# Bacteria-to-human protein networks reveal origins of endogenous DNA damage

Jun Xia[1-5,19], Li-Ya Chiu[6,19], Ralf B. Nehring[1-4], María Angélica Bravo Núñez†[1-4], Qian Mei[1-4,7], Mercedes Perez[6], Yin Zhai[2,4], Devon M. Fitzgerald[1-4], John P. Pribis[1-5], Yumeng Wang[8-10], Chenyue W. Hu[11], Reid T. Powell[12], Sandra A. LaBonte[13], Ali Jalali[4,14], Meztli L. Matadamas Guzmán‡[1-4], Alfred M. Lentzsch[6], Adam T. Szafran[15], Mohan C. Joshi[1,3,4,§], Megan Richters[1-4], Janet L. Gibson[1-4], Ryan L. Frisch¶[1-4], P.J. Hastings[1,4], David Bates[1,3,4], Christine Queitsch[16], Susan G. Hilsenbeck[4], Cristian Coarfa[4,15], James C. Hu[8], Deborah A. Siegele[8], Kenneth L. Scott[1,4,5], Han Liang[8-10], Michael A. Mancini[4,15], Christophe Herman[1,3,5,17], Kyle M. Miller[4,6,17] & Susan M. Rosenberg[1-5,7,18]

[1]Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, Texas 77030, USA.

[2]Department of Biochemistry and Molecular Biology, Baylor College of Medicine, Houston, Texas 77030, USA.

[3]Department of Molecular Virology and Microbiology, Baylor College of Medicine, Houston, Texas 77030, USA.

[4]Dan L Duncan Comprehensive Cancer Center, Baylor College of Medicine, Houston, Texas 77030, USA.

[5]Graduate Program in Integrative Molecular and Biomedical Sciences, Baylor College of Medicine, Houston, Texas 77030, USA.

[6]Department of Molecular Biosciences, Institute for Cellular and Molecular Biology, University of Texas at Austin, Austin, Texas 78712 USA.

[7]Systems, Synthetic and Physical Biology Program, Rice University, Houston, Texas 77030, USA.

[8]Graduate Program in Structural and Computational Biology and Molecular Biophysics, Baylor College of Medicine, Houston, Texas 77030, USA.

[9]Department of Bioinformatics and Computational Biology, The University of Texas MD Anderson Cancer Center, Houston, Texas 77030, USA.

[10]Department of Systems Biology, The University of Texas MD Anderson Cancer Center, Houston, Texas 77030, USA.

[11]Department of Bioengineering, Rice University, Houston, Texas 77030, USA.

[12]Institute of Biosciences and Technology, Texas A&M University, Houston, Texas 77030, USA.

[13]Department of Biochemistry and Biophysics, Texas A&M University and Texas AgriLife Research, College Station, TX 77843, USA.

[14]Department of Neurosurgery, Baylor College of Medicine, Houston, 77030, Texas USA.

[15]Department of Molecular and Cellular Biology, Baylor College of Medicine, Houston, Texas 77030, USA.

[16]Department of Genome Sciences, University of Washington, Seattle, Washington 98195.

[17]Senior author

[18]Lead author

[19]These authors contributed equally to this work.

41    †Present address: Graduate School of the Stowers Institute for Medical Research,1000 East 50th
42    Street, Kansas City, MO 64110, USA.

43    ‡Present address:  Doctorate in Biomedical Science, Universidad Nacional Autónoma de México,
44    México.

45    § Present address: Multidisciplinary Centre for Advance Research and Studies (MCARS), Jamia
46    Millia Islamia, New Delhi 110025, India.

47    ¶Present address:  DuPont Industrial Biosciences, 200 Powder Mill Road. Wilmington, DE 19803.

48

49    **Correspondence:    herman@bcm.edu    (CH),    kyle.miller@austin.utexas.edu    (KMM),**
50    **smr@bcm.edu (SMR)**
51

52    **SUMMARY**

53

54    DNA damage provokes mutations and cancer, and results from external carcinogens or
55    endogenous cellular processes. Yet, the intrinsic instigators of DNA damage are poorly
56    understood. Here we identify proteins that promote endogenous DNA damage when
57    overproduced: the DNA-damaging proteins (DDPs). We discover a large network of DDPs
58    in *Escherichia coli* and deconvolute them into six DNA-damage-causing function clusters,
59    demonstrating DDP mechanisms in three: reactive-oxygen increase by transmembrane
60    transporters, chromosome loss by replisome binding, and replication stalling by transcription
61    factors. Their 284 human homologs are over-represented among known cancer drivers, and their
62    expression in tumors predicts heavy mutagenesis and poor prognosis. Half of tested human
63    homologs, when overproduced in human cells, promote DNA damage and mutation, with DNA-
64    damaging mechanisms like those in *E. coli*. Together, our work reveals DDP networks that
65    provoke endogenous DNA damage and may indicate functions of many human known and newly
66    implicated cancer-promoting proteins.

67

## INTRODUCTION

DNA damage often underlies "spontaneous" mutations (Hastings et al., 1976; Tubbs and Nussenzweig, 2017), which drive cancer, genetic diseases, pathogen drug resistance and evasion of immune responses, and evolution generally. DNA damage can be caused by exogenous agents such as radiation or tobacco smoke, and indeed the vast majority of known carcinogens are DNA-damaging agents, and because of that, mutagens (Chatterjee and Walker, 2017). However, most DNA damage is generated endogenously within cells (Jackson and Loeb, 2001; Tubbs and Nussenzweig, 2017), by intrinsic cellular processes that involve macromolecule components including proteins. Presumably, proteins that promote endogenous DNA damage are required by cells, but cause DNA damage as side effects of their necessary functions and/or when dysregulated. The identities and functions of the proteins that promote endogenous DNA damage in cells in any organism are poorly understood or unknown (Figure 1A). We sought to identify these proteins systematically, and to understand how they might promote endogenous DNA damage, here.

One way to identify the proteins that *promote* endogenous DNA damage is by overproduction, which is a natural event that occurs frequently by copy-number alteration and other routes, and is a major source of cancer-driving functions (Zack et al., 2013). Given the conservation of DNA biology across the tree of life (Aravind et al., 1999; Makarova and Koonin, 2013), identification of proteins that promote spontaneous endogenous DNA damage carried out in any organism could potentially inform strategies for prevention, diagnosis, and treatment of disease, including therapeutic inhibition of cancer development, aging, and evolution of pathogens (Fitzgerald et al., 2017).

Some proteins that *prevent* or *reduce* endogenous DNA damage levels in cells have been identified by loss-of-function mutations/knock-downs that increase DNA damage (Alvaro et al., 2007; Lovejoy et al., 2009; Paulsen et al., 2009). DNA-repair proteins are among this category because they *reduce* levels of endogenous DNA damage. Similarly, proteins that *prevent/reduce* endogenous DNA damage, along with other kinds of proteins, will be among genes the loss-of-function of which promotes mutagenesis or genome instability (Putnam et al., 2016; Yuen et al., 2007). By contrast, no unbiased screen has been reported for proteins that actively *promote* endogenous DNA damage in cells. Though a limited nucleus-specific screen identified some (Lovejoy et al., 2009), the range of functions and numbers of proteins, processes, and mechanisms that cause endogenous DNA damage are unknown.

Here we report the comprehensive discovery in *Escherichia coli* of a large, diverse network of proteins that promote endogenous DNA damage when overproduced: the "DNA-damaging" proteins (DDPs). We show that their upregulation promotes mutagenesis, and use massive function-based assays to identify kinds, causes, and consequences of intrinsic DNA damage provoked by overproduction of the 208 *E. coli* DDPs. We found that they group into six discreet function clusters, and determine molecular mechanisms of DNA-damage generation from three of these. We identify their human homologs, also a large network, and find that they are overrepresented among known cancer instigators, and that their expression in human cancers predicts poor outcomes and high mutation loads. We show that overproduction of human homologs in human cells also promotes endogenous DNA damage and mutation. We determine the mechanisms of DNA-damage instigation for two cancer-associated human DDPs, both of which mimic bacterial mechanisms and suggest unexpected roles in cancer. The identities and functions of the proteins in bacterial and human DDP networks provide an important general model for illuminating mechanisms of genesis of endogenous DNA damage, and may inform

114   cancer-promoting function discovery of many known and newly discovered cancer-driving
115   proteins.
116

117   **RESULTS**
118   **Large Diverse Protein Network Promotes DNA Damage**
119   We screened an inducible overproduction library of all >4000 *E. coli* proteins in two steps to
120   identify clones with increased endogenous DNA-damage levels (STAR Methods, Figure 1B). For
121   both the primary and secondary screens, we measured fluorescence of cells that carry a
122   fluorescence-reporter gene driven by an SOS DNA-damage-response-activated promoter at a non-
123   genic chromosomal site (Nehring et al., 2016) (Figure 1B). The promoter fusion reports DNA-
124   damage-response induction, and not spurious promoter firing (Pennington and Rosenberg, 2007)
125   (Figures S1A and S1B). First, in the primary screen, we used a fluorescence plate-reader, which is
126   high-throughput but low-resolution, to identify potential clones with increased DNA damage and
127   so increased fluorescence (Figures 1B and 1C). This identified potential positive candidates (414
128   proteins). Second, we eliminated false positives from the primary plate-reader screen using a
129   sensitive flow-cytometry secondary screen in the same strains (Figures 1B and 1D). Flow
130   cytometry, though low-throughput, is highly sensitive and reports DNA damage at the single-cell
131   level (Pennington and Rosenberg, 2007) (STAR Methods, Figure 1D). The stringent, high-
132   sensitivity flow-cytometry secondary screen validated 208 of the proteins identified in the primary
133   screen as genuine DDPs that cause increased DNA damage when overproduced (Figures 1D-G;
134   Table S1).
135       Further, we tested a representative sample of 66 of the DDP-overproduction clones (Table
136   S1) for whether their increased fluorescence requires the SOS-response-activator protein RecA
137   and a functioning (inducible) LexA SOS repressor, as expected if the fluorescence results from
138   activation of the SOS DNA-damage response (Pennington and Rosenberg, 2007). All 66 showed
139   RecA- and LexA-dependent high fluorescence, demonstrating the presence of DNA damage and
140   a genuine SOS response (Figures 1H and S1C; Table S1). We also ruled out the possibility that
141   DDP overproduction might increase mCherry protein fluorescence itself, independently of DNA
142   damage, by showing that a separate representative sample of 40 of the 208 DDPs did not cause
143   increased fluorescence from the same mCherry reporter gene under the control of a non-SOS
144   promoter (Figure S1D; Table S1).
145       DNA-damage-promotion by overproduction of the 208 DDPs is observed additionally in
146   three ways. First, 95% of the 208 validated *E. coli* DDPs displayed an independent DNA-damage-
147   related phenotype in at least one of nine assays that are either less sensitive or more DNA-damage-
148   type specific than the SOS-flow cytometric assay in the secondary screen. For example, we tested
149   a representative sample of 67 of the DDP-overproduction clones for the presence of damaged,
150   single-stranded DNA using a less sensitive assay for microscopically visible foci of a fluorescent
151   DNA-damage-sensor protein RecA*GFP (Renzette et al., 2005), which is less sensitive because it
152   uses a partially functional RecA protein (Renzette et al., 2005). The data nevertheless show a
153   significant association of RecA foci with SOS-positive (DDP) clones in that 32 of the 67 tested
154   showed increased foci ($r = 0.7$, $p = 1.3 \times 10^{-10}$, Pearson's correlation) (Figure 1I; Table S1 for clone-
155   by-clone results). Also, later in this paper, we explored the kinds and causes of endogenous DNA
156   damage promoted by the 208 *E. coli* DDPs, using seven assays for specific kinds of DNA
157   damage—all more DNA-damage-type-specific than the SOS-response assay used here, and we
158   additionally assayed DDPs for mutagenesis, described below (summarized Table S1; Figure 1J).
159   All but 12 of the 208 clones were positive in either the RecA*GFP-focus assay, and/or one of the

