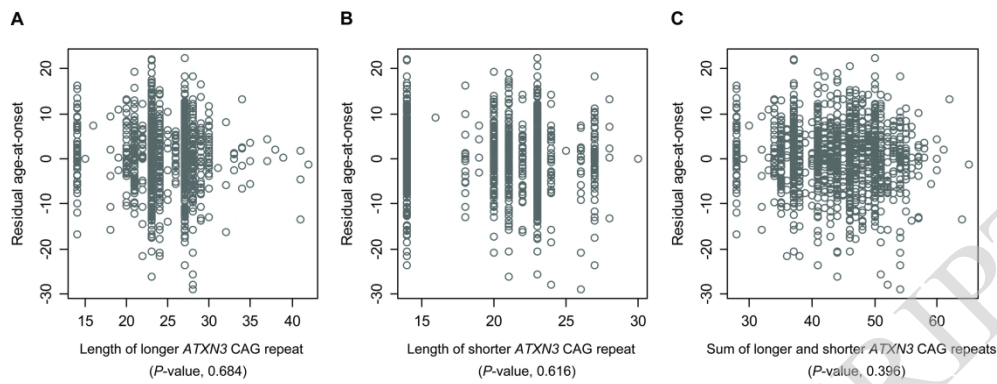
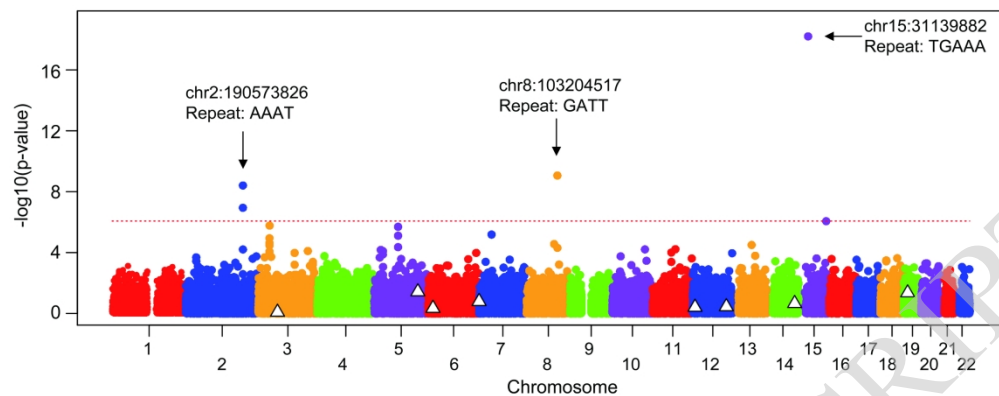


Figure 1

231x87mm (300 x 300 DPI)

ACCEPTED MANUSCRIPT



Figuer 2

228x89mm (300 x 300 DPI)

ACCEPTED MANUSCRIPT

Table 1. Statistical analysis to test association between residual age-at-onset of HD and CAG repeat size of other polyglutamine disease-associated genes.

Gene	Disease	Subjects (<i>n</i>)	Range of repeat length	<i>P</i> -value		
				Longer repeat	Shorter repeat	Sum of repeats
<i>ATNI</i>	DRPLA	551	4-19	0.4691	0.1653	0.0635
<i>ATXN1</i>	SCA1	502	21-37	0.3091	0.3281	0.8652
<i>ATXN2</i>	SCA2	604	14-33	0.7241	0.9873	0.8670
<i>ATXN3</i>	SCA3	503	14-41	0.5038	0.9361	0.5517
<i>CACNA1A</i>	SCA6	497	4-18	0.5760	0.1205	0.5201
<i>TBP</i>	SCA17	483	29-42	0.837	0.452	0.3904

