

Widespread impact of DNA replication on mutational mechanisms in cancer

Tomkova, M.^a, Tomek, J.^b, Kriaucionis, S.^a and Schuster-Böckler, B.^{a*}

AFFILIATION

- a. Ludwig Cancer Research Oxford
University of Oxford
Old Road Campus Research Building
Oxford OX3 7DQ
United Kingdom
- b. Department of Physiology, Anatomy and Genetics
University of Oxford
Oxford OX1 3PT
United Kingdom

KEYWORDS

Mutagenesis, DNA replication, Mutational signatures, DNA repair

MANUSCRIPT

Although mutagens can attack DNA at any time, at least one round of replication is required before damage becomes a fixed mutation. DNA replication therefore plays an important role in mutagenesis, yet little is known about how replication and various mutagenic mechanisms interact. Here, we present the first pan-cancer analysis of the relationship between mutagenic mechanisms, represented by their sequence signatures¹, and DNA replication. Using whole-genome sequencing data from 3056 patients spanning 19 cancer types, we observe a significant impact of replication on 22 out of 29 detected mutational signatures. Association with replication timing and asymmetry around replication origins shed new light on several mutagenic processes, such as suggesting that oxidative damage to the nucleotide pool substantially contributes to the mutational landscape of esophageal adenocarcinoma. Together, our results indicate an involvement of DNA replication and the associated damage repair in most mutagenic processes.

Understanding the mechanisms of mutagenesis in cancer is important for the prevention and treatment of the disease^{2,3}. Mounting evidence suggests replication itself contributes to cancer risk⁴. Copying of DNA is intrinsically asymmetrical, with leading and lagging strands being processed by distinct sets of enzymes⁵, and different genomic regions replicating at defined times during S phase⁶. Previous analyses have focused either on the genome-wide distribution of mutation rate or on the strand specificity of individual base changes. These studies revealed that the average mutation frequency is increased in late-replicating regions^{7,8}, and that the asymmetric synthesis of DNA during replication leads to strand-specific frequencies of base changes⁹⁻¹². However, the extent to which DNA replication influences distinct mutational mechanisms, with their manifold possible causes, remains incompletely understood.

DNA replication in eukaryotic cells is initiated around replication origins (ORI), from where it proceeds in both directions, synthesizing the leading strand continuously and the lagging strand discontinuously (Fig. 1a). We used two independent data sets to describe replication direction relative to the reference sequence, one derived from high-resolution replication timing data¹² and the other from direct detection of ORIs¹³, corrected for technical artefacts¹⁴ (see Methods). The

former provides information for more genomic loci, while the latter is of higher resolution. As a third measure of DNA replication, we compared regions replicating early during S phase to regions replicating late¹².

A mutational signature is a unique combination of frequencies of all 96 possible mutation types (a base-pair mutation, annotated on the pyrimidine, and 5' and 3' flanking nucleotides)¹. Here, we calculated *strand-specific* signatures¹⁵ that add strand information to each mutation type, based on the direction of DNA replication¹² (Fig. 1b). We further condensed the strand-specific signatures into *directional signatures* consisting of 96 mutation types, each assigned either “leading” or “lagging” direction depending on the frequency in the strand-specific signature (Fig. 1c). These directional signatures can be used to separately compute the exposure to the signature on the leading and lagging strand in individual samples (Fig. 1d). We applied this novel algorithm to somatic mutations in 3056 whole-genome sequences across 19 cancer types (Supplementary Table 1), excluding genes from the analysis in order to prevent potential confounding of the results by transcription strand asymmetry^{1,12} or selection. Samples with microsatellite-instability (MSI) and POLE mutations were treated as separate groups, since they are associated with specific mutational processes. In total, we detected 25 mutational signatures that each corresponded to one of the COSMIC signatures¹⁶ and 4 novel signatures, which were primarily found in samples that had not been previously used for signature extraction (myeloid blood, skin, MSI, and ovarian cancers) (Fig. S1, S14–19).

In total, 22 out of 29 signatures exhibited significant replication strand asymmetry or significant correlation with replication timing (Fig. 2a–b, S1). Such widespread replication bias across the mutational landscape is surprising, considering that previous reports documented strand bias for only a few mutational processes such as activity of the APOBEC class of enzymes that selectively edit exposed single-stranded cytosines on the lagging strand^{12,15,17–19}. Our observations confirm that both APOBEC signatures (2 and 13) exhibit clear strand asymmetry, with signature 13 being the most significantly asymmetric signature ($p < 8e^{-100}$). We also observe differences in these signatures with respect to replication timing: signature 2 shows enrichment in late replicating regions, whereas signature 13 appears uncorrelated with replication timing (Fig. 3), which is consistent with previous reports¹⁵. These results validate that our approach is able to correctly identify strand and timing asymmetries of mutagenic processes. Consequently, we next tried to interpret the replication biases we observed in other mutational signatures.

Amongst the better understood mutational mechanisms, several involve replicative processes and DNA repair, such as mismatch-repair deficiency (MMR)²⁰ or mutations in the proofreading domain of Pol ϵ (“POLE-M samples”)^{9,21}. We first analyzed the signatures representing these mechanisms, since they can be directly attributed to a known molecular process. All 5 signatures previously associated with MMR and the novel MSI-linked signature N4 exhibit replication strand asymmetry, generally with enrichment of C>T mutations on the leading strand template and C>A and T>C mutations on the lagging strand template (Fig. 4, S2). It has previously been proposed that the correlation of overall mutation rate with replication timing (as shown in Fig 2b) is a direct result of the activity of MMR²². In contrast, we observed a more complex relationship. Some MMR signatures in MMR deficient patients do not correlate with replication timing (sig. 15, 21, 26) or do so only in one direction (sig. 20), whereas others show clear timing asymmetry (sig. 6 and N4, Fig. S2), indicating that MMR might be only one of several factors influencing mutagenesis in a timing-dependent manner.