160 7 assays described below, or mutagenesis (summarized Table S1). This equates to 95% of the 208
161 proteins showing DNA-damage-related phenotypes in an independent assay. Finally, all of the
162 twelve not validated in a non-SOS-based assay (Table S1) were shown to increase fluorescence
163 SOS-response dependently, showing only RecA- LexA-dependent fluorescence increase (Figures
164 1H and S1C; Table S1) and not general fluorescence increase (Figure S1D; Table S1). Collectively,
165 these data demonstrate independently that all 208 DDPs promote DNA damage when
166 overproduced.
167 　　　　The 208 proteins span many different classes that function in diverse cellular and metabolic
168 processes (Figure 1F; Table S1), and only 8% encode known DNA-repair proteins (blue font,
169 Figure 1F; Table S1). Although DNA-repair proteins *reduce* DNA damage when expressed
170 normally, their overproduction, here, increased DNA damage, which might occur by perturbing
171 undamaged DNA, by titrating repair partner proteins away from DNA damage and/or inhibiting
172 DNA repair. We call all of these proteins DNA-damaging proteins (DDPs).
173 　　　　The DDPs constitute a network both functionally and by protein-protein associations.
174 Functionally, they cause the same phenotype—increased DNA damage on overproduction—but
175 their reported protein functions are remarkably diverse (Figure 1F; Table S1). We used STRING—
176 a database that contains known and predicted associations between protein pairs—to examine
177 whether these diverse proteins had any other known or predicted connections indicative of a
178 network. STRING measures protein-protein associations or interactions of many kinds including
179 indirect associations (e.g., co-occurrence in different organisms' genomes, or in papers in the
180 literature) and direct protein-protein interactions, among others (STAR Methods). Using STRING
181 with an interaction score cut-off of ≥0.6 (medium-to-high confidence, STAR Methods), we found
182 that the 208 DDPs form a significant network via protein-protein interaction data (Figure 1G,
183 specific interactions, Figure S2A) with more interactions than random sets of 208 *E. coli* proteins
184 ($p = 2.0 \times 10^{-31}$, hypergeometric test, Supplemental Discussion 1). When known DNA-repair
185 proteins are removed, the STRING network is still significant compared with random sets of the
186 same number of *E. coli* proteins ($p = 9 \times 10^{-7}$, hypergeometric test), indicating that this association
187 network is not solely via DNA-repair-protein associations. When both DNA-repair and DNA-
188 replication proteins are removed—both known to interact directly with DNA—the STRING
189 network is no longer significant compared with random sets of the same number of *E. coli* proteins
190 ($p = 0.08$, hypergeometric test). These data suggest that, as might be expected from the highly
191 diverse protein functions (Figure 1F), these proteins work in many different cellular processes that
192 may share only their various effects on DNA damage, seen by significant association as a network
193 only when DNA-repair and replication proteins are included, which we identify as the hubs of the
194 DDP STRING interaction network (Figure 1G).
195 　　　　The DDP network is estimated to be larger than the 208 proteins identified (Supplemental
196 Discussion 2, Figure S1E). Although the premier overproduction (Mobile) library was used (Saka
197 et al., 2005), its composition of some native genes and some genes encoding five additional amino
198 acids is indicated by our data to have prevented detection of some additional DDPs (Supplemental
199 Discussion 2).
200
201 **Endogenous DNA Damage Increases Mutations**
202 We tested the hypothesis that triggering endogenous DNA damage would increase mutation rates
203 (Figure 1A) using DDPs:  a useful test because they were discovered based on DNA damage,
204 rather than genome instability/mutation rate.  Mutation rates of a sample of 32 representative *E.*
205 *coli* DDP clones were assayed by a modified forward-mutation fluctuation-test assay (Figure 1J,

6

206    STAR Methods). We chose 10 DDP clones from the low-damage group (<5-fold increase in
207    endogenous DNA-damage levels) and 22 from the high-damage group with a >5-fold increase in
208    endogenous DNA-damage levels (Supplemental Discussion 3; Table S1). The data in Figure 1J
209    show that increased endogenous DNA-damage levels are associated significantly with elevated
210    mutation rates. We confirmed that the mutagenesis assay reported genuine loss-of-function
211    mutations in the *c*I mutation-reporter gene (Figure S1F), rather than gene-regulatory or epigenetic
212    changes, by sequencing mutations in a sample of independent isolated mutants of 10 representative
213    DDP clones (Figure S1F). The sequence analyses of selected DDP clones revealed additionally
214    that high DNA damage led to various kinds of mutations including base substitutions, indels,
215    transposition events, and gross chromosomal rearrangements (GCRs, including large deletions,
216    Figure S1F, right). These mutations mimic the increased small mutations and GCRs seen in various
217    cancers (Stratton, 2011). Although the SOS response induces mutations by upregulation of low-
218    fidelity DNA polymerases (Pols) V and IV (Kobayashi et al., 2002; Maor-Shoshani et al., 2000;
219    Wagner and Nohmi, 2000), some of the mutations (Figure S1F, blue font) differ from common Pol
220    V and Pol IV errors (Figure S1F, red font). The data suggest that the type of DNA damage, and
221    not merely induction of the SOS response, may influence the kinds and rates of mutations made
222    (Figures 1J and S1F, Table S1). The data show that overproduction of many functionally diverse
223    *E. coli* proteins (Figure 1F; Table S1) causes increased DNA-damage loads (Figures 1C-E; Table
224    S1), and genome instability with mutations of essentially all kinds (Figures 1J and S1F; Table S1).
225
226    **Human Homologs of Bacterial DDPs a Network Associated with Cancers**
227    As an unbiased quantitative way to find human DDPs, we identified 284 human homologs of the
228    *E. coli* DDPs via BLASTp and deltaBLAST searches (STAR Methods, Figure 2A; Table S2).
229    "Homologs" are defined here as proteins with amino-acid similarity that may result from possible
230    evolutionary relatedness (STAR Methods). The 284 human homologs are used here as candidate
231    human DDPs (hDDPs), and are homologs of 68 of the *E. coli* DDPs (shown, Table S2). The
232    remaining *E. coli* DDPs are mostly analogs of human proteins, which function similarly but are
233    not homologous (Serres et al., 2001), and many are of unknown function. The hDDP candidate
234    proteins also constitute a protein-protein interaction network (Figure 2B, specific interactions
235    Figure S2B), with significantly more interactions than sets of 284 random human proteins ($p = 1.2$
236    $\times 10^{-327}$, hypergeometric test), or random human homologs of *E. coli* proteins, which also differ
237    from random human proteins, but less so ($p = 1.8 \times 10^{-49}$, hypergeometric test, discussed
238    Supplemental Discussion 1). Only 5.6% of the human homologs are known DNA-repair proteins
239    (blue font, Figure 2A), again indicating a different class of candidate genome-integrity-affecting
240    proteins. Like the *E. coli* network (Figure 1G), DNA-repair and -replication proteins are central
241    hubs (Figure 2B).
242    We tested whether the human homologs of *E. coli* DDPs are both relevant to human
243    cancers, and behave as a network, by examining their associations with various kinds of data from
244    human cancers. We observe strong associations of the 284-protein network in cancer data of
245    several kinds. First, the human homologs (Figure 2A; Table S2) of *E. coli* DDPs are significantly
246    overrepresented among known (Forbes et al., 2015) and predicted (D'Antonio and Ciccarelli,
247    2013) cancer drivers in a curated consensus in the Sanger Institute's Catalogue of Somatic
248    Mutations In Cancer (COSMIC) (Forbes et al., 2015) and the database of D'Antonio and Ciccarelli
249    (D'Antonio and Ciccarelli, 2013) ($p = 0.0002$, Fisher's exact test; Figure 2C; Table S3), which
250    contain gain- and loss-of-function drivers. Human homologs of random *E. coli* proteins are not
251    overrepresented ($p = 0.48$, Fisher's exact test). Thus, the cancer association is specific to DDP

7

252 homologs, not conserved proteins generally. The human homologs remain overrepresented among
253 known and predicted drivers when homologs of *E. coli* DNA-repair proteins and other known
254 human DNA-repair proteins (Figure 2A) are excluded ($p$ = 0.05, Figure 2C). No other
255 comprehensive overexpression screen for DNA damage has been reported; however, we analyzed
256 cancer association in published data from a limited, selected-candidate overexpression screen in
257 human cells (Lovejoy et al., 2009), and found that these also show cancer association
258 (Supplemental Discussion 4). These overlap with our 284 hDDP candidate network by only one
259 protein—FIGNL—indicating that many new candidates were revealed by the *E. coli* screen.
260     Additionally, we found that candidate human DDP genes show increased copy numbers in
261 26 cancer types in the cohort of patients in The Cancer Genome Atlas (Gao et al., 2013) (TCGA)
262 (Figures 2D and S3A-C; Table S4). About 40% of the 284 human homolog genes have increased
263 copy numbers in cancers (GISTIC threshold copy-number gain ≥ 1), either cancer-specifically or
264 across cancers, compared with fewer than 20% of non-DDP genes amplified in those cancers
265 (Figure 2D). The fractions of patients with increased copy numbers of each of the genes encoding
266 the 284 human homologs of *E. coli* DDPs are shown for all 284 in Table S4 (examples, Figure
267 S3A-C). The human homologs are enriched as copy-number increases in cancers compared with
268 non-DDP human genes (Figure 2D, $p$ = 0.04, one-way Fisher's exact test), suggesting that their
269 overexpression is associated with cancers.
270     We next examined the outcomes for cancer survival relative to mRNA levels of the 284
271 candidate human DDPs using cancer-patient and RNA data in TCGA (Gao et al., 2013). We found
272 that in at least four cancer types, increased levels of the 284 RNAs, relative to the total RNAs, is
273 associated with decreased overall survival (Figure 2E). This association results not just from the
274 known (Forbes et al., 2015) and predicted (D'Antonio and Ciccarelli, 2013) cancer driver genes in
275 the network, but is seen also in the network genes not known previously to drive cancers (Figure
276 S3D-F). These data indicate an association of candidate hDDP gene overexpression with poor
277 survival in these cancers, and further highlight the network properties/predictive power of the 284
278 DDP-homolog protein/gene set. Moreover, increase of the 284 human DDP-homolog RNAs,
279 relative to all RNAs, is also associated with total genomic mutation burden in cancers (relative to
280 the patient normal tissue) in at least 12 cancer types in TCGA (Gao et al., 2013) (Figure 2F). Even
281 stronger association is seen for the subset of the 284 homologs that are known (Forbes et al., 2015)
282 or predicted (D'Antonio and Ciccarelli, 2013) cancer-driving genes (Figure 2F). These data support
283 the possibility that overexpression in the candidate hDDP network is associated with mutagenesis
284 in human cancers.
285

