

**Deficiency in DNA mismatch repair of methylation damage is a major
mutational process in cancer**

Hu Fang¹, Xiaoqiang Zhu¹, Jieun Oh¹, Jayne A. Barbour¹, Jason W. H. Wong^{1,*}

¹School of Biomedical Sciences, Li Ka Shing Faculty of Medicine, The University of
Hong Kong, Hong Kong Special Administrative Region

*Correspondence: jwhwong@hku.hk

1 **Abstract**

2 DNA mismatch repair (MMR) is essential for maintaining genome integrity with its
3 deficiency predisposing to cancer¹. MMR is well known for its role in the post-
4 replicative repair of mismatched base pairs that escape proofreading by DNA
5 polymerases following cell division². Yet, cancer genome sequencing has revealed that
6 MMR deficient cancers not only have high mutation burden but also harbour multiple
7 mutational signatures³, suggesting that MMR has pleiotropic effects on DNA repair. The
8 mechanisms underlying these mutational signatures have remained unclear despite
9 studies using a range of *in vitro*^{4,5} and *in vivo*⁶ models of MMR deficiency. Here, using
10 mutation data from cancer genomes, we identify a previously unknown function of
11 MMR, showing that the loss of non-canonical replication-independent MMR activity
12 is a major mutational process in human cancers. MMR is comprised of the MutS α
13 (MSH2/MSH6) and MutL α (MLH1/PMS2) complexes⁷. Cancers with deficiency of
14 MutS α exhibit mutational signature contributions distinct from those deficient of
15 MutL α . This disparity is attributed to mutations arising from the unrepaired
16 deamination of 5-methylcytosine (5mC), i.e. methylation damage, as opposed to
17 replicative errors by DNA polymerases induced mismatches. Repair of methylation
18 damage is strongly associated with H3K36me3 chromatin but independent of binding
19 of MBD4, a DNA glycosylase that recognise 5mC and can repair methylation damage.
20 As H3K36me3 recruits MutS α , our results suggest that MutS α is the essential factor in
21 mediating the repair of methylation damage. Cell line models of MMR deficiency
22 display little evidence of 5mC deamination-induced mutations as their rapid rate of

1 proliferation limits for the opportunity for methylation damage. We thus uncover a non-
 2 canonical role of MMR in the protection against methylation damage in non-dividing
 3 cells.

4

5 Keywords: Mismatch repair; Mutational signature; Microsatellite instability cancers;

6 H3K36me3; DNA methylation

7

1 Introduction

2 DNA mismatch repair (MMR) is a highly conserved DNA repair pathway, well
3 known for its ability to recognise and remove errors from the newly synthesised DNA
4 strand during replication⁷. The general mechanism of MMR in the correction of DNA
5 replication errors in humans is well established. The process is initiated by the MutS α
6 heterodimer, comprised of MSH2 and MSH6, which recognise mis-incorporated bases
7 in double stranded DNA behind the replication fork. Subsequently, the MutL α
8 heterodimer, which is comprised of MLH1 and PMS2, is recruited to excise the
9 sequence surrounding the mutated base. The mismatched section of the daughter strand
10 is digested by EXO1 exonuclease and the gap filled by a DNA polymerase⁷. MMR is
11 also known to have non-canonical roles outside the context of DNA replication⁸. One
12 of the better-known functions of non-canonical MMR (ncMMR) is facilitating somatic
13 hyper-mutation of the immunoglobulin locus in lymphoid cells⁹. NcMMR has also been
14 shown to be activated by DNA lesions resulting in an error-prone repair process that
15 leads to the formation of A:T mutations¹⁰⁻¹². While these studies have found that
16 ncMMR is generally associated with increased mutagenesis through the recruitment of
17 error-prone DNA polymerase eta (POLH), more recently, it has also emerged that
18 ncMMR is capable of protecting actively transcribed genes by removing DNA lesions
19 in a transcription dependent manner¹³. The mechanistic details of how ncMMR
20 achieves error-free repair remains unclear, but there is also evidence that ncMMR can
21 facilitate active DNA demethylation¹⁴, suggesting that there are multiple ncMMR
22 pathways and a pathway for high-fidelity ncMMR may exist.

The use of somatic mutational signatures has contributed significantly to our understanding of the underlying mutational processes in cancer¹⁵. Currently seven single based substitution (SBS) mutational signatures have been associated with MMRd across various cancer types³. MMRd mutational signatures, SBS6 and SBS15 are characterised by high frequency of C>T mutations, SBS21 and SBS26 have a dominant mutation spectrum of T>C, while SBS44 has relatively high contributions from C>A, C>T and T>C mutations. SBS14 and SBS20 are associated with MMRd concurrent with polymerase epsilon (POLE) and polymerase delta (POLD1) exonuclease mutations, respectively. Recently, it was shown that two mutational processes largely underlie the MMRd specific mutational signatures¹⁶. Both of these mutational processes are believed to be associated with DNA polymerase associated replication errors, yet, only one has been reproducible *in vitro* using clonal models of MMRd cell lines^{5,16}. Thus, the mechanisms underlying the mutational processes and the mutational signatures observed in human MMRd cancers remain unclear.

In this study, we provide evidence that ncMMR is required for the repair of 5-methylcytosine (5mC) deamination induced G:T mismatches (we will refer to this as 5mC deamination damage) outside of the context of DNA replication. We show that mutations arising from unrepaired 5mC deamination events are prevalent in MMRd cancers, particularly those deficient of MutS α . Therefore, the deficiency of ncMMR activity is a major mutational process in MMRd cancers.

Results

1 **Profile of mismatch repair genes in MSI-H tumors**

2 We first investigated three tumor types (CRC, colorectal cancer; STAD, stomach
3 adenocarcinoma; UCEC, uterine corpus endometrial carcinoma) as they have been
4 recognised as MSI prone and contain the majority of MSI events¹⁷. A set of 316 cancer
5 samples with MSI-High (MSI-H) phenotype was obtained from TCGA Pan-cancer
6 cohort for these tumor types. In order to focus on samples where the major mutational
7 processes is MMRd, we excluded samples with concurrent POLE and POLD1
8 exonuclease domain mutations as defined by high contributions from SBS14 and
9 SBS20, respectively¹⁸. The remaining samples were stratified into MutL α and MutS α
10 mutants after careful review of the underlying mutations, RNA expression and DNA
11 methylation status of the four canonical mismatch repair genes (*MLH1*, *PMS2*, *MSH2*
12 and *MSH6*) (see Methods, **Extended Data Fig 1, Supp Table 1**). In the end, 197 MutL α
13 and 18 MutS α deficient cancer samples were identified. In addition, 14 samples were
14 found with both MutL α and MutS α defects and 37 samples were unable to be
15 conclusively classified based on the available data.

16 As expected, all the MSI-H samples showed high mutation load for both indels and
17 single base substitutions with high proportions of C>T and T>C (**Fig 1A**). *MLH1*
18 promoter methylation (defined as beta value > 0.25) was observed in 78% (149/191) of
19 MSI-H samples with available methylation data and this is in line with previous
20 studies^{19,20}. For the curated MutL α deficient cancers, 93.4% (184/197) had aberrant
21 expression of *MLH1* while the majority with MutS α deficiency harboured truncating
22 mutations in either *MSH2* or *MSH6* (**Fig 1A**). In order to compare the mutation

spectrum for MutL α and MutS α deficient samples, principal component analysis was performed based on trinucleotide context point mutation frequencies. Samples with MutS α deficiency were clustered together with a relatively high frequency of C>T mutations, while MutL α deficient cancers were more distributed with a broader mutation spectrum (**Fig 1B**). These results suggest that there may be differences underlying the mutational process of MutL α and MutS α deficiency.

MutS α and MutL α deficient cancers display differential mutational signatures

To determine if MutS α and MutL α deficient cancers have different mutational processes, we used Sigfit²¹ to assign the somatic mutations of each sample to the five MMRd associated SBS mutational signatures along with the age-associated SBS1, which is present in most cancers. SBS1 contributed to a surprisingly high proportion of mutations in many MMR samples (**Fig 2A**), however, this may reflect difficulty in resolving SBS1 and SBS6 as the two signatures show a high degree of similarity. Generally, MutS α deficient cancers had the highest contribution from SBS1+SBS6 (**Fig 2A**). To simplify the representation of mutational processes in MMRd cancers, we adopted the use of two signatures proposed by Nemeth et al.¹⁶. Using the non-negative matrix factorization (NMF) algorithm, two *de novo* signatures were decomposed from the mutations of the MMRd samples. In line with Nemeth et al.¹⁶, two distinctive *de novo* signatures were obtained with signature A (SigA), characterised by a high frequency of C>T mutations, particularly in CpG sites while signature B (SigB) showed a broader spectrum with C>A, C>T and T>C mutations (**Fig 2B**). We then calculated

1 the cosine similarity between the reconstructed spectrum derived from these two
2 signatures and the real spectrum for each sample genome. Most of the samples had high
3 cosine similarity above 0.85 (**Extended Data Fig 2A**). Next, we compared the newly
4 decomposed signatures with the previously reported MMRd related signatures³. SigA
5 showed relatively high cosine similarity with SBS6 and SBS15 while SigB was more
6 alike with SBS21 and SBS26 (**Extended Data Fig 2B**).

