1    Discovery of novel DNA cytosine deaminase activities enables a nondestructive

2    single-enzyme methylation sequencing method for base resolution high-coverage

3    methylome mapping of cell-free and ultra-low input DNA

4

5    Romualdas Vaisvila[1]#*, Sean R. Johnson[1]*, Bo Yan[1], Nan Dai[1], Billal M. Bourkia[1], Minyong Chen[1], Ivan R.

6    Corrêa Jr.[1], Erbay Yigit[1], and Zhiyi Sun[1]#

7        1.   New England Biolabs Inc., 240 County Road, Ipswich, MA 01938, United States

8    * These authors contributed equally

9    # Corresponding authors: vaisvila@neb.com and sunz@neb.com

# Abstract

11    Cytosine deaminases have important uses in the detection of epigenetic modifications and in genome

12    editing. However, the range of applications of deaminases is limited by a small number of well

13    characterized enzymes. To expand the toolkit of deaminases, we developed an in-vitro approach that

14    bypasses a major hurdle with their severe toxicity in expression hosts. We systematically assayed the

15    activity of 175 putative cytosine deaminases on an unprecedented variety of substrates with

16    epigenetically relevant base modifications. We found enzymes with high activity on double- and single-

17    stranded DNA in various sequence contexts including novel CpG-specific deaminases, as well as enzymes

18    without sequence preference. We also report, for the first time, enzymes that do not deaminate

19    modified cytosines. The remarkable diversity of cytosine deaminases opens new avenues for

20    biotechnological and medical applications. Using a newly discovered non-specific, modification-sensitive

21    double-stranded DNA deaminase, we developed a nondestructive single-enzyme 5-methylctyosine

22    sequencing (SEM-seq) method. SEM-seq enables accurate, high-coverage, base-resolution methylome

23    mapping of scarce biological material including clinically relevant cell-free DNA (cfDNA) and single-cell

24    equivalent 10 pg input DNA. Using SEM-seq, we generated highly reproducible base-resolution 5mC

25    maps, accounting for nearly 80% of CpG islands for a low input human cfDNA sample offering valuable

26    information for identifying potential biomarkers for detection of early-stage cancer and other diseases.

27    This streamlined protocol will enable robust, high-throughput, high-coverage epigenome profiling of

28    challenging samples in research and clinical settings.

# Introduction

30    Cytosine deaminases are widespread enzymes that are involved in numerous important cellular

31    processes. In eukaryotes, the APOBEC (Apolipoprotein B mRNA Editing Catalytic Polypeptide-like) family

32    of proteins plays important roles in antibody diversification and innate immunity against retroviruses[1].

33    In bacteria, deaminases are found in polymorphic toxin systems (PTS), which are multi-domain proteins

34    involved in intra- and interspecies competition[2–4]. Deaminases which act on nucleic acids have been

35    used in many biotechnological applications, including as key components of base editing tools[5–7] and in

36    assays to detect various DNA and RNA modifications[8–12].

37 Recently, APOBEC3A deaminase was featured as a non-destructive enzymatic alternative to bisulfite
38 conversion to detect 5-methylcytosine (5mC) and 5-hydroxymethylcytosine (5hmC) DNA modifications
39 (EM-seq[9] , LR-EM-seq[10] and ACE-seq[11]). APOBEC3A deaminates cytosine (C) to uridine (U) in single-
40 stranded DNA (ssDNA). APOBEC3A also deaminates 5mC and 5hmC, albeit with reduced activity, but
41 does not deaminate 5-carboxylcytosine (5caC) and glucosylated 5hmC (5gmC). In ABOBEC3A-based
42 modification detection assays, 5mC can be protected through conversion to 5caC and 5gmC by
43 combined Tet methylcytosine dioxygenase 2 (TET2) and T4-phage beta-glucosyl transferase (T4-BGT)
44 activity.

45 Compared to bisulfite conversion-based sequencing methods[13], enzymatic deamination protocols do not
46 damage DNA, require lower amounts of input DNA, produce less biased data, and are more compatible
47 with long read sequencing and enrichment of long amplicons[9–11]. We envisioned that the discovery of
48 deaminases with new properties would open the possibility of leaping beyond some of the limitations of
49 current enzymatic methods. Specifically, sequence-agnostic robust activity on double-stranded DNA
50 (dsDNA) combined with a lack of activity on 5mC and 5hmC would enable a streamlined, one-step, one-
51 enzyme protocol for 5mC mapping.

52 About a decade ago, Iyer, Zhang, and colleagues reported extensive computational analyses of the
53 diversity and phylogenetic distribution of enzymes with the deaminase fold, including a large variety of
54 putative cytosine deaminases found in bacterial polymorphic toxin systems[2,3]. For many years, most of
55 that sequence space remained unexplored experimentally. Recently, Mok, de Moraes, and colleagues
56 described the characterization, engineering, and application to base editing of DddA, of a cytosine
57 deaminase from bacterial toxin systems acting on dsDNA substrates[5–7]. We considered DddA unsuitable
58 for epigenetics modification detection due to its overall low activity and strong sequence preference.
59 Nevertheless, inspired by these earlier works, we endeavored to discover and characterize a broader
60 range of putative cytosine deaminases and assess their suitability for cytosine modification detection
61 and other applications.

62 In this study, we report an extensive survey of the enzymatic activity of cytosine deaminases from
63 bacterial polymorphic toxin systems, phages, and gene cassettes. We expressed deaminases using an in-
64 vitro system that circumvents their well-documented host toxicity and assayed 175 enzymes from 13
65 deaminase families. Combining Liquid Chromatography–Mass Spectrometry (LC-MS) and Next
66 Generation Sequencing (NGS) methods, we measured their deaminase activity and substrate selectivity
67 on ssDNA, dsDNA, using unmodified and modified substrates incorporating 5-methylcytosine (5mC), 5-
68 hydroxymethylcytosine (5hmC), 5-formylcytosine (5fC), 5-carboxycytosine (5caC), 5-
69 glucosyloxymethylcytosine (5gmC), and N4-methylcytosine (N4mC). Our work uncovered bacterial
70 deaminases with diverse activities, including enzymes with a wide range of sequence preferences and
71 modification sensitivities on ssDNA and dsDNA. We identified a subset of enzymes with properties
72 desirable for improving epigenetic modification detection methods. Most notably, we found a
73 modification-sensitive deaminase that is active on dsDNA without sequence constraints. We
74 demonstrated the application of this enzyme in a new one-tube deaminase-mediated sequencing
75 method, SEM-seq, for human methylome profiling.

# Results

## In-vitro screening of 175 deaminase sequences

78  We developed a bioinformatics strategy to look for putative DNA deaminase sequences from public
79  databases. Then, we devised an experimental pipeline to test deamination activity in-vitro (Fig. 1a).
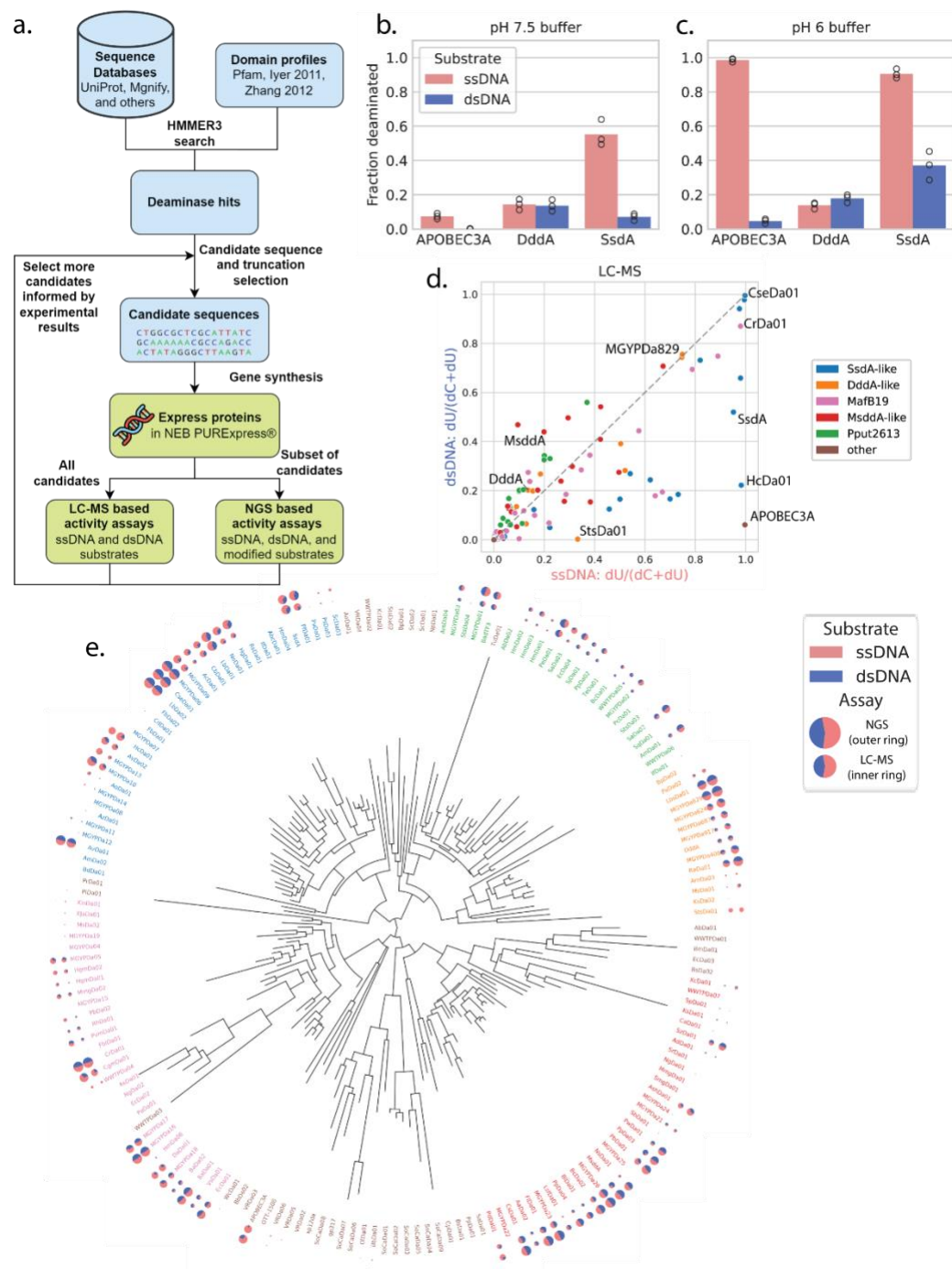
80  Using as queries, HMMER3[14] deaminase family profiles from Pfam[15] and reports by Iyer, Zhang and
81  colleagues[2,3] (Supplementary Fig. 1), we searched for new deaminases in six different databases, with
82  most candidates coming from UniProt[16] or Mgnify[17]. We picked hits from diverse deaminase families
83  with the intent to cover a broad range of the sequence space and thus catalog deaminase activities.
84  Candidate selection was performed over many iterations. Active enzymes in our initial screen primarily
85  derived from bacterial polymorphic toxin systems. Later rounds focused on sequences in the same
86  families as those of the active enzymes from previous selections.

