

ExpansionHunter Denovo: A computational method for locating known and novel repeat expansions in short-read sequencing data

Egor Dolzhenko^{1#} (edolzhenko@illumina.com)

Mark F. Bennett^{2,3,4#} (mark.bennett@wehi.edu.au)

Phillip A. Richmond^{5#} (prichmond@cmmt.ubc.ca)

Brett Trost^{6,7} (brett.trost@sickkids.ca)

Sai Chen¹ (schen6@illumina.com)

Joke J.F.A. van Vugt⁸ (J.F.A.vanVugt-2@umcutrecht.nl)

Charlotte Nguyen^{6,7,9} (charlotte.nguyen@sickkids.ca)

Giuseppe Narzisi¹⁰ (gnarzisi@nygenome.org)

Vladimir G. Gainullin¹ (vgainullin@illumina.com)

Andrew Gross¹ (agross@illumina.com)

Bryan Lajoie¹ (blajoie@illumina.com)

Ryan J. Taft¹ (rtaft@illumina.com)

Wyeth W. Wasserman⁵ (wyeth@cmmt.ubc.ca)

Stephen W. Scherer^{6,7,9,11} (stephen.scherer@sickkids.ca)

Jan H. Veldink⁸ (J.H.Veldink@umcutrecht.nl)

David R. Bentley¹² (DBentley@illumina.com)

R K.C. Yuen^{6,7,9#} (ryan.yuen@sickkids.ca)

Melanie Bahlo^{2,3#} (bahlo@wehi.edu.au)

Michael A. Eberle^{1#} (meberle@illumina.com)

¹Illumina Inc., 5200 Illumina Way, San Diego, CA 92122, USA

²Population Health and Immunity Division, The Walter and Eliza Hall Institute of Medical Research, 1G Royal Parade, Parkville VIC 3052, Australia

³Department of Medical Biology, The University of Melbourne, 1G Royal Parade, Parkville VIC 3052, Australia

⁴Epilepsy Research Centre, Department of Medicine, The University of Melbourne, Austin Health, 245 Burgundy Street, Heidelberg VIC 3084, Australia

⁵Centre for Molecular Medicine and Therapeutics, BC Children's Hospital, University of British Columbia, Vancouver, BC V5Z 4H4, Canada

⁶Genetics and Genome Biology, The Hospital for Sick Children, University of Toronto, 686 Bay Street, Toronto, ON M5G 0A4, Canada

⁷The Center for Applied Genomics, The Hospital for Sick Children, University of Toronto, 686 Bay Street, Toronto, ON M5G 0A4, Canada

⁸Department of Neurology, Brain Center Rudolf Magnus, University Medical Center Utrecht, Utrecht, The Netherlands

⁹Department of Molecular Genetics, University of Toronto, 1 King's College Circle, Toronto, ON M5S 2E5, Canada

¹⁰New York Genome Center, 101 Avenue of the Americas, New York 10013, USA

¹¹The McLaughlin Centre, University of Toronto, 686 Bay Street, Toronto, ON M5G 0A4, Canada

¹²Illumina Cambridge Ltd, Illumina Centre, 19 Granta Park, Great Abington, Cambridge CB21 6DF, UK

equal contribution

Running head: Discovery of repeat expansions with ExpansionHunter Denovo

Abstract

Expansions of short tandem repeats are responsible for over 40 monogenic disorders, and undoubtedly many more pathogenic repeat expansions (REs) remain to be discovered. Existing methods for detecting REs in short read sequencing data require predefined repeat catalogs. However recent discoveries have emphasized the need for detection methods that do not require candidate repeats to be specified in advance. To address this need, we introduce ExpansionHunter Denovo, an efficient catalog-free method for genome-wide detection of REs. Analysis of real and simulated data shows that our method can identify large expansions of 41 out of 43 pathogenic repeats, including nine recently reported non-reference REs not discoverable via existing methods.

ExpansionHunter Denovo is freely available at

<https://github.com/Illumina/ExpansionHunterDenovo>

Keywords

Repeat expansions, short tandem repeats, whole genome sequencing data, genome-wide analysis, Friedreich's ataxia, Myotonic Dystrophy type 1, Huntington's disease, fragile X syndrome

Background

High-throughput whole-genome sequencing (WGS) has experienced rapid reductions in per-genome costs over the past ten years¹ driving population-level sequencing projects and precision medicine initiatives at an unprecedented scale²⁻⁷. The availability of large sequencing datasets now allows researchers to perform comprehensive genome-wide searches for disease-associated variants. The primary limitations of these studies are the completeness of the reference genome and the ability to identify putative causal variations against the reference background. A wide variety of software tools can identify variations relative to the reference genome such as single nucleotide variants (SNVs) and short (1-50 bp) insertions and deletions (indels)⁸⁻¹³, copy number variants (CNVs)^{14,15} and structural variants (SVs)¹⁵⁻¹⁷. A common feature of these variant callers is their reliance on sequence reads that at least partially align to the reference genome. However, because some variants include large amounts of inserted sequence relative to the reference, methods that can analyze reads that do not align to the reference are also needed.

A particularly important category of variants that involve long insertions are repeat expansions (REs) such as the expansions in *C9orf72* repeat associated with amyotrophic lateral sclerosis (ALS). This repeat consists of three copies of CCGGGG motif in the reference (18 bp total) whereas the pathogenic mutations are comprised of at least 30 copies of the motif (180 bp total) and may reach into thousands of bases^{18,19}. REs are known to be responsible for dozens of monogenic disorders^{20,21}.

Several recently-developed tools can detect REs longer than the standard short read sequencing read length of 150 bp²²⁻²⁷. These tools have all been demonstrated to be capable of accurately detecting pathogenic expansions of simple short tandem repeats (STRs). However, recent discoveries have shown that many pathogenic repeats have complex structure and hence require more flexible methods. For instance: (a) REs causing spinocerebellar ataxia types 31 and 37, familial adult myoclonic epilepsy types 1, 2, 3, and 4, and Baratela-Scott syndrome²⁸⁻³⁴ occur within an inserted sequence relative to the reference; (b) expanded repeats recently shown to cause spinocerebellar ataxia, familial adult myoclonic epilepsy, and cerebellar ataxia with neuropathy and bilateral vestibular areflexia syndrome have different composition relative to the reference STR²⁸⁻³³; (c) Unverricht-Lundborg disease, a type of progressive myoclonus epilepsy, is caused by an expansion of a dodecamer (12-mer) repeat³⁵. None of the existing methods are capable of discovering all of these REs.

