



Published in final edited form as:

*Nat Genet.* 2019 January ; 51(1): 36–41. doi:10.1038/s41588-018-0285-7.

## Error-prone bypass of DNA lesions during lagging strand replication is a common source of germline and cancer mutations

Vladimir B. Seplyarskiy<sup>#1,2,3</sup>, Evgeny E. Akkuratov<sup>#4,5</sup>, Natalia Akkuratova<sup>#5</sup>, Maria A. Andrianova<sup>6</sup>, Sergey I. Nikolaev<sup>7,8</sup>, Georgii A. Bazykin<sup>3,6</sup>, Igor Adameyko<sup>9,10</sup>, and Shamil R. Sunyaev<sup>1,2,\*</sup>

<sup>1</sup>Division of Genetics, Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA

<sup>2</sup>Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA

<sup>3</sup>Institute for Information Transmission Problems of the Russian Academy of Sciences (Kharkevich Institute), Moscow, Russia

<sup>4</sup>Present address: Science for Life Laboratory, Department of Applied Physics, Royal Institute of Technology, Stockholm, Sweden

<sup>5</sup>Institute of Translational Biomedicine, Saint Petersburg State University, St. Petersburg 199034, Russia

<sup>6</sup>Skolkovo Institute of Science and Technology, Skolkovo, Russia

<sup>7</sup>UMR8200 - CNRS, Stabilité Génétique et Oncogénèse, France; Gustave Roussy Cancer Campus, F-94805, Villejuif, France; Université Paris Saclay, Paris Sud - Orsay, F-91400, France.

<sup>8</sup>Department of Dermatology and Venereology, Université Paris 7, Saint Louis Hospital, FR-75010 Paris, France.

<sup>9</sup>Department of Physiology and Pharmacology, Karolinska Institutet, 171 77, Stockholm, Sweden.

<sup>10</sup>Center for Brain Research, Medical University Vienna, 1090 Vienna, Austria.

# These authors contributed equally to this work.

### Abstract

Users may view, print, copy, and download text and data-mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use:[http://www.nature.com/authors/editorial\\_policies/license.html#terms](http://www.nature.com/authors/editorial_policies/license.html#terms)

\* Correspondence should be addressed to [ssunyaev@rics.bwh.harvard.edu](mailto:ssunyaev@rics.bwh.harvard.edu).

Authors Contributions

V.B.S., G.A.B. and S.R.S. designed the study. V.S.B. performed the data analyses. V.B.S., E.E.A., N.V.A. and I.A. designed and performed experiments. M.A.A. performed data preprocessing and helped with results presentation. S.I.N. retrieved genomic data for squamous cell carcinoma. V.B.S. and S.R.S. drafted the manuscript. All authors contributed to the final version of the paper.

Competing Financial Interests

Authors declare no financial interest

Accession Codes

Experimental data generated in this study deposited in sequence read archive under accession number SRP151915 (<https://www.ncbi.nlm.nih.gov/sra/SRP151915>)

Studies in experimental systems have identified a multitude of mutational mechanisms including DNA replication infidelity and DNA damage followed by inefficient repair or replicative bypass. However, the relative contributions of these mechanisms to human germline mutation remain unknown. Here, we show that error-prone damage bypass on the lagging strand plays a major role in human mutagenesis. Transcription-coupled DNA repair removes lesions on the transcribed strand; lesions on the non-transcribed strand are preferentially converted into mutations. In human polymorphism we detect a striking similarity between mutation types predominant on non-transcribed strand and on the strand lagging during replication. Moreover, damage-induced mutations in cancers accumulate asymmetrically with respect to the direction of replication, suggesting that DNA lesions are resolved asymmetrically. We experimentally demonstrate that replication delay greatly attenuates the mutagenic effect of UV-irradiation confirming that replication converts DNA damage into mutations. We estimate that at least 10% of human mutations arise due to DNA damage.

---

Experiments in model organisms have uncovered that DNA polymerases make errors and resulting mismatches become mutations<sup>1</sup>. An alternative mechanism of mutagenesis due to mis-repaired DNA damage or damage bypassed by translesion (TLS) polymerases has been extensively studied in experimental systems exposed to exogenous mutagens<sup>2,3</sup>. Although this sheds light on the mechanistic details of mutagenesis, experimental systems provide little information on the relative contributions of these mechanisms to naturally occurring human mutations. Recently, computational genomics approaches have revealed statistical properties of mutations occurring in the germline<sup>4-7</sup>, in tumors<sup>8</sup> and in embryo during early stages of development<sup>9-11</sup>. In cancer, many types of mutations have been successfully attributed to the action of specific mutagens<sup>12</sup>. A number of studies have explored how cancer mutations scale with age at diagnosis<sup>13,14</sup> and how human germline mutations scale with paternal age<sup>15-17</sup>. It was hypothesized that the dependency of the number of accumulated mutations on the number of cell divisions may reflect the replicative origin of mutations<sup>18,19</sup>. However, a quantitative model suggests that accumulation of both damage-induced and co-replicative mutations may scale with the number of cell divisions<sup>20</sup>. Therefore, we still do not know whether DNA damage substantially contributes to human mutations or whether they mostly arise due to errors in replication.

To discriminate between co-replicative mutations and damage-induced mutations, we rely on statistical properties of mutations unequivocally associated with DNA damage. Both germline and cancer mutations leave footprints in the form of mutational asymmetry with respect to the direction of transcription (T-asymmetry). T-asymmetry reflects the prevalence of mutations that originate from lesions on the non-transcribed strand that could not be repaired by transcription-coupled repair (TC-NER)<sup>21,22</sup>. Thus, the analysis of T-asymmetry may be used to quantify mutations arising from DNA lesions. Genomic data on cancers in which most mutations are caused by specific, well-understood, damage-inducing agents provide an additional perspective on properties of damage-induced mutations. Notably, the level of T-asymmetry in these cancers is exceptionally high.

The most obvious statistical feature associated with replication is asymmetry with respect to the direction of the replication fork (R-asymmetry). R-asymmetry may reflect differential

fidelity of replication between the leading and lagging strands. Alternatively, R-asymmetry may be caused by the strand-specific bypass of DNA damage. Bulky DNA lesions not repaired prior to replication can either lead to fork regression followed by error-free repair or be bypassed by TLS polymerases<sup>23,24</sup>. It does not remove the lesion and commonly introduces mutations on the newly synthesized strand. It has been asserted that the error-prone bypass process has different properties on leading and lagging strands<sup>23,27</sup> that would lead to R-asymmetry.

As a starting point, we compare R-asymmetry with T-asymmetry. To avoid the interference of statistical signals between the two types of asymmetries, R-asymmetry is estimated exclusively in intergenic regions and T-asymmetry only in genic regions. We calculate R-asymmetries for the 92 types of single-nucleotide mutations in each trinucleotide context. CpG>TpG mutations are excluded because they usually arise via conversion of methylated cytosines directly into thymines by deamination<sup>28</sup> (Supplementary Note 1). Figure 1 shows data for rare (allele frequency below 0.1%) SNPs from the gnomAD dataset. Supplementary Figure 1 shows that R-asymmetries across different contexts are concordant between rare SNPs and *de novo* mutations.