87  Because of their inherent toxicity to expression hosts, deaminases are often obtained in low yields, or
88  are subjected to mutations that inactivate enzymatic activity. As such, putative deaminase proteins were
89  expressed using an in-vitro protein synthesis system. An LC-MS-based assay was developed to quantify
90  deaminase activities on unmodified cytosines in ssDNA, dsDNA, and RNA substrates. The single-stranded
91  ΦX174 Virion DNA and the double-stranded ΦX174 RF I DNA, sharing the same template sequence were
92  used as DNA substrates. These 5,386-nucleotide (nt) length substrates present C in a rich diversity of
93  sequence contexts, thus allowing the detection of deaminase activities regardless of the sequence
94  specificities. Isotope labeled firefly luciferase (Fluc) mRNA was used as the substrate for testing
95  deaminase activity on RNA for a subset of deaminases.

96  We validate our screening assay on the previously published SsdA and DddA enzymes. Considering the
97  potential mutagenicity of DNA deaminases, we compared the use of the PURExpress in-vitro protein
98  synthesis system with different transcription templates. These were composed of either unmodified
99  DNA or of DNA with 5hmdC in place of dC to potentially block template mutation during the
100  transcription and translation reactions. Of note, DNA templates constructed from the 5hmdCTP and
101  dCTP generated similar results (Supplementary Fig. 2a). Further investigation showed that deaminase
102  activity is inhibited by the 1X PURExpress buffer (Supplementary Fig. 2b) but can be rescued in the
103  diluted buffer. Our data suggested this is a viable approach to efficiently produce active proteins of toxic
104  deaminases for further characterization. As a case in point, we showed by LC-MS analysis that SsdA, like
105  APOBEC3A, prefers ssDNA, and DddA is active on dsDNA (Fig. 1b and Supplementary Fig. 2 c,d), which is
106  in agreement with the published results. We further uncovered that all these three enzymes have equal
107  or better activity in a lower pH buffer (pH 6)[18] than that of a previously published buffer (pH 7.5)[7] (Fig.
108  1b,c).

109  We screened a total of 175 new enzymes across 13 cytosine deaminase families (Fig. 1d,e,
110  Supplementary Fig. 1). Our assay revealed that many bacterial DNA deaminases act on both ssDNA and
111  dsDNA, with some deaminating both types of substrates with equal efficiency. Five deaminase families,
112  DddA-like, SsdA-like, Pput2613, MafB19, and MsddA-like, accounted for nearly all enzymes with DNA
113  activity detectable by LC-MS (Fig. 1d,e). In four of the five families, at least one deaminase was found
114  with higher RNA deaminase activity than APOBEC3A  (Supplementary Fig. 3), which has previously been
115  reported to have activity on RNA substrates[19]. Examples of enzymes highly active on DNA were CseDa01,
116  MGYPDa06, and LbDa02, which deaminated close to 100% of cytosines in dsDNA and ssDNA substrates;

117    CrDa01, which deaminated 87% and 97% of cytosines in dsDNA and ssDNA, respectively; and
118    MGYPDa829 and LbsDa01, which deaminated about 75% of cytosines in both dsDNA and ssDNA
119    substrates. When we tested at lower concentrations, MGYPDa829 (DddA-like) showed equal activity on
120    ssDNA and dsDNA, whereas CseDa01 (SsdA-like) preferred ssDNA (Supplementary Fig. 4). Overall, most
121    of the active bacterial deaminases display a preference for ssDNA substrates. While we found a few
122    cytosine deaminases that strongly prefer single-stranded substrates, such as HcDa01 and StsDa01, we
123    did not find a single bacterial deaminase that only acts on dsDNA. For simplicity, herein we will use the
124    term "dsDNA deaminase" for enzymes that act on dsDNA, even though they also deaminate ssDNA to
125    some extent.

126

**Figure 1. Screen design and results overview.** a) Schematic of candidate selection and screening strategy. b) Activity of control enzymes expressed in PURExpress and assayed in pH 7.5 buffer. c) Activity of control enzymes expressed in PURExpress and assayed in pH 6 buffer. d) Activities of all screened enzymes, as measured by the LC-MS assay. e) Maximum likelihood tree of screened enzymes with activities on unmodified ssDNA (pink) and dsDNA (blue), as measured by the LC-MS assay (inner ring) and NGS assay (outer ring) indicated by area on pie charts. Enzymes that showed low activity in the LC-MS assay were not assayed by NGS. DNA emoji designed by OpenMoji – the open-source emoji and icon project. License: CC BY-SA 4.0.
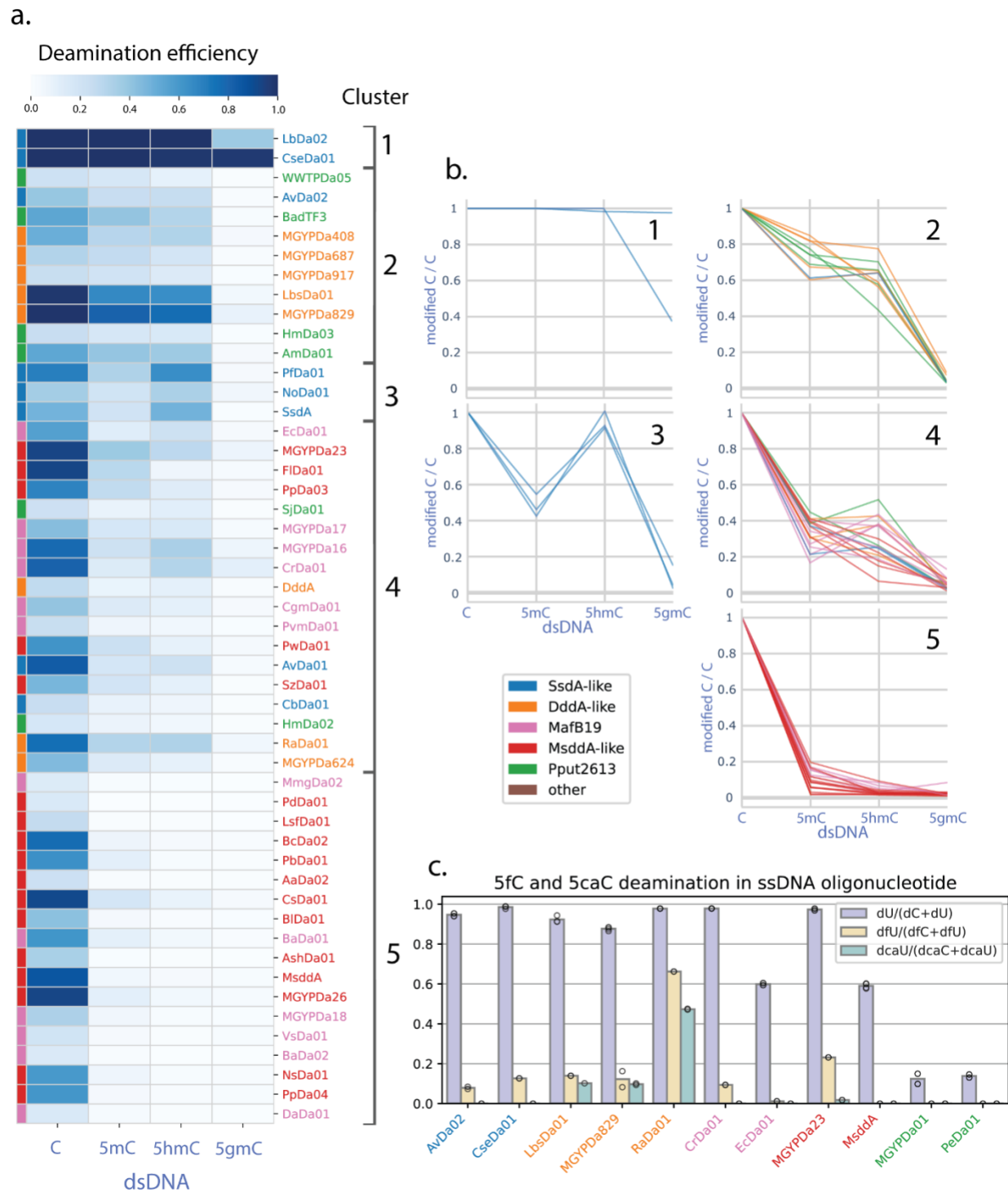
5

## A wide spectrum of sensitivity to DNA cytosine modifications

We next measured deaminase activity on a variety of physiologically relevant cytosine modifications. We developed an NGS assay to characterize deaminase sequence preference and sensitivity to 5C-methylation, 5C-hydroxymethylation, 5C-glucosylation, and N4-methylation (N4mC) in both single-stranded and double-stranded DNA. In addition, we used the LC-MS assay to examine activity on 5C-formylated, and 5C-carboxylated substrates. To validate these methods, we showed that the deamination efficiency on unmodified substrates measured by NGS was broadly consistent with that of LC-MS (Pearson r = 0.85 and 0.86 for ssDNA and dsDNA respectively) (Supplementary Fig. 5).