We have developed ExpansionHunter Denovo (EHdn), a novel method for performing genome-wide search for expanded repeats, to address the limitations of the existing approaches. Our method (a) does not require prior knowledge of the genomic coordinates of the REs, (b) can detect nucleotide composition changes within the expanded repeats, and (c) is applicable to both short and long motifs. EHdn scans the existing alignments of short reads from one or many sequencing libraries, including the unaligned and misaligned reads, to identify approximate locations of long repeats and their nucleotide composition. EHdn is computationally efficient because it does not

re-align reads. Depending on the sensitivity settings, EHdn can analyze a single 30-40x WGS sample in about 30 minutes to 2 hours on a single CPU thread.

In this study, we demonstrate that EHdn can be used to rediscover the REs associated with fragile X syndrome (FXS), Friedreich Ataxia (FRDA), Myotonic Dystrophy type 1 (DM1), and Huntington's disease (HD) using case-control analysis to compare a small number of affected individuals (N=14-35) to control samples (N=150) and also by using outlier analysis to compare individual cases to controls. We then characterize large (longer than the read length) repeats in our control cohort to get a sense of baseline variability of these long repeats. Finally, we demonstrate the capabilities of our method by analyzing simulated expansions of various classes of tandem repeats known to play an important role in human disease. Taken together, our findings demonstrate that EHdn is a robust tool for identifying previously-unknown pathogenic repeat expansions in both cohort and single-sample outlier analysis, capable of identifying a new, previously inaccessible class of REs.

Results

ExpansionHunter Denovo

Overview

The length of disease-causing REs tends to exceed the read length of modern short-read sequencing technologies³⁶. Thus pathogenic expansions of many repeats can be detected by locating reads that are completely contained inside the repeats. As in our previous work^{25,27}, we call these reads in-repeat reads (IRRs). We implemented a method, ExpansionHunter Denovo (EHdn), for performing a genome-wide search for IRRs in BAM/CRAM files containing read alignments. EHdn computes genome-wide STR profiles containing locations and counts of all identified IRRs. Subsequent

comparisons of STR profiles across multiple samples can reveal the locations of the expanded repeats.

Genome-wide STR profiles

Genome-wide STR profiles computed by EHDn contain information about two types of IRRs: anchored IRRs and paired IRRs. Anchored IRRs are IRRs whose mates are confidently aligned to the genomic sequence adjacent to the repeat (Methods). Paired IRRs are read pairs where both mates are IRRs with the same repeat motif. Repeats exceeding the read length generate anchored IRRs (Figure 1, middle panel). Repeats that are longer than the fragment length of the DNA library produce paired IRRs in addition to anchored IRRs (Figure 1, right panel). The genomic coordinates where the anchored reads align correspond to the approximate locations of loci harboring REs and the number of IRRs is indicative of the overall RE length.

The information about anchored IRRs is summarized in an STR profile for each repeat motif (e.g. CGG) by listing regions containing anchored IRRs in close proximity to each other together with the total number of anchored IRRs identified there (Figure 2, middle). Note that the mapping positions of anchored IRRs correspond to the positions of anchor reads; mapping positions of IRRs themselves are not used because their alignments are often unreliable.

Contrary to anchored IRRs, the origin of paired IRRs cannot be determined if a genome contains multiple long repeats with the same motif. Due to this, STR profiles only contain the overall count of paired IRRs for each repeat motif.

Comparing STR profiles across multiple samples

To compare STR profiles across multiple samples, the profiles must first be merged together. During this process, the nearby regions, across multiple samples, with anchored IRRs are merged and the associated counts are depth-normalized and tabulated for each sample (Figure 2, right; Methods). The total counts of paired IRRs

are also normalized and tabulated for each sample. The resulting per-sample counts can be compared in two ways: If the samples can be partitioned into cases and controls where a significant subset of cases is expected to contain expansions of the same repeat then a case/control analysis can be performed using a Wilcoxon rank-sum test (Methods). Alternatively, if no enrichment for any specific expansion is expected, an outlier analysis (Methods) can be used to flag repeats that are expanded in a small subgroup of cases compared to the rest of the dataset. Case-control and outlier analyses can be performed on either anchored IRRs or paired IRRs, which we call locus and motif methods, respectively (Methods). Thus, the locus method can reveal locations of repeat expansions while the motif method can reveal the overall enrichment for long repeats with a given motif.

