

SHORT REPORT

Open Access



Association analysis of repetitive elements and R-loop formation across species

Chao Zeng^{1,2*} , Masahiro Onoguchi^{1,2}  and Michiaki Hamada^{1,2,3,4*} 

Abstract

Background: Although recent studies have revealed the genome-wide distribution of R-loops, our understanding of R-loop formation is still limited. Genomes are known to have a large number of repetitive elements. Emerging evidence suggests that these sequences may play an important regulatory role. However, few studies have investigated the effect of repetitive elements on R-loop formation.

Results: We found different repetitive elements related to R-loop formation in various species. By controlling length and genomic distributions, we observed that satellite, long interspersed nuclear elements (LINEs), and DNA transposons were each specifically enriched for R-loops in humans, fruit flies, and *Arabidopsis thaliana*, respectively. R-loops also tended to arise in regions of low-complexity or simple repeats across species. We also found that the repetitive elements associated with R-loop formation differ according to developmental stage. For instance, LINEs and long terminal repeat retrotransposons (LTRs) are more likely to contain R-loops in embryos (fruit fly) and then turn out to be low-complexity and simple repeats in post-developmental S2 cells.

Conclusions: Our results indicate that repetitive elements may have species-specific or development-specific regulatory effects on R-loop formation. This work advances our understanding of repetitive elements and R-loop biology.

Keywords: R-loop, Repetitive element, Transposable element

Introduction

An R-loop is a three-stranded structure composed of a DNA:RNA hybrid with a displaced single-stranded DNA (ssDNA). Although R-loops were initially considered to be rare by-products of transcription, recent reports have suggested that R-loops are widely distributed in eukaryotic genomes [1–3] and are involved in gene regulation and genome integrity [4–9]. R-loops regulate gene expression through a variety of molecular mechanisms. For instance,

in promoter regions, R-loops promote or inhibit gene transcription by decreasing methylation [10] or enhancing polycomb-mediated gene silencing [11]. Surprisingly, a recent report indicated that R-loops can shield the ribosomal gene expression by RNA polymerases II from the transcription conflicts caused by other RNA polymerases [12]. In terminator regions, R-loops can induce RNA polymerase stalling to improve the efficiency of transcription termination [13]. Simultaneously, these R-loops will promote nascent RNA cleavage and 3' transcript degradation [14]. With respect to genome integrity, R-loops can affect genome dynamics [15] and telomere stability [16]. In centromeric regions, R-loops are reported to maintain genome stability by promoting chro-

*Correspondence: chao.zeng@hamadalab.com; mhamada@waseda.jp

¹AIST-Waseda University Computational Bio Big-Data Open Innovation Laboratory (CBBDOIL), National Institute of Advanced Industrial Science and Technology, 63-520, 3-4-1, Okubo Shinjuku-ku, 169-8555 Tokyo, Japan

²Faculty of Science and Engineering, Waseda University, 55N-06-10, 3-4-1 Okubo Shinjuku-ku, 169-8555 Tokyo, Japan

Full list of author information is available at the end of the article



© The Author(s). 2021 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

matin condensation [15] or chromosome segregation [17]. In telomeric regions, telomere repeat-containing RNAs (TERRAs) preferentially accumulate in short telomeres to form telomere-repairing R-loops [16]. Previous reviews described other details of R-loops and genome integrity [4–9]. Additionally, with an increasing number of R-loop-related diseases (such as neurological diseases and cancers) being reported and studied [18–22], a deep understanding of the biology of R-loop structures is critical.

Genome-wide mapping of R-loop structures, performed with immunoprecipitation-based high-throughput sequencing [3, 23], provides a global approach for quantitative analysis and systematic characterization of R-loops. This technique involves immunoprecipitating R-loop structures using a DNA-RNA hybrid antibody (S9.6) and then subjecting those R-loop-forming DNAs to direct DNA-sequencing, known as DRIP-seq (DNA-RNA immunoprecipitation followed by high-throughput DNA sequencing) [23]. Additionally, other R-loop profiling methods have been developed in the past decade by modifying antibodies and/or protocols [2, 24–29].

Previous studies have uncovered several characteristics of R-loop formation. For example, within the genome, R-loops are prone to be formed in GC-skew (asymmetric strand distribution of guanines and cytosines) [23, 30] or AT-skew (asymmetric strand distribution of adenines and thymines) [1, 2] regions. The stability of G quadruplex (G4) structures is also related to R-loop formation [31]. Notably, multiple reports have indicated that R-loop formation is associated with short tandem repeats, especially trinucleotide repeat (TNR) expansion [32–35]. An R-loop predictive model was designed with the feature of short tandem repeats in mind [36]. Additionally, limited studies have shown that transposable elements (TEs), including Ty1 [37] and LINEs [38], may play a role in the formation of R-loops. Specific epigenetic signatures also connect with R-loop structures [3, 38]. Although R-loops are thought to form co-transcriptionally or in cis (transcription and R-loop formation at the same locus) [6], accumulating evidence has indicated that R-loop structures can also form post-transcriptionally or in trans (RNA transcribed from a locus hybridizes to a distal locus) [39–43]. However, our understanding of R-loop formation and their sequence characteristics in the genome is still limited.

In this study, we focused on the effects of repetitive elements on R-loop formation. By comparing the frequency of repetitive elements in R-loops and the controls, we derive the repeat class or family that is associated with R-loop formation. For this purpose, we separately prepared three groups of controls, corresponding to uniform distribution in the genome, specific length and genomic distributions, and co-transcriptional formation. The three controls may not be appropriate for all the data studied;

for example, we found that more than 60% of R-loops in fruit fly do not overlap with the transcribed region (implying that the third control is not preferred). We observed different repetitive elements related to R-loop formation in various species. Based on the second control mentioned above, we discovered that satellites, LINEs, and DNA transposons were each enriched for R-loops in humans, fruit flies, and *Arabidopsis thaliana*, respectively. Additionally, R-loops were mainly found in regions of low-complexity or simple repeats across these three species. Interestingly, we also found that the repetitive elements associated with R-loop formation differ according to the organism's developmental stage. For fruit flies, LINEs and LTRs are more likely to contain R-loops in embryos, which changes to R-loops being more prevalent in areas of low-complexity and simple repeats in S2 cells. To our knowledge, this is the first comprehensive analysis of the association between repetitive elements and R-loop formation. This work improves our insights into the potential functions of repetitive elements and our understanding of the biological mechanisms underlying R-loop formation.

Methods

Genomic sequences and gene annotations of human (hg38) and fruit fly (dm6) were downloaded from the UCSC genome browser [44]. The *A. thaliana* genome (TAIR10) and gene annotation were obtained from Ensembl [45]. Repetitive sequences were downloaded from the UCSC repeatmasker track [44]. We extracted genome-wide R-loop regions and transcribed regions from DRIP-seq and GRO-seq (global run-on sequencing) data, respectively, following the steps mentioned below. Bedtools (v2.25.0) [46] were used to extract (intersect sub-command) or remove (subtract sub-command) overlaps between regions. Statistics and enrichment analysis were implemented in home-made Python scripts.

