

**Genome-wide genetic screening with chemically-mutagenized haploid  
embryonic stem cells**

Josep V. Forment<sup>1,2</sup>, Mareike Herzog<sup>1,2</sup>, Julia Coates<sup>1</sup>, Tomasz Konopka<sup>3</sup>, Bianca V.  
Gapp<sup>3</sup>, Sebastian M. Nijman<sup>3,4</sup>, David J. Adams<sup>2</sup>, Thomas M. Keane<sup>2</sup> and Stephen  
P. Jackson<sup>1,2</sup>

<sup>1</sup>The Wellcome Trust and Cancer Research UK Gurdon Institute, and Department of  
Biochemistry, University of Cambridge, Cambridge, UK

<sup>2</sup>The Wellcome Trust Sanger Institute, Hinxton, Cambridge, UK

<sup>3</sup>Ludwig Institute for Cancer Research Ltd. and Target Discovery Institute, Nuffield  
Department of Medicine, University of Oxford, Oxford, UK

<sup>4</sup>Research Center for Molecular Medicine of the Austrian Academy of Sciences  
(CeMM), Vienna, Austria

**Authors for correspondence:**

Josep V. Forment [j.forment@gurdon.cam.ac.uk](mailto:j.forment@gurdon.cam.ac.uk)

Stephen P. Jackson [s.jackson@gurdon.cam.ac.uk](mailto:s.jackson@gurdon.cam.ac.uk)

24    **Abstract**

25    In model organisms, classical genetic screening via random mutagenesis provides  
26    key insights into the molecular bases of genetic interactions, helping defining  
27    synthetic-lethality, synthetic-viability and drug-resistance mechanisms. The limited  
28    genetic tractability of diploid mammalian cells, however, precludes this approach.  
29    Here, we demonstrate the feasibility of classical genetic screening in mammalian  
30    systems by using haploid cells, chemical mutagenesis and next-generation  
31    sequencing, providing a new tool to explore mammalian genetic interactions.

32

33 Classical genetic screens with chemical mutagens assign functionality to genes in  
34 model organisms<sup>1,2</sup>. Since most mutagenic agents yield single-nucleotide variants  
35 (SNVs), mutation clustering provides information on the functionality of protein  
36 domains, and defines key amino acid residues within them<sup>3</sup>. RNA interference  
37 (RNAi) allows forward-genetic screening in human cell cultures<sup>3</sup>, and insertional  
38 mutagenesis in near-haploid human cancer cells<sup>4</sup> and whole-genome CRISPR/Cas9  
39 small-guide RNA (sgRNA) libraries have also been used for this purpose<sup>5,6</sup>. Although  
40 powerful, such loss-of-function (LOF) approaches miss phenotypes caused by  
41 separation-of-function or gain-of-function SNV mutations, are less informative on  
42 defining functional protein regions, and are not well suited to studying functions of  
43 essential genes<sup>7</sup>. Here, we describe the generation of chemically mutagenized  
44 mammalian haploid cell libraries, and establish their utility to identify recessive  
45 suppressor mutations by using resistance to 6-thioguanine (6-TG) as a proof-of-  
46 principle.

47  
48 Comprehensive libraries of homozygous SNV-containing mutant clones are not  
49 feasible to obtain in cells with diploid genomes. To circumvent this, we used H129-3  
50 haploid mouse embryonic stem cells (mESCs)<sup>8</sup> that we had mock-treated or treated  
51 with varying doses of the DNA-alkylating agent ethylmethanesulfonate (EMS), a  
52 chemical inducer of SNVs<sup>9</sup> (**Fig. 1a; Supplementary Results, Supplementary Fig.**  
53 **1a**). For comparison, the same procedure was performed on diploid H129-3 mESCs  
54 (**Supplementary Fig. 1b**). Haploid and diploid mutant libraries were then screened  
55 for suppressors of cellular sensitivity to 6-TG (**Fig. 1b**). Ensuing analyses revealed  
56 EMS-dose dependent induction of 6-TG resistance, with more clones arising in  
57 haploid than in diploid cells (**Fig. 1c**), thus highlighting the advantage of identifying  
58 suppressor mutations in a haploid genetic background.

59  
60 Next, we isolated 196 6-TG resistant clones from EMS-generated haploid cell  
61 libraries. To assess the feasibility of identifying causative suppressor mutations, we  
62 subjected DNA samples from seven resistant clones, and from control mESCs not  
63 treated with EMS, to whole-exome DNA sequencing. Ensuing analyses, comparing  
64 sequences from EMS-resistant clones with control mESCs and the 129S5 mouse  
65 genome (see Methods), identified homozygous base insertions/deletions (INDELs)  
66 and SNVs. Only 11.3% of these affected coding sequences and were non-

67 synonymous (**Fig. 1d**). Thus, while each resistant clone had ~370 INDEL/SNV  
68 mutations (**Supplementary Fig. 1c**), on average only ~40 of these were in coding  
69 sequences and non-synonymous.

70

71 We then identified candidate suppressor genes by analyzing this set of non-  
72 synonymous mutations. We defined suppressor gene candidates as those being  
73 mutated in multiple independent clones and harboring multiple potential deleterious  
74 mutations as assigned by prediction software (see Methods and **Supplementary**  
75 **Data Set 1**). *Hprt*, the gene encoding the sole 6-TG target<sup>10</sup> (**Fig. 1b**), was mutated  
76 in five of the seven sequenced clones (**Supplementary Data Set 1**). Moreover, it  
77 was the only gene mutated in multiple clones that carried likely deleterious mutations  
78 in all cases (**Fig. 2a**). Furthermore, these *Hprt* mutations affected different residues  
79 of the coding sequence (**Supplementary Data Set 1**). By contrast, only three non-  
80 synonymous mutations in other genes mutated in more than one clone were  
81 predicted to be deleterious, and no other gene contained a likely deleterious  
82 mutation in more than one clone (**Fig. 2a, Supplementary Data Set 1**). This  
83 analysis established that, without using any previous knowledge regarding the nature  
84 of suppressor loci, sequencing just a few clones identified *Hprt* as the top  
85 suppressor-gene candidate.