Strikingly, there is a high concordance between T-asymmetry and R-asymmetry across mutation types in both tri- and penta-nucleotide contexts (Figure 1, Supplementary Figure 2). Mutation types that are predominant on the lagging strand are also more common on the non-transcribed strand (Figure 1a;  $R^2=0.84$ ;  $p\text{-value}=5.6*10^{-37}$ ). This association holds even for six basic mutation classes separately (Figure 1b).

T-asymmetry arises from mutations induced by DNA damage on the non-transcribed strand that is invisible to TC-NER repair<sup>6</sup>. As a result, level of T-asymmetry scales with the proportion of damage-induced mutations. Figure 1 suggests that R-asymmetry may be due to similarly differential resolution of DNA damage between leading and lagging strands. DNA lesions on the lagging strand would be more frequently converted into mutations, probably due to error-prone damage bypass.

To follow up on this hypothesis, we analyze R-asymmetry in genomes of cancers exposed to specific mutagens. Four cancer types in PCAWG datasets contain samples with high levels of T-asymmetry in specific mutation contexts: melanoma, predominated by UV-induced C>T mutations (signature 7)<sup>8</sup>; two lung cancers (LUAD and LUSC), predominated by smoking-induced G>T mutations (signature 4); and liver cancer, with a high prevalence of A>G mutations (signatures 12 and 16). All of these processes reflect the action of DNA-damaging mutagens. We find that about 95% of these samples demonstrate a weak but usually significant excess of mutations on the lagging strand (Figure 2, Supplementary Table 1). The mutagens acting primarily outside of replication also cause R-asymmetry, strongly suggesting that error-prone bypass on the lagging strand happens frequently. Levels of T-asymmetry and R-asymmetry are correlated across samples in lung and liver cancers (Supplementary Figure 3). Melanoma is the exception, possibly because variation in T-asymmetry across samples is primarily due to variation in TC-NER activity rather than damage intensity (Supplementary Figure 3). A recent study also found an excess of damage-induced mutations corresponding to COSMIC signatures 23 (unknown etiology) and 24

(aflatoxin) on the lagging strand<sup>29</sup>. Consistently with our interpretation, samples lacking damage-induced signatures do not exhibit a lagging strand bias (Supplementary Table 2).

The observed R-asymmetry is limited to samples with signatures of bulky damage rather than any type of damage. Tumors with *MUTYH* deficiency have a high load of mutations from oxo-guanine lesions that do not block progression of RNA and DNA polymerases. Neither T- nor R-asymmetry is detectable in these samples (Supplementary Figure 4). In contrast, R-asymmetry is significantly enhanced in cutaneous squamous cell carcinoma tumors from patients with congenital *XPC* deficiency (Xeroderma Pigmentosum)<sup>30</sup>. These tumors lack the global genome repair (GG-NER) activity and have elevated levels of bulky damage (See Supplementary Note 2 and Supplementary Figure 5).

If DNA lesions are more frequently bypassed by TLS on the lagging strand directly across the lesion, they will persist on this strand through replication. Therefore, mutational asymmetry caused by the bypass in turn causes the asymmetry of unrepaired DNA damage. We utilize time series XR-seq data<sup>31</sup> to test whether the activity of the NER system is biased with respect to the replication fork direction. Indeed, repair is more frequently observed on the lagging strand (Figure 3). Moreover, the difference between leading and lagging strands sharply increases with time after UV irradiation as more and more cells complete a round of replication.

To test whether the differential activity of the NER system reflects the preferential bypass of DNA damage, we analyze the Damage-seq dataset<sup>32</sup>. Damage-seq was used to detect DNA damage (cyclobutane pyrimidine dimers) over a series of time points following the exposure of fibroblasts to UV radiation. The data show a clear dependency on transcription and preferential retention of damage on the non-transcribed strand (Supplementary Figure 6). A lagging strand bias of DNA damage progressively increases with time, mirroring the trend in XR-seq data (Figure 3).

Collectively, the above observations support the differential replication bypass hypothesis suggesting that many damage-induced mutations do not arise from mis-repair; instead bulky lesions are converted to mutation during replication. Thus, replication delay should reduce mutation rate in cells exposed to damaging agents, because it would provide more time for cells to complete repair. To test this directly, we compared UV-irradiated fibroblasts exposed and not exposed to roscovitine, which reversibly arrests replication (Figure 4). Colonies grown from fibroblasts not treated with the chemical have ~14,000 mutations with spectra matching the UV-signature (Figure 5). These mutations demonstrated both T- and R-asymmetries quantitatively similar to cancer data ( $\log_2(\text{T-asym.})=0.50$ ,  $\log_2(\text{R-asym.})=0.17$ ,  $p<0.01$  for both). In sharp contrast, colonies derived from UV-irradiated cells that experienced replication delay (51 hours of roscovitine treatment) possessed just ~2000 mutations with no evident UV-signature. Control cells that were treated by roscovitine but not exposed to UV irradiation have a highly similar spectrum of mutations and only ~400 fewer mutations. Therefore, replication delay decreased UV-induced mutation load by more than 30 fold. This provides a strong support for the error-prone replication bypass of bulky lesions being the major source of mutations, at least in our experimental system.

Interestingly, it also suggests that mutations in melanoma are primarily accumulating in dividing cells.

R-asymmetry not related to error-prone bypass was previously detected in several cancers and in experimental systems. It was attributed to differences in fidelity between Polymerase  $\epsilon$  and Polymerase  $\delta$ <sup>1,33–35</sup> or differential efficiency of mismatch repair between leading and lagging strands<sup>33–36</sup>. APOBEC deaminates cytosines on the lagging strand<sup>33,35,37,38</sup> and misincorporation of oxo-guanine in esophageal cancer is highly asymmetric<sup>39</sup>. However, these processes neither match patterns observed for human germline mutations nor explain the strong association between R- and T-asymmetries and experimental data on UV-irradiated cells. A mechanism alternative to error-prone bypass may be responsible for R-asymmetry of CpG>TpG mutations in the human germline that are not caused by bulky lesions (Supplementary Figure 7 and Supplementary Note 1).

One possible alternative explanation for the similarity between R-asymmetry and T-asymmetry involves the exposure of DNA to a single-stranded conformation (ssDNA): the lagging strand stays in the single-stranded state during replication for a longer period, while the non-transcribed strand may occasionally adopt the single-stranded state because of R-loop formation<sup>40,41</sup>. We have tested the effect of R-loops on T-asymmetry and found that, in the germline, the asymmetry does not increase in regions prone to R-loops compared to flanking regions within the same transcript (Supplementary Figure 8a). Additional clues to the role of ssDNA may be provided by APOBEC-induced mutations because APOBEC mutations have a strong affinity for ssDNA<sup>42,43</sup>. Again, we do not find that R-loops substantially affect the distribution of APOBEC-induced mutations in cancers (Supplementary Figure 8b). Hence, it is unlikely that ssDNA is the cause of T-asymmetry and of the association between T- and R-asymmetries.