DRIP-seq analysis

Reads were aligned to the corresponding genome using BWA-MEM (0.7.17-r1188) [47] with default parameters. For paired-end reads (fruit fly and *A. thaliana*), reads aligned in a proper pair (SAMtools view -f 2)[48] were considered as mapped reads. For single-end reads (human), we extracted mapped reads (SAMtools view -F -4) for subsequent analysis. Mapped reads were sorted by SAMtools, and polymerase chain reaction (PCR) duplicates were marked with Picard MarkDuplicates (v2.18.1) [49] using default parameters. Subsequently, we discarded the read duplicates from the alignments (SAMtools view -h -F 1024).

We used DRIP-seq data from IP (immunoprecipitation), Input, and RnaseH-treated samples, separately. Using the input sample as a control, we extracted peaks in the IP and

RnaseH-treated samples. MACS (v2.2.7.1) was applied to detect peaks from the alignments. In addition to the same MACS parameters (-q 0.001 -broad -broad-cutoff 0.001 -keep-dup all), we used -g hs -f BAM, -g dm -f BAMPE, and -g 1.36e8 -f BAMPE for human, fruit fly, and *A. thaliana*, respectively. We selected an output file (.broadPeak) for the final peak detection results. After removing the peaks in the IP sample that overlapped with the peaks in the RnaseH-treated sample, we obtained the final R-loop peaks.

GRO-seq analysis

For short reads with lengths of less than 100 nt (human and fruit fly), we sequentially utilized BWA-ALN and BWA-SAMSE with default parameters for mapping. Alternatively, for *A. thaliana*, we used BWA-MEM using default parameters. Homer (v4.11) [50] (findPeaks -style groseq) was applied to detect peaks (indicating transcripts). In cases where a tissue or cell line corresponded to multiple samples (e.g., biological replicates), we used BEDtools (merge sub-command) to combine peaks from all samples.

Enrichment analysis

To investigate the enrichment of repetitive elements in the R-loop-forming regions, we calculated the percentage of bases in the repetitive elements in the R-loops. For this purpose, we also prepared three control groups based on different hypotheses. First, assuming that R-loops are randomly distributed in the genome, we calculated the content of the repetitive elements in the genome as a control (referred to as “genome control”). Second, assuming that R-loops have a preferential distribution of length and genomic locations, we randomly selected 1000 groups of sequences from the genome while maintaining the same number of R-loops as well as length and genomic location (referred to as “sampling control”). Finally, assuming that R-loops are overwhelmingly formed directly where they are transcribed, we used the transcribed regions defined in the GRO-seq data as a control (referred to as “GRO control”).

For genome and GRO controls, we calculated the percentage of the bases of the repetitive elements in R-loops and control sequences x_k and y_k , respectively, and then computed the

$$\text{enrichment}_k = \log_{10} \left(\frac{x_k}{y_k} \right), \quad (1)$$

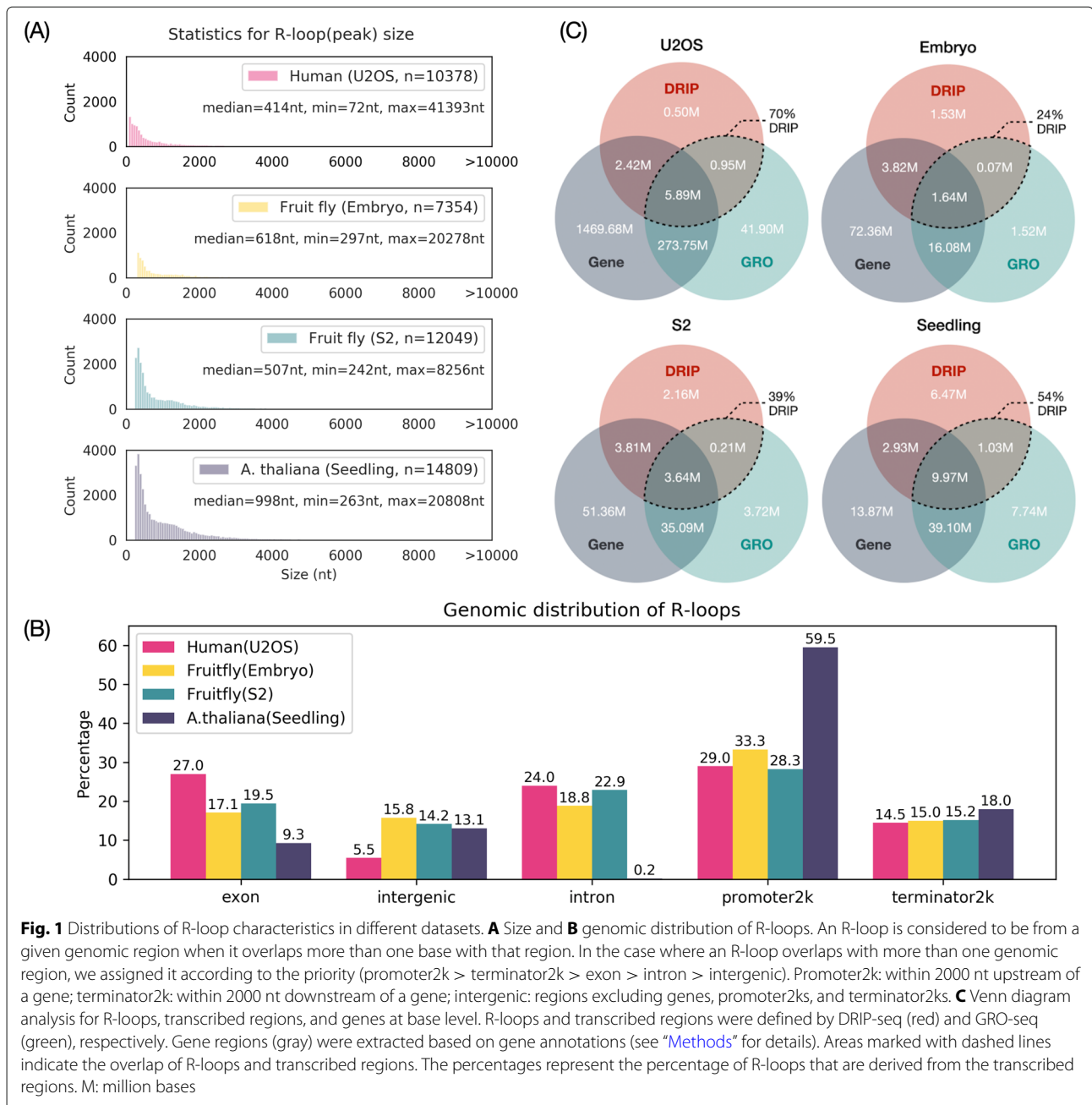
where k denotes a repetitive element (repeat class or family). For the sampling control, we computed the percentage of the bases of repetitive elements in R-loops and control groups as p_k and $\{q_{k,j}\}$ ($j \in \{1, \dots, 1000\}$), respectively. Then, we calculated the Z score of p_k in $\{q_{k,j}\}$ ($j \in \{1, \dots, 1000\}$) as the enrichment metric.

