



Multisample motif discovery and visualization for tandem repeats

Yaran Zhang, Marc Hulsman, Alex Salazar, et al.

Genome Res. published online November 13, 2024
Access the most recent version at doi:[10.1101/gr.279278.124](https://doi.org/10.1101/gr.279278.124)

P<P	Published online November 13, 2024 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .



The NEW Vortex Mixer

USC
SCIENTIFIC
A THERMOFISHER COMPANY

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 Title:

2 **Multi-sample motif discovery and visualization for tandem repeats**

3

4 Authors:

5 Yaran Zhang¹, Marc Hulsman^{1,2}, Alex Salazar¹, Niccolò Tesi^{1,2}, Lydian Knoop¹, Sven van der
6 Lee^{1,2,4}, Sanduni Wijesekera¹, Jana Krizova¹, Erik-Jan Kamsteeg⁵, and Henne Holstege^{1,2,3,4,*}

7 1. Section Genomics of Neurodegenerative Diseases and Aging, Department of Clinical

8 Genetics, Vrije Universiteit Amsterdam, Amsterdam UMC, Amsterdam, The Netherlands

9 2. Delft Bioinformatics Lab, Delft University of Technology, Delft, The Netherlands

10 3. Alzheimer Center Amsterdam, Neurology, Vrije Universiteit Amsterdam, Amsterdam UMC

11 location VUmc, Amsterdam, The Netherlands

12 4. Amsterdam Neuroscience, Neurodegeneration, Amsterdam, The Netherlands

13 5. Department of Human Genetics, Radboud University Medical Center, Nijmegen, The
14 Netherlands

15

16 *Corresponding author:

17 Henne Holstege: h.holstege@amsterdamumc.nl

18

19 Running title: A tool to characterize tandem repeats

20

21 Keywords: tandem repeats; repeat-motifs; motif discovery; multi-sample; long-read
22 sequencing; repeat visualization tool; LRS special issue

23

24

25 **Abstract**

26 Tandem Repeats (TR) occupy a significant portion of the human genome and are the source of
27 polymorphism due to variations in sizes and motif compositions. Some of these variations have been
28 associated with various neuropathological disorders, highlighting the clinical importance of assessing
29 the motif structure of TRs. Moreover, assessing the TR motif variation can offer valuable insights into
30 evolutionary dynamics and population structure. Previously, characterizations of TRs have been
31 limited by short-read sequencing technology, which lacks the ability to accurately capture the full TR
32 sequences. As long-read sequencing becomes more accessible and can capture the full complexity of
33 TRs, there is now also a need for tools to characterize and analyze TRs using long-read data across
34 multiple samples. In this study, we present MotifScope, a novel algorithm for characterization and
35 visualization of TRs based on a *de novo* *k*-mer approach for motif discovery. Comparative analysis
36 against established tools reveals that MotifScope can identify a greater number of motifs and more
37 accurately represent the underlying repeat sequence. Moreover, MotifScope has been specifically
38 designed to enable motif composition comparisons across assemblies of different individuals, as well
39 as across long-read sequencing reads within an individual, through combined motif discovery and
40 sequence alignment. We showcase potential applications of MotifScope in diverse fields, including
41 population genetics, clinical settings, and forensic analyses.

42

43 Introduction

44 A large part of the human genome consists of repetitive elements. One such class of repeats are
45 tandem repeats (TRs), which are DNA sequences characterized by the contiguous repetition of at least
46 one nucleotide, accounting for approximately 6-8% of the human genome (Cui et al. 2024; English et
47 al. 2023; Rajan-Babu et al. 2024). TRs are broadly classified based on the size of their repetitive motif:
48 TRs with motif size ≤ 6 bp are referred to as short tandem repeats (STRs), while those with larger
49 motifs and variability in copy numbers are categorized as variable number tandem repeats (VNTRs)
50 (Tautz 1993; Eslami Rasekh et al. 2021).

51 TRs are highly polymorphic, making them a major source of diversity in human genomes (Jeffreys et
52 al. 1985). In fact, 13-17 STRs are currently used in North America and in the United Kingdom to
53 uniquely identify a person (Hammond et al. 1994; Opel et al. 2007; Glynn 2022; Mallinder et al. 2022).
54 Due to the high genetic variability and lack of linkage disequilibrium with each other, the probability
55 of two unrelated individuals sharing a perfect match of this set of STRs is less than 1 in 1 billion (Reilly
56 2001). This variability has led to the adoption of these TRs as standard tools in forensics analyses,
57 where they play a crucial role in DNA profiling and identifying individuals with a high degree of
58 accuracy (Moretti et al. 2001; Jobling and Gill 2004).

59 TRs have been associated with a range of neurological disorders (Chintalaphani et al. 2021; Hannan
60 2018; Tang et al. 2017). For instance, Friedreich ataxia (FRDA) can be caused by homozygous
61 expansion of a GAA repeat in the first intron of frataxin (*FXN*) gene, while the expansion of a GGGGCC
62 repeat intronic of *C9orf72* gene has been linked to increased risk of amyotrophic lateral sclerosis (ALS)
63 and frontotemporal dementia (FTD) (Pandolfo 2009; DeJesus-Hernandez et al. 2011). VNTR
64 expansions have also been implicated in several diseases (De Roeck et al. 2018; Hannan 2018; Song et
65 al. 2018). For example, the expansion of a 25 bp repeat in the intronic region of the ATP binding
66 cassette subfamily A member 7 (*ABCA7*) gene has been linked to an increased risk of Alzheimer's

67 disease (AD) (De Roeck et al. 2018). However, it is not just the repeat length that is associated with
68 disease. In most pathogenic TRs, the motif composition of the repeat is also important (Chen et al.
69 2020; Cortese et al. 2019; Ishiura et al. 2018; Seixas et al. 2017; Wright et al. 2020). For instance, in
70 patients with cerebellar ataxia with neuropathy and vestibular areflexia syndrome (CANVAS),
71 expansions of the repeat in replication factor C subunit 1 (*RFC1*) gene are composed primarily of
72 AAGGG or GACAG, variations from the more common AAAAG motif found in non-expanded alleles
73 (Cortese et al. 2019; Scriba et al. 2020).

74 Furthermore, TRs can provide insight into (the evolution of) population structure due to their high
75 mutability (Course et al. 2021; Ellegren 2004; Rosenberg et al. 2002; Lu et al. 2023). Course et al.
76 demonstrated significant differences in repeat length among VNTRs in genes including ADP-
77 ribosyltransferase 1 (*ART1*), PROP paired-like homeobox 1 (*PROP1*), and dynein 2 intermediate chain
78 1 (*DYNC2I1*), as well as substantial differences in motif organization in poly(rC) binding protein 3
79 (*PCBP3*) between superpopulations (Course et al. 2021). Using a pangenome-based approach, Lu et al.
80 discovered that more than 8,000 VNTRs show differential motif usage across populations (Lu et al.
81 2021). These findings emphasize the importance of considering both repeat length and motif
82 organization in TR genotyping.

83 With new sequencing technologies emerging that combine longer read length and higher accuracy, it
84 is now possible to deeply characterize TRs. As such, multiple methods have been developed to
85 genotype and profile TRs. One class of methods has been proposed which relies on databases of TRs,
86 describing for each TR the motifs that are to be considered. This has an advantage that discovered
87 motifs are more homogenous across individuals. For instance, Ren et al. developed a toolkit, *vamos*,
88 that generated a representative set of motifs for over 460,000 TRs in the human genome (Ren et al.
89 2023). This was achieved by selecting motifs for each TR in a reference set of genomes, while allowing
90 for some sequence divergence. Then, *vamos* uses this motif database to annotate TRs in the query
91 genome. Similarly, Dolzhenko et al. developed the Tandem Repeat Genotyping Tool (TRGT),

92 specifically for PacBio HiFi sequencing data (Dolzhenko et al. 2024). It uses pre-specified motifs to
93 genotype simple TR regions, while more complex repeats are defined using hidden Markov models
94 (HMMs). Additionally, TRGT comes with a visualization module, the TRVZ tool, which displays read-
95 level evidence supporting the genotype calls made by TRGT as well as the TR motifs. While the use of
96 a motif database can work well for relatively stable TRs, it may result in the loss of important motifs
97 for more variable TRs, for instance, rare motifs that are relevant to diseases. In addition, relying on a
98 database hinders the application to novel repeats, or application to species not covered by the
99 database (currently databases are only available for the human genome). Therefore, there remains a
100 need for more versatile methods that can accurately capture the complexity of TR sequences across
101 diverse genomic contexts.

102 Another class of methods detects motifs in a given sequence using a *de novo* approach. This allows
103 them to handle extensive genomic diversity, including complex TRs in which motifs can be highly
104 variable. A drawback of these methods is that in a setting in which multiple individuals are analyzed
105 together, it can be hard to canonicalize the discovered motifs across individuals. For instance, the
106 widely used Tandem Repeats Finder (TRF) program generates a separate annotation for each motif it
107 identifies, leaving it to the user to select the optimal motif representation, and to canonicalize these
108 representations across different individuals (Benson 1999). Recently, Masutani et al. proposed an
109 algorithm, uTR, to decompose TRs after selecting a better set of motifs according to maximum
110 parsimony that minimizes replication slippage events (Masutani et al. 2023). Still, this method will
111 analyze each allele sequence individually. Furthermore, these methods are not sensible to small
112 mutations (e.g., single nucleotide polymorphisms, SNPs). These small variations can however be
113 biologically important, for instance in clinical settings in which some pure repeats are considered more
114 pathogenic than interrupted repeats (Rafehi et al. 2023).

115 Here, we present MotifScope, a flexible toolkit for motif annotation and visualization of TRs from
116 sequencing data, that uses a *de novo* *k*-mer-based approach for motif discovery. To evaluate

117 MotifScope performances, we compared it to the three existing tools for motif discovery: uTR, TRF
118 and vamos. Our findings indicate that MotifScope identified a greater number of motifs and reflected
119 the actual repeat sequence more accurately than other tools. Additionally, we show potential
120 applications of MotifScope in population genetics to explore population stratification due to TRs, in
121 clinical studies to study pathogenic TRs, and in a forensic setting.

122 Results

123 MotifScope is a tool for characterizing and visualizing the motif composition of TRs. The input for
124 MotifScope is a FASTA-formatted file containing the sequence(s) to annotate (Figure 1). This can be a
125 single repetitive sequence, individual reads from one individual, or a collection of assembled alleles
126 from multiple individuals. There are three major algorithmic steps in MotifScope: (1) iteratively
127 identifying and annotating motifs using a k -mer-based approach; (2) color mapping of motifs using a
128 dimensional reduction technique, and (3) clustering and aligning sequences based on their motif
129 composition. MotifScope iteratively discovers and annotates motifs that make up long-continuous
130 sequences within a single sequence without error correction, but also across sequences to enable joint
131 discovery and annotation across multiple genomes. The output consists of a FASTA-formatted file
132 reporting each input sequence with motif information, representing motifs as sequences followed by
133 their respective counts. Additionally, MotifScope provides a tab-separated file summarizing the
134 amount of sequence covered by each motif and their corresponding count in each input sequence.
135 Finally, a graphic representation of the motif composition for each TR in each sample can be generated.