7 As we found that the mutation spectrum of MutS α deficient samples had generally
8 higher proportion of C>T mutations and SBS1+SBS6 relative to MutL α deficient
9 samples (**Fig 1A-B and Fig 2A**), we sought to examine the contribution of the two *de*
10 *novo* signatures in each MMR sample. Samples with MutS α deficiency had
11 significantly higher contribution of SigA relative to those with MutL α deficiency ($p <$
12 0.001, Student's t-test, **Fig 2C**), with samples deficient of both complexes having SigA
13 contribution less than MutS α but not significantly different to MutL α deficient samples
14 ($p < 0.01$ and $p = 0.41$, respectively, Student's t-test, **Extended Data Fig 2C**). To
15 exclude any potential cancer specific effect, we examined the signature contribution in
16 CRC, STAD and UCEC separately, and found that SigA is significantly more enriched
17 in MutS α compared with MutL α deficient samples across all cancer types ($p < 0.01$,
18 Student's t-test, **Extended Data Fig 2D-F**). We next expanded our data to three
19 independent cohorts to validate this observation. Due to limited availability of data
20 types for these cohorts, the approach to classify MutS α and MutL α deficiency status is
21 slightly different from the TCGA dataset (see Methods). MSK-CRC²² and MSK-
22 UCEC²³ cohorts contain 99 CRC and 22 UCEC MSI-H samples, respectively, with only

1 targeted sequencing data available. After fitting the mutations from the samples to SigA
2 and SigB, both the cohorts showed significant enrichment of SigA in MutSα compared
3 with MutLα deficient samples ($p < 0.001$ and $p = 0.04$ for CRC and UCEC respectively,
4 Student's t-test, **Fig 2D and E**). Furthermore, the analysis was also performed on the
5 Depmap²⁴ cohort comprising of 99 MMRd cell lines across 16 cancer types. Again,
6 MutSα deficient samples had significant enrichment in SigA compared with MutLα
7 deficient samples ($p = 0.001$, Student's t-test, **Fig 2F**). Finally, as complex MSH2 and
8 MSH6 mutations are a frequent mechanism of MSI in prostate cancer²⁵, we identified
9 a further 4 MutSα mutant samples in the TCGA prostate adenocarcinoma (PRAD)
10 cohort and confirmed them all to have high SigA contribution (>0.845 , **Extended Data**
11 **Fig 2G**). Together, these results suggest that these two MMR associated mutational
12 processes contribute to different extent to the overall mutation spectrum of MutSα and
13 MutLα deficient cancers.

14

CpG C>T mutations in MutSα mutants show no replication strand bias compared with non-CpG C>T mutations

17 Mutation density varies across cancer genomes²⁶. Due to differential MMR efficiency,
18 strong correlation is found between DNA replication timing and mutation density where
19 late replicating regions have higher mutation density compared with early replicating
20 regions of the genome²⁷. In MMRd cancers, the association of mutation density and
21 replication timing becomes less apparent for MMR dependent mutational processes.
22 Given that we found different contributions of mutational processes in MutSα and

1 MutL α deficient cancers, we sought to determine how the processes are influenced by
 2 replication timing. As expected, higher mutation density was observed in late
 3 replicating regions compared with early replicating regions in MMR proficient
 4 microsatellite stable (MSS) cancers but the difference was reduced in MutS α and
 5 MutL α deficient MSI samples (**Extended Data Fig 3A**). We observed significant
 6 difference in the dependence of the mutation load and replication time between MutS α
 7 and MutL α (**Extended Data Fig 3B**) and ascribe this difference to CpG C>T mutations
 8 but not non-CpG C>T or T>C mutations (**Extended Data Fig 3C-O**). These data
 9 suggest that there may be differences in the dependence of CpG C>T mutations on
 10 replication compared with other types of substitution mutations.

11 To further examine the relationship between mutation formation in MMRd cancers
 12 and DNA replication, we next investigated the replication strand bias of mutations.
 13 Mutations in MMRd samples have generally been attributed to unrepaired errors that
 14 have escaped from polymerase proofreading during DNA replication²⁸. Due to the
 15 differential fidelity of polymerases during DNA synthesis, mutation load in the leading
 16 strand and lagging strand are expected to be asymmetric²⁹. We examined the replication
 17 asymmetry of C>T/G>A and A>G/T>C mutations in the leading and lagging strands.
 18 Both mutation types showed strand bias with the leading strand, generating more C>T
 19 and A>G ($p < 0.05$ for both, Chi-squared test, **Fig 3A**). The lack of strand bias for
 20 exome-wide simulated mutations confirmed that this bias is related to replication rather
 21 than sequence composition (**Fig 3B**). We next compared the strand bias for each
 22 individual MutS α and MutL α deficient samples. While there was no difference in

1 replication bias in A>G/T>C mutations between the MutL α and MutS α deficient
2 samples ($p = 0.849$, Student's t-test, **Fig 3C**), the bias of C>T/G>A mutation is
3 significantly different, with MutS α showing less asymmetry than MutL α deficient
4 samples ($p=0.002$, Student's t-test, **Fig 3D**). Further, we compared CpG C>T and non-
5 CpG C>T replication bias individually for MutS α and MutL α deficient samples. CpG
6 C>T mutations showed significantly less bias compared with non-CpG C>T mutations
7 ($p<0.001$ for both MutS α and MutL α , Student's t-test, **Fig 3E**). We further observed
8 significant CpG C>T strand bias difference between MutS α and MutL α deficient
9 samples ($P<0.001$, Student's t-test) while there was no significant difference for non-
10 CpG C>T mutations (**Fig 3E**). Although there was insufficient whole genome
11 sequenced (WGS) samples with MutS α deficiency ($n = 1$), with increased number of
12 mutations we were able to assess strand bias across all mutation types (**Extended Data**
13 **Fig 4A-F**) and consistent results were observed in the strand bias difference between
14 CpG C>T and non-CpG C>T mutations (MutL α , $p < 0.001$, Student's t-test, **Extended**
15 **Data Fig 4G**). Reduced strand bias of CpG C>T mutations was also observed in WGS
16 MSS samples ($p < 0.001$, **Extended Data Fig 4H**). Furthermore, as some MutL α
17 deficient samples had relatively high SigA contribution, we directly correlated SigA
18 contribution with the degree of CpG C>T replication bias. Negative correlation was
19 observed with samples with high proportions of SigA showing less bias ($R = -0.29$, $p <$
20 0.0001 , **Extended Data Fig 4I**). Together these results suggest that CpG C>T mutations
21 associated with SigA in human MMRd cancers may not have arisen from unrepaired
22 DNA polymerase errors as they do not exhibit the characteristic replication asymmetry

found for other types of substitution mutations.

Most CpG C>T mutations in MMRd cancers are caused by the deamination of 5-methylcytosine

5mC has the tendency to undergo spontaneous deamination into thymine and is a major mutagenic process in the human genome³⁰. Methyl-binding domain 4 (MBD4) is a key DNA glycosylase responsible for the removal of 5mC deamination damage³¹. Cancer patients with biallelic germline MBD4 deficiency present with extremely high frequency of CpG C>T mutations, providing a model for mutations induced by 5mC deamination³². CpG C>T mutations in MBD4 deficient cancers are likely to arise outside the context of DNA replication thus replication asymmetry would not be expected. To test if this is the case, we obtained somatic mutations from whole genome sequenced MBD4-deficient cancers³² and compared the CpG C>T replication strand bias with MMRd mutants. We also included POLE exonuclease domain mutant cancers as their somatic mutations are known to be predominantly leading strand biased³³ and it has been proposed that POLE mutants are particularly erroneous when replicating 5mC³⁴. MBD4 mutants show little strand bias for all mutation types while there was substantial strand bias for POLE mutants (**Fig 4A**), an observation not present in simulated mutations (**Fig 4B**). The strand bias of CpG C>T mutations was close to zero for MBD4 mutants while, as expected, POLE mutants presented strong leading strand bias for both CpG and Non-CpG C>T mutations (**Fig 4C-D**). Although CpG C>T mutations for POLE mutants were highly enriched in TCG trinucleotide context, we

observed consistent CpG C>T mutation strand bias in other contexts (**Extended Data Fig 5**). These results suggest that the pattern of MMRd CpG C>T mutations is more similar to MBD4 mutants where the mutations arise from the deamination of 5mC and are less dependent on process of DNA replication.

As CpG C>T mutation rate increases with methylation level in MMRd samples³⁵, we compared replication strand bias of CpG C>T mutations in lowly and highly methylated sites in MMRd and POLE mutant cancers. Mutations at highly methylated regions showed significantly less strand bias compared with lowly methylated sites for MMRd samples ($p = 0.0233$, Student's t-test, **Fig 4E**) while POLE mutants presented comparably strong strand bias at both lowly and highly methylated regions ($p = 0.80$, Student's t-test, **Fig 4F**). This further supports our hypothesis that CpG C>T mutations in MMR deficient cancers arise from replication independent deamination of 5mC.

MMR repairs 5-methylcytosine deamination damage in a H3K36me3 dependent manner

The histone mark H3K36me3 is an important epigenetic modification involved in the recruitment of MMR to chromatin³⁶. One of the hallmarks of H3K36me3 dependent MMR is differential repair of exons and introns where exons show significantly less mutations than expected due to increased H3K36me3 and MMR activity compared with introns³⁷. To examine if MMR might play a role in the repair of 5mC deamination damage, we compared the observed and expected CpG C>T mutation densities in exons and introns in MBD4 deficient cancers and compared this to MMRd (i.e. MSI-H), MSS,

1 and POLE mutant cancers (**Fig 5A-D**). Due to the proximity of introns and exons in the
2 genome, comparison of their mutation density automatically controls for transcriptional
3 activity and replication timing, both of which are also correlated with H3K36me3³⁸.
4 Exons have more observed and expected number of CpG C>T mutations compared with
5 introns due to their generally higher GC content and CpG methylation levels³⁹.
6 Meanwhile, compared to introns, MSS and POLE showed substantial and significant
7 decrease in observed exonic CpG C>T compared to expected (37.3% and 31.2%, $p <$
8 0.0001, one sample t-test, **Fig 5 A-B**), while this decrease was substantially less in MSI
9 cancers due to the loss of MMR (2.74%, **Fig 5C**). Surprisingly, substantial and
10 significant decrease in the observed exonic mutation density compared with expected
11 exonic mutation density was also observed in MBD4 mutants (21.6%, $p <$ 0.0001, one
12 sample t-test, **Fig 5D**). As MBD4 is responsible for the repair for 5mC deamination
13 damage, the decrease in observed exonic mutations in MBD4 mutants suggests that
14 MMR may also be playing a role in the repair of G:T mismatches.