Results and discussion

Different controls provide multiple perspectives on the association between repetitive elements and R-loop formation

To investigate the association between repetitive elements and R-loop formation, we asked whether there was an overrepresented or underrepresented repetitive element in the R-loops. For this purpose, we considered three controls. First, R-loops are randomly and uniformly distributed in the genome. Second, R-loops of specific lengths are distributed in specific regions of the genome. Finally, most R-loops are co-transcriptional and are more likely to form where nascent transcripts are produced [6]. Thus, we prepared a genome control, sampling control, and GRO control to calculate the enrichment of repetitive elements in R-loops under the corresponding conditions (see “Methods” for details). We observed from the DRIP-seq data that the median R-loop lengths in animals (human and fruit fly) ranged from 414 to 618 nt, whereas those in plants (*A. thaliana*) were longer, with a median of 998 nt (Fig. 1A). In addition, only 5.5%–15.8% of R-loops were detected in the intergenic region (Fig. 1B), which implies that the majority of R-loops preferentially form surrounding gene regions, which is consistent with previous reports [3, 24, 51]. In addition, R-loops formed most frequently in promoter regions across species (Fig. 1B). For *A. thaliana*, in particular, up to 60% of R-loops were distributed in promoter regions, which reduces the occurrence of R-loops in other regions of the genome, especially in introns, which account for only 0.2% of R-loops formed. To this end, a sampling control was randomly generated in different datasets, simulating the distribution of R-loop lengths and genomic locations.

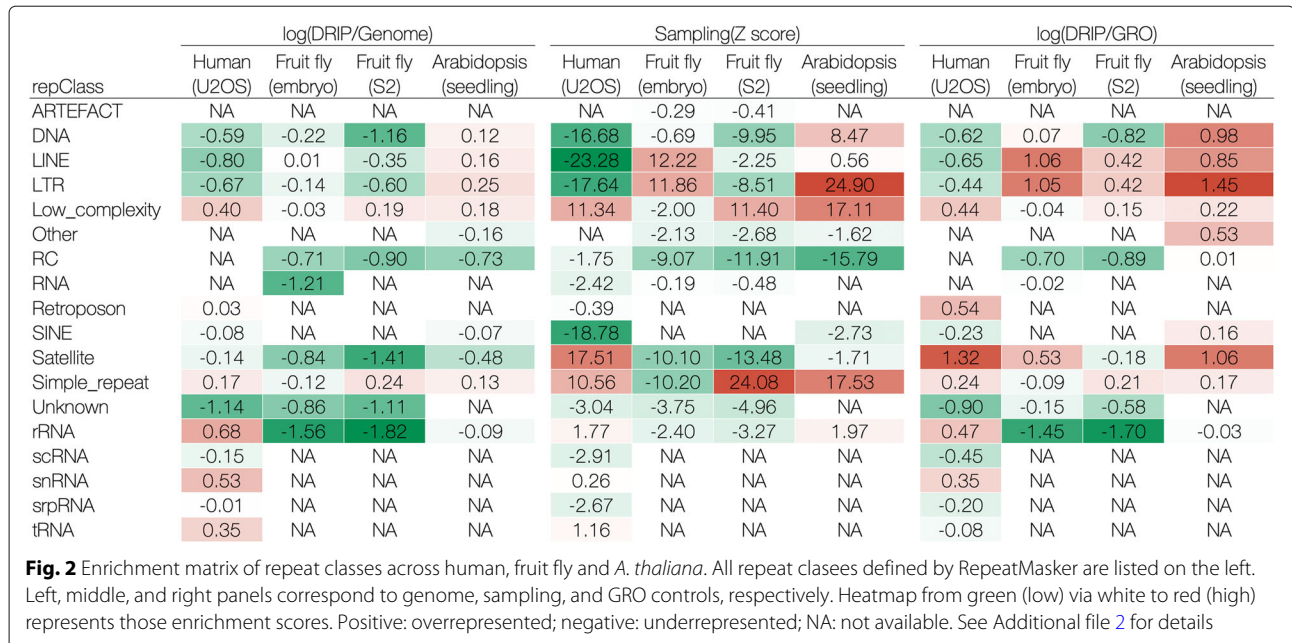
Considering that R-loops may form in irregularly transcribed regions (e.g., enhancers, promoters) that are not detectable by ordinary RNA-seq, we used GRO-seq data that can identify nascent transcripts to define transcriptional regions (including genes) in the genome for different tissues or cell lines. We then randomly generated the GRO control from real-time transcriptional regions. Notably, according to the GRO-seq data, the proportion of nascent transcripts detected in the intergenic region varied widely among species (3.18%–10.43%, Additional file 1), and was highest in humans. We speculate that the reason for this is the presence of a fraction of unannotated genes, and U2OS is a cell line that may have cancer-specific transcripts. To verify whether the majority of R-loops are co-transcriptional, we checked the percentage of overlaps between R-loops (as defined by DRIP-seq) and transcribed regions (as defined by GRO-seq). Although humans have 10.43% of transcribed regions in intergenic areas (see Additional file 1), the overlap between R-loops and these transcribed regions is higher (70%) compared to the other two species (Fig. 1C). The lowest overlap



between R-loops and transcribed regions was observed in fruit flies (embryos, 24%). It can be observed that all the transcribed regions defined in the GRO-seq data had a high match with the annotated gene regions (84.85%–91.8%, Additional file 1). Therefore, we speculate that the R-loops in fruit flies (embryos) may be formed by transcripts derived from distant or other chromosomal regions.

Repetitive elements are associated with r-loop formation across species and developmental stages

By first comparing the enrichment of repetitive elements in humans, fruit flies (S2 cells), and *A. thaliana*, we observed some patterns of relative consistency among species (Fig. 2, middle panel). In the sampling control, we observed that low-complexity and simple repeats tended to be enriched in R-loops among species, while



rolling circle (RC) appeared less frequently in R-loops than in random cases. Consistent with the previous report [52], DNA, LINEs, and LTRs were underrepresented in animal (human and fruit fly) R-loops; however, they were overrepresented in plant (*A. thaliana*) R-loops. In human and *A. thaliana* R-loops, short interspersed nuclear elements (SINEs) were underrepresented, whereas rRNA was mildly enriched. We also found some species-specific patterns. For example, satellites were significantly enriched in human R-loops, but significantly under-enriched in fruit fly R-loops. Notably, when we switched the controls to the transcriptional region (GRO control), we found that retrotransposons and satellites were enriched in human and *A. thaliana* R-loops, respectively, suggesting a positive correlation between the formation of cis R-loops and these two repetitive elements. Additionally, in humans and *A. thaliana*, the high degree of overlap between R-loops and transcriptional regions (Fig. 1C) enhances the confidence of the cis formation of R-loops in these two species.

We asked whether the association between repetitive elements and R-loop formation varies during development. For this purpose, we compared embryos and S2 cells of the fruit fly. Surprisingly, in the sampling control, we found that R-loops were more likely to form in regions of the embryo containing LINEs and LTRs (Fig. 2, middle panel). In post-developmental S2 cells, R-loops were highly aggregated in regions with low complexity or simple repeats. However, when we applied the GRO control (Fig. 2, right panel), the above results were relatively attenuated (i.e., low-complexity and simple repeats), as even

LINEs and LTRs were not separated in the embryo and S2 cells showing enrichment in the R-loops. Note that, in fruit flies, a small fraction (up to 39%) of R-loops overlapping with the transcribed regions might be a source of this inconsistency (Fig. 1C).