286 **Human Homologs Promote DNA Damage and Mutation**
287 We validated a sample of human candidate DDPs as genuine DNA-damage instigators in human
288 cells (Figure 3). We tested our hypothesis that overproduction of these proteins, which can result
289 from gene amplification, can promote DNA damage relevant to cancers by testing a sample in
290 which about half the homologs were known to be amplified in cancers in TCGA (Gao et al., 2013)
291 and the other half were not (Supplemental Discussion 5). We were also limited by availability in
292 human cDNA-clone collections (Yang et al., 2011) (Supplemental Discussion 5). Because many
293 genes in those collections are not full length (Table S5), we cloned several de novo (STAR
294 Methods) to create 70 full-length sequence-verified overexpression GFP-fusion genes encoding
295 human homologs of *E. coli* DDPs, 3 human homologs of *E. coli* damage-down proteins, as possible
296 negative controls (Table S5), and 20 control random non-DDPs (Table S5, ~half of which are
297 random human homologs of *E. coli* proteins). Using transient transfection of these human

298  overexpression clones, we performed three flow-cytometric assays, which we developed (Figure
299  3A), to screen for increased DNA damage at the single-cell level in human cells that produce the
300  candidate hDDPs, shown by GFP (Supplemental Discussion 6 for the infeasibility of stable clones).
301  We screened for—(i) increased γH2AX levels, a marker for DNA double-strand breaks (DSBs)
302  (Kinner et al., 2008); (ii) increased γH2AX in a sensitized screen in cells treated with a
303  nonhomologous-break-repair (DNA-PK) inhibitor; and (iii) increased levels of the DNA-damage
304  marker protein phospho-p53, which indicates activation of the DNA-damage response (Sakaguchi
305  et al., 1998). The data show that the human homologs are enriched for genuine DDPs that increase
306  DNA damage upon overproduction (Figure 3B). Among the 73 human homologs, we found that
307  45% (33 of the 73) showed increased DNA damage (Figure 3B). This is highly significant ($p <$
308  $0.0001$ one-way Fisher's exact test) compared with the 20 random human proteins (Figures S3G-
309  I; Table S6). Of the 33 validated hDDPs, only one (FIGNL1) was known previously to increase
310  DNA damage upon overproduction (Lovejoy et al., 2009). Thus, we identified and validated 33
311  genuine human DDPs.
312      We note, however, that the human-cell DNA-damage assays used here favor detection of
313  DSBs, not all DNA-damage types comprehensively. Thus, many more of the human homologs
314  may be DNA-damage instigating for other kinds of DNA damage than is estimated here.
315      As in *E. coli*, we found that overproduction of validated hDDPs increased mutation rates
316  in human cells. Using a forward-mutation fluctuation-test assay for hypoxanthine-guanine
317  phosphoribosyl transferase (HPRT) deficiency (STAR Methods), we found increased mutation
318  rates for 4 out of 4 overproduced high-DNA-damage hDDP clones compared with cell-only,
319  vector-only, and GFP-tubulin-overproducing negative controls (Figure 3C and Supplemental
320  Discussion 8). Thus, increased mutation rates result from overproduction of validated hDDPs in
321  human cells. These results support the hypothesis that hDDPs may drive cancers based on their
322  ability to increase genome instability—a known cancer-driving phenotype (Hanahan and
323  Weinberg, 2011).
324      As shown in Figure 3E, 32 of the 33 validated human DDPs are *E. coli* DDP homologs
325  from the following categories: (i) 16% that are both known (Forbes et al., 2015) or predicted
326  (D'Antonio and Ciccarelli, 2013) cancer drivers and amplified in TCGA cancers; (ii) 53% that are
327  amplified in cancers and not known or predicted drivers; (iii) 6% known/predicted cancer drivers
328  that are not known to be amplified in cancers; and (iv) 25% that are neither gene-amplified in
329  cancers nor previously known or predicted drivers. None of these classes was predicted previously
330  to be DNA-damage promoting, and some might not have been hypothesized to potentially promote
331  cancer via overproduction (e.g., classes iii and iv). In Supplemental Discussion 7, we use the rates
332  of validation in each class tested to estimate that there are likely to be many additional hDDPs
333  among the 284-protein candidate hDDP network, that would test positive in our assays.
334      Bioinformatically, ~75% of the validated hDDP genes show cancer-associated copy-
335  number increases in the TCGA patient-cohort data (Gao et al., 2013) (GISTIC threshold copy-
336  number gain ≥ 1, Figure 3D). The fraction of cancer-specific or across-cancer copy-number-
337  increased genes among the validated hDDP genes is higher than the candidates we tested that were
338  not validated (Figure 3D, $p = 0.02$, one-way Fisher's exact test). The data provide support for our
339  hypothesis that validated hDDP genes may drive cancer via DNA damage when overexpressed,
340  and imply that our cell-based DNA-damage assays relate to human cancer biology.
341
342  **Functional Systems Biology**
343  Having identified DDPs in *E. coli* and human, we used the tractable *E. coli* model for further

9

344     function discovery. We sought to create a multi-parameter, minable data set of phenotypes to bin
345     the 208 *E. coli* DDPs into function clusters that reflect the kinds, causes, and consequences of
346     DNA damage provoked by their overproduction. Our phenotypes are based on seven quantitative
347     functional assays, many at the single-cell level (Figure 4). Two of these employ synthetic proteins
348     that trap, fluorescently label, and allow quantification, as well as genomic mapping, of specific
349     DNA-damage-intermediate structures in single living cells (Shee et al., 2013; Xia et al., 2016). We
350     use the data to predict, then demonstrate, mechanisms by which dysregulation of diverse conserved
351     proteins increase endogenous DNA damage.
352

353     **Proteins that Instigate DNA Double-Strand Breaks**
354     We used the engineered double-strand-break (DSB)-end-specific binding protein GamGFP (Shee
355     et al., 2013) to quantify DSBs in single living cells. GamGFP "traps" DSB ends and prevents their
356     repair in *E. coli* and mammalian cells (Shee et al., 2013). GamGFP produced from a chromosomal
357     regulatable gene cassette labels one-ended and two-ended DSBs as fluorescent foci in *E. coli* at an
358     estimated 70% efficiency (Shee et al., 2013) (i.e., 30% of DSBs present are not seen as foci). We
359     quantified GamGFP foci in each of the 208 DDP-overproducing clones by automated microscopy
360     (STAR Methods). We found that 87 of the 208 DDP-overproducing clones displayed significantly
361     more GamGFP foci than vector-only controls (Figures 4A-B and S4A, Table S1, Supplemental
362     Discussion 9) and than 25 random SOS-negative (non-DDP) clones, none of which had increased
363     GamGFP foci (Figure S4A). Our finding of 121 DDP-overproducing clones without increased
364     GamGFP foci suggests that single-stranded (ss)DNA, the inducing signal for the SOS DNA-
365     damage response, frequently accumulates at sites other than DSBs. The 41% of *E. coli* DDP clones
366     with elevated GamGFP DSB foci are not enriched in any gene-function category (Table S1),
367     implying that DSBs—a common result of many DNA-damaging mechanisms (Merrikh et al.,
368     2012)—result from various cellular processes.
369

370     **Stalled Reversed Replication Forks**
371     Stalling of DNA replication leads to DNA damage (Jackson and Bartek, 2009) and can create four-
372     way DNA junctions when stalled replication forks "reverse" such that the new DNA strands
373     basepair with each other (Figure 4C, reversed fork, RF) (Seigneur et al., 1998). Reversed forks
374     (RFs) block resumption of DNA replication, lead to replication-fork breakage (Seigneur et al.,
375     1998; Yeeles et al., 2013), and so are both a kind of DNA damage, and reflect a cause of DNA
376     damage—stalled replication. We quantified stalled, RFs as fluorescent foci of the engineered 4-
377     way-junction-specific DNA-binding protein RuvCDefGFP (RDG) in cells lacking homology-
378     directed-repair protein RecA ($\Delta recA$), in which essentially all RDG foci represent RFs (Xia et al.,
379     2016). RDG labels 4-way junctions with about 50% efficiency in live *E. coli* (Xia et al., 2016).
380     We found that 106 of the 208 DDP-overproducing clones showed increased RDG (RF) foci
381     relative to the vector-only control (Figures 4D and S4B, Table S1) and to 30 control, SOS-negative
382     overproducing clones (Figure S4B; Table S1). Among the 106 clones with increased RDG foci,
383     49 also show increased DSBs (Table S1), detected as GamGFP foci, showing significant
384     correlation ($p = 0.03$, r=0.15 Spearman's correlation). Overall, at least 51% of DDPs promote
385     replication stalling on overproduction, indicating the importance of DNA replication to DNA-
386     damage generation by many of the overproduced proteins, but also suggesting that mechanisms in
387     addition to replication may underlie DNA damage induced by dysregulation of many other DDPs.
388