Taken together, the observed mutation patterns in the germline and in cancer, XR-seq and Damage-seq data and our experiments point to differential damage bypass as a likely source of R-asymmetry. This suggests that DNA damage substantially contributes to spontaneous mutations. T-asymmetry allows us to conservatively quantify its contribution. Assuming that DNA damage is uniform, TC-NER is completely error-free and is the only cause of the T-asymmetry (see Methods), we compute the minimal fraction of damage-induced mutations in highly transcribed genes. Extrapolation of this estimate to the whole genome suggests that 10% of human germline mutations, 51–52% of mutations in melanoma, 40–44% of mutations in lung cancer, and 25–27% of mutations in liver cancer are due to DNA damage. The estimates are high for cancers affected by known mutagens, although still lower than existing estimates<sup>8</sup>, attesting to the conservative nature of our analysis.

From the biochemical perspective, a higher conversion rate of damage due to mutations on the lagging strand is unsurprising, as replication of the leading strand is less tolerant to damage. Helicase is attached to the leading strand and is therefore more sensitive to damage on this strand<sup>23,27</sup>. Damage on the leading strand blocks Polymerase  $\epsilon$ , potentially causing fork uncoupling and stalling. This, in turn, may cause fork regression with lesion repair, template switch or homologous repair<sup>23</sup> – all these processes are error-free. With the exception of break-induced replication producing highly complex mutations, fork stalling is

usually resolved by error-free mechanisms. Meanwhile, lesions on the lagging strand are unlikely to cause fork stalling and instead often only result in a short gap downstream from the lesion<sup>23,27</sup>. Consequently, damage on the lagging strand is rarely removed during replication and is instead simply bypassed by error-prone mechanisms (TLS). Our results corroborate earlier findings in the yeast system, where as much as 90% of spontaneous mutations have been attributed to TLS through DNA lesions<sup>44,45</sup>.

Our experimental results show that the number of damage-induced mutations reduces with the increasing timespan between introduction of DNA damage and cell division. The computational analysis suggests that mutations statistically associated with replication do not necessarily arise as a result of replication errors alone. Earlier studies have demonstrated the dependency of the number of accumulated mutations on the number of cell divisions. In line with theoretical models, we note that the mutation rate scaling with the number of replications does not establish the mechanistic origin of mutations<sup>20</sup>.

## Methods

### Human polymorphism and cancer mutation data

To analyze mutational patterns reflected in human DNA polymorphism, we extracted SNPs with derived allele frequency <0.1% from gnomAD data<sup>48</sup>. Cancer somatic mutations were extracted from PCAWG dataset<sup>49</sup>. Cancer somatic mutations identified in XPC<sub>wt</sub> and XPC<sub>-/-</sub> skin SCC samples were downloaded from dbGap (phs000830). Samples with MUTYH deficiency were chosen according to annotation from Scarpa et al<sup>50</sup>.

### Experimental data on DNA damage and repair

The XR-seq dataset for cyclobutane pyrimidine dimers (CPD) reported in Adar *et al.*<sup>31</sup> allowed us to estimate the amount of DNA damage actively repaired by NER following UV irradiation. To directly assess the presence of unrepaired DNA damage, we used the Damage-seq data for CPDs provided by Hu *et al.*<sup>32</sup>. We did not use XR-seq and Damage-seq data for pyrimidinepyrimidone (6–4) photoproducts because these lesions are repaired too quickly to permit an accurate analysis of the effect of damage bypass over successive rounds of replication.

### R-asymmetry

As described previously<sup>37</sup>, the “derivative” (normalized rate of change) of replication timing may serve as a predictor of the preferential replication fork direction. This approach was proposed by Chen *et al.*<sup>5</sup> and has been used in recent cancer genomics studies<sup>33,35</sup>.

We focused on genomic regions showing a strong preference for a specific fork direction as evident from the replication timing “derivative”. For the analysis, XR-seq, and Damage-seq (Figures 1a and b, Figures 3a and b), we used a conservative threshold corresponding to 10% of genomic regions with the highest absolute values of the replication timing “derivative”. However, this threshold appeared too restrictive for cancer genome analyses because many individual tumors have insufficient numbers of mutations within the 10% of the genome, so we relaxed the threshold to 40% for these analyses. Both of these thresholds have been used



in previous studies<sup>7,33</sup>, and the results have generally been robust with respect to the threshold chosen.

For each individual analysis, we selected the most relevant available replication timing dataset: IMR-90 for lung cancers, HepG2 for liver cancer, and NHEK for melanoma and squamous carcinoma. For germline mutations, there is no relevant cell and we decided to consider regions with replication direction conserved across tissue types requiring that all 7 tissues have the same sign of the replication timing “derivative”; and at least in half of the tissues (4 out of 7) have value of the “derivative” exceeding 40% threshold. We also used replication timing data obtained from NHEK cell line to predict the preferential fork direction in the analysis of XR-seq and Damage-seq data and our experimental dataset (matching the tissue but not the exact cell type).

For each mutation type, we calculated R-asymmetry as the ratio of mutation density on the lagging strand to the mutation density on the leading strand. Samples with fewer than 100 mutations on each strand were excluded from the analysis to reduce sampling noise.

XPC knockouts have a distinct mutational spectrum that is dominated by TpCpT>TpTpT mutations (Supplementary Figure 9) and we restrict our test to this mutation type. Supplementary Figure 5 focuses on the magnitude of the effect in each tumor rather than on the presence of the effect. We therefore excluded samples with fewer than 500 mutations on each strand. The relaxation of the threshold to 100 mutations does not change the conclusions (data not shown).

In order to exclude the impact of T-asymmetry on the R-asymmetry estimation, we restricted the analysis of R-asymmetry to intergenic regions.

### T-asymmetry

For each mutation type, we calculated T-asymmetry as a ratio of mutation density on the transcribed strand to mutation density on the non-transcribed strand. Gene annotations and transcription direction were determined according to the knownGene track of the UCSC genome browser. Tumors with T-asymmetry >1.2 (for any of the six major mutation classes) were considered to have high level of T-asymmetry. Even with this lenient criterion, only four cancer types (melanoma, LUAD, LUSC, and liver cancer) had more than 20 tumor samples in this category. To order the genes by their expression levels, we selected the most relevant tissues from Gtex<sup>51</sup>: testis for SNPs from gnomAD, sun-exposed skin for melanoma, liver for liver cancer, and lung for lung cancers.