Various repeat families associated with r-loop formation

Further, we investigated the relationship between R-loop formation and the repetitive elements in each species at the repeat family level. For humans, we consistently observed an enrichment of R-loops, in the sampling and the GRO controls, containing cent, telo, low-complexity, simple repeat, snRNA and rRNA in the genome (Fig. 3). Interestingly, for the GRO control, we found that SVA repeat elements (belonging to the retroposon class) and satellites preferentially occurred in the R-loops.

For fruit fly embryos, we consistently observed the enrichment of specific repeat families in R-loops in both the sampling control and the GRO control (Fig. 4). For example, Gypsy and Pao, which are both LTRs; Jockey, R1, CR1, and LOA, all belonging to the LINEs; TcMar-Tc1 and PiggyBac, both of which are in the GRO control, I and R2 belonging to LINEs, copia (belonging to the LTRs), and satellites also tended to be enriched in R-loops. For S2 cells, the number of enriched repeat families in R-loops was significantly reduced compared to that in embryos (Fig. 4). In addition to the enrichment of Pao and CR1 that can still be observed in R-loops (which is consistent with the results in embryos), we also observed an increase in simple repeats and low-complexity in the R-loops of S2 cells.

repFamily	%Genome	%Sampling	%GRO	%DRIP	log(DRIP/Genome)	Sampling(Z score)	log(DRIP/GRO)
L1	17.06	11.78	9.89	1.39	-1.09	-21.29	-0.85
Alu	10.22	13.99	15.28	8.38	-0.09	-17.66	-0.26
ERVL-MaLR	3.62	2.67	2.11	0.64	-0.75	-14.56	-0.52
L2	3.36	3.40	4.11	1.77	-0.28	-12.76	-0.37
ERV1	2.80	2.47	1.62	0.75	-0.57	-9.30	-0.33
MIR	2.73	3.03	3.15	2.38	-0.06	-8.05	-0.12
centr	2.31	0.08	0.02	1.14	-0.31	18.96	1.75
ERVL	1.86	1.42	1.10	0.29	-0.80	-10.23	-0.57
hAT-Charlie	1.49	1.54	1.76	0.55	-0.43	-13.41	-0.51
Simple_repeat	1.23	1.29	1.05	1.83	0.17	10.56	0.24
TcMar-Tigger	1.18	1.04	1.35	0.15	-0.90	-10.20	-0.95
CR1	0.39	0.36	0.40	0.10	-0.57	-7.49	-0.59
ERVK	0.31	0.33	0.22	0.16	-0.30	-2.56	-0.15
hAT-Tip100	0.28	0.25	0.16	0.09	-0.49	-5.19	-0.24
Low_complexity	0.20	0.25	0.18	0.51	0.40	11.34	0.44
SVA	0.14	0.17	0.04	0.15	0.03	-0.39	0.54
Satellite	0.13	0.15	0.05	0.08	-0.22	-1.26	0.23
Gypsy	0.12	0.09	0.04	0.02	-0.69	-3.52	-0.28
RTE-X	0.11	0.10	0.10	0.02	-0.77	-3.92	-0.73
hAT-Blackjack	0.11	0.09	0.07	0.02	-0.71	-3.96	-0.50
TcMar-Mariner	0.09	0.08	0.07	0.01	-1.17	-3.61	-1.03
TcMar-Tc2	0.05	0.05	0.06	0.01	-0.56	-2.68	-0.59
RTE-BovB	0.04	0.04	0.04	0.01	-0.89	-3.40	-0.92
hAT	0.04	0.04	0.03	0.01	-0.52	-3.29	-0.37
Unknown	0.02	0.02	0.01	0.00	-1.14	-3.04	-0.90
MULE-MuDR	0.02	0.01	0.01	0.00	-0.68	-0.90	-0.47
tRNA-RTE	0.02	0.02	0.03	0.00	-0.75	-3.07	-0.81
LTR	0.02	0.02	0.01	0.00	NA	-2.33	NA
PiggyBac	0.02	0.02	0.02	0.01	-0.48	-1.19	-0.62
hAT-Ac	0.01	0.01	0.01	0.01	-0.16	-0.86	-0.09
tRNA	0.01	0.01	0.02	0.01	-0.01	-0.39	-0.17
Helitron	0.01	0.01	0.01	0.00	NA	-1.75	NA
snRNA	0.01	0.04	0.02	0.04	0.53	0.26	0.35
DNA	0.01	0.01	0.01	0.00	NA	-1.93	NA
telo	0.01	0.03	0.02	0.55	1.76	20.87	1.53
srpRNA	0.01	0.03	0.01	0.01	-0.01	-2.67	-0.20
5S-Deu-L2	0.01	0.01	0.01	0.00	-0.40	-1.20	-0.48
rRNA	0.01	0.02	0.01	0.04	0.68	1.77	0.47
scRNA	0.00	0.01	0.01	0.00	-0.15	-2.91	-0.45
RNA	0.00	0.01	0.00	0.00	NA	-2.42	NA
Dong-R4	0.00	0.00	0.00	0.00	NA	-0.88	NA
Penelope	0.00	0.00	0.00	0.00	-0.34	-0.75	-0.31
acro	0.00	0.00	0.00	0.00	NA	-0.26	NA
tRNA-Deu	0.00	0.00	0.00	0.00	NA	-1.35	NA
hAT-Tag1	0.00	0.00	0.00	0.00	NA	-0.64	NA
TcMar	0.00	0.00	0.00	0.00	NA	-0.63	NA
Merlin	0.00	0.00	0.00	0.00	NA	-0.34	NA
TcMar-Pogo	0.00	0.00	0.00	0.00	NA	-0.11	NA
PIF-Harbinger	0.00	0.00	0.00	0.00	NA	-0.28	NA

