

1 **New, easy, quick and efficient DNA replication timing analysis by high-**
2 **throughput approaches.**

3 Jihad Hadjadj^{1¶}, Thomas Denecker^{2¶}, Eva Guérin¹, Su-Jung Kim¹, Fabien Fauchereau¹,
4 Giuseppe Baldacci¹, Chrystelle Maric^{1\$} and Jean-Charles Cadoret^{1\$*}

5

6 ¹Pathologies de la Réplication de l'ADN, Institut Jacques-Monod, UMR7592, CNRS,
7 Université de Paris, F-75013, Paris, France

8 ²Institut de Biologie Intégrative de la Cellule UMR9198, CNRS, Université Paris-Saclay,
9 Université Paris-Sud, F-91405, Orsay, France

10 [¶]These authors contributed equally to this work.

11 ^{\$} co-last authors

12 * corresponding author

13 E-mail: jean-charles.cadoret@ijm.fr

14

15

16 **Running Title:** Easy replication timing analysis with START-R

17

18

19 **Keywords:** Replication Timing, Bioinformatic method, genomic organization

1 **Abstract**

2 DNA replication must be faithful and follow a well-defined spatio-temporal program closely
3 linked to transcriptional activity, epigenomic marks, intra-nuclear structures, mutation rate
4 and cell fate determination. Among the readouts of the DNA replication spatio-temporal
5 program, replication timing (RT) analyses require complex, precise and time-consuming
6 experimental procedures, and the study of large-size computer files. We improved the RT
7 protocol to speed it up and increase its quality and reproducibility. Also, we partly automated
8 the RT protocol and developed a user-friendly software: the START-R suite (Simple Tool for
9 the Analysis of the Replication Timing based on R). START-R suite is an open source web
10 application using an R script and an HTML interface to analyze DNA replication timing in a
11 given cell line with microarray or deep-sequencing results. This novel approach can be used
12 by every biologist without requiring specific knowledge in bioinformatics. It also reduces the
13 time required for generating and analyzing simultaneously data from several samples.
14 START-R suite detects constant timing regions (CTR) but also, and this is a novelty, it
15 identifies temporal transition regions (TTR) and detects significant differences between two
16 experimental conditions. The informatic global analysis requires less than 10 minutes.

1 **Introduction**

2 DNA replication is a highly regulated process involved in the maintenance of genome stability
3 (Hanahan and Weinberg, 2011; Macheret and Halazonetis, 2015; Técher et al., 2017). Its
4 accuracy relies partly on a spatio-temporal program that regulates timing and location of
5 origin firing (Dileep et al., 2015; Rivera-Mulia and Gilbert, 2016). Based on this program,
6 replication is organized into large-scale domains that replicate at different times in S phase
7 (Ryba et al., 2010; Cornacchia et al., 2012). During the last decade, different groups
8 including our laboratory, showed that the replication-timing program (RT) is finely tuned and
9 maintained from an S phase to the following one (Hadjadj et al., 2016; Brustel et al., 2017;
10 Almeida et al., 2018). It also appeared that this program is modified during cell differentiation
11 (Hiratani et al., 2010; Gilbert 2012; Hadjadj et al., 2016). However, it remains unclear how
12 this program is established and maintained. Protocols developed to study the RT in specific
13 cell lines have been established in different labs (Hansen et al., 2010; Ryba et al., 2011;
14 Dileep et al., 2012; Marchal et al., 2018). Differences between RT protocols may produce
15 different results, sometimes devoid of biological relevance. DNA-Immunoprecipitation (DNA-
16 IP) is a critical step of the RT protocol. We optimized duration and reproducibility of this step
17 by using the SX-8G IP-Star® Compact Automated System (Diagenode®). Thus, it now lasts
18 only 1 day (instead of 2-3 days before) and produces highly reproducible results regardless
19 of the experimenter. In order to make the analysis of experimental results more accurate and
20 reproducible, we also implemented two web-based softwares: START-R Analyzer and
21 START-R Viewer, showing user-friendly interfaces (HTML and simple-click controls) that can
22 be used by any biologist. The START-R Analyzer was initially based on a script developed in
23 2011 by David Gilbert's laboratory (Ryba et al., 2011) which is not anymore currently working
24 as it is, due to different software updates. The script was improved by implementing new
25 tools for the detection of temporal transition regions (TTR) and for the fast identification of
26 differential results between two experiments. These softwares are available online on our

1 GitHub group website (<https://github.com/thomasdenecker/START-R>), they are free and
2 each developer can improve them according to specific needs. We validated the START-R
3 suite with Drosophila, zebrafish, mouse and human RT data obtained with microarrays or
4 high-throughput sequencing. Using this automated DNA-IP protocol followed by analysis with
5 the START-R Suite, it becomes easier for a large number of laboratories to carry out studies
6 on the RT, thus opening up to new research perspectives.

1 **Results**

2 **An improved RT protocol using the IP-Star robot**

3 The protocol developed to analyze genome-wide replication-timing program (RT) in
4 mammalian cells lasts 3 weeks (Ryba et al., 2011). It includes pulse-labeling of cells with
5 nucleotide analog 5-bromo-2-deoxyuridine (BrdU) followed by flow cytometry cell sorting
6 (FACS) of labeled cells into two S-phase fractions. Then, immunoprecipitation (IP) targeting
7 the BrdU-labeled DNA is performed. IP is a time-consuming step and needs to be carefully
8 monitored in order to get precise and specific signals. Thus, we optimized the BrdU IP
9 protocol by introducing an automated step using the SX-8G IP-Star® Compact Automated
10 System. It is now possible to simultaneously perform IP of 16 samples that correspond to 8
11 early and 8 late fractions, which means 8 RT experiments. The automated system allows to
12 perform DNA-IP overnight. Its high standardization improves the reproducibility of RT profiles
13 from one experiment to another, whoever the experimenter. To prove this point, RT analyses
14 were performed by four different experimenters with four independent RKO cell line cultures;
15 two experimenters performed handmade DNA-IPs independently, while two others
16 independently used the IP-Star robot. Then, we compared the percentage of differences
17 between each handmade experiment and each “IP-Star” experiment. We found a difference
18 of 4.25% between both handmade experiments while only a difference of 0.11% was
19 observed between each independent experiment performed with the “IP-Star” robot
20 (Supplemental Fig. S2 and Supplemental Table S1).

21

22 Once immunoprecipitated, newly synthesized DNA is amplified with a SeqPlex™ enhanced
23 DNA Amplification kit designed for microarray or deep-sequencing experiments. For
24 microarray experiments, labeled DNA is hybridized to a whole genome comparative
25 hybridization microarray (CGH microarray, 180,000 probes, one every 13Kb). We showed
26 that a microarray with only 60,000 probes is not sufficient to produce a detailed RT profile

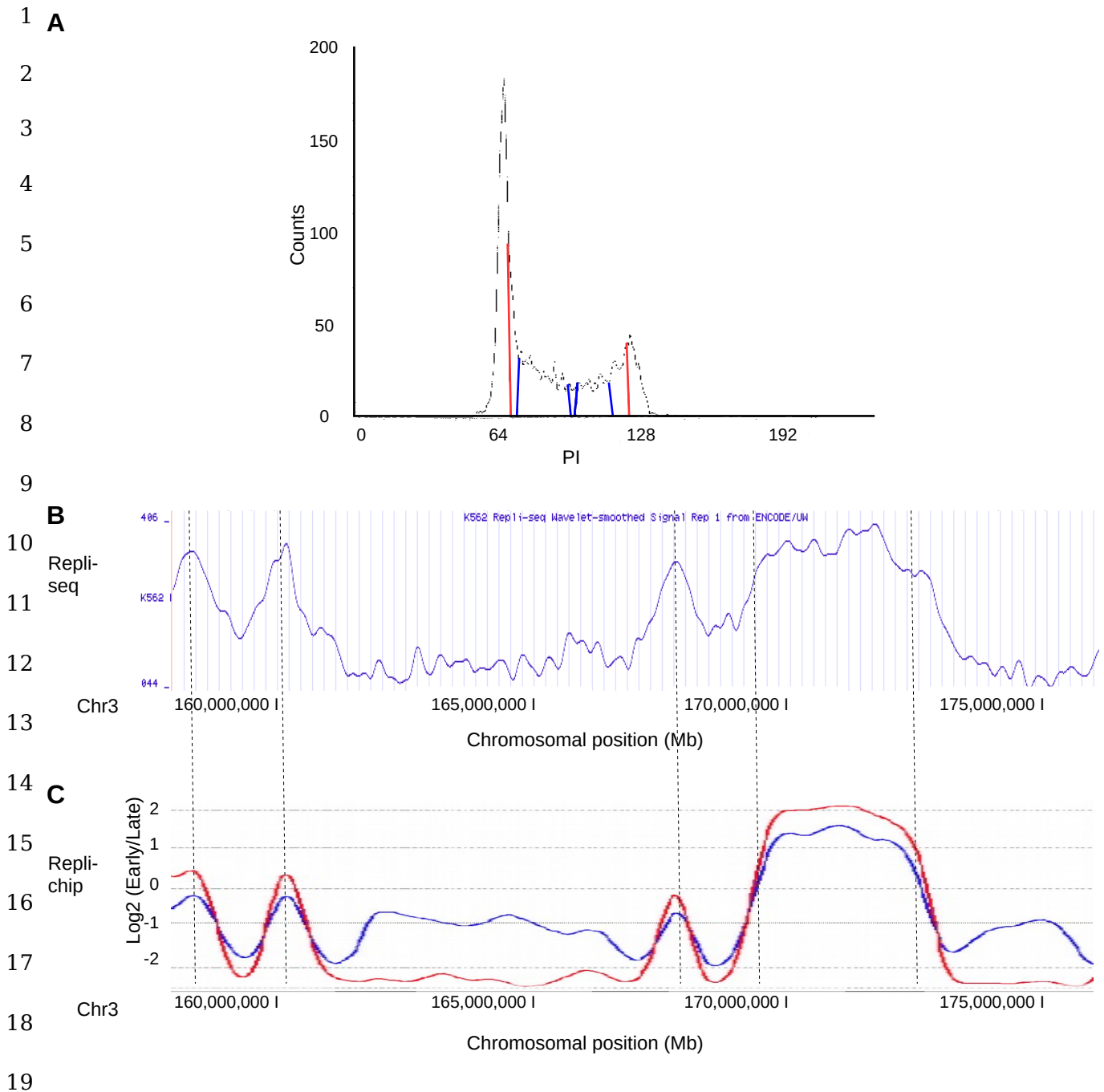
1 (Supplemental Fig. S3). After scanning, the generated picture is analyzed through the feature
2 extraction software (Feature Extraction 9.1 Agilent) that measures Cy3 and Cy5 intensity
3 values for each of the 180,000 probes of the microarray. This generates a table containing
4 the measures of processed signals that are used by START-R. “SystematicName”,
5 “gProcessedSignal and rProcessedSignal” columns are the ones used by default by START-
6 R Analyzer to generate the whole-genome RT profile and its analysis. Other column names
7 could also be used but experimenters have to be sure that such names are well matched
8 when the file is downloaded to START-R Analyzer.