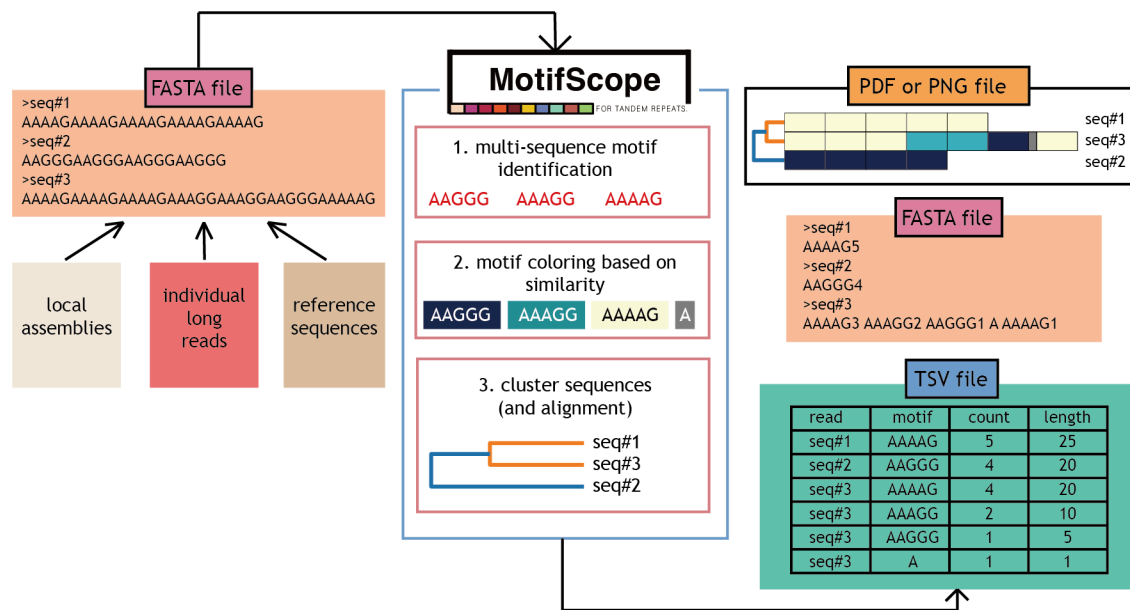


Figure 1. Overview of MotifScope. Input to MotifScope consists of a set of sequences. MotifScope first identifies motifs by evaluating the length of consecutive sequence stretches formed by each k -mer present in the sequences. It then iteratively annotates consecutive occurrences of k -mers as motifs and masks them in the sequences. Subsequently, single occurrences of identified motifs are annotated and masked. To visualize the motif composition, motifs are colored based on their sequence similarity to each other. The sequences are then clustered based on their motif composition. It also offers the option to perform multiple sequence alignment based on motif compositions.

136 Accuracy and efficiency

137 To evaluate the performance of MotifScope in characterizing TRs, we conducted a comparative
 138 analysis alongside recently developed TR-analysis methods: TRF, uTR, and vamos (with both the
 139 original motif set, “vamos original”, and an efficient motif set, “vamos efficient”) on the HG002
 140 genome sequenced with PacBio HiFi technology (see Methods). We used a set of TRs from the PacBio
 141 repeat catalog, which contains 171,146 TRs and benchmarked the tools using long-read whole-
 142 genome sequencing of the HG002 genome with PacBio HiFi technology (see Methods).

143 It’s important to note that vamos can only be applied to locations in its own repeat catalog, consisting
 144 of 467,104 VNTRs. The size distribution of TRs within the PacBio catalog differed from that in the
 145 vamos VNTR catalog. Specifically, repeats were smaller in the PacBio catalog, with 99.97% of repeats

146 being ≤ 100 bp in the GRCh38 human reference genome, while in the vamos VNTR catalog, 80.42% of
147 repeats are ≤ 100 bp (Supplementary Figure 1A). The motif sizes are also larger for the TRs in the
148 vamos catalog compared to the PacBio catalog (Supplementary Figure 1B). The PacBio catalog also
149 showed reduced sequence complexity, as illustrated by the distinct k -mer ($k = 10$) counts observed in
150 each repeat in the GRCh38 human reference genome (Supplementary Figure 1C).

151 Hence, we evaluated the performances of the tools based on a subset of 5,486 TRs that exactly
152 overlapped between these two catalogs (i.e., same start and end coordinates for each TR). This set of
153 TRs is more similar to the TRs in PacBio catalog based on length, motif size and sequence complexity.
154 Additionally, to provide a comprehensive evaluation, we also randomly sampled 5,000 repeats from
155 the vamos VNTR catalog and compared the performances of the tools on this set of repeats.
156 MotifScope identified a greater number of motifs compared to other tools (Figure 2A). To evaluate
157 the quality of motif identification, we assessed the normalized edit distance between the
158 concatenation of the motif representations generated by these tools and the true repeat sequences
159 (Figure 2B). By design, MotifScope consistently achieves an edit distance of 0, indicating an exact
160 match between its motif description and the underlying repeat sequence. In contrast, other methods
161 exhibited increasing edit distances, for example, by sorting the edit distance in ascending order, at
162 90th percentile, the edit distance was 0 for MotifScope, 0.022 for TRF, 0.032 for vamos original, 0.029
163 for vamos efficient and 0.098 for uTR. This shows that MotifScope's description of TRs reflects the
164 actual repeat sequences more accurately compared to the other three tools.

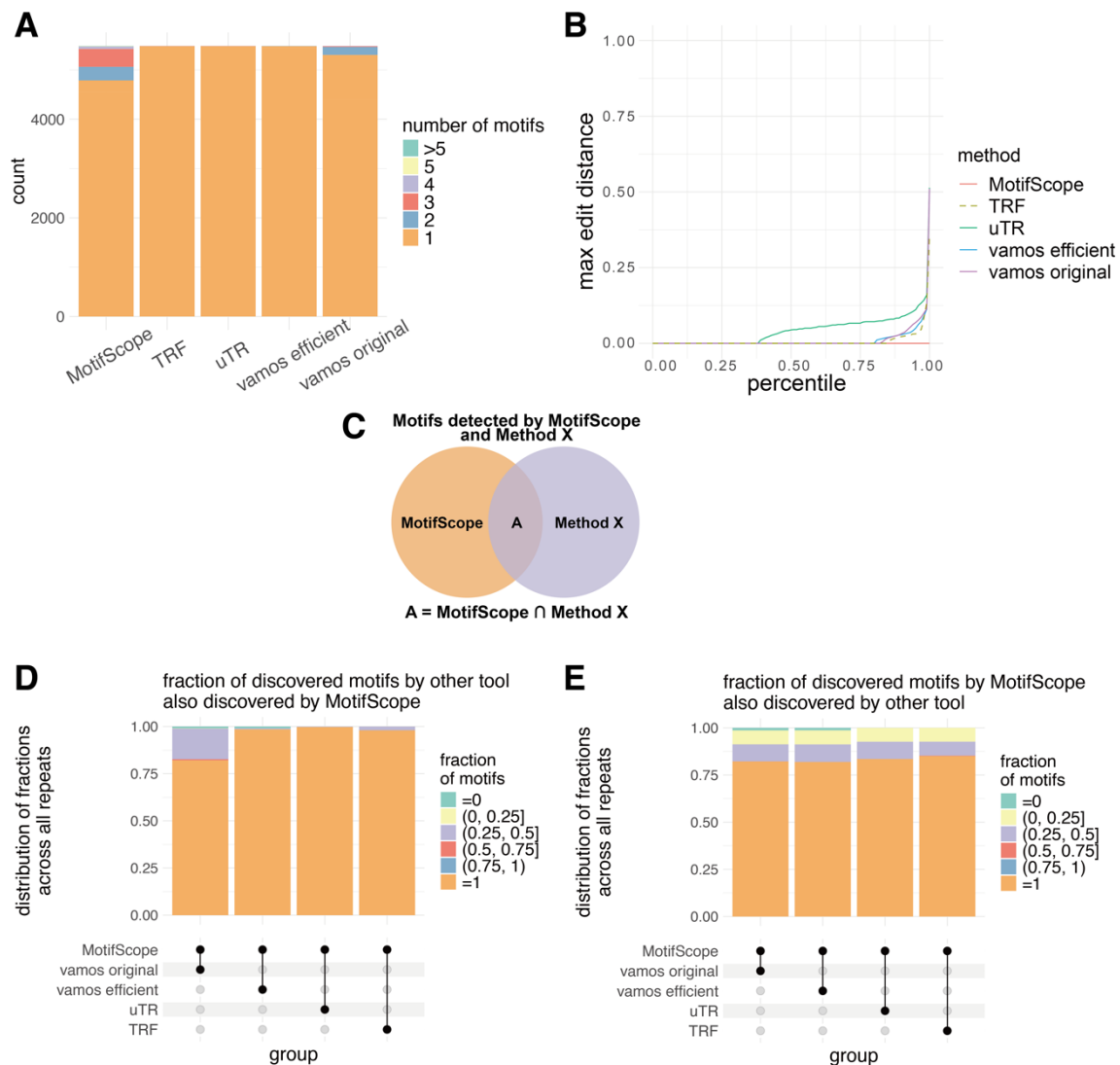


Figure 2. Comparative analysis of tandem repeat characterization. Four tools were tested in the analysis, MotifScope, TRF, uTR, and vamos. For vamos, here we used both the original motif set (vamos original) and the efficient motif set (vamos efficient). (A) shows the number of motifs discovered by each of the four tools. (B) shows the edit distance between the actual sequence and the results obtained from the four tools, normalized with respect to repeat length. (D) and (E) shows the intersection of motifs between MotifScope and other tools (dots connected by lines below the X axis): in the bar plot, each column shows the results obtained from both MotifScope and the respective other tools (Method X) for all 5,486 loci; for (D), the stacked bar plot shows the fraction of intersected motifs over the total number of motifs found by the other tool for each loci, as shown in (C), $A / \text{Method X}$, and the segments in the bar denotes the distribution of the ratio among all these loci; for (E) the stacked bar plot shows the fraction of intersected motifs over the total number of motifs found by MotifScope, as shown in (C), $A / \text{MotifScope}$, and the segments in the bar denotes the distribution of the ratio among all these loci.

165 We further evaluated the extent to which different tools captured the same motifs within the studied
166 TRs: we found that MotifScope identified all motifs found by TRF in 97.89% of loci, by uTR in 99.75%,
167 by vamos original in 79.64%, and by vamos efficient in 98.34% of the loci (Figure 2D). The relatively
168 higher fraction of motifs identified by vamos original but not by MotifScope is attributed to the
169 presence of non-repeated motifs that occur only once in vamos original motif sets, which are
170 characterized differently by MotifScope as single nucleotides. Moreover, MotifScope was also able to
171 pick up motifs that were not detected by the other tools. For example, in 15.30% of TRs, MotifScope
172 found motifs not identified by TRF, in 16.49% not identified by uTR, in 17.70% and 17.94% not
173 identified by vamos original and vamos efficient, respectively (Figure 2E). These additional motifs
174 found by MotifScope are largely composed of single-nucleotide motifs. In the 3,544 comparisons
175 between MotifScope and other methods where new motifs were detected by MotifScope, 87.4%
176 contained at least one single-nucleotide motif. However, these single-nucleotide motifs make up a
177 small percentage of the repeat sequences (1.74%). Altogether, 86.4% of sequences had no single-
178 nucleotide motifs, and the percentage increased with sequence complexity, as demonstrated by
179 distinct *k*-mer count (Supplementary Figure 2). These additional motifs allow MotifScope to more
180 accurately represent the underlying repeat sequence. When using 5,000 VNTR randomly sampled
181 from the vamos catalog, MotifScope identified more motifs and reflected the sequence more
182 accurately compared to other tools. Additionally, the overlap of motifs identified by MotifScope and
183 other tools decreased in this subset, reflecting the higher complexity of this set of TR (Supplementary
184 Figure 3).

185 We evaluated computational performances in terms of running time and memory usage using 48
186 known pathogenic TR from PacBio catalog on 1, 5, 10, 20 and all 47 genomes from the Human
187 Pangenome Reference Consortium (HPRC), respectively (see Methods). We found that running time
188 of MotifScope was mainly driven by figure generation and startup time. The motif discovery step took
189 on average 5 seconds, and only marginally increased when multiple genomes were considered
190 (Supplementary Figure 4). When coupled with the figure generation, this took 60.94 seconds for 48

191 TRs for a single genome, increasing to 538.13 seconds for 47 genomes. The maximum memory usage
192 was 0.3 GB when coupled with the figure generation on a single genome, increasing to 4.70 GB for 47
193 genomes.