15 Recently, the MMR system has been reported to preferentially protect actively
16 transcribed genes from mutation during transcription¹³. Consistent with this, we found
17 that there was a negative correlation between CpG C>T mutation density and gene
18 expression level for MSS and MBD4 mutants while there were more mutations in
19 highly expressed genes for MMRd samples (**Fig 5E**). To determine if transcription
20 coupled nucleotide excision repair (TC-NER) may also have a role in repairing 5mC
21 deamination damage, we examined the transcription strand bias in MSI, MSS and
22 MBD4 mutant cancers. We found that the transcription strand bias of CpG C>T

1 mutations was close to zero in all cases (**Fig 5F-H**) suggesting that TC-NER is not
2 involved in its repair. In MBD4 mutant cancers, this lack of transcription strand bias,
3 further suggests that MMR is likely to be playing a role in the repair of the G:T
4 mismatches. We next examined the association of CpG C>T mutation and different
5 epigenetics marks for MBD4 mutants by multivariable logistic regression. Since CpG
6 C>T mutations are highly dependent on methylation (**Extended Data Fig 6A**), we only
7 selected mutations with highly methylated CpG (> 0.9) to ensure that we delineate the
8 impact of histone modifications from DNA methylation. Apart from replication timing,
9 we identified histone mark H3K36me3 had the lowest hazard ratio (HR) for CpG C>T
10 mutation formation (HR = 0.88, $p < 0.001$, **Fig 5I, Supp Table 2**). There were fewer
11 mutations in the regions of high H3K36me3 signal suggesting that in the absence of
12 MBD4, the repair of CpG C>T mutations is dependent on H3K36me3 (**Fig 5J**).
13 Interestingly, H3K36me3 also positively correlates with methylation level (**Extended**
14 **Data Fig 6B**). This suggests that in MBD4 mutants, although CpG methylation is the
15 strongest determinant of mutation density, H3K36me3 activity is also an important
16 factor for accounting for CpG mutation density. Taken together, these results suggest
17 that even in the absence of MBD4, MMR has some capacity to facilitate the repair of
18 5mC deamination damage in a H3K36me3 dependent manner.

19

20 **MMR rather than MBD4 is essential for the repair of 5-methylcytosine**
21 **deamination damage**

22 While purified MBD4 can excise mismatched bases from DNA *in vitro*⁴⁰, it is

unclear whether MBD4 can repair 5mC deamination damage in the absence of MMR proteins *in vivo*. Previous studies have shown that MBD4 binding in the genome is strongly associated with DNA methylation density but is only weakly associated with H3K36me3⁴¹. To determine if MBD4 is able to facilitate the repair of 5mC deamination damage, we identified regions of the genome which are enriched in MBD4 and H3K36me3 based on ChIP-seq data from ENCODE. As the MBD4 ChIP-seq dataset is only available for the HepG2 cell line, we also used H3K36me3 ChIP-seq data from HepG2 to avoid cell type specific bias. After removing regions of low mappability, we identified 119,237 1kb windows in the genome that have either high (top 20%) MBD4 or H3K36me3 (**Fig 6A**). We also identified windows with low MBD4 or low H3K36me3. Using the respective regions, we compared the observed/expected mutation density across the windows in MBD4 mutants, MSS and MMRd (MSI) cancers. In high MBD4 regions, the observed/expected mutation rate was generally lower than the other regions (**Fig 6B**). This is due to the lower methylation level of CpGs in these regions despite overall high density of CpG, as well as the enrichment in early replicating regions (**Extended Data Fig 7A**). Importantly, the observed/expected mutation load was broadly similar across the three cancer types suggesting that MBD4 alone has little impact of observed mutation density. By contrast, for high H3K36me3 regions, in line with the importance of H3K36me3 for MMR, MSI had significantly more observed/expected mutations than MBD4 mutants and MSS cancers. The over 38% decrease in observed/expected mutation load in MBD4 reinforces our earlier results that MMR can repair 5mC deamination damage in the absence of MBD4. To further

1 quantify the impact of MBD4 binding on 5mC deamination damage, we generated
2 multivariable regression models and found that in MMRd cancers, MBD4 signal was
3 not associated with lower likelihood of mutation (HR = 1.06, Figure 6C, see MutL α
4 and MutS α separately in **Extended Data Fig 7B-C**). In MSS cancers, although a small
5 decrease in HR was present for MBD4, H3K36me3 was a larger contributor to mutation
6 likelihood (HR = 0.89 versus HR = 0.96, **Figure 6D**).

7 It has been observed that truncating mutations in MBD4 are common in MSI
8 cancers and truncated MBD4 can exert a dominant negative effect⁴². We have
9 previously shown that MSI cancers with and without MBD4 truncating mutations does
10 not show any difference in C>T mutation density at methylation CpG sites³⁵. To
11 examine this further in the context of MutL α and MutS α deficient MSI cancers, we
12 found that 13.7% (27/197) of MutL α and 11.1% (2/18) of MutS α deficient cancers
13 harboured MBD4 truncating mutations (**Extended Data Fig 7D, Supp Table3**). We did
14 not find any difference in the mutational signature contributions in MBD4 wild-type
15 and mutants for both MutL α and MutS α deficient cancers (**Extended Data Fig 7E**).
16 We also examined the expression of MBD4 in TCGA MSI and MSS samples. While
17 MBD4 expression was significantly lower in MutL α deficient cancers compared with
18 MSS cancers (**Extended Data Fig 7F**), no association between MBD4 expression and
19 SigA contribution was observed for MutL α deficient cancers (**Extended Data Fig 7G**).
20 This further suggests that MBD4 is not a major factor in limiting CpG C>T mutations
21 in MMRd cancers.

22 As MutS α deficient cancers have the highest SigA contribution and is responsible

1 for the recognition of mismatched bases in DNA, our data suggests that it is the essential
2 component for the repair of 5mC deamination damage. Interestingly, we found TDG to
3 be upregulated in MutL α deficient cancers (**Extended Data Fig 7H**). This upregulation
4 suggests that TDG may be the alternative glycosylase that can collaborate with MMR
5 in the repair of 5mC deamination damage independently of MutL α and MBD4 (**Figure**
6 **6E**).

7 8 **Somatic mutations in MutS α mutant cell lines are largely caused by DNA** 9 **replication errors**

10 As MutS α deficient human cancers are highly enriched in SigA with a high
11 proportion of CpG C>T mutations, we were intrigued that clonal mutations derived
12 from two independent cultured MutS α mutant cell lines^{43,44} had a broader distribution
13 of mutation types (**Fig 7A-B**), with both presenting a high contribution of SigB (**Fig**
14 **7C**). One is based on the human HAP1 cell line with knockout of MSH6 mediated by
15 CRISPR-Cas9⁴³ and the other one is based on human DLD-1 colon adenocarcinoma
16 cells with MSH6 deficiency⁴⁴. In both cases, somatic mutations were acquired in culture
17 following clonal isolation and expansion. To determine if the CpG C>T mutations were
18 likely to have formed in a similar way as in human cancers, the replication strand bias
19 of CpG C>T and non-CpG C>T mutations for these cell lines was examined.
20 Interestingly, unlike the human MMRd cancers, both MutS α deficient cell lines showed
21 strong replication strand bias with no significant difference between CpG and non-CpG
22 C>T mutations ($p > 0.05$, Student's t-test, **Fig 7D-E**). As cell lines replicate almost

continuously and the mutations are acquired over the course of just one month, the result is consistent with the CpG C>T mutations forming during DNA replication. In contrast, cancer samples or pre-malignant cells replicate much more slowly than cell lines with mutations accumulating over many years, consistent with most CpG C>T mutations occurring via the deamination of 5mC, independent of DNA replication.

6

7 **Discussion**

8 A number of previous studies have sought to use genomics to determine the
9 mechanisms underlying the different mutational signatures in MMRd cancers^{4-6,16,45}.
10 However, these efforts have been either restricted by cell line/organoid models^{4,5}, non-
11 mammalian models⁴⁵, a lack of samples⁶ or incomplete assessment of all the data types
12 available¹⁶. In this study, we carefully determined the status of the canonical MMR
13 genes and classified each sample as being MutS α or MutL α deficient. In doing so, we
14 found significant differences in the mutational signatures of MutS α and MutL α across
15 four independent cohorts. Specifically, MutS α deficiency presents a high CpG C>T
16 mutation spectrum, while MutL α mutants have a broader mutational spectrum
17 including C>A, C>T and T>C mutations. These results are consistent with a recent
18 published study of whole exome sequenced MSH2/MSH6-deficient gliomas which
19 were all found to have a high frequency of CpG C>T mutation similar to SigA⁴⁶.
20 Another study based on the MSH6 null mouse model also reported elevated mutation
21 frequency and predominance of G:C to A:T transition in MSH6 deficient small
22 intestinal epithelium⁴⁷.

23 Due to the susceptibility of cytosine (at CpG sites) to a variety of chemical
24 modifications, CpG C>T mutations are common in cancer genomes and are generally

1 recognised to be the result of 5mC deamination. Nevertheless, 5mC can also lead to
2 CpG mutations in other ways^{30,35,48}. While MMR is known to correct G:T and other
3 mismatches that result from DNA replication errors, the repair of 5mC deamination
4 should require excision of thymine from the G:T mismatch to restore the correct G:C
5 pair³¹. It has been well established that MBD4 plays a critical role in the repair of
6 mutations caused by 5mC deamination³². C>T mutations that arise from unrepaired
7 5mC deamination induced G:T should show no replication strand asymmetry as the
8 deamination process should be largely independent of DNA replication (**Fig 4C**). By
9 contrast, CpG C>T mutations that arise as a result of DNA replication errors, such as
10 those in POLE mutant cancers³⁴, display a high degree of asymmetry (**Fig 4D**). As the
11 CpG C>T mutations observed in MMRd cancers, particularly those deficient of MutS α ,
12 lack strong replication strand asymmetry, this suggests that most of these mutations
13 likely arise from replication independent 5mC deamination. As further evidence, we
14 found that CpG C>T mutations generated by clonal expansion of MMRd cell lines show
15 significant replication strand bias (**Fig 7D and E**), reflecting the rapid rate of cell
16 division and lack of time for 5mC deamination induced mutations to accumulate. Thus,
17 we propose that the two main mutational processes operational in MMRd cancers are a
18 replication independent 5mC deamination induced C>T mutational signature (i.e. SigA)
19 and a replication dependent DNA polymerase error driven signature (i.e. SigB) (**Fig**
20 **6E**). As SigA contributes to over 50% of mutations observed in most MMRd cancers
21 (**Fig2 C-F**), the inability to repair 5mC deamination damage can be considered the
22 major mutational process in MMRd cancers.