86

87 In addition to HPRT inactivation, mutations in genes for DNA mismatch repair (MMR)  
88 proteins confer 6-TG resistance<sup>11</sup>, as does inactivation of the DNA methyltransferase  
89 DNMT1<sup>12</sup>. Notably, the two whole-exome sequenced clones that did not carry *Hprt*  
90 mutations contained nonsense mutations in MMR genes (**Supplementary Data Set**  
91 **1, Supplementary Fig. 1d**). To further analyze coverage of our mutant libraries, we  
92 subjected the 189 additional suppressor clones we retrieved to targeted exon  
93 sequencing of the six known suppressor genes (**Fig. 1b**). With the exception of  
94 *Dnmt1* (see below), we identified predicted deleterious mutations in all known  
95 suppressor genes in homozygosis in two or more resistant clones (**Fig. 2b** top  
96 panels, **Supplementary Data Set 2**). Importantly, introducing wild-type versions of  
97 *Hprt* or *Mlh1* into resistant clones containing mutations in these genes restored 6-TG  
98 sensitivity (**Supplementary Figure 2**), confirming them as phenotypic drivers. Thus,  
99 if the non-targeted whole-exome sequence approach that we carried out in the initial  
100 analysis of seven clones had been applied to all 196 suppressor clones, *Hprt*, *Msh2*,

101 *Msh6*, *Mlh1* and *Pms2* would have been identified as suppressor gene candidates,  
102 confirming the feasibility of the approach to identify most or all resistance loci.

103  
104 Interestingly, ~20% (40) of clones presented two or more heterozygous deleterious  
105 mutations in the same suppressor gene (**Supplementary Data Set 2**). We note that  
106 haploid cell cultures cannot be maintained indefinitely and become diploid over  
107 time<sup>8,13</sup>. Accordingly, identified heterozygous mutations could have arisen after  
108 diploidization of the original EMS-treated haploid populations, or could have occurred  
109 in the small proportion of diploid H129-3 cells in the EMS-treated enriched haploid  
110 populations (**Fig. 1a**). Regardless of their origin, deleterious heterozygous mutations  
111 could only generate 6-TG resistance if each affected one allele of the gene,  
112 effectively inactivating both copies. Heterozygous mutations that we observed in  
113 *Dnmt1* occurred in such close proximity that they could be analyzed from the same  
114 sequencing reads. As we observed no co-occurrence in the same reads  
115 (**Supplementary Fig. 3a**), we concluded that *Dnmt1* mutants were compound  
116 heterozygotes, and confirmed this through Sanger sequencing (**Supplementary Fig.**  
117 **3b**). Furthermore, as these mutations all scored as potentially deleterious for DNMT1  
118 protein function (**Supplementary Data Set 2**), it is likely that they caused 6-TG  
119 resistance (see below). *Dnmt1* would thus be included in the list of suppressor gene  
120 candidates when considering deleterious heterozygous mutations. Furthermore, this  
121 analysis increased the numbers of clones identified with mutations in other  
122 suppressor loci (**Fig. 2b**, lower panels).

123  
124 Highlighting the applicability of our methodology to identify functionally important  
125 protein regions, we retrieved variants linked to *Hprt* mutations causative of Lesch-  
126 Nyhan syndrome<sup>14</sup>, as well as mutations in MMR genes linked to Lynch syndrome<sup>15</sup>  
127 (**Fig. 2c**). Partially reflecting the mutational preferences of EMS (see below), we  
128 found mRNA splicing variant mutations potentially affecting total protein levels  
129 (**Supplementary Data Set 2**). These were particularly prevalent in *Hprt* (**Fig. 2b**),  
130 and a detailed analysis confirmed their impacts on reducing HPRT protein levels  
131 (**Supplementary Figure 4**). These results highlight how production of aberrant  
132 mRNA splicing and associated reduction of protein product is an important  
133 consequence of EMS mutagenesis.

We also identified mutations that had not been previously reported, the majority of which were predicted to have deleterious effects on protein function (**Supplementary Fig. 5a, Supplementary Data Set 2**). To verify their impacts, we introduced newly identified MLH1 (A612T) and DNMT1 (G1157E) mutations into wild-type mESCs by CRISPR/Cas9 gene editing (**Supplementary Fig. 5b,c**). H129-3 mESCs carrying these mutations were more resistant to 6-TG than their wild-type counterparts (**Supplementary Fig. 5d**), supporting these mutations being causative of the suppressor phenotype. mESCs carrying targeted mutations in *Dnmt1* and *Mlh1* also allowed examination of their effects on cell proliferation. As observed under non-selective conditions, mutations in *Mlh1*, and especially in *Dnmt1*, impaired cell proliferation (**Supplementary Fig. 5e**), potentially helping to explain the low proportion of *Dnmt1* mutant suppressors arising from our screen. DNMT1-deficient cells exhibit 6-TG resistance, but the mechanism for this is not completely understood<sup>12,16</sup>. Our results point to an important role of Dnmt1 methyltransferase activity in mediating 6-TG sensitivity, as suppressor mutations identified in our screen localized to that domain (**Fig. 2c**). Collectively, these results further validated our pipeline to identify suppressor mutations.