We found that the screened enzymes generally had low activity on N4mC (Supplementary Fig. 6) but displayed a wide spectrum of activities on different modification types at the cytosine 5-position (Fig. 2, Supplementary Fig. 7). We clustered the deaminases into five functional clusters based on their specificity for C, 5mC, 5hmC and 5gmC in dsDNA. Most deaminases showed a moderate decrease of activity toward 5mC and 5hmC modifications, and a strong decrease of activity toward 5gmC compared to unmodified dsDNA substrates (Fig. 2a,b, clusters #2 and 4). A subset of deaminases displayed similar activities across all the cytosine types (Fig. 2a,b, cluster #1). This group includes CseDa01, which efficiently deaminates C, 5mC, 5hmC, and even gmC to nearly 100%. Another interesting subset showed decreased activity on 5mC and 5gmC, but not on 5hmC (Fig. 2a,b, cluster #3). This subset, which includes PfDa01 and SsdA, can potentially be used to distinguish 5hmC from 5mC. Enzymes from clusters #1-4 tended to have more relaxed modification sensitivity to ssDNA compared to dsDNA (Supplementary Fig. 7). The most exciting discovery was of a unique functional cluster with very low or no activity on modified substrates (Fig. 2a,b, cluster #5). The enzymes from cluster #5 with the highest deamination activity of unmodified C were found to be phylogenetically related (Fig. 1e). We referred to the most active enzyme of this group (92% of unmodified cytosine deaminated and less than 3% of the modified cytosines), originating from human mouth metagenome (European Nucleotide Archive accession ERZ773077)[20], as "Modification-Sensitive DNA Deaminase A" (MsddA) and the family of deaminases related to it as MsddA-like.

We also examined deamination activities on 5fC and 5caC of representative enzymes from each family using ssDNA oligonucleotide substrates containing four cytosines in a 'Cg' context in which all cytosines were either unmodified, 5-formylated, or 5-carboxylated (Fig. 2c, Supplementary Fig. 8). We observed that MsddA, EcDa01 (MafB19 family), and the two enzymes from the Pput2613 family did not deaminate either modified substrate. In contrast, all other tested deaminases displayed some activity on 5fC, and the three enzymes from the DddA-like family had activity on both 5fC and 5caC.
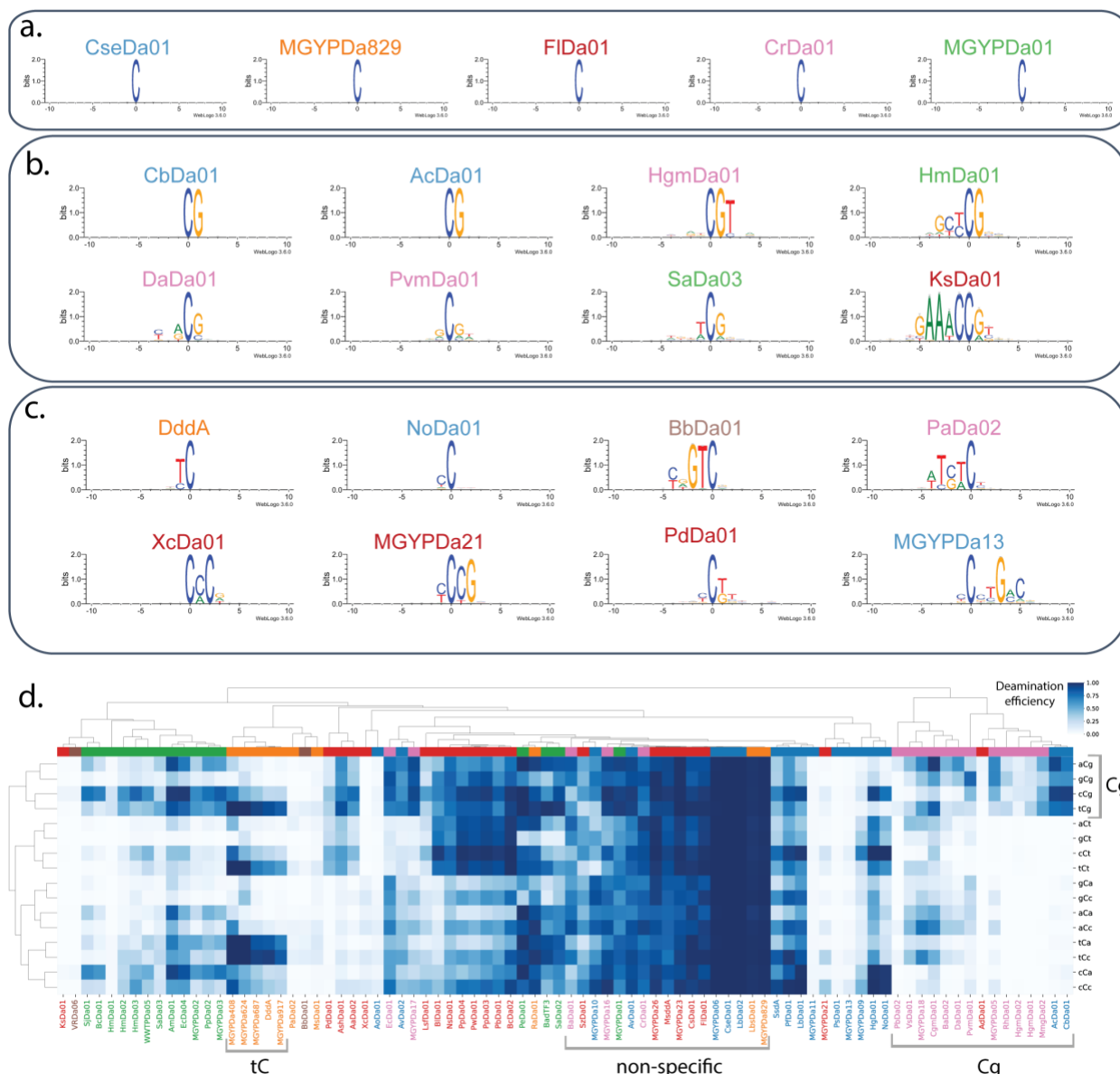
**Figure 2. Deamination efficiency of representative enzymes on C-5 modified substrates in dsDNA.** a) Deamination efficiency on C, 5mC, 5hmC, and 5gmC in dsDNA, measured by NGS assay. The enzymes were grouped into five clusters, using average linkage clustering of cosine distances between deamination efficiencies. b) Deamination efficiency on dsDNA C, 5mC, 5hmC, and 5gmC divided by deamination efficiency on C of enzymes in the five clusters. c) Deamination efficiency on 5fC and 5caC in ssDNA oligonucleotides by representative deaminases from each family, measured by LC-MS assay.

## A wide variety of DNA sequence preferences

The single-base resolution of our NGS assay and the use of a large size and sequence complexity of substrates enabled us to accurately survey the sequence preference landscape for our deaminase library. Sequence Logo plotting[22] of 10-base context on either side of sites deaminated in unmodified dsDNA uncovered a diversity of context preferences (Fig. 3a-c, Supplementary Fig. 9), including a subset of enzymes that display no sequence constraints (Fig. 3a). The non-specific enzymes also tend to have high overall activity (Fig. 3d). In addition to displaying different preferences for sequences at the 5' end position(s) to the deaminated cytosine, reminiscent of enzymes from the eukaryotic APOBEC deaminase family[18], many bacterial deaminases display sequence selectivity for the immediate 3' position, some specifically recognizing CpG dinucleotides (Fig. 3b,d). Deaminase recognition sequence sites can extend beyond the nCn context, with preferences for sequences of various lengths and compositions (Fig. 3c).

As a more systematic analysis of sequence preference across different substrate types, we calculated the average deamination efficiency of all nCn sequences in unmodified and modified dsDNA and ssDNA substrates and conducted clustering analysis across diverse enzymes and specificities (Fig. 3d, Supplementary Fig. 10). We observed that enzymes from the same family tend to cluster together. For example, MafB19 deaminases tend to prefer Cg, DddA-like deaminases tend to prefer tC, and MsddA-like deaminases tend to have reduced activity on Ca. However, all five families contained examples of non-specific deaminases (Fig. 3a). Specificities were generally consistent across substrates, noting that enzymes tend to display more context specific preferences on substrates on which they had lower total activity (Supplementary Fig. 10).

**Figure 3. Sequence context preference of representative enzymes.** a) Logos of example non-specific deaminases from each family. b) Logos of example deaminases with preference for G at the +1 position, including CpG-specific deaminases. c) Logos of example deaminases with diverse sequence specificities. d) Activity on representative deaminases in the nCn contexts of unmodified dsDNA. Rows and columns are sorted based on average linkage clustering of cosine distances. All logos are of sites with >=50% deamination efficiency in unmodified dsDNA. Deamination efficiencies measured by NGS assay. Names are colored according to family: SsdA-like (blue), DddA-like (orange), MafB19 (pink), MsddA-like (red), Pput2613 (green), other (brown).

9

## SEM-seq, a nondestructive streamlined single-enzyme methylation sequencing method for accurate base-resolution methylome analysis

The discovery of the non-specific methylation-sensitive dsDNA deaminase MsddA enabled us to perform SEM-seq, a single-enzyme method for methylation sequencing at single base resolution. SEM-seq eliminated the need for TET2/T4-BGT protection and denaturing steps that are required of APOEC3A-based protocols (Fig. 4a).