9

10 **The Repli-chip protocol leads to results similar to the “6 fractions Repli-Seq” protocol**

11 To test the accuracy of our Repli-chip approach, we compared it with the previously used “6
12 fractions-Replication-sequencing (Repli-seq)” method (Hansen et al., 2010). Both
13 approaches were applied to the same K562 cell line. We retrieved these Repli-seq data (6
14 fractions corresponding to G1, S1, S2, S3, S4 and G2) from the UCSC genome browser
15 website (Kent et al., 2002).

16 Initially, we chose the position of cell sorting windows in S phase in our Repli-chip
17 experiment on the basis of previous validated RT protocols (blue lines, Fig. 1A). When
18 S1(Early) and S2 (Late) fractions are limited only to the S-Phase during Repli-chip
19 experiments (blue lines, Fig. 1A), some regions finally appeared as replicating in the middle
20 of S phase (blue smooth RT profile, Fig. 1C), whilst they had been shown to replicate very
21 late in the “6 fractions Repli-seq” experiment (Fig. 1B). This could be explained by the low
22 amounts of nascent DNA fragments present in the Early and Late sorted fractions restricted
23 to S phase. The analysis of this Early/Late ratio results in an artefactual mid S phase
24 replicating domain (blue smooth RT profile, Fig. 1C). Thus, these regions appear to replicate
25 in the middle of S phase while corresponding to domains replicated during the very late S-
26 phase and the beginning of G2.



20 **Figure 1. Repli-chip and Repli-seq protocols lead to similar results.**

21 **A)** Cell cycle profile of K562 cells after propidium iodide labeling. The blue lines indicate
22 windows used to sort cells in the first (S1) and the second part of S phase (S2). The red lines
23 indicate the new wider delimitations of Early and Late sorting windows. **B)** Replication profile
24 of K562 cell line obtained with Repli-seq approach. **C)** Replication profiles of K562 cell line
25 obtained with our Repli-chip protocol. The blue line depicts the replication profile obtained
26 after using the sorting windows limited to the S phase. The red line depicts replication
27 profiles obtained using wider sorting windows overlapping slightly with G1 (on the left) and
28 G2 (on the right) phases of cell cycle (as shown in A). The x-axis displays chromosomal
29 positions (Mb) and the y-axis log₂ (Early/Late) intensities. Dashed vertical lines show
30 common RT regions between the B and C profiles.

1 To fix this problem, we have expanded the cell sorting windows to the end of G1 for the Early
2 fraction, and into the beginning of G2 for the Late fractions (red lines, Fig. 1A). Using these
3 cell-sorting parameters RT profiles are highly similar to the “6 fractions-Repli-seq”
4 experiments despite the gap left between S1 and S2 fractions to avoid cross-contaminations
5 between both fractions, as shown for chromosome 19 and all other chromosomes (red
6 smooth RT profile, Figs. 1B, 1C and Supplemental Fig. S4, respectively). Therefore, an
7 accurate Repli-chip protocol with correct cell-sorting parameters for two fractions provides
8 similar results to those produced by the “6 fractions Repli-Seq” approach, in a less expensive
9 and less time-consuming way.

10

11 **START-R suite for automation of RT analysis allows robust statistical analysis with a** 12 **user-friendly interface**

13 We developed a software suite starting from a script created by David Gilbert's group in 2011
14 (Ryba et al., 2011), that we improved and updated with different current versions of tools and
15 algorithms. We also implemented new functions as TTR detection and differential analysis.
16 The START-R suite, which stands for Simple Tool for the Analysis of the Replication Timing
17 based on R, is implemented into an HTML interface for more efficient and easier installation,
18 use and sharing by biologists. START-R is built-in with Docker that packages START-R into
19 a virtual container (Supplemental Fig. S1). Thus, START-R can be easily deployed at a
20 personal computer or on a server, and can run independently of any library updating. This
21 was not the case in the script developed in 2011 (Ryba et al., 2011), which makes that script
22 much less easy to use (Supplemental Fig. S1). As indicated by its acronym, START-R has
23 the strength of being based on a statistical approach using R, allowing researchers non-
24 initiated in R programming to analyze their data. While the usability of the program would
25 tend to limit its adaptability, START-R provides as many parameters as possible for a
26 comprehensive analysis of the RT program (Supplemental Fig. S5A to 5K). Furthermore, we
27 added new scaling, normalization and smoothing methods (Supplemental Figs. S6, S7) and

1 also novel statistical approaches to detect differences between two samples (Supplemental
2 Fig. S8, S9). A classical differential analysis performed with START-R takes only 5-6 min
3 (compared to several hours without using START-R), with a personal computer containing an
4 intel® Xeon(R) CPU E5-1620-3.60GHz × 8 core and 32GB of memory with the 18.04 Ubuntu
5 version. START-R Analyzer runs by default with the hg18 human genomic annotations
6 including the position of the centromeres. It can also be used without the centromere
7 positions or with the centromere positions uploaded in a corresponding chromosome
8 coordinates file for all organisms and all annotations (>hg18 for human genome). This
9 flexibility is one of the new aspects of the START-R suite that allows to analyze RT program
10 in every organisms (Supplemental Fig. S5B).

11

12 **A large panel of new settings and tools for RT analysis**

13 In our method, we based our script on four major steps: normalisation (between Early and
14 Late fractions, between two replicates, and between two independent experiments,
15 Supplemental Fig. S6), smoothing (LOESS, Simple, Weighted, Modified, Triangular,
16 Exponential and Running methods including limma, Ritchie et al., 2015, Supplemental Fig.
17 S7), identification of transition timing regions (TTRs), and segmentation. The originality of our
18 approach is to first detect TTRs in order to better identify Constant Timing Regions (CTRs,
19 Supplemental Figs. S8A, S8B). The identification of TTRs is based on their intrinsic
20 properties: regions that include more than three consecutive probes with significantly
21 different Early/Late intensity log ratios are considered as TTRs (Supplemental Fig. S8A). The
22 statistical significance of differences between probes is calculated by the outlier boxplot
23 method (Supplemental Fig. S9, Krzywinski and Altman, 2013). Following TTRs detection,
24 START-R Analyzer localizes CTRs: TTRs are subtracted from the genome (Supplemental
25 Fig. S8B) and the remaining regions are considered as CTRs. If TTRs were not excluded
26 initially, the CTRs would overlap the adjacent TTRs (Supplemental Fig. S8B). This overlap

1 would result in a less precise segmentation because the algorithm would take into account
2 the probes positioned in adjacent TTRs for calculating the segment value. After subtraction
3 of TTRs, START-R Analyzer scans the remaining regions through a sliding window-based
4 algorithm. It will potentially divide these regions into different CTRs if the standard deviation
5 values of intensities reach a chosen threshold. At the end of these steps, START-R Analyzer
6 automatically generates a BED file for CTRs and TTRs making easier further bioinformatics
7 analyses and the display of the RT domains via a genome browser (Supplemental Fig. S5F).
8 It also produces descriptive statistics of the normalization step and RT statistical elements for
9 each chromosome (domain size distribution, summary of segmentation process). Finally, a
10 codebook is generated to ensure the traceability of options chosen for each analysis. The
11 user-friendly interface facilitates the choice among the different analysis parameters
12 (Supplemental Figs. S5A to S5K).