194 **Merit of *de novo* motif discovery**

195 MotifScope uses a *de novo* motif discovery approach to enable the discovery of rare, possibly
196 pathogenic motifs. In Figure 3A, we present the motif characterization results of a TR within intron 2
197 of *RFC1* (Chr4:39348424-39348483), where biallelic AAGGG or GACAG repeat expansions have been
198 associated with an autosomal recessive neurological disorder, CANVAS. We employed MotifScope to
199 characterize this repeat on HG002 and a Dutch CANVAS patient (Wang et al. 2022; van de Pol et al.
200 2023). The patient is a compound heterozygous carrier of two different AAAAG expansions of 6.28
201 and 7.69 kb. In HG002, all four tools annotated the repeat using the motif AAAAG or its cyclic shifts.
202 In the CANVAS patient, MotifScope identified three major motifs: GACAG, GACAA, and AAAAG. Using
203 these motifs, the TR can be characterized as $(GACAG)_{1136}(GACAA)_{118}GAAG(AAAAG)_{14}$ with 4 indels for
204 haplotype 1 and $(GACAG)_{1420}(GACAA)_{114}GAAG(AAAAG)_{15}$ with 20 indels for haplotype 2. Note that the
205 results of MotifScope also represent the indels through additional motifs (total of 6 motifs), resulting
206 in a fully accurate representation of the actual repeat sequences (edit distance = 0.0). In contrast, uTR
207 annotated the sequences using GACAG and GACAA, $(GACAG)_{1137}(GACAA)_{133}$ for haplotype 1 and
208 $(GACAG)_{1422}(GACAA)_{130}$ for haplotype 2, with an average normalized edit distance of 0.005, mislabeling
209 the AAAAG motifs at the end of the alleles. TRF explained the sequences with a single GACAG motif,
210 $(GACAG)_{1255}$ for haplotype 1 and $(GACAG)_{1535}$ for haplotype 2, resulting in an average normalized edit
211 distance of 0.019. Finally, vamos original primarily characterized the sequence with motifs GGGAC and
212 AAAAG, $(GGGAC)_{1177}(AAAAG)_{130}$ for haplotype 1 and $(GGGAC)_{1418}(AAAAG)_{130}$ for haplotype 2 generating
213 an average normalized edit distance of 0.212, while vamos efficient predominantly used motifs
214 GGAAA, GGCAA, and AAAAG for annotation, resulting in $(GGAAA)_{1181}(GGCAA)_{115}(AAAAG)_{14}$ for

221 **Joint motif discovery and annotation across multiple** 222 **genomes**

223 MotifScope offers the capability to analyze multiple samples simultaneously, which can improve the
224 characterization and comparison of TR haplotypes. For example, when MotifScope was applied to the
225 *RFC1* repeat in a single genome HG01175 from the HPRC, it initially failed to identify the starting motif
226 AAAAG and motif GACGG immediately after the AAAGG repeat in haplotype 1. Instead, these motifs
227 were represented as single nucleotides (e.g., A, A, A, A, G for AAAAG) (highlighted in the red box in
228 Figure 4A). As the GACGG and AAAAG motifs are in the original motif set, vamos original was able to
229 identify these two single motifs on HG01175 (Supplementary Figure 5). However, when jointly
230 analyzing HG01175 with HG01109 and HG00733, MotifScope managed to also identify the single
231 copies of these motifs in HG01175, as these were recurring motifs across samples (Figure 4B, C, D). In
232 addition, joint analysis also revealed that HG01175 haplotype 1 and HG01109 haplotype 1 share the
233 same motif structure: AAAAG(AAGGG)₇(GACGG)_{1/2}(AAAGGG)_n(AAAGGGAAGG)₂AAAG(GAAA)₂AAG
234 (Figure 4D).

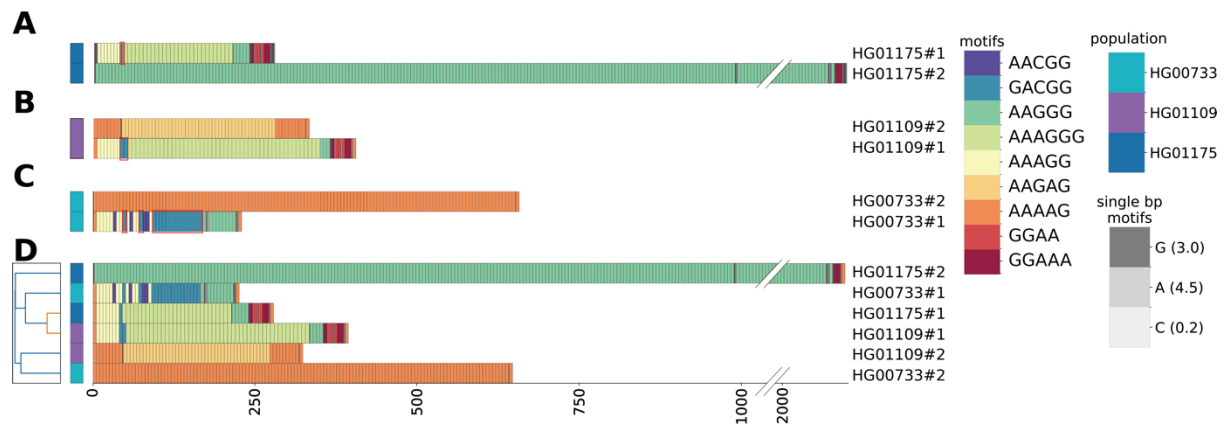


Figure 4. Motif characterization of the *RFC1* repeat in three HPRC genomes. The results from MotifScope on the assemblies of three genomes, (A) HG01175, (B) HG01109 and (C) HG00733. These three genomes are Admixed Americans. (D) shows the joint analysis result of these three genomes. The sequence “GACGG” in these sequences are highlighted in red boxes. The left panel shows the clustering of sequences along with genome identifiers, represented by the corresponding color bar on the second-to-right side. The right panel shows the motif composition of the repeat, with distinct motifs represented in different colors, as indicated by the color bar on the right side of the figure.

235 We assessed the stability of the set of discovered motifs for known pathogenic TRs by sampling
 236 subsets of 1 to 47 genomes from the HPRC samples, and comparing the similarity of discovered motifs
 237 between subsets. The motifs identified were generally stable, with increased stability observed as
 238 more genomes were included (Supplementary Figure 6A). In fact, this can be seen as a form of joint
 239 calling approach, that ensures a consistent motif representation across different sequences and
 240 samples, enhancing the reliability and comparability of the results. However, the stability varies with
 241 locus sequence complexity. For example, the notch 2 N-terminal like C (*NOTCH2NLC*) TR locus
 242 (Chr1:149390803-149390842) showed high stability across genomes with one single motif present in
 243 all genomes, while the *RFC1* TR locus showed high diversity across individuals due to presence of 8
 244 unique motifs in the HPRC genomes. As a result, for *RFC1*, a larger number of individuals had to be
 245 analyzed to accurately capture the extensive motif diversity in the population (Supplementary Figure
 246 6B).

247 **Profiling clinically relevant loci in the population**

248 TRs are known to have a wide variability in motif sequence, motif size, and repeat size across
249 populations, suggesting the importance of examining multiple samples from different populations
250 collectively. For example, a pathogenic repeat in the gene brain expressed associated with NEDD4 1
251 (*BEAN1*) (Chr16:66490397-66490466) is associated with an autosomal dominant neurological disorder,
252 Spinocerebellar Ataxia Type 31 (SCA31), and coincides with a specific GAATG repeat insertion found
253 solely in the Japanese population (Ishikawa and Nagai 2019). We applied MotifScope to this locus
254 using 47 genome assemblies from the HPRC, and identified a cluster of African alleles with distinct
255 motif structure, with A, T and ATT insertions in the TAAAA repeat (in red box in Figure 5A). Additionally,
256 two expanded alleles were observed: a 2.47 kb expansion with CAATA motifs in an admixed American
257 allele and a 0.99 kb TAAAA expansion in an East Asian allele. Notably, we observed expansions with
258 CAATA and TAAAA motifs (Figure 5B) as well as another motif TAACA that was found in a South Asian
259 individual (Figure 5A).

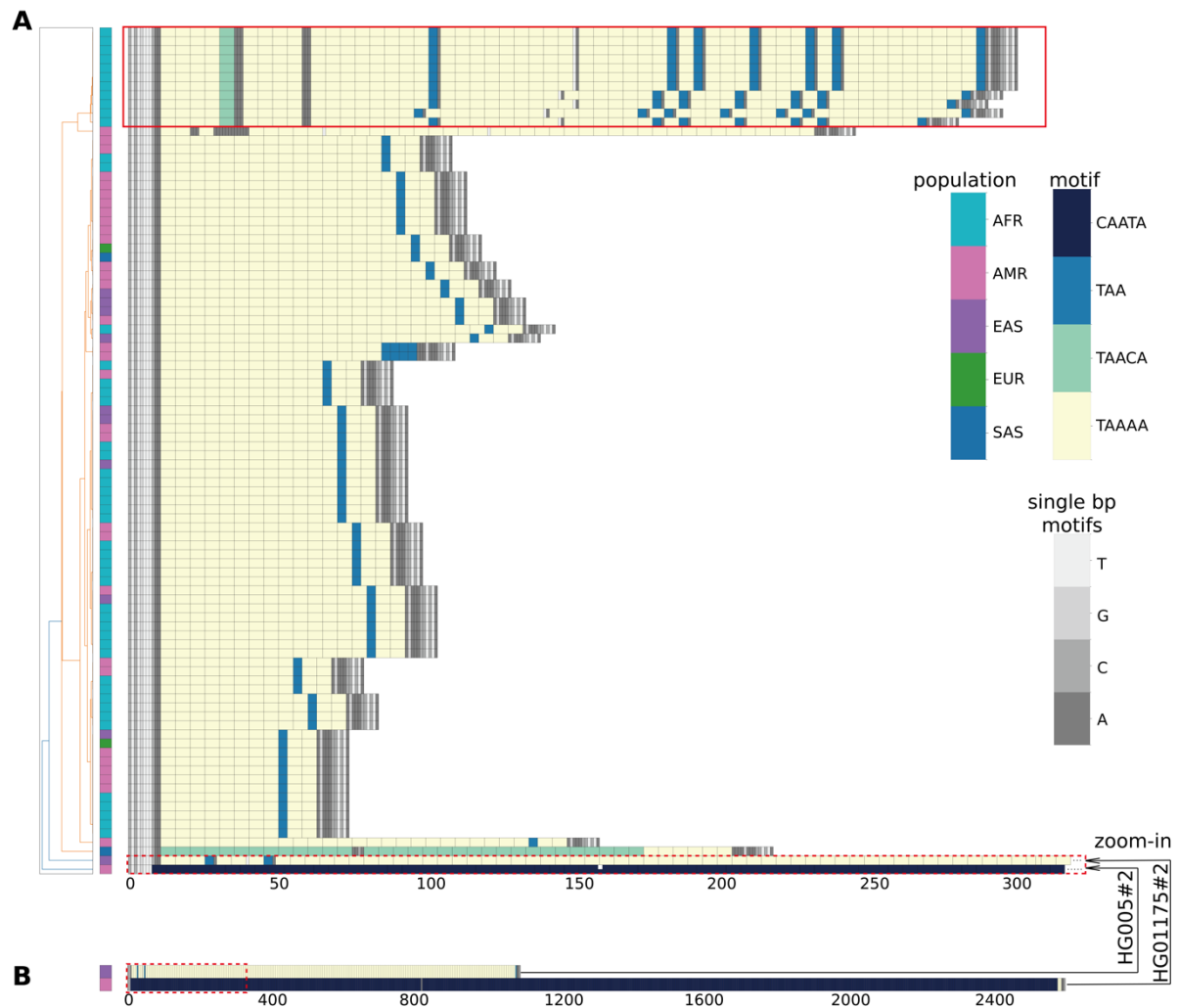


Figure 5. *BEAN1* repeat in HPRC samples. For (A), the leftmost panel presents the clustering of assembly sequences from HPRC sample ($n = 47$) with 10 bp flanking the repeat, along with population origin of the alleles in the adjacent column. The color code, denoting population origin, is in the second-to-right color bar (SAS: South Asian; EUR: European; EAS: East Asian; AMR: Admixed American; AFR: African). The right panel visually represents repeat composition, with distinct colors signifying different motifs. The two expanded alleles are truncated in this figure. The full motif compositions of these two alleles are shown in (B).