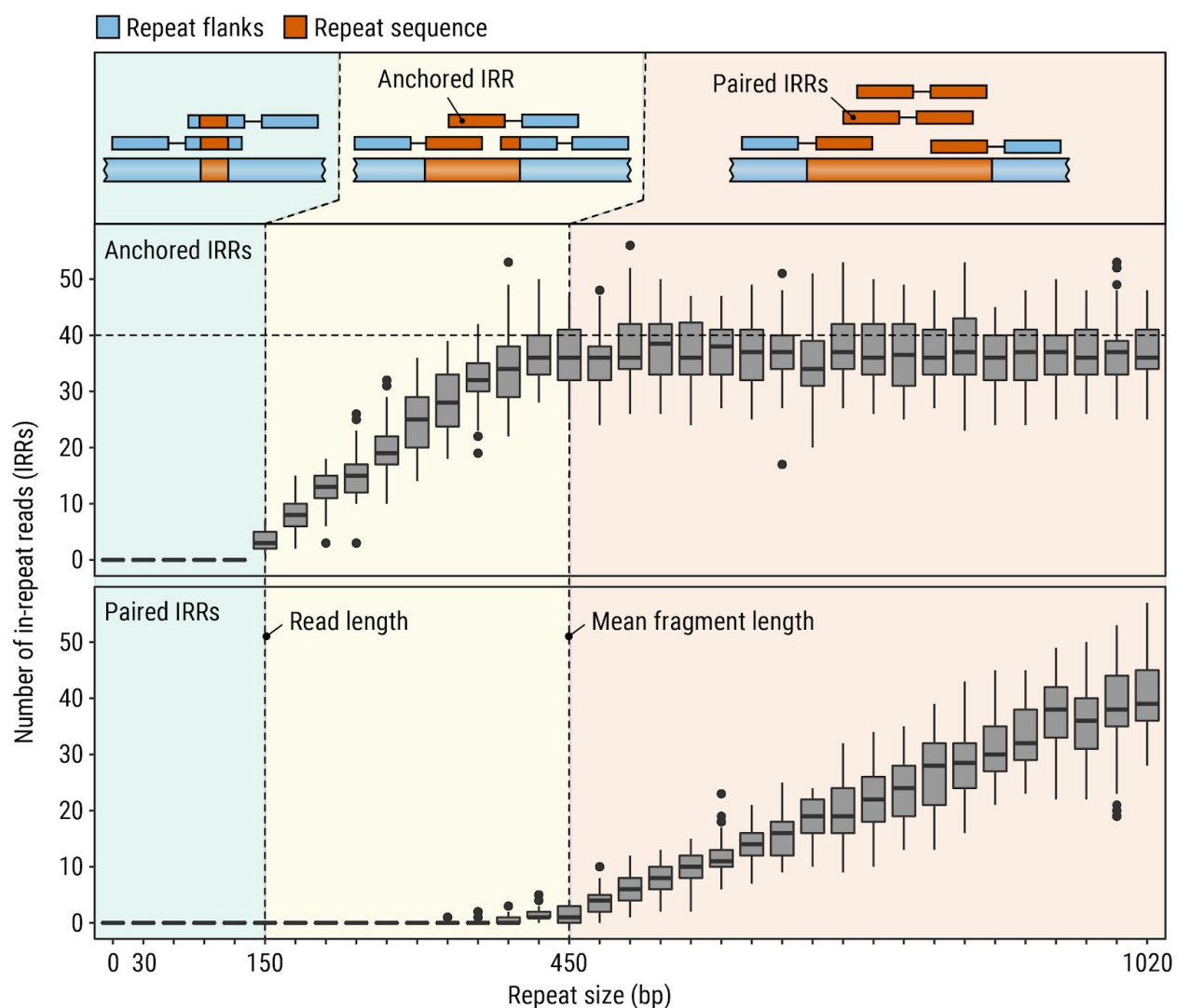


Figure 1. Diagram illustrating the types and counts of reads generated by simulating repeats of different lengths. When the repeat is shorter than the read length (left panels), there are no IRRs associated with the repeat. When a repeat is longer than the read length but shorter than the fragment length (middle panels), anchored IRRs but no paired IRRs are present. As the repeat length approaches and exceeds the fragment length (right panels), paired IRRs are generated in addition to anchored IRRs.

Baseline simulations

To demonstrate the baseline expectation of how the numbers of anchored and paired IRRs vary with repeat length, we simulated 2x150 bp reads at 20x coverage with 450 bp mean fragment length for the repeat associated with Huntington's disease and varied the repeat length from 0 to 340 CAG repeats (0 to 1,020 bp; Supplemental Information). No IRRs occur when the repeat was shorter than the read length (Figure 1, left panel). When the repeat is longer than the read length but shorter than the fragment length (Figure 1, middle panel), the number of anchored IRRs increases proportionally to the length of the repeat. As the length of the repeat approaches and exceeds the mean fragment length (Figure 1, right panel), the number of paired IRRs increases linearly with the length of the repeat. Because anchored IRRs require one of the reads to "anchor" outside of the repeat region, the number of anchored IRRs is limited by the fragment length and remains constant as the repeat grows beyond the mean fragment length. It is important to note that real sequence data may introduce additional challenges compared to the simulated data. For example, sequence quality in low complexity regions or interruptions in the repeat may impact the ability to identify some IRRs.

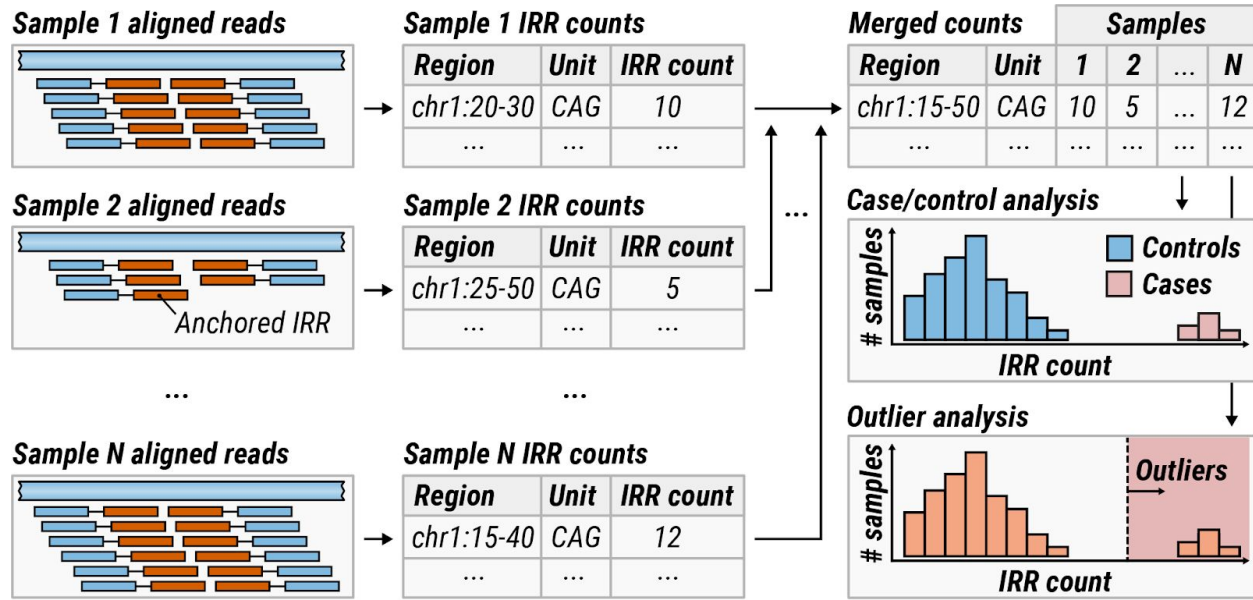


Figure 2. (Left) A search for anchored IRRs is performed across all aligned reads. (Middle) The IRR counts are summarized into STR profiles. (Right) The resulting STR profiles are merged across all samples. If the dataset can be partitioned into cases and controls, IRR counts in these groups are compared for each locus. Alternatively, if no such partition is possible, an outlier analysis is performed.