Around 12% of resistant clones (23) did not present mutations in any of the known suppressor genes (**Fig. 2b**). We subjected these clones to whole-exome DNA and RNA sequencing. DNA sequencing of the unassigned clones and control samples allowed an unprecedented description of EMS mutagenic action, confirming its preference for producing SNVs and transition rather than transversion mutations (**Supplementary Fig. 6**). Although whole-exome sequencing retrieved causative mutations in all control 6-TG resistant samples, no other gene candidate could be identified from the remaining orphan suppressors (**Supplementary Data Set 3**). RNA sequencing, however, revealed reduced expression levels of *Hprt*, *Mlh1* or *Msh6* as likely causes of suppression in several such clones (**Fig. 2d; Supplementary Data Set 4**). Further studies will be required to define whether epigenetic alterations or mutations outside of exon regions, and hence not covered by exome-targeted DNA sequencing, could explain the nature of remaining orphan suppressor clones.

Collectively, our findings establish that classical genetic screening can be effectively performed in mammalian systems by combining use of haploid cells, chemical SNV induction, and next-generation sequencing. The use of haploid cells when creating SNV mutant libraries identifies recessive suppressor point-mutations, in contrast to diploid cell screening where only dominant mutations are retrieved<sup>17</sup>. Furthermore, EMS induction of SNVs generates complex mutant libraries, increasing the probability of identifying suppressor loci compared to isolation of rare, spontaneous suppressor events<sup>18</sup>. Through screening for cellular resistance to 6-TG, we identified point mutations in all described suppressor genes. This highlights the power of our approach to comprehensively identify suppressor loci with low error rates, as no false positive suppressor candidate genes were found. Moreover, as we have established for 6-TG suppressor loci, our methodology has value in delineating key amino-acid residues required for protein function, thus helping to explain molecular mechanisms of suppression. We note that SNV-based mutagenesis will be useful to identify separation-of-function and gain-of-function mutations, including those in essential genes. Also, through studies performed in cells bearing mutations in another gene, our approach has the potential to investigate gene-gene interactions in a comprehensive manner. In addition, we envisage the applicability of this approach in human haploid cells<sup>19,20</sup>. Chemical mutagenesis of haploid cells, either alone or in combination with LOF screens, has the potential to bring functional genomics in mammalian systems to a hitherto unachieved comprehensive level.

## **Methods**

Methods and associated references are available in the online version of this paper.

## Acknowledgements

We thank all S.P.J. lab members for discussions, especially A. Blackford, F. Puddu, C. Schmidt and P. Marco-Casanova for critical reading of the manuscript, and C. Le Sage and T.-W. Chiang for advice with CRISPR/Cas9 gene editing. We thank M. Leeb for H129-3 cells and advice on haploid ES cell culture conditions, and J. Hackett for advice in generating stable ES cell lines. We thank C.D. Robles-Espinoza for helping designing the array of baits for the exon-capture experiment, and J. Hewinson for technical support. Research in the S.P.J. laboratory is funded by Cancer Research UK (CRUK; programme grant C6/A11224), the European Research Council and the European Community Seventh Framework Programme (grant agreement no. HEALTH-F2-2010-259893; DDResponse). Core funding is provided by Cancer Research UK (C6946/A14492) and the Wellcome Trust (WT092096). S.P.J. receives salary from the University of Cambridge, supplemented by CRUK. J.V.F. was funded by Cancer Research UK programme grant C6/A11224 and the Ataxia Telangiectasia Society. J.C. was funded by Cancer Research UK programme grant C6/A11224. D.J.A. is supported by CRUK. Research leading to these results has received funding from the European Research Council under the European Union's Seventh Framework Programme (FP7/2007-2013)/ERC grant agreement no. [311166]. B.V.G. is supported by a Boehringer Ingelheim Fonds PhD fellowship.



#### **Author contributions**

J.V.F. and S.P.J. designed the project. J.V.F. mutagenized haploid cells, performed 6-TG selection and isolated suppressor clones. J.V.F. and J.C. expanded suppressor clones, isolated gDNA and prepared samples for sequencing. M.H. analyzed DNA sequencing data, supervised by T.M.K. and D.J.A. J.V.F. and J.C. produced stable cell lines and CRISPR/Cas9 knock-ins. J.V.F. and J.C. isolated RNA from suppressor clones and prepared samples for sequencing. B.V.G. produced RNA sequencing libraries and T.K. analyzed RNA sequencing data, supervised by S.M.N. J.V.F. and S.P.J. wrote the manuscript, with input from all authors.