### Exclusion of replica B2 at 48h from Damage-seq

T-asymmetry and the difference between genic and non-genic regions are the main results of the Damage-seq experiments<sup>32</sup> that support the utility of the data for the genome-wide analysis of bulky DNA damage and repair by the NER system. Thus, for quality control of the Damage-seq data, we calculated T-asymmetry and the ratio of reads in intergenic and genic regions separately for all replicas. T-asymmetry and the ratio of reads in intergenic and genic regions were normalized using the corresponding values for naked DNA. We found that the replicates were generally concordant at each time point with the exception of the

48h point, where we found substantial T-asymmetry and prevalence of mutations in intergenic regions in replica A but essentially no signal in replica B2 (Supplementary Figure 10). At other time points, we observed a clear, time-dependent increase in T-asymmetry and decrease in the fraction of damages in genic regions, as expected. Based on these observations, we argue that the absence of the signal in replica B2 at 48h is an artifact. Therefore, this data point was excluded. As shown in Supplementary Figure 10c, this replica is also a clear outlier in the analysis of R-asymmetry.

### Estimate of the proportion of mutations arising due to DNA damage in human cancers and the germline

To conservatively estimate the proportion of damage-induced mutations, we capitalized on the statistical signal of T-asymmetry that is associated with DNA damage. The T-asymmetry introduced by co-transcriptional processes cannot be a consequence of replication infidelity. Therefore, mutations responsible for the T-asymmetry must be damage-induced. Since transcribed and non-transcribed regions can have different susceptibilities to DNA damage, we conservatively compared the levels of mutations between transcribed strand and immediately adjacent flanking sequences rather than between transcribed and non-transcribed strands:

$$t = \frac{\mu_{intergenic}}{\mu_{transcribed\_strand}},$$

where  $\mu_{transcribed\_strand}$  is the mutation density on the transcribed strand and  $\mu_{intergenic}$  is the mutation density in flanking intergenic regions.

To estimate  $t$ , we used the 10% of genes with the highest expression levels. We conservatively assumed that all damage on transcribed strands is efficiently repaired. Thus, the fraction of damage-induced mutations in transcribed regions and in intergenic regions is expressed as:

$$f_{genic} = \frac{t-1}{t+1}$$

$$f_{intergenic} = \frac{t-1}{t}$$

If  $a$  denotes the fraction of mutations in genic regions, and  $b$  is the fraction of mutations in intergenic regions, the fraction of damage-induced mutations for the whole genome ( $f_{genome}$ ) is expressed as:

$$f_{genome} = af_{genic} + bf_{intergenic}$$



The conservative nature of this estimate is evident in the cancer data. Although nearly all mutations in melanoma are caused by UV irradiation, our estimate attributes only 50% of mutations to DNA damage (Supplementary table 3).

Confidence intervals have been obtained by sampling mutations with replacement 200 times for human polymorphism and by resampling tumors 200 times for cancers

### R-loops

We used data on strand-specific R-loops from Sanz *et al.*<sup>40</sup>. Most R-loops were on the template strand, and we considered only such R-loops. For control regions, we used intronic regions within the same gene that were 500 nucleotides apart from the R-loop peak and 500 nucleotides long.

### CpG islands

Annotation of CpG islands was downloaded from the UCSC genome browser (cpgIslandExt).

### Experimental procedures

Human fibroblast cells from skin (GM00637) were purchased from the National Institute of General Medical Sciences Human Genetic Cell Repository (Coriell Institute). They were maintained with Minimum Essential Medium (M5650, Sigma Aldrich) supplemented with 10% fetal bovine serum (10270–106, Gibco) and 2mM L-Glutamine at 37°C in the 5% CO<sub>2</sub>.

We generated genetically homogenous colonies via two successive passages starting from a single cell.

Cells were irradiated with a lamp (112537, Merck) emitting 254 nm UV light (2 J/(m<sup>2</sup>\*sec)) during 10 seconds, resulting in 20 J/m<sup>2</sup> irradiation. For a subset of colonies, we added 30µM roscovitine (R7772, Sigma Aldrich). For cells to be UV irradiated, roscovitine was added 3 hours prior to the UV treatment.

After 48 hours of incubation without changing the medium we split cells with low density in order to select individual colonies and subsequently cultivate them to achieve 1\*10<sup>6</sup> cells (approximately for 4 weeks). DNA was isolated with PureLink Genomic DNA Mini Kit (K182000, ThermoFisher) and then sequenced by MacroGen Inc on Illumina's HiSeq X Ten with the average coverage of 30X. Overall, we produced six colonies including two treated with roscovitine and UV-irradiation; two UV-irradiated, but with no roscovitine in the medium; one colony incubated with roscovitine, but not irradiated; and one control colony that was not treated (Figure 4).

To quantify the change in proliferation after treatment with UV-light and/or roscovitine cells on coverslips were incubated with 5 µg/ml EdU for 24 hours. Then for each condition we made 15 measurements (5 different regions on 3 coverslips; at average 20 cells per region) of the fraction of cells that incorporated EdU. Cells were stained with the EdU detection kit (Click-iT EdU Imaging kit C10337, Thermo Fisher) to count divided cells and stained with Hoechst to count the total number of cells.

In each sample, we measured the proliferation rate via EdU incorporation during the first 24 hours (adding EdU 5 minutes after UV-irradiation and staining cells after 24 hours) and second 24 hours (adding EdU after 24 hours and staining the cells after 48 hours). Examples of the EdU staining are shown in Supplementary Figure 11 and Supplementary Figure 12. Adding roscovitine to the medium decreases proliferation rate by 2–3 fold compared to the control population (Supplementary Figure 13). UV-irradiation itself decreased the proliferation rate by 5 fold, followed by a substantial recovery on day 2. Combination of the UV-irradiation and roscovitine almost completely halted cell proliferation both on days one and two. Moreover, we observed that during the colony selection, cells treated with roscovitine grew slower than non-treated cells.

### Mutation calling

To obtain the set of mutations from sequenced reads, we performed following steps: first we trimmed reads with TrimGalore-0.4.5 in paired mode, then we mapped reads with bwa-0.7.12 according GATK best practice, then we call mutations from bam files with MuTect2 using the control colony (no roscovitine treatment or UV-irradiation) as “normal” and other colonies (treated with roscovitine, UV or both) colonies as “tumor”. Finally, we filtered out all the mutations observed in more than one colony. Mutation spectra for all replicates are shown on Supplementary Figure 14.

### Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

### Acknowledgments

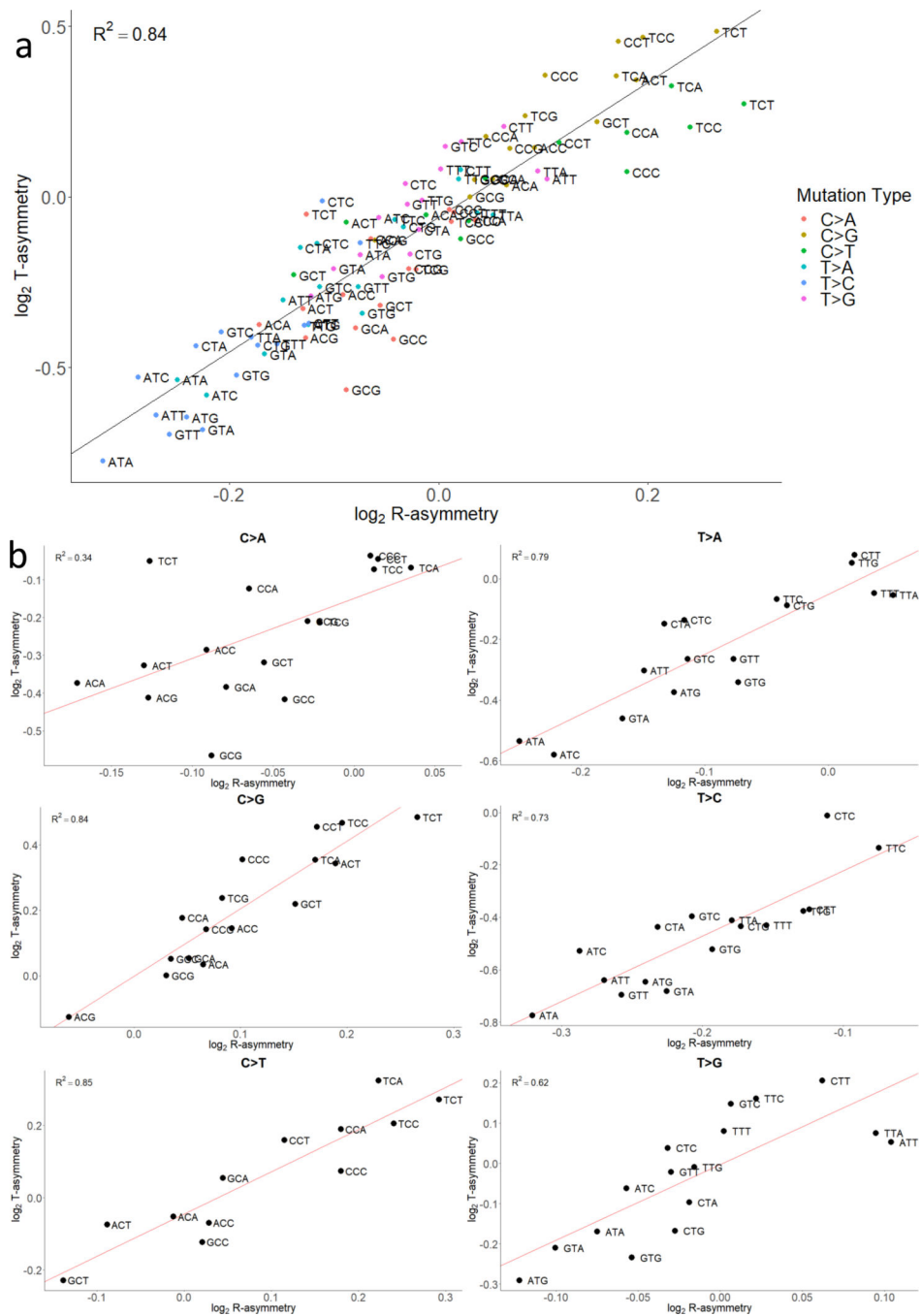
We thank Sergei Mirkin, Dmitry Gordenin, Cristopher Cassa and Donate Weghorn for useful comments on the manuscript, Lionel Sanz and Frédéric Chédin for help with R-loop data and Blake Boulerice for proofreading. This study was supported by grants: R35GM127131, R01MH101244 and U01HG009088; N.A. was supported from the Russian Science Foundation grant №16–15–10273

### References

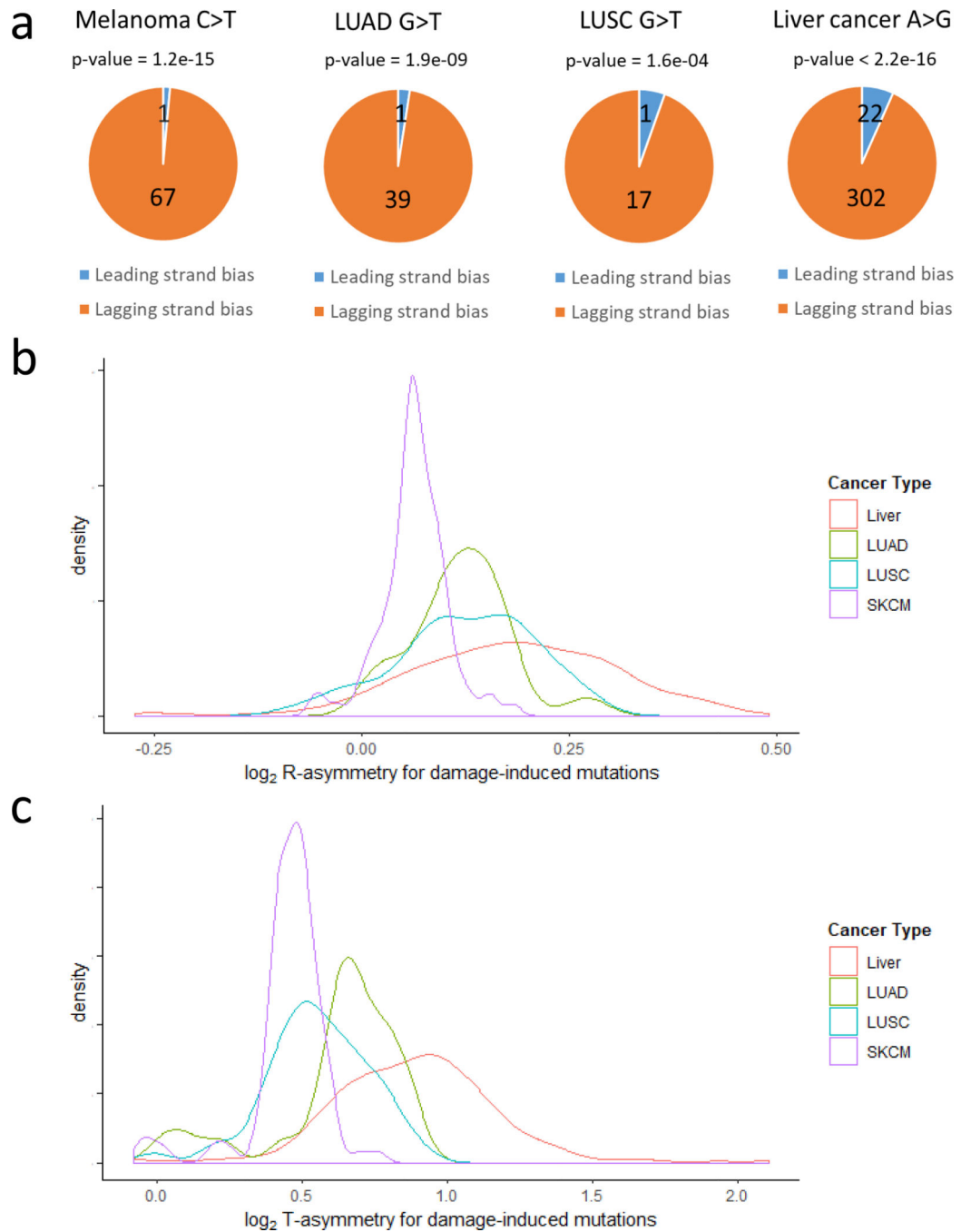
1. Lujan SA et al. Heterogeneous polymerase fidelity and mismatch repair bias genome variation and composition. *Genome Res.* 24, 1751–1764 (2014). [PubMed: 25217194]
2. Boiteux S & Jinks-Robertson S DNA Repair Mechanisms and the Bypass of DNA Damage in *Saccharomyces cerevisiae*. *Genetics* 193, 1025–1064 (2013). [PubMed: 23547164]
3. Cohen IS et al. DNA lesion identity drives choice of damage tolerance pathway in murine cell chromosomes. *Nucleic Acids Res.* 43, 1637–1645 (2015). [PubMed: 25589543]
4. Baker A et al. Replication fork polarity gradients revealed by megabase-sized U-shaped replication timing domains in human cell lines. *PLoS Comput. Biol.* 8, e1002443 (2012). [PubMed: 22496629]
5. Chen C-L et al. Replication-associated mutational asymmetry in the human genome. *Mol. Biol. Evol.* 28, 2327–2337 (2011). [PubMed: 21368316]
6. Polak P & Arndt PF Transcription induces strand-specific mutations at the 5' end of human genes. *Genome Res.* 18, 1216–1223 (2008). [PubMed: 18463301]
7. Seplyarskiy VB, Andrianova MA & Bazykin GA APOBEC3A/B-induced mutagenesis is responsible for 20% of heritable mutations in the TpCpW context. *Genome Res.* 27, 175–184 (2017). [PubMed: 27940951]
8. Alexandrov LB et al. Signatures of mutational processes in human cancer. *Nature* 500, 415–421 (2013). [PubMed: 23945592]