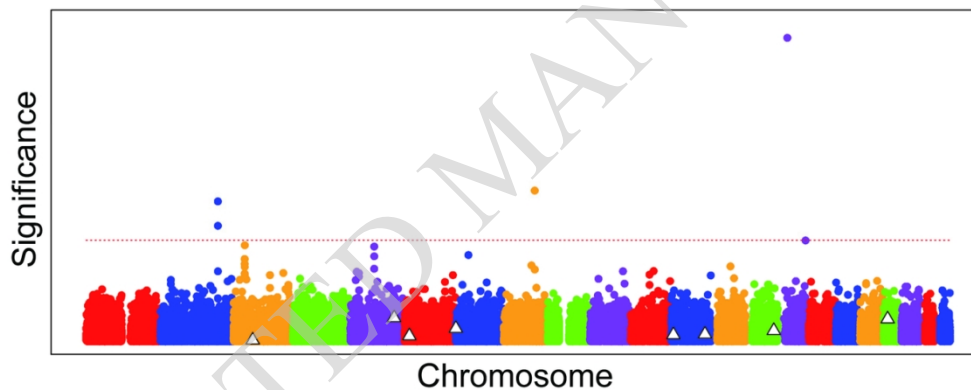
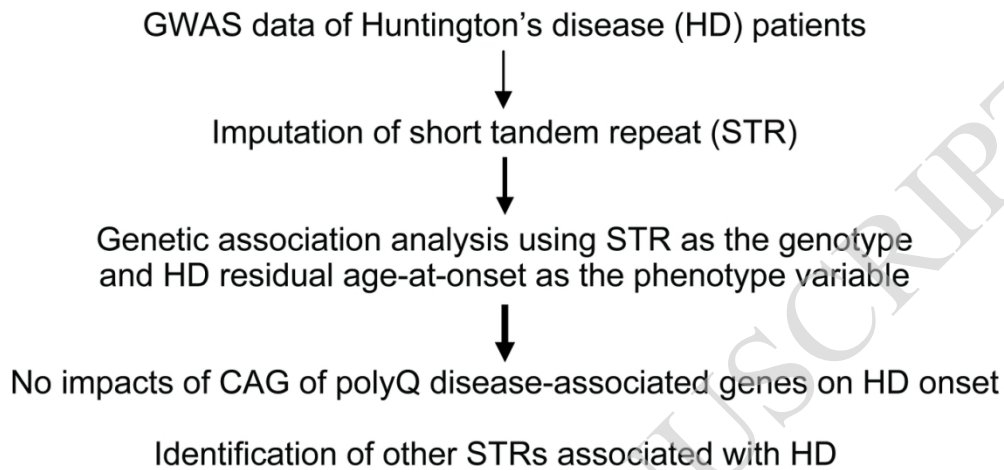
To test whether residual age-at-onset of HD is significantly associated with CAG repeat lengths of other polyglutamine disease-associated genes, we performed linear regression analyses. For each of the test genes, either the longer, shorter, or the sum of both repeat alleles was used as a continuous predictor variable to explain residual age-at-onset of HD. Sample size and *p*-value are shown. R-squared values were smaller than 1% for all tested alleles.

Table 2. Significance of CAG repeat in imputed STR association analysis.

Polyglutamine disease	Chromosomal location (GRCh37/hg19)	Gene	<i>P</i> -value
DRPLA (Dentatorubropallidolusian atrophy)	Chr12:7045880	<i>ATN1</i>	0.5785
SCA1 (Spinocerebellar ataxia type 1)	Chr6:16327865	<i>ATXN1</i>	0.4584
SCA2 (Spinocerebellar ataxia type 2)	Chr12:112036754	<i>ATXN2</i>	0.3602
SCA3 (Spinocerebellar ataxia type 3)	Chr14:92537353	<i>ATXN3</i>	0.2198
SCA6 (Spinocerebellar ataxia type 6)	Chr19:13318673	<i>CACNA1A</i>	0.0442
SCA7 (Spinocerebellar ataxia type 7)	Chr3:63898361	<i>ATXN7</i>	0.8243
SCA12 (Spinocerebellar ataxia type 12)	Chr5:146258291	<i>PPP2R2B</i>	0.0377
SCA17 (Spinocerebellar ataxia type 17)	Chr6:170870996	<i>TBP</i>	0.1638

Significance of CAG repeat of other polyglutamine disease genes extracted from the genome wide STR association analysis using an additive model. Since STRs were imputed only for autosomes, the SBMA CAG repeat on the X chromosome was not assessed.