389     **Reactive Oxygen via Transporters and Metabolism**

390 We quantified intracellular levels of reactive oxygen species (ROS) in single cells using the
391 peroxide-indicator dye, dihydrorhodamine (DHR) (Gutierrez et al., 2013) and flow cytometry
392 (Figures 4E, F and S4C, D). We found that 56 of the overproduced DDPs caused increased ROS
393 levels (Figures 4F and S4C, D). We show that the high ROS levels contribute to DNA damage in
394 at least 16 of the 56 ROS-elevated DDP-producing clones, in which the DNA damage was
395 significantly reduced by the ROS-quenching agent thiourea (Keren et al., 2013), compared with
396 the vector-only control (Figure 4G). Thus, high endogenous ROS levels underlie DNA damage in
397 a subset of the DDP-producing clones. These (Figure 4G) comprise five membrane-spanning
398 transporters (investigated below), an excess compared with the prevalence of transporters among
399 *E. coli* proteins ($p$=0.002, hypergeometric test). The other eleven proteins relate to metabolism
400 processes (Table S1), implying that perturbation of metabolic pathways can cause DNA damage
401 by increasing ROS.
402
403 **DNA Loss**
404 Loss of DNA in cells can result from various problems including chromosome-segregation failure
405 (Joshi et al., 2013), for example, from incomplete DNA replication or incomplete homology-
406 directed repair between chromosomes, either of which can leave two chromosomes attached at cell
407 division (Hendricks et al., 2000). We identified 67 DDP clones with increased frequencies of
408 DNA-depleted ("anucleate") cells (Figures 4H, I and S4E, F; Table S1) using flow cytometry of
409 single cells with DNA and cell membranes stained separately. Overproduced DNA-repair and
410 replication proteins are enriched among these clones ($p$ = 0.04, one-way Fisher's exact test, Table
411 S1), implying that excessive DNA-repair and replication proteins promote DNA damage that leads
412 to DNA erosion or chromosome-segregation failure. Their overproduction might alter
413 stoichiometry of and hinder DNA-repair or replication complexes, and/or perturb DNA directly.
414
415 **Reduced DNA-repair Capacity via Specific Proteins**
416 Reduction of DNA-repair capacity could increase DNA-damage levels, and could result either
417 from saturation of repair pathways with excessive DNA damage, or from overproduction of
418 proteins that interfere directly with DNA-repair. Either would mimic repair-deficiency (and
419 mutator phenotype), without a DNA-repair-gene mutation. We assayed DNA-repair capacity
420 indirectly, as sensitivity to DNA-damaging agents that produce damage repaired by specific DNA-
421 repair mechanisms: DSB-, ssDNA-break-, and ROS-instigator phleomycin (Steighner and Povirk,
422 1990); base-oxidizing agent hydrogen peroxide ($H_2O_2$); and DNA cross-linking agent mitomycin-
423 C (MMC). These cause DNA damage repaired by homology-directed repair (HR, phleomycin),
424 base-excision repair ("BER", $H_2O_2$), and, for MMC, both nucleotide-excision repair and HR
425 (Friedberg et al., 2005). We found that 106, 75, and 10 of the 208 DDP clones were sensitive to
426 phleomycin (Figures 4J and S5A-B), $H_2O_2$ (Figures 4L and S5C), and MMC (Figures 4K and S5D,
427 E), respectively, shown by reduced cell densities in cultures (normalized for any effects of protein
428 overproduction on overall growth rates, STAR Methods). Collectively, 140 DDP-overproducing
429 strains were sensitive to at least one DNA-damaging agent (Figure S5G; Table S1), and 45 to
430 multiple drugs (Figure S5G; Table S1). Non-DDP overproduction clones were not enriched for
431 sensitivity to DNA-damaging agents and differ from the DDP-network group for sensitivity to
432 phleomycin and $H_2O_2$, but not for MMC (Figure S5). The data suggest that overproduction of
433 various DDPs provoke different kinds of DNA damage that overwhelm or inhibit distinct DNA-
434 repair mechanisms.
435         We excluded the possibility that most DNA-damage sensitivities resulted from DDP-

11

436   induced heritable mutation, by showing no sensitivity *after* removal from DNA damage (Figure
437   S5H, Supplemental Discussion 10). We also excluded transcriptional downregulation of DNA-
438   repair genes, by RNA-seq of DNA repair genes (Figure S5F, Supplemental Discussion 10). The
439   data suggest that specific DNA-repair pathways are either inhibited directly or saturated by DNA
440   damage caused by overproduction of specific DDPs and imply that dysregulating diverse proteins,
441   such as by gene copy-number alteration, can create DNA-repair deficiency without mutation of
442   DNA-repair genes.
443
444   **Clustering *E. coli* Function Data Implicates Mechanisms**
445   We grouped the quantitative data from the functional assays using stability-based clustering [(Hu
446   et al., 2015), Progeny clustering] (Figures 4M,N and S5G). The quantitative data on RDG (RF)
447   foci analyzed with three other quantitative parameters measured in single-cells—ROS levels, DNA
448   loss, and DSBs (all discussed above)—revealed that high RF loads are enriched in a specific cluster
449   (Figure 4M). The RF-dense cluster is significantly enriched for DNA-binding transcription factors
450   (examined below), with 29%, compared with 12% among the network as a whole ($p = 0.002$, one-
451   way Fisher's exact test, Figure 4M; Table S1). The data indicate that distinct protein functions
452   preferentially stall replication.
453       Grouping all quantitative data sets revealed six discreet function clusters (Figure 4N; Table
454   S1), which may indicate at least six different potential mechanisms and cellular consequences of
455   DNA-damage promotion by DDPs (reduced DNA-damage classes discussed Supplemental
456   Discussion 11). We compared the function clusters with protein-protein-association data in the
457   DDP network (Figures 4N and S2A, Table S1). Whereas the entire DDP network shows
458   significantly more protein-protein associations than sets of 208 random *E. coli* proteins,
459   superimposition of the function clusters onto the protein-protein interaction network indicates that
460   cluster 2 shows even more protein-protein interactions than the DDP network as a whole ($p =$
461   $0.0007$, one-was Fisher's exact test). These data support associations of function clusters with
462   particular biological mechanisms, three examined below.
463
464   **Transcription Factor Binding to DNA Promotes Replication Stalls**
465   Clusters 5 and 6 of Figures 4N (listed Table S1) show increased replication stalling/reversed forks
466   (RFs, per Figure 4M) and are most enriched for DNA-binding transcription factors: transcriptional
467   activators and repressors. We hypothesized that persistent binding of a protein to DNA might
468   create a replication "roadblock", stall forks, and cause RFs. RFs can be regarded as DNA damage
469   and also cause additional DNA damage when cleaved by endonucleases (Seigneur et al., 1998). In
470   support of this hypothesis, we found, first, that mutational ablation of the DNA-binding-domains
471   (DBDs) of three of the transcription factors—CsgD, HcaR and MhpR—abolished both their
472   abilities to promote SOS-inducing DNA damage (Figures 5A-C), and RDG (RF) foci (Figures 5D-
473   E and S6A). The data indicate that these transcription factors must bind DNA to provoke DNA
474   damage and RFs upon overproduction. Second, we created an mCherry (red) fusion of the CsgD
475   transcription factor (Ogasawara et al., 2011), and see that it forms foci DBD-dependently (Figure
476   5F and S6B), suggesting that foci reflect the DNA-bound transcription factor. We found that most
477   of the CsgD-mCherry foci, and also HcaR-mCherry foci, co-localized with RDG (RF) foci
478   (Figures 5F-H), suggesting that RFs form near the DNA-bound transcription factors. Foci of DNA-
479   bound proteins are distinguishable at ~50kb apart on DNA, e.g., (Shee et al., 2013); thus, these
480   data indicate that RDG/RF foci accumulate in the vicinity of the sites of transcription-factor
481   binding to DNA (Figures 5G and 5H). High resolution mapping of RDG (RFs) by ChIP-seq in the

12

482   genome of CsgD-overproducing cells showed RDG (RFs) enriched near the transcription factor's
483   target DNA-binding sites CsgD-DBD-dependently (Figure 5I), supporting the hypothesis that the
484   bound transcription factor stalls replication causing fork reversal nearby. CsgD has 10
485   experimentally well-characterized binding sites (Ogasawara et al., 2011), and we found that the
486   CsgD-DBD-dependent RDG (RF) ChIP-seq peaks are very significantly enriched in 10kb regions
487   surrounding known CsgD DNA-binding sites (representative peaks, Figure 5I; Supplemental
488   Discussion 12, rest of known sites; Figure S7). CsgD-DBD-dependent RDG (RF) ChIP-seq peaks
489   occurred both upstream and downstream of the binding sites in the replication paths (Figure S7,
490   discussed Supplemental Discussion 12, and Figure 5J). Our data support a model (Figure 5J) in
491   which overproduced DNA-bound transcription factors can create roadblocks to replication, which
492   leads to increased fork stalling and reversal near where the transcription factors bind, causing DNA
493   damage.

494

495   ### *E. coli* and Human Transporter Overproduction Elevates ROS
496   Membrane-spanning transporters are the largest category of human homologs of the *E. coli* DDPs,
497   and several are both overrepresented among known cancer drivers and also provoke DNA damage
498   on overproduction (Figures 2A, Tables S2 and S3). We found that *E. coli* membrane transporters
499   are overrepresented at 26% in the high-ROS cluster in Figure 4M compared with 11% over the
500   whole network ($p = 0.004$, one-way Fisher's exact test, Figure 6A-C; Table S1). Further, sixteen
501   DDP clones with high ROS caused DNA damage ROS-dependently in that the damage was
502   reduced by ROS quenching (Figure 4I). These include five transporters, the increased ROS and
503   ROS-dependent DNA damage of which are shown in Figures 6B-D. Three of the five are $H^+$
504   symporters, a significant enrichment compared with the frequency of $H^+$ symporters encoded in
505   the genome (Keseler et al., 2017) ($p = 2.7x10^{-5}$, hypergeometric test), one transports polypeptides,
506   and the remaining one $Mg^{2+}$.
507         Proton ($H^+$) symporters import molecules concurrently with $H^+$. We found that
508   overproduction of each of the three $H^+$ symporters conferred reduced intracellular pH (increased
509   $H^+$) (Figures 6E-G), implying that overproduction increased their symporter activities. However,
510   their induction of ROS was not well correlated with their reduction of pH (Figure 6H), suggesting
511   that other cargos that they import may cause the increased ROS and DNA damage, or that simply
512   compromising membrane integrity and cellular boundaries may provoke ROS and DNA damage.
513   A specific model for XanQ, the strongest ROS-promoter among them, is illustrated in (Figure 6I),
514   discussed Supplemental Discussion 13. Overall, the data reveal that DNA-damage induction can
515   result from increased transporter activity, leading to high levels of DNA-damaging ROS (Figures
516   6A-C). The data suggest that disturbing cellular boundaries can cause DNA damage via ROS.
517         CorA, an inner membrane $Mg^{2+}$ transporter, which transports $Co^{2+}$ and $Ni^{2+}$ less efficiently
518   (Kehres and Maguire, 2002), elevates ROS (Figures 6C-D) and DNA damage (Figure 6B) when
519   overproduced. Both might occur by increasing the usually minimal import of $Co^{2+}$ and $Ni^{2+}$ (Figure
520   6I). $Ni^{2+}$ is toxic and induces DNA damage via oxidative stress (Cameron et al., 2011) because
521   $Ni^{2+}$ binds sulfhydryl groups commonly found in anti-oxidative enzymes (Schmidt et al., 2009).
522   Increased $Mg^{2+}$ import is unlikely to underlie the ROS and DNA damage because $Mg^{2+}$ is the most
523   abundant metal ion in cells and excessive $Mg^{2+}$ does not seem to affect the activities of the many
524   $Mg^{2+}$-utilizing enzymes (Hartwig, 2001).
525         In human cells, multiple DNA-damaging mechanisms are implicated. First, a survey of the
526   subcellular localization of all 33 overproduced, validated hDDPs showed that 16 were cytoplasmic;
527   10 were nuclear, and 7 were found throughout the cell (Figures 6J), suggesting that direct contact

13

528    with DNA is not needed for many overproduced hDDPs' instigation of DNA damage.