23 Activities of MMR outside the context of DNA replication have been termed
24 ncMMR and have been generally associated with the (error-prone) repair of DNA
25 lesions through the recruitment of error-prone DNA polymerases¹⁰⁻¹². It is therefore

1 surprising that a form of ncMMR appears to participate in the error-free repair of 5mC
2 deamination damage. Nevertheless, a recent study showed evidence of ncMMR
3 dependent error-free repair of oxidative damage in actively transcribed genes¹³. Our
4 data suggests that for the repair of 5mC deamination damage MMR is only performing
5 the function of mismatch/lesion recognition to promote recruitment of other DNA
6 repair factors pathways. Interaction of MMR and base excision repair has been studied
7 *in vitro* in the context of active DNA demethylation¹⁴ and their results support a model
8 where the two pathways collaborate to facilitate error-free removal of 5mC induced
9 lesions including G:T mismatches in a MSH2 and TDG dependent manner. Our
10 findings support this model. MutS α is known to have the role in DNA mismatch
11 recognition and its MSH6 subunit contains the PWWP domain that binds H3K36me3
12 to facilitate its recruitment³⁶. The high contribution of SigA for MutS α deficient cancers
13 suggests that MutS α is essential for efficient mismatch recognition to recruit the
14 MutL α -MBD4 complex, and potentially TDG, for the repair of 5mC deamination
15 damage. Future studies will be required to fully elucidate the mechanisms of this form
16 of ncMMR in *in vivo*.

17 In summary, we demonstrate that replication independent 5mC deamination
18 contributes to most CpG C>T mutations in MMRd cancers. We find that MMR is in
19 fact essential for the repair of 5mC deamination induced G:T mismatches. This non-
20 canonical MMR function is likely to be MutS α dependent as MutS α deficient cancers
21 are highly enriched in CpG C>T mutations. These results provide new insights of
22 mutational processes in MMRd cancers and further our understanding of the ever-
23 important MMR pathway.

24

1 **Acknowledgments**

2 The authors would like to thank Ian Majewski and Rebecca Poulos for valuable
 3 feedback on the manuscript. This research is funded by Seed Funding from the
 4 University Grants Council, The University of Hong Kong to JWHW.

5

6 **Author Contributions**

7 J.W.H.W conceived and designed the research; H.F, J.O. and J.W.H.W developed
 8 methodology and performed research; H.F, X.Z., J.O. J.A.B. and J.W.H.W analysed
 9 data; H.F. and J.W.H.W wrote the manuscript.

10

11 **Declaration of Interests**

12 The authors have declared that no competing interests exist.

13

1 **Methods**

2 **Data collection**

3 A list of 316 MSI-H cancer samples including colorectal cancer (CRC), stomach
4 adenocarcinoma (STAD) and uterine corpus endometrial carcinoma (UCEC) were
5 obtained from TCGA Pan-Cancer dataset⁴⁹. Other three independent cohorts were
6 obtained for validation: Depmap cohort comprises 99 MSI-H samples across 16 cancer
7 types²⁴. MSK-CRC and MSK-UCEC cohorts contain 99 MSI-H colorectal cancers and
8 22 MSI-H endometrial cancers respectively, both with targeted sequencing data^{22,23}. All
9 these data are summarised in **Supp Table 4**.

11 **MutS α and MutL α classification**

12 For the 316 samples from TCGA Pan-Cancer cohort, signature contributions were
13 assigned by fitting COSMIC mutation signature v3 via the Sigfit R package²¹. Samples
14 with high contributions (>10%) of signature SBS10a/b, SBS14 and SBS20 were
15 excluded to avoid the effect of mutational processes that are caused by POLE and
16 POLD1. DNA mutation, RNA expression and methylation data were applied to the
17 remaining 266 samples for classification by the steps below: (1) Linear regression
18 analysis was performed based on MLH1 methylation and expression. The regression
19 equation was obtained as: $y=9.050-4.996*x$. Hypermethylated MLH1 is defined as
20 $\beta>0.25$ and the low MLH1 expression cutoff value was obtained as 7.8 based on the
21 equation. (2) Then MutS α and MutL α were determined based on the RNA expression
22 and truncating mutation of the MMR genes that are elaborated in **Extended Data Fig**
23 **1**. For 99 MSI-H samples from the Depmap cohort, we first classified samples with
24 truncating mutations in MSH2/MSH6 as MutS α if they have no aberration in
25 MLH1/PMS2. Then the remaining samples with no aberration in MSH2/MSH6 were

classified as MutL α (**Supp Table 5**). Due to the availability of data for MSK-CRC cohort, the classification of MutS α and MutL α is based on DNA mutations of MMR genes. Samples with truncating mutations in MSH2/MSH6 and without truncating mutations in MLH1/PMS2 are classified as MutS α . The remaining samples are classified as MutL α (**Supp Table 56**). For samples from the MSK-UCEC cohort the classification is based on immunohistochemistry of the four MMR genes (**Supp Table 7**).

Mutation simulation at tri-nucleotide resolution

Mutation simulation was performed by SigProfilerSimulator⁵⁰ to preclude the bias of tri-nucleotide composition which could affect the mutation distribution in local regions. Briefly, the total number of simulated mutations for each sample is equal to the observed mutations, but the position of the mutation is relocated according to the frequency of tri-nucleotide context along the given region. Each sample was simulated 100 times and all the mutations are combined as expected mutations for subsequent local mutation density analysis.

Mutational signature analysis

The profile of each signature was represented by six substitution subtypes: C>A, C>G, C>T, T>A, T>C and T>G. For signatures generated by tri-nucleotide context, each substitution on the cancer genome was examined by incorporating information on the bases immediately 5' and 3' to each mutated base to generate 96 possible mutational types. *De novo* signatures were extracted by Sigfit which applies a Bayesian inference algorithm²¹. Mutational signatures were displayed and reported based on the observed tri-nucleotide frequency of the human genome. For validation cohorts, contribution of

1 *de novo* signatures for each sample was calculated by fitting the mutations to the
2 extracted *de novo* signatures.

3

4 **Replication timing and mutation density**

5 The replication time of different chromosome regions was obtained for the HepG2 cell
6 line from the ENCODE data portal⁵¹. Exome sequence with known replication time was
7 divided into five bins from late to early: [-3.88, 51.98), [51.98, 66.30), [66.30, 74.95),
8 [74.95, 80.74), [80.74, 87.95]. The counts of mutations within each bin were calculated
9 as observed mutation. Similarly, the expected mutation counts were also computed for
10 each bin based on simulated data. The slope was obtained from the linear regression
11 model based on the correlation of mutation ratio (obs/exp) and replication timing.

12

13 **Gene expression and mutation density**

14 The general gene expression data were obtained from GTEx Portal and all expressed
15 genes were integrated into four bins according to the expression quartile. For each bin,
16 only sites located within early replicated regions are adopted. The size of each bin was
17 calculated based on the length of exons of each gene. The count of observed and
18 expected mutations was calculated for each bin to determine mutation density.

19

20 **Calculating strand asymmetries**

21 Replication direction was defined using replication timing profiles from a previously
22 published paper⁵². Left (leading)- and right (lagging)-replicating regions were
23 determined by the derivative of the profile. For a given mutation type in the right
24 replication direction, the mutation counts (N_l) in that region were calculated, and its
25 complementary mutation was calculated as n_l . Correspondingly, the mutation counts of

1 this mutation type in left replication direction were calculated as N_2 , and its
2 complementary mutation was calculated as n_2 . Then, asymmetry (A) was calculated as:

$$3 \quad A = \log_2((N_1 + n_2)/(N_2 + n_1))$$

4 For the transcription strand asymmetry, coding and template strand were obtained
5 from a published study⁵³ and the asymmetry is reported as log2 ratio of (mutation count
6 within template regions) / (mutation counts within coding regions).

7

8 **Computing mutation density in exonic and intronic regions**

9 All gene coordinates were obtained from the UCSC table browser. Middle exons and
10 middle introns were extracted for each gene. Then, the mutation density was calculated
11 as mutation counts per megabase for both exonic and intronic regions.

12

13 **Associations between MBD4 mutant mutation density, histone marks and CpG** 14 **methylation**

15 As MBD4 mutants are derived from acute myelocytic leukemia, we obtained whole
16 genome bisulfite sequencing data for E050 (Mobilized CD34 Primary Cells). Histone
17 marks including H3K36me3, H3K27me3, H3K9me3, H3K4me3, H3K4me1, H3K27ac
18 as well as DNase I hypersensitive site are derived from common myeloid progenitor
19 and CD34-positive samples. For the data to estimate regression model for MSS and
20 MMRd (MSI) samples, the mutations are from TCGA CRC cancer. Histone marks as
21 well as DNase I hypersensitive site are obtained from Homo sapiens large intestine
22 male embryo (108 days). All these data are downloaded from the Roadmap
23 Epigenomics Atlas⁵⁴. MBD4 ChIP-seq data obtained from HepG2 cells from
24 ENCODE⁵¹. Only sites with methylation value >0.9 are adopted for fitting the logistic
25 model. For the correlation of CpG methylation, MBD4 mutants mutation and

1 H3K36me3 signal, all cytosines in the CpG di-nucleotide were merged into 12 bins
 2 according to their methylation level as: [0], (0, 0.1], ..., (0.9, 1.0), [1]. These bins were
 3 then used as intersected regions to calculate mutation density in each methylation level.
 4 H3K36me3 bins were set based on the H3K36me3 signal. For the grouping of MBD4
 5 and H3K36me occupied regions, the H3K36me data were also obtained from HepG2
 6 from ENCODE⁵¹. MBD4 and H3K36me3 signal/input were calculated across 1kb
 7 windows across the genome. Regions that had an average mappability score of <1 based
 8 on UCSC *Duke Uniqueness* 35 bp and those that overlapped *DAC blacklisted regions*
 9 were removed from the analysis.