We validated SEM-seq using genomic DNA from the human GM12878 cell line and benchmarked it against two published datasets from standard EM-seq[9] and WGBS (ENCODE accession ENCSR890UQO)[23] protocols. Unmethylated lambda and CpG-methylated pUC19 DNA were spiked in as controls to estimate deamination rates. Two technical replicates of SEM-seq were sequenced on the Illumina Nova-seq platform. 277 million paired-end reads from two replicates of each method were used for the comparative analyses.
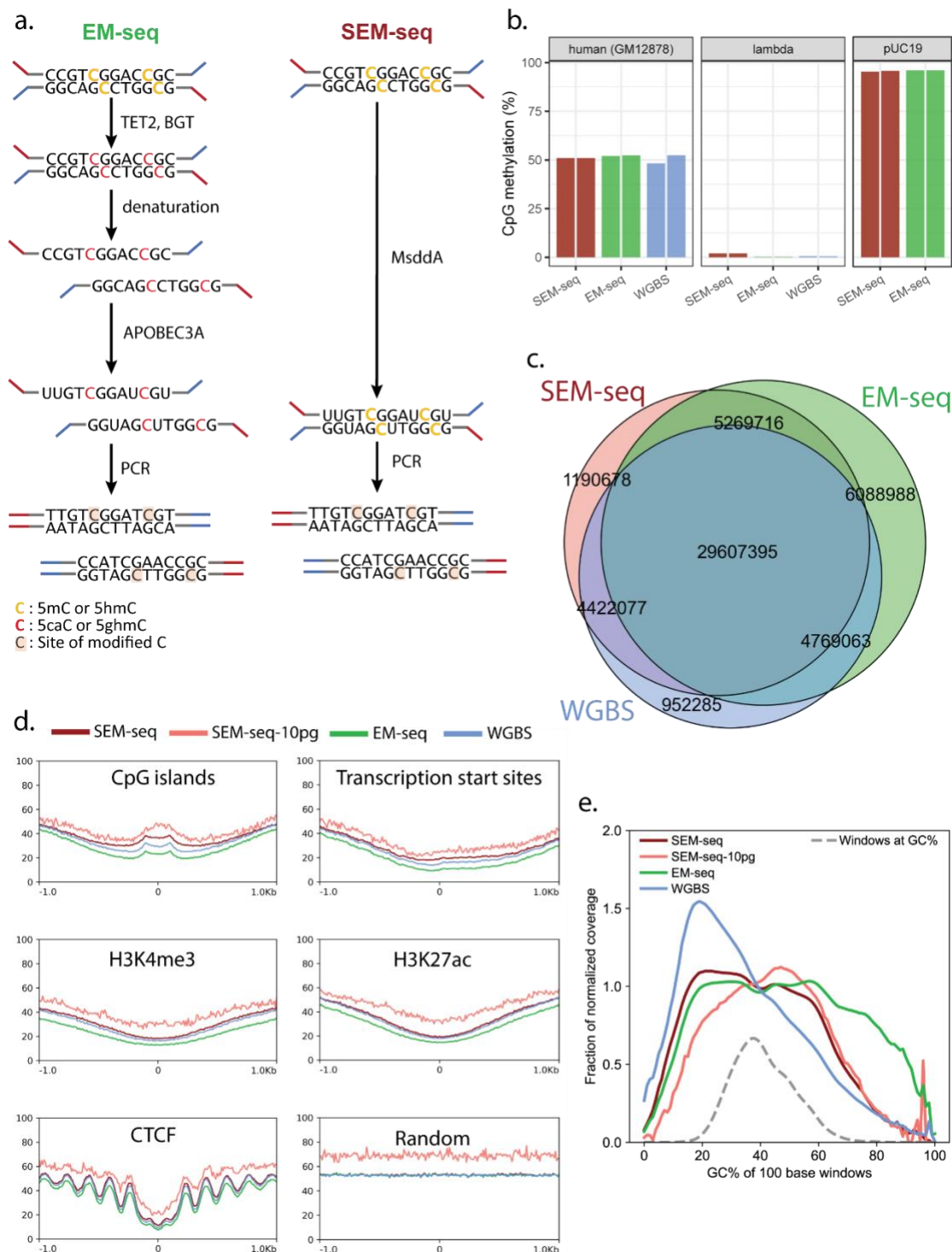
The library quality metrics and methylation results were highly consistent between technical replicates, suggesting a high reproducibility and robustness of SEM-seq (Supplementary Fig. 11). SEM-seq libraries produced very similar CpG methylation results for human GM12878 as those of EM-seq and WGBS (Fig. 4b-d), despite a slightly lower CpG deamination rate on unmodified cytosines for lambda phage DNA (SEM-seq 98.0%, EM-seq 99.8%; WGBS 99.4%). The 5mC non-conversion rates calculated from the CpG-methylated pUC19 DNA were 95.4% and 95.7% for each of the two SEM-seq replicates. These data were nearly identical to that of EM-seq (~96%) (Fig. 4b), but without requiring protection by TET2 and BGT. The conversion rates of unmodified pUC19 CH sites were 96.3% and 95.7% for the two SEM-seq replicates. Due to generally low methylation level of CH sites in the human genomes, we focused our methylome analysis on CpG sites. After combining the data of the two technical replicates, SEM-seq identified 40.4 million high-confidence methylated CpG sites, of which 86% and 84% agreed with EM-seq and WGBS datasets, respectively (Fig. 4c). The methylation levels of commonly methylated CpG sites are also highly correlated between SEM-seq and the other two methods (Pearson correlations: SEM-seq vs. EM-seq 0.91, SEM-seq vs. WGBS 0.88, EM-seq vs. WGBS 0.90). SEM-seq accurately produced the expected CpG methylation profiles for key epigenomic features, such as an enrichment of methylation in CpG islands, a reduced methylation levels near transcription start sites and at active chromatin markers and enhancers, and a typical regularly spaced oscillation pattern surrounding CTCF binding sites with depleted methylation at the center[24] (Fig. 4d). Furthermore, SEM-seq, like EM-seq, is non-destructive and gives more even read coverage across genomic regions of variable GC contents than the bisulfite-based WGBS method (Fig. 4e).
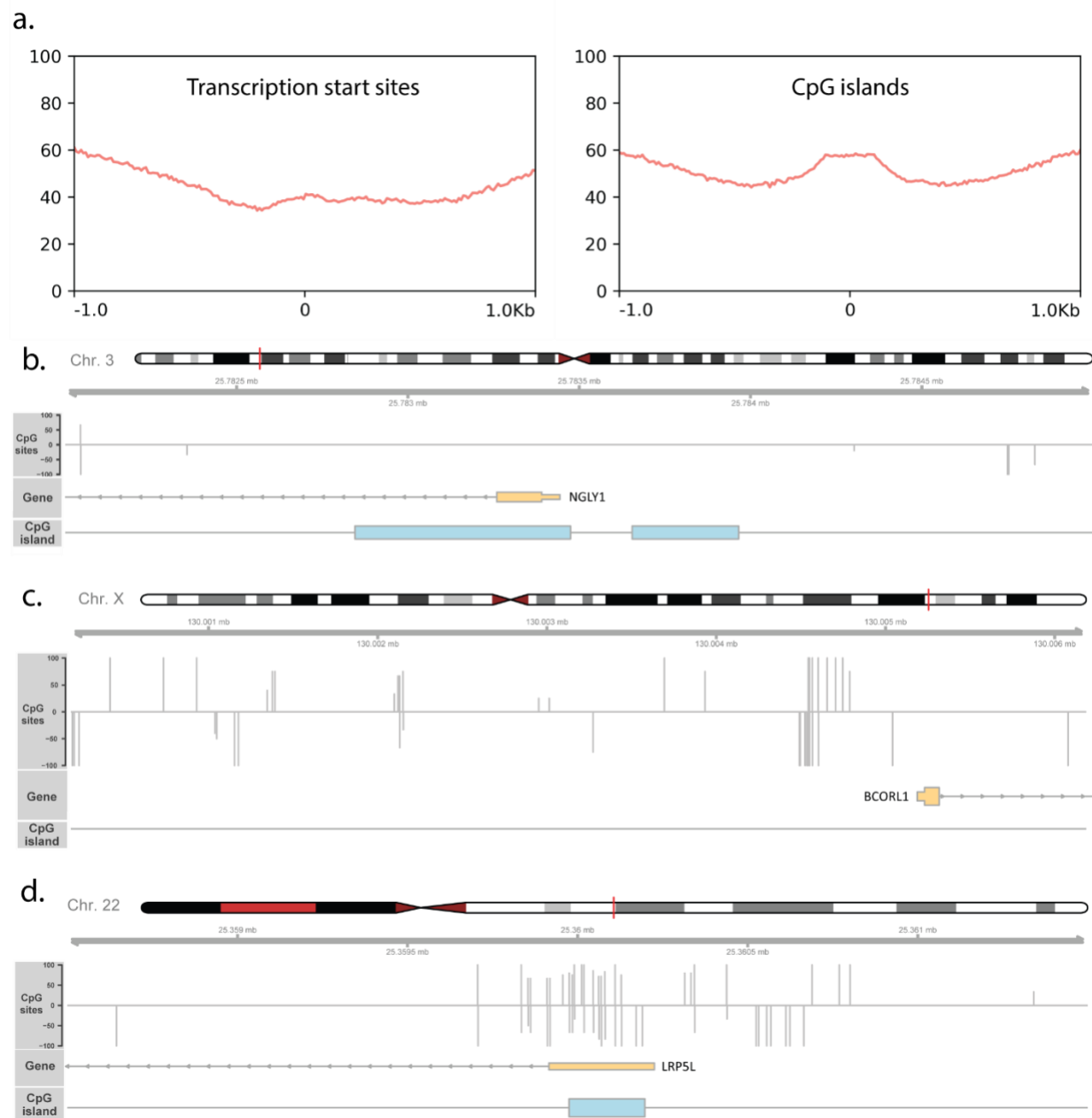
## SEM-seq for cfDNA and 10 pg input

The streamlined, robust protocol of SEM-seq makes it an advantageous method for obtaining accurate and highly reproducible methylome information from scarce biological samples such as cell free DNA (cfDNA), which has important clinical applications for noninvasive prenatal diagnosis and early cancer detection and monitoring[25]. We applied SEM-seq to human cell free DNA (cfDNA) and sequenced two replicate SEM-seq libraries each using 3 ng of cfDNA, generating 187.0 million and 204.5 million paired-end reads. The two cfDNA libraries covered 26.6 million and 30.6 million CpG sites respectively with at least 3X coverage. Among them, 25.6 and 29.4 million methylated CpG sites were identified with 18.5 million agreeing between the two replicates (Supplementary Fig. 12a). The methylation levels of individual CpG sites also correlated well between the replicates (Pearson correlation=0.71)

247    demonstrating accurate, base resolution quantification. A single cfDNA SEM-seq library generated
248    accurate 5mC profiles in epigenetic important genomic features resembling those of the human
249    GM12878 genomic DNA (Fig. 5a, Supplementary Fig. 12b). Remarkably, SEM-seq was able to provide
250    single-base-resolution 5mC quantification of 21536 CpG islands accounting for 78% of annotated CpG
251    islands in the human genome, leading to the detection of 11426 hypomethylated (methylation level
252    <20%) and 6712 hypermethylated (methylation level >70%) CpG islands and their associated genes
253    (Supplementary Fig. 12c, Supplementary Table S1). As expected, gene promoter regions are depleted in
254    5mC globally. SEM-seq allowed us to further examine the distribution and methylation level of individual
255    cytosine sites of each gene at a finer resolution. Consistent with the global trend, many promoter
256    regions are hypomethylated, for example the N-glycanase 1 (NGLY1) promoter (Fig. 5b). However, a
257    subset of genes are heavily modified in the promoter, such as BCL6 corepressor like 1 (BCORL1) (Fig.
258    5c)—a transcriptional corepressor that was suggested a prognostic factor that promotes tumor
259    metastasis[26], and LDL receptor related protein 5 like (LRP5L), of which a CpG island near the TSS is also
260    hypermethylated (Fig. 5d). Such high-resolution methylation information is very valuable for detecting
261    aberrant DNA methylation in cancer and other diseased samples and identify potential biomarkers for
262    detection and classification of early-stage cancer.