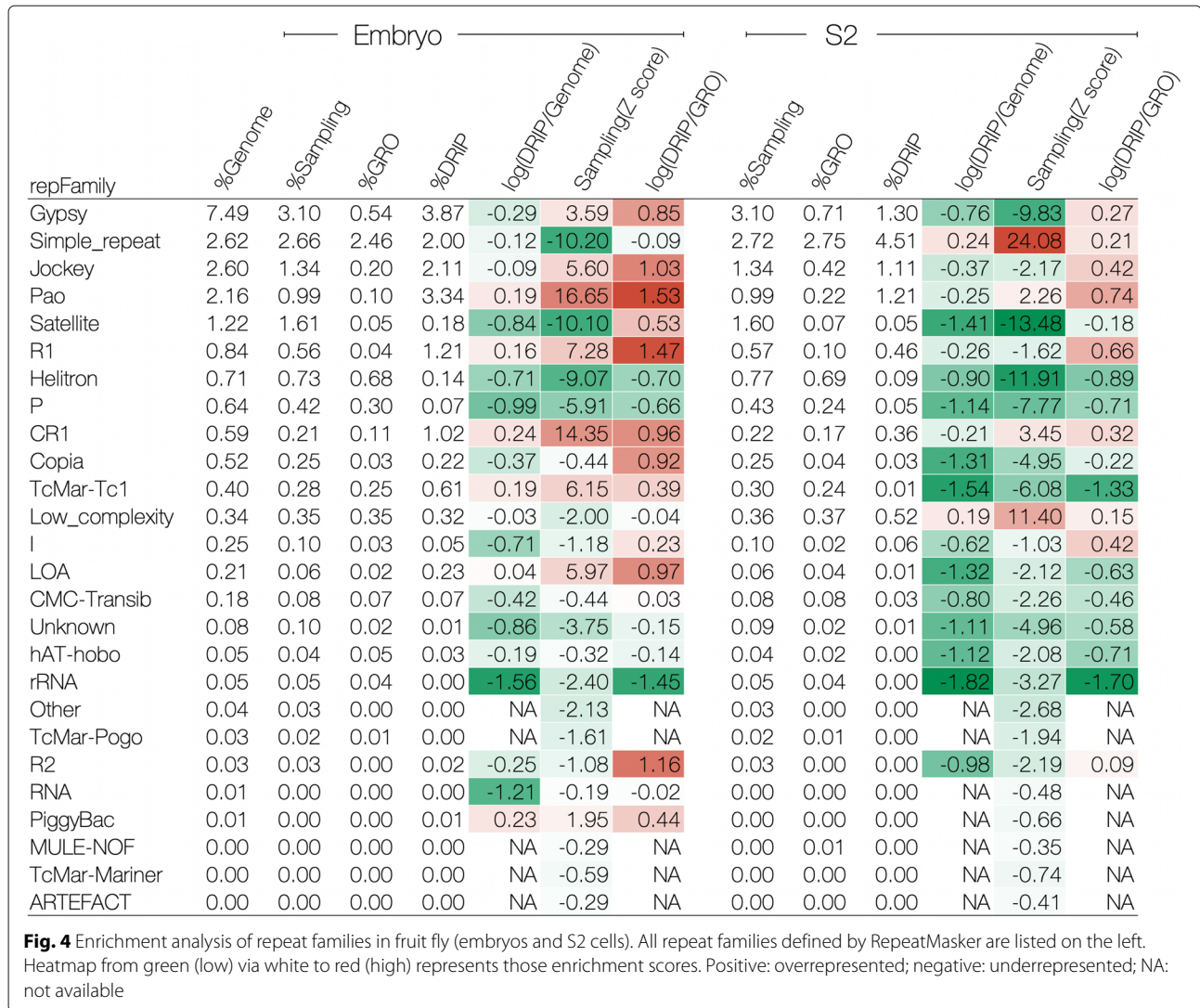
Fig. 3 Enrichment analysis of repeat families in human U2OS cells. All repeat families defined by RepeatMasker are listed on the left. Heatmap from green (low) via white to red (high) represents those enrichment scores. Positive: overrepresented; negative: underrepresented; NA: not available

For *A. thaliana*, nearly half of the repeat families were enriched in R-loops in both the sampling and GRO controls (Fig. 5). These repeat families were Gypsy, Copia, LTR, all belonging to the LTRs; MULE-MuDR, CMC-EnSpm, En-Spm, DNA, which are DNA transposons; L1 belonging to the LINES; centr; simple repeat and low-complexity. However, PIF-Harbinger, all belonging to DNA, satellite, composite, and tRNA, were only observed to be enriched in R-loops in the GRO control.

Repetitive elements differently contribute to r-loop formation and function

We found that repetitive elements contribute differently to R-loop formation among the samples investigated

in this study. Human U2OS cells showed that 21.19% of the DRIP-seq signals contained repetitive elements, while 43.19% of the GRO-seq data contained these signatures (Additional file 2A). Further analysis revealed that some repetitive sequences, especially TEs, including LINES, SINES, LTRs, and DNA families, did not tend to form R-loops in U2OS cells. On the other hand, low-complexity, satellite, simple repeat, retroposon, snRNA, and rRNA sequences were enriched in R-loop regions compared to non-repetitive sequences in the GRO-seq control. These results are consistent with previous reports [3, 24, 51, 53, 54]. Notably, a recent report has shown that low-complexity and simple repeat sequences are strongly associated with promoter regions [55], as are R-loop



structures [3, 24, 51]. These results suggest that repetitive elements, such as low-complexity and simple repeats, are the key features of R-loop formation in promoter regions. Interestingly, low-complexity sequences have also been shown to be associated with Ezh2 binding, which is a component of polycomb repressive complex 2 (PRC2), and have methyltransferase activity for histone H3 lysine 27 [55]. Another report has shown that R-loop formation is required for the recruitment of PRC2 and repression of a subset of polycomb target genes [11]. These results suggest that R-loop formation involving low-complexity elements could be important for the recruitment of PRC2 and epigenetic regulation of target genes. Therefore, we hypothesize that repetitive elements in R-loop regions might contribute differently to the subsequent function of R-loop formation.

In contrast to human U2OS cells, *A. thaliana* seedlings showed that 22.25% of the DRIP-seq signals contained

repetitive elements, while only 3.08% of the GRO-seq data contained these elements (Additional file 2C). In addition to simple repeats, low-complexity, and satellites, which are prone to form R-loops in human U2OS cells, TEs, including LTRs, DNA transposons, and LINES, were more preferentially enriched in R-loop regions in *A. thaliana* seedlings. These results imply that R-loop formation does not simply depend on genomic sequence features but depends highly on the species (or biological contexts). Given that R-loop formation is essential for epigenetic regulation [3, 24, 51], TEs that form R-loops could be critical regulatory elements for gene regulation in *A. thaliana* seedlings. Further analysis of such factors will reveal the functional significance of R-loop formation in TEs.

To investigate the contribution of repetitive elements in R-loop formation at different developmental stages, we compared the distribution of repetitive elements in R-loop regions between fly embryos and S2 cell lines.