## References

1. Forsburg, S. L. The art and design of genetic screens: yeast. *Nat Rev Genet* **2**, 659–668 (2001).
2. St Johnston, D. The art and design of genetic screens: *Drosophila melanogaster*. *Nat Rev Genet* **3**, 176–188 (2002).
3. Boutros, M. & Ahringer, J. The art and design of genetic screens: RNA interference. *Nat Rev Genet* **9**, 554–566 (2008).
4. Carette, J. E. *et al.* Haploid genetic screens in human cells identify host factors used by pathogens. *Science* **326**, 1231–1235 (2009).
5. Koike-Yusa, H., Li, Y., Tan, E.-P., Velasco-Herrera, M. D. C. & Yusa, K. Genome-wide recessive genetic screening in mammalian cells with a lentiviral CRISPR-guide RNA library. *Nat Biotechnol* **32**, 267–273 (2014).
6. Shalem, O. *et al.* Genome-scale CRISPR-Cas9 knockout screening in human cells. *Science* **343**, 84–87 (2014).
7. Rolef Ben-Shahar, T. *et al.* Eco1-dependent cohesin acetylation during establishment of sister chromatid cohesion. *Science* **321**, 563–566 (2008).
8. Leeb, M. & Wutz, A. Derivation of haploid embryonic stem cells from mouse embryos. *Nature* **479**, 131–134 (2011).
9. Munroe, R. R. & Schimenti, J. J. Mutagenesis of mouse embryonic stem cells with ethylmethanesulfonate. *Methods Mol. Biol.* **530**, 131–138 (2009).
10. LePage, G. A. & Jones, M. Purinethiols as feedback inhibitors of purine synthesis in ascites tumor cells. *Cancer Res* **21**, 642–649 (1961).
11. Swann, P. F. *et al.* Role of postreplicative DNA mismatch repair in the cytotoxic action of thioguanine. *Science* **273**, 1109–1111 (1996).
12. Guo, G., Wang, W. & Bradley, A. Mismatch repair genes identified using genetic screens in Blm-deficient embryonic stem cells. *Nature* **429**, 891–895 (2004).
13. Elling, U. *et al.* Forward and reverse genetics through derivation of haploid mouse embryonic stem cells. *Cell Stem Cell* **9**, 563–574 (2011).
14. Jinnah, H. A., De Gregorio, L., Harris, J. C., Nyhan, W. L. & O'Neill, J. P. The spectrum of inherited mutations causing HPRT deficiency: 75 new cases and a review of 196 previously reported cases. *Mutat Res* **463**, 309–326 (2000).
15. Jiricny, J. Postreplicative mismatch repair. *Cold Spring Harb Perspect Biol* **5**, a012633 (2013).

- 259 16. Loughery, J. E. P. *et al.* DNMT1 deficiency triggers mismatch repair defects in  
260 human cells through depletion of repair protein levels in a process involving  
261 the DNA damage response. *Hum Mol Genet* **20**, 3241–3255 (2011).
- 262 17. Kasap, C., Elemento, O. & Kapoor, T. M. DrugTargetSeqR: a genomics- and  
263 CRISPR-Cas9-based method to analyze drug targets. *Nat Chem Biol* **10**, 626–  
264 628 (2014).
- 265 18. Smurnyy, Y. *et al.* DNA sequencing and CRISPR-Cas9 gene editing for target  
266 validation in mammalian cells. *Nat Chem Biol* **10**, 623–625 (2014).
- 267 19. Blomen, V. A. *et al.* Gene essentiality and synthetic lethality in haploid human  
268 cells. *Science* **350**, 1092–1096 (2015).
- 269 20. Sagi, I. *et al.* Derivation and differentiation of haploid human embryonic stem  
270 cells. *Nature* **532**, 107–111 (2016).

271

## Figure legends

**Figure 1. Generation of mutagenized libraries. (a)** Experimental workflow. **(b)** Schematic of 6-TG metabolism and genotoxicity. Inactivating mutations in genes highlighted in red have been shown to confer 6-TG resistance. **(c)** Suppressor frequencies to 6-TG treatment of different EMS-mutagenized libraries, represented as number of suppressor clones isolated per 10,000 plated cells. **(d)** Locations and consequences of identified mutations.

**Figure 2. Identification of suppressor mutations. (a)** Genes identified through whole-exome sequencing of seven 6-TG resistant clones, harboring at least two independent mutations in different clones. Mutations were assigned as deleterious or neutral according to PROVEAN and SIFT software (see Methods). **(b) Top panels.** Distribution of mutations identified in suppressor gene candidates; numbers of independent clones are in brackets and types of *Hprt* mutations are shown in more detail on the pie-chart to the right. **Bottom panels.** Distribution of all suppressor gene candidate mutations identified, including heterozygous deleterious mutations. **(c)** Distribution of point mutations on DNMT1, HPRT and MMR proteins; each square represents an independent clone. Asterisks (\*) denote STOP-codon gains. Except for HPRT, all proteins are shown at a proportional scale. **(d) *Hprt*, *Mlh1* and *Msh6*** mRNA expression levels (fragments per kilobase per million reads). Numbers next to dots are clone identifiers (see Supplementary Data Set 2). Black dots indicate wild-type (WT) samples, red dots represent clonal samples whose mutations were identified via targeted exon capture sequencing (controls; see Supplementary Data Set 2), and white dots represent samples for which no causative mutations were identified. Error bars represent uncertainties on expression estimates. **Lower panel.** Reduced *Hprt* mRNA levels correspond to reduced protein production as detected by western blot. Uncut gel images are available in Supplementary Fig. 7.

## **Online Methods**

### **Cell lines and culture conditions**

H129-3 haploid mouse embryonic stem cells (mESCs)<sup>8</sup> were used for the experiments described in this paper. When pure haploid content was required, cells were grown in chemically defined 2i medium plus LIF as described previously<sup>8</sup>. In all other cases, cells were grown in DMEM high glucose (Sigma) supplemented with glutamine, streptomycin, penicillin, non-essential amino acids, sodium pyruvate,  $\beta$ -mercaptoethanol and LIF. All plates and flasks were gelatinized prior to cell seeding. All cells used in this study were mycoplasma free.

### **Cell sorting**

Cell sorting for DNA content was performed after staining with  $15\mu\text{g ml}^{-1}$  Hoechst 33342 (Invitrogen) on a MoFlo flow sorter (Beckman Coulter). The haploid 1n peak was purified. Analytic flow profiles of DNA content were recorded after fixation of the cells in ethanol, RNase digestion and staining with propidium iodide (PI) on a Fortessa analyzer (BD Biosciences). Cell cycle profiles were produced using FlowJo software (Tree Star).