Unexpectedly, three MMR signatures (sig. 6, 26, and N4) showed increased exposures around ORIs (Fig. 4, S2–3, S13). Based on experiments in yeast, it has been suggested that MMR is involved in balancing the differences in fidelity of the leading and lagging polymerases¹⁰, in particular repairing errors made by Pol α ¹⁰, which primes the leading strand at ORIs and each lagging strand Okazaki fragment²³ and lacks intrinsic proofreading capabilities²⁴. It has been recently shown that error-prone Pol α -synthesized DNA is retained *in vivo*, causing an increase of mutations on the lagging strand¹¹, and that regions around ORIs have a high density of Pol α -synthesized DNA. It is therefore possible that increased exposure to signatures 6 and 26 around ORIs is caused by incomplete repair of Pol α -induced errors. The most common Pol α -induced mismatches normally repaired by MMR are G-dT and C-dT, leading to C>T mutations on the leading and C>A mutations on the lagging strand²⁵, matching our observations in the MMR-linked signatures. Notably, we also detected weaker but still significant exposure to MMR signatures in samples with seemingly intact mismatch repair (Fig. S3). Replication strand asymmetry in these samples was substantially smaller, but the higher exposure to signatures 6 and 26 around ORIs remained. These findings are compatible with a model in which mismatch repair balances the effect of mis-incorporation of nucleotides by Pol α .

POLE-M samples were previously reported to be “ultra-hypermutated” with excessive C>A and C>T mutations on the leading strand^{9,12,21}. Two mutational signatures (10 and 14) have been associated with Pol ϵ , the main leading strand polymerase^{23,26}. As expected, we observe very strong strand asymmetry for these two signatures in all POLE-M samples, with an increase of C>A, C>T, and T>G mutations on the leading strand (Fig. 4, S4). As with MMR signatures, we also found weak but significant evidence of signature 10 and 14 in samples without Pol ϵ defects (POLE-WT). Strikingly, however, in these samples the strand asymmetry was in the inverse orientation compared to the POLE-M samples, *i.e.*, increased C>A, C>T, and T>G mutations on the lagging strand (Fig. 4, S5, S12). Conversely, we detected presence of the non-POLE signatures 18 and 28 in POLE-M samples, but in the inverse orientation compared to POLE-WT samples. All four signatures (10, 14, 18 and 28) exhibited a stronger correlation with replication timing and distance from ORI in POLE-WT samples than in POLE-M samples. We therefore hypothesize that POLE-linked signatures are originally caused by a process that affects both strands, and under normal circumstances is slightly enriched on the lagging strand. In POLE-M samples the lack of replication-associated proofreading would lead to a strong relative increase in these mutations on the leading strand, explaining the flipped orientation of signatures.

We next focused on signatures that have not previously been reported to be connected to replication, or for which the causal mechanism is unknown. Our data show a link between DNA replication and exogenous mutagens such as UV light (signature 7), tobacco smoke (signature 4) or aristolochic acid (AA) (signature 22)²⁷. In these signatures, we observed marked correlation with replication timing (Fig. 4, S6–7). Higher mutation frequency late in replication has been observed in mouse embryonic fibroblast (MEFs) treated with AA or Benzo[a]pyrene (B[a]P, a mutagen in tobacco smoke)²⁸. This increased mutagenicity might be attributed to differences in DNA damage tolerance between early and late replication. Translesion synthesis (TLS), an error-prone DNA damage tolerance mechanism, has been observed to increase in activity and mutagenicity later in the cell cycle when replicating DNA damaged by B[a]P²⁹, leading to more mutations later during the cell cycle. We also observed weak but significant replication strand asymmetry in the mutagen-induced signatures. This matches a previously observed lower efficiency of bypass of DNA damage on the lagging strand³⁰ and strong mutational strand asymmetry in cells lacking Pol η , the main TLS polymerase responsible for the replication of UV-induced photolesions³¹. Altogether, our data

highlight the importance of replication in converting DNA damage into actual mutations and suggest that bypass of DNA damage occurring on the lagging template results in detectably lower fidelity on this strand.

Signature 17 had the largest median strand asymmetry (p value < $1e^{-59}$) and also is one of the signatures with the strongest correlations with replication timing (Fig. 2, 4). The mutational process causing this signature is unclear. We noted that the timing asymmetry and exposure distribution around ORIs to signature 17 closely resembled that of signatures 4 and 7, suggesting a possible link to DNA damage. Signature 17 is most prominent in gastric cancers and esophageal adenocarcinoma (EAC), where it appears early during disease development³², and it is also present in Barrett's esophagus (BE), a precursor to EAC³³. Due to the importance of gastro-esophageal and duodenogastric reflux in the development of BE and EAC³⁴⁻³⁶ and the resulting oxidative stress³⁷⁻⁴⁰, it has been speculated that oxidative damage could cause the mutation patterns characteristic for Signature 17^{41,42}. Oxidative stress affects not only bases in the DNA, but also the nucleotide pool, such as the oxidation of dGTP to 8-oxo-dGTP. This oxidized dGTP derivative has been shown to induce T>G transversions⁴³⁻⁴⁵ through incorporation by TLS polymerases into DNA opposite A on the template strand⁴⁶. Importantly, the resulting mismatch of 8-oxo-G and A has been shown in yeast to be more efficiently repaired into G:C when 8-oxo-G is on the lagging strand template^{47,48}, resulting in an enrichment of T>G mutations on the lagging strand template. Our data show strong lagging-strand bias of T>G mutations and overall higher exposure to signature 17 on the lagging strand, supporting the hypothesis that signature 17 is a by-product of oxidative damage.