repFamily	%Genome	%Sampling	%GRO	%DRIP	log(DRIP/Genome)	Sampling(Z_score)	log(DRIP/GRO)
Gypsy	5.21	2.98	0.18	9.44	0.26	25.70	1.72
MULE-MuDR	2.31	2.47	0.24	3.39	0.17	7.23	1.14
Helitron	1.78	1.97	0.32	0.33	-0.73	-15.79	0.01
Copia	1.43	2.05	0.25	2.43	0.23	3.52	0.99
Simple_repeat	1.18	1.23	1.08	1.60	0.13	17.53	0.17
L1	1.03	1.46	0.21	1.50	0.16	0.56	0.85
CMC-EnSpm	0.76	0.66	0.01	1.23	0.21	7.59	2.08
Satellite	0.47	0.26	0.02	0.09	-0.69	-3.28	0.58
centr	0.40	0.14	0.00	0.19	-0.32	1.34	2.61
Low_complexity	0.37	0.36	0.33	0.55	0.18	17.11	0.22
MuDR	0.27	0.29	0.03	0.52	0.28	5.24	1.29
DNA	0.27	0.21	0.07	0.22	-0.09	0.38	0.46
hAT-Ac	0.24	0.28	0.09	0.19	-0.10	-2.99	0.32
hAT	0.14	0.16	0.03	0.04	-0.57	-5.76	0.05
PIF-Harbinger	0.13	0.16	0.04	0.08	-0.20	-3.30	0.28
TcMar-Pogo	0.12	0.13	0.05	0.04	-0.48	-6.37	-0.09
En-Spm	0.09	0.09	0.00	0.13	0.16	1.66	1.52
LTR	0.06	0.05	0.00	0.13	0.30	8.55	1.41
TcMar-Stowaway	0.06	0.07	0.02	0.02	-0.45	-6.50	-0.09
SINE	0.05	0.06	0.04	0.04	-0.05	-1.63	0.06
Composite	0.05	0.05	0.01	0.04	-0.13	-1.55	0.53
tRNA	0.04	0.05	0.02	0.03	-0.08	-2.53	0.35
hAT-Tag1	0.03	0.03	0.01	0.01	-0.41	-2.25	-0.05
rRNA	0.01	0.00	0.01	0.01	-0.09	1.97	-0.03
Harbinger	0.01	0.01	0.00	0.00	-0.43	-1.08	0.02
TcMar-Mariner	0.01	0.01	0.00	0.00	-0.67	-2.27	-0.31
centromeric	0.00	0.00	0.00	0.00	NA	-0.49	NA

Fig. 5 Enrichment analysis of repeat families in *A. thaliana* seedlings. All repeat families defined by RepeatMasker are listed on the left. Heatmap from green (low) via white to red (high) represents those enrichment scores. Positive: overrepresented; negative: underrepresented; NA: not available

In fly embryos, 15.5% of the DRIP-seq signals contained repetitive elements, as compared to only 5.35% of the GRO-seq data (Additional file 2B). In S2 cells, 9.81% of the DRIP-seq signals contained some repetitive elements, while 6.28% of the GRO-seq data contained those elements (Additional file 2B). These results show that repetitive element contribution to R-loop formation is more prominent in embryos than in S2 cells, suggesting that the impact of repetitive elements on R-loop formation remarkably changes in different developmental stages or cell lineages. We also observed that LTRs, LINEs, and satellites were highly enriched in embryo R-loops and were less enriched in S2 R-loops. Conversely, simple repeats and low complexity were relatively enriched in S2 cells and less enriched in embryos. We speculate that repetitive elements could change their function through R-loop formation, along with the developmental context. For example, gypsy, which is known as one of the major insulator elements in flies [56], is more highly enriched in embryo R-loops than S2 R-loops. R-loop formation on gypsy may alter the function of the insulator or protein complex on insulator bodies,

resulting in the downstream regulation of the chromatin compartment. This case is consistent with the recent observation that R-loop formation is associated with an enhancer- and insulator-like state [3]. Further investigation is required to reveal the relationship between R-loop formation and the insulator function of gypsy elements.

R-loop formation might be derived from TE regulation

Our results highlight the impact of TE elements on R-loop formation, especially at different developmental stages. This suggests that the TE sequence itself could tend to form an R-loop. Because TEs originate from exogenous viruses, they are the target of gene silencing by multiple layers of defense mechanisms to prevent the harmful effects of TE activity. Therefore, R-loop formation involving TEs might be one such mechanism by which cells mitigate the effects of TEs. It has been shown that R-loop formation can stimulate transcription of an antisense sequence, resulting in the formation of heterochromatin [57, 58]. This mechanism is suitable if R-loop formation has a role in silencing TE elements. Similarly, it is

reasonable that R-loops have a role in regulatory signals of epigenetic regulations if their functional origin is derived from TE regulation. Moreover, chromatin loosening following the depletion of histone H1 induces the accumulation of R-loops in heterochromatic regions enriched with repetitive elements, including several types of TEs [59]. This result suggests that TE elements could preferentially form R-loop structures, when their silencing by heterochromatin is resolved. This is consistent with the notion that transcribing TE sequences increase the likelihood of R-loop formation. Taken together, R-loop formation might be intimately correlated with TE sequences, although further experimental studies are required to confirm this hypothesis.

Concluding remarks

In this study, we reanalyzed DRIP-seq data to investigate the impact of repetitive elements on R-loop formation. We found that satellites, LINEs, and DNA transposons were enriched for R-loops in humans, fruit flies, and *A. thaliana*, respectively. Consistently, we observed that R-loops preferred to form in regions of low-complexity or simple repeats across species. Additionally, we also found that the repetitive elements associated with R-loop formation differ according to the developmental stage. LINEs and LTRs are more likely to promote R-loop formation in embryos (fruit fly). However, R-loop formation changes in S2 cells to being more prevalent in low complexity and simple repeat areas of the genome. These results imply that repetitive elements may have species-specific or development-specific regulatory effects on R-loop formation. To our knowledge, this is the first study to analyze the association between repetitive elements and R-loop formation across species and developmental stages. Our results show that various repetitive elements may distinctly contribute to R-loop formation in a biological context-dependent manner. This work advances our understanding of repetitive elements and R-loop biology; future research should aim to determine the mechanism of R-loop formation on each repetitive element and its biological function.

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13100-021-00231-5>.

Additional file 1: Genomic distribution of DRIP- and GRO-seq peaks. Percentages of peaks overlapping with gene, promoter2k and intergenic regions shown in the last three columns. Peaks overlapping by at least 1nt with other regions are counted according to the priority (promoter2k > gene > intergenic). Promoter2k means 2000nt upstream of a gene region.

Additional file 2: Enrichment analyses of repeat classes across (A)human, (B)fruit Fly, and (C)*A. thaliana*.

Additional file 3: Public DRIP- and GRO-seq datasets used in this study.

Abbreviations

LINE: long interspersed nuclear element, LTR: long terminal repeat retrotransposon, ssDNA: single-stranded DNA, TERRA: telomere repeat-containing RNA, DRIP-seq: DNA–RNA immunoprecipitation followed by high-throughput DNA sequencing, G4: G quadruplex, TNR: trinucleotide repeat, TE: transposable element, PCR: polymerase chain reaction, IP: immunoprecipitation, GRO-seq: global run-on sequencing, RC: rolling circle, SINE: short interspersed nuclear element, PRC2: polycomb repressive complex 2

Acknowledgements

We thank Dr. Martin Frith and Dr. Yutaka Saito for their helpful discussions. We thank Ms. Risa Maemura for her participation in the initial survey of R-loop-related studies. Computations were partially performed on the NIG supercomputer at ROIS National Institute of Genetics.

Authors' contributions

MH conceived and supervised this study. CZ designed and performed the experiments. All authors contributed to analysis and interpretation of the data. CZ and MO wrote the manuscript. MH revised the manuscript critically. All authors read and approved the final manuscript.

Funding

This work was supported by JSPS KAKENHI [grant numbers JP17K20032, JP16H05879, and JP20H00624 to MH; JP20K15784 to CZ].

Availability of data and materials

DRIP-seq and GRO-seq datasets are available from public repositories. See Additional file 3 for the full list of the datasets.