### **Ethylmethanesulfonate (EMS) treatment**

Mutagenesis with EMS and measurement of killing and suppression frequency was performed as described previously<sup>9</sup>, with the following modifications. After cell sorting, haploid cells were grown in 2i medium plus LIF and changed to DMEM plus LIF for the overnight EMS treatment. After EMS treatment, cells were cultured for 5 passages in DMEM plus LIF and plated into 6-well plates at a density of  $5 \times 10^5$  cells per well. Cells were treated with  $2\mu\text{M}$  6-thioguanine (6-TG; Sigma) for 6 days, supplying new media with drug daily. Cells were then grown in medium without 6-TG until mESC colonies could be picked.

### **DNA isolation and exome sequencing**

mESC clones were grown into 12-well plates. Genomic DNA was extracted from confluent wells using QIAamp DNA Blood Mini Kit (QIAGEN) and cleaned performing a proteinase K (QIAGEN) digestion step. Genomic DNA (approximately  $1\mu\text{g}$ ) was

fragmented to an average size of 150 bp and subjected to DNA library creation using established Illumina paired-end protocols. Adapter-ligated libraries were amplified and indexed via PCR. A portion of each library was used to create an equimolar pool comprising 8 indexed libraries. For whole-exome sequencing, each pool was hybridized to SureSelect RNA baits (Mouse\_all\_exon; Agilent Technologies). Whole-exome sequencing was performed with 8 DNA samples per sequencing lane (first 7 suppressors plus control) or 15 DNA samples per sequencing lane (subsequent 66 suppressors analysed). For the exon-capture experiment, samples were hybridized with a specific array of RNA baits (Agilent) covering the exonic sequences of *Dnmt1*, *Hprt*, *Mlh1*, *Mlh3*, *Msh2*, *Msh3*, *Msh4*, *Msh5*, *Msh6*, *Pms1*, *Pms2* and *Setd2* genes. Sequence targets were captured and amplified in accordance with manufacturer's recommendations. Enriched libraries were subjected to 75 base paired-end sequencing (HiSeq 2500; Illumina) following manufacturer's instructions. A single sequencing library was created for each sample, and the sequencing coverage per targeted base per sample is given in Supplementary Data Set 5. All raw sequencing data is available from ENA under accession numbers ERP003577 and ERP005179.

#### **DNA sequence analysis**

Sequencing reads were aligned to the *Mus musculus* GRCm38 (mm10) assembly (Ensembl version release 68) using BWA (v0.5.10-tpx). All lanes from the same library were merged into a single BAM file with Picard tools (<http://broadinstitute.github.io/picard>), and PCR duplicates were marked by using 'MarkDuplicates'<sup>21</sup>. SNVs and INDELs were called using SAMtools (v1.3) mpileup followed by BCFtools (v1.3)<sup>22</sup>. The following parameters were used for Samtools mpileup: -g -t DP,AD -C50 -pm3 -F0.2 -d10000. BCFtools call parameters were: -vm -f GQ. The variants were annotated using the Ensembl Variant Effect Predictor<sup>23</sup>. Variants were filtered to remove any variants detected outside the bait regions and any heterozygous variants where appropriate. Additionally, variants were filtered using VCFtools (v0.1.12b) vcf-annotate with options -H -f +/-q=25/SnpGap=7/d=5 and custom filters were written to exclude variants with a GQ score of less than 10<sup>24</sup>. INDELs were left aligned using BCFtools norm. VCFtools vcf-isec was used to remove variants present in the control sample from all other samples as well as variants present in sequencing of a mouse strain from the 129S5 background<sup>25</sup>. INDELs called from whole exome sequencing data were further verified using the

microassembly based caller Scalpel<sup>26</sup> and discarded from the data if not identified by both callers. All remaining variants were used to generate a visualization of mutational patterns. All SNVs were assigned to one of 96 possible triplet channels using the GRCm38 assembly to identify flanking bases.

## **Antibodies**

Rabbit anti-HPRT (Abcam ab10479, 1: 10 000 dilution), mouse anti-MSH6 (BD Biosciences 610919, 1: 2 000), mouse anti-PMS2 (BD Biosciences 556415, 1: 1 000), rabbit anti-MRE11 (Abcam ab33125, 1: 10 000) and mouse anti-MLH1 (BD Biosciences 554073, 1: 1 000) were used for western blot analysis.

## **Complementation assays**

Human *MLH1* was amplified from pEGFP-MLH1<sup>27</sup> and cloned into pPB-CMV-HA-pA-IN<sup>28</sup> using *EcoRI* and *MluI* sites to generate pPB-Tet-MLH1. Cells from the SC\_6TG5758127 *Mlh1* mutant clone (see Supplementary Data Set 2) were transfected with a combination of pCMV-HyPBase<sup>29</sup>, pPB-CAG-rtTAM2-IP (a derivative of pPBCAG-rtTAIRESNeo<sup>28</sup> where the neomycin resistance cassette was replaced by a puromycin resistance one, gift from J. Hackett) and pPB-CMV-HA-pA-IN or pPB-Tet-MLH1, in a 1:1:10 ratio using TransIT-LT1 transfection reagent (Mirus) and following manufacturer's instructions. 48 h after transfection, selection was applied with 3 µg/ml puromycin for 6 days. Resistant cell populations were plated into 6-well plates (125 000 cells per well) and *MLH1* expression was induced by the addition of 1 µg/ml doxycycline. 24 h after doxycycline induction, cells were left untreated or treated with 2 µg/ml 6-TG for 6 days. Surviving cells were stained using crystal violet.

Cells from SC\_6TG5758069 and SC\_6TG5758117 *Hprt* mutant clones (see Supplementary Data Set 2) were transfected with pEGFP-C1 (Clontech) or pCMV6-AC-Hprt-GFP (OriGene MG202453) using TransIT-LT1 transfection reagent (Mirus) and following manufacturer's instructions. 48 h after transfection, selection was applied with 175 µg/ml G418 for several days, until GFP-positive colonies were picked. Cells were left untreated or treated with 2 µg/ml 6-TG for 6 days. Surviving cells were stained using crystal violet. Microscopy images were obtained from an Olympus IX71 microscope using Cell<sup>F</sup> imaging software (Olympus).