260 We also applied MotifScope to a pathogenic repeat recently identified in the gene fibroblast growth
 261 factor 14 (*FGF14*) (Chr13:102161575-102161726) (Supplementary Figure 7) (Rafehi et al. 2023;
 262 Pellerin et al. 2023). We used the HPRC assemblies as well as the otter-assembled allele sequences of
 263 246 Dutch AD patients and 238 Dutch centenarians (see Methods and Supplementary Figure 8). The
 264 expansion of this repeat has been associated with an autosomal dominant adult-onset ataxia SCA27B.
 265 Despite sharing the same GAA repeat backbone, 10 different repeat structures were identified across

266 all evaluated individuals. Similar to the GRCh38 human reference genome, the majority of the alleles
267 carried a non-expanded GAA repeat. However, the other assemblies carried different insertions inside
268 the GAA repeat: insertions of single As; an GAAGAG repeat insertion; an GAG repeat insertion
269 immediately followed by a GAGAAG repeat insertion; an [(GAA)₁(CAG)₂] repeat insertion; a GAGAAG
270 repeat insertion; insertions of different single nucleotides in slightly expanded GAA repeat; and an
271 East Asian individual with an 1.85 kb expansion composed of an [(GAG)₁(GAA)₄] repeat, an
272 [(GCA)₂(GAA)₂] repeat, an [(GCA)₁(GAA)₂] repeat, CAGAAG repeats, and CAG repeat insertions
273 (Supplementary Figure 7). Notably, MotifScope revealed that all African individuals had the same
274 starting sequence of the repeat as the non-expanded GAA alleles whereas individuals of other ancestry
275 that did not carry the pure non-expanded GAA allele had a different starting sequence (Supplementary
276 Figure 7). This case illustrates that the motifs and motif structures of TRs can be highly variable in the
277 population and MotifScope is able to identify them and cluster sequences with the same motif
278 structure together.

279 MotifScope was also applied to TRs in the ataxin 8 (*ATXN8*) gene (Chr13:70139383-70139428), where
280 interruptions in the CAG repeat with CCG repeat are believed to increase pathogenicity and are more
281 likely to cause Spinocerebellar Ataxia Type 8 (SCA8) (Koob et al. 1999). However, no such interruptions
282 were found in the Dutch centenarians, AD patients, or HPRC individuals (Supplementary Figure 9).
283 Similarly, in the ataxin 2 (*ATXN2*) TR (Chr12:111598950-111599019), where interruptions in the CAG
284 repeat with CAA motifs are associated with Spinocerebellar Ataxia Type 2 (SCA2), no reported
285 pathogenic interruption patterns were detected in these samples (Supplementary Figure 10) (Charles
286 et al. 2007).

287 We also applied MotifScope to an intronic repeat in gene *ABCA7* (Chr19:1049437-1050066) on HG002,
288 where the motif size is 25 bp (Supplementary Figure 11). MotifScope was able to identify 24 different
289 motifs in this sequence including the 4 single nucleotides. However, a substantial portion of the

290 sequence was annotated with single-nucleotide motifs: 21.4% and 25.3% for the two allele sequences,
291 respectively.

292 **Analyzing TRs at read-level**

293 Whereas MotifScope can jointly analyze multiple individuals, it can also jointly characterize and
294 visualize all mapped sequencing reads that span a TR region, which can unveil somatic and technical
295 variations. For instance, in Figure 6A, we display all spanning reads from the blood of a Dutch CANVAS
296 patient on the *RFC1* repeat, revealing a motif structure characterized by $(GACAG)_n(GACAA)_n(AAAAG)_n$.
297 However, these reads varied in size, ranging from 7.0 kb to 9.5 kb, and contained different SNVs at
298 different positions. These reads showed varying copies of different motifs suggesting the presence of
299 somatic mutation and/or technical errors (Supplementary Figure 12A).

300

301

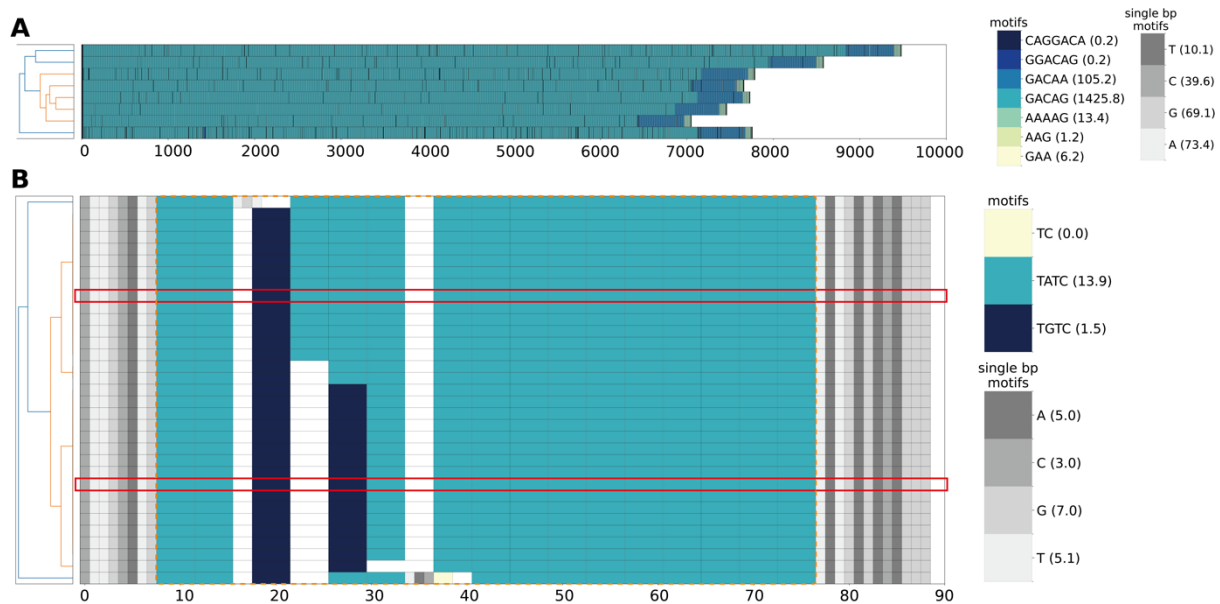


Figure 6. Motif characterization with MotifScope on long-reads. (A) shows motif characterization of the spanning reads of the *RFC1* repeat in the blood of a CANVAS patient. The sequences contain the *RFC1* repeat and 10 bp sequences flanking both sides of this region. (B) shows motif characterization of the spanning reads of forensic loci *D3S135*. The two assembly alleles of the HG002 genome are highlighted in red boxes and the repeat is highlighted in the orange box with 10 bp flanking this region. In each figure, the clustering of the sequences is shown in the left panel, the right panel shows the motif composition of the repeat, with distinct motifs represented in different colors, as indicated by the color bar on the right side of the figure. The number following the motif on the color bar indicates the average number of occurrences of the motif per read.

302 Another example is shown in Figure 6B, where we analyzed all reads as well as the phased assemblies
 303 for the forensic locus *D3S1358* (Chr3:45540738-45540802) in HG002, using the multiple sequence
 304 alignment (MSA) feature provided by MotifScope. Our analysis revealed that 14 out of 31 reads
 305 matched one assembly allele, $(TCTA)_1(TCTG)_2(TCTA)_{12}$, 12 out of 31 reads matched the other
 306 assembled allele, $(TCTA)_1(TCTG)_1(TCTA)_{14}$, while 5 out of 31 reads did not exactly match to any of the
 307 two assembled alleles. Of these 5 reads, 3 reads contained a whole motif gain or loss compared to the
 308 two assemblies, 1 read contained a single C deletion, and 1 read contained a C insertion in addition to
 309 a whole TCTA motif loss. This is also evident from the count of different motifs across different reads
 310 (Supplementary Figure 12B). This analysis allows one to assess somatic stability and/or the propensity

311 for technical sequencing errors in a region. Additionally, we randomly selected 8, 15, and 30 reads
312 from forensic locus *D2S1338* (Supplementary Figure 13). AAGG, CAGG and AGG motifs were
313 consistently identified by MotifScope with similar counts per read. Additional AGGA motifs were
314 discovered in one of the 15 and 30 reads, which contained single nucleotide deletion. This suggests
315 that the influence of sequencing coverage on motif discovery is minimal.

316 MotifScope can also be applied to reads sequenced with Oxford Nanopore Technologies. We applied
317 MotifScope to the forensic locus *D8S1179* using HG002 reads sequenced with PacBio Sequel II, PacBio
318 Revio, Nanopore R9.4.1 (simplex) chemistry, and Nanopore duplex, in reference motifs guided mode,
319 as this locus has well-defined known motifs. As shown in Supplementary Figure 14, MotifScope
320 produced similar results for reads from PacBio Sequel II, PacBio Revio, and Nanopore duplex, while
321 Nanopore R9.4.1 chemistry displayed less repeat purity. This likely reflects the lower sequencing
322 quality of Nanopore R9.4.1 chemistry compared to the other technologies. Local assemblies from
323 these technologies consistently showed the same motif compositions.

324

325 Discussion

326 Advancements in long read sequencing technologies have significantly enhanced our ability to explore
327 TRs within the human genome. These variants have emerged as crucial elements in genetic studies
328 due to their implications in evolution and diseases (Perry et al. 2008; Weischenfeldt et al. 2013). They
329 often display multiallelic patterns in the human population, and therefore, a detailed examination of
330 the composition of these repeats is important for unraveling their functional implications and
331 evolutionary history. To address the complexities of TRs, we developed MotifScope, a tool designed
332 to characterize and visually represent TR compositions. In a benchmarking study against existing tools,

333 MotifScope outperformed by identifying more motifs and more accurately reflecting sequence
334 composition.

335 By performing motif discovery *de novo*, MotifScope does not have to rely on a predetermined set of
336 motifs to characterize TR sequences, which makes it able to detect motifs that are rare or unseen in
337 the general population. This is especially important in a clinical setting, as studies have shown that
338 motif alterations in repeats can be the cause of disease. For instance, in CANVAS patients, pathogenic
339 alleles of the *RFC1* repeat have been identified as expansions of GACAG and AAGGG sequences, rather
340 than the AAAAG repeat observed in the GRCh38 human reference genome (Cortese et al. 2019; Scriba
341 et al. 2020; van de Pol et al. 2023). Similarly, in five familial adult myoclonic epilepsy (FAME) subtypes,
342 expanded TTTC segments were found inside the TTTTA repeat in sterile alpha motif domain
343 containing 12 (*SAMD12*), StAR related lipid transfer domain containing 7 (*STARD7*), membrane
344 associated ring-CH-type finger 6 (*MARCHF6*), trinucleotide repeat containing adaptor 6A (*TNRC6A*)
345 and Rap guanine nucleotide exchange factor 2 (*RAPGEF2*) genes (Bennett et al. 2020; Florian et al.
346 2019; Corbett et al. 2019; Ishiura et al. 2018). Given that both the original and efficient set of motifs
347 from vamos were constructed with the HPRC and Human Structural Variation Consortium (HGSVC)
348 samples, some of these pathogenic motifs are not represented in the default motif sets (Figure 3A)
349 (Ren et al. 2023). Consequently, there is a risk of overlooking these rare, biologically relevant motifs,
350 which could have significant implications for disease diagnosis and understanding.