Ethics approval and consent to participate

None declared.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Author details

¹AIST-Waseda University Computational Bio Big-Data Open Innovation Laboratory (CBBD-OIL), National Institute of Advanced Industrial Science and Technology, 63-520, 3-4-1, Okubo Shinjuku-ku, 169-8555 Tokyo, Japan. ²Faculty of Science and Engineering, Waseda University, 55N-06-10, 3-4-1 Okubo Shinjuku-ku, 169-8555 Tokyo, Japan. ³Institute for Medical-oriented Structural Biology, Waseda University, 2-2, Wakamatsu-cho Shinjuku-ku, 162-8480 Tokyo, Japan. ⁴Graduate School of Medicine, Nippon Medical School, 1-1-5, Sendagi, Bunkyo-ku, 113-8602 Tokyo, Japan.

Received: 10 November 2020 Accepted: 8 January 2021

Published online: 20 January 2021

References

1. Wahba L, Costantino L, Tan FJ, Zimmer A, Koshland D. S1-drip-seq identifies high expression and polyA tracts as major contributors to r-loop formation. *Gene Dev.* 2016;30(11):1327–38.
2. Xu W, Xu H, Li K, Fan Y, Liu Y, Yang X, Sun Q. The r-loop is a common chromatin feature of the arabidopsis genome. *Nat plants.* 2017;3(9):704–14.
3. Sanz LA, Hartono SR, Lim YW, Steyaert S, Rajpurkar A, Ginno PA, Xu X, Chédin F. Prevalent, dynamic, and conserved r-loop structures associate with specific epigenomic signatures in mammals. *Mol cell.* 2016;63(1):167–78.
4. Aguilera A, García-Muse T. R loops: from transcription byproducts to threats to genome stability. *Mol cell.* 2012;46(2):115–24.
5. Sollier J, Cimprich KA. Breaking bad: R-loops and genome integrity. *Trends Cell Biol.* 2015;25(9):514–22.
6. Chédin F. Nascent connections: R-loops and chromatin patterning. *Trends Genet.* 2016;32(12):828–38.
7. García-Muse T, Aguilera A. R loops: from physiological to pathological roles. *Cell.* 2019;179(3):604–18.