263    We also made SEM-seq libraries from only 10 pg of human GM12878 DNA, equivalent to single cell
264    input[27]. With 28 million paired-end reads, a single 10 pg library covered 9.3 M CpG sites and 8219 (30%)
265    CpG islands. It revealed similar methylation patterns to the 50 ng EM-seq and SEM-seq libraries, with
266    low GC coverage bias (Fig. 4e), and the expected 5mC distributions at various key genomic features (Fig.
267    4d). This result suggests that SEM-seq is suitable for single-cell methylome studies.

**Figure 4. SEM-seq Human GM12878 DNA compared to EM-seq, and whole genome bisulfite sequencing (WGBS).** a) Schematic of EM-seq and SEM-seq protocols. b) Total CpG methylation in GM12878, lambda (negative internal control) and CpG-methylated pUC19 (positive internal control) DNA. c) Counts of methylated CpG sites called in common by the three methods. d) Comparison of CpG methylation profiles in neighborhoods around specific genomic features. e) Effects of GC-content on read coverage from the three methods.

12

**Figure 5. SEM-seq on human cfDNA sample**. a) CpG methylation profiles in neighborhoods around transcription start sites (TSS) and CpG islands. b) 1.5 Kb upstream and 1 Kb downstream flanking region of the TSS of NGLY1 (hypomethylated). c) 5 Kb upstream and 1Kb downstream flanking region of the TSS of BCORL1 (hypomethylated). d) 1.5 Kb upstream and 1 Kb downstream flanking region of the TSS of LRP5L (hypomethylated). For b), c), and d), the methylation levels of CpG sites are shown as vertical grey lines. Positive and negative values indicate CpG sites on the top strand and bottom strand, respectively. Coordinates correspond to positions on respective chromosomes (GRch38).

13

## Discussion

284

285 This work presents a systematic experimental screening of 175 new cytosine deaminases spanning 13
286 distinct deaminase families. Screened enzymes cover most deaminase families that have been
287 hypothesized to act on cytosines[2,3]. Nearly all the active bacterial deaminases we characterized are
288 found in polymorphic toxin systems (PTS) and come from five of the 13 families: SsdA-like, DddA-like,
289 MafB19, Pput2613, and MsddA-like. The last of which we christened in reflection of the sensitivity to
290 modified cytosines displayed by many of its members.

291 Our screen reveals that bacterial deaminases are a versatile group of enzymes with diverse and
292 previously unknown properties. We show that many bacterial deaminases act on dsDNA, contrasting
293 with earlier observations that most DNA cytosine deaminases strongly prefer ssDNA substrates. While
294 some bacterial DNA deaminases are strictly single-strand specific enzymes, which parallels all the known
295 eukaryotic DNA deaminases, we did not find any deaminases that only accept double-strand DNA
296 substrates. This suggests that the deamination activity on single-stranded substrates may be the
297 ancestral state of all the bacterial deaminases.

298 Several recent studies have reported a functional characterization of enzymes from the DddA-like
299 family[28–30]. One study[28] showed that deletion of SPKK-related motifs at the C-terminus of DddA
300 abrogates its activity on dsDNA. In agreement with those results, we found that StsDa01, a deaminase
301 from the DddA-like family, lacks a C-terminal SPKK-related motif and strongly prefers ssDNA substrates.
302 Furthermore, we found that by truncating the C-terminus of MGYPDa829, this enzyme was converted
303 from having similar activity on both dsDNA and ssDNA to strongly preferring ssDNA (Supplementary Fig.
304 13). The enzyme resulting from a swap of the C-terminus from MGYPDa829 onto StsDa01 retained a
305 strong preference for ssDNA. We observed that dsDNA deaminases outside the DddA-like family lack
306 SPKK-related motifs, therefore the association between dsDNA activity and the SPKK-related motif does
307 not seem to generalize beyond the DddA-like family.

308 Bacterial deaminase activities also vary widely across DNA modification types including 5mC, 5hmC, 5fC,
309 5caC, 5gmC, and N4mC. In contrast to a wide spectrum of activities on modifications at C-5 (Fig. 2),
310 including the first report of deaminases with strong activity on 5caC and 5gmC, the screened enzymes
311 generally had low activity on N4mC (Supplementary Fig. 6). Finally, we found that bacterial deaminases
312 display a broad spectrum of sequence specificities, including enzymes with no apparent biases, enzymes
313 with diverse sequence preferences at the 5' end position relative to the target C, and enzymes with
314 strong sequence preferences at the 3' end positions, particularly a strong bias toward CpG
315 dinucleotides.

316 The diverse properties of bacterial deaminases lend themselves to various applications, including as
317 base editors and for epigenetic modification mapping. Information on substrate preference serves as a
318 key for selecting suitable enzymes for different applications. In base editing applications, their small size
319 (100-200 amino acids) (Supplementary Fig. 14) could provide advantages for therapeutic delivery. Also, a
320 strong ssDNA preference coupled with different context specificities will likely help reduce off-target
321 editing. For detecting the important epigenetic mark 5mC, a sequence independent cytosine deaminase
322 that can discriminate between methylated and unmethylated sites would be ideally suited for broad
323 methylome analysis. Indeed, we demonstrate the successful application in human whole methylome
324 mapping of the non-specific, modification-sensitive dsDNA deaminase MsddA. Our streamlined
325 methylation sequencing protocol, SEM-seq, not only is free of the damage-inducing bisulfite treatment,

326     but also does not require any accessory protection proteins (TET2 and T4-BGT) nor harsh denaturing
327     steps. We leveraged the advantages of SEM-seq by successfully applying it to methylome analysis in
328     clinically relevant cfDNA and in challenging, single-cell equivalent, 10 pg input DNA samples.

329     Our study unravels another example of the great potential for developing biotechnological tools from
330     previously unexploited bacterial and bacteriophage enzymes, in our case DNA methylation mapping
331     using bacterial DNA deaminase enzymes with diverse substrate preferences and modification
332     sensitivities.

333

## Limitations of the Study

335     Deaminases screened in this study or used for the SEM-seq experiments were all produced using the
336     PURExpress in-vitro transcription/translation (IVTT) system without further purification. The lack of
337     purified and quantified deaminases limited our ability to optimize SEM-seq conditions, for example by
338     fine-tuning enzyme concentration. The current SEM-seq conversion rate of unmodified CH is around
339     96%, which is lower than unmodified CH conversion rates typically observed in EM-seq and WGBS
340     experiments. This CH conversion rate needs to be improved in order to accurately detect low level non-
341     CpG methylations in most mammalian samples. Future work with purified deaminases will focus on
342     optimizing SEM-seq conditions for improved sensitivity of 5mC detection in both CpG and non-CpG
343     contexts.

## Methods

### Candidate selection

To guide the search for new enzymes, we first curated a list of HMMER3[14] cytosine deaminase sequence profiles. 29 profiles came from the CDA clan (CL0109) from the Pfam[15] database (version 34) (excluding the TM1506, LpxI_C, FdhD-NarQ, and AICARFT_IMPCHas, which are thought to not encode deaminases), 17 profiles were built from multiple sequence alignments (MSAs) of deaminase families defined by Iyer et al. (2011)[2], and one profile was built from a multiple sequence alignment found in Zhang et al. (2012)[3]. The profiles from Iyer largely overlap with the Pfam profiles, with a few differences, despite some similar names (Supplementary Fig. 1). The Pfam MafB19-deam (PF14437) profile and the similarly named profile from Iyer were found to be different from each other. None of the enzymes we screened had a best profile match to the Pfam MafB19-deam profile, so our usage of MafB19 corresponds to the Iyer profile.

The Pfam DYW_deaminase (PF14432) profile is biased towards eukaryotic members of that family, whereas the similarly named profile from Iyer captures both eukaryotic and bacterial DYW deaminases. We therefore split the Iyer DYW profile into three separate profiles, the one comprising mostly bacterial enzymes we called SsdA-like, as it contains the previously described SsdA deaminase[7], the other, which more closely resembles the DYW profile from Pfam, we called Iyer2011_DYW, the combined profile we called Iyer2011_DYW_combined (giving a new total of 18 profiles from Iyer). We found experimentally that the enzymes with strong matches towards the profile from Zhang et al. 2012[3] tended to have poor deamination activity on modified cytosines. We therefore renamed this profile MsddA-like, for "modification sensitive DNA deaminase A-like".

Some candidate sequences were selected directly from the MSAs listed in Iyer et al. (2011) and Zhang et al. (2012). Others were selected from hmmsearch hits of the profiles described above against six different databases: UniProt[16], Mgnify[17], IMG/VR[31], IMG/M[32], gene cassette metagenomes[33], wastewater treatment plant metagenomes[34], and GenBank[35]. Candidate selection and experimental screening were performed over many iterations. Candidates were selected from the hmmsearch hits with the intent to cover a broad range of sequence space, but with a focus on sequences similar to those that were shown to be active in earlier rounds. Mok et al. (2020)[5] and de Moraes et al. (2021)[7] previously reported on the active bacterial DNA deaminases, SsdA, DddA, and BadTF3, so similarity to those enzymes was also used as a criterion for candidate selection.

Most of the deaminases we tested were found as fusions to larger proteins, for example as parts of polymorphic toxin systems. In our assays we only expressed the deaminase domain rather than the full proteins. To determine the boundaries of the deaminase domain, we ran AlphaFold2[36] via the LocalColabFold[37] package on the deaminase domain plus up to about 1000 amino acids upstream and downstream. We visualized the AlphaFold2 predicted structures in PyMOL. N-terminal truncation sites were generally selected at several amino acids before helix 1 of the deaminase domain. The deaminase domains were typically found at the C-terminus of the fusion protein, so for most sequences C-terminal truncations were not necessary. For cases where the boundaries of the deaminase domain were difficult to distinguish, we either tried several different truncations (Supplementary Fig. 15) or relied on the Predicted Aligned Error (PAE) metric reported by AlphaFold2. We found that residues with a low PAE to residues in the core of the deaminase domain corresponded to our intuitive notions of the boundaries of the deaminase domain. To aid in the visualization of PAE, we developed a Colab notebook called

16

386 PAEView
387 (https://github.com/seanrjohnson/proteinotes/blob/master/colab_notebooks/PAEView.ipynb).

## Sequence naming

389 For convenience, each screened sequence was given a short name. The names are related to the
390 database or species of origin for the sequence. Da = deaminase, MGYP = Mgnify protein, Hm = hot
391 metagenome, VR = IMG/VR, WWTP = wastewater treatment plant, SoCa = soil gene cassette, chimera =
392 chimeric sequence. Other prefixes are mostly two or three letters drawn from the name of the source
393 organism or the source environment of the metagenome data. Some sequences also have prefixes or
394 suffixes of the form extN#, extC#, d#, Cd#, which indicate, respectively, N-terminal extensions, C-
395 terminal extensions, N-terminal deletions, and C-terminal deletions of the indicated number of residues,
396 compared to the candidate with the un-affixed name.