Analysis of sequencing data

Detection of expanded repeats in case-control studies

Given a sufficient number of samples with the same phenotype, pathogenic REs may be identified by searching for regions with significantly longer repeats in cases compared to controls (see Figure 2). To demonstrate the feasibility of such analyses, we analyzed 91 Coriell samples with experimentally-confirmed expansions in repeats associated with Friedreich's ataxia (FRDA; N=25), Myotonic Dystrophy type 1 (DM1; N=17), Huntington disease (HD; N=14), and fragile X syndrome (FXS; N=35). This dataset has been previously used to benchmark the performance of existing methods²³⁻²⁵.

The pathogenic cutoffs for FRDA, DM1, and FXS repeats are greater than the read length, so our analysis of simulated data suggests that anchored IRRs are likely to be present in each sample with one of these expansions (Figure 1). The pathogenic cutoff for the HD repeat (120 bp) is less than the read length (150 bp) used in this study, so a subset of samples with Huntington's disease may not contain relevant IRRs making this expansion harder to detect, however it is easily detectable with existing methods.

We separately compared samples with expansions in *FXN* (FRDA), *DMPK* (DM1), *HTT* (HD), or *FMR1* (FXS) genes (cases) against a control cohort of 150 unrelated Coriell samples of African, European, and East Asian ancestry³⁷. Each case-control comparison revealed a clear enrichment of anchored IRRs from the corresponding repeat region (Figure 3). This analysis demonstrated that ExpansionHunter Denovo can re-identify known pathogenic repeat expansions without prior knowledge of the location or repeat motif when the pathogenic repeat length is longer than (or nearly the size of) the read length and the repeat is highly penetrant.

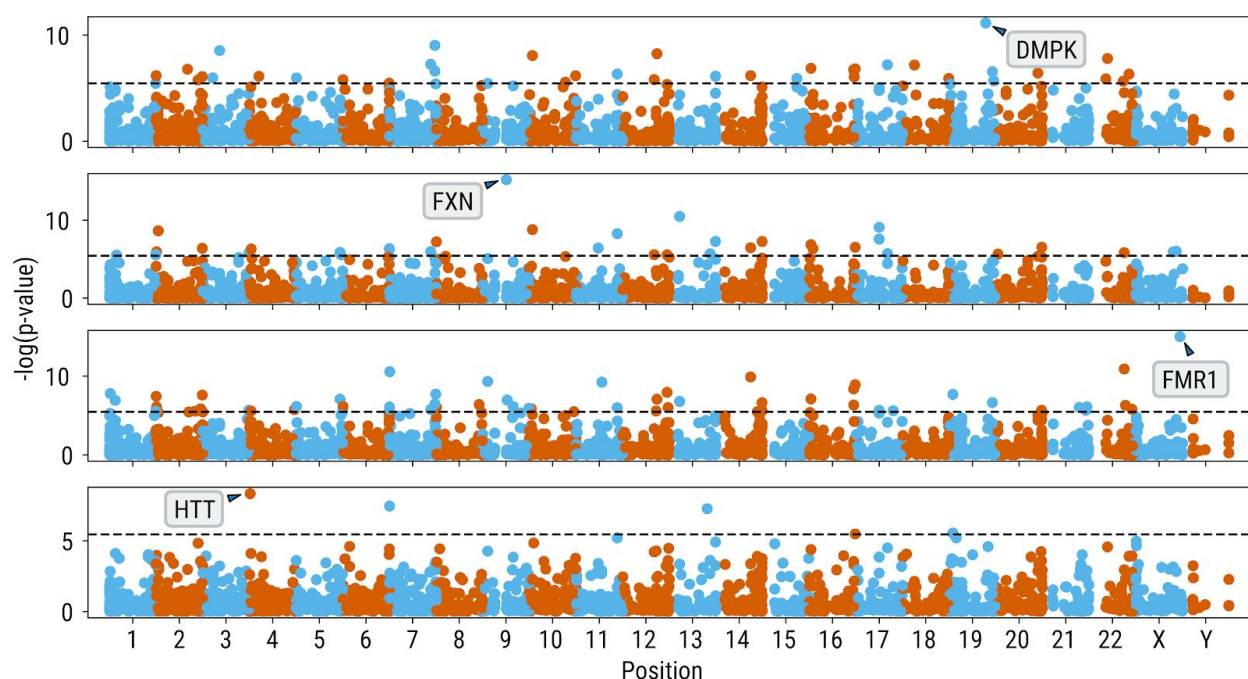


Figure 3: Genome-wide analysis of anchored IRRs comparing cases with known pathogenic expansions in *DMPK*, *FXN*, *FMR1* and *HTT* genes (top to bottom) to 150 controls.

Detection of expanded repeats in mixed sample cohorts

In many discovery projects it can be difficult to isolate patients that harbor the same repeat expansion based on the phenotype alone. For instance, (a) the repeat expansion in the *C9orf72* gene is present in fewer than 10% of ALS patients and (b) many ataxias can be caused by expansions of a variety of repeats. Such problems call for methods that are suitable for heterogeneous disease cohorts.