351 MotifScope also offers the capability for joint motif analysis across samples. Similar to vamos, which
352 utilizes motifs present in the HPRC and HGSVC samples, MotifScope thereby uses sequence
353 information from other samples to enrich motif characterization (Figure 4). This allows it to correctly
354 annotate single occurrences of a motif in a haplotype, by leveraging information from other genomes
355 in which the motif is expanded, highlighting patterns that reflect evolutionary proximity. Furthermore,
356 joint analysis in combination with motif canonicalization ensures consistent annotation of motifs

357 across all sequences. This prevents the occurrence of different motifs and/or cyclic shifted forms of
358 the same motifs in different sequences, enhancing comparison across sequences.

359 Such comparisons will also highlight population differences. Several studies have demonstrated that
360 while overall patterns of repeat variations are highly similar across populations, there are notable
361 exceptions with population-specific patterns (Ziaei Jam et al. 2023). For instance, common CAG
362 expansions in an intronic repeat within the gene carbonic anhydrase 10 (*CA10*) were found
363 predominantly in African individuals, and the motif usage of an intronic repeat in gene *PCBP3* showed
364 substantial differences across modern superpopulations (Ziaei Jam et al. 2023; Course et al. 2021).
365 Another example is the *BEAN1* repeat, which showed population-specific pathogenic expansions
366 (Figure 5) (Ishikawa and Nagai 2019). The ataxia SCA31, found exclusively in the Japanese population,
367 aligns with a specific GAATG repeat insertion in this repeat, which is also exclusive to the Japanese
368 population, suggesting a strong founder effect. This insertion ranges 2.5–3.8 kb in size and was found
369 to inversely correlate with disease age of onset. Our analysis of the *FGF14* repeat in the HPRC samples,
370 Dutch AD patients, and Dutch cognitively healthy centenarians (Supplementary Figure 7) also revealed
371 highly population specific patterns and revealed that only certain haplotype clusters showed
372 expansions.

373 Repeat sequences might also differ at the read level. Somatic instabilities within TRs have been widely
374 observed, including in forensic STRs (Figure 6B) and pathological conditions such as cancers and repeat
375 expansion disorders like Huntington's disease, myotonic dystrophy type 1 and CANVAS (Figure 6A),
376 contributing to variability at the individual read level (dos Santos et al. 2012; Monckton 2021; Veitch
377 et al. 2007; Ciosi et al. 2019; Chintalaphani et al. 2021). Simultaneously, despite the high accuracy of
378 long read sequencing technologies, such as PacBio HiFi sequencing with an error rate lower than 0.05%
379 and Oxford Nanopore duplex with an error rate lower than 0.09%, errors can still occur during
380 sequencing, leading to variations between individual reads (Tesi et al. 2024). This can make it
381 challenging to separate somatic from technical variation. However, we find that variations in reads

382 also include whole-motif gain or loss in addition to SNVs. Given that homopolymer errors (i.e., indels)
383 are the primary source of long-read sequencing systematic errors when using the PacBio technology,
384 whole-motif gain and loss might be reflective of somatic variations (Au et al. 2012). With MotifScope,
385 we can examine these different classes of variations between reads while annotating them with a
386 consistent motif set. MotifScope also uses grayscale to color single-nucleotide motifs to visually
387 separate them from multibase motifs.

388 Alignment of sequences can further facilitate repeat sequence comparisons. However, traditional
389 alignment methods may struggle to accurately align highly repetitive sequences, in which motif gains
390 and losses occur, while only subtle differences can be used as alignment markers. Standard alignment
391 approaches therefore do not always yield the most meaningful results. Here, we addressed this issue
392 by aligning sequences based on their motif composition, performing alignment on so-called “motif
393 sequences”, which lead to a more biologically relevant alignment. This approach also enhances the
394 visual representation and comparability of repeat loci, particularly for analyzing differences between
395 reads (Figure 6B) in complex loci and identifying variations among haplotypes with similar motif
396 structures.

397 One element which sets MotifScope apart is that motifs are annotated exactly as they appear in the
398 input sequences, providing a precise representation of repeat structures. In contrast, other tools
399 usually allow for some flexibility in describing repeat sequences. This leads to simpler and more
400 condensed motif patterns, but can also hide crucial information. In particular, small variations can
401 highlight evolutionary proximity and expansion locations. Moreover, small variations can also be
402 highly relevant in clinical settings, in which certain repeats have been found to only have pathogenic
403 effects when the repeat is pure, i.e. is not interrupted by small sequence variations. One example is
404 an intronic repeat in gene *FGF14*. Expansion of this repeat (> 250 repeats) can cause adult-onset ataxia
405 SCA50/ATX-FGF14 (Rafehi et al. 2023). Previous studies have found that only individuals carrying pure
406 GAA repeat expansions developed the disease and it was hypothesized that it is because these long

407 GAA repeats are known to form secondary structures, inhibiting the transcription of the gene (Rafehi
408 et al., 2023). With MotifScope, the structure of the repeat can be easily checked and visualized. As
409 shown in Supplementary Figure 7, three individuals had > 250 GAA repeats, yet they all have SNVs
410 within the repeat sequence, and none were given ataxia diagnosis. However, it's important to note
411 that these samples have low coverage, with 3, 2, and 1 reads supporting the assemblies respectively,
412 suggesting the possibility of sequencing errors affecting both the purity and size of the repeats.

413 Due to biological variations in TRs and technical errors introduced during sequencing, TR sequences
414 often deviate from a perfect repetition of a single motif. Despite these variations, *k*-mers that form
415 longer continuous sequence stretches are still more likely to explain more of the sequences and are
416 consequently the frequently observed motifs. While this approach is effective for active repeats in
417 which motif copies have not substantially diverged, it can also present a limitation in analyzing highly
418 diverged sequences (Supplementary Figure 3). Exact motif discovery approaches can struggle with
419 sequences in which repeats are hard to identify due to the buildup of mutations. This is more likely to
420 occur in VNTRs with large motifs. For instance, Supplementary Figure 11 shows an example of an
421 intronic *ABCA7* repeat, where the sequence is highly complex. While the concatenation of results from
422 MotifScope produces the true sequence, these additional motifs do not enhance the interpretation of
423 the repeat's structure, as much of the sequence is labeled with single nucleotides. MotifScope
424 prioritizes motifs that can annotate long consecutive sequences, but when nearby motifs are all
425 different, it fails to identify them. In such cases, tools like uTR, which allow distances to the true
426 sequence, or vamos, which utilizes an established motif set to decompose repeat sequences, may be
427 more suitable alternatives. Future iterations of MotifScope might therefore benefit from
428 incorporating a more robust initial repeat structure identification step. Exact motif descriptions of
429 highly complex repeats in combination with the here proposed motif color mapping could thereby
430 reveal both the high-level structure of the repeat locus as well as more recent patterns of motif
431 divergence.

432 With the continued adoption of long-read data in research and clinical settings, characterization and
433 visualization of tandem repeats will remain an active research area, as TRs constitute a major source
434 of biological variation (Jeffreys et al. 1985). MotifScope offers a new unique angle to this, which might
435 lead to new insights into motif structure and pathogenic mechanisms. Joint motif analysis thereby
436 facilitates sequence comparisons, which is leveraged in assembly-pileups for comparative analyses of
437 motif structures across individuals, as well as read-pileups for analyzing somatic differences within an
438 individual. The joint analysis of TR motif structures might also open up new avenues for case/control
439 studies to not only take into account repeat length but also repeat structure in a biologically
440 meaningful manner. This could shed new light on the biological and phenotypical impact of tandem
441 repeats.

442 **Methods**

443 MotifScope aims to characterize and visualize motif organization of TRs. It is designed to specifically
444 target TRs that are variable among individuals. Although it works best with genome assemblies, it can
445 be applied to any collection of genomic sequences (e.g., individual long-reads) from different
446 technologies (e.g., PacBio, Nanopore).

447 It should be noted that MotifScope operates on sequences provided by the user, and therefore it is
448 important for users to provide the target sequences for analysis. MotifScope is primarily for long-read
449 sequencing data as long-read sequencing enables the characterization of the majority of TR sequences.
450 However, it's worth noting that the tool accepts FASTA-formatted files as input, allowing for the
451 analysis of various sequence data types, including reference sequences, genome assemblies, long
452 reads, and potentially short read sequencing data. It can be run in three modes: *assembly* mode (by
453 providing sample information) for assessing variations between individuals, *reads* mode for comparing
454 (somatic) variations within sequencing reads obtained from a single individual, and *single-sequence*

455 mode (by disabling sequence clustering) for analyzing the motif structure of a single sequence.
456 Optionally, motif discovery can be guided by providing an expected set of motifs (in a TSV file). Finally,
457 MotifScope can be used for performing MSA based on motif composition. MotifScope always outputs
458 the motif composition in a FASTA-formatted file (including each motif and the relative number of
459 copies). Additionally, it provides a tab-separated file reporting the fraction of sequence covered by
460 each motif. Optionally, MotifScope can generate a visual representation of motif composition across
461 sequences.

462 **Algorithm**

463 **Motif discovery and annotation**

464 MotifScope identifies repeat motifs across a given set of input sequences $S = \{s_1, s_2, \dots, s_n\}$ based on
465 highly occurring k -mers (see Supplementary Algorithm 1 for implementation details). Due to the
466 repetitive nature of TRs, a single k -mer k_x in a tandemly repeated sequence can yield a set of distinct
467 k -mers (K_x) as the repeat motif can be circularly permuted, for example for k -mer k_{AAG} , $K_{AAG} =$
468 {"AAG", "AGA", "GAA"}.

469 k -mer frequencies for varying lengths of k (the user can specify the maximum size of k -mer to screen)
470 are first computed across all sequences. To that end, input sequences are concatenated into a string
471 ($s_{combined} = \{s_1\}\{s_2\}\dots\{s_n\}$), and a suffix array and LCP array is computed for $s_{combined}$ using the
472 libsaïs library (available at <https://github.com/IlyaGrebnov/libsaïs>) (Nong et al. 2011) . We iteratively
473 walk through the suffix array, while keeping a set of active k -mers and their count. At each position i
474 in the suffix array, the count of active k -mers with a $size \leq LCP[i]$ is increased by one. Active k -mers
475 with a $size > LCP[i]$ are stored in a result list, accompanied by the count with which they occurred in
476 the sequence $s_{combined}$.

477 k -mers that can be represented as repetition of shorter sequences are excluded, e.g., AGAG is
 478 removed because it can be viewed as $(AG)_2$. Also, k -mers are not considered if they contain the
 479 sequence separation symbol “\$” or occur only once.

480 The final list of these k -mers are then sorted based on $k \times count$ so that k -mers that can mask longer
 481 sequences will be considered first. Given the set of sorted k -mers $K = \{k_a, k_b, k_c, \dots\}$, MotifScope then
 482 iteratively identifies and annotates them across the input sequences (Algorithm 1). In brief, for each
 483 iteration, a k -mer k_j is selected as the candidate motif, and the maximum continuous masked
 484 sequence length $l_{mcs}(k_j)$ is determined. This is done for all k -mers, until $k \times count$ of the next k -mer
 485 is smaller than the maximum value of l_{mcs} that has already been observed for previously considered
 486 k -mers. The k -mer with the largest l_{mcs} value is subsequently selected, and masked from the
 487 sequence. The masking operation is detailed in Supplementary Algorithm 2.