The example of signature 17 demonstrates how the characteristic relationship between mutational signatures and DNA replication can lead to experimentally testable hypotheses and thus help to reveal the mechanisms of many currently unexplained mutational processes. In summary, our results provide evidence that DNA replication interacts with most mutational mechanisms, suggesting that differences among DNA polymerases and repair enzymes might play a larger part in the accumulation of mutations than previously appreciated.

FIGURES

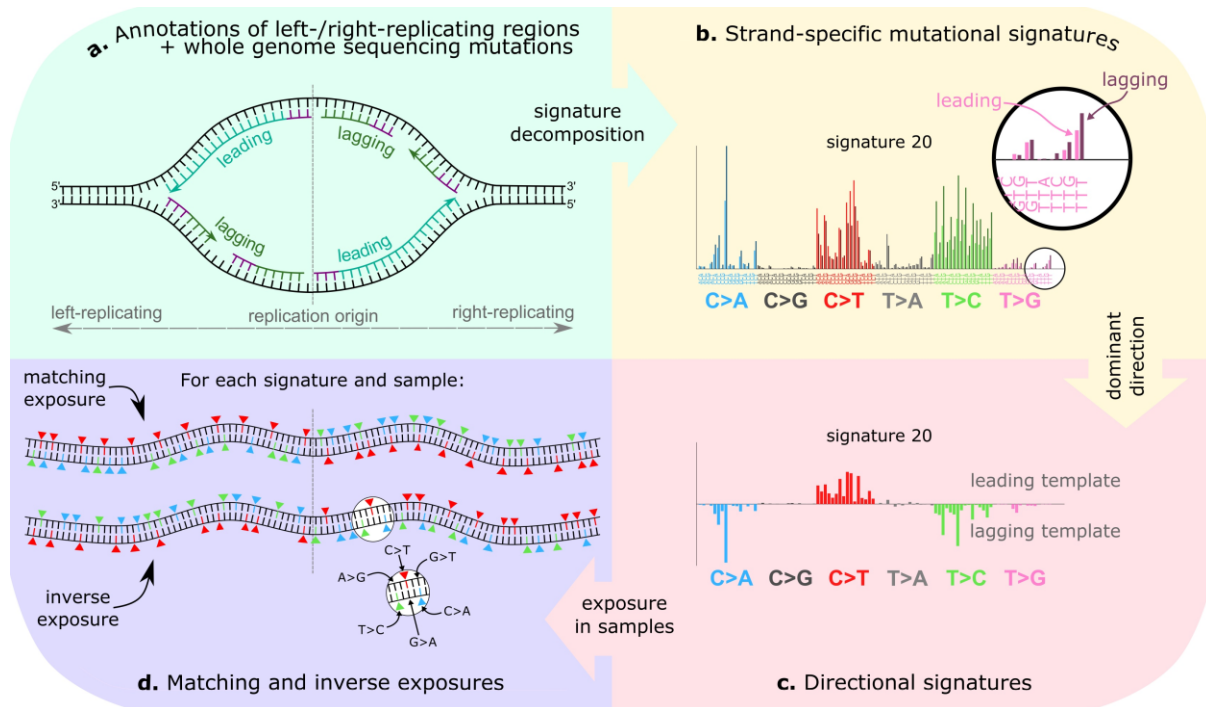


Figure 1: Methods overview. **a.** Mutation frequency on the leading and lagging strand is computed using annotated left/right-replicating regions and somatic single-nucleotide mutations oriented according to the strand of the pyrimidine in the base-pair. **b.** Leading and lagging strand-specific mutational signatures are extracted (signature 20 is shown as an example). **c.** Each of the 96 mutation types is annotated according to its dominant direction (upwards-facing bars for leading, downwards-facing bars for lagging template preference). **d.** Exposures to the directional signatures are separately quantified for the leading and lagging strand of each patient. The exposure in the *matching orientation* reflects the extent to which mutations in pyrimidines on the leading (and lagging) strand can be explained by the leading (and lagging) component of the signature, respectively. Conversely, the exposure in the *inverse orientation* reflects how mutations in pyrimidines on the leading strand can be explained by the lagging component of the signature (or vice-versa) (Methods). Top part of 1d shows an example of a sample with completely matching exposure, given the signature in 1c, with C>T mutations on the leading template and C>A and T>C mutations on the lagging template, whereas bottom part of 1d shows an example of a sample with completely inverse exposure.