To solve this problem, we follow the approach taken previously by others and compare each case sample against the control cohort to identify outliers^{24,26} (Methods). To demonstrate the efficacy of this approach, we combined each sample from the pool of samples with expansions in *FXN*, *DMPK*, *HTT*, and *FMR1* genes with 150 controls to generate a total of 91 datasets, each containing 151 samples. We then performed an outlier analysis on the counts of anchored IRRs (Methods) in each dataset.

In 77% of the datasets, the expanded repeat ranked within the top 10 repeats based on the outlier score (Figure 4). This number increased to 82% when the analysis was restricted to short motifs between 2 and 6 bp. EHdn performed well for *DMPK* and *FXN* repeats, identifying these REs within the top 10 ranks for 41 out of 42 cases. The *FMR1* expansion was only ranked in the top 10 for 22 out of 35 cases known to have the expansion. This result is consistent with a previous comparison, which found this locus had poorest performance across all RE detection tools²⁴. The performance for the *HTT* repeat was the worst demonstrating an acknowledged limitation that EHdn was not designed to detect REs shorter than the read length. These results demonstrate that for expansions exceeding the read length EHdn has comparable performance to previously published catalog-based methods designed to specifically target known pathogenic REs^{23–25,27}.

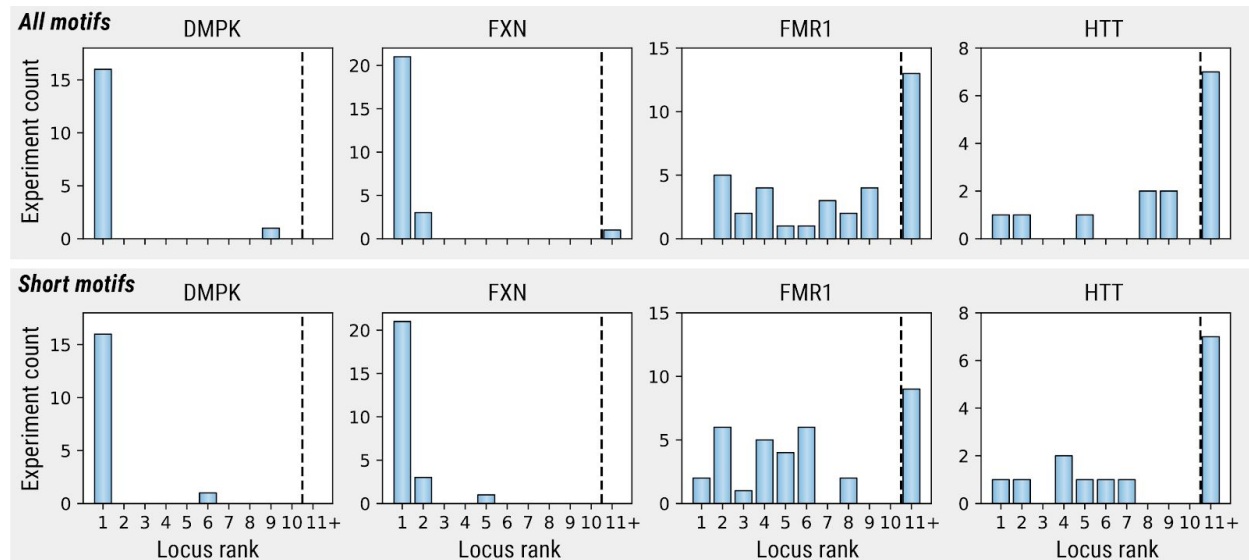


Figure 4: Ranking of known expansions based on the outlier score computed for anchored IRRs. Each rank originates from a genome-wide analysis of a dataset consisting of a single sample with a known expansion and 150 controls. (Top row) Ranks for all identified motifs. (Bottom row) Ranks for motifs 6bp and shorter.

The landscape of long repeats within a control population

To explore the landscape of large repeats in the general population, we applied EHdn to 150 unrelated Coriell samples of African, European, and East Asian ancestry³⁷. To limit this analysis to higher confidence repeats we considered loci where EHdn identified at least five anchored IRRs and motifs supported by at least five paired IRRs in a single sample. Altogether, EHdn identified 1,572 unique motifs spanning between two and 20 bp, 94% of which were longer than 6 bp. Of these, 19% were found in at least half of the samples and 23% were found in just one sample. On average, each person had 658 loci with long repeats. As expected, the telomeric motif AACCT is particularly abundant. It was found in about ~23,000 IRRs per sample. Similarly, the centromeric motif AATGG was found in ~5,000 IRRs per sample.

Exploring the limitations of catalog-based RE detection methods

Genome wide catalogs are limiting

Existing methods for repeat expansion detection rely upon predefined repeat catalogs^{22–26,38}. Although these methods can analyze user-defined catalogs, defining the genomic locations of repeats is a complex task and there is a risk of missing many potentially-pathogenic loci.

To evaluate the limitations of a catalog-based approach, we curated a set of 53 pathogenic or potentially pathogenic repeats (Table S1) and checked if they were present in two commonly-used catalogs: (a) STRs with up to 6bp motifs) from the UCSC genome browser simple repeats track^{39,40} utilized by STRetch and exSTRa and (b) the GangSTR catalog.

Nine of the pathogenic repeats are not present in the reference genome and hence are absent from both catalogs (Figure S3). For the remaining 44 loci, we considered a locus to be present in a catalog if it overlapped with one of the catalog's repeats by at least one base pair. Our comparison showed that 22 loci are present in both catalogs, 12 are missing from the GangSTR catalog and present in the UCSC catalog, five are missing from the UCSC catalog and present in the the GangSTR catalog, and five are missing from both catalogs (Figure S2). While it is possible to update the catalogs to include these known pathogenic repeats, the number of missing, potentially pathogenic, REs remains unknown. The lack of concordance between the two catalogs is also worth noting (Figure S2). EHdn has been designed to overcome the limitations inherent to catalog-based methods. To demonstrate the efficacy of our method we show that it can detect expansions of pathogenic repeats without *a priori* knowledge of their location or nucleotide composition through a series of simulations.