13 We added a step allowing the differential analysis of RT programs from two experiments.
14 Thus, we can now compare RT profiles obtained in different conditions and/or with different
15 cell lines to identify loci and elements that can modify the RT program. Our differential
16 analysis includes three different methods of comparison: the Mean method, the Euclidean
17 method and the Segment comparison method.

18 The Mean method compares the means of log ratio intensities (Early/Late) obtained in two
19 different experiments. Mean is calculated for a sliding window of 30 successive probes
20 corresponding to a 300 kb genomic domain, which is consistent with the size of replication-
21 timing domains already described (Rivera-Mulia and Gilbert, 2016). The overlapping
22 parameter, defining the number of probes overlapping successive windows was initially set to
23 15. After the calculation of nominal and adjusted p-values (t-tests for mean comparison), the
24 user can choose the p-value thresholds to distinguish significant differences between two
25 conditions. The p-value adjustment method is chosen among a list of classical procedures
26 (such as Bonferroni, Holm and Benjamini and Hochberg methods).

1 The Euclidean method computes the squared differences of the log ratio intensities
2 (Early/Late) for the same probes in two different experiments. The squared differences of
3 every probes are plotted in a boxplot and the outliers are considered as significantly different
4 following the threshold chosen by users (Supplemental Fig. S9). Regions of differential
5 intensities are defined as more than three consecutive outlier probes.

6 The segment approach uses the TTR and CTR values generated by START-R in the TTR
7 and segmentation steps. The user can choose any parameter as this approach employs an
8 empiric method. First, the method compares the presence or the absence of TTRs in two
9 experiments (Supplemental Fig. S8C). If TTRs are detected in both experiments, the
10 segment approach compares their slopes(Supplemental Fig. S8D). Finally, means of probe
11 intensities are compared between CTRs common to two different experiments, by t-test. The
12 threshold for the t-test significance is calculated according to the following procedure: two
13 empiric CTRs are generated by picking randomly intensity values in both CTR for each
14 experiment and a p-value is calculated by t-test. This permutation process is repeated
15 10.000 times and results to the estimation of an empiric p-value that is the threshold p-value
16 for CTR comparison. The last major implementation is START-R Viewer (Supplemental Fig.
17 S5K). This web-based interface allows the visualization of the RT profile generated by
18 START-R Analyzer in dynamic charts obtained with the Plotly library (Sievert et al., 2017). It
19 creates figures from a specific file generated by START-R Analyzer , integrating all the
20 analyses performed by START-R analyzer that another genome browser cannot display
21 optimally. One can easily identify CTRs, TTRs (Fig. 2A) and significantly advanced or
22 delayed regions (Fig. 2B). The user can also choose color options for the RT display and
23 take automatically a screenshot of the RT profile to create a figure. We therefore developed
24 a genome browser to optimally display the maximum of resources generated by START-R
25 Analyzer .
26

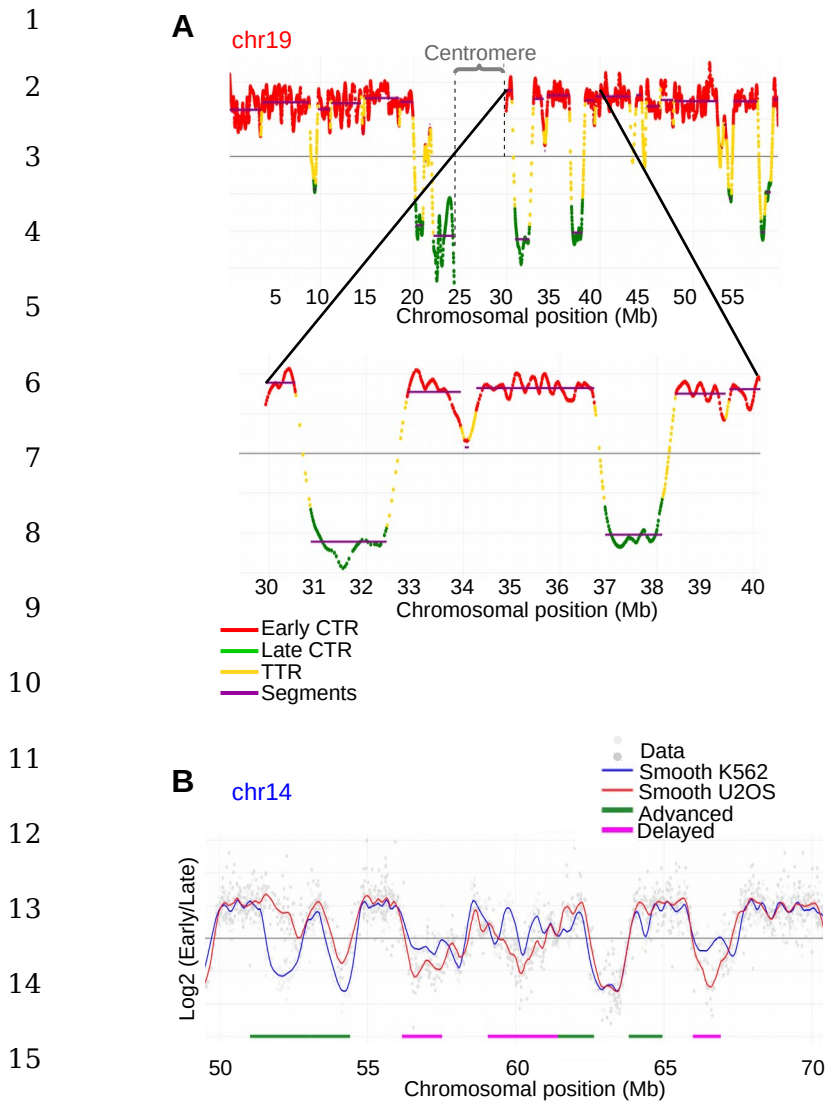
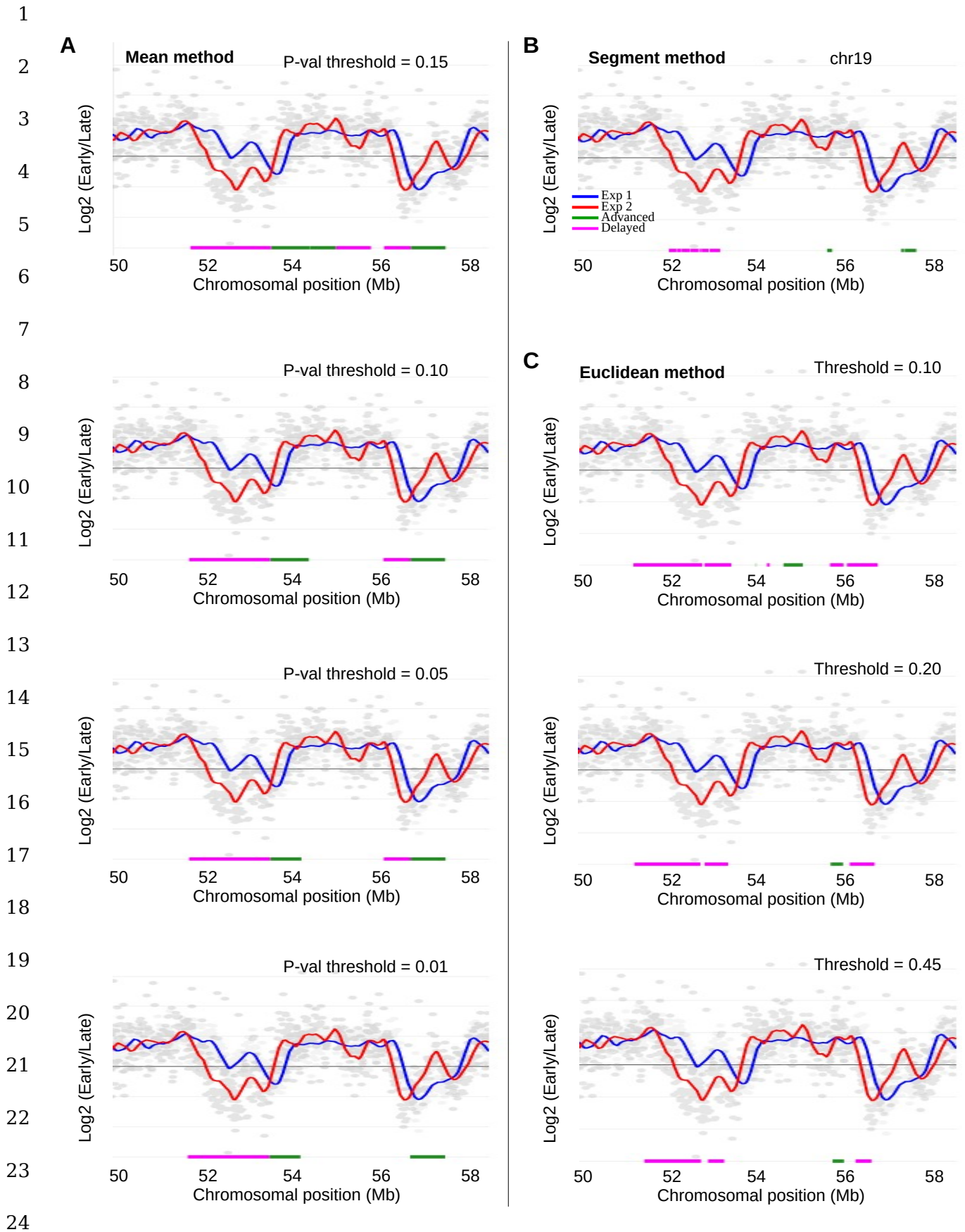


Figure 2. Examples of data generated with our RT protocol visualized by START-R-Viewer.