10

11 **Supplementary tables**

12 Supp Table 1. Sample information for TCGA cohort.

13 Supp Table 2. Multivariable logistic regression model.

14 Supp Table 3. MBD4 truncating mutation annotation in MSI samples.

15 Supp Table 4. Data cohorts summary.

16 Supp Table 5. Sample information for Depmap cohort.

17 Supp Table 6. Sample information for MSK-CRC cohort.

18 Supp Table 7. Sample information for MSK-UCEC cohort.

19

20 **References**

- 21 1 Peltomaki, P. Role of DNA mismatch repair defects in the pathogenesis of human cancer.
22 *J Clin Oncol* **21**, 1174-1179, doi:10.1200/JCO.2003.04.060 (2003).
- 23 2 Kunkel, T. A. & Erie, D. A. J. A. R. B. DNA mismatch repair. **74**, 681-710 (2005).
- 24 3 Alexandrov, L. B. *et al.* The repertoire of mutational signatures in human cancer. *Nature*
25 **578**, 94-+, doi:10.1038/s41586-020-1943-3 (2020).
- 26 4 Drost, J. *et al.* Use of CRISPR-modified human stem cell organoids to study the origin of
27 mutational signatures in cancer. *Science* **358**, 234-+, doi:10.1126/science.aao3130 (2017).
- 28 5 Zou, X. *et al.* Dissecting mutational mechanisms underpinning signatures caused by
29 replication errors and endogenous DNA damage. 2020.2008.2004.234245,
30 doi:10.1101/2020.08.04.234245 %J bioRxiv (2020).

1 6 Davies, H. *et al.* Whole-Genome Sequencing Reveals Breast Cancers with Mismatch Repair
2 Deficiency. *Cancer Res* **77**, 4755–4762, doi:10.1158/0008-5472.CAN-17-1083 (2017).

3 7 Kunkel, T. A. & Erie, D. A. Eukaryotic Mismatch Repair in Relation to DNA Replication. *Annu*
4 *Rev Genet* **49**, 291–313, doi:10.1146/annurev-genet-112414-054722 (2015).

5 8 Crouse, G. F. Non-canonical actions of mismatch repair. *DNA Repair (Amst)* **38**, 102–109,
6 doi:10.1016/j.dnarep.2015.11.020 (2016).

7 9 Cascalho, M., Wong, J., Steinberg, C. & Wabl, M. Mismatch repair co-opted by
8 hypermutation. *Science* **279**, 1207–1210, doi:10.1126/science.279.5354.1207 (1998).

9 10 Pena-Diaz, J. *et al.* Noncanonical Mismatch Repair as a Source of Genomic Instability in
10 Human Cells. *Mol Cell* **47**, 669–680, doi:10.1016/j.molcel.2012.07.006 (2012).

11 11 Zlatanou, A. *et al.* The hMsh2–hMsh6 Complex Acts in Concert with Monoubiquitinated
12 PCNA and Pol eta in Response to Oxidative DNA Damage in Human Cells. *Mol Cell* **43**,
13 649–662, doi:10.1016/j.molcel.2011.06.023 (2011).

14 12 Supek, F. & Lehner, B. Clustered Mutation Signatures Reveal that Error-Prone DNA Repair
15 Targets Mutations to Active Genes. *Cell* **170**, 534–+, doi:10.1016/j.cell.2017.07.003 (2017).

16 13 Huang, Y. P., Gu, L. Y. & Li, G. M. H3K36me3-mediated mismatch repair preferentially
17 protects actively transcribed genes from mutation. *J Biol Chem* **293**, 7811–7823,
18 doi:10.1074/jbc.RA118.002839 (2018).

19 14 Grin, I. & Ishchenko, A. A. An interplay of the base excision repair and mismatch repair
20 pathways in active DNA demethylation. *Nucleic Acids Res* **44**, 3713–3727,
21 doi:10.1093/nar/gkw059 (2016).

22 15 Volkova, N. V. *et al.* Mutational signatures are jointly shaped by DNA damage and repair.
23 *Nat Commun* **11**, doi:10.1038/s41467-020-15912-7 (2020).

24 16 Nemeth, E. *et al.* Two main mutational processes operate in the absence of DNA mismatch
25 repair. *DNA Repair* **89**, doi:ARTN 102827
26 10.1016/j.dnarep.2020.102827 (2020).

27 17 Cortes-Ciriano, I., Lee, S., Park, W. Y., Kim, T. M. & Park, P. J. A molecular portrait of
28 microsatellite instability across multiple cancers. *Nat Commun* **8**, doi:ARTN 15180
29 10.1038/ncomms15180 (2017).

30 18 Haradhvala, N. J. *et al.* Distinct mutational signatures characterize concurrent loss of
31 polymerase proofreading and mismatch repair. *Nat Commun* **9**, doi:ARTN 1746
32 10.1038/s41467-018-04002-4 (2018).

33 19 Simpkins, S. B. *et al.* MLH1 promoter methylation and gene silencing is the primary cause
34 of microsatellite instability in sporadic endometrial cancers. *Hum Mol Genet* **8**, 661–666,
35 doi:DOI 10.1093/hmg/8.4.661 (1999).

36 20 Haraldsdottir, S. *et al.* Patients with colorectal cancer associated with Lynch syndrome and
37 MLH1 promoter hypermethylation have similar prognoses. *Genet Med* **18**, 863–868,
38 doi:10.1038/gim.2015.184 (2016).

39 21 Gori, K. & Baez-Ortega, A. J. b. sigfit: flexible Bayesian inference of mutational signatures.
40 372896 (2018).

41 22 Yaeger, R. *et al.* Clinical Sequencing Defines the Genomic Landscape of Metastatic
42 Colorectal Cancer. *Cancer Cell* **33**, 125–+, doi:10.1016/j.ccell.2017.12.004 (2018).

43 23 Soumerai, T. E. *et al.* Clinical Utility of Prospective Molecular Characterization in Advanced
44 Endometrial Cancer. *Clin Cancer Res* **24**, 5939–5947, doi:10.1158/1078-0432.Ccr-18-

0412 (2018).

24 Ghandi, M. *et al.* Next-generation characterization of the Cancer Cell Line Encyclopedia. *Nature* **569**, 503–+, doi:10.1038/s41586-019-1186-3 (2019).

25 Pritchard, C. C. *et al.* Complex MSH2 and MSH6 mutations in hypermutated microsatellite unstable advanced prostate cancer. *Nat Commun* **5**, 4988, doi:10.1038/ncomms5988 (2014).

26 Supek, F. & Lehner, B. Scales and mechanisms of somatic mutation rate variation across the human genome. *DNA Repair (Amst)* **81**, 102647, doi:10.1016/j.dnarep.2019.102647 (2019).

27 Lawrence, M. S. *et al.* Mutational heterogeneity in cancer and the search for new cancer-associated genes. *Nature* **499**, 214–218, doi:10.1038/nature12213 (2013).

28 Preston, B. D., Albertson, T. M. & Herr, A. J. DNA replication fidelity and cancer. *Semin Cancer Biol* **20**, 281–293, doi:10.1016/j.semcancer.2010.10.009 (2010).

29 Haradhvala, N. J. *et al.* Mutational Strand Asymmetries in Cancer Genomes Reveal Mechanisms of DNA Damage and Repair. *Cell* **164**, 538–549, doi:10.1016/j.cell.2015.12.050 (2016).

30 Tomkova, M. & Schuster-Bockler, B. DNA Modifications: Naturally More Error Prone? *Trends Genet* **34**, 627–638, doi:10.1016/j.tig.2018.04.005 (2018).

31 Bellacosa, A. & Drohat, A. C. Role of base excision repair in maintaining the genetic and epigenetic integrity of CpG sites. *DNA Repair* **32**, 33–42, doi:10.1016/j.dnarep.2015.04.011 (2015).

32 Sanders, M. A. *et al.* MBD4 guards against methylation damage and germ line deficiency predisposes to clonal hematopoiesis and early-onset AML. *Blood* **132**, 1526–1534, doi:10.1182/blood-2018-05-852566 (2018).

33 Fang, H. *et al.* Mutational processes of distinct POLE exonuclease domain mutants drive an enrichment of a specific TP53 mutation in colorectal cancer. *Plos Genet* **16**, doi:ARTN e1008572 10.1371/journal.pgen.1008572 (2020).

34 Tomkova, M., Tomek, J., Kriaucionis, S. & Schuster-Bockler, B. Mutational signature distribution varies with DNA replication timing and strand asymmetry. *Genome Biol* **19**, doi:ARTN 129 10.1186/s13059-018-1509-y (2018).

35 Poulos, R. C., Olivier, J. & Wong, J. W. H. The interaction between cytosine methylation and processes of DNA replication and repair shape the mutational landscape of cancer genomes. *Nucleic Acids Res* **45**, 7786–7795, doi:10.1093/nar/gkx463 (2017).

36 Li, F. *et al.* The Histone Mark H3K36me3 Regulates Human DNA Mismatch Repair through Its Interaction with MutS alpha. *Cell* **153**, 590–600, doi:10.1016/j.cell.2013.03.025 (2013).

37 Frigola, J. *et al.* Reduced mutation rate in exons due to differential mismatch repair. *Nat Genet* **49**, 1684–+, doi:10.1038/ng.3991 (2017).

38 Huang, Y. P. & Li, G. M. DNA mismatch repair preferentially safeguards actively transcribed genes. *DNA Repair* **71**, 82–86, doi:10.1016/j.dnarep.2018.08.010 (2018).

39 Saxonov, S., Berg, P. & Brutlag, D. L. A genome-wide analysis of CpG dinucleotides in the human genome distinguishes two distinct classes of promoters. *Proc Natl Acad Sci U S A* **103**, 1412–1417, doi:10.1073/pnas.0510310103 (2006).

1 40 Morera, S. *et al.* Biochemical and structural characterization of the glycosylase domain of
2 MBD4 bound to thymine and 5-hydroxymethyluracil-containing DNA. *Nucleic Acids Res*
3 **40**, 9917-9926, doi:10.1093/nar/gks714 (2012).