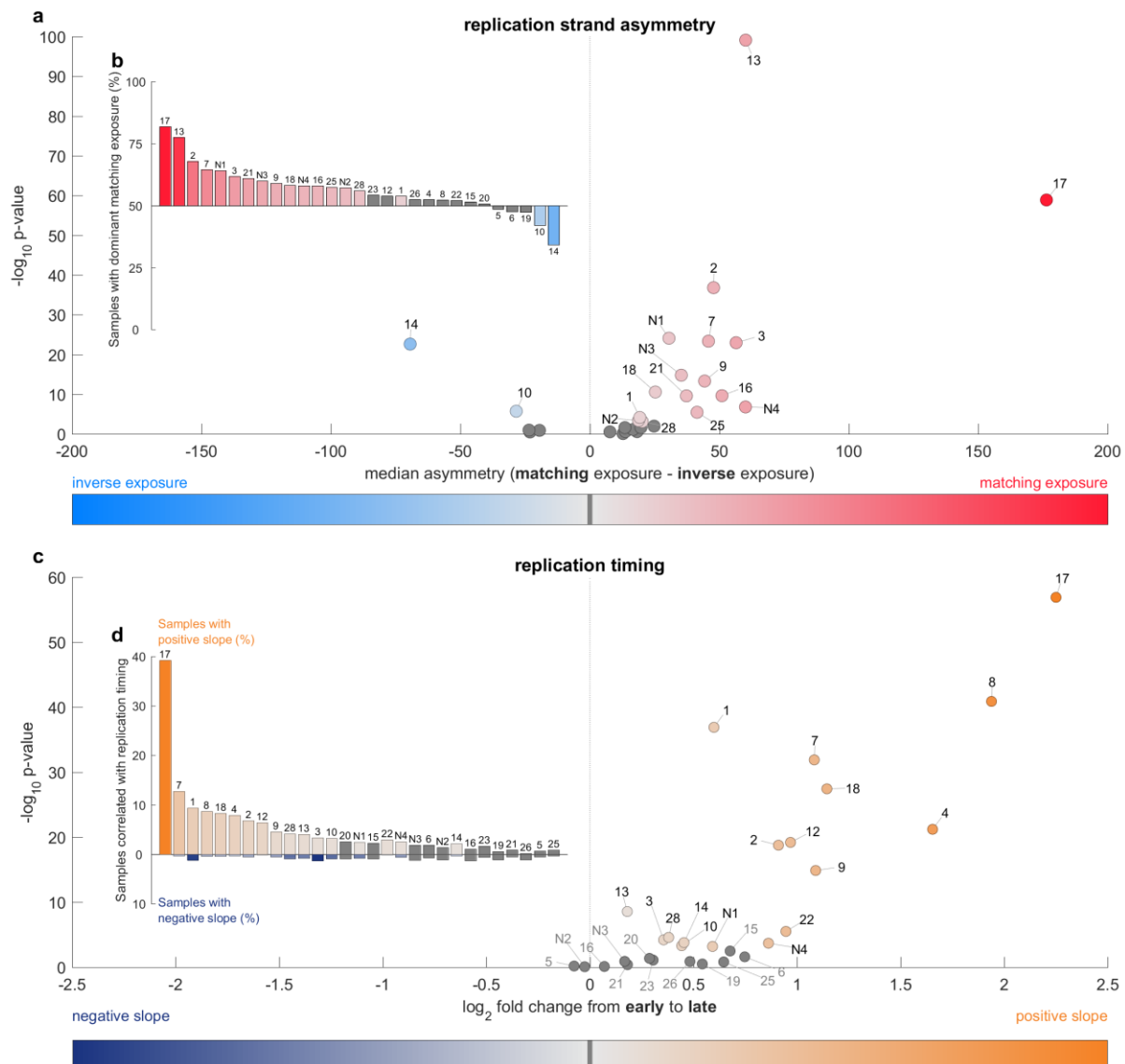


Figure 2: Most mutational signatures exhibit a significant replication strand asymmetry and/or correlation with replication timing. **a.** The difference of matching and inverse exposure is computed for each sample and signature. For each signature, the median value of these differences (in samples exposed to this signature) is plotted against $-\log_{10}$ p-value (signtest of strand asymmetry per sample; with Bonferroni correction). **b.** Percentage of samples that have higher matching than inverse exposure to the signature denoted above/below each bar. **c.** Correlation of exposures with replication timing. A line is fitted to average exposure in four quartiles of replication timing. \log_2 -transformed fold change from average exposure in early (bottom quartile) to late (top quartile) is plotted on the x-axis. The y-axis represents significance of the direction of the slope in individual samples (signtest of slope sign per sample: 0 for non-significant correlation, -1 for negative, 1 for positive; with Bonferroni correction). **d.** Percentage of samples with a significantly positive and negative correlation with exposure, respectively.

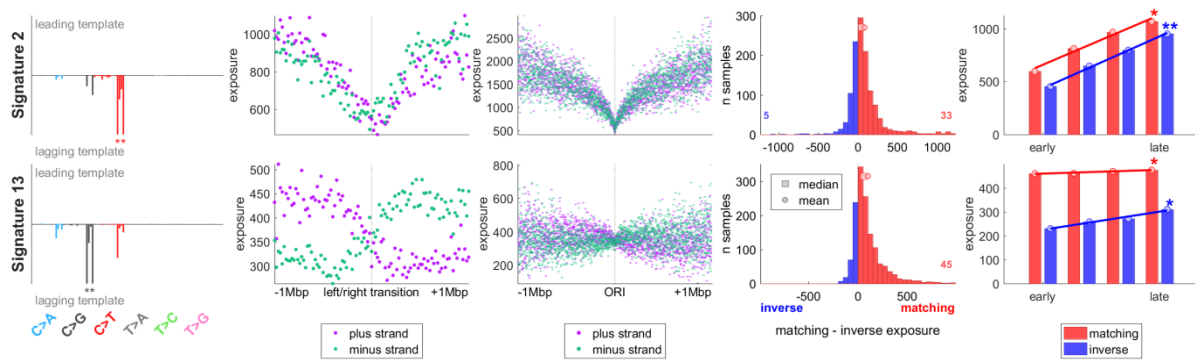


Figure 3: APOBEC signatures show strong but distinct effects of replication. Column 1: directional signatures for the two APOBEC signatures. Column 2: mean exposure on the plus (Watson) and minus (Crick) strand around transitions between left- and right-replicating regions. The transition corresponds to a region enriched for replication origins. Column 3: mean exposure on the plus and minus strand around directly ascertained replication origins. Column 4: distribution of differences between matching and inverse exposure amongst patients with sufficient exposure. Number of outliers is denoted by the small numbers on the sides. Column 5: mean matching and inverse exposure in four quartiles of replication timing; asterisks represent significance of the fit (F-test for coefficient of deviation from 0). The leading and lagging strand annotations used in columns 4 and 5 are based on the direction of replication derived from replication timing data.