A) START-R Viewer allows visualizing RT data with many features. The top panel displays the distribution of early and late constant timing regions (CTR, in red and green, respectively) and of transition timing regions (TTR, in yellow) on a portion of human chromosome 19. Segments corresponding to regions of constant timing are shown in purple. Chromosome 19 centromere is indicated by grey dashed lines and a curly bracket. The bottom panel displays a zoom of a smaller region of chromosome 19 where timing profile can be seen through the zoom option of START-R-Viewer. B) Differential analyses are done on a portion of human chromosome 14 comparing RT profiles of two cell lines: K562 in blue and U2OS in red. Advanced (green) and Delayed (pink) regions are identified with START-R Analyzer using the mean comparison analysis with the Holm's p-value correction and a limit corrected p-value of 0.05. Light grey and grey spots indicate data from both RT experiments.

1 **How to choose adapted thresholds for differential analyses?**

2 As described above, START-R Analyzer proposes different parameters depending on the
3 chosen method. The central point for users is to select the correct method for the differential
4 analysis and to choose the adapted parameters for this process. To test the impact of
5 parameters on the detection of differences between two conditions, we compared RT data
6 from the K562 and U2OS cell lines. We used the same normalization, smoothing, TTR and
7 segment detection methods for both cell lines. We tested the Mean, Segment and Euclidean
8 methods for differential analyses (Fig. 3A, 3B and 3C). Only Mean and Euclidean methods
9 offer the possibility to change threshold parameters. The segment method does not use p-
10 value thresholds, or thresholds, since it is based on empirical parameters.

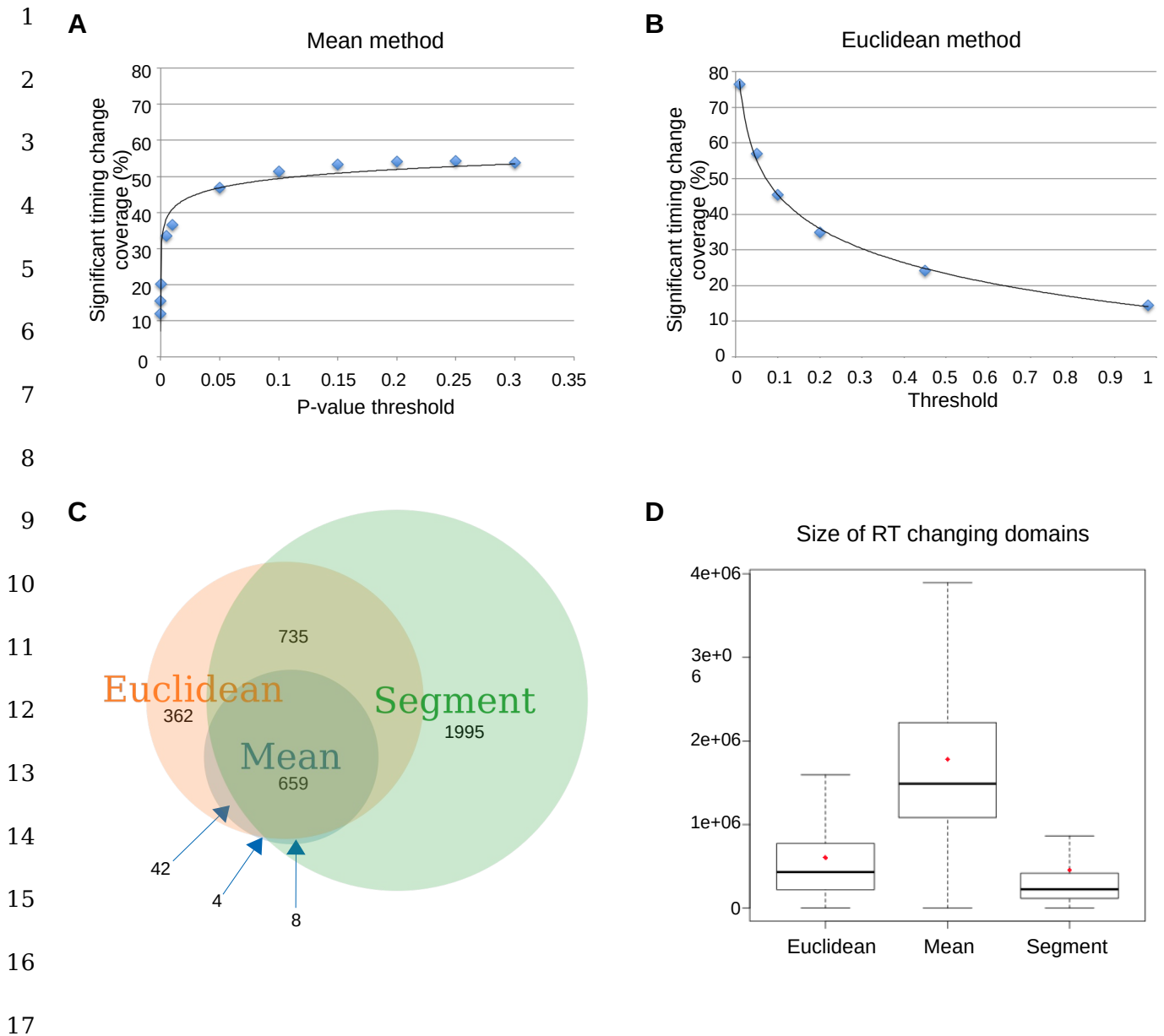


1 **Figure 3. Comparison of differential analyses of RT profiles with START-R Analyzer**
2 **using the Mean, Segment and Euclidean methods.**

3 **A)** Differential analyses allow the comparison of RT profiles of a portion of chromosome 19
4 for two cell lines, K562 (blue) and U2OS (red). Light grey and grey spots indicate data from
5 both RT experiments. The left panels show the identification of Advanced (green) and
6 Delayed (pink) regions using the Mean method with a corrected p-value threshold ranging
7 from 0.15 to more stringent p-values of 0.10, 0.05 and 0.01, respectively. **B)** Differential
8 analysis of the same chromosomal region using the Segment method based on empirical
9 parameters. **C)** Differential analyses of the same chromosomal region with the Euclidean
10 method with a threshold varying from 0.10 to 0.45.
11

12

13 To further explore the relationship between the p-value threshold or threshold and the
14 detection sensitivity of true RT changes with the Mean and Euclidean methods, we examined
15 the significant timing domain changes when the p-value threshold, or the threshold, increase
16 (Fig. 3A and 3C). We used graphs depicting the significant timing change coverage related
17 with the p-value threshold, or threshold (Figs. 4A and 4B and Supplemental Fig. S10). We
18 defined an optimized p-value threshold, or threshold, as the parameter at which the gain of
19 additional RT change coverage is minimal. For each method, we draw the chord of the curve
20 (Figs. 4A and 4B and Supplemental Fig. S10). The perpendicular and longer segment
21 between the chord and the curve was defined, indicating the optimized parameter
22 (Supplemental Fig. S10). These tests showed optimized p-value threshold and threshold for
23 the Mean and Euclidean methods of 0.025 and 0.192, respectively (Supplemental Fig. S10A
24 and S10B). Using these parameters, we found 713 RT changing regions with the Mean
25 method, 1798 with the Euclidean method, and 3397 with the Segment method (Fig. 4C).
26 Each method shows its particularities, however, 659 common RT changing regions were
27 detected by the three methods. While the number of regions with RT changes is different for
28 each method, the global genome coverage is the same. As START-R Analyzer works
29 quickly, we propose that experimenters use the same procedure in order to determine their
30 own optimized parameters suitable for their experiments.



1 Then, we compared the three methods with different combinations of parameters
2 (Supplemental Fig. S11) in order to evaluate the capacity of each method to detect
3 significant RT changes. As part of a comparison between the K562 and U2OS cell lines, the
4 Segment method allows the detection of the highest number of RT changes (3397 regions
5 with an average size of \approx 450Mb). The Mean Method shows the lowest RT changes (around
6 710 regions with an average size of \approx 1,800Mb), representing around 18-21% of RT changes
7 detected by the Segment method (Supplemental Fig. S11). The Euclidean method reveals
8 an intermediate number of RT changes (around 1932 regions with an average size of \approx 607
9 Mb). Thus, each method shows a specific range concerning the length of regions with RT
10 changes (Fig. 4D).