Simulated expansions of known pathogenic repeats

We simulated various pathogenic REs at multiple sizes larger than the read length (Methods). For each of these REs, we simulated reads from the expanded locus and then spiked them into a real WGS sample to keep our simulations as realistic as possible. The resulting samples were analyzed with EHdn and STRetch to benchmark the ability of each method to prioritize pathogenic expansions. Our comparisons were focused on STRetch because this method was specifically designed to search for novel expansions using a genome-wide catalog.

Our simulations show that EHdn ranks 29 out of 34 pathogenic repeats in the top 10 when their size exceeds the read length (Table S3-S6). STRetch prioritizes 25 out of 29 repeats in the top 10 and the five remaining repeats are missing from its catalog. Two expansions missed by both methods, located in the genes *NUTM2B-AS1* and *NOTCH2NLC*, have flanking sequences that are highly similar to other genomic regions⁴¹ making it challenging to uniquely align reads flanking the repeat locus. However, EHdn can detect expansions of these repeats when their size exceeds the fragment length through motif-based analysis, though it cannot place the expansions in any specific region. One of REs detected by EHdn and missed by STRetch is the pathogenic *CSTB* repeat with a motif length of 12bp (Supplemental Information). This is because STRetch is limited to detection of motifs with length up to six base pairs. To further highlight this strength of EHdn, we confirmed that it can detect other REs with long motifs (Supplemental Information).

Thus, EHdn offers similar performance to catalog-based methods on large expansions of known repeats despite performing a catalog-free genome-wide search, hence allowing detection of an important new class of REs.

Complex repeat expansions and insertions

Recent discoveries of pathogenic complex expansion/insertion events demonstrate the need for methods capable of detecting such expansions. One such example is the

recently-discovered repeat expansion causing cerebellar ataxia with neuropathy and bilateral vestibular areflexia syndrome (CANVAS)^{42,43}. Rafehi et al⁴³ demonstrated that EHdn is the only computational method capable of discovering this expansion.

Motivated by these findings, we sought to benchmark EHdn's ability to detect other recently-published REs with complex structure (Figure S3). We simulated nine such REs with four distinct configurations reported previously. For eight out of nine REs, including a simulated version of the CANVAS expansion, EHdn was able to detect one or both of the expanded repeats in each locus (Table S7). The expansion of the remaining locus, *SAMD12*, is also detectable by setting a more permissive mapping quality threshold for anchor reads (Table S7).

Discussion

Here, we introduced a new software tool, EHdn, that can identify novel REs using high-throughput WGS data. We tested EHdn by comparing samples with known REs against a control group of 150 diverse individuals and performed simulation studies across a range of pathogenic or potentially pathogenic REs. These analyses show that EHdn offers comparable performance to targeted methods on known pathogenic repeats while also being able to detect repeats absent from existing catalogs.

Recent discoveries have highlighted the importance of complex pathogenic repeat expansions involving non-reference insertions²⁸⁻³³. EHdn is currently the only method capable of discovering these expansions from BAM or CRAM files without the need for re-alignment of the supporting reads. Additionally, we anticipate that EHdn can replace existing more manual and less computationally efficient discovery pipelines, such as the TRhist-based pipeline⁴¹, where identification of enriched repeat motifs is followed by ad-hoc realignment of relevant reads to the reference genome and manual evaluation of loci where these reads align.

EHdn has some limitations and areas for further improvement. It is limited to the detection of repetitive sequence longer than the read length and cannot, in general, detect shorter expansions. However, detection of these shorter expansions are feasible with the existing catalog-based methods, or SV detection methods. It may be possible to extend the detection limit to shorter repeat expansions, however increasing the search space will lead to increased runtime and reduced power to detect outlier expansions.

In many previous studies, identification of pathogenic REs required years of work and involved linkage studies to isolate the region of interest followed by targeted sequencing to identify the likely causative mutations. EHdn can be used as a front-line tool in such studies to rapidly identify candidate REs. The loci flagged by EHdn can be further studied by defining custom input files describing these novel REs for analysis with targeted methods, as well as other molecular assays. The benefits of this approach were demonstrated in a recent study, where EHdn successfully identified a novel complex pathogenic RE⁴⁴.

Conclusions

We presented ExpansionHunter Denovo, a new genome-wide and catalog-free method to search for REs in WGS data. We demonstrated that EHdn consistently detects REs in real and simulated data. Given the widespread adoption of WGS for rare disease diagnosis, we expect that EHdn will enable further RE discoveries that will likely resolve the cause of disease in many individuals.

Methods

Identification of IRRs

To determine if a read r is an in-repeat read, we first check the read for periodicity. We define $I_k(i) = 1$ if $r_i = r_{i+k}$ and 0 otherwise, where r_i and r_{i+k} are the i th and $(i+k)$ th

bases of the read r . We then let $S(k) = \sum_{i=0}^{L-k-1} I_k(i) / (L-k)$ where L is the read length.

Note that if a read consists of a perfect stretch of repeat units of length k then $S(k) = 1$.

We search across of range of motif lengths (by default $k \in \{2, 3, \dots, 20\}$) for the

smallest k such that $S(k) \geq t$ where t is a set threshold (we use $t = 0.8$ in all our

analyses). If such a value of k is found, we extract the putative repeat unit using the

most frequent bases at each offset $0 \leq i \leq k - 1$. Since the orientation of the repeat

where a given IRR originated is unknown in general, the unit of the repeat is

ambiguous. To remove the ambiguity, we select the smallest repeat unit in

lexicographical order under circular permutation and reverse complement operations.

We then use this putative repeat unit to calculate a weighted-purity (WP) score of a

read²⁵. We assume that a read is an IRR if it achieves WP score of at least 0.9. The WP

score lowers the penalty for low-quality mismatches in order to account for the

possibility of an increased base-call error rate that may occur in highly repetitive regions

of the genome.

EHdn searches for IRRs among unaligned reads and reads whose mapping quality

(MAPQ) is below a set threshold which in the analysis presented here was set to 40.

For this study, we limited our analysis to motif lengths between two and 20 base pairs.