Genome-wide association study (GWAS) to identify short tandem repeats that modify Huntington's disease



Graphical Abstract

118x118mm (300 x 300 DPI)

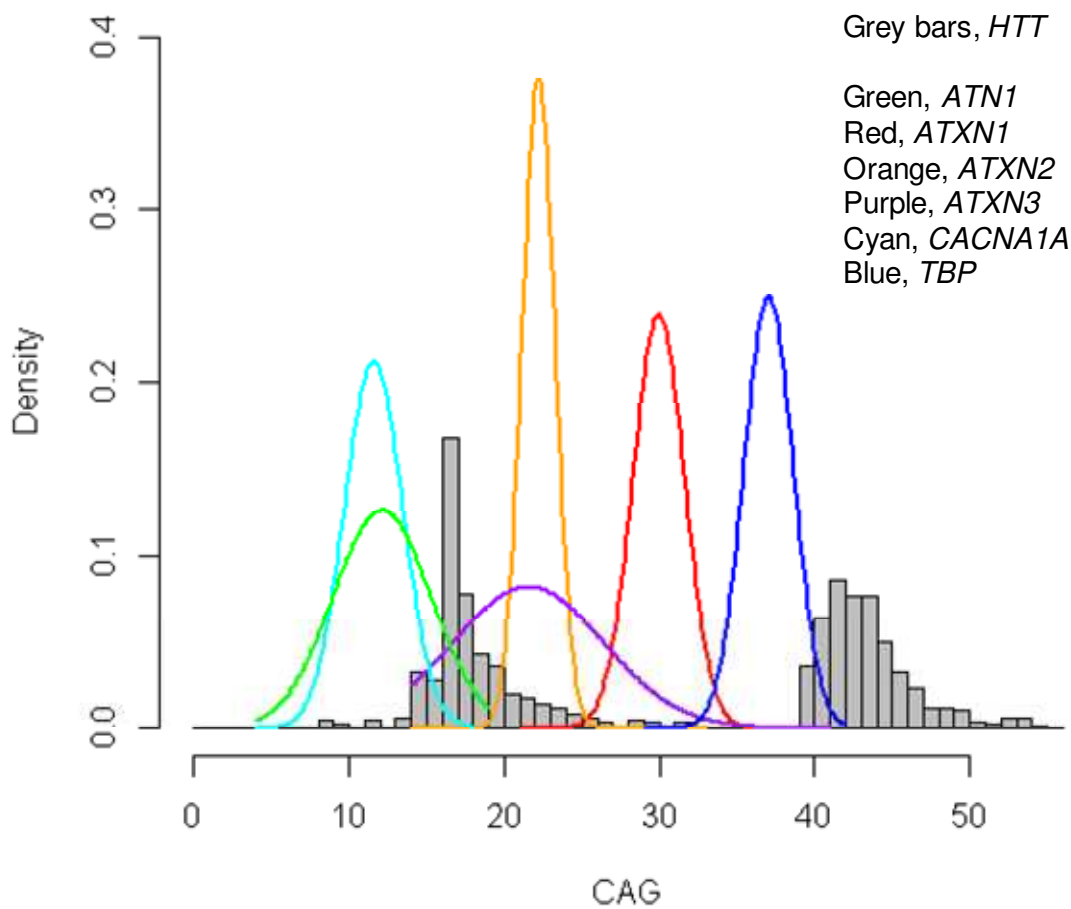
Supplementary Table 1. SNP association with residual age-at-onset in HD.

For each of other CAG repeat expansion disease genes, we evaluated the levels of association based on our recent SNP association for HD modification. A RefSeq select transcript as a representative transcript was used for a given region to identify the SNP with the highest significance in the GeM Euro 9K website. To correct the top SNP *P*-value for the gene size and number of SNPs we applied modified Sidak method to obtain corrected *P*-value.