529        We screened a sample of thirteen validated hDDP-producing clones for those with DNA
530    damage that could be suppressed by ROS quencher N-acetyl cysteine (NAC), to identify those that
531    instigate ROS-dependent DNA damage. We found that overproduction of a membrane-spanning
532    transporter promotes DNA damage via ROS. KCNAB1/2 promote high DNA damage (Figure 6K),
533    are cytoplasmic when overproduced (Figure 6J), and are subunits of intracellular voltage-gated $K^+$
534    channels that function in redox transformations of xenobiotics (Hlavac et al., 2014). Increased
535    *KCNAB2* mRNA is found in breast cancer (Hlavac et al., 2014), but how KCNAB1/2
536    overproduction might promote cancer is unknown. We found that DNA-damage promotion by
537    KCNAB1/2 relies at least partly on ROS, in that ROS-quenching NAC treatment reduced DNA-
538    damage induction by KCNAB1 or KCNAB2 overproduction (Figure 6K). Cells overproducing
539    several other validated hDDPs showed no reduction in DNA damage with NAC treatment (Figure
540    6K), indicating that these DDPs promote DNA damage by other or additional mechanisms.

541

542    ***E. coli* Pol IV and Human DNMT1 Promote DNA Damage via Replisome-Clamp Interaction**
543    DNA polymerase (Pol) IV, encoded by *dinB*, is among the highest DNA-damage generators in the
544    *E. coli* DDP network (Figures 7A-C; Table S1), generating both DSBs seen by increased GamGFP
545    foci (Table S1) and ssDNA gaps, inferred as follows.  SOS DNA-damage-response induction by
546    overproduced Pol IV was partially RecB- (DSB) and partially RecF- (ssDNA-gap) dependent
547    (Figures 7D-E); and RecB and RecF are required for SOS induction by DSBs and ssDNA gaps,
548    respectively (McPartland et al., 1980), implicating DSBs and ssDNA gaps as Pol IV-induced
549    damage.  Overproduced Pol IV also promotes chromosome loss (Figure 4E). We show, that Pol
550    IV-induced DNA damage occurs via its interaction with the replisome sliding clamp as follows.

551        Pol IV is a poorly processive, low-fidelity DNA polymerase that traverses specific
552    damaged bases that more efficient DNA polymerases cannot copy (Wagner et al., 2000). Pol IV
553    competes with more processive DNA polymerases (Frisch et al., 2010; Hastings et al., 2010), by
554    competition for binding the replisome sliding clamp protein, beta (Dohrmann et al., 2016; Heltzel
555    et al., 2012). We hypothesized that the documented ability of Pol IV to slow replication-fork
556    progression (Heltzel et al., 2012; Uchida et al., 2008) might underlie its generation of DNA damage
557    when overproduced. In support of this hypothesis, we found, first, that overproducing Pol IV
558    simultaneously with DNA Pol II, its competitor for the replisome (Frisch et al., 2010), caused a
559    roughly 50% reduction in the Pol IV-dependent DNA damage (Figures 7B-C). Second, cells
560    producing a mutant replisome clamp-loader protein with reduced Pol IV loading and increased
561    loading of the major replicative DNA polymerase, Pol III (Dohrmann et al., 2016), also showed
562    reduced DNA damage from Pol IV overproduction (Figures 7D-E, *dnaX* tau-only), but no
563    reduction of SOS on its own (Figure 7F). Third, a C-terminal Pol IV deletion that abolishes Pol IV
564    interaction with the replisome beta sliding clamp (Uchida et al., 2008) also abolished DNA damage
565    upon Pol IV overproduction (ΔBBD, Figures 7B-C).  We conclude that Pol IV interaction with the
566    replisome clamp is required for its generation of DNA damage on overproduction.

567        Perhaps surprisingly, Pol IV catalytic activity (DNA synthesis) was not required for all of
568    its DNA-damage induction; the catalytically-inactive Pol IV R49F mutant (Wagner et al., 1999)
569    reduced only half the DNA damage caused by overproduction (Pol IV cat⁻, Figures 7B and 7C),
570    without reducing Pol IV protein levels (Uchida et al., 2008; Wagner et al., 1999) (Figure S7D).
571    Thus, mere binding of Pol IV to the replisome, not only its poorly processive catalytic activity,
572    appears sufficient to cause DNA damage. The synthesis-dependent component of DNA damage
573    might have resulted from the ability of Pol IV to incorporate oxidized guanine (8-oxo-dG) into

14

574  DNA, per (Foti et al., 2012), which leads to two different strand-breaking BER processes that
575  begin with base removal by MutM and MutY DNA glycosylases (Foti et al., 2012). However, loss
576  of neither glycosylase diminished DNA damage caused by Pol IV overproduction (Figures 7E and
577  7G). The data imply that BER following 8-oxo-dG incorporation is not how excess Pol IV
578  promotes most DNA damage. Overall, Pol IV promotes DNA damage dependently on replisome-
579  clamp interaction, and only partly dependently on catalysis (model, Figure S7E).
580      We found that human DNMT1 overproduction also induces DNA damage in human cells
581  based on binding the replisome clamp, and independently of its catalytic activity:  a non-canonical
582  potential cancer-driving role. DNMT1 is the major human DNA methyltransferase that methylates
583  DNA upon replication (Jin and Robertson, 2013). Hypomorphic mutations in *DNMT1* promote
584  microsatellite instability (Jin and Robertson, 2013). Increased DNMT1 is common to several
585  cancer types, and causes hypermethylation, proposed to downregulate tumor-suppressor genes
586  (Biniszkiewicz et al., 2002), which would constitute a cell-biological, rather than DNA-based,
587  tumor-promoting role. Surprisingly, we found that, in addition, DNMT1 promotes DNA damage
588  independently of its DNA-methylation activity, in that overproduction of DNMT1 catalytically-
589  dead mutant proteins (Figure 7H) increased DNA damage similarly to wild-type DNMT1 (Figures
590  7I and 7J). Overproduction of two other DNA methyltransferases did not increase DNA-damage
591  levels (Figure 7I). DNMT1 truncations (Figure 7H) revealed that DNMT1 promotion of DNA
592  damage required its PCNA-binding domain (PBD), which binds the replisome sliding clamp:
593  PCNA (Figures 7I, S3L and S3M). Rad18-mediated monoubiquitination of PCNA, a DNA damage
594  response (Mortusewicz et al., 2005), also resulted from DNMT1 overproduction, also methylase-
595  independently and requiring the DNMT1 PBD (Figures 7J and S3M). These data suggest that mere
596  binding of overproduced DNMT1 to the replisome clamp promotes DNA damage, and a resulting
597  DNA-damage response, independently of methylation. The finding that both *E. coli* DNA
598  polymerase IV and human DNMT1 promote DNA damage dependently on replisome-clamp
599  binding and independently of their catalytic activities (Figure 7A-G) indicates generality of this
600  mode of generation of DNA damage. The data suggest that promotion of DNA damage by
601  DNMT1-PCNA complexes (Figures 7H-J), and resulting mutagenesis (Figure 3C), may promote
602  cancers other than or in addition to via the known cell-biological/regulatory function of DNMT1
603  in DNA methylation, and that many clamp-binding proteins may act similarly when their genes
604  are overexpressed/amplified.
605
606
607  **DISCUSSION**
608
609  The identities and functions of proteins in the *E. coli* and human DNA-damaging-protein networks
610  reveal that multiple diverse proteins, cellular processes, and molecular mechanisms underlie
611  genesis of endogenous DNA damage. Although obtained in an overexpression screen, these are
612  likely to represent natural causes of spontaneous endogenous DNA damage, first, because gene
613  overexpression in cells in populations is remarkably common, on the order of tens of percents of
614  bacterial cells for any given gene (Elowitz et al., 2002), with copy-number gains of any
615  chromosomal region occurring in $10^{-3}$ of cells (Reams et al., 2010), and expected to be comparable
616  in human cells (Hastings et al., 2009). Second, the association of the 284 human DDP homologs
617  with four aspects of human cancer data indicates natural biological relevance (reviewed below).
618
619  **Endogenous DNA Damage in Cancer**

15

620 The 284 human homologs of *E. coli* DDP genes are overrepresented among known (Forbes et al.,
621 2015) and predicted (D'Antonio and Ciccarelli, 2013) cancer drivers (Figure 2C), overrepresented
622 in cancers as amplified (Figures 2D and S3A-C, Table S3), and their increased expression in
623 human cancers associated with poor outcomes (Figures 2E and S3D-F) and heavy mutation loads
624 (Figure 2F). These associations support their overexpression being both biologically relevant, and
625 relevant to cancer biology specifically.
626 The DDPs appear likely to represent a new broad function class of cancer-promoting
627 proteins, and the earliest in cancer. Cancer-gene functions have been grouped into multiple specific
628 categories (Hanahan and Weinberg, 2011) that fit into two broad classes of function (Kinzler and
629 Vogelstein, 1997): the cancer-cell-biology-altering "gatekeepers", mutations in which make cell
630 biology more cancer-like, and the genomic "caretakers"—the DNA-repair genes, mutations in
631 which elevate mutation rate and so drive cancer by increasing gatekeeper mutations (Figure 7K).
632 The DDPs are expected to act *before* and upstream of DNA-repair functions promoting the
633 endogenous DNA damage that necessitates repair (Figure 7K), and so instigating some of the
634 earliest events in cancer development.
635 The large number of DDPs span diverse protein functions, the cancer-driving functions of
636 many of which may be obscure or mis-assigned. Some of the mechanisms of DDP action may
637 necessitate reëvaluation of their cancer driving roles, and also of the drugs designed to inhibit them.
638 For example, we found that human DNA methyltransferase DNMT1 causes DNA damage
639 independently of its methylation activity, via its interaction with the replisome sliding clamp
640 (Figure 7H-J). Current cancer drugs against DNMT1 target the methylase activity (Jones et al.,
641 2016), and not replisome binding. It is unclear which activities of DNMT1 promote cancer, and so
642 which should be drugged. Our finding may inform the development of and use of DNMT1-
643 targeting strategies, taking into account its multifunctionality, broadly across cancer types. These
644 results underscore the importance of determining all functions of a protein that are cancer
645 promoting.
646
647 **The *E. coli*-to-Cancer Gene-function Atlas of Bacterial and Human DDP Phenotypes**
648 Our data from seven quantitative assays for kinds, causes, and consequences of endogenous DNA
649 damage promoted by *E. coli* DDPs constitute a rich resource and framework for within- and cross-
650 species discovery of conserved DNA-damage-generating mechanisms. We used them to identify
651 six main function or phenotype clusters (Figure 4N), and implicate, then test and demonstrate,
652 three *E. coli* DDP mechanisms (Figures 5-7), two also identified in human cells (Figures 6 and 7).
653 We created a minable web-based resource for searching the complete *E. coli* function data, and the
654 functional data from the validated human DDPs: the *E. coli*-to-Cancer Gene-function Atlas
655 (ECGA) (https://microbialphenotypes.org/wiki/index.php/Special:ECGA). The ECGA data can be
656 searched via the bacterial proteins' or their human homologs' names or by function key words.
657 ECGA can be used for querying/generation of hypotheses for *E. coli* and potential conserved
658 human-protein functions, and as it develops in future, for other organisms.
659
660 **Mechanisms that Cause Endogenous DNA Damage**
661 In *E. coli* DNA-binding transcription factors caused replication-fork stalling and reversal (Figure
662 5) by apparent blocking of replication forks by the bound transcription factors (Figure 5J). Though
663 replication-transcription conflicts have been engineered by reversing chromosome segments
664 (Tehranchi et al., 2010), engineering multiple transcription-factor binding sites into long arrays
665 (Magnan et al., 2015), or knock out of RNA-removal proteins (Wahba et al., 2016), the kinds of