## Phylogenetic Tree

398 Amino acid sequence alignments were all calculated using MUSCLE[38]. Trees were generated using raxml-
399 ng (v. 1.1)[39]. The tree was rooted at the midpoint and rendered using ETE3[40].

## Hmm profile comparison

401 Similarity between deaminase hmm profiles was computed using hmmer_compare.py, a new python
402 implementation of the algorithm described by Söding[41]. Scores were calculated for each pair of profiles.
403 Scores were converted to distances using the formula:

405 distance = 1 - (pairwise_score / min(profile1_self_score, profile2_self_score))

407 A UPGMA tree was generated from the distance matrix using scikit-learn[42] and rendered with Geneious
408 Prime (https://www.geneious.com/) to visualize distances between profiles.
409 Code for generating hmmer3 profile trees and alignments is available at:
410 https://github.com/seanrjohnson/hmmer_compare

## Sources of DNA and RNA substrates

413 Unmodified DNA oligonucleotides were purchased from IDT Coralville, IA), 5caC and 5fC containing
414 oligonucleotides were synthesized by NEB (Ipswich, MA). dsDNA oligonucleotides were annealed in 10
415 mM Tris–HCl, pH 8.0 buffer. *E. coli* C2566 genomic DNA was purified using Monarch Genomic DNA
416 Purification Kit (NEB Ipswich, MA), GM12878 genomic DNA was obtained from Coriell Cell Repositories
417 (Camden, NJ). Single donor human plasma was obtained from Innovative Research (Novi, MI) and cfDNA
418 was extracted using the BioChain (Newark, CA) cfPure MAX V2 Cell-Free DNA Extraction Kit.
419 Unmethylated cl857 Sam7 Lambda DNA (Promega, Fitchburg, WI), fully C-methylated XP12 phage DNA
420 (Dr. Yan-Jiun Lee, NEB Ipswich MA), fully C-hydroxymethylated T4147[43] and fully C-hydroxy-methyl-beta-
421 glucosylated T4 alpha-glucosyl-transferase knockout (AGT-) genomic DNAs (Dr. Lidija Truncaite,  Vilnius
422 University Life Sciences Center, Vilnius Lithuania), pRSSM1.PleII plasmid (containing N4mC at cacCgc
423 sites) (Dr. Iain Murray, NEB Ipswich MA) were used as spike in controls in NGS deamination assay.

424 Isotope labeled firefly luciferase (Fluc) mRNA was synthesized by HiScribe T7 High Yield RNA Synthesis
425 Kit (E2040S, NEB, Ipswich) using stable heavy isotopes of ATP and CTP (NLM-3987-CA and CNLM-4267-
426 CA-20, Cambridge Isotope Laboratory, Tewksbury, MA). Resulting Fluc mRNA was purified twice by
427 Monarch RNA cleanup kit (T2040S, NEB, Ipwsich, MA) to ensure that unincorporated nucleotides were

428  completely removed. RNA was eluted in nuclease-free water, quantified by Qubit RNA broad range
429  assay (Q10211, ThermoFisher, Waltham, MA) and stored at -20 °C.

## In-vitro expression of deaminases

431  The candidate DNA deaminase genes first were codon-optimized then added flanking sequences
432  containing T7 promoter at 5' end (GCGAATTAATACGACTCACTATAGGGCTTAAGTATAAGGAGGAAAAAAT)
433  and T7 terminator at 3' end (CTAGCATAACCCCTCTCTAAACGGAGGGGTTTATTTGG) and ordered as linear
434  gBlocks from IDT (Coralville, IA, USA). Template DNA for in vitro protein synthesis was generated with
435  Phusion® Hot Start Flex DNA Polymerase (NEB. Ipswich MA) using gBlocks as template and flanking
436  primers (Supplementary Table S2). The PCR products were purified using Monarch PCR and DNA
437  Cleanup kit (NEB, Ipswich, MA). DNA concentration was quantified using a NanoDrop
438  spectrophotometer (Thermo Fisher Scientific, Inc., Waltham, MA). 100 - 400 ng of PCR fragments were
439  used as template DNA to synthesize analytic amounts of DNA deaminases using PURExpress In Vitro
440  Protein Synthesis kit (NEB, Ipswich, MA) following manufacturer's recommendations.

## Deamination activity assay on single and double stranded ΦX174 DNA substrates

442  To test the activity of in-vitro expressed DNA deaminases, 2 µL of PURExpress sample was mixed with
443  300 ng of ΦX174 Virion DNA (ssDNA substrate, NEB Ipswich MA) or ΦX174 RF I DNA (dsDNA substrate,
444  NEB, Ipswich MA) into buffer containing 50 mM Bis-Tris pH 6.0, 0.1% Triton X-100 and incubated for 1 h
445  at 37°C for a total volume of 50 µL. The deaminated ΦX174 DNA was purified using Monarch PCR and
446  DNA Cleanup kit (NEB, Ipswich, MA). DNA concentration was quantified using a NanoDrop
447  spectrophotometer (Thermo Fisher Scientific, Inc., Waltham, MA). 150 ng of deaminated DNAs were
448  digested to nucleosides with the Nucleoside Digestion Mix (NEB, Ipswich, MA) following manufacturer's
449  recommendations. LC-MS/MS analysis was performed by injecting digested DNAs on an Agilent 1290
450  Infinity II UHPLC equipped with a G7117A diode array detector and a 6495C triple quadrupole mass
451  detector operating in the positive electrospray ionization mode (+ESI). UHPLC was carried out on a
452  Waters XSelect HSS T3 XP column (2.1 × 100 mm, 2.5 µm) with a gradient mobile phase consisting of
453  methanol and 10 mM aqueous ammonium acetate (pH 4.5). MS data acquisition was performed in the
454  dynamic multiple reaction monitoring (DMRM) mode. Each nucleoside was identified in the extracted
455  chromatogram associated with its specific MS/MS transition: dC [M+H]$^+$ at m/z 228.1→112.1; dU
456  [M+H]$^+$ at m/z 229.1→113.1; d$^m$C [M+H]$^+$ at m/z 242.1→126.1; and dT [M+H]$^+$ at m/z 243.1→127.1.
457  External calibration curves with known amounts of the nucleosides were used to calculate their ratios
458  within the samples analyzed.

## Deamination activity assay on single and double-stranded oligonucleotide DNA substrates

461  To test the activity of in vitro expressed DddA and SsdA DNA deaminases on oligonucleotide substrates
462  (Supplementary Table S3), a 2 µL of PURExpress sample was mixed with 300 ng of 44 bp single-stranded
463  or double-stranded DNA substrates (annealed in 10 mM TRIS ph=8.0) with single cytosine in the contexts
464  TC, AC, CC, and GC in a buffer consisting of 20 mM Tris-HCl pH 7.5, 200 mM NaCl, 1 mM DTT for a total
465  volume of 50 µL. Cytosine deamination to uracil percentage was measured with LC-MS/MS as described
466  above.

## Synthesis of 5fC and 5caC modified DNA substrates

Oligonucleotides with 5fC and 5caC modified bases were synthesized in-house using standard phosphoramidite chemistry and supplier deprotection protocols (Glen Research, Sterling, VA). Control oligonucleotide (4CpG_contr) was purchased from Integrated DNA Technologies (IDT, Coralville, IA).

## Deaminase activity assay on 5fC and 5caC modified DNA

To test activity on 5fC and 5caC modified DNA, we designed a set of substrates (40 bp) containing 4 modified cytosines (Supplementary Table S4). The set of oligonucleotides included preferable deamination sites for 11 DNA deaminases from 5 representative clades. In each reaction, we mixed each modified oligonucleotide (dcaC or dfC) with the control oligonucleotide (C only) in a ratio of 1:1 (800 ng+800 ng) to monitor deamination of cytosine to uracil. After incubation for 5 h at 37°C in - reaction buffer containing 50 mM Bis-Tris pH 6.0, 0.1% Triton X-100 with different DNA deaminases, DNA was purified using Monarch PCR and DNA Cleanup kit, digested to nucleosides with the Nucleoside Digestion Mix (NEB, Ipswich MA) and the reaction products were quantified with LC-MS/MS.

## Deamination activity assay on RNA substrate

To test the activity of in-vitro expressed deaminases on RNA, 2 µL of freshly made PURExpress sample was mixed with 200 ng of heavy isotope labeled Fluc mRNA into buffer containing 50 mM Bis-Tris pH 6.0, 0.1% Triton X-100 and incubated for 1 h at 37°C for a total volume of 50 µL. Then, deaminase treated RNA was purified with NEBNext sample purification beads (E7104S, NEB, Ipswich, MA) using 1.2x beads to sample volume ratio, eluted in 40 µL nuclease free water. RNA sample was filtrated using a 0.22 µm centrifugal filter (UFC30GV00, MilliporeSigma, Darmstadt, Germany) to ensure all the beads were removed from the sample. Purified and filtrated RNA was then digested to single nucleoside by Nucleoside Digestion Mix (M0649S, NEB, Ipswich, MA). LC-MS/MS analysis was performed by injecting digested stable isotope labeled RNAs on an Agilent 1290 Infinity II UHPLC equipped with a G7117A diode array detector and a 6495C triple quadrupole mass detector operating in the positive electrospray ionization mode (+ESI). UHPLC was carried out on a Waters XSelect HSS T3 XP column (2.1 × 100 mm, 2.5 µm) with a gradient mobile phase consisting of methanol and 10 mM aqueous ammonium acetate (pH 4.5). MS data acquisition was performed in the dynamic multiple reaction monitoring (DMRM) mode. Each nucleoside was identified in the extracted chromatogram associated with its specific MS/MS transition: rC* [M+H]+ at m/z 256.1→119.1; rU* [M+H]+ at m/z 256.1→119.1  (*: Stable isotope labeled nucleosides) External calibration curves with known amounts of the nucleosides were used to calculate their ratios within the samples analyzed. The conversion rate was calculated using the following formula:

labeled Uridine / (labeled Uridine + labeled C)

## DNA deamination NGS assay

50 ng of unmodified *E. coli* C2566 genomic DNA was combined with control DNAs with various DNA modification types (Supplementary Table S5) in 50 µL of 10 mM Tris pH 8.0.