Gene	Disease	Top SNP	MAF (%)	<i>P</i> -value	Corrected <i>P</i> -value
<i>ATN1</i>	DRPLA	rs181318837	1.539	0.036365	0.296653
<i>AR</i>	SBMA	rs5918762	15.58	0.372432	1
<i>ATXN1</i>	SCA1	rs80281835	1.677	0.003874	0.937191
<i>ATXN2</i>	SCA2	rs77838113	2.609	0.016352	0.599493
<i>ATXN3</i>	SCA3	rs55961283	2.499	0.103693	0.999995
<i>ATXN7</i>	SCA7	rs77203794	1.473	0.025268	0.952429
<i>CACNA1A</i>	SCA6	rs145803932	1.749	0.017769	0.99987
<i>PPP2R2B</i>	SCA12	rs4705448	1.114	0.000421	0.207854
<i>TBP</i>	SCA17	rs73256671	6.509	0.039786	0.598874
<i>DMPK</i>	DM1	rs183029748	1.269	0.072635	0.492709

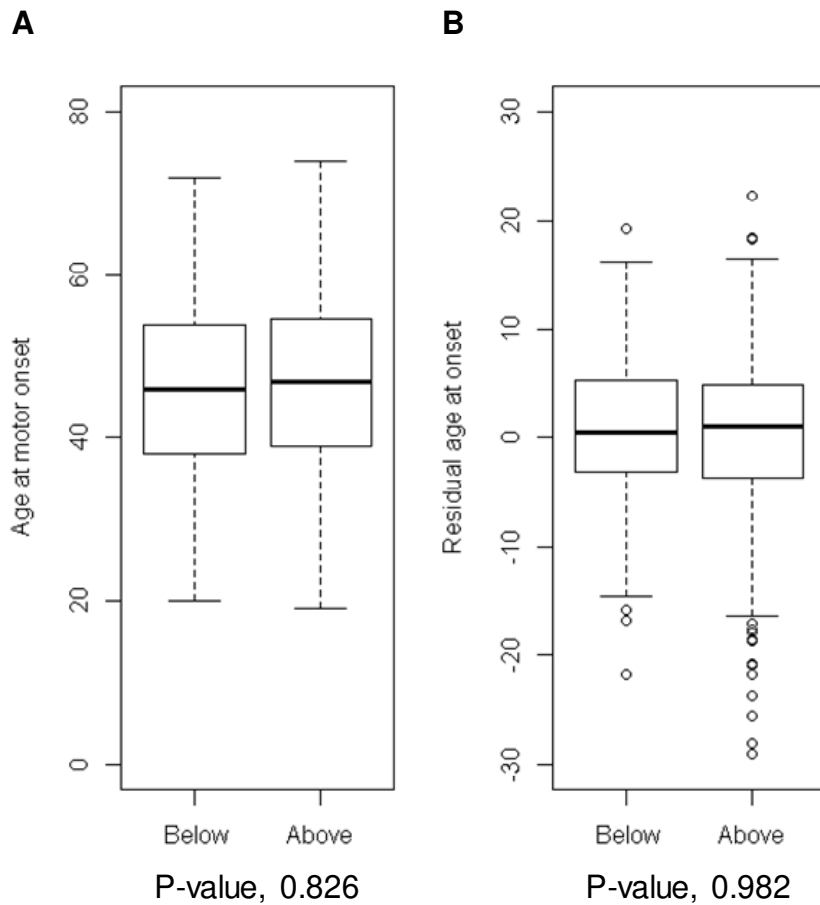
Supplementary Figure 1. Distributions of CAG repeats of polyglutamine disease-causing genes in HD subjects.

Distributions of CAG repeats at both alleles for each polyglutamine disease-causing gene are plotted. The sizes of CAG repeats of other polyglutamine disease-causing genes were determined in 483 to 604 HD individuals (Table 1). For each of other polyglutamine disease gene, mean and standard deviation was used to generate a theoretical normal distribution to plot the patterns of CAG repeat length distribution. The grey background histogram represents *HTT* CAG repeat lengths.