11 Nevertheless, 92 to 95% RT changes detected by the Mean method overlap those detected
12 with the other two methods (Supplemental Fig. S11). 40-48% RT changes detected by the
13 Euclidean method overlap those detected by the Segment method, and 29-40% overlap
14 those detected by the Mean method (Supplemental Fig. S11). With the optimized
15 parameters, only 0.5% of RT changes found by the Mean method are unique, compared to
16 20% for the Euclidean method and 58.7% for the Segment method (Fig. 4C). These
17 observations show that a large part of RT changes, but not all, are detected by the Mean
18 method. The Segment method appears to be more sensitive and able to detect more new
19 regions than the other two methods. Each method has its own detection characteristics. Each
20 user, depending on the asked biological questions, has to choose the most suitable one for
21 the analysis. When the goal is to identify most regions with a real RT change, we
22 recommend using the Mean method. When the goal is to find all regions with a RT change,
23 we recommend using Segment or Euclidean method but with increased risks of obtaining
24 false positives. Since the analyses with START-R are fast, we also recommend performing
25 the analysis using the three methods and keeping the common results. Therefore, it is
26 important that START-R Analyzer proposes these three tools.

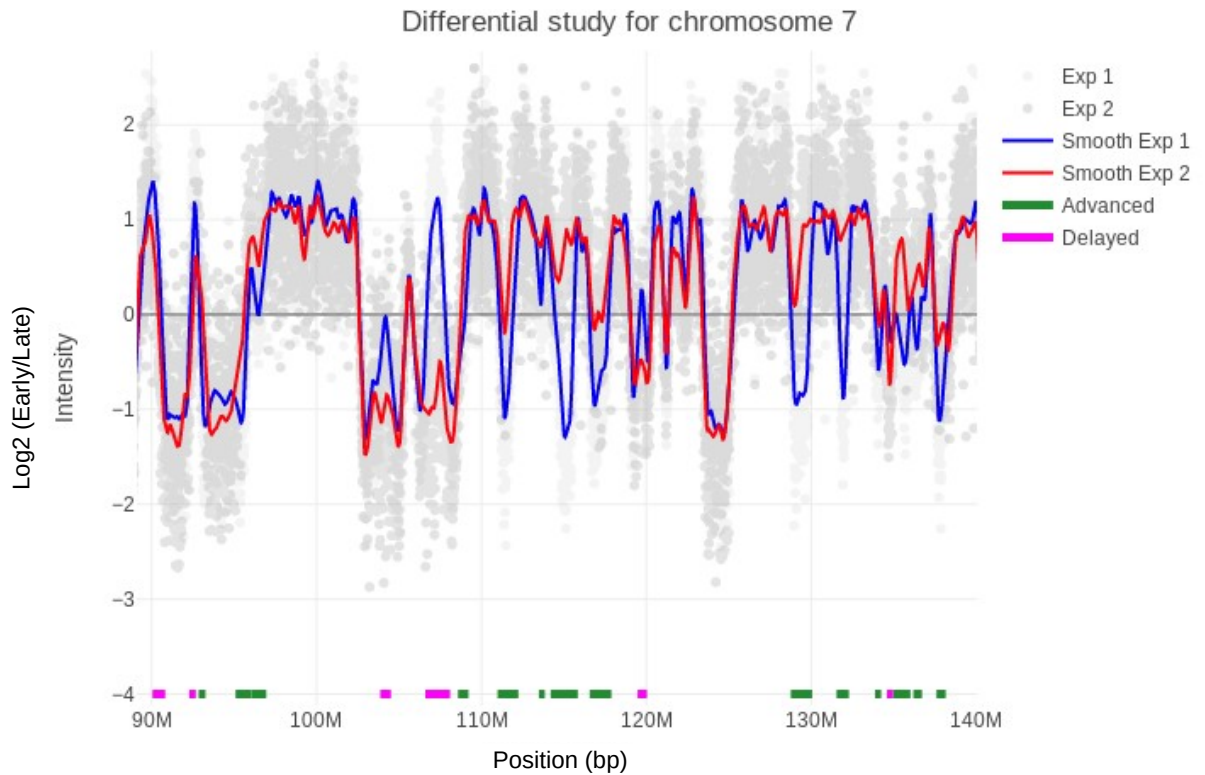
27

1 **START-R analysis of replication-timing programs during differentiation in mouse: a**
2 **new analysis of previous data**

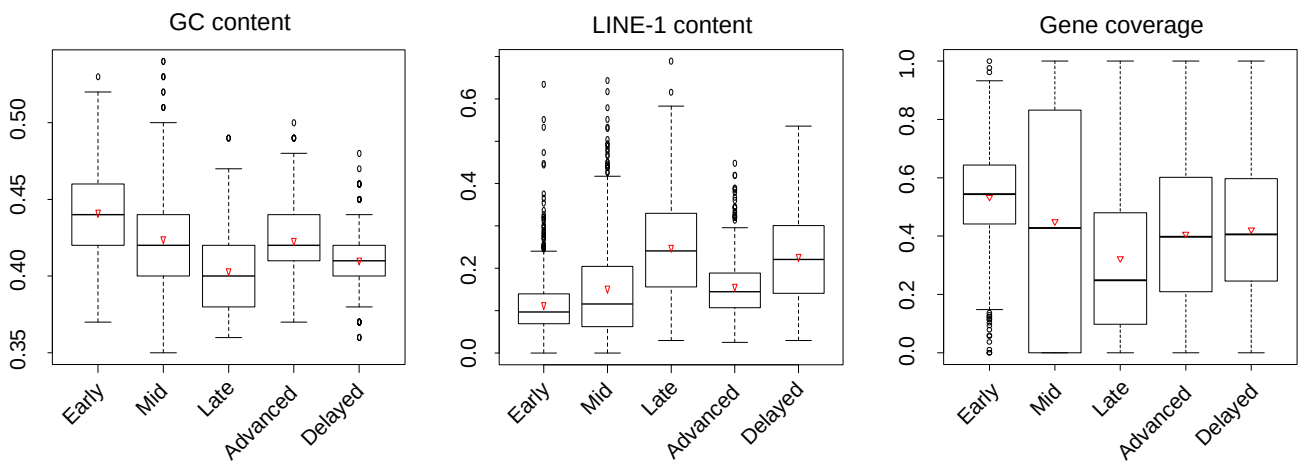
3 To validate our START-R based-approach without *a priori* consideration, we decided to re-
4 analyze the data obtained by the Gilbert's group concerning the changes of replication-timing
5 program during cell differentiation in mouse (Hiratani et al., 2008). They found that 20% of
6 replication domains change between the D3esc and D3npc9 cell lines. There are two types
7 of changes: Early-to-Late (EtoL or delayed) and Late-to-Early (LtoE or advanced). Each
8 modified timing region had a particular molecular signature: LtoE regions show a GC/LINE-1
9 density and gene coverage similar to constant early regions, while EtoL regions showed GC/
10 LINE-1 density and gene coverage similar to constant late regions. We used the same raw
11 data for START-R analysis. First, we converted the raw data with the convertPair.R script to
12 be in the correct format for START-R Analyzer. Then, we used START-R Analyzer with the
13 standard option: Loess Early/Late normalisation, scale inter-replica normalisation, inter-
14 experiments standardization, Loess method for smoothing (span=300kb), 2.5 for SD
15 difference between two segments, and Holm's method with a p-value=0.05 for the differential
16 analysis. With these parameters, 2,066 CTRs are detected in the genome and 910 regions
17 show a different replication timing between D3esc and D3npc9 (Fig. 5A). Thanks to START-
18 R Analyzer that automatically generates BED files, it is easy to import files into a GALAXY
19 session (Afgan et al., 2018) in order to continue the molecular characterization and to
20 generate complementary results. Advanced and delayed regions show the aforementioned
21 specific molecular signatures for the GC/LINE-1 content and gene coverage (Fig. 5B;
22 Hiratani et al., 2008). Unfortunately, we cannot compare the regions that we have identified
23 with START-R with those previously discovered, since the article of Hiratani *et al* does not
24 mention the genome coordinates. However, our results confirm that START-R Analyzer is
25 robust, whatever the mammalian replication timing program studied and whatever the type of

1 microarrays used (from Agilent or Affimetrix or Nimblegen). In fact, we detected identical
2 molecular signatures in the regions showing RT changes.