4 41 Baubec, T., Ivanek, R., Lienert, F. & Schubeler, D. Methylation-dependent and -
5 independent genomic targeting principles of the MBD protein family. *Cell* **153**, 480-492,
6 doi:10.1016/j.cell.2013.03.011 (2013).

7 42 Bader, S. A., Walker, M. & Harrison, D. J. A human cancer-associated truncation of MBD4
8 causes dominant negative impairment of DNA repair in colon cancer cells. *Br J Cancer* **96**,
9 660-666, doi:10.1038/sj.bjc.6603592 (2007).

10 43 Zou, X. Q. *et al.* Validating the concept of mutational signatures with isogenic cell models.
11 *Nat Commun* **9**, doi:ARTN 1744
12 10.1038/s41467-018-04052-8 (2018).

13 44 ádám *et al.* Long-term treatment with the PARP inhibitor niraparib does not increase the
14 mutation load in cell line models and tumour xenografts. (2018).

15 45 Meier, B. *et al.* Mutational signatures of DNA mismatch repair deficiency in *C. elegans* and
16 human cancers. *Genome Res* **28**, 666-675, doi:10.1101/gr.226845.117 (2018).

17 46 Touat, M. *et al.* Mechanisms and therapeutic implications of hypermutation in gliomas.
18 *Nature* **580**, 517-+, doi:10.1038/s41586-020-2209-9 (2020).

19 47 Mark, S. C. *et al.* Elevated mutant frequencies and predominance of G : C to A : T transition
20 mutations in Msh6(-/-) small intestinal epithelium. *Oncogene* **21**, 7126-7130,
21 doi:10.1038/sj.onc.1205861 (2002).

22 48 Nabel, C. S., Manning, S. A. & Kohli, R. M. The Curious Chemical Biology of Cytosine:
23 Deamination, Methylation, and Oxidation as Modulators of Genomic Potential. *Acs Chem*
24 *Biol* **7**, 20-30, doi:10.1021/cb2002895 (2012).

25 49 Wilks, C. *et al.* The Cancer Genomics Hub (CGHub): overcoming cancer through the power
26 of torrential data. *Database-Oxford*, doi:ARTN bau093
27 10.1093/database/bau093 (2014).

28 50 Bergstrom, E. N., Barnes, M., Martincorena, I. & Alexandrov, L. B. J. B. Generating realistic
29 null hypothesis of cancer mutational landscapes using SigProfilerSimulator. (2020).

30 51 Sloan, C. A. *et al.* ENCODE data at the ENCODE portal. *Nucleic Acids Res* **44**, D726-D732,
31 doi:10.1093/nar/gkv1160 (2016).

32 52 Koren, A. *et al.* Differential relationship of DNA replication timing to different forms of
33 human mutation and variation. *Am J Hum Genet* **91**, 1033-1040,
34 doi:10.1016/j.ajhg.2012.10.018 (2012).

35 53 Vohringer, H. & Gerstung, M. J. b. Learning mutational signatures and their
36 multidimensional genomic properties with TensorSignatures. 850453 (2019).

37 54 Roadmap Epigenomics, C. *et al.* Integrative analysis of 111 reference human epigenomes.
38 *Nature* **518**, 317-330, doi:10.1038/nature14248 (2015).

39

Figure 1

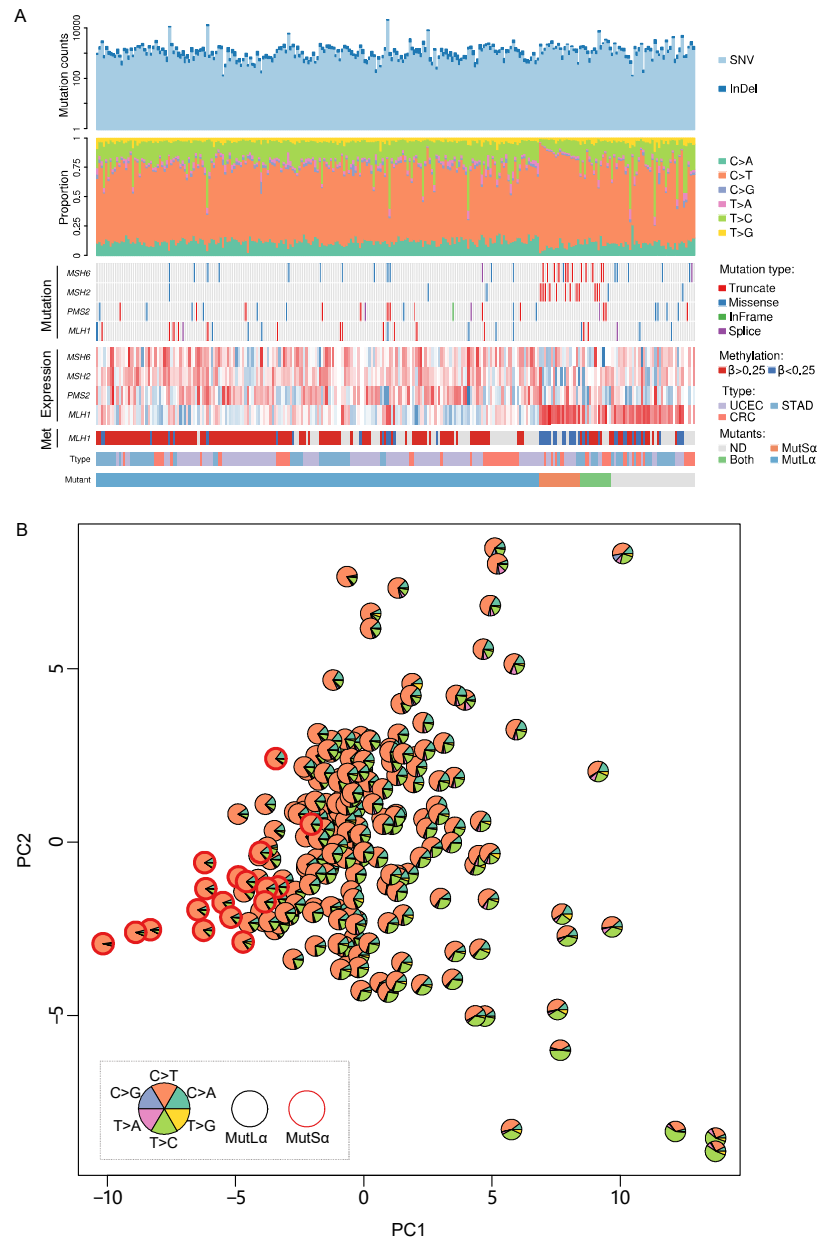


Fig 1. Landscape of MSI-H samples. (A) Profile of mutation burden and six types of mutation frequency, as well as the aberrant status of mismatch repair genes including DNA mutation, RNA expression and methylation. Cancer types and mutants classification are also indicated. **(B)** Principal component analysis of MSI-H cancer samples based on the frequency of 96 types of mutational contexts. The fractions of the six types of mutations are represented by the area of the sectors and MutS α mutants are highlighted in red.

Figure 2

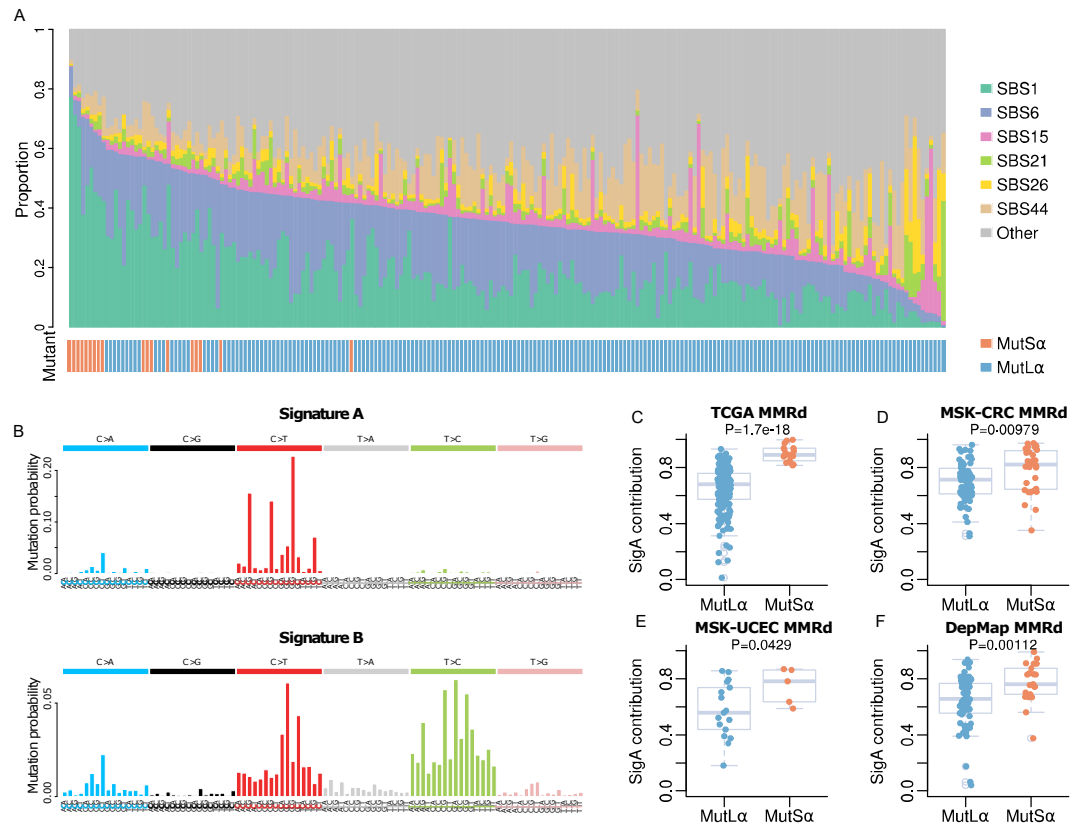


Fig 2. Mutation signature contribution in MutSa and MutLa mutants. (A) Fraction of MMRd associated signatures and age-related signature SBS1 contribution in MutSa and MutLa mutants. (B) The spectrum of *de novo* signatures extracted from TCGA MSI-H cancer samples. (C-F) The boxplot of SigA contribution for MutSa and MutLa mutants in TCGA-MMRd, MSK-CRC, MSK-UEC and Depmap MMRd cohort. P-values were calculated by two-tailed Student's t-test.