Algorithm 1 Motif annotation

Input: a set of TR sequences of one region, $S = s_1, s_2, s_3, \dots, s_n$, and parameters $kmin$ and $kmax$, defining the range for screening kmers

Output: motif annotation of S

```

1: function ANNOTATESEQUENCES( $S, kmin, kmax$ )
2:    $i \leftarrow 0$ 
3:    $M \leftarrow$  an empty set
4:    $T \leftarrow$  an empty set  $\triangleright T$  is the set of positions that are tagged by motifs
5:    $S_{cur} \leftarrow \text{join}(S, '\$')$ 
6:   while  $\max\{|s| \text{ for } s \in S\} > 1$  do
7:      $k_{best} \leftarrow \text{SELECTBESTKMER}(S, kmin, kmax)$ 
8:      $m_i \leftarrow \text{CANONICALIZEKMER}(i, M, k_{best})$ 
9:      $M \leftarrow M \cup \{m_i\}$ 
10:    if  $|k_{best}| = 1$  then
11:      break
12:    else
13:       $i \leftarrow i + 1$ 
14:    end if
15:     $i \leftarrow i + 1$ 
16:     $P \leftarrow$  the set of start positions of all uninterrupted sequences of at
    least two copies of  $m_i$  in  $S_{cur}$ 
17:     $R, S_{cur} \leftarrow \text{MASKMOTIF}(S_{cur}, m_i, P)$   $\triangleright R$  is the set of positions
    that are masked with  $m_i$ 
18:     $T \leftarrow T \cup R$ 
19:  end while
20:  for  $m \in M$  do
21:     $P \leftarrow$  the set of start positions of all single occurrences of  $m$  in  $S_{cur}$ 
22:     $R, S_{cur} \leftarrow \text{MASKMOTIF}(S_{cur}, m, P)$ 
23:     $T \leftarrow T \cup R$ 
24:  end for
25:  return  $T$ 

```

488 The unmasked sequences are then used as the input for the subsequent iteration to discover the next
 489 candidate motif. For the i^{th} iteration (where $i > 1$), an additional step is taken to provide a canonicalized

490 description of candidate motifs (Supplementary Algorithm 3). For example, if the set of already
491 selected k -mers $M = \{\text{"TGAGA"}\}$, the next candidate motif, m_2 , is canonicalized towards TGAGC instead
492 of one of its cyclical rotations GAGCT, AGCTG, GCTGA or CTGAG. This aims to ensure that these motifs
493 are in a comparable representation, enabling clearer comparison between them. To achieve this,
494 MotifScope considers all cyclical rotations of candidate k -mer m_i , and selects the rotation that
495 produces the maximum sum of pairwise alignment scores compared to all previously identified
496 candidate motifs in M . This k -mer is then considered canonicalized and is added to M .

497 This k -mer selection process stops when the longest remaining sequence ≤ 1 bp in length or the length
498 of the identified motif is 1 bp. All single occurrences of previously identified candidate motifs in the
499 remaining sequences are then tagged with the corresponding motif as well. The bases that remain
500 uncharacterized are then tagged with the single nucleotide at that position. In this way, each base in
501 the sequences is assigned to one motif.

502 **Clustering and alignment**

503 To effectively compare the patterns of motifs in sequences and to enable further downstream analysis,
504 MotifScope clusters sequences based on their motif organization and length. Hierarchical clustering is
505 subsequently performed on the pairwise distance matrix of these motif sequences.

506 By default, sequences are first translated into motif sequences, in which each unique motif is assigned
507 a character. Nucleotide sequences are then translated into a "motif sequence" with these characters
508 according to the motif assignments. Next, edit distances between pairs of sequences are calculated
509 using the Levenshtein algorithm.

510 Alternatively, full multiple sequence alignment is performed. This allows for an aligned representation
511 in the figure, and also provides clustering distances. MotifScope makes use of the partial order
512 alignment algorithm and dual affine gap penalties, as implemented in abPOA library (Gao et al. 2021).
513 MSA can both be performed at the nucleotide level, as well at the level of motif sequences. For motif

514 sequences, in which individual letters represent the motif occurrences, match and mismatch costs are
515 set according to pairwise alignment scores of the motif sequences.

516 **Dimensional reduction of motifs to a color map**

517 MotifScope accurately represents the underlying sequences, and uses a color-based visualization to
518 display TR composition. It supports reflecting motif similarity in a color spectrum such that similar
519 colors correspond to similar motifs. This is achieved by projecting the pairwise alignment score matrix
520 of motifs into a one-dimensional space using Uniform Manifold Approximation and Projection (UMAP)
521 or multidimensional scaling (MDS). It also supports using random RGB colors to represent motifs.
522 Single-nucleotide motifs are colored with grayscale to ensure they are well-separated from multi-
523 nucleotide motifs.

524 **Benchmarking**

525 We benchmarked MotifScope's ability to identify motifs and accurately represent TR sequences in the
526 context of recently developed methods: uTR, vamos, and TRF. To do so, we used HG002 genome
527 assembly, and an overlapping set of 5,486 TRs between the PacBio repeat catalog (version 0.3.0,
528 available at <https://github.com/PacificBiosciences/trgt/tree/main/repeats>) and the vamos repeat
529 catalog (based on the exact same start and end coordinates). We also randomly sampled 5,000 repeats
530 from the vamos VNTR catalog and compared the performances between these four tools on this set
531 of repeats.

532 For vamos, several motif databases have been made available by the authors: “vamos original”, which
533 uses motifs identified in samples from the HPRC and HGSCV; and “vamos efficient”, in which rare
534 motifs have been replaced with more common ones while ensuring a bounded total replacement cost
535 (compression strength $q = 0.2$) (Ren et al. 2023).

536 The number of motifs discovered was calculated for the 5,486 TRs on HG002. MotifScope, uTR, and
537 vamos, generate a single representation that can consist of multiple motifs for each sequence: in these
538 cases, all the motifs included in the representation were included. For TRF, which can produce multiple
539 representations each with a different motif, all the motifs reported were included. For the
540 comparisons, all motifs were corrected for cyclic shifts and all motifs were represented with the
541 shortest unit possible: for example, AGAG would be represented as AG.

542 We calculated the edit distance for MotifScope, uTR, and vamos, between the concatenation of the
543 motif representation of the repeat and the true sequence of the repeat. For example, if a tool
544 annotated a TR as (AGG)₃, then the relative motif-derived sequence would be AGG AGG AGG. This
545 sequence was compared to the true underlying sequence. Edit distance was then normalized by the
546 length of the true repeat sequence. Because MotifScope annotates TR sequences through exact
547 matching, the edit distance for MotifScope is always 0 by definition. In the case of TRF, where multiple
548 results were sometimes provided for a single repeat, for each characterization, the edit distances were
549 calculated using the result from TRF (i.e., motif * copy number) and the true underlying sequence of
550 the characterization. These distances were further normalized by dividing by the length of the
551 corresponding parts of the true repeat sequence. The average of these values was used to represent
552 the normalized edit distance for TRF. For the motif overlap between MotifScope and the other tools,
553 the fraction of intersected motifs between MotifScope and another tool over the total number of
554 motifs found by the other tool ($A / \text{Method X}$ in Figure 2C), and the fraction of intersected motifs
555 between MotifScope and another tool over the total number of motifs found by MotifScope ($A /$
556 MotifScope in Figure 2C) were calculated for each locus.

557 Sequencing data

558 **Public sequencing data:** Individual PacBio HiFi reads as well as publicly available whole genome
559 assemblies of the paternal and maternal haplotypes of HG002 genome were used for benchmarking

560 and read level analysis (Wang et al. 2022). PacBio-based whole genome assemblies of the publicly
561 available HPRC samples were also used to assess repeats across individuals (Wang et al. 2022; Liao et
562 al. 2023).

563 **CANVAS patients:** The PacBio HiFi sequencing data of a blood sample of a Dutch CANVAS patient was
564 additionally used to evaluate the performance of different tools in a clinical setting based on the *RFC1*
565 repeat (van de Pol et al. 2023). The genome of the Dutch CANVAS patient was assembled with hifiasm
566 (version 0.16). In addition, PacBio HiFi sequencing data of a blood sample of another Dutch CANVAS
567 patient was used to show read variability (van de Pol et al. 2023).

568 **Dutch AD patients and cognitively healthy centenarians:** PacBio HiFi sequencing data of 246 AD
569 patients and 238 Dutch cognitively healthy centenarians were additionally used for multi-genome
570 comparisons of TR motifs. The sequencing data are available through Alzheimer Genetics Hub
571 (<https://www.alzheimergenetics.org/>). Formal data requests can be submitted via the contact form at
572 <https://alzheimergenetics.org/contact/>. Additional information about the sequencing and data
573 processing can be found in Salazar et al., 2023 (Salazar et al. 2023). Targeted local assembly of the
574 regions of interest was done on these genomes using otter on HiFi reads (available at
575 <https://github.com/holstegelab/otter>) (Tesi et al. 2024).

576 Software availability

577 MotifScope has been written in Python (version \geq 3.10). The analyses were done with MotifScope
578 version 1.0.0. The code, documentation, example files, a conda environment, a packaged Docker
579 image and scripts used in this manuscript are publicly available at
580 <https://github.com/holstegelab/MotifScope> and as Supplemental Code. Additionally, MotifScope is
581 also available as a web server at <https://motifscope.holstegelab.eu>.

582 Competing interest statement

583 HH has a collaboration contract with Muna Therapeutics, PacBio, Neurimmune and Alchemab. She
584 serves in the scientific advisory boards of Muna Therapeutics and is an external advisor for Retromer
585 Therapeutics.

586 **Acknowledgements**

587 The authors are grateful to all study participants, their family members, the participating medical staff,
588 general practitioners, pharmacists and all laboratory personnel involved in patient diagnosis, blood
589 collection, blood biobanking, DNA preparation and sequencing. Part of the work in this manuscript
590 was carried out on the Cartesius supercomputer, which is embedded in the Dutch national e-
591 infrastructure with the support of SURF Cooperative. Computing hours were granted to H. H. by the
592 Dutch Research Council ('100plus': project# vuh15226, 15318, 17232, and 2020.030; 'Role of VNTRs
593 in AD'; project# 2022.31, 'Alzheimer's Genetics Hub' project# 2022.38). This work is supported by a
594 VIDI grant from the Dutch Scientific Counsel (#NWO 09150172010083) and a public-private
595 partnership with TU Delft and PacBio, receiving funding from ZonMW and Health~Holland, Topsector
596 Life Sciences & Health (PPP-allowance), and by Alzheimer Nederland WE.03-2018-07. H.H., S.L., are
597 recipients of ABOARD, a public-private partnership receiving funding from ZonMW (#73305095007)
598 and Health~Holland, Topsector Life Sciences & Health (PPP-allowance; #LSHM20106). S.L. is recipient
599 of ZonMW funding (#733050512). H.H. was supported by the Hans und Ilse Breuer Stiftung (2020),
600 Dioraphte 16020404 (2014) and the HorstingStuit Foundation (2018). Acquisition of the PacBio Sequel
601 II long read sequencing machine was supported by the ADORE Foundation (2022).

602 Research of Alzheimer center Amsterdam is part of the neurodegeneration research program of
603 Amsterdam Neuroscience. Alzheimer Center Amsterdam is supported by Stichting Alzheimer
604 Nederland and Stichting Steun Alzheimercentrum Amsterdam. The clinical database structure was
605 developed with funding from Stichting Dioraphte.

606 Author contributions: Conceived the study: HH; Wrote the manuscript: YZ, AS, MH, NT, HH; Patient
607 selection: NT, HH, E-JK, Patient blood collection and sequencing: SW, JK LK, E-JK; Data management:
608 NT, MH, SvdL, Bioinformatic analysis: YZ, MH.