3 **A**



16 **B**



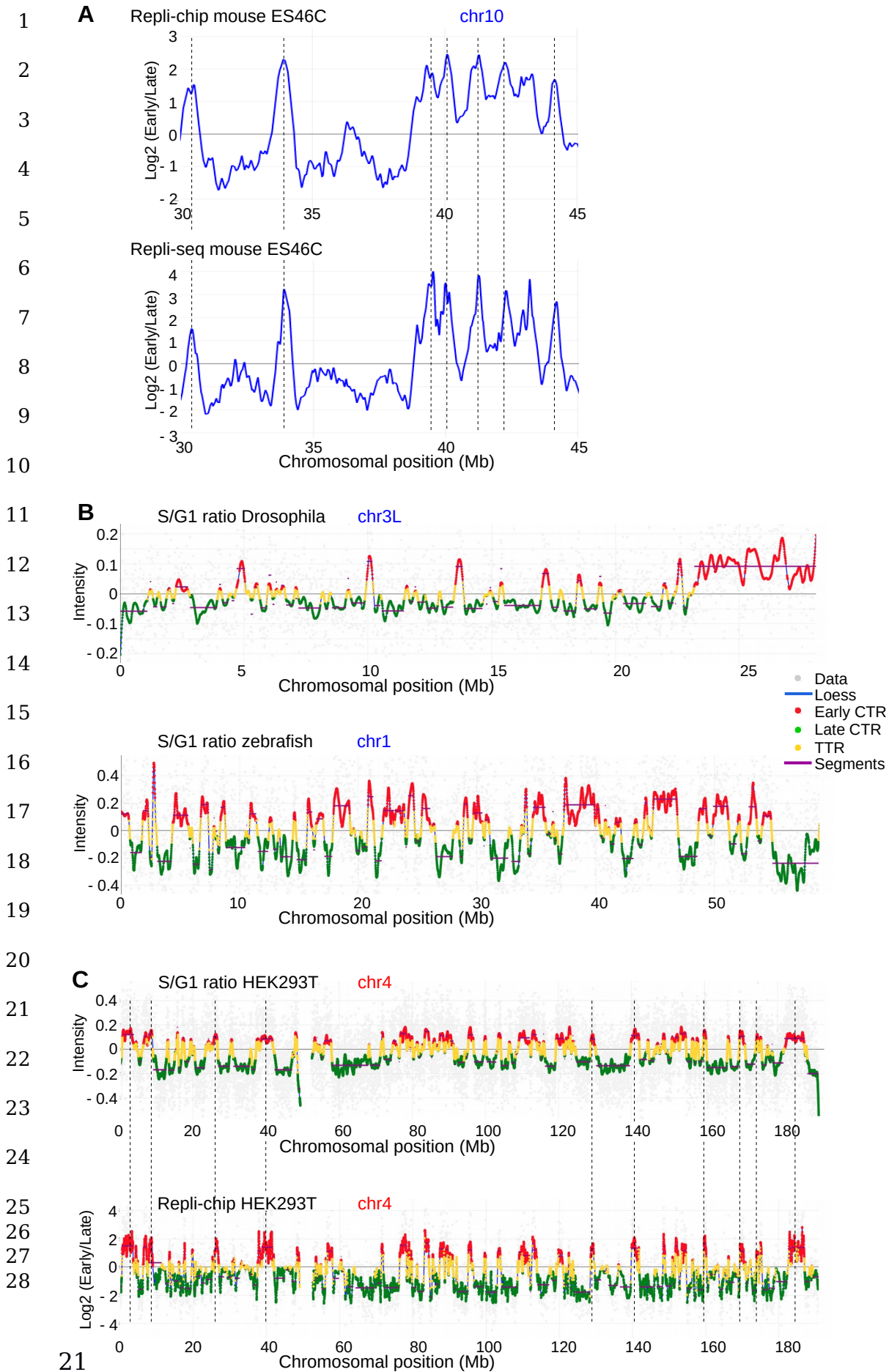
23 **Figure 5. Genomic characteristics of regions harboring different replication timing**
24 **programs.**

25 **A)** START-R Differential analysis of RT profiles is shown for a portion of chromosome 7 in
26 mouse D3esc (blue) and D3npc9 (red) cells. Light grey and grey spots indicate data from
27 both RT experiments. **B)** Boxplots illustrate differences in GC content, LINE-1 content and
28 gene coverage between Early, Mid and Late replicating regions. The two other categories

1 show the characteristics of Advanced and Delayed regions. For each category, the mean
2 value is indicated by an open red triangle. The band at the middle of the box indicates the
3 median value. The bottom and top of the box are the 25th and 75th percentiles. Bottom and
4 top whiskers represent the limits with exclusion of outliers (open circles).
5

6 **Validation of START-R with Early-Late Repli-seq data from mouse**

7 Many data of replication-timing program can be obtained with Repli-seq experiments, but
8 their analysis is time-consuming and often require bioinformatics skills. We analyzed the
9 Early/Late repli-seq data from Marchal and co-workers (Marchal et al., 2018). With different
10 genomic tools (using GALAXY in our case), we obtained the alignment of reads and a BAM
11 coverage file essential for the integration in the START-R pipeline. We specifically developed
12 a supplemental script to convert the BAM coverage file to a log Early/Late file
13 (`convert_bamcoverage_file.R`) to be sure that the integration into the START-R pipeline was
14 correct. Then, we compared the RT smooth profile from E/L Repli-seq with similar data
15 obtained with microarrays (Fig. 6A). The profiles are almost identical, exactly as described by
16 Marchal and co-workers (Marchal *et al.*, 2018). Thus, START-R Analyzer and Viewer can be
17 easily used to analyze E/L Repli-seq data, showing its versatility and its simplicity of use.



1 **Figure 6. The START-R suite allows analysis and visualisation of both Repli-chip and**
2 **Repli-seq data from different model systems.**

3 **A)** RT profiles of a portion of mouse chromosome 10 from ES46C cell line are generated
4 using Repli-chip (top panel) and Repli-seq data (bottom panel) with START-R software.
5 Dashed vertical lines show common RT regions between both profiles. **B)** RT profiles
6 obtained by S/G1 ratios are shown for the left part of Drosophila chromosome 3 (3L) and for
7 zebrafish chromosome 1 (blue lines). The profiles display distribution of early and late CTRs,
8 in red and green, respectively and of TTRs, in yellow. Segments corresponding to regions of
9 constant timing are shown in purple. Grey spots indicate data from RT experiments. **C)** RT
10 profiles of human HEK293T chromosome 4 are generated using S/G1 ratio and Repli-chip
11 data. The empty space inside the RT profiles represent the centomere region.
12

13

14 **Validation of START-R with S-G1 Repli-seq data from Drosophila, zebrafish and**
15 **human.**

16 Other laboratories use the ratio of DNA content between G1 and S phases to analyze the RT
17 program. We wanted to know if START-R suite can run the correct analyses with this type of
18 data and also with other organisms than mouse and human. We performed exactly the same
19 pipeline used for early-late Repli-seq data described above for Drosophila, zebrafish and
20 human S/G1 data (Armstrong et al., 2018; Siefert et al., 2017; Massey et al., 2019), in order
21 to be sure that the integration into the START-R pipeline was correct. Then and as expected,
22 START-R can be run with S/G1 log ratio data for Drosophila, zebrafish and human (Fig. 6B
23 and 6C). We observed similar profiles as the ones already observed for these different
24 organisms. However, our analysis of HEK293T RT changes with the Repli-chip method gave
25 higher differences between distant extreme values, which optimized the detection of RT
26 changes.

1 **Discussion**

2 In this study, we show a new automated protocol for generating and analyzing RT profiles in
3 human and mouse genomes. This approach relies on both the automation of the IP step and
4 on new web-based softwares, START-R Analyzer and Viewer (see graphical abstract in
5 Supplemental Material). The IP-Star® robot reduces the length of the IP step from 3 days to
6 an overnight experiment allowing the user to test 16 samples at the same time. This protocol
7 is very interesting because samples could then be treated either by Repli-seq (2-fractions) or
8 by microarrays. In addition, since the main experiment is carried out by a robot, there are few
9 differences due to different experimenters. As demonstrated by our results, the degree of
10 reproducibility of experiments using the IP-Star® robot is very high (Supplemental Fig. S2).
11 In addition, the choice of the cell-sorting window is primordial. The window extension
12 overlapping G1 and G2/M (Fig. 1A) gives exactly the same profile obtained with the “6-
13 fractions Repli-seq” method (Fig. 1B and 1C), in shorter time and reduced financial costs.
14 The “6-fractions Repli-seq” approach, which is long and expensive, can discourage a number
15 of labs. This approach makes these experiments much more affordable with the same level
16 of precision.

17 The START-R suite facilitates the analysis by making it more accessible to non-
18 bioinformatician researchers. User-friendly interfaces integrate all used steps to generate RT
19 profiles (Fig. 2) and users can choose different parameters at every step. Compared to the
20 previous, no longer in use method (Ryba et al., 2011), START-R Analyzer detects TTR
21 regions and better refines and improves the CTRs detection. In addition, START-R Analyzer
22 contains new calculation methods for the identification of differences between two conditions
23 or two cell lines (Fig. 3). This flexibility gives the users the opportunity to choose the
24 differential analysis method and different parameters according to their questions (Fig. 4,
25 Supplemental Figs. S10 and S11). START-R also processes data from different organisms
26 (Figs. 5 and 6) and those obtained with different methods, such as two fractions Repli-seq,

1 two fractions Repli-chip and S/G1 fractions data. It also automatically generates files with
2 different output formats essential for further molecular characterizations and compatible for
3 classical bioinformatic tools and/or for GALAXY genomic tools. START-R Analyzer is not
4 exclusively developed for mammalian genome as we also generated RT analyses for
5 *Drosophila* and zebrafish genomes (Fig. 6).