Figure 3

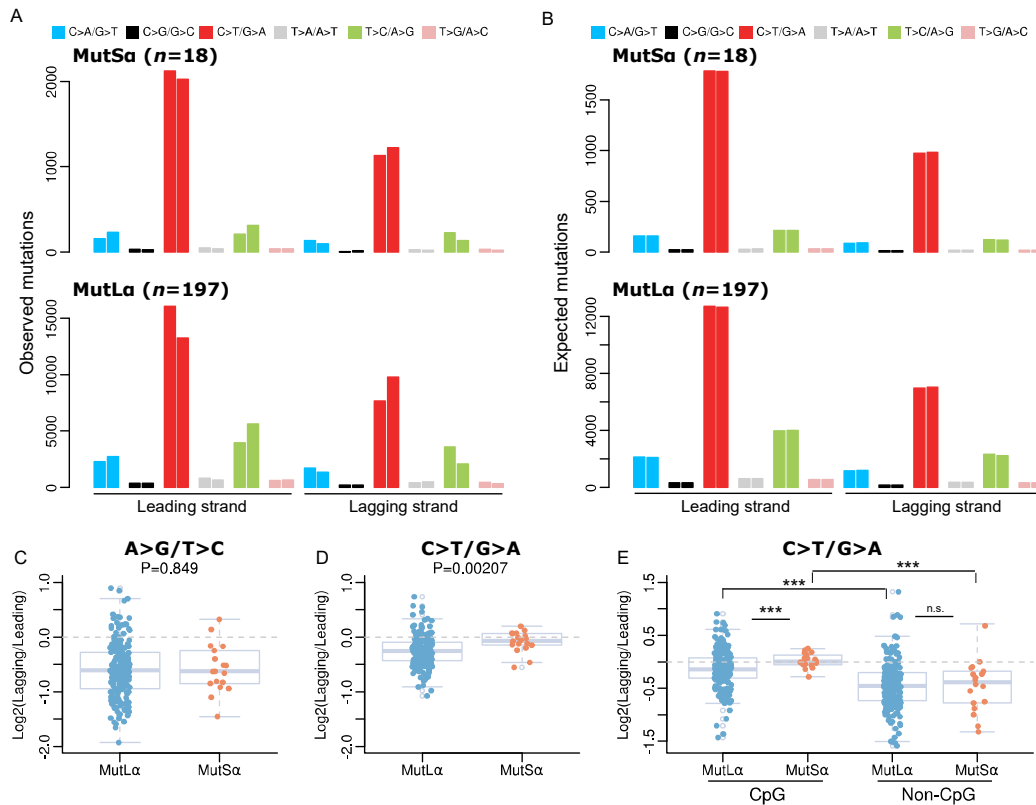


Fig 3. Replication asymmetry for MMR deficiency samples. Landscape of replication asymmetry for all observed mutations (A) and expected mutations (B) in MutSa and MutLa mutants. The expected mutations were obtained from simulation data that consider the abundance of tri-nucleotide mutational contexts. (C-D) Boxplot of replication stand bias for A>G/T>C and C>T/G>A mutations in MutSa and MutLa mutants. (E) Boxplot of replication stand bias for CpG C>T and non-CpG C>T mutations in MutSa and MutLa mutants. *** <0.001, n.s. >0.05, two-tailed Student's t-test.

Figure 4

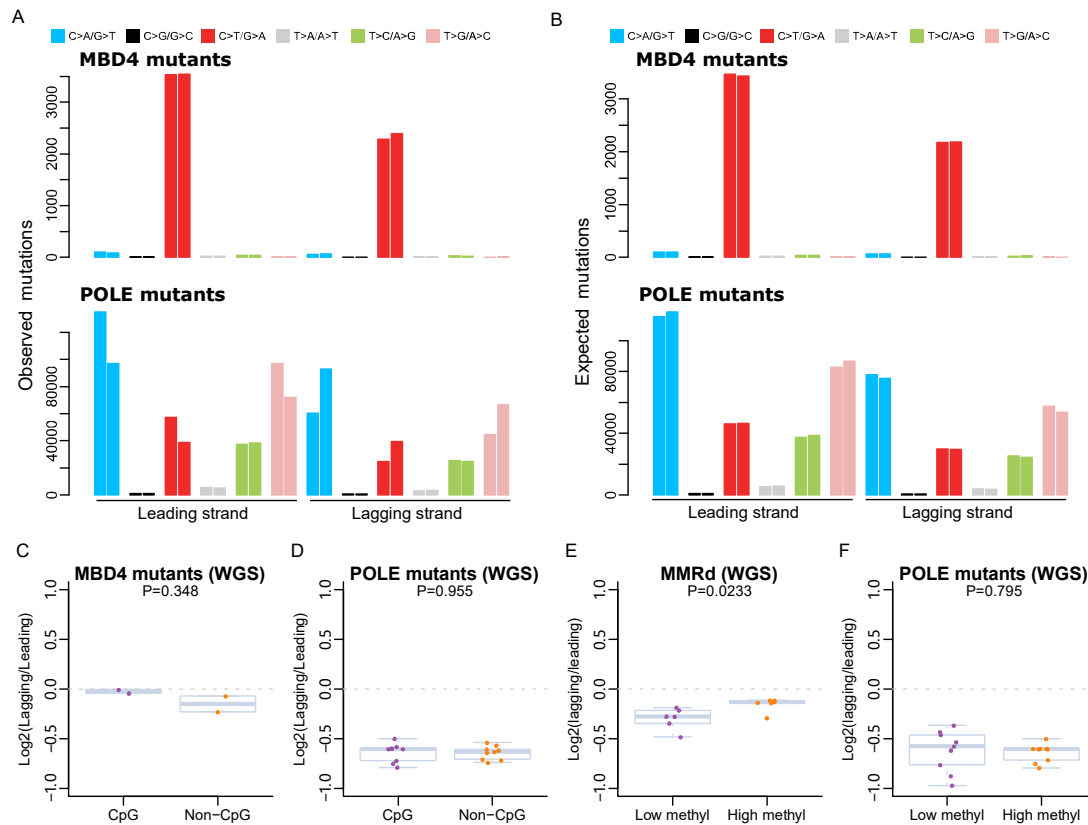


Fig 4. Replication asymmetry for MBD4 and POLE mutants. Landscape of replication asymmetry for all observed mutations (A) and expected mutations (B) in MBD4 and POLE mutants. The expected mutations are obtained from simulation data that consider the abundance of tri-nucleotide mutational contexts. (C-D) Boxplot of replication stand bias for CpG C>T and non-CpG C>T mutations in MBD4 and POLE mutants. (E-F) Boxplot of replication stand bias for CpG C>T mutations in highly methylated and lowly methylated regions for MMR deficiency samples and POLE mutants. Sites with beta value >0.3 are defined as highly methylated while <0.3 as lowly methylated. The range of mutation counts in lowly and highly methylated sites for calculating strand bias in MMRd samples were (52-156) and (2,347-6,145) respectively, and for POLE mutants (72-703) and (4,032-39,580) respectively.

Figure 5

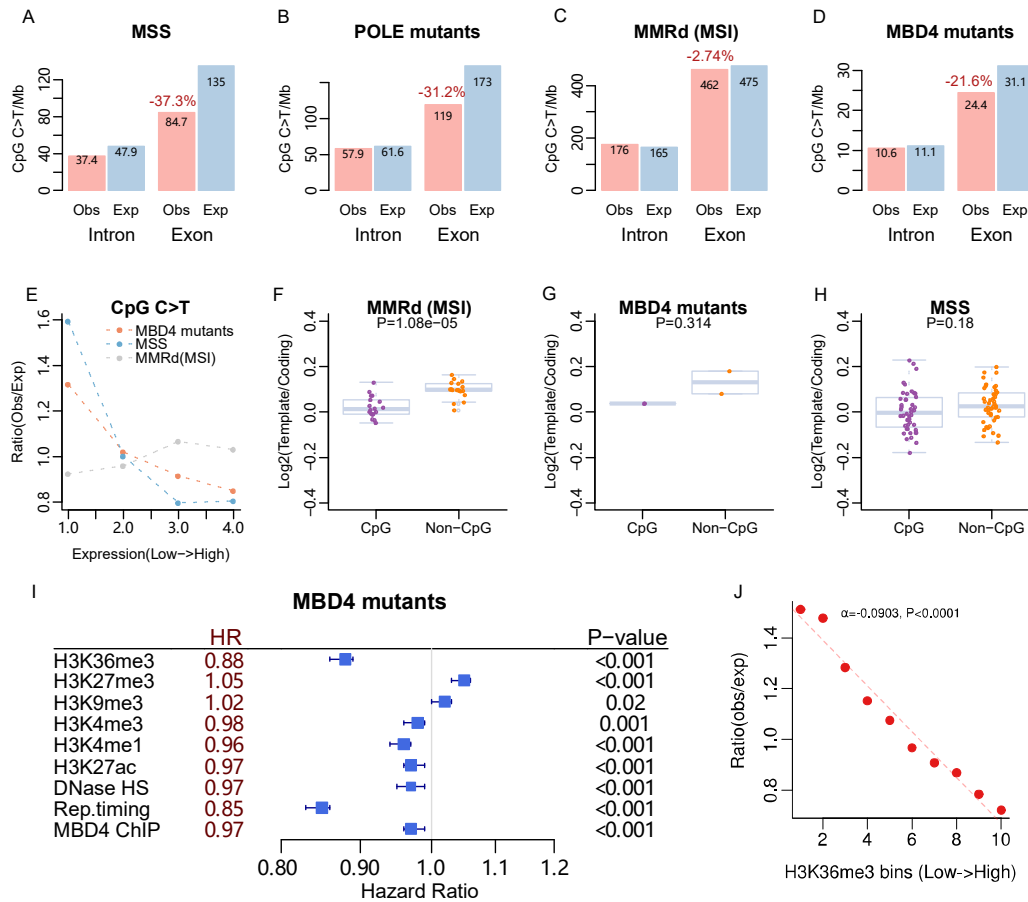


Fig 5. Association of mutation frequency with local determinants for different samples. (A-D) Observed and expected mutation densities in exons and introns across MSS, POLE mutants, MMR deficiency samples and MBD4 mutants. The expected mutations are obtained from simulation data that consider the abundance of trinucleotide mutational contexts. The decrease of observed and expected mutation density in exonic regions is indicated and calculated as (obs-exp)/exp. **(E)** Correlation of CpG C>T mutation ratio (obs/exp) with gene expression for MMR deficiency samples, MSS and MBD4 mutants. The P-values of the correlation are $7.7e-4$, $<2.2e-16$ and 0.167 for MBD4 mutants, MMR deficiency samples and MSS respectively, and they were obtained from the linear regression model by fitting observed mutation density with unbinned gene expression. **(F-H)** Boxplot of transcription strand bias for