16

666    DNA damage generated were not identified, nor was it known that any of these mechanisms could
667    occur in natural genomes with only an endogenous protein upregulated, as often occurs in cells.
668    Our results indicate that fork reversal is common and protein-function specific. Nearly 10% of
669    human genes encode transcription factors (Levine and Tjian, 2003), including many cancer-driving
670    overproduction (onco-)proteins, some known to promote DNA damage when overproduced, e.g.,
671    c-Myc, and E2F1 (Pickering and Kowalik, 2006; Vafa et al., 2002). Thus, based on our observation
672    of transcription-factor binding and DNA-damage-production in *E. coli*, many onco-protein
673    transcription factors might promote cancer similarly to *E. coli* transcription factors—by causing
674    genome-destabilizing DNA damage, potentially RFs. Moreover, in *E. coli* and human cells, we
675    discovered that increased transmembrane transporter activities of several different kinds elevate
676    ROS levels causing DNA damage (Figures 4M and 6A-F). This previously unknown mechanism,
677    also apparent with the human KCNAB1/2 transporter, might explain the KCNAB2 association
678    with cancers (Hlavac et al., 2014)—a hypothesis that remains to be tested. Further, we found that
679    both *E. coli* DNA polymerase IV (Figure 7A-G) and human DNMT1 (Figures 7H-J, and S3L and
680    M) provoke DNA damage via binding their respective replisome sliding clamps when
681    overproduced, independently of their catalytic activities. Disruption of the replisome leading to
682    replication-fork collapse (or other means Figure S7E), is thought to be an important source of DNA
683    damage (Kuzminov, 1995) likely to apply to dysregulation of many kinds of proteins.
684

685    **DNA Damage as Potential Cancer Biomarker**
686    The existence of many diverse means of increasing endogenous DNA damage, and the predicted
687    large sizes and diversity of both bacterial and human DDP networks indicate that dysregulation of
688    any of many proteins is likely to be mutagenic via DNA damage (Figures 1J, 2F and S1F). Because
689    many different proteins, processes, and mechanisms instigate DNA damage, DNA damage itself
690    might be a robust predictor of cancer and genetic-disease susceptibility. The ability to detect high
691    DNA damage could potentially make DNA-damage screening attractive for early identification of
692    at-risk individuals, at a time before genome-sequencing would identify disease-associated
693    mutations. Additionally, the success of cancer immune therapy "checkpoint inhibitors" is limited
694    to high-mutagenesis cancers, apparently because mutagenesis creates diverse tumor antigens that
695    can be attacked by the stimulated immune system (Germano et al., 2017). Thus, immune therapy
696    is currently aimed primarily at some DNA-repair-defective cancers (Germano et al., 2017). Our
697    data suggest that DNA damage, or upregulation of DDP network genes, may predict additional
698    susceptibilities in various cancers.
699

700    **REFERENCES**

701

702    Alvaro, D., Lisby, M., and Rothstein, R. (2007). Genome-wide analysis of Rad52 foci reveals
703    diverse mechanisms impacting recombination. PLoS Genet *3*, e228.

704    Aravind, L., Walker, D.R., and Koonin, E.V. (1999). Conserved domains in DNA repair proteins
705    and evolution of repair systems. Nucleic Acids Res *27*, 1223-1242.

706    Asad, N.R., Asad, L.M.B.O., Almeida, C.E.B.d., Felzenszwalb, I., Cabral-Neto, J.B., and Leitão,
707    A.C. (2004). Several pathways of hydrogen peroxide action that damage the E. coli genome.
708    Genetics and Molecular Biology *27*, 291-303.

709     Benjamini, Y., and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and
710     powerful approach to multiple testing. Journal of the royal statistical society Series B
711     (Methodological), 289-300.

712     Berkopec, A. (2007). HyperQuick algorithm for discrete hypergeometric distribution. Journal of
713     Discrete Algorithms *5*, 341-347.

714     Biniszkiewicz, D., Gribnau, J., Ramsahoye, B., Gaudet, F., Eggan, K., Humpherys, D.,
715     Mastrangelo, M.A., Jun, Z., Walter, J., and Jaenisch, R. (2002). Dnmt1 overexpression causes
716     genomic hypermethylation, loss of imprinting, and embryonic lethality. Mol Cell Biol *22*, 2124-
717     2135.

718     Bolger, A.M., Lohse, M., and Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina
719     sequence data. Bioinformatics *30*, 2114-2120.

720     Bonocora, R.P., and Wade, J.T. (2015). ChIP-seq for genome-scale analysis of bacterial DNA-
721     binding proteins. Methods Mol Biol *1276*, 327-340.

722     Boratyn, G.M., Schaffer, A.A., Agarwala, R., Altschul, S.F., Lipman, D.J., and Madden, T.L.
723     (2012). Domain enhanced lookup time accelerated BLAST. Biol Direct *7*, 12.

724     Brombacher, E., Dorel, C., Zehnder, A.J., and Landini, P. (2003). The curli biosynthesis
725     regulator CsgD co-ordinates the expression of both positive and negative determinants for
726     biofilm formation in Escherichia coli. Microbiology *149*, 2847-2857.

727     Cameron, K.S., Buchner, V., and Tchounwou, P.B. (2011). Exploring the molecular mechanisms
728     of nickel-induced genotoxicity and carcinogenicity: a literature review. Rev Environ Health *26*,
729     81-92.

730     Carrasco, B., Cozar, M.C., Lurz, R., Alonso, J.C., and Ayora, S. (2004). Genetic recombination
731     in Bacillus subtilis 168: contribution of Holliday junction processing functions in chromosome
732     segregation. J Bacteriol *186*, 5557-5566.

733     Chatterjee, N., and Walker, G.C. (2017). Mechanisms of DNA damage, repair, and mutagenesis.
734     Environ Mol Mutagen *58*, 235-263.

735     D'Antonio, M., and Ciccarelli, F.D. (2013). Integrated analysis of recurrent properties of cancer
736     genes to identify novel drivers. Genome Biol *14*, R52.

737     Dohrmann, P.R., Correa, R., Frisch, R.L., Rosenberg, S.M., and McHenry, C.S. (2016). The
738     DNA polymerase III holoenzyme contains gamma and is not a trimeric polymerase. Nucleic
739     Acids Res *44*, 1285-1297.

740     Dudin, O., Geiselmann, J., Ogasawara, H., Ishihama, A., and Lacour, S. (2014). Repression of
741     flagellar genes in exponential phase by CsgD and CpxR, two crucial modulators of Escherichia
742     coli biofilm formation. J Bacteriol *196*, 707-715.

743    Elowitz, M.B., Levine, A.J., Siggia, E.D., and Swain, P.S. (2002). Stochastic gene expression in
744    a single cell. Science *297*, 1183-1186.

745    Fitzgerald, D.M., Hastings, P., and Rosenberg, S.M. (2017). Stress-induced mutagenesis:
746    implications in cancer and drug resistance. Annu Rev Cancer Biol *1*, 119-140.

747    Forbes, S.A., Beare, D., Gunasekaran, P., Leung, K., Bindal, N., Boutselakis, H., Ding, M.,
748    Bamford, S., Cole, C., Ward, S.*, et al.* (2015). COSMIC: exploring the world's knowledge of
749    somatic mutations in human cancer. Nucleic Acids Res *43*, D805-811.

750    Foster, P.L. (2006). Methods for determining spontaneous mutation rates. Methods Enzymol
751    *409*, 195-213.

752    Foti, J.J., Devadoss, B., Winkler, J.A., Collins, J.J., and Walker, G.C. (2012). Oxidation of the
753    guanine nucleotide pool underlies cell death by bactericidal antibiotics. Science *336*, 315-319.

754    Friedberg, E.C., Walker, G.C., Siede, W., and Wood, R.D. (2005). DNA repair and mutagenesis
755    (American Society for Microbiology Press).

756    Frisch, R.L., Su, Y., Thornton, P.C., Gibson, J.L., Rosenberg, S.M., and Hastings, P.J. (2010).
757    Separate DNA Pol II- and Pol IV-dependent pathways of stress-induced mutation during double-
758    strand-break repair in Escherichia coli are controlled by RpoS. J Bacteriol *192*, 4694-4700.

759    Galhardo, R.S., Almeida, C.E., Leitao, A.C., and Cabral-Neto, J.B. (2000). Repair of DNA
760    lesions induced by hydrogen peroxide in the presence of iron chelators in Escherichia coli:
761    participation of endonuclease IV and Fpg. J Bacteriol *182*, 1964-1968.

762    Gao, J., Aksoy, B.A., Dogrusoz, U., Dresdner, G., Gross, B., Sumer, S.O., Sun, Y., Jacobsen, A.,
763    Sinha, R., Larsson, E.*, et al.* (2013). Integrative analysis of complex cancer genomics and
764    clinical profiles using the cBioPortal. Sci Signal *6*, pl1.

765    Gao, R., Davis, A., McDonald, T.O., Sei, E., Shi, X., Wang, Y., Tsai, P.C., Casasent, A., Waters,
766    J., Zhang, H.*, et al.* (2016). Punctuated copy number evolution and clonal stasis in triple-negative
767    breast cancer. Nat Genet *48*, 1119-1130.

768    Germano, G., Lamba, S., Rospo, G., Barault, L., Magri, A., Maione, F., Russo, M., Crisafulli, G.,
769    Bartolini, A., Lerda, G.*, et al.* (2017). Inactivation of DNA repair triggers neoantigen generation
770    and impairs tumour growth. Nature *552*, 116-120.

771    Gutierrez, A., Laureti, L., Crussard, S., Abida, H., Rodriguez-Rojas, A., Blazquez, J., Baharoglu,
772    Z., Mazel, D., Darfeuille, F., Vogel, J.*, et al.* (2013). beta-Lactam antibiotics promote bacterial
773    mutagenesis via an RpoS-mediated reduction in replication fidelity. Nat Commun *4*, 1610.

774    Hanahan, D., and Weinberg, R.A. (2011). Hallmarks of cancer: the next generation. Cell *144*,
775    646-674.

776    Hanzelmann, S., Castelo, R., and Guinney, J. (2013). GSVA: gene set variation analysis for
777    microarray and RNA-seq data. BMC Bioinformatics *14*, 7.