Motif lengths equal to one were excluded to eliminate the large number of homopolymer

repeats from the downstream analyses since we identified over 30 times as many

homopolymer IRRs as IRRs with longer repeat motifs.

EHdn designates a read pair as a paired IRR if both mates are IRRs with the same repeat motif. A read is designated as an anchored IRR if it is an IRR and its mate is not an IRR and has MAPQ above a set threshold which was set to 50 for this study. Parameters such as the maximum allowed MAPQ for an IRR, the minimum allowed MAPQ for an anchor, and the range of repeat unit lengths for which to search are all tunable with EHdn. For example, setting the anchor read MAPQ threshold to 0 and the IRR MAPQ threshold to 60 ensures that every read pair in the alignment file is analyzed (assuming that the MAPQ values range from 0 to 60) at the cost of a corresponding increase in runtime.

Merging IRRs

Because an anchored IRR is assigned to the location of the aligned anchor read and not the position of the actual repeat (whose exact location may be unknown), a single repeat may produce anchored IRRs at a variety of locations centered around the repeat. To account for this, anchored IRRs with the same repeat motif are merged if their anchors are aligned within 500 bp of one another. When multiple samples are analyzed, the anchor regions are also merged across all samples and the counts of anchored IRRs (normalized to 40x read depth) are tabulated for each merged region and sample. Additionally, the depth-normalized counts of paired IRRs are tabulated for each repeat motif and sample.

Prioritization of expanded repeats

EHdn supports case-control and outlier analyses of the underlying dataset. The case-control analysis is based on a one-sided Wilcoxon rank-sum test. It is appropriate for situations where a significant subset of cases is expected to contain expansions of the same repeat.

The outlier analysis is appropriate for heterogeneous cohorts where enrichment for any specific expansion is not expected. The outlier analysis bootstraps the sampling

distribution of the 95% quantile and then calculates the z-scores for cases that exceed the mean of this distribution. The z-scores are used for ranking the repeat regions. Similar outlier-detection frameworks were also developed for exSTRa²⁴ and STRetch²⁶.

Both the case-control and the outlier analyses can be applied either to the counts of anchored IRRs or to the counts of paired IRRs. We refer to these as locus or motif methods, respectively. The high-ranking regions flagged by the analysis of anchored IRRs correspond to approximate locations of putative repeat expansions. The high-ranking motifs flagged by the analysis of paired IRRs correspond to the overall enrichment for repeats with that motif.

Defining relevant repeat expansions

A catalogue of pathogenic or potentially pathogenic repeat expansions was collated from the literature. We supplemented this catalog with recently reported STRs linked with gene expression⁴⁵, and repeats with longer motifs overlapping with disease genes (Supplemental Information).

Simulated repeat expansions

Expanded repeats were simulated using a strategy similar to that taken by BamSurgeon⁴⁶. Briefly, we simulated reads in a 2Kb region around an expanded repeat and then aligned the reads to the reference genome. We then removed reads in the same region from a WGS control sample and merged alignments of real and simulated data together (Figure S2; Supplemental Information).

Human WGS data

The control WGS samples are from Illumina Polaris dataset³⁷. All samples were sequenced on an Illumina HiSeqX instrument using TruSeq DNA PCR-free sample

prep. The 91 Coriell samples with experimentally-confirmed repeat expansions in *DMPK*, *FMR1*, *FXN* and *HTT* were introduced in our earlier publication²⁵.

Funding

BT was funded by the Canadian Institutes for Health Research Banting Postdoctoral Fellowship and the Canadian Open Neuroscience Platform Research Scholar Award. MB was supported by an Australian National Health and Medical Research Council (NHMRC) Program Grant (GNT1054618) and an NHMRC Senior Research Fellowship (GNT1102971). This work was made possible through Victorian State Government Operational Infrastructure Support and Australian Government NHMRC IRIISS.

References

1. Muir, P. *et al.* The real cost of sequencing: scaling computation to keep pace with data generation. *Genome Biol.* **17**, 53 (2016).
2. Erikson, G. A. *et al.* Whole-Genome Sequencing of a Healthy Aging Cohort. *Cell* **165**, 1002–1011 (2016).
3. Telenti, A. *et al.* Deep sequencing of 10,000 human genomes. *Proc. Natl. Acad. Sci. U. S. A.* **113**, 11901–11906 (2016).
4. Gudbjartsson, D. F. *et al.* Large-scale whole-genome sequencing of the Icelandic population. *Nat. Genet.* **47**, 435–444 (2015).
5. Nagasaki, M. *et al.* Rare variant discovery by deep whole-genome sequencing of 1,070 Japanese individuals. *Nat. Commun.* **6**, 8018 (2015).

6. 1000 Genomes Project Consortium *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
7. Consortium, P. M. A. S. & Project MinE ALS Sequencing Consortium. Project MinE: study design and pilot analyses of a large-scale whole-genome sequencing study in amyotrophic lateral sclerosis. *European Journal of Human Genetics* **26**, 1537–1546 (2018).
8. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
9. Raczy, C. *et al.* Isaac: ultra-fast whole-genome secondary analysis on Illumina sequencing platforms. *Bioinformatics* **29**, 2041–2043 (2013).
10. Rimmer, A. *et al.* Integrating mapping-, assembly- and haplotype-based approaches for calling variants in clinical sequencing applications. *Nat. Genet.* **46**, 912–918 (2014).
11. Poplin, R. *et al.* Creating a universal SNP and small indel variant caller with deep neural networks. *bioRxiv* 092890 (2016). doi:10.1101/092890
12. Garrison, E. & Marth, G. Haplotype-based variant detection from short-read sequencing. *arXiv [q-bio.GN]* (2012).
13. Poplin, R. *et al.* A universal SNP and small-indel variant caller using deep neural networks. *Nat. Biotechnol.* **36**, 983–987 (2018).
14. Roller, E., Ivakhno, S., Lee, S., Royce, T. & Tanner, S. Canvas: versatile and scalable detection of copy number variants. *Bioinformatics* **32**, 2375–2377 (2016).
15. Abyzov, A., Urban, A. E., Snyder, M. & Gerstein, M. CNVnator: an approach to discover, genotype, and characterize typical and atypical CNVs from family and population genome sequencing. *Genome Res.* **21**, 974–984 (2011).
16. Chen, X. *et al.* Manta: rapid detection of structural variants and indels for germline and cancer sequencing applications. *Bioinformatics* **32**, 1220–1222 (2016).