6 START-R Viewer produces a nice interface to visualize all the data generated by START-R
7 Analyzer. It facilitates the analysis of RT. In addition, it makes easier the navigation along the
8 genome to take screenshots suitable for future figures (Supplemental Fig. S5K).

9 START-R Analyzer and START-R viewer freewares are available on GitHub and their source
10 codes are open to anyone who wants to improve and to integrate them into a personal
11 computer or server, whatever the operating system. Finally, the START-R suite is very
12 flexible since it can use data from different microarray platforms, from Repli-seq experiments
13 (with 2 fractions or S/G1 fractions), and from different organisms (Figs. 2, 3, 5 and 6). In
14 conclusion, it is now possible for any biologist or laboratory to readily explore new or
15 previous replication timing data simply and quickly. Thus, a large number of laboratories can
16 today use our approach and our softwares to find out if their experimental conditions are
17 affecting the replication timing process or are correlated with other molecular mechanisms.

18 START-R also allows to determine what parts of the genome are impacted and to
19 characterize further those loci. Thanks to the accessibility of our approaches and softwares,
20 their speed and efficiency, new research perspectives can be efficiently envisaged.

1 **Methods**

2 **Experimental procedures to obtain early and late DNA replicating fractions**

3 BrdU incorporation and cell fixation

4 $3 \cdot 10^7$ exponentially growing mammalian cells were incubated with 0.5 mM BrdU (Abcam,
5 #142567), protected from light, at 37°C for 90 minutes. Cells fixed in 75% final cold EtOH can
6 be stored at -20°C.

7

8 Cell sorting

9 10^7 BrdU labeled cells were incubated with 80 $\mu\text{g}/\text{mL}$ Propidium Iodide (Invitrogen, P3566)
10 and with 0,4 mg/ml RNaseA (Roche, 10109169001) for 1h at room temperature with orbital
11 shaking at 180 rpm. 10^5 cells were sorted in early and late S phase fractions using a
12 Fluorescence Activated Cell Sorting system (INFLUX 500 Cytopeia, BD Biosciences) in Lysis
13 Buffer (50mM Tris pH=8, 10mM EDTA, 0.5% SDS, 300mM NaCl).

14

15 DNA extraction and sonication

16 DNA from sorted cells was extracted using Proteinase K treatment (200 $\mu\text{g}/\text{ml}$, Thermo
17 Scientific, EO0491) followed by phenol-chloroform extraction and sonicated to a size of 500-
18 1,000 base pair (bp), as previously described (Hadjadj et al., 2016).

19

20 Immunoprecipitation using SX-8G IP-Star® Compact Automated System (Diagenode)

21 Immunoprecipitations from 10^5 cells were performed using IP star robot at 4°C (indirect 200 μl
22 method, Diagenode) with an anti-BrdU antibody (10 μg , purified mouse Anti-BrdU, BD
23 Biosciences, #347580). Denatured DNA was incubated 5 hours with anti-BrdU antibodies in
24 IP buffer (10mM Tris pH=8, 1mM EDTA, 150mM NaCl, 0.5% Triton X-100, 7mM NaOH)
25 followed by 5 hours incubation with Dynabeads Protein G (Invitrogen, 10004D) (Hadjadj et
26 al., 2016). Beads were then washed with Wash Buffer (20mM Tris pH=8, 2mM EDTA,

1 250mM NaCl, 1% Triton X-100). Reversion was performed at 37°C during 2 hours with a
2 solution containing 1% SDS and 0.5mg Proteinase K followed, after the beads removal, by
3 an incubation at 65°C during 6 hours in the same solution.

4

5 DNA purification and quantitative PCR

6 Immunoprecipitated BrdU labeled DNA fragments were extracted with phenol-chloroform and
7 precipitated with cold ethanol. Control quantitative PCRs (qPCRs) were performed using
8 oligonucleotides specific of mitochondrial DNA, early (*BMP1* gene) or late (*DPPA2* gene)
9 replicating regions (Ryba et al., 2011; Hadjadj et al., 2016).

10

11 Amplification

12 Whole genome amplification was performed using SeqPlex™ Enhanced DNA Amplification kit
13 as described by the manufacturer (Sigma-Aldrich, SEQXE). Amplified DNA was purified
14 using PCR purification product kit as described by the manufacturer (Macherey-Nagel,
15 740609.50). DNA amount was measured using a Nanodrop. Quantitative PCRs using the
16 oligonucleotides described above were performed to check whether the ratio between early
17 and late replication regions was still maintained after amplification.

18

19 Universal Linkage System (ULS™) labeling, chip loading and scanning for Repli-chip 20 experiment (microarray)

21 Early and late nascent DNA fractions were labelled with Cy3-ULS and Cy5-ULS,
22 respectively, using the ULS arrayCGH labeling Kit (Kreatech, EA-005).

23 Same amounts of early and late-labeled DNA were loaded on mouse or human DNA
24 microarrays (SurePrint G3 Human CGH arrays, Agilent Technologies, G4449A).

25 Hybridization was performed as previously described (Hadjadj et al., 2016). The following
26 day microarrays were scanned using an Agilent C-scanner with Feature Extraction 9.1

1 software (Agilent technologies). RT datasets are available in the Gene Expression Omnibus
2 (GEO) database under the accession numbers GSM2111308 for U2OS, GSM2111313 for
3 K562 and GSM2111310 for HEK293T cell lines.

4

5 **START-R suite**

6 START-R Analyzer and START-R Viewer are open source web-based applications (doi:
7 10.5281/zenodo.3243339), developed using the Shiny R package (Chang et al., 2018).
8 START-R suite was concatenated using Docker (Supplemental Fig. S1) in order to install,
9 use and share it easily. START-R softwares can be used with different operating systems:
10 Windows, Mac OS and Linux. The source code and the installation procedure are available
11 on GitHub (<https://github.com/thomasdenecker/START-R>) with information on how to use
12 the software.

13 To install START-R, users should read and follow the Readme file available on GitHub web
14 site (<https://github.com/thomasdenecker/START-R/blob/master/README.md>). Briefly, users
15 should install Docker and then follow installation procedures described for each OS in the
16 Readme file. Finally, in order to run the START-R suite, the user should double-click on the
17 START-R file (Windows) or launch the command line (Linux / MacOS X), followed by
18 opening an internet browser at the following URLs: <http://localhost:3838/> for START-R
19 Analyzer and <http://localhost:3839/> for START-R Viewer (Supplemental Fig. S1).

20

21 **Early/Late Repli-seq data and conversion**

22 In order to validate our softwares, we used data from the GEO database, with the accession
23 numbers GSM2496038 and GSM2496039, corresponding to Repli-seq 46C mouse cells -
24 Early S fraction or Late S fraction, respectively. Data were managed using different tools
25 from GALAXY server (Afgan et al., 2018) but data can also be processed manually using
26 specific command lines of the algorithms used on a local computer. Read mapping was

1 obtained using Bowtie2 (2.3.4.2 version) with the very sensitive end-to-end option. Then,
2 PCR duplicates were removed by RmDUP from SAMTools (2.0.1 version). We used
3 BamCoverage (3.1.2.0.0 version with default parameters) with a bin size of 10kb
4 (corresponding to the spacing of probes on microarrays) and a Reads Per Kilobase Million
5 (RPKM) normalisation to generate a bedGraph file. A headline was added to the file to name
6 the 4 columns (chr, start, end, gProcessedSignal for early file or rProcessedSignal for late
7 file, respectively). Then, a script to convert and merge the bedGraph files from early and late
8 samples to a format compatible with START-R analyzer was developed. This script is
9 available on GitHub (<https://github.com/thomasdenecker/START-R>) in “supplement script”
10 file as [convert_bamcoverage_file.R](#).

11

12 **Validation of START-R suite using microarray data from other laboratories**

13 We compared microarray data with Repli-seq data obtained with mouse ES46C cell line
14 (GEO accession numbers: GSM2496037 and GSM2496038-039, respectively). Then, we
15 analyzed with the START-R suite the microarrays data obtained by Hiratani *et al.* (2008) of
16 D3esc and D3npc9 cell lines during mouse cells differentiation (GEO accession numbers:
17 GSM450273 and GSM450285, respectively). As data extracted from the Nimblegen platform
18 are in PAIR format, we used a script to convert data into a valid format for START-R
19 Analyzer (convertPair.R, available on GitHub <https://github.com/thomasdenecker/START-R>
20 in “supplement script” file).

21

22 **Validation of START-R suite using S/G1 data from multiple species**

23 Different laboratories analyze variations of DNA copy number between G1 and S phase cells
24 (S/G1 ratio) to study the replication timing program with Repli-seq. We used data obtained
25 from different organisms such as *Drosophila*, zebrafish and human (Armstrong *et al.*, 2018;
26 Siefert *et al.*, 2017; Massey *et al.*, 2019), to validate the START-R suite (GEO accession
27 numbers: GSM3154888 and GSM3154890 for HWT *Drosophila melanogaster* female larvae

1 wing disc cells in S and G1 phase, respectively; GSM2282090 for 28hpf *Danio rerio*
2 embryos; SRX3413939-40 for HEK293T human cells in S and G1 phase, respectively). As
3 previously, reads from G1 and S fractions are mapped with Bowtie2, then PCR duplicates
4 are removed by RmDUP tool, and Bamcoverage is used to obtain the coverage with RPKM.
5 Then, as above, the bedGraph files are converted to a format compatible with START-R
6 analyzer with "[convert_bamcoverage_file.R](#)." script.