778    Hartwig, A. (2001). Role of magnesium in genomic stability. Mutat Res *475*, 113-121.

779    Hastings, P.J., Hersh, M.N., Thornton, P.C., Fonville, N.C., Slack, A., Frisch, R.L., Ray, M.P.,
780    Harris, R.S., Leal, S.M., and Rosenberg, S.M. (2010). Competition of Escherichia coli DNA
781    polymerases I, II and III with DNA Pol IV in stressed cells. PLoS One *5*, e10862.

782    Hastings, P.J., Lupski, J.R., Rosenberg, S.M., and Ira, G. (2009). Mechanisms of change in gene
783    copy number. Nat Rev Genet *10*, 551-564.

784    Hastings, P.J., Quah, S.K., and von Borstel, R.C. (1976). Spontaneous mutation by mutagenic
785    repair of spontaneous lesions in DNA. Nature *264*, 719-722.

786    Heckman, K.L., and Pease, L.R. (2007). Gene splicing and mutagenesis by PCR-driven overlap
787    extension. Nat Protoc *2*, 924-932.

788    Heltzel, J.M., Maul, R.W., Wolff, D.W., and Sutton, M.D. (2012). Escherichia coli DNA
789    polymerase IV (Pol IV), but not Pol II, dynamically switches with a stalled Pol III* replicase. J
790    Bacteriol *194*, 3589-3600.

791    Hendricks, E.C., Szerlong, H., Hill, T., and Kuempel, P. (2000). Cell division, guillotining of
792    dimer chromosomes and SOS induction in resolution mutants (dif, xerC and xerD) of
793    Escherichia coli. Mol Microbiol *36*, 973-981.

794    Hlavac, V., Brynychova, V., Vaclavikova, R., Ehrlichova, M., Vrana, D., Pecha, V., Trnkova,
795    M., Kodet, R., Mrhalova, M., Kubackova, K*., et al.* (2014). The role of cytochromes p450 and
796    aldo-keto reductases in prognosis of breast carcinoma patients. Medicine (Baltimore) *93*, e255.

797    Hu, C.W., Kornblau, S.M., Slater, J.H., and Qutub, A.A. (2015). Progeny Clustering: A Method
798    to Identify Biological Phenotypes. Sci Rep *5*, 12894.

799    Hu, C.W., and Qutub, A.A. (2016). progenyClust: an R package for Progeny Clustering. R
800    JOURNAL *8*, 328-338.

801    Jackson, A.L., and Loeb, L.A. (2001). The contribution of endogenous sources of DNA damage
802    to the multiple mutations in cancer. Mutat Res *477*, 7-21.

803    Jackson, S.P., and Bartek, J. (2009). The DNA-damage response in human biology and disease.
804    Nature *461*, 1071-1078.

805    Jin, B., and Robertson, K.D. (2013). DNA methyltransferases, DNA damage repair, and cancer.
806    Adv Exp Med Biol *754*, 3-29.

807    Johnson, S.C. (1967). Hierarchical clustering schemes. Psychometrika *32*, 241-254.

808    Jones, P.A., Issa, J.P., and Baylin, S. (2016). Targeting the cancer epigenome for therapy. Nat
809    Rev Genet *17*, 630-641.

810   Joshi, M.C., Magnan, D., Montminy, T.P., Lies, M., Stepankiw, N., and Bates, D. (2013).
811   Regulation of sister chromosome cohesion by the replication fork tracking protein SeqA. PLoS
812   Genet *9*, e1003673.

813   Kehres, D.G., and Maguire, M.E. (2002). Structure, properties and regulation of magnesium
814   transport proteins. Biometals *15*, 261-270.

815   Kelley, E.E., Khoo, N.K., Hundley, N.J., Malik, U.Z., Freeman, B.A., and Tarpey, M.M. (2010).
816   Hydrogen peroxide is the major oxidant product of xanthine oxidase. Free Radic Biol Med *48*,
817   493-498.

818   Keren, I., Wu, Y., Inocencio, J., Mulcahy, L.R., and Lewis, K. (2013). Killing by bactericidal
819   antibiotics does not depend on reactive oxygen species. Science *339*, 1213-1216.

820   Keseler, I.M., Mackie, A., Santos-Zavaleta, A., Billington, R., Bonavides-Martinez, C., Caspi,
821   R., Fulcher, C., Gama-Castro, S., Kothari, A., Krummenacker, M.*, et al.* (2017). The EcoCyc
822   database: reflecting new knowledge about Escherichia coli K-12. Nucleic Acids Res *45*, D543-
823   D550.

824   Kim, S.R., Matsui, K., Yamada, M., Gruz, P., and Nohmi, T. (2001). Roles of chromosomal and
825   episomal dinB genes encoding DNA pol IV in targeted and untargeted mutagenesis in
826   Escherichia coli. Mol Genet Genomics *266*, 207-215.

827   Kinner, A., Wu, W., Staudt, C., and Iliakis, G. (2008). Gamma-H2AX in recognition and
828   signaling of DNA double-strand breaks in the context of chromatin. Nucleic Acids Res *36*, 5678-
829   5694.

830   Kinzler, K.W., and Vogelstein, B. (1997). Cancer-susceptibility genes. Gatekeepers and
831   caretakers. Nature *386*, 761, 763.

832   Kobayashi, S., Valentine, M.R., Pham, P., O'Donnell, M., and Goodman, M.F. (2002). Fidelity
833   of Escherichia coli DNA polymerase IV. Preferential generation of small deletion mutations by
834   dNTP-stabilized misalignment. J Biol Chem *277*, 34198-34207.

835   Kuzminov, A. (1995). Collapse and repair of replication forks in Escherichia coli. Mol Microbiol
836   *16*, 373-384.

837   Kuzminov, A. (2011). Homologous Recombination-Experimental Systems, Analysis, and
838   Significance. EcoSal Plus *4*.

839   Levine, M., and Tjian, R. (2003). Transcription regulation and animal diversity. Nature *424*, 147-
840   151.

841   Li, H. (2013). Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
842   arXiv preprint arXiv:13033997.

843   Loiselle, F.B., and Casey, J.R. (2010). Measurement of Intracellular pH. Methods Mol Biol *637*,
844   311-331.

845  Lovejoy, C.A., Xu, X., Bansbach, C.E., Glick, G.G., Zhao, R., Ye, F., Sirbu, B.M., Titus, L.C.,
846  Shyr, Y., and Cortez, D. (2009). Functional genomic screens identify CINP as a genome
847  maintenance protein. Proc Natl Acad Sci U S A *106*, 19304-19309.

848  Magnan, D., Joshi, M.C., Barker, A.K., Visser, B.J., and Bates, D. (2015). DNA Replication
849  Initiation Is Blocked by a Distant Chromosome-Membrane Attachment. Curr Biol *25*, 2143-
850  2149.

851  Makarova, K.S., and Koonin, E.V. (2013). Archaeology of eukaryotic DNA replication. Cold
852  Spring Harb Perspect Biol *5*, a012963.

853  Maor-Shoshani, A., Reuven, N.B., Tomer, G., and Livneh, Z. (2000). Highly mutagenic
854  replication by DNA polymerase V (UmuC) provides a mechanistic basis for SOS untargeted
855  mutagenesis. Proc Natl Acad Sci U S A *97*, 565-570.

856  McClure, R., Balasubramanian, D., Sun, Y., Bobrovskyy, M., Sumby, P., Genco, C.A.,
857  Vanderpool, C.K., and Tjaden, B. (2013). Computational analysis of bacterial RNA-Seq data.
858  Nucleic Acids Res *41*, e140.

859  McPartland, A., Green, L., and Echols, H. (1980). Control of recA gene RNA in E. coli:
860  regulatory and signal genes. Cell *20*, 731-737.

861  Merrikh, H., Zhang, Y., Grossman, A.D., and Wang, J.D. (2012). Replication-transcription
862  conflicts in bacteria. Nat Rev Microbiol *10*, 449-458.

863  Michaels, M.L., Pham, L., Cruz, C., and Miller, J.H. (1991). MutM, a protein that prevents G.C--
864  --T.A transversions, is formamidopyrimidine-DNA glycosylase. Nucleic Acids Res *19*, 3629-
865  3632.

866  Miller, J. (1993). A short course in bacterial genetics: a laboratory manual and handbook for
867  Escherichia coli and related bacteria. Trends in Biochemical Sciences-Library Compendium *18*,
868  193.

869  Mortusewicz, O., Schermelleh, L., Walter, J., Cardoso, M.C., and Leonhardt, H. (2005).
870  Recruitment of DNA methyltransferase I to DNA repair sites. Proc Natl Acad Sci U S A *102*,
871  8905-8909.

872  Nehring, R.B., Gu, F., Lin, H.Y., Gibson, J.L., Blythe, M.J., Wilson, R., Bravo Nunez, M.A.,
873  Hastings, P.J., Louis, E.J., Frisch, R.L.*, et al.* (2016). An ultra-dense library resource for rapid
874  deconvolution of mutations that cause phenotypes in Escherichia coli. Nucleic Acids Res *44*,
875  e41.

876  Ogasawara, H., Yamamoto, K., and Ishihama, A. (2011). Role of the biofilm master regulator
877  CsgD in cross-regulation between biofilm formation and flagellar synthesis. J Bacteriol *193*,
878  2587-2597.

879     Paulsen, R.D., Soni, D.V., Wollman, R., Hahn, A.T., Yee, M.C., Guan, A., Hesley, J.A., Miller,
880     S.C., Cromwell, E.F., Solow-Cordero, D.E., *et al.* (2009). A genome-wide siRNA screen reveals
881     diverse cellular processes and pathways that mediate genome stability. Mol Cell *35*, 228-239.

882     Pennington, J.M., and Rosenberg, S.M. (2007). Spontaneous DNA breakage in single living
883     Escherichia coli cells. Nat Genet *39*, 797-802.

884     Pickering, M.T., and Kowalik, T.F. (2006). Rb inactivation leads to E2F1-mediated DNA
885     double-strand break accumulation. Oncogene *25*, 746-755.

886     Putnam, C.D., Srivatsan, A., Nene, R.V., Martinez, S.L., Clotfelter, S.P., Bell, S.N., Somach,
887     S.B., de Souza, J.E., Fonseca, A.F., de Souza, S.J., *et al.* (2016). A genetic network that
888     suppresses genome rearrangements in Saccharomyces cerevisiae and contains defects in cancers.
889     Nat Commun *7*, 11256.

890     Rahman, M., Jackson, L.K., Johnson, W.E., Li, D.Y., Bild, A.H., and Piccolo, S.R. (2015).
891     Alternative preprocessing of RNA-Sequencing data in The Cancer Genome Atlas leads to
892     improved analysis results. Bioinformatics *31*, 3666-3672.

893     Ramirez, F., Dundar, F., Diehl, S., Gruning, B.A., and Manke, T. (2014). deepTools: a flexible
894     platform for exploring deep-sequencing data. Nucleic Acids Res *42*, W187-191.

895     Reams, A.B., Kofoid, E., Savageau, M., and Roth, J.R. (2010). Duplication frequency in a
896     population of Salmonella enterica rapidly approaches steady state with or without recombination.
897     Genetics *184*, 1077-1094.