## **Prediction of mutation consequences on protein function**

Amino acid mutations were analysed using PROVEAN<sup>30</sup> and SIFT<sup>31</sup> software. Scores below -2.5 for PROVEAN and 0.05 for SIFT indicate likely deleterious effects.

## **Sanger sequencing**

PCR amplifications from genomic DNA were performed using the following oligonucleotides: Dnmt1-1157F 5'- CGAGATGCCTGGTAGACACA -3', Dnmt1 1157R 5'- GAGTAGGCCTGAGGAGAGCA -3', Dnmt1 1477F 5'- GCTACAAAACCCCAGGAAGC -3', Dnmt1 1477R 5'- CAGGATCAGATTGGCGTGAC -3'. PCR products from SC\_6TG5758159 and SC\_6TG5758161 *Dnmt1* mutant suppressors (see Supplementary Data Set 2) were cloned using Zero Blunt TOPO PCR cloning kit (Thermo Fisher Scientific) and following manufacturer's instructions.

## **Gene editing**

Sequences for DNA templates for small guide RNAs were generated using CRISPR Design (<http://crispr.mit.edu>) and cloned into pAiO-Cas9 D10A<sup>32</sup>. Sequences of the guides were the following: Dnmt1-1 5'- TCGGAAGGATTCCACCAAGC -3', Dnmt1-2 5'-ACATCCAGGGTCCGGAGCTT -3', Mlh1-1 5'- AGGACGACGGCCCGAAGGAA -3'; Mlh1-2 5'- GCCACTTTCAGGACTGTCTA -3'. H129-3 cells were transfected with Dnmt1 or Mlh1 targeting plasmids and single-stranded DNA oligonucleotides (200 nt, IDT Technologies) containing the desired mutations using *TransIT*-LT1 transfection reagent (Mirus) and following manufacturer's instructions. 48 h after transfection GFP-positive cells were sorted on a MoFlo flow sorter (Beckman Coulter) and seeded into gelatinized plates. Colonies forming after 5-6 days were picked into 96-well plates, DNA was isolated using QuickExtract DNA extraction solution (Epicentre Biotechnologies) and PCR amplifications of the edited regions were performed. Sequences of the oligonucleotides used were as follows: Dnmt1-F 5'- CGAGATGCCTGGTAGACACA -3', Dnmt1-R 5'- GAGTAGGCCTGAGGAGAGCA -3', Mlh1-F 5'- TGTCCCAACCTAGGGACTTG -3', Mlh1-R 5'- TGCTGGCCTTAGACAGTCCT -3'. PCR products (358 bp for *Dnmt1*, 287 bp for *Mlh1*) were digested with *EcoRI* restriction enzyme and run on a 3% agarose 1xTAE gel for 1.5 h at 150 V. Positive clones (those producing two DNA fragments after



*EcoRI* digestion of approx. 180 bp (*Dnmt1*) or 200 and 80 bp (*Mlh1*) were confirmed by Sanger sequencing of the PCR products, and tested for resistance to 6-TG as described for the screen.

### **Population doublings**

Each cell line was seeded in duplicate into 2 rows of a 24-well plate at a density of 25 000 cells/well. Cells were collected daily and cell counts were measured using a Countess II Automated Cell Counter (ThermoFisher Scientific) using Trypan Blue staining to discard dead cells.

### **RNA isolation and sequencing**

mESC clones were grown in 24-well plates. Total RNA was extracted from confluent wells using RNeasy Mini Kit (QIAGEN). Libraries for RNA-seq were prepared from 500 ng total RNA using the QuantSeq 3' mRNA-Seq kit (Lexogen) according to manufacturer's instructions. An exception to the instruction was the application of 13 instead of the recommended 12 PCR cycles for library amplification. Libraries were pooled in equal concentrations. Prior to sequencing, a T-fill reaction was performed on a cBot as described previously<sup>33</sup>, providing the T-fill solution in a primer tube strip. Finally, sequencing was carried out using an Illumina HiSeq-2500 using 50 bp single read v3 chemistry. Raw sequencing data is available from ENA under accession number ERP014134.

### **RNA sequence analysis**

Reads were trimmed of adapter sequences using Cutadapt (v.1.2.1). High-quality reads were extracted using TriageTools<sup>34</sup> (v0.2.2, long reads –length 35, high-quality bases –quality 9, and complex sequences –lzw 0.33). Alignments onto the mm10 genome were carried out using GSNAP<sup>35</sup> (v2014-02-28) with Gencode gene splice junctions. Expression levels were obtained using Exp3p (github.com/tkonopka/Exp3p v0.1) and then processed with custom R scripts (Supplementary Data Set 6).

### **Statistical analyses**

All groups analysed showed comparable variances.

### **Accession codes**

DNA sequencing data is available from ENA under accession numbers ERP003577 and ERP005179. RNA sequencing data is available from ENA under accession number ERP014134.