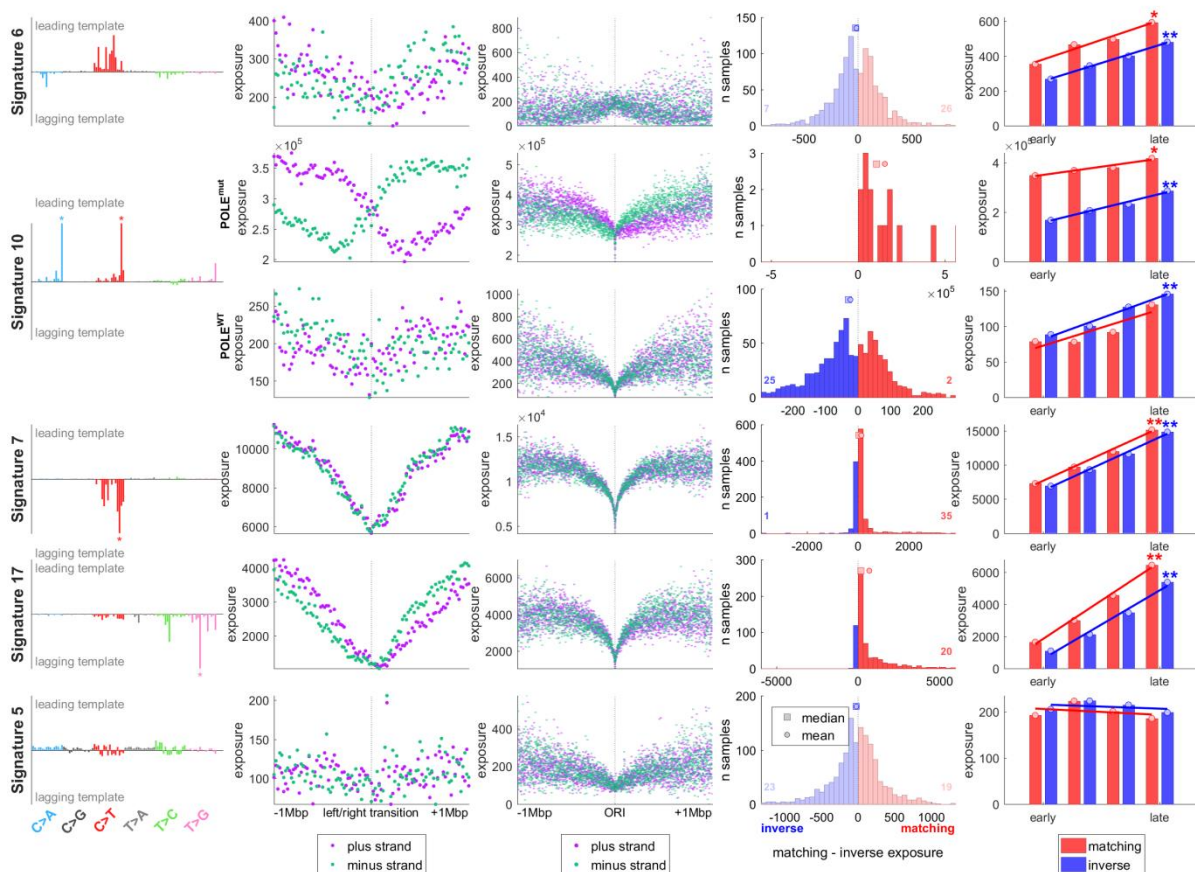


Figure 4: Different mutational signatures exhibit characteristic timing and strand asymmetry profiles. Columns show directional signature (column 1), distribution around timing transition regions (column 2) and around replication origins (column 3), per-patient mutation strand asymmetry (column 4; non-significant asymmetry is shown in light-coloured histogram) and

correlation with replication timing (column 5), as described in Fig 3. Row 1: Signature 6, associated with mismatch-repair deficiency. Row 2–3: signature 10, associated with POLE errors, shown for patients with known POLE mutations (row 2), and those without (row 3). Row 4: signature 7, representing UV-induced damage. Row 5: signature 17, characteristic of gastric and oesophageal cancers. Row 6: Signature 5, of unknown aetiology, is not discernibly affected by replication.

ACKNOWLEDGMENTS

We thank Dr. Mary Muers for comments on the manuscript. S.K. and B.S.-B. are funded by Ludwig Cancer Research. S.K. received funding from BBSRC grant BB/M001873/1. M.T. and J.T. are funded by EPSRC (EP/F500394/1) and Bakala Foundation.

AUTHOR CONTRIBUTIONS

B.S.-B. and M.T. designed the study. M.T. performed the analysis with contributions from J.T. B.S.-B. and M.T. wrote the manuscript with contributions from S.K. and J.T.