8. Crossley MP, Bocek M, Cimprich KA. R-loops as cellular regulators and genomic threats. *Mol cell*. 2019;73(3):398–411.
9. Niehrs C, Luke B. Regulatory r-loops as facilitators of gene expression and genome stability. *Nat Rev Mol Cell Biol*. 2020;21(3):167–78.
10. Grunseich C, Wang IX, Watts JA, Burdick JT, Guber RD, Zhu Z, Bruzel A, Lanman T, Chen K, Schindler AB, et al. Senataxin mutation reveals how r-loops promote transcription by blocking dna methylation at gene promoters. *Mol cell*. 2018;69(3):426–37.
11. Skourti-Stathaki K, Triglia ET, Warburton M, Voigt P, Bird A, Pombo A. R-loops enhance polycomb repression at a subset of developmental regulator genes. *Mol cell*. 2019;73(5):930–45.
12. Abraham KJ, Khosraviyani N, Chan JNY, Gorthi A, Samman A, Zhao DY, Wang M, Bokros M, Vidya E, Ostrowski LA, Oshidari R, Pietrobon V, Patel PS, Algouneh A, Singhania R, Liu Y, Yerlici VT, De Carvalho DD, Ohh M, Dickson BC, Hakem R, Greenblatt JF, Lee S, Bishop AJR, Mekhail K. Nucleolar rna polymerase ii drives ribosome biogenesis. *Nat*. 2020;585(7824):298–302.
13. Proudfoot NJ. Transcriptional termination in mammals: Stopping the rna polymerase ii juggernaut. *Science*. 2016;352(6291):aad9926.
14. Skourti-Stathaki K, Proudfoot NJ, Gromak N. Human senataxin resolves rna/dna hybrids formed at transcriptional pause sites to promote xrn2-dependent termination. *Mol cell*. 2011;42(6):794–805.
15. Castellano-Pozo M, Santos-Pereira JM, Rondón AG, Barroso S, Andújar E, Pérez-Alegre M, García-Muse T, Aguilera A. R loops are linked to histone h3 s10 phosphorylation and chromatin condensation. *Mol cell*. 2013;52(4):583–90.
16. Graf M, Bonetti D, Lockhart A, Serhal K, Kellner V, Maicher A, Jolivet P, Teixeira MT, Luke B. Telomere length determines terra and r-loop regulation through the cell cycle. *Cell*. 2017;170(1):72–85.
17. Kabeche L, Nguyen HD, Buisson R, Zou L. A mitosis-specific and r loop-driven atr pathway promotes faithful chromosome segregation. *Science*. 2018;359(6371):108–14.
18. Skourti-Stathaki K, Proudfoot NJ. A double-edged sword: R loops as threats to genome integrity and powerful regulators of gene expression. *Gene Dev*. 2014;28(13):1384–96.
19. Groh M, Gromak N. Out of balance: R-loops in human disease. *PLoS Genet*. 2014;10(9):1004630.
20. Costantino L, Koshland D. The yin and yang of r-loop biology. *Curr Opin Cell Biol*. 2015;34:39–45.
21. Richard P, Manley JL. R loops and links to human disease. *J Mol Biol*. 2017;429(21):3168–80.
22. Khristich AN, Mirkin SM. On the wrong dna track: Molecular mechanisms of repeat-mediated genome instability. *J Biol Chem*. 2020;295(13):4134–70.
23. Ginno PA, Lott PL, Christensen HC, Korf I, Chédin F. R-loop formation is a distinctive characteristic of unmethylated human cpG island promoters. *Mol cell*. 2012;45(6):814–25.
24. Chen L, Chen J-Y, Zhang X, Gu Y, Xiao R, Shao C, Tang P, Qian H, Luo D, Li H, et al. R-chip using inactive rna polymerase ii reveals dynamic coupling of r-loops with transcriptional pausing at gene promoters. *Mol cell*. 2017;68(4):745–57.
25. Dumelie JG, Jaffrey SR. Defining the location of promoter-associated r-loops at near-nucleotide resolution using bisdrp-seq. *Elife*. 2017;6:28306.
26. Tan-Wong SM, Dhir S, Proudfoot NJ. R-loops promote antisense transcription across the mammalian genome. *Mol cell*. 2019;76(4):600–16.
27. Yan Q, Shields EJ, Bonasio R, Sarma K. Mapping native r-loops genome-wide using a targeted nuclease approach. *Cell Rep*. 2019;29(5):1369–80.
28. Malig M, Hartono SR, Giapfagnione JM, Sanz LA, Chedin F. Ultra-deep coverage single-molecule r-loop footprinting reveals principles of r-loop formation. *J Mol Biol*. 2020;432(7):2271–88.
29. Crossley MP, Bocek MJ, Hamperl S, Swigut T, Cimprich KA. qdrp: a method to quantitatively assess rna-dna hybrid formation genome-wide. *Nucleic Acids Res*. 2020;48(14):e84.
30. Ginno PA, Lim YW, Lott PL, Korf I, Chédin F. Gc skew at the 5' and 3' ends of human genes links r-loop formation to epigenetic regulation and transcription termination. *Genome Res*. 2013;23(10):1590–600.
31. De Magis A, Manzo SG, Russo M, Marinello J, Morigi R, Sordet O, Capranico G. Dna damage and genome instability by g-quadruplex ligands are mediated by r loops in human cancer cells. *Proc Natl Acad Sci*. 2019;116(3):816–25.
32. Lin Y, Dent SY, Wilson JH, Wells RD, Napierala M. R loops stimulate genetic instability of ctg·cag repeats. *Proc Natl Acad Sci*. 2010;107(2):692–7.
33. Groh M, Lufino MM, Wade-Martins R, Gromak N. R-loops associated with triplet repeat expansions promote gene silencing in friedreich ataxia and fragile x syndrome. *PLoS Genet*. 2014;10(5):1004318.
34. Loomis EW, Sanz LA, Chédin F, Hagerman PJ. Transcription-associated r-loop formation across the human fmr1 cgg-repeat region. *PLoS Genet*. 2014;10(4):1004294.
35. Su XA, Freudenreich CH. Cytosine deamination and base excision repair cause r-loop-induced cag repeat fragility and instability in *saccharomyces cerevisiae*. *Proc Natl Acad Sci*. 2017;114(40):8392–401.
36. Wongsurawat T, Jenjaroenpun P, Kwok CK, Kuznetsov V. Quantitative model of r-loop forming structures reveals a novel level of rna-dna interactome complexity. *Nucleic Acids Res*. 2012;40(2):16.
37. El Hage A, Webb S, Kerr A, Tollervey D. Genome-wide distribution of rna-dna hybrids identifies rna polymerase h targets in trna genes, retrotransposons and mitochondria. *PLoS Genet*. 2014;10(10):1004716.
38. Nadel J, Athanasiadou R, Lemetre C, Wijetunga NA, Broin PO, Sato H, Zhang Z, Jeddeloh J, Montagna C, Golden A, et al. Rna: Dna hybrids in the human genome have distinctive nucleotide characteristics, chromatin composition, and transcriptional relationships. *Epigenetics chromatin*. 2015;8(1):1–19.
39. Toriumi K, Tsukahara T, Hanai R. R-loop formation in trans at an aggag repeat. *J Nucleic Acids*. 2013;2013:629218.
40. Chu H-P, Cifuentes-Rojas C, Kesner B, Aeby E, Lee H-g, Wei C, Oh HJ, Boukhali M, Haas W, Lee JT. Terra rna antagonizes atrx and protects telomeres. *Cell*. 2017;170(1):86–101.
41. Zhou Z, Giles KE, Felsenfeld G. Dna·rna triple helix formation can function as a cis-acting regulatory mechanism at the human β -globin locus. *Proc Natl Acad Sci*. 2019;116(13):6130–9.
42. Ariel F, Lucero L, Christ A, Mammarella MF, Jegu T, Veluchamy A, Mariappan K, Latrasse D, Blein T, Liu C, et al. R-loop mediated trans action of the apolo long noncoding rna. *Mol cell*. 2020;77(5):1055–65.
43. Feretzaki M, Pospisilova M, Valador Fernandes R, Lunardi T, Krejci L, Lingner J. RAD51-dependent recruitment of TERRA lncRNA to telomeres through R-loops. *Nat*. 2020;587(7833):303–8.
44. Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, Haussler D. The human genome browser at ucsc. *Genome Res*. 2002;12(6):996–1006.
45. Yates AD, Achuthan P, Akanni W, Allen J, Allen J, Alvarez-Jarreta J, Amode MR, Armean IM, Azov AG, Bennett R, et al. Ensembl 2020. *Nucleic Acids Res*. 2020;48(D1):682–8.
46. Quinlan AR, Hall IM. Bedtools: a flexible suite of utilities for comparing genomic features. *Bioinforma*. 2010;26(6):841–2.
47. Li H, Durbin R. Fast and accurate short read alignment with burrows-wheeler transform. *Bioinforma*. 2009;25(14):1754–60.
48. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R. The sequence alignment/map format and samtools. *Bioinforma*. 2009;25(16):2078–9.
49. Picard toolkit. <https://doi.org/broadinstitute.github.io/picard>. Accessed 14 Oct 2020.
50. Wang D, Garcia-Bassets I, Benner C, Li W, Su X, Zhou Y, Qiu J, Liu W, Kaikkonen MU, Ohgi KA, et al. Reprogramming transcription by distinct classes of enhancers functionally defined by ena. *Nat*. 2011;474(7351):390–4.
51. Yan Q, Shields EJ, Bonasio R, Sarma K. Mapping native r-loops genome-wide using a targeted nuclease approach. *Cell Rep*. 2019;29(5):1369–80.
52. Nadel J, Athanasiadou R, Lemetre C, Wijetunga NA, Ó Broin P, Sato H, Zhang Z, Jeddeloh J, Montagna C, Golden A, Seoighe C, Grealley JM. RNA:DNA hybrids in the human genome have distinctive nucleotide characteristics, chromatin composition, and transcriptional relationships. *Epigenetics chromatin*. 2015;8:46.
53. Johnson WL, Yewdell WT, Bell JC, McNulty SM, Duda Z, O'Neill RJ, Sullivan BA, Straight AF. RNA-dependent stabilization of SUV39H1 at constitutive heterochromatin. *Elife*. 2017;6:e25299.
54. Velazquez Camacho O, Galan C, Swist-Rosowska K, Ching R, Gamalinda M, Karabiber F, De La Rosa-Velazquez I, Engist B, Koschorz B, Shukeir N, Onishi-Seebacher M, van de Nobelen S, Jenuwein T. Major satellite repeat RNA stabilize heterochromatin retention of Suv39h enzymes by

RNA-nucleosome association and RNA:DNA hybrid formation. *Elife*. 2017;6:e25293.

55. Lu JY, Shao W, Chang L, Yin Y, Li T, Zhang H, Hong Y, Percharde M, Guo L, Wu Z, Liu L, Liu W, Yan P, Ramalho-Santos M, Sun Y, Shen X. Genomic repeats categorize genes with distinct functions for orchestrated regulation. *Cell Rep*. 2020;30(10):3296–3115.
56. Geyer PK, Corces VG. DNA position-specific repression of transcription by a *Drosophila* zinc finger protein. *Genes Dev*. 1992;6(10):1865–73.
57. Skourti-Stathaki K, Kamieniarz-Gdula K, Proudfoot NJ. R-loops induce repressive chromatin marks over mammalian gene terminators. *Nat*. 2014;516(7531):436–9.
58. Nakama M, Kawakami K, Kajitani T, Urano T, Murakami Y. DNA-RNA hybrid formation mediates RNAi-directed heterochromatin formation. *Genes Cells*. 2012;17(3):218–33.
59. Bayona-Feliu A, Casas-Lamesa A, Reina O, Bernués J, Azorín F. Linker histone H1 prevents R-loop accumulation and genome instability in heterochromatin. *Nat Commun*. 2017;8(1):283.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