9. Harland C et al. Frequency of mosaicism points towards mutation-prone early cleavage cell divisions. *bioRxiv* 079863 (2016). doi:10.1101/079863
10. Lindsay SJ, Rahbari R, Kaplanis J, Keane T & Hurles M Striking differences in patterns of germline mutation between mice and humans. *bioRxiv* 082297 (2016). doi:10.1101/082297
11. Ju YS et al. Somatic mutations reveal asymmetric cellular dynamics in the early human embryo. *Nature* 543, 714–718 (2017). [PubMed: 28329761]
12. Helleday T, Eshtad S & Nik-Zainal S Mechanisms underlying mutational signatures in human cancers. *Nat. Rev. Genet* 15, 585–598 (2014). [PubMed: 24981601]
13. Alexandrov LB et al. Clock-like mutational processes in human somatic cells. *Nat. Genet* 47, 1402–1407 (2015). [PubMed: 26551669]
14. Podolskiy DI, Lobanov AV, Kryukov GV & Gladyshev VN Analysis of cancer genomes reveals basic features of human aging and its role in cancer development. *Nat. Commun* 7, (2016).
15. Kong A et al. Rate of de novo mutations and the importance of father's age to disease risk. *Nature* 488, 471–475 (2012). [PubMed: 22914163]
16. Francioli LC et al. Genome-wide patterns and properties of de novo mutations in humans. *Nat. Genet* 47, 822–826 (2015). [PubMed: 25985141]
17. Wong WSW et al. New observations on maternal age effect on germline de novo mutations. *Nat. Commun* 7, 10486 (2016). [PubMed: 26781218]
18. Moorjani P, Gao Z & Przeworski M Human Germline Mutation and the Erratic Evolutionary Clock. *PLoS Biol.* 14, e2000744 (2016). [PubMed: 27760127]
19. Tomasetti C, Li L & Vogelstein B Stem cell divisions, somatic mutations, cancer etiology, and cancer prevention. *Science* 355, 1330–1334 (2017). [PubMed: 28336671]
20. Gao Z, Wyman MJ, Sella G & Przeworski M Interpreting the Dependence of Mutation Rates on Age and Time. *PLoS Biol.* 14, e1002355 (2016). [PubMed: 26761240]
21. Fousteri M & Mullenders LHF Transcription-coupled nucleotide excision repair in mammalian cells: molecular mechanisms and biological effects. *Cell Res.* 18, 73–84 (2008). [PubMed: 18166977]
22. Martejn JA, Lans H, Vermeulen W & Hoeijmakers JHJ Understanding nucleotide excision repair and its roles in cancer and ageing. *Nat. Rev. Mol. Cell Biol* 15, 465–481 (2014). [PubMed: 24954209]
23. Yeeles JTP, Poli J, Marians KJ & Pasero P Rescuing stalled or damaged replication forks. *Cold Spring Harb. Perspect. Biol* 5, a012815 (2013). [PubMed: 23637285]
24. Roberts SA & Gordenin DA Hypermutation in human cancer genomes: footprints and mechanisms. *Nat. Rev. Cancer* 14, 786–800 (2014). [PubMed: 25568919]
25. Supek F & Lehner B Clustered Mutation Signatures Reveal that Error-Prone DNA Repair Targets Mutations to Active Genes. *Cell* 170, 534–547.e23 (2017). [PubMed: 28753428]
26. Rogozin IB et al. DNA polymerase  $\eta$  mutational signatures are found in a variety of different types of cancer. *Cell Cycle* 17, 348–355 (2018). [PubMed: 29139326]
27. Hedglin M & Benkovic SJ Eukaryotic Translesion DNA Synthesis on the Leading and Lagging Strands: Unique Detours around the Same Obstacle. *Chem. Rev* (2017). doi:10.1021/acs.chemrev.7b00046
28. Shen JC, Rideout WM & Jones PA The rate of hydrolytic deamination of 5-methylcytosine in double-stranded DNA. *Nucleic Acids Res.* 22, 972–976 (1994). [PubMed: 8152929]
29. Letouzé E et al. Mutational signatures reveal the dynamic interplay of risk factors and cellular processes during liver tumorigenesis. *Nat. Commun* 8, 1315 (2017). [PubMed: 29101368]
30. Zheng CL et al. Transcription restores DNA repair to heterochromatin, determining regional mutation rates in cancer genomes. *Cell Rep.* 9, 1228–1234 (2014). [PubMed: 25456125]
31. Adar S, Hu J, Lieb JD & Sancar A Genome-wide kinetics of DNA excision repair in relation to chromatin state and mutagenesis. *Proc. Natl. Acad. Sci. U. S. A* 113, E2124–2133 (2016). [PubMed: 27036006]
32. Hu J, Adebali O, Adar S & Sancar A Dynamic maps of UV damage formation and repair for the human genome. *Proc. Natl. Acad. Sci. U. S. A* (2017). doi:10.1073/pnas.1706522114