REFERENCES

1. Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature* **500**, 415–21 (2013).
2. Secrier, M. *et al.* Mutational signatures in esophageal adenocarcinoma define etiologically distinct subgroups with therapeutic relevance. *Nat. Genet.* **2016**, (2016).
3. Stenzinger, A. *et al.* Mutations in POLE and survival of colorectal cancer patients--link to disease stage and treatment. *Cancer Med.* **3**, 1527–1538 (2014).
4. Tomasetti, C. *et al.* Variation in cancer risk among tissues can be explained by the number of stem cell divisions. *Science* **347**, 78–81 (2015).
5. Lujan, S. A. *et al.* DNA Polymerases Divide the Labor of Genome Replication. *Trends Cell Biol.* **26**, 640–654 (2016).
6. Fragkos, M. *et al.* DNA replication origin activation in space and time. *Nat. Rev. Mol. Cell Biol.* **16**, 360–74 (2015).
7. Stamatoyannopoulos, J. a *et al.* Human mutation rate associated with DNA replication timing. *Nat. Genet.* **41**, 393–395 (2009).
8. Lawrence, M. S. *et al.* Mutational heterogeneity in cancer and the search for new cancer-associated genes. *Nature* **499**, 214–8 (2013).
9. Shinbrot, E. *et al.* Exonuclease mutations in DNA Polymerase epsilon reveal replication strand specific mutation patterns and human origins of replication. *Genome Res.* 1740–1750 (2014).
10. Lujan, S. A. *et al.* Mismatch Repair Balances Leading and Lagging Strand DNA Replication Fidelity. *PLoS Genet.* **8**, (2012).
11. Reijns, M. A. M. *et al.* Lagging-strand replication shapes the mutational landscape of the genome. *Nature* **518**, 502–506 (2015).
12. Haradhvala, N. J. *et al.* Mutational Strand Asymmetries in Cancer Genomes Reveal Mechanisms of DNA Damage and Repair. *Cell* **164**, 538–549 (2016).
13. Besnard, E. *et al.* Unraveling cell type-specific and reprogrammable human replication origin signatures associated with G-quadruplex consensus motifs. *Nat. Struct. Mol. Biol.* **19**, 837–844 (2012).
14. Foulk, M. S. *et al.* Characterizing and controlling intrinsic biases of lambda exonuclease in nascent strand sequencing reveals phasing between nucleosomes and G-quadruplex motifs around a subset of human replication origins. *Genome Res.* **25**, 725–735 (2015).
15. Morganella, S. *et al.* The topography of mutational processes in breast cancer genomes. *Nat. Commun.* **7**, 11383 (2016).
16. Wellcome Trust Sanger Institute. COSMIC: Signatures of Mutational Processes in Human Cancer. Available at: <http://cancer.sanger.ac.uk/cosmic/signatures>.
17. Hoopes, J. I. *et al.* APOBEC3A and APOBEC3B Preferentially Deaminate the Lagging Strand Template during DNA Replication. *Cell Rep.* 1–10 (2016).
18. Green, A. M. *et al.* APOBEC3A damages the cellular genome during DNA replication. *Cell Cycle* **15**, 998–1008 (2016).

19. Seplyarskiy, V. B. *et al.* APOBEC-induced mutations in human cancers are strongly enriched on the lagging DNA strand during replication. *Genome Res.* **26**, 174–182 (2016).
20. Zhao, H. *et al.* Mismatch repair deficiency endows tumors with a unique mutation signature and sensitivity to DNA double-strand breaks. *Elife* **3**, e02725 (2014).
21. Shlien, A. *et al.* Combined hereditary and somatic mutations of replication error repair genes result in rapid onset of ultra-hypermuted cancers. *Nat. Genet.* **47**, 257–262 (2015).
22. Supek, F. *et al.* Differential DNA mismatch repair underlies mutation rate variation across the human genome. *Nature* **521**, 81–84 (2015).
23. Stillman, B. DNA Polymerases at the Replication Fork in Eukaryotes. *Mol. Cell* **30**, 259–260 (2008).
24. McCulloch, S. D. *et al.* The fidelity of DNA synthesis by eukaryotic replicative and translesion synthesis polymerases. *Cell Res.* **18**, 148–161 (2008).
25. Nick McElhinny, S. a *et al.* Differential correction of lagging-strand replication errors made by DNA polymerases α and δ . *Proc. Natl. Acad. Sci. U. S. A.* **107**, 21070–21075 (2010).
26. Georgescu, R. E. *et al.* Reconstitution of a eukaryotic replisome reveals suppression mechanisms that define leading/lagging strand operation. *Elife* **2015**, 1–20 (2015).
27. Helleday, T. *et al.* Mechanisms underlying mutational signatures in human cancers. *Nat. Rev. Genet.* **15**, 585–598 (2014).
28. Nik-zainal, S. *et al.* The genome as a record of environmental exposure. *Mutagenesis* **30**, 763–770 (2015).
29. Diamant, N. *et al.* DNA damage bypass operates in the S and G2 phases of the cell cycle and exhibits differential mutagenicity. *Nucleic Acids Res.* **40**, 170–180 (2012).
30. Cordeiro-Stone, M. *et al.* Asymmetry of DNA replication and translesion synthesis of UV-induced thymine dimers. *Mutat. Res.* **510**, 91–106 (2002).
31. McGregor, W. G. *et al.* Abnormal, Error-Prone Bypass of Photoproducts by Xeroderma Pigmentosum Variant Cell Extracts Results in Extreme Strand Bias for the Kinds of Mutations Induced by UV Light. *Mol. Cell. Biol.* **19**, 147–154 (1999).
32. Murugaesu, N. *et al.* Tracking the genomic evolution of esophageal adenocarcinoma through neoadjuvant chemotherapy. *Cancer Discov.* **5**, 821–832 (2015).
33. Ross-Innes, C. S. *et al.* Whole-genome sequencing provides new insights into the clonal architecture of Barrett’s esophagus and esophageal adenocarcinoma. *Nat. Genet.* **47**, 1–11 (2015).
34. Souza, R. F. The role of acid and bile reflux in oesophagitis and Barrett’s metaplasia. *Biochem. Soc. Trans.* **38**, 348–52 (2010).
35. Erichsen, R. *et al.* Erosive Reflux Disease Increases Risk for Esophageal Adenocarcinoma, Compared With Nonerosive Reflux. *Clin. Gastroenterol. Hepatol.* **10**, 475–480.e1 (2012).
36. Fein, M. *et al.* Importance of duodenogastric reflux in gastro-oesophageal reflux disease. *Br. J. Surg.* **93**, 1475–1482 (2006).
37. Kauppi, J. *et al.* Increased Oxidative Stress in the Proximal Stomach of Patients with Barrett’s Esophagus and Adenocarcinoma of the Esophagus and Esophagogastric Junction. *Transl.*