## Methods References

21. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
22. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
23. McLaren, W. *et al.* Deriving the consequences of genomic variants with the Ensembl API and SNP Effect Predictor. *Bioinformatics* **26**, 2069–2070 (2010).
24. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158 (2011).
25. Keane, T. M. *et al.* Mouse genomic variation and its effect on phenotypes and gene regulation. *Nature* **477**, 289–294 (2011).
26. Narzisi, G. *et al.* Accurate de novo and transmitted indel detection in exome-capture data using microassembly. *Nat Methods* **11**, 1033–1036 (2014).
27. Hong, Z. *et al.* Recruitment of mismatch repair proteins to the site of DNA damage in human cells. *J Cell Sci* **121**, 3146–3154 (2008).
28. Murakami, K. *et al.* NANOG alone induces germ cells in primed epiblast in vitro by activation of enhancers. *Nature* **529**, 403–407 (2016).
29. Yusa, K., Zhou, L., Li, M. A., Bradley, A. & Craig, N. L. A hyperactive piggyBac transposase for mammalian applications. *Proc Natl Acad Sci USA* **108**, 1531–1536 (2011).
30. Choi, Y. & Chan, A. P. PROVEAN web server: a tool to predict the functional effect of amino acid substitutions and indels. *Bioinformatics* **31**, 2745–2747 (2015).
31. Kumar, P., Henikoff, S. & Ng, P. C. Predicting the effects of coding non-synonymous variants on protein function using the SIFT algorithm. *Nat Protoc* **4**, 1073–1081 (2009).
32. Chiang, T.-W. W., le Sage, C., Larrieu, D., Demir, M. & Jackson, S. P. CRISPR-Cas9(D10A) nickase-based genotypic and phenotypic screening to enhance genome editing. *Sci Rep* **6**, 24356 (2016).
33. Wilkening, S. *et al.* An efficient method for genome-wide polyadenylation site

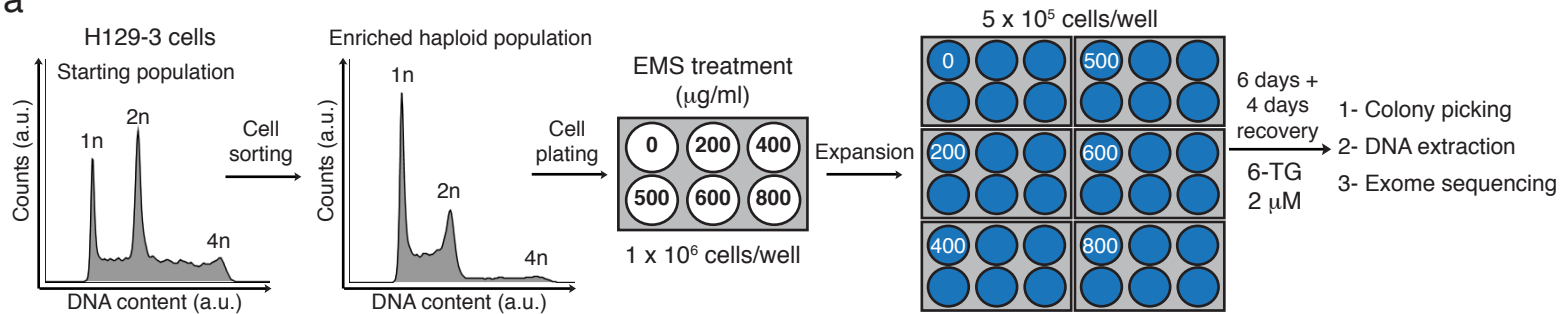
- 503 mapping and RNA quantification. *Nucleic Acids Res* **41**, e65–e65 (2013).
- 504 34. Fimereli, D., Detours, V. & Konopka, T. TriageTools: tools for partitioning and  
505 prioritizing analysis of high-throughput sequencing data. *Nucleic Acids Res* **41**,  
506 e86–e86 (2013).
- 507 35. Wu, T. D. & Nacu, S. Fast and SNP-tolerant detection of complex variants and  
508 splicing in short reads. *Bioinformatics* **26**, 873–881 (2010).