17. Layer, R. M., Chiang, C., Quinlan, A. R. & Hall, I. M. LUMPY: a probabilistic framework for structural variant discovery. *Genome Biol.* **15**, R84 (2014).
18. DeJesus-Hernandez, M. *et al.* Expanded GGGGCC hexanucleotide repeat in noncoding region of C9ORF72 causes chromosome 9p-linked FTD and ALS. *Neuron* **72**, 245–256 (2011).
19. Renton, A. E. *et al.* A hexanucleotide repeat expansion in C9ORF72 is the cause of chromosome 9p21-linked ALS-FTD. *Neuron* **72**, 257–268 (2011).
20. La Spada, A. R. & Paul Taylor, J. Repeat expansion disease: progress and puzzles in disease pathogenesis. *Nat. Rev. Genet.* **11**, 247–258 (2010).
21. Hannan, A. J. Tandem repeats mediating genetic plasticity in health and disease. *Nat. Rev. Genet.* **19**, 286–298 (2018).
22. Tang, H. *et al.* Profiling of Short-Tandem-Repeat Disease Alleles in 12,632 Human Whole Genomes. *Am. J. Hum. Genet.* **101**, 700–715 (2017).
23. Mousavi, N., Shleizer-Burko, S. & Gymrek, M. Profiling the genome-wide landscape of tandem repeat expansions. *bioRxiv* 361162 (2018). doi:10.1101/361162
24. Tankard, R. M. *et al.* Detecting Expansions of Tandem Repeats in Cohorts Sequenced with Short-Read Sequencing Data. *Am. J. Hum. Genet.* **103**, 858–873 (2018).
25. Dolzhenko, E. *et al.* Detection of long repeat expansions from PCR-free whole-genome sequence data. *Genome Res.* **27**, 1895–1903 (2017).
26. Dashnow, H. *et al.* STRetch: detecting and discovering pathogenic short tandem repeat expansions. *Genome Biol.* **19**, 121 (2018).
27. Dolzhenko, E. *et al.* ExpansionHunter: a sequence-graph-based tool to analyze variation in short tandem repeat regions. *Bioinformatics* **35**, 4754–4756 (2019).
28. Sato, N. *et al.* Spinocerebellar ataxia type 31 is associated with ‘inserted’ penta-nucleotide

- repeats containing (TGGAA)*n*. *Am. J. Hum. Genet.* **85**, 544–557 (2009).
29. Seixas, A. I. *et al.* A Pentanucleotide ATTTTC Repeat Insertion in the Non-coding Region of DAB1, Mapping to SCA37, Causes Spinocerebellar Ataxia. *Am. J. Hum. Genet.* **101**, 87–103 (2017).
 30. Ishiura, H. *et al.* Expansions of intronic TTTC A and TTTTA repeats in benign adult familial myoclonic epilepsy. *Nat. Genet.* **50**, 581–590 (2018).
 31. Corbett, M. A. *et al.* Intronic ATTTTC repeat expansions in STARD7 in familial adult myoclonic epilepsy linked to chromosome 2. *Nat. Commun.* **10**, 4920 (2019).
 32. Florian, R. T. *et al.* Unstable TTTTA/TTTCA expansions in MARCH6 are associated with Familial Adult Myoclonic Epilepsy type 3. *Nat. Commun.* **10**, 4919 (2019).
 33. Yeetong, P. *et al.* TTTCA repeat insertions in an intron of YEATS2 in benign adult familial myoclonic epilepsy type 4. *Brain* **142**, 3360–3366 (2019).
 34. LaCroix, A. J. *et al.* GGC Repeat Expansion and Exon 1 Methylation of XYLT1 Is a Common Pathogenic Variant in Baratela-Scott Syndrome. *Am. J. Hum. Genet.* **104**, 35–44 (2019).
 35. Lalioti, M. D. *et al.* Dodecamer repeat expansion in cystatin B gene in progressive myoclonus epilepsy. *Nature* **386**, 847–851 (1997).
 36. Ashley, E. A. Towards precision medicine. *Nature Reviews Genetics* **17**, 507–522 (2016).
 37. Illumina. Illumina/Polaris. *GitHub* Available at: <https://github.com/Illumina/Polaris>. (Accessed: 20th November 2019)
 38. Bahlo, M. *et al.* Recent advances in the detection of repeat expansions with short-read next-generation sequencing. *F1000Res.* **7**, (2018).
 39. Benson, G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res.* **27**, 573–580 (1999).

40. Kent, W. J. *et al.* The human genome browser at UCSC. *Genome Res.* **12**, 996–1006 (2002).
41. Ishiura, H. *et al.* Noncoding CGG repeat expansions in neuronal intranuclear inclusion disease, oculopharyngodistal myopathy and an overlapping disease. *Nat. Genet.* **51**, 1222–1232 (2019).
42. Cortese, A. *et al.* Biallelic expansion of an intronic repeat in RFC1 is a common cause of late-onset ataxia. *Nat. Genet.* **51**, 649–658 (2019).
43. Rafehi, H. *et al.* Bioinformatics-Based Identification of Expanded Repeats: A Non-reference Intronic Pentamer Expansion in RFC1 Causes CANVAS. *Am. J. Hum. Genet.* **105**, 151–165 (2019).
44. Rafehi, H. *et al.* Validation of new tools to identify expanded repeats: an intronic pentamer expansion in RFC1 causes CANVAS. doi:10.1101/597781
45. Fotsing, S. F. *et al.* The impact of short tandem repeat variation on gene expression. *Nat. Genet.* **51**, 1652–1659 (2019).
46. Ewing, A. D. *et al.* Combining tumor genome simulation with crowdsourcing to benchmark somatic single-nucleotide-variant detection. *Nat. Methods* **12**, 623–630 (2015).