- Oncol.* **9**, 336–339 (2016).
38. Rasanen, J. V. *et al.* The expression of 8-hydroxydeoxyguanosine in oesophageal tissues and tumours. *Eur. J. Surg. Oncol.* **33**, 1164–1168 (2007).
 39. Jimenez, P. *et al.* Free radicals and antioxidant systems in reflux esophagitis and Barrett's esophagus. *World J. Gastroenterol.* **11**, 2697–2703 (2005).
 40. Dvorak, K. *et al.* Bile acids in combination with low pH induce oxidative stress and oxidative DNA damage: relevance to the pathogenesis of Barrett's oesophagus. *Gut* **56**, 763–771 (2007).
 41. Dulak, A. M. *et al.* Exome and whole-genome sequencing of esophageal adenocarcinoma identifies recurrent driver events and mutational complexity. *Nat. Genet.* **45**, 478–86 (2013).
 42. Nones, K. *et al.* Genomic catastrophes frequently arise in esophageal adenocarcinoma and drive tumorigenesis. *Nat. Commun.* **5**, 1–9 (2015).
 43. Inoue, M. *et al.* Induction of chromosomal gene mutations in *Escherichia coli* by direct incorporation of oxidatively damaged nucleotides: New evaluation method for mutagenesis by damaged dna precursors in vivo. *J. Biol. Chem.* **273**, 11069–11074 (1998).
 44. Satou, K. *et al.* Mutagenic effects of 8-hydroxy-dGTP in live mammalian cells. *Free Radic. Biol. Med.* **42**, 1552–1560 (2007).
 45. Satou, K. *et al.* Involvement of specialized DNA polymerases in mutagenesis by 8-hydroxy-dGTP in human cells. *DNA Repair (Amst)*. **8**, 637–642 (2009).
 46. Kamiya, H. Mutations Induced by Oxidized DNA Precursors and Their Prevention by Nucleotide Pool Sanitization Enzymes. *Genes Environ.* **29**, 133–140 (2007).
 47. Pavlov, Y. I. *et al.* Evidence for Preferential Mismatch Repair of Lagging Strand DNA Replication Errors in Yeast. *Curr. Biol.* **13**, 744–748 (2003).
 48. Mudrak, S. V *et al.* The polymerase eta translesion synthesis DNA polymerase acts independently of the mismatch repair system to limit mutagenesis caused by 7,8-dihydro-8-oxoguanine in yeast. *Mol. Cell. Biol.* **29**, 5316–26 (2009).
 49. Bass, A. J. *et al.* Genomic sequencing of colorectal adenocarcinomas identifies a recurrent VTI1A-TCF7L2 fusion. *Nat. Genet.* **43**, 964–8 (2011).
 50. Wang, K. *et al.* Whole-genome sequencing and comprehensive molecular profiling identify new driver mutations in gastric cancer. *Nat. Genet.* **46**, 573–82 (2014).
 51. Saunders, C. T. *et al.* Strelka: Accurate somatic small-variant calling from sequenced tumor-normal sample pairs. *Bioinformatics* **28**, 1811–1817 (2012).
 52. Encode Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
 53. Rosenthal, R. *et al.* deconstructSigs: delineating mutational processes in single tumors distinguishes DNA repair deficiencies and patterns of carcinoma evolution. *Genome Biol.* **17**, 31 (2016).

METHODS

Somatic mutations. Cancer somatic mutations in 3056 whole-genome sequencing samples (Supplementary Table 1) were obtained from the data portal of The Cancer Genome Atlas (TCGA), the data portal of the International Cancer Genome Consortium (ICGC), and previously published data in peer-review journals^{1,21,41,49,50}. For the TCGA samples, aligned reads of paired tumor and normal samples were downloaded from the UCSC CGHub website under TCGA access request #10140 and somatic variants were called using Strelka (version 1.0.14)⁵¹ with default parameters.

Direction of replication. Left- and right-replicating domains were taken from¹². Each domain (called territory in the original source code and data) is 20kbp wide and annotated with the direction of replication and with replication timing.

Excluded regions. The following regions were excluded: regions with low unique mappability of sequencing reads (positions with mean mappability in 100bp sliding windows below 0.99 from UCSC mappability track), gencode protein coding genes, and blacklisted regions defined by Anshul Kundaje⁵² (Anshul_Hg19UltraHighSignalArtifactRegions.bed, Duke_Hg19SignalRepeatArtifactRegions.bed, and wgEncodeHg19ConsensusSignalArtifactRegions.bed from <http://mitra.stanford.edu/kundaje/akundaje/release/blacklists/hg19-human/>).

Mutation frequency analysis. All variants were classified by the pyrimidine of the mutated Watson-Crick base pair (C or T), strand of this base pair (C or T), and the immediate 5' and 3' sequence context into 96 possible mutation types as described by Alexandrov *et al.*¹. The frequency of trinucleotides on each strand was computed for each replication domain. Then the mutation frequency of each mutation type in each replication domain on the leading (plus=Watson strand in left replicating domains; minus=Crick strand in right replicating domains) and lagging strand (vice versa) was computed for each sample.