509

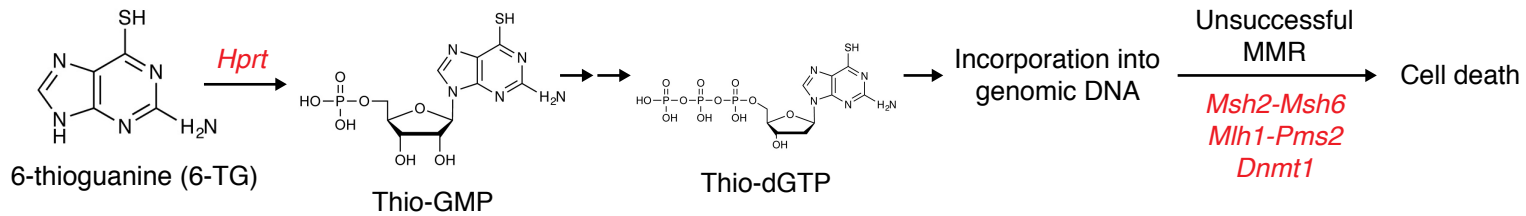
510 **Competing financial interests**

511 The authors declare no competing financial interests.

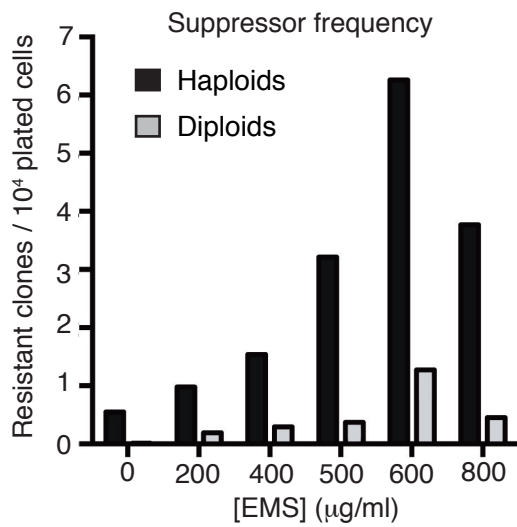
a



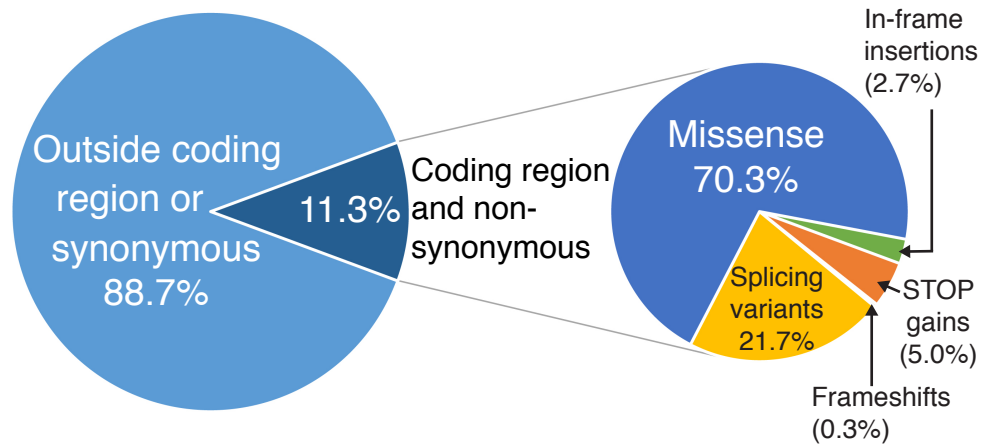
b



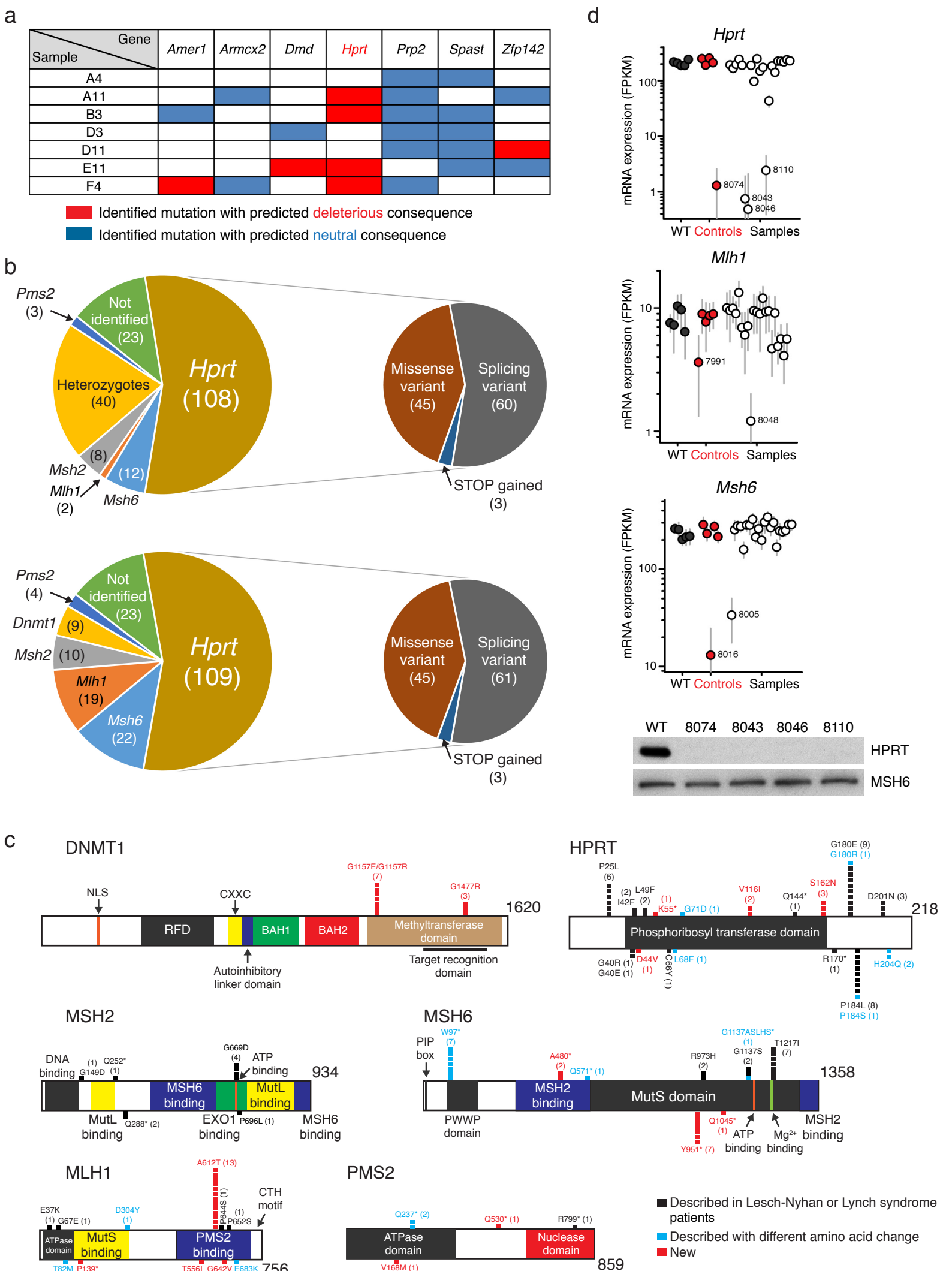
c



d



Forment et al, Figure 1



Forment et al, Figure 2

Mouse ES cells



Sorting



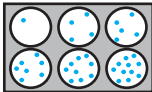
Haploid ES cells



EMS ( $\mu\text{g/ml}$ )

Toxic drug  
selection

6-TG



Next-generation  
sequencing