898     Renzette, N., Gumlaw, N., Nordman, J.T., Krieger, M., Yeh, S.P., Long, E., Centore, R.,
899     Boonsombat, R., and Sandler, S.J. (2005). Localization of RecA in *Escherichia coli* K-12 using
900     RecA-GFP. Mol Microbiol *57*, 1074-1085.

901     Saito, Y., Uraki, F., Nakajima, S., Asaeda, A., Ono, K., Kubo, K., and Yamamoto, K. (1997).
902     Characterization of endonuclease III (nth) and endonuclease VIII (nei) mutants of Escherichia
903     coli K-12. J Bacteriol *179*, 3783-3785.

904     Saka, K., Tadenuma, M., Nakade, S., Tanaka, N., Sugawara, H., Nishikawa, K., Ichiyoshi, N.,
905     Kitagawa, M., Mori, H., Ogasawara, N., *et al.* (2005). A complete set of Escherichia coli open
906     reading frames in mobile plasmids facilitating genetic studies. DNA Res *12*, 63-68.

907     Sakaguchi, K., Herrera, J.E., Saito, S., Miki, T., Bustin, M., Vassilev, A., Anderson, C.W., and
908     Appella, E. (1998). DNA damage activates p53 through a phosphorylation-acetylation cascade.
909     Genes Dev *12*, 2831-2841.

910     Schmidt, K., Wolfe, D.M., Stiller, B., and Pearce, D.A. (2009). Cd2+, Mn2+, Ni2+ and Se2+
911     toxicity to Saccharomyces cerevisiae lacking YPK9p the orthologue of human ATP13A2.
912     Biochem Biophys Res Commun *383*, 198-202.

913     Seigneur, M., Bidnenko, V., Ehrlich, S.D., and Michel, B. (1998). RuvAB acts at arrested
914     replication forks. Cell *95*, 419-430.

23

915   Serres, M.H., Gopal, S., Nahum, L.A., Liang, P., Gaasterland, T., and Riley, M. (2001). A
916   functional update of the Escherichia coli K-12 genome. Genome Biol *2*, RESEARCH0035.

917   Shee, C., Cox, B.D., Gu, F., Luengas, E.M., Joshi, M.C., Chiu, L.Y., Magnan, D., Halliday, J.A.,
918   Frisch, R.L., Gibson, J.L.*, et al.* (2013). Engineered proteins detect spontaneous DNA breakage
919   in human and bacterial cells. Elife *2*, e01222.

920   Sottoriva, A., Kang, H., Ma, Z., Graham, T.A., Salomon, M.P., Zhao, J., Marjoram, P.,
921   Siegmund, K., Press, M.F., Shibata, D.*, et al.* (2015). A Big Bang model of human colorectal
922   tumor growth. Nat Genet *47*, 209-216.

923   Stanton, R.C. (2012). Glucose-6-phosphate dehydrogenase, NADPH, and cell survival. IUBMB
924   Life *64*, 362-369.

925   Steighner, R.J., and Povirk, L.F. (1990). Bleomycin-induced DNA lesions at mutational hot
926   spots: implications for the mechanism of double-strand cleavage. Proc Natl Acad Sci U S A *87*,
927   8350-8354.

928   Stern, M.J., Ames, G.F., Smith, N.H., Robinson, E.C., and Higgins, C.F. (1984). Repetitive
929   extragenic palindromic sequences: a major component of the bacterial genome. Cell *37*, 1015-
930   1026.

931   Stratton, M.R. (2011). Exploring the genomes of cancer cells: progress and promise. Science
932   *331*, 1553-1558.

933   Sun, G., Chung, D., Liang, K., and Keles, S. (2013). Statistical analysis of ChIP-seq data with
934   MOSAiCS. Methods Mol Biol *1038*, 193-212.

935   Szklarczyk, D., Franceschini, A., Wyder, S., Forslund, K., Heller, D., Huerta-Cepas, J.,
936   Simonovic, M., Roth, A., Santos, A., Tsafou, K.P.*, et al.* (2015). STRING v10: protein-protein
937   interaction networks, integrated over the tree of life. Nucleic Acids Res *43*, D447-452.

938   Tehranchi, A.K., Blankschien, M.D., Zhang, Y., Halliday, J.A., Srivatsan, A., Peng, J., Herman,
939   C., and Wang, J.D. (2010). The transcription factor DksA prevents conflicts between DNA
940   replication and transcription machinery. Cell *141*, 595-605.

941   Thomason, L.C., Costantino, N., and Court, D.L. (2007). E. coli genome manipulation by P1
942   transduction. Curr Protoc Mol Biol *Chapter 1*, Unit 1 17.

943   Tipparaju, S.M., Liu, S.Q., Barski, O.A., and Bhatnagar, A. (2007). NADPH binding to beta-
944   subunit regulates inactivation of voltage-gated K(+) channels. Biochem Biophys Res Commun
945   *359*, 269-276.

946   Torkelson, J., Harris, R.S., Lombardo, M.J., Nagendran, J., Thulin, C., and Rosenberg, S.M.
947   (1997). Genome-wide hypermutation in a subpopulation of stationary-phase cells underlies
948   recombination-dependent adaptive mutation. EMBO J *16*, 3303-3311.

949  Tubbs, A., and Nussenzweig, A. (2017). Endogenous DNA Damage as a Source of Genomic
950  Instability in Cancer. Cell *168*, 644-656.

951  Uchida, K., Furukohri, A., Shinozaki, Y., Mori, T., Ogawara, D., Kanaya, S., Nohmi, T., Maki,
952  H., and Akiyama, M. (2008). Overproduction of Escherichia coli DNA polymerase DinB (Pol
953  IV) inhibits replication fork progression and is lethal. Mol Microbiol *70*, 608-622.

954  Vafa, O., Wade, M., Kern, S., Beeche, M., Pandita, T.K., Hampton, G.M., and Wahl, G.M.
955  (2002). c-Myc can induce DNA damage, increase reactive oxygen species, and mitigate p53
956  function: a mechanism for oncogene-induced genetic instability. Mol Cell *9*, 1031-1044.

957  Wagner, J., Fujii, S., Gruz, P., Nohmi, T., and Fuchs, R.P. (2000). The beta clamp targets DNA
958  polymerase IV to DNA and strongly increases its processivity. EMBO Rep *1*, 484-488.

959  Wagner, J., Gruz, P., Kim, S.R., Yamada, M., Matsui, K., Fuchs, R.P., and Nohmi, T. (1999).
960  The dinB gene encodes a novel E. coli DNA polymerase, DNA pol IV, involved in mutagenesis.
961  Mol Cell *4*, 281-286.

962  Wagner, J., and Nohmi, T. (2000). Escherichia coli DNA polymerase IV mutator activity:
963  genetic requirements and mutational specificity. J Bacteriol *182*, 4587-4595.

964  Wahba, L., Costantino, L., Tan, F.J., Zimmer, A., and Koshland, D. (2016). S1-DRIP-seq
965  identifies high expression and polyA tracts as major contributors to R-loop formation. Genes
966  Dev *30*, 1327-1338.

967  Xia, J., Chen, L.T., Mei, Q., Ma, C.H., Halliday, J.A., Lin, H.Y., Magnan, D., Pribis, J.P.,
968  Fitzgerald, D.M., Hamilton, H.M.*, et al.* (2016). Holliday junction trap shows how cells use
969  recombination and a junction-guardian role of RecQ helicase. Sci Adv *2*, e1601605.

970  Yang, X., Boehm, J.S., Yang, X., Salehi-Ashtiani, K., Hao, T., Shen, Y., Lubonja, R., Thomas,
971  S.R., Alkan, O., Bhimdi, T.*, et al.* (2011). A public genome-scale lentiviral expression library of
972  human ORFs. Nat Methods *8*, 659-661.

973  Yeeles, J.T., Poli, J., Marians, K.J., and Pasero, P. (2013). Rescuing stalled or damaged
974  replication forks. Cold Spring Harb Perspect Biol *5*, a012815.

975  Yuen, K.W., Warren, C.D., Chen, O., Kwok, T., Hieter, P., and Spencer, F.A. (2007). Systematic
976  genome instability screens in yeast and their potential relevance to cancer. Proc Natl Acad Sci U
977  S A *104*, 3925-3930.

978  Zack, T.I., Schumacher, S.E., Carter, S.L., Cherniack, A.D., Saksena, G., Tabak, B., Lawrence,
979  M.S., Zhsng, C.Z., Wala, J., Mermel, C.H.*, et al.* (2013). Pan-cancer patterns of somatic copy
980  number alteration. Nat Genet *45*, 1134-1140.
981

982  **AUTHOR CONTRIBUTIONS**

1017 **Figure 1**. **Comprehensive Discovery of *E. coli* DNA-Damaging-Protein Network**

1018 **Figure 1**. **Comprehensive Discovery of *E. coli* DNA-Damaging-Protein Network**
1019 (A) We sought comprehensive identification of gain-of-function *E. coli* "DNA-damaging"
1020 proteins (DDPs): proteins that provoke DNA damage when overproduced, modeling frequent
1021 protein overproduction via various mechanisms common in cancers and bacteria. Although
1022 hypotheses for the origins of endogenous DNA damage have been suggested (shown), how these
1023 may arise, whether they are common naturally/spontaneously, and what proteins cause them are
1024 unclear.
1025 (B) Scheme of *E. coli* comprehensive overproduction screen for DDPs. (1) Primary screen: the
1026 complete *E. coli* overexpression Mobile library was screened by plate reader for increased
1027 fluorescence from an SOS-DNA-damage-response-reporter gene, P$_{sulA}$*mCherry* (Nehring et al.,
1028 2016). (2) Secondary screen: potential positive clones from the primary screen were validated,
1029 and false-positives eliminated, by sensitive flow-cytometric assay, which reports fluorescence per
1030 cell at the single-cell level.
1031 (C) Representative results of primary plate-reader screen: afu, arbitrary fluorescence units (SOS
1032 activity), per OD$_{600}$ unit, indicating biomass. Red, potential "damage-up" DDP hits with fold
1033 change >30%.
1034 (D) Representative flow-cytometric validation of SOS-positive (DDP) overexpression clones from
1035 plate-reader screens. Dashed line, flow cytometry "gate" above which cells are scored as SOS-
1036 positive (STAR Methods). Validated DDP clones have significantly higher frequencies of SOS-
1037 positive cells than the vector control. Blue, vector control; red, DDP-overproduction clones.
1038 (E) Quantification of increased DNA damage measured as % SOS-positive cells in the 208
1039 validated *E. coli* DDP clones (identities with data, Table S1).
1040 (F) *E. coli* DDP network summary; proteins of many different functions are DDPs. Functions of
1041 the 208 *E. coli* DDPs (Table S1). Few (8%) are known DNA-repair proteins (blue, Table S1).
1042 (G) Protein-protein associations of the *E. coli* over-production DDP network, CytoScape software
1043 generated from STRING 10.0 database. Other, defined Table S1; specific protein associations,