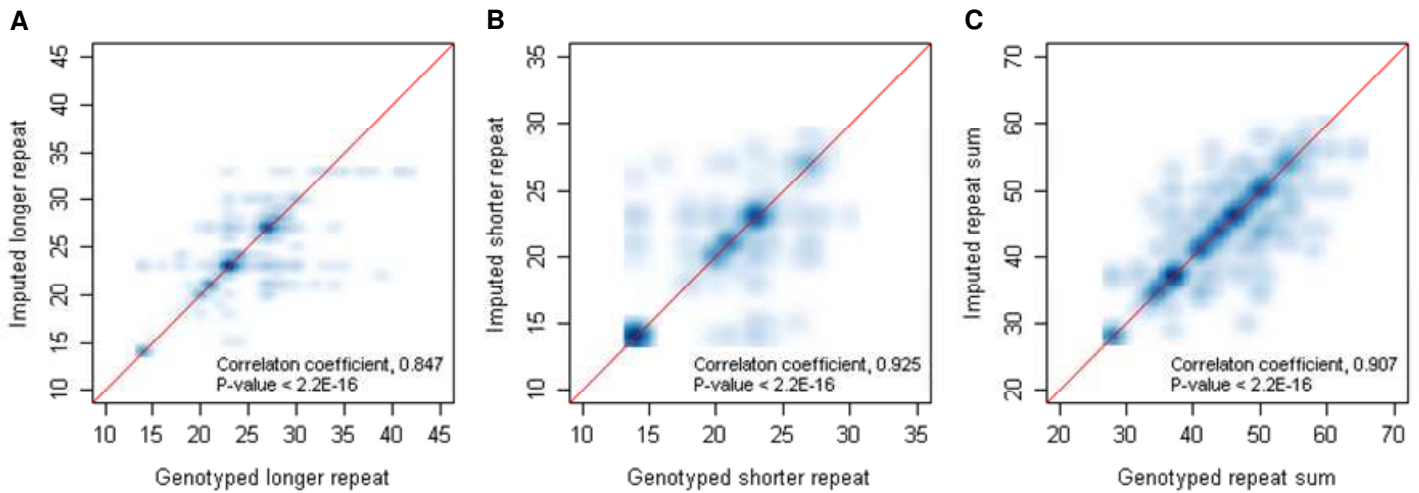
Supplementary Figure 2. Same age-at-onset and residual age-at-onset between HD subjects carrying below and above CAG repeat of *ATXN3*.

HD subjects were grouped based on the median of longer allele of the *ATXN3* CAG repeat (i.e., 23), generating a below the median group (*ATXN3* CAG < 23) and an above the median group (*ATXN3* CAG > 23). Subsequently, age-at-onset (A) and residual age-at-onset (B) were compared between two groups using Mann-Whitney *U*-test. *P*-values are shown at the bottom of each plot.



Supplementary Figure 3. Correlation between genotyped and imputed *ATXN3* CAG repeats.

To evaluate the accuracy of STR imputation, we compared *ATXN3* CAG repeats that were experimentally determined and imputed. A total of 1388 samples had both typed and imputed *ATXN3* CAG repeats. We compared longer repeat (A), shorter repeat (B), and the sum of two repeats (C) in imputed data (Y-axis) and typed data (X-axis). We also calculated percentage of identical or similar (difference smaller than 5) were calculated for longer, shorter and sum (D).



D

	Longer repeat	Shorter repeat	Sum
Identical between genotyped and imputed repeats	74.6%	87.2%	69.2%
Difference between genotyped and imputed repeats < 5	93.8%	95.4%	93.8%

Supplementary Figure 4. The lack of association between HD residual age-at-onset and imputed *ATXN3* CAG repeats

Genome-wide STR association analysis used sum of the two alleles for each STR, arguing against modification of HD by *ATXN3* repeat. We further tested whether imputed longer allele (A) or imputed shorter allele (B) is associated with residual age-at-onset. Also, we compared age-at-onset (C) and residual age-at-onset (D) between HD subjects with below and above median repeat length. Respective *P*-values are shown below of each plot.