33. Haradhvala NJ et al. Mutational Strand Asymmetries in Cancer Genomes Reveal Mechanisms of DNA Damage and Repair. *Cell* 164, 538–549 (2016). [PubMed: 26806129]
34. Andrianova MA, Bazykin GA, Nikolaev SI & Seplyarskiy VB Human mismatch repair system balances mutation rates between strands by removing more mismatches from the lagging strand. *Genome Res.* 27, 1336–1343 (2017). [PubMed: 28512192]
35. Morganello S et al. The topography of mutational processes in breast cancer genomes. *Nat. Commun* 7, 11383 (2016). [PubMed: 27136393]
36. Lujan SA et al. Mismatch repair balances leading and lagging strand DNA replication fidelity. *PLoS Genet* 8, e1003016 (2012). [PubMed: 23071460]
37. Seplyarskiy VB et al. APOBEC-induced mutations in human cancers are strongly enriched on the lagging DNA strand during replication. *Genome Res.* 26, 174–182 (2016). [PubMed: 26755635]
38. Hoopes JI et al. APOBEC3A and APOBEC3B Preferentially Deaminate the Lagging Strand Template during DNA Replication. *Cell Rep.* 14, 1273–1282 (2016). [PubMed: 26832400]
39. Tomkova M, Tomek J, Kriaucionis S & Schuster-Boeckler B Widespread impact of DNA replication on mutational mechanisms in cancer. *bioRxiv* 111302 (2017). doi:10.1101/111302
40. Sanz LA et al. Prevalent, Dynamic, and Conserved R-Loop Structures Associate with Specific Epigenomic Signatures in Mammals. *Mol. Cell* 63, 167–178 (2016). [PubMed: 27373332]
41. Skourti-Stathaki K & Proudfoot NJ A double-edged sword: R loops as threats to genome integrity and powerful regulators of gene expression. *Genes Dev.* 28, 1384–1396 (2014). [PubMed: 24990962]
42. Roberts SA et al. Clustered mutations in yeast and in human cancers can arise from damaged long single-strand DNA regions. *Mol. Cell* 46, 424–435 (2012). [PubMed: 22607975]
43. Burns MB et al. APOBEC3B is an enzymatic source of mutation in breast cancer. *Nature* 494, 366–370 (2013). [PubMed: 23389445]
44. Quah S-K, Borstel RC von & Hastings PJ. The Origin of Spontaneous Mutation in *Saccharomyces Cerevisiae*. *Genetics* 96, 819–839 (1980). [PubMed: 7021317]
45. Lawrence CW & Maher VM Mutagenesis in eukaryotes dependent on DNA polymerase zeta and Rev1p. *Philos. Trans. R. Soc. B Biol. Sci* 356, 41–46 (2001).
46. Goldmann JM et al. Parent-of-origin-specific signatures of de novo mutations. *Nat. Genet.* 48, 935–939 (2016). [PubMed: 27322544]
47. Moorjani P, Amorim CEG, Arndt PF & Przeworski M Variation in the molecular clock of primates. *Proc. Natl. Acad. Sci. U. S. A.* 113, 10607–10612 (2016). [PubMed: 27601674]
48. Lek M et al. Analysis of protein-coding genetic variation in 60,706 humans. *Nature* 536, 285–291 (2016). [PubMed: 27535533]
49. Pan-cancer analysis of whole genomes | *bioRxiv*. Available at: <https://www.biorxiv.org/content/early/2017/07/12/162784>. (Accessed: 21st February 2018)
50. Scarpa A et al. Whole-genome landscape of pancreatic neuroendocrine tumours. *Nature* 543, 65–71 (2017). [PubMed: 28199314]
51. Consortium T Gte. The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue gene regulation in humans. *Science* 348, 648–660 (2015). [PubMed: 25954001]



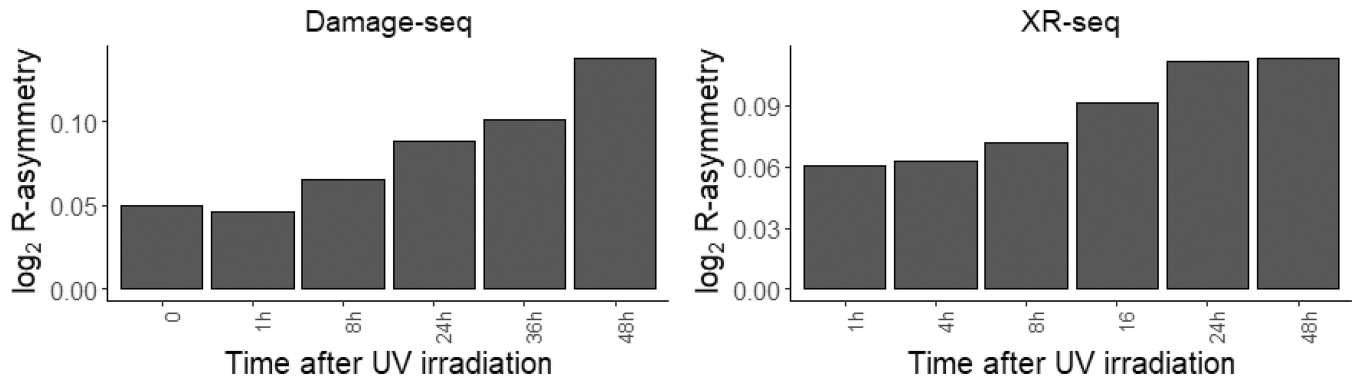
**Figure 1. R-asymmetry and T-asymmetry patterns in human polymorphism.**  
**a**, Relationship between R-asymmetry and T-asymmetry for 92 mutation types (NpCpG>T mutations excluded). **b**, Relationship between R-asymmetry and T-asymmetry shown separately for the six types of single-nucleotide mutations to highlight the effects of adjacent nucleotides.



**Figure 2. Damage-induced mutations preferentially reside on the lagging strand.**

**a**, Number of tumor samples among melanomas, lung adeno carcinomas (LUAD), lung squamous carcinomas (LUSC), and liver cancers that have more damage-induced mutations on the leading than on the lagging strand (p-values shown for the goodness-of-fit chi-square test). **b,c** distribution of R-asymmetry (**b**) and T-asymmetry (**c**) values. Samples with T-asymmetry less than 1.2 were excluded from panel **b**.