609

610 References

- 611 Au KF, Underwood JG, Lee L, Wong WH. 2012. Improving PacBio Long Read Accuracy by Short Read
612 Alignment. *PLoS One* **7**: e46679.
613 <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0046679> (Accessed March
614 5, 2024).
- 615 Bennett MF, Oliver KL, Regan BM, Bellows ST, Schneider AL, Rafehi H, Sikta N, Crompton DE,
616 Coleman M, Hildebrand MS, et al. 2020. Familial adult myoclonic epilepsy type 1 SAMD12
617 TTCA repeat expansion arose 17,000 years ago and is present in Sri Lankan and Indian
618 families. *European Journal of Human Genetics* **2020 28:7 28**: 973–978.
619 <https://www.nature.com/articles/s41431-020-0606-z> (Accessed March 5, 2024).
- 620 Benson G. 1999. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res* **27**:
621 573–580. <https://academic.oup.com/nar/article/27/2/573/1061099> (Accessed March 20,
622 2023).
- 623 Charles P, Camuzat A, Benammar N, Sellal F, Destée A, Bonnet AM, Lesage S, Le Ber I, Stevanin G,
624 Dürr A, et al. 2007. Are interrupted SCA2 CAG repeat expansions responsible for parkinsonism?
625 *Neurology* **69**: 1970–1975.
- 626 Chen Z, Xu Z, Cheng Q, Tan YJ, Ong HL, Zhao Y, Lim WK, Teo JX, Foo JN, Lee HY, et al. 2020.
627 Phenotypic bases of NOTCH2NLC GGC expansion positive neuronal intranuclear inclusion
628 disease in a Southeast Asian cohort. *Clin Genet* **98**: 274–281.
629 <https://onlinelibrary.wiley.com/doi/full/10.1111/cge.13802> (Accessed March 5, 2024).
- 630 Chintalaphani SR, Pineda SS, Deveson IW, Kumar KR. 2021. An update on the neurological short
631 tandem repeat expansion disorders and the emergence of long-read sequencing diagnostics.
632 *Acta Neuropathologica Communications* **2021 9:1 9**: 1–20.
633 <https://actaneurocomms.biomedcentral.com/articles/10.1186/s40478-021-01201-x> (Accessed
634 September 27, 2022).
- 635 Ciosi M, Maxwell A, Cumming SA, Hensman Moss DJ, Alshammari AM, Flower MD, Durr A, Leavitt
636 BR, Roos RAC, Holmans P, et al. 2019. A genetic association study of glutamine-encoding DNA
637 sequence structures, somatic CAG expansion, and DNA repair gene variants, with Huntington
638 disease clinical outcomes. *EBioMedicine* **48**: 568–580.
- 639 Corbett MA, Kroes T, Veneziano L, Bennett MF, Florian R, Schneider AL, Coppola A, Licchetta L,
640 Franceschetti S, Suppa A, et al. 2019. Intronic ATTTC repeat expansions in STARD7 in familial

- 641 adult myoclonic epilepsy linked to chromosome 2. *Nature Communications* 2019 10:1 **10**: 1–10.
642 <https://www.nature.com/articles/s41467-019-12671-y> (Accessed March 5, 2024).
- 643 Cortese A, Simone R, Sullivan R, Vandrovцова J, Tariq H, Yan YW, Humphrey J, Jaunmuktane Z,
644 Sivakumar P, Polke J, et al. 2019. Biallelic expansion of an intronic repeat in RFC1 is a common
645 cause of late-onset ataxia. *Nature Genetics* 2019 51:4 **51**: 649–658. [https://www-nature-](https://www-nature-com.vu-nl.idm.oclc.org/articles/s41588-019-0372-4)
646 [com.vu-nl.idm.oclc.org/articles/s41588-019-0372-4](https://www-nature-com.vu-nl.idm.oclc.org/articles/s41588-019-0372-4) (Accessed October 18, 2022).
- 647 Course MM, Sulovari A, Gudsnuk K, Eichler EE, Valdmanis PN. 2021. Characterizing nucleotide
648 variation and expansion dynamics in human-specific variable number tandem repeats. *Genome*
649 *Res* **31**: gr.275560.121. <https://genome.cshlp.org/content/early/2021/07/09/gr.275560.121>
650 (Accessed May 17, 2022).
- 651 Cui Y, Ye W, Li JS, Li JJ, Vilain E, Sallam T, Li W. 2024. A genome-wide spectrum of tandem repeat
652 expansions in 338,963 humans. *Cell* **187**: 2336–2341.e5.
653 <http://www.cell.com/article/S0092867424002526/fulltext> (Accessed June 10, 2024).
- 654 De Roeck A, Duchateau L, Van Dongen J, Cacace R, Bjerke M, Van den Bossche T, Cras P,
655 Vandenberghe R, De Deyn PP, Engelborghs S, et al. 2018. An intronic VNTR affects splicing of
656 ABCA7 and increases risk of Alzheimer’s disease. *Acta Neuropathol* **135**: 827–837.
657 <https://link.springer.com/article/10.1007/s00401-018-1841-z> (Accessed June 4, 2024).
- 658 DeJesus-Hernandez M, Mackenzie IR, Boeve BF, Boxer AL, Baker M, Rutherford NJ, Nicholson AM,
659 Finch NCA, Flynn H, Adamson J, et al. 2011. Expanded GGGGCC Hexanucleotide Repeat in
660 Noncoding Region of C9ORF72 Causes Chromosome 9p-Linked FTD and ALS. *Neuron* **72**: 245–
661 256. <http://www.cell.com/article/S0896627311008282/fulltext> (Accessed March 20, 2023).
- 662 Dolzhenko E, English A, Dashnow H, De Sena Brandine G, Mokveld T, Rowell WJ, Karniski C,
663 Kronenberg Z, Danzi MC, Cheung WA, et al. 2024. Characterization and visualization of tandem
664 repeats at genome scale. *Nature Biotechnology* 2024 1–9.
665 <https://www.nature.com/articles/s41587-023-02057-3> (Accessed March 3, 2024).
- 666 dos Santos GC, de Souza Góoes AC, de Vitto H, Moreira CC, Avvad E, Rumjanek FD, de Moura Gallo
667 CV. 2012. Genomic instability at the 13q31 locus and somatic mtDNA mutation in the D-loop
668 site correlate with tumor aggressiveness in sporadic Brazilian breast cancer cases. *Clinics* **67**:
669 1181–1190.
- 670 Ellegren H. 2004. Microsatellites: simple sequences with complex evolution. *Nature Reviews*
671 *Genetics* 2004 5:6 **5**: 435–445. <https://www-nature-com.vu-nl.idm.oclc.org/articles/nrg1348>
672 (Accessed March 20, 2023).
- 673 English A, Dolzhenko E, Jam HZ, Mckenzie S, Olson ND, Coster W De, Park J, Gu B, Wagner J, Eberle
674 MA, et al. 2023. Benchmarking of small and large variants across tandem repeats. *bioRxiv*
675 2023.10.29.564632. <https://www.biorxiv.org/content/10.1101/2023.10.29.564632v1>
676 (Accessed June 10, 2024).
- 677 Eslami Rasekh M, Hernández Y, Drinan SD, Fuxman Bass JI, Benson G. 2021. Genome-wide
678 characterization of human minisatellite VNTRs: population-specific alleles and gene expression
679 differences. *Nucleic Acids Res* **49**: 4308–4324. [https://dx-doi-org.vu-](https://dx-doi-org.vu-nl.idm.oclc.org/10.1093/nar/gkab224)
680 [nl.idm.oclc.org/10.1093/nar/gkab224](https://dx-doi-org.vu-nl.idm.oclc.org/10.1093/nar/gkab224) (Accessed March 3, 2024).
- 681 Florian RT, Kraft F, Leitão E, Kaya S, Klebe S, Magnin E, van Rootselaar AF, Buratti J, Kühnel T,
682 Schröder C, et al. 2019. Unstable TTTTA/TTTCA expansions in MARCH6 are associated with

- 683 Familial Adult Myoclonic Epilepsy type 3. *Nature Communications* 2019 10:1 **10**: 1–14.
684 <https://www.nature.com/articles/s41467-019-12763-9> (Accessed March 5, 2024).
- 685 Gao Y, Liu Y, Ma Y, Liu B, Wang Y, Xing Y. 2021. abPOA: an SIMD-based C library for fast partial order
686 alignment using adaptive band. *Bioinformatics* **37**: 2209–2211.
687 <https://dx.doi.org/10.1093/bioinformatics/btaa963> (Accessed August 15, 2024).
- 688 Glynn CL. 2022. Bridging Disciplines to Form a New One: The Emergence of Forensic Genetic
689 Genealogy. *Genes* 2022, Vol 13, Page 1381 **13**: 1381. <https://www.mdpi.com/2073-4425/13/8/1381/htm> (Accessed March 3, 2024).
- 691 Hammond HA, Jin L, Zhong Y, Thomas Caskey C, Chakraborty R. 1994. Evaluation of 13 short tandem
692 repeat loci for use in personal identification applications. *Am J Hum Genet* **55**: 175.
693 </pmc/articles/PMC1918216/?report=abstract> (Accessed March 3, 2024).
- 694 Hannan AJ. 2018. Tandem repeats mediating genetic plasticity in health and disease. *Nature Reviews*
695 *Genetics* 2018 19:5 **19**: 286–298. <https://www.nature.com/articles/nrg.2017.115> (Accessed
696 March 20, 2023).
- 697 Ishikawa K, Nagai Y. 2019. Molecular Mechanisms and Future Therapeutics for Spinocerebellar
698 Ataxia Type 31 (SCA31). *Neurotherapeutics* **16**: 1106–1114.
699 <https://link.springer.com/article/10.1007/s13311-019-00804-6> (Accessed November 21, 2023).
- 700 Ishiura H, Doi K, Mitsui J, Yoshimura J, Matsukawa MK, Fujiyama A, Toyoshima Y, Kakita A, Takahashi
701 H, Suzuki Y, et al. 2018. Expansions of intronic TTCA and TTTA repeats in benign adult familial
702 myoclonic epilepsy. *Nature Genetics* 2018 50:4 **50**: 581–590.
703 <https://www.nature.com/articles/s41588-018-0067-2> (Accessed March 5, 2024).
- 704 Jeffreys AJ, Wilson V, Thein SL. 1985. Hypervariable ‘minisatellite’ regions in human DNA. *Nature*
705 *1985 314:6006* **314**: 67–73. <https://www-nature-com.vu-nl.idm.oclc.org/articles/314067a0>
706 (Accessed March 20, 2023).
- 707 Jobling MA, Gill P. 2004. Encoded evidence: DNA in forensic analysis. *Nat Rev Genet* **5**: 739–751.
- 708 Koob MD, Moseley ML, Schut LJ, Benzow KA, Bird TD, Day JW, Ranum LPW. 1999. An untranslated
709 CTG expansion causes a novel form of spinocerebellar ataxia (SCA8). *Nature Genetics* 1999 21:4
710 **21**: 379–384. https://www.nature.com/articles/ng0499_379 (Accessed August 15, 2024).
- 711 Liao WW, Asri M, Ebler J, Doerr D, Haukness M, Hickey G, Lu S, Lucas JK, Monlong J, Abel HJ, et al.
712 2023. A draft human pangenome reference. *Nature* 2023 617:7960 **617**: 312–324.
713 <https://www.nature.com/articles/s41586-023-05896-x> (Accessed September 21, 2024).
- 714 Lu TY, Munson KM, Lewis AP, Zhu Q, Tallon LJ, Devine SE, Lee C, Eichler EE, Chaisson MJP. 2021.
715 Profiling variable-number tandem repeat variation across populations using repeat-pangenome
716 graphs. *Nature Communications* 2021 12:1 **12**: 1–12. <https://www.nature.com/articles/s41467-021-24378-0> (Accessed August 23, 2022).
- 718 Lu TY, Smaruj PN, Fudenberg G, Mancuso N, Chaisson MJP. 2023. The motif composition of variable
719 number tandem repeats impacts gene expression. *Genome Res* **33**: 511–524.
720 <https://genome.cshlp.org/content/33/4/511.full> (Accessed September 22, 2024).
- 721 Mallinder B, Pope S, Thomson J, Beck LA, McDonald A, Ramsbottom D, Court DS, Vanhinsbergh D,
722 Barber M, Evett I, et al. 2022. Interpretation and reporting of mixed DNA profiles by seven
723 forensic laboratories in the UK and Ireland. *Forensic Sci Int Genet* **58**: 102674.