7

8 **GC content, Long Interspersed Nuclear Elements-1 (LINE-1) and gene coverage** 9 **calculation**

10 LINE-1 elements coordinates were extracted from the University of California Santa Cruz
11 (UCSC) Genome Browser Repeat Masker track (Smit et al., 1996-2010). Two steps are
12 required to calculate GC content: (i) DNA sequence extraction from Early, Mid, Late,
13 Advanced, Delayed regions coordinates with the Extract Genomic DNA tool (2.2.4 version
14 with default parameters); (ii) calculation of GC content with the GeeCee EMBOSS tool (5.0.0
15 version with default parameters) present in the GALAXY website (Afgan et al., 2018).
16 Coverage tools from GALAXY website were used to determine LINE-1 content and gene
17 coverage. All boxplots were made using R program (3.5.1 version, R Core Team, 2019).

18

19 **Data Access**

20 All raw data generated in this study have been submitted to the NCBI Gene Expression
21 Omnibus (GEO; <https://www.ncbi.nlm.nih.gov/geo/>) under accession number GSE141122.
22 The START-R suite and the script used to perform the analysis in this study are available on
23 GitHub at <https://github.com/thomasdenecker/START-R>,

24

25

1 **Acknowledgments**

2 This project was supported by La Ligue Nationale Contre le Cancer (RS16/75-108 and
3 RS17/75-135), the Groupement des Entreprises Françaises en Lutte Contre le Cancer
4 (GEFLUC), the Institut National du Cancer INCa-10493, the IdEx Université de Paris ANR-
5 18-IDEX-0001 and by the generous legacy from Ms Suzanne Larzat to our group.

6 We thank Gaëlle Lelandais for discussions during START-R elaboration and for allowing
7 Thomas Denecker to work on this project during his PhD work. We also acknowledge the
8 ImagoSeine core facility of the Institut Jacques-Monod, member of the France BioImaging
9 (ANR-10-INBS-04) supported by the Region Île-de-France (E539).

10

11 **Disclosure Declaration**

12 The authors declare that they have no conflict on interest.

13

14 **References**

15 Afgan E, Baker D, Batut B, Van Den Beek M, Bouvier D, Ech M, Chilton J, Clements D,
16 Coraor N, Grüning BA, et al. 2018. The Galaxy platform for accessible, reproducible and
17 collaborative biomedical analyses: 2018 update. *Nucleic Acids Res* **46**: W537–W544.

18 Almeida R, Fernández-Justel JM, Santa-María C, Cadoret J-C, Cano-Aroca L, Lombraña R,
19 Herranz G, Agresti A, Gómez M. 2018. Chromatin conformation regulates the
20 coordination between DNA replication and transcription. *Nat Commun* **9(1)**: 1590-

21 Armstrong RL, Penke TJR, Strahl BD, Matera AG, McKay DJ, MacAlpine DM, Duronio RJ.
22 2018. Chromatin conformation and transcriptional activity are permissive regulators of
23 DNA replication initiation in *Drosophila*. *Genome Res.* **(11)**:1688-1700.

24 Brustel J, Kirstein N, Iazard F, Grimaud C, Prorok P, Cayrou C, Schotta G, Abdelsamie AF,
25 Déjardin J, Méchali M, et al. 2017. Histone H4K20 trimethylation at late-firing origins
26 ensures timely heterochromatin replication. *EMBO J* **36(18)**: 2726-2741.

- 1 Chang W, Cheng J, Allaire J, Xie Y, McPherson J. 2018. shiny: Web Application Framework
2 for R. R package version 1.2.0.No Title. <https://CRANR-project.org/package=shiny>.
- 3 Cornacchia D, Dileep V, Quivy JP, Foti R, Tili F, Santarella-Mellwig R, Antony C, Almouzni
4 G, Gilbert DM, Buonomo SBC. 2012. Mouse Rif1 is a key regulator of the replication-
5 timing programme in mammalian cells. *EMBO J* **31**: 3678–3690.
- 6 Dileep V, Didier R, Gilbert DM. 2012. Genome-wide analysis of replication timing in
7 mammalian cells: Troubleshooting problems encountered when comparing different cell
8 types. *Methods* **57**: 165–169.
- 9 Dileep V, Rivera-Mulia JC, Sima J, Gilbert DM. 2015. Large-Scale Chromatin Structure–
10 Function Relationships during the Cell Cycle and Development: Insights from
11 Replication Timing. *Cold Spring Harb Symp Quant Biol* **80**: 53–63.
- 12 Gilbert DM. 2010. Cell fate transitions and the replication timing decision point. *J Cell Biol*
13 **191**: 899–903.
- 14 Hadjadj D, Denecker T, Maric C, Fauchereau F, Baldacci G, Cadoret JC. 2016.
15 Characterization of the replication timing program of 6 human model cell lines.
16 *Genomics Data* **9**: 113-117.
- 17 Hanahan D. and Weinberg RA. 2011 Hallmarks of cancer: The next generation. *Cell* **144**:
18 646–674.
- 19 Hansen RS, Thomas S, Sandstrom R, Canfield TK, Thurman RE, Weaver M, Dorschner MO,
20 Gartler SM, Stamatoyannopoulos JA. 2010. Sequencing newly replicated DNA reveals
21 widespread plasticity in human replication timing. *Proc Natl Acad Sci* **107**: 139–144.
- 22 Hiratani I, Ryba T, Itoh M, Yokochi T, Schwaiger M, Chang CW, Lyou Y, Townes TM,
23 Schübeler D, Gilbert DM. 2008. Global reorganization of replication domains during
24 embryonic stem cell differentiation. *PLoS Biol* **6**: 2220–2236.
- 25 Kent WJ, Sugnet CW, Furey TS, Roskin KM, Kent WJ, Sugnet CW, Furey TS, Roskin KM.
26 2002. The Human Genome Browser at UCSC. *Genome Res* **19**: 1228–31.

- 1 Krzywinski M, Altman N. 2013. Error bars. *Nat Methods* **10**: 921-922.
- 2 Macheret M, Halazonetis TD. 2015. DNA Replication Stress as a Hallmark of Cancer. *Annu*
3 *Rev Pathol Mech Dis* **10**: 425–448.
- 4 Marchal C, Sasaki T, Vera D, Wilson K, Sima J, Rivera-Mulia JC, Trevilla-García C, Nogues
5 C, Nafie E, Gilbert DM. 2018. Genome-wide analysis of replication timing by next-
6 generation sequencing with E/L Repli-seq. *Nat Protoc* **13**: 819–839.
- 7 Massey DJ, Kim D, Brooks KE, Smolka MB, Koren A. 2019. Next-Generation Sequencing
8 Enables Spatiotemporal Resolution of Human Centromere Replication Timing. *Genes*
9 (Basel) **10(4)**: pii: E269
- 10 R Core Team. 2019. R: A language and environment for statistical computing. *R Found Stat*
11 *Comput Vienna, Austria URL <https://www.R-project.org/>.*
- 12 Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, Smyth GK. 2015. limma powers
13 differential expression analyses for RNA-sequencing and microarray studies. *Nucleic*
14 *Acids Res* **43(7)**: e47.
- 15 Rivera-Mulia JC, Gilbert DM. 2016. Replicating Large Genomes: Divide and Conquer. *Mol*
16 *Cell* **62**: 756–765.
- 17 Ryba T, Hiratani I, Lu J, Itoh M, Kulik M, Zhang J, Schulz TC, Robins AJ, Dalton S, Gilbert
18 DM. 2010. Evolutionarily conserved replication timing profiles predict long-range
19 chromatin interactions and distinguish closely related cell types. *Genome Res* **20(6)**:
20 761-770.
- 21 Ryba T, Battaglia D, Pope BD, Hiratani I, Gilbert DM. 2011. Genome-scale analysis of
22 replication timing: From bench to bioinformatics. *Nat Protoc* **6**: 870–895.
- 23 Siefert JC, Georgescu C, Wren JD, Koren A, Sansam CL. 2017. DNA replication timing
24 during development anticipates transcriptional programs and parallels enhancer
25 activation. *Genome Res* **8**:1406-1416.
- 26 Sievert C, Parmer C, Hocking T, Chamberlain S, Ram K, Corvellec M, Despouy P. 2017.
27 plotly: Create Interactive Web Graphics via “plotly.js”. R package version 4.7.1.

1 <https://CRANR-project.org/package=plotly>.

2 Smit A, Hubley R, Green P. 1996-2010. RepeatMasker Open-3.0.

3 <http://www.repeatmasker.org>.

4 Técher H, Koundrioukoff S, Nicolas A, Debatisse M. 2017. The impact of replication stress
5 on replication dynamics and DNA damage in vertebrate cells. *Nat Rev Genet* **18(9)**:
6 535-550.