**Figure 3. R-asymmetry in UV-irradiated cells.**

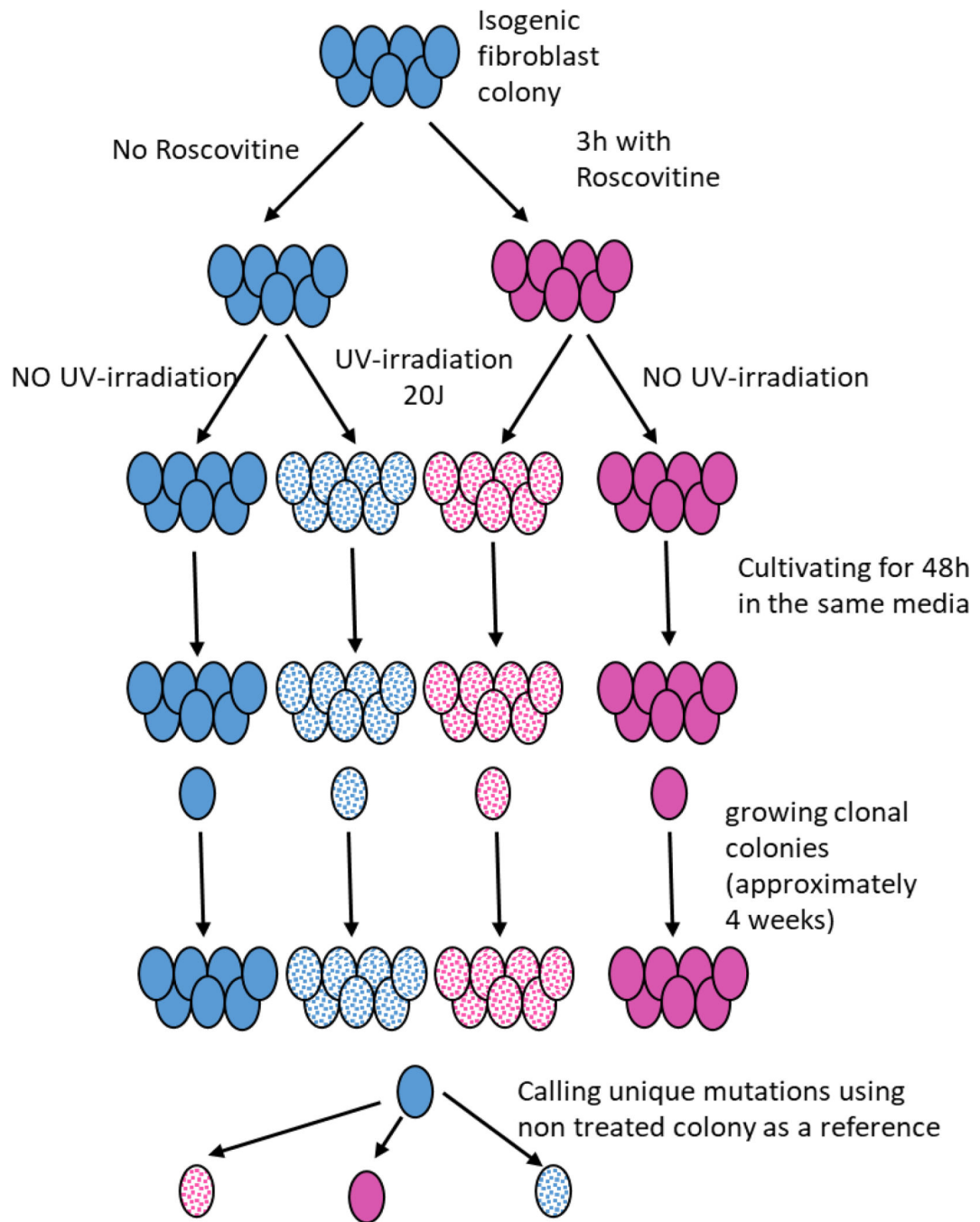
R-asymmetry of repaired CPD damage (left) and CPD damage remaining in DNA (right) as a function of time since UV irradiation.

Author Manuscript

Author Manuscript

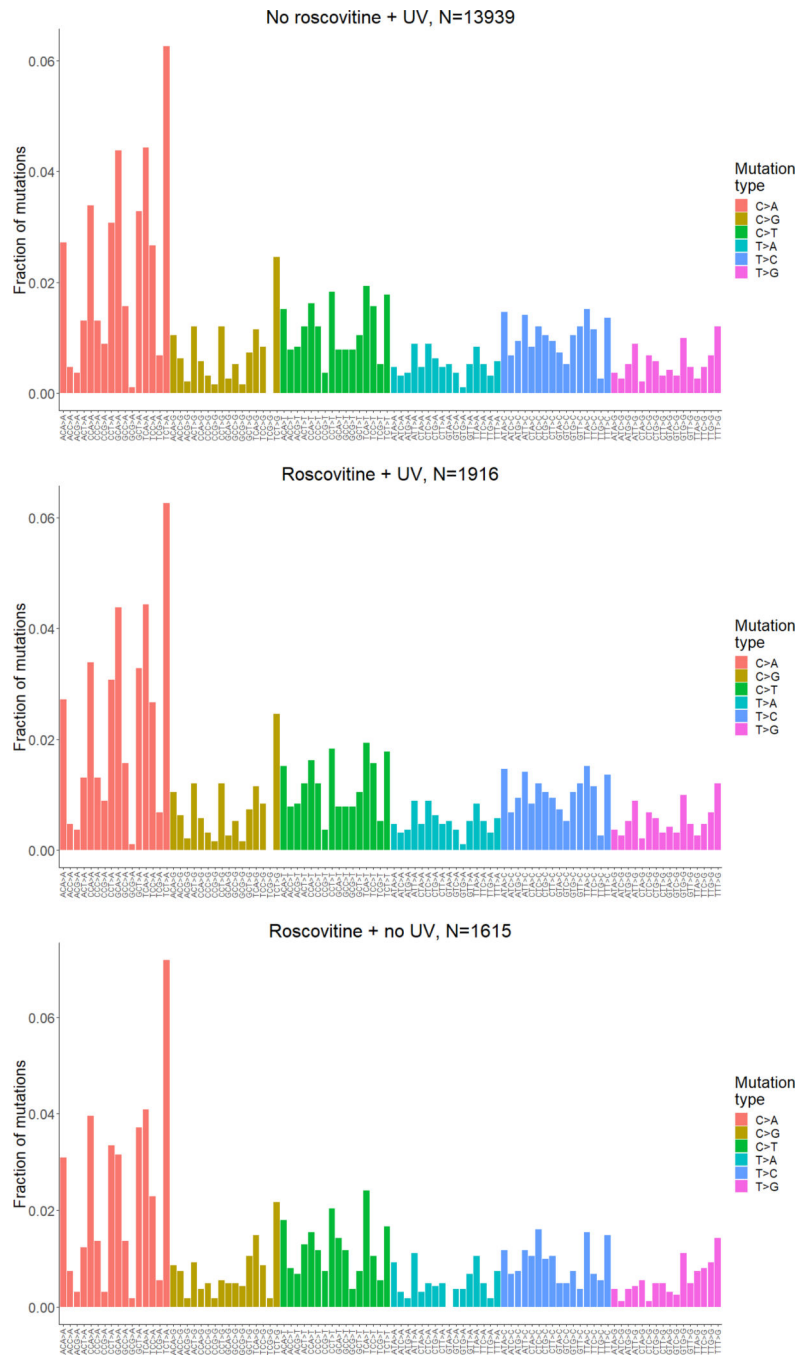
Author Manuscript

Author Manuscript



**Figure 4. Experimental design to test effect of replication delay on the rate of UV-induced mutations.**

Clonal colonies of fibroblast cells shown in pink were treated with roscovitine for 3 hours in advance of the UV-irradiation. Colonies shown in blue were not treated by roscovitine. Half of the colonies were irradiated with UV (20J) (dotted), and the other half were used as control. Randomly chosen cells from each colony were used to start new genetically homogeneous colonies.



**Figure 5.** Quantity and spectra of mutations in fibroblast colonies identified by whole genome sequencing