- 724 Masutani B, Kawahara R, Morishita S. 2023. Decomposing mosaic tandem repeats accurately from
725 long reads ed. T. Marschall. *Bioinformatics* **39**.
726 <https://academic.oup.com/bioinformatics/article/39/4/btad185/7114028> (Accessed May 16,
727 2023).
- 728 Monckton DG. 2021. The Contribution of Somatic Expansion of the CAG Repeat to Symptomatic
729 Development in Huntington's Disease: A Historical Perspective. *J Huntingtons Dis* **10**: 7–33.
730 <https://orcid.org/0000-0002-8298-8264> (Accessed March 5, 2024).
- 731 Moretti TR, Baumstark AL, Defenbaugh DA, Keys KM, Smerick JB, Budowle B. 2001. Validation of
732 Short Tandem Repeats (STRs) for Forensic Usage: Performance Testing of Fluorescent Multiplex
733 STR Systems and Analysis of Authentic and Simulated Forensic Samples. *J Forensic Sci* **46**: 647–
734 660. <https://dx.doi.org/10.1520/JFS15018J> (Accessed March 3, 2024).
- 735 Nong G, Zhang S, Chan WH. 2011. Two efficient algorithms for linear time suffix array construction.
736 *IEEE Transactions on Computers* **60**: 1471–1484.
- 737 Opel KL, Chung DT, Drábek J, Butler JM, McCord BR. 2007. Developmental Validation of Reduced-Size
738 STR Miniplex Primer Sets*. *J Forensic Sci* **52**: 1263–1271.
739 <https://onlinelibrary.wiley.com/doi/full/10.1111/j.1556-4029.2007.00584.x> (Accessed March 3,
740 2024).
- 741 Pandolfo M. 2009. Friedreich ataxia: the clinical picture. *J Neurol* **256 Suppl 1**: 3–8.
742 <https://pubmed.ncbi.nlm.nih.gov/19283344/> (Accessed March 3, 2024).
- 743 Pellerin D, Danzi MC, Wilke C, Renaud M, Fazal S, Dicaire M-J, Scriba CK, Ashton C, Yanick C, Beijer D,
744 et al. 2023. Deep Intronic FGF14 GAA Repeat Expansion in Late-Onset Cerebellar Ataxia . *New*
745 *England Journal of Medicine* **388**: 128–141. [https://www-nejm-org.vu-](https://www-nejm-org.vu-nl.idm.oclc.org/doi/full/10.1056/NEJMoa2207406)
746 [nl.idm.oclc.org/doi/full/10.1056/NEJMoa2207406](https://www-nejm-org.vu-nl.idm.oclc.org/doi/full/10.1056/NEJMoa2207406) (Accessed March 5, 2024).
- 747 Perry GH, Yang F, Marques-Bonet T, Murphy C, Fitzgerald T, Lee AS, Hyland C, Stone AC, Hurles ME,
748 Tyler-Smith C, et al. 2008. Copy number variation and evolution in humans and chimpanzees.
749 *Genome Res* **18**: 1698–1710. [https://genome-cshlp-org.vu-](https://genome-cshlp-org.vu-nl.idm.oclc.org/content/18/11/1698.full)
750 [nl.idm.oclc.org/content/18/11/1698.full](https://genome-cshlp-org.vu-nl.idm.oclc.org/content/18/11/1698.full) (Accessed March 5, 2024).
- 751 Rafehi H, Read J, Szmulewicz DJ, Davies KC, Snell P, Fearnley LG, Scott L, Thomsen M, Gillies G, Pope
752 K, et al. 2023. An intronic GAA repeat expansion in FGF14 causes the autosomal-dominant
753 adult-onset ataxia SCA27B/ATX-FGF14. *The American Journal of Human Genetics* **110**: 105–119.
- 754 Rajan-Babu IS, Dolzhenko E, Eberle MA, Friedman JM. 2024. Sequence composition changes in short
755 tandem repeats: heterogeneity, detection, mechanisms and clinical implications. *Nature*
756 *Reviews Genetics* 2024 1–24. <https://www.nature.com/articles/s41576-024-00696-z> (Accessed
757 June 10, 2024).
- 758 Reilly P. 2001. Legal and public policy issues in DNA forensics. *Nature Reviews Genetics* 2001 2:4 **2**:
759 313–317. <https://www.nature.com/articles/35066091> (Accessed March 5, 2024).
- 760 Ren J, Gu B, Chaisson MJP. 2023. vamos: variable-number tandem repeats annotation using efficient
761 motif sets. *Genome Biology* 2023 24:1 **24**: 1–18.
762 <https://genomebiology.biomedcentral.com/articles/10.1186/s13059-023-03010-y> (Accessed
763 September 15, 2023).

- 764 Rosenberg NA, Pritchard JK, Weber JL, Cann HM, Kidd KK, Zhivotovsky LA, Feldman MW. 2002.
765 Genetic structure of human populations. *Science (1979)* **298**: 2381–2385. <https://www-science->
766 [org.vu-nl.idm.oclc.org/doi/10.1126/science.1078311](https://www-science-org.vu-nl.idm.oclc.org/doi/10.1126/science.1078311) (Accessed March 20, 2023).
- 767 Salazar A, Tesi N, Knoop L, Pijnenburg Y, Lee S van der, Wijesekera S, Krizova J, Hiltunen M, Damme
768 M, Petrucelli L, et al. 2023. An AluYb8 retrotransposon characterises a risk haplotype of
769 TMEM106B associated in neurodegeneration. *medRxiv* 2023.07.16.23292721.
770 <https://www.medrxiv.org/content/10.1101/2023.07.16.23292721v3> (Accessed June 4, 2024).
- 771 Scriba CK, Beecroft SJ, Clayton JS, Cortese A, Sullivan R, Yau WY, Dominik N, Rodrigues M, Walker E,
772 Dyer Z, et al. 2020. A novel RFC1 repeat motif (ACAGG) in two Asia-Pacific CANVAS families.
773 *Brain* **143**: 2904–2910. <https://academic.oup.com/brain/article/143/10/2904/5939924>
774 (Accessed December 13, 2022).
- 775 Seixas AI, Loureiro JR, Costa C, Ordóñez-Ugalde A, Marcelino H, Oliveira CL, Loureiro JL, Dhingra A,
776 Brandão E, Cruz VT, et al. 2017. A Pentanucleotide ATTC Repeat Insertion in the Non-coding
777 Region of DAB1, Mapping to SCA37, Causes Spinocerebellar Ataxia. *Am J Hum Genet* **101**: 87–
778 103. <http://www.cell.com/article/S0002929717302422/fulltext> (Accessed March 20, 2023).
- 779 Song JHT, Lowe CB, Kingsley DM. 2018. Characterization of a Human-Specific Tandem Repeat
780 Associated with Bipolar Disorder and Schizophrenia. *Am J Hum Genet* **103**: 421–430.
781 <http://www.cell.com/article/S0002929718302386/fulltext> (Accessed March 20, 2023).
- 782 Tang H, Kirkness EF, Lippert C, Biggs WH, Fabani M, Guzman E, Ramakrishnan S, Lavrenko V,
783 Kakaradov B, Hou C, et al. 2017. Profiling of Short-Tandem-Repeat Disease Alleles in 12,632
784 Human Whole Genomes. *Am J Hum Genet* **101**: 700–715.
785 <http://www.cell.com/article/S0002929717303828/fulltext> (Accessed March 25, 2022).
- 786 Tautz D. 1993. Notes on the definition and nomenclature of tandemly repetitive DNA sequences. *EXS*
787 **67**: 21–28. https://link.springer.com/chapter/10.1007/978-3-0348-8583-6_2 (Accessed March
788 20, 2023).
- 789 Tesi N', Salazar A, Zhang Y, Van Der Lee S, Hulsman M, Knoop L, Wijesekera S, Krizova J, Schneider A-
790 F, Pennings M, et al. 2024. Characterising tandem repeat complexities across long-read
791 sequencing platforms with TREAT. *bioRxiv* 2024.03.15.585288.
792 <https://www.biorxiv.org/content/10.1101/2024.03.15.585288v1> (Accessed June 4, 2024).
- 793 van de Pol M, O'Gorman L, Corominas-Galbany J, Cliteur M, Derks R, Verbeek NE, van de
794 Warrenburg B, Kamsteeg EJ. 2023. Detection of the ACAGG Repeat Motif in RFC1 in Two Dutch
795 Ataxia Families. *Movement Disorders* **38**: 1555–1556. <https://onlinelibrary-wiley-com.vu->
796 [nl.idm.oclc.org/doi/full/10.1002/mds.29441](https://onlinelibrary-wiley-com.vu-nl.idm.oclc.org/doi/full/10.1002/mds.29441) (Accessed March 3, 2024).
- 797 Veitch NJ, Ennis M, McAbney JP, Shelbourne PF, Monckton DG. 2007. Inherited CAG·CTG allele
798 length is a major modifier of somatic mutation length variability in Huntington disease. *DNA*
799 *Repair (Amst)* **6**: 789–796.
- 800 Wang T, Antonacci-Fulton L, Howe K, Lawson HA, Lucas JK, Phillippy AM, Popejoy AB, Asri M, Carson
801 C, Chaisson MJP, et al. 2022. The Human Pangenome Project: a global resource to map
802 genomic diversity. *Nature* 2022 **604**: 7906 **604**: 437–446. <https://www-nature-com.vu->
803 [nl.idm.oclc.org/articles/s41586-022-04601-8](https://www-nature-com.vu-nl.idm.oclc.org/articles/s41586-022-04601-8) (Accessed March 3, 2024).

804 Weischenfeldt J, Symmons O, Spitz F, Korbel JO. 2013. Phenotypic impact of genomic structural
805 variation: insights from and for human disease. *Nature Reviews Genetics* 2013 14:2 14: 125–
806 138. <https://www-nature-com.vu-nl.idm.oclc.org/articles/nrg3373> (Accessed March 5, 2024).

807 Wright GEB, Black HF, Collins JA, Gall-Duncan T, Caron NS, Pearson CE, Hayden MR. 2020.
808 Interrupting sequence variants and age of onset in Huntington’s disease: clinical implications
809 and emerging therapies. *Lancet Neurol* 19: 930–939.
810 <http://www.thelancet.com/article/S1474442220303434/fulltext> (Accessed March 20, 2023).

811 Ziaei Jam H, Li Y, DeVito R, Mousavi N, Ma N, Lujumba I, Adam Y, Maksimov M, Huang B, Dolzhenko
812 E, et al. 2023. A deep population reference panel of tandem repeat variation. *Nature*
813 *Communications* 2023 14:1 14: 1–15. [https://www-nature-com.vu-](https://www-nature-com.vu-nl.idm.oclc.org/articles/s41467-023-42278-3)
814 [nl.idm.oclc.org/articles/s41467-023-42278-3](https://www-nature-com.vu-nl.idm.oclc.org/articles/s41467-023-42278-3) (Accessed March 5, 2024).

815