Then the DNA was transferred to a Covaris microTUBE (Covaris, Woburn, MA) and sheared to 300 bp using the Covaris S2 instrument. 50 µL of sheared material was transferred to a PCR strip tube to begin library construction. NEBNext DNA Ultra II Reagents (NEB, Ipswich, MA) were used according to the manufacturer's instructions for end repair, A-tailing, and adaptor ligation. The custom made Pyrollo-dC adaptor (NEB Organic Synthesis Division, Ipswich MA), where all dCs are replaced with Pyrollo-dC, was

508    used: (ACACTCTTTCCCTACACGACGCTCTTCCGATC*T and
509    [Phos]GATCGGAAGAGCACACGTCTGAACTCCAGTCA). The ligated samples were mixed with 110 µL of
510    resuspended NEBNext Sample Purification Beads and cleaned up according to the manufacturer's
511    instructions. The library was eluted in 17 µL of water. For ssDNA libraries, prior to deamination the DNA
512    was denatured by heating at 90°C for 10 minutes followed by cooling for 2 minutes on ice. The DNA was
513    then deaminated in 50 mM Bis-Tris pH 6.0, 0.1% Triton X-100, using 1 µL of dsDNA deaminase
514    synthesized as described above with an incubation time of 1 hour at 37°C. After deamination reaction, 1
515    µL of Thermolabile Proteinase K (NEB, Ipswich, MA) was added and incubated additional 30 min at 37°C
516    followed by 10 min at 60°C. 5µM of NEBNext Unique Dual Index Primers and 25 µL NEBNext Q5U Master
517    Mix (New England Biolabs, Ipswich, MA, USA) were added to the DNA and PCR amplified. The PCR
518    reaction samples were mixed with 50 µL of resuspended NEBNext Sample Purification Beads and
519    cleaned up according to the manufacturer's instructions. The library was eluted in 15 µL of water. The
520    libraries were analyzed and quantified by High sensitivity DNA analysis using a chip inserted into an
521    Agilent Bioanalyzer 2100. The libraries were sequenced using the Illumina NextSeq and NovaSeq
522    platforms. Paired-end sequencing of 150 cycles (2 x 75 bp) was performed for all the sequencing runs.
523    Base calling and demultiplexing were carried out with the standard Illumina pipeline.

## Measure deaminase activity on N4mC

525    Deaminase activity on N4mC was measured from all the cacCgc sites in the genome of pRSSM1.PleII
526    plasmid, of which the capitalized C is N4 methylated. To avoid the confounding effect of sequence
527    selectivity, we only considered enzymes that have a minimum 10% activity on unmodified cacCgc sites in
528    the *E. coli* DNA and calculated a relative deamination efficiency value, which is the percentage of one
529    enzyme's activity on N4mC compared to that on unmodified C in the same sequence context.

## SEM-seq library preparation of human genomic DNA

531    50 ng of human GM12878 genomic DNA spiked with 0.15 ng CpG-methylated pUC19, 1 ng unmethylated
532    lambda DNA and 1 ng methylated XP12 DNA were sonicated with the Covaris S2 instrument, end
533    repaired and ligated to the Pyrollo-dC adaptor as described above. The deamination of the resulting
534    ligation product was performed with 1 µL of dsDNA deaminase (MsddA) in Low pH buffer (50 mM Bis-
535    Tris pH 6.0, 0.1% Triton X-100) with the addition of 50mM NaCl and 1mM DTT in 20 µL total reaction
536    volume for 3 h at 37 °C. The reaction was stopped by adding 1 µL of Thermolabile Proteinase K (NEB,
537    Ipswich, MA) and incubating for 30 min at 37°C then an additional 10 min at 60°C. 15 µL of reaction
538    without purification was directly used for PCR amplification (6 cycles) in 50 µL total volume and cleaned
539    up as mentioned above. The whole genome libraries were sequenced on the Illumina NovaSeq 6000
540    platform (100 bp paired-end).

## Low input SEM-seq libraries

542    50 ng of human GM12878 genomic DNA spiked with 0.15 ng CpG-methylated pUC19, 1 ng unmethylated
543    lambda DNA and 1 ng methylated XP12 DNA were sonicated with the Covaris S2 instrument. Sheared
544    DNA was diluted to a concentration of 0.2 pg/µL, from which 50 µL (10 pg) of were used for end repair
545    and adaptor ligation as described above. For adaptor ligation Pyrollo-dC adaptor was diluted 25 times.
546    The deamination of the resulting ligation product was performed with 1 µL of dsDNA deaminase
547    (MsddA) in Low pH buffer (50 mM Bis-Tris pH 6.0, 0.1% Triton X-100) with the addition 15 ng of Lambda
548    DNA (Promega, Fitchburg, WI) in 20 µL total reaction volume for 3 h at 37 °C. The reaction was stopped
549    by adding 1 µL of Thermolabile Proteinase K (NEB, Ipswich, MA) and incubating for 30 min at 37°C then

550    an additional 10 min at 60°C. 15 μL of reaction without purification was directly used for PCR

551    amplification (12 and 14 cycles) in 50 μL total volume and cleaned up as mentioned above two times.

552    The whole genome libraries were sequenced on the Illumina NovaSeq 6000 platform (100 bp paired-

553    end).

554

## Data analysis

### Read processing and mapping

557    Raw Illumina reads were first trimmed by the Trim Galore software

558    (https://github.com/FelixKrueger/TrimGalore) to remove adapter sequences and low-quality bases from

559    the 3' end. Unpaired reads due to adapter/quality trimming were also removed during this process. The

560    trimmed read sequences were C to T converted and were then mapped to the corresponding reference

561    using the Bismark program[44] with the default Bowtie2 settings[45]. For the NGS deamination assay, a

562    composite reference was created by combining the complete genome sequences of *E. coli* C2566, phage

563    lambda, phage XP12, phage T4, plasmid pRRSlac, and Adenovirus. For human libraries made by SEM-seq,

564    EM-seq and WGBS, a composite reference included the human genome (GRCh38) and the complete

565    sequences of unmethylated lambda and CpG-methylated pUC19 controls.

566    The first 5 bp at the 5' end of R2 reads were removed to reduce end-repair errors and aligned read pairs

567    that shared the same alignment ends were regarded as PCR duplicates and were discarded. For human

568    libraries, aligned reads that contained excessive cytosines in non-CpG context (more than 4 in 100bp)

569    were removed before calculating methylation level for individual sites because these reads likely contain

570    high conversion errors.

571    The numbers of Ts (converted not methylated) and Cs (unconverted modified) of each covered cytosine

572    position were then calculated from the remaining good quality alignments using Bismark methylation

573    extractor.

### Analysis of deamination efficiency and sequence preference from the NGS data

575    Deamination events were inferred from C to T changes at each covered cytosine position in the

576    genomes based on Bismark methylation extractor output. The 20 bp flanking sequences (10 bp

577    upstream and 10 bp downstream) were extracted from all the covered cytosines that had at least 10X

578    read coverage from the individual genomes and divided the cytosines sites into different groups based

579    on their deamination rates (e.g., >=90%, >=50%). Flanking sequences of each cytosine group were used

580    to make sequence logo using WebLogo 3[22] to infer deamination sequence preference.

### Determination of methylated CpG sites in human genome by Binomial correction

582    For all samples except the 10 pg SEM-seq libraries, to avoid false identification of methylated positions

583    due to the incomplete conversion of unmodified cytosines, we used Binomial statistics with Benjamini-

584    Hochberg correction to determine the methylated CpG sites that have significantly higher methylation

585    levels compared with the background. The P-value of the one-sided binomial test ($H_0: P \leq P_0$, $H_a: P >$

586    $P_0$) can be calculated as $P = 1 - F(k; n, P_0)$, where $F$ is the cumulative distribution function of the

587    binomial distribution, $k$ is the number of unmethylated cytosine, $n$ is the coverage (number of

588    methylated and unmethylated cytosine), $P$ is the sample methylation level for each region, and $P_0$ is the

589    background methylation level which is the non-conversion rate estimated from unmethylated lambda

590    genome. The P-value was further adjusted using Benjamini-Hochberg method. Only the CpG positions

591    with coverage above a certain cutoff and FDR (false discovery rate) < 0.05 were defined as methylated

592 CpG sites. For a single replicate and replicates combined data, we require a minimum of 3X and 5X
593 coverage, respectively.
594

595 **Genome-wide coverage analysis**
596 We combined the filtered bam files from the two technical replicates to investigate the sequencing
597 coverage across the genome. The GC content bias statistics were calculated using Picard (version
598 2.26.11) CollectGcBiasMetrics (http://broadinstitute.github.io/picard) and was plotted using a custom
599 python script.

600 **Determination of coverage around CpG islands and transcription start sites (TSS)**
601 We combined the filtered bam files from the two technical replicates to examine the sequencing
602 coverage of different methods around the CpG islands and transcription start sites. The reference bed
603 files of human CpG islands and TSS were downloaded from UCSC Table Browser[46].  and The coverage
604 was computed using bamCoverage from deepTools (version 3.3.1)[47] with the following parameters: --
605 binSize 50 --normalizeUsing RPGC --effectiveGenomeSize 2913022398 –ignoreDuplicates. Then the
606 coverage was calculated for per genome regions using computeMatrix from deepTools with -a 1000 -b
607 1000; and scale-regions and reference-point option were used for CpG islands and TSS, respectively.

608 **Determination of methylation status of CpG sites at epigenomic regions**
609 We studied the methylation status of defined CpG sites at key epigenetic regions. We used GM12878
610 ChIP-seq data sets (processed hg38 bed narrowPeak files) from ENCODE portal[23]: ENCFF023LTU
611 (H3K27ac), ENCFF188SZS (H3K4me3), and ENCFF796WRU (CTCF). The methylation of CpG sites around
612 these epigenetic regions was calculated using computeMatrix reference-point from deepTools with the
613 following parameters: -a 1000 -b 1000 --skipZeros --referencePoint center.