1 CpG C>T and non-CpG C>T mutations in MMR deficiency samples, MSS samples and
 2 MBD4 mutants. **(I)** The hazard ratio of different epigenetics marks for CpG C>T
 3 mutation formation from multi-variable logistic regression model. 95% confidence
 4 level is indicated. P-value is calculated by Wald's test. **(J)** Correlation between
 5 mutations in MBD4 mutants and H3K36me3 signal from mobilised CD34+ primary
 6 cells.

7

Figure 6

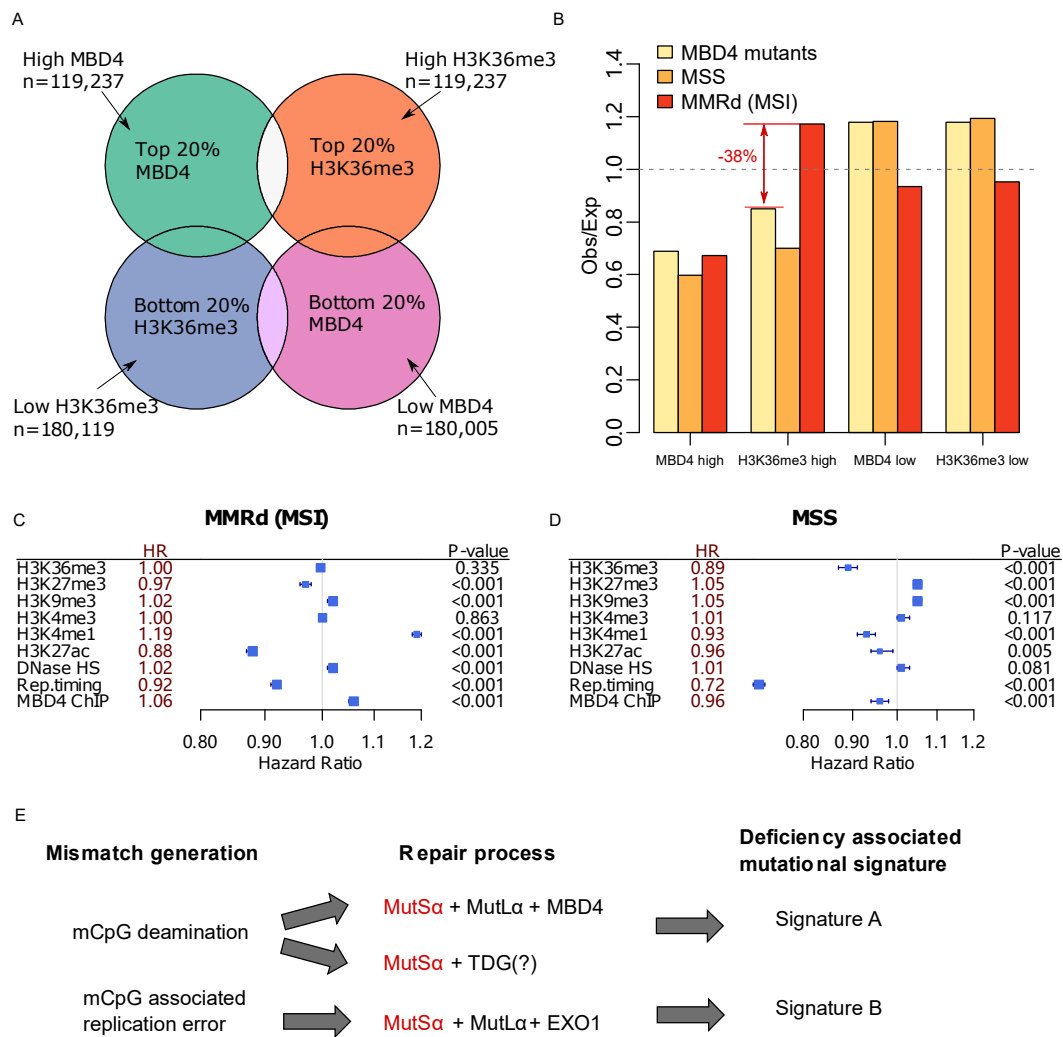


Fig 6. Association of MBD4 binding sites, histone mark H3K36me3 and mutations.

(A) Venn diagram indicating the number of regions classified as top and bottom MBD4 and H3K36me3 signal based on the HepG2 cell line. (B) The ratio of observed and expected CpG C>T mutations in different regions for MBD4 mutants, MSS and MMRd cancers. The hazard ratio of different epigenetics marks for CpG C>T mutation formation from multivariable logistic regression model for MMRd (C) and MSS (D) cancers. 95% confidence level is indicated. P-value was calculated by Wald's test. (E) Schematic of proposed mechanism of mismatch formation, repair and mutations.

Figure 7

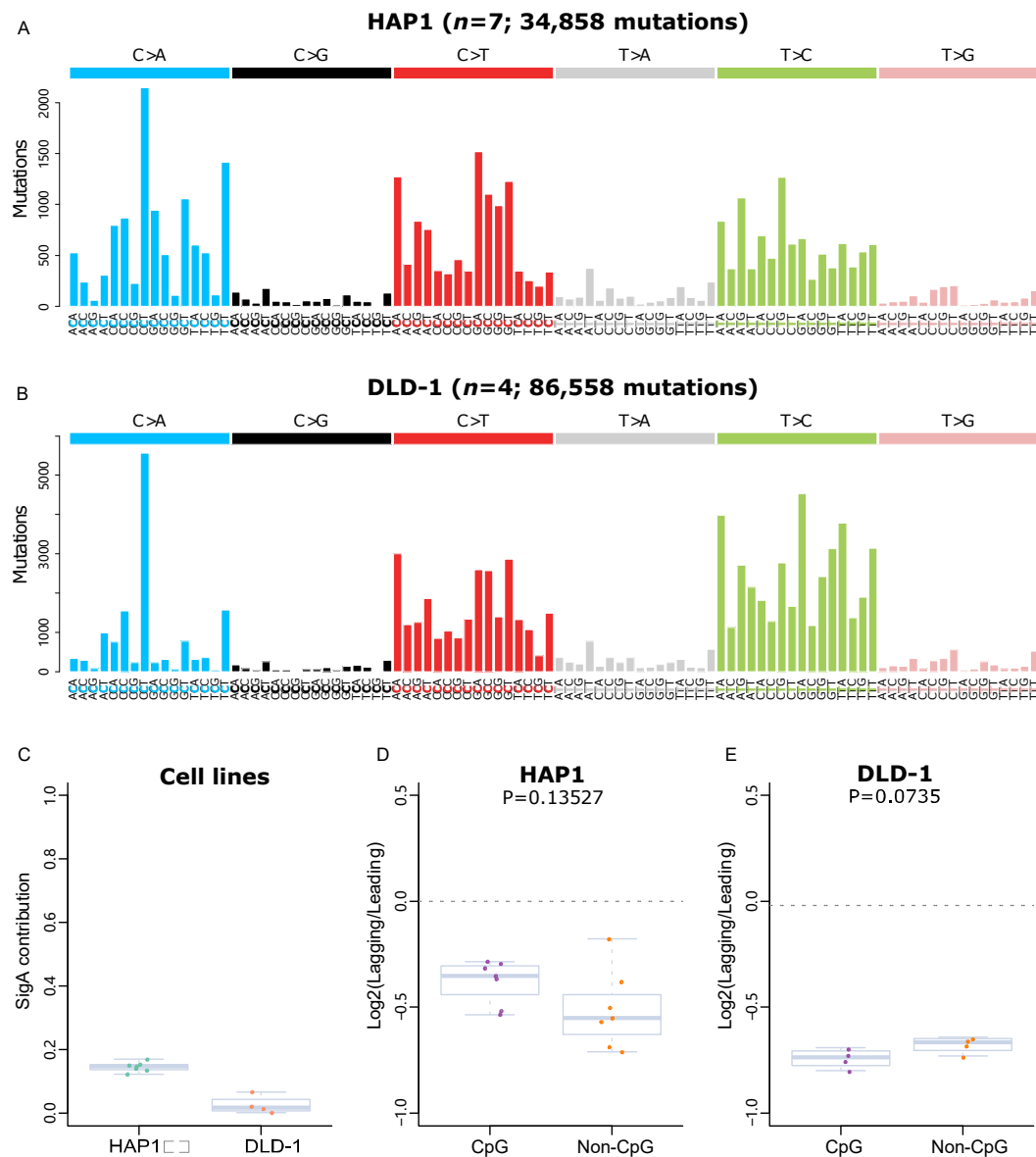


Fig 7. Mutation spectrum and replication asymmetry for MutS α mutant cell lines. Mutation spectrum for cultured HAP1 (A) and DLD-1 (B) cell lines. (C) SigA contribution for HAP1 and DLD-1 cell lines. Boxplot of replication stand bias for CpG C>T and non-CpG C>T mutations in HAP1 (D) and DLD-1 (E) cell lines.