Extraction of mutational signatures. Matlab code¹ was used for extraction of strand-specific mutational signatures. The input data were the mutation counts on the leading and lagging strands (summed from all replicating domains together, but without the excluded regions) in each sample. The 192-elements-long mutational signatures (example in Fig. 1b) were extracted in each cancer type separately (for K number of signatures between 2 and 7). The best K with minimal error and maximal stability (minimizing $\text{error}_K / \max(\text{error}) + (1 - \text{stability}_K)$ and with stability of at least 0.8) was selected for each cancer type. Signatures present in only a small number of samples with very low exposures were excluded ($(95^{\text{th}} \text{ percentile of exposures of this signature}) / (\text{mean total exposure per samples}) < 0.2$). The remaining signatures were then normalized by the frequency of trinucleotides in the leading and lagging strand and subsequently multiplied by the frequency of trinucleotides in the genome. This made them comparable with the 30 previously identified whole-genome-based COSMIC signatures (<http://cancer.sanger.ac.uk/cosmic/signatures>). Signatures extracted in each cancer type and COSMIC signatures were all pooled together (with equal values in the leading and lagging part in the COSMIC signatures) and were clustered using unsupervised hierarchical clustering (with cosine distance and complete linkage). A threshold was selected to identify clusters of similar signatures. Mis-clustering was avoided by manual examination (and whenever necessary re-assignment) of all signatures in all clusters. The resulting 29 signatures (representing the detected clusters) contained 25 previously observed (COSMIC) and 4 new signatures. For the subsequent analysis, the signatures were converted back to 96 values: the 25 previously observed signatures

were used in their original form and average of the leading and lagging part were used for the 4 newly identified signatures.

Annotation of signatures with leading and lagging direction. Each signature was annotated with strand direction: which of the 96 mutation types were higher on the leading strand and which on the lagging strand (Fig. 1c). This was based on the dominant strand direction within the signature's cluster. Types with unclear direction and small values were assigned according to the predominant direction of other trinucleotides of the same mutation group, such as C>T.

Calculating strand-specific exposures in individual samples. Exposures to leading and lagging parts of the signatures on the leading and lagging strands in individual samples were quantified using non-negative least squares regression using the Matlab function $e = lsqnonneg(S, m)$, where

$$S = \begin{pmatrix} S_{LD} & S_{LG} \\ S_{LG} & S_{LD} \end{pmatrix}, m = \begin{pmatrix} m_{LD} \\ m_{LG} \end{pmatrix}, e = \begin{pmatrix} e_{matching} \\ e_{inverse} \end{pmatrix}.$$

The matrix S_{LD} has 96 rows and 29 columns and represents the leading parts of the signatures, *i.e.* the elements of the lagging parts contain zeros in this matrix. Similarly, S_{LG} has the same size, but contains zeros in the leading parts. The vector m_{LD} of length 96 contains mutations on the leading strand (again normalized by trinucleotides in leading strand/whole genome), and similarly m_{LG} contains mutations from the lagging strand. Finally, *lsqnonneg* finds a non-negative vector of exposures e such that it minimizes a function $|m - C \cdot e|$. A similar approach has been used in⁵³ for finding exposures to a given set of signatures. Our extension includes the strand-specificity of the signatures. The interpretation of the model is that the *matching exposure* $e_{matching}$ represents exposure of the leading part of the signature on the leading strand and exposure of the lagging part of the signature on the lagging strand, whereas $e_{inverse}$ represents the two remaining options. It is important to note that the direction of the mutation is relative to the nucleotide in the base pair chosen as the reference, *i.e.*, mutations of a pyrimidine on the leading strand correspond to mutations of a purine on the lagging strand. In order to minimize the number of spurious signature exposures, the least exposed signature was incrementally removed (in both leading and lagging parts) while the resulting error did not exceed the original error by 0.5%. The resulting reported values in each sample and signature were the difference (or fold change) of $e_{matching}$ and $e_{inverse}$. In each signature, the signtest was used to compare matching and inverse exposures across samples with sufficient minimal exposure (at least 10) to the signature. Bonferroni correction was applied to correct for multiple testing.

Replication origins. The left/right transitions of the replication domains represent regions with on average higher density of replication origins. In order to get better resolution of the replication origins, and to validate the results using an independent estimates of left- and right-replicating domains, genome-wide maps of human replication origins from NS-seq by¹³ were used. Eight fastq files (HeLa, iPS, hESC, IMR; each with two replicates) were downloaded and mapped to hg19 using bowtie2 (version 2.1.0). To control for the inefficient digestion of λ -exo step of NS-seq, reads from non-replicating genomic DNA (LexoG0) were used as a control¹⁴. Peaks were called using "macs callpeak" with parameters --gsize=hs --bw=200 --qvalue=0.05 --mfold 5 50 and LexoG0 mapped reads as a control. Only peaks covered in at least seven of the eight samples were used. 1000 1kbp bins were generated to the left and right of each origin, as long as they did not reach half the distance to the next origin. We then used these replication direction annotations in the 1kbp bins to calculate strand-specific exposures in individual samples as above and ascertained that both

approaches lead to qualitatively very similar mutational strand asymmetries in individual signatures (Fig. S20).

Quantification of exposures with respect to replication timing, left/right transitions, and replication origins. Replication domains were divided into four quartiles by their average replication timing. The entire exposure quantification was computed separately in each quartile, or bin around left/right transition or bin around replication origin. In replication timing plots, a linear regression model (function `fitlm` in MatLab) was fitted to the mean exposure in each quartile (separately for matching and inverse exposures) and the significance of the linear coefficient was tested using F-test for the hypothesis that the regression coefficient is zero (function `coefTest` in MatLab).