# 614 Supplementary information
615 Supplementary Table 1: Supplementary_Table_1.xlsx

616 Supplementary Figures 1-16, and Supplementary tables S2-S6: Supplementary_Figures_and_tables.pdf

# 617 Data Availability
618 NGS substrate specificity libraries and SEM-seq libraries are available at NCBI GEO accession: GSE233932

# 619 Acknowledgements

628

## 629 Competing interests

630 The authors are employees of New England Biolabs, a manufacturer and vendor of molecular biology
631 reagents.

## 632 References

633 1. Cervantes-Gracia, K., Gramalla-Schmitz, A., Weischedel, J. & Chahwan, R. APOBECs orchestrate

634     genomic and epigenomic editing across health and disease. *Trends Genet.* **37**, 1028–1043 (2021).

635 2. Iyer, L. M., Zhang, D., Rogozin, I. B. & Aravind, L. Evolution of the deaminase fold and multiple origins

636     of eukaryotic editing and mutagenic nucleic acid deaminases from bacterial toxin systems. *Nucleic*

637     *Acids Res.* **39**, 9473–9497 (2011).

638 3. Zhang, D., de Souza, R. F., Anantharaman, V., Iyer, L. M. & Aravind, L. Polymorphic toxin systems:

639     Comprehensive characterization of trafficking modes, processing, mechanisms of action, immunity

640     and ecology using comparative genomics. *Biol. Direct* **7**, 18 (2012).

641 4. Ruhe, Z. C., Low, D. A. & Hayes, C. S. Polymorphic Toxins and Their Immunity Proteins: Diversity,

642     Evolution, and Mechanisms of Delivery. *Annu. Rev. Microbiol.* **74**, 497–520 (2020).

643 5. Mok, B. Y. *et al.* A bacterial cytidine deaminase toxin enables CRISPR-free mitochondrial base editing.

644     *Nature* **583**, 631–637 (2020).

645 6. Mok, B. Y. *et al.* CRISPR-free base editors with enhanced activity and expanded targeting scope in

646     mitochondrial and nuclear DNA. *Nat. Biotechnol.* (2022) doi:10.1038/s41587-022-01256-8.

647 7. de Moraes, M. H. *et al.* An interbacterial DNA deaminase toxin directly mutagenizes surviving target

648     populations. *eLife* **10**, e62967 (2021).

649 8. Meyer, K. D. DART-seq: an antibody-free method for global m6A detection. *Nat. Methods* **16**, 1275–

650     1280 (2019).

651 9. Vaisvila, R. *et al.* Enzymatic methyl sequencing detects DNA methylation at single-base resolution

652     from picograms of DNA. *Genome Res.* **31**, 1280–1289 (2021).

653    10.    Sun, Z. *et al.* Nondestructive enzymatic deamination enables single-molecule long-read amplicon

654          sequencing for the determination of 5-methylcytosine and 5-hydroxymethylcytosine at single-base

655          resolution. *Genome Res.* **31**, 291–300 (2021).

656    11.    Schutsky, E. K. *et al.* Nondestructive, base-resolution sequencing of 5-hydroxymethylcytosine

657          using a DNA deaminase. *Nat. Biotechnol.* **36**, 1083–1090 (2018).

658    12.    Xiao, Y.-L. *et al.* Transcriptome-wide profiling and quantification of N6-methyladenosine by

659          enzyme-assisted adenosine deamination. *Nat. Biotechnol.* 1–11 (2023) doi:10.1038/s41587-022-

660          01587-6.

661    13.    Frommer, M. *et al.* A genomic sequencing protocol that yields a positive display of 5-

662          methylcytosine residues in individual DNA strands. *Proc. Natl. Acad. Sci. U. S. A.* **89**, 1827–1831

663          (1992).

664    14.    Eddy, S. R. Accelerated Profile HMM Searches. *PLOS Comput. Biol.* **7**, e1002195 (2011).

665    15.    Mistry, J. *et al.* Pfam: The protein families database in 2021. *Nucleic Acids Res.* **49**, D412–D419

666          (2021).

667    16.    The UniProt Consortium. UniProt: the universal protein knowledgebase in 2021. *Nucleic Acids*

668          *Res.* **49**, D480–D489 (2021).

669    17.    Mitchell, A. L. *et al.* MGnify: the microbiome analysis resource in 2020. *Nucleic Acids Res.* **48**,

670          D570–D578 (2020).

671    18.    Ito, F., Fu, Y., Kao, S.-C. A., Yang, H. & Chen, X. S. Family-Wide Comparative Analysis of Cytidine

672          and Methylcytidine Deamination by Eleven Human APOBEC Proteins. *J. Mol. Biol.* **429**, 1787–1799

673          (2017).

674    19.    Barka, A. *et al.* The Base-Editing Enzyme APOBEC3A Catalyzes Cytosine Deamination in RNA with

675          Low Proficiency and High Selectivity. *ACS Chem. Biol.* **17**, 629–636 (2022).

676  20.    Mitchell, A. L. *et al.* EBI Metagenomics in 2017: enriching the analysis of microbial communities,

677       from sequence reads to assemblies. *Nucleic Acids Res.* **46**, D726–D735 (2018).

678  21.    Füllgrabe, J. *et al.* Simultaneous sequencing of genetic and epigenetic bases in DNA. *Nat.*

679       *Biotechnol.* 1–8 (2023) doi:10.1038/s41587-022-01652-0.

680  22.    Crooks, G. E., Hon, G., Chandonia, J.-M. & Brenner, S. E. WebLogo: A Sequence Logo Generator.

681  23.    Davis, C. A. *et al.* The Encyclopedia of DNA elements (ENCODE): data portal update. *Nucleic*

682       *Acids Res.* **46**, D794–D801 (2018).

683  24.    Sun, Z. *et al.* High-Resolution Enzymatic Mapping of Genomic 5-Hydroxymethylcytosine in

684       Mouse Embryonic Stem Cells. *Cell Rep.* **3**, 567–576 (2013).

685  25.    Lianidou, E. Detection and relevance of epigenetic markers on ctDNA: recent advances and

686       future outlook. *Mol. Oncol.* **15**, 1683–1700 (2021).

687  26.    Yin, G. *et al.* BCORL1 is an independent prognostic marker and contributes to cell migration and

688       invasion in human hepatocellular carcinoma. *BMC Cancer* **16**, 103 (2016).

689  27.    Gillooly, J. F., Hein, A. & Damiani, R. Nuclear DNA Content Varies with Cell Size across Human

690       Cell Types. *Cold Spring Harb. Perspect. Biol.* **7**, a019091 (2015).

691  28.    Mi, L. *et al.* DddA homolog search and engineering expand sequence compatibility of

692       mitochondrial base editing. *Nat. Commun.* **14**, 874 (2023).

693  29.    Guo, J. *et al.* A Novel Double-Stranded DNA Deaminase-Based and Transcriptional Activator-

694       Assisted Nuclear and Mitochondrial Cytosine Base Editors with Expanded Target Compatibility and

695       Enhanced Activity. SSRN Scholarly Paper at https://doi.org/10.2139/ssrn.4227259 (2022).

696  30.    Huang, J. *et al.* Discovery of new deaminase functions by structure-based protein clustering.

697       2023.05.21.541555 Preprint at https://doi.org/10.1101/2023.05.21.541555 (2023).

698  31.    Paez-Espino, D. *et al.* IMG/VR: a database of cultured and uncultured DNA Viruses and

699       retroviruses. *Nucleic Acids Res.* **45**, gkw1030 (2017).

700    32.     Chen, I.-M. A. *et al.* The IMG/M data management and analysis system v.6.0: new tools and

701          advanced capabilities. *Nucleic Acids Res.* **49**, D751–D763 (2021).

702    33.     Ghaly, T. M., Geoghegan, J. L., Alroy, J. & Gillings, M. R. High diversity and rapid spatial turnover

703          of integron gene cassettes in soil. *Environ. Microbiol.* **21**, 1567–1574 (2019).

704    34.     Singleton, C. M. *et al.* Connecting structure to function with the recovery of over 1000 high-

705          quality metagenome-assembled genomes from activated sludge using long-read sequencing. *Nat.*

706          *Commun.* **12**, 2009 (2021).

707    35.     Da, B. *et al.* GenBank. *Nucleic Acids Res.* **41**, (2013).

708    36.     Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* 1–11 (2021)

709          doi:10.1038/s41586-021-03819-2.

710    37.     Mirdita, M. *et al.* ColabFold: making protein folding accessible to all. *Nat. Methods* **19**, 679–682

711          (2022).

712    38.     Edgar, R. C. High-accuracy alignment ensembles enable unbiased assessments of sequence

713          homology and phylogeny. 2021.06.20.449169 Preprint at

714          https://doi.org/10.1101/2021.06.20.449169 (2022).

715    39.     Kozlov, A. M., Darriba, D., Flouri, T., Morel, B. & Stamatakis, A. RAxML-NG: a fast, scalable and

716          user-friendly tool for maximum likelihood phylogenetic inference. *Bioinformatics* **35**, 4453–4455

717          (2019).

718    40.     Huerta-Cepas, J., Serra, F. & Bork, P. ETE 3: Reconstruction, Analysis, and Visualization of

719          Phylogenomic Data. *Mol. Biol. Evol.* **33**, 1635–1638 (2016).

720    41.     Söding, J. Protein homology detection by HMM–HMM comparison. *Bioinformatics* **21**, 951–960

721          (2005).

722    42.     Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *Mach. Learn. PYTHON*.

723    43.    Bair, C. L. & Black, L. W. A type IV modification dependent restriction nuclease that targets

724    glucosylated hydroxymethyl cytosine modified DNAs. *J. Mol. Biol.* **366**, 768–778 (2007).

725    44.    Krueger, F. & Andrews, S. R. Bismark: a flexible aligner and methylation caller for Bisulfite-Seq

726    applications. *Bioinforma. Oxf. Engl.* **27**, 1571–1572 (2011).

727    45.    Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–

728    359 (2012).

729    46.    Karolchik, D. *et al.* The UCSC Table Browser data retrieval tool. *Nucleic Acids Res.* **32**, D493-496

730    (2004).

731    47.    Ramírez, F. *et al.* deepTools2: a next generation web server for deep-sequencing data analysis.

732    *Nucleic Acids Res.* **44**, W160–W165 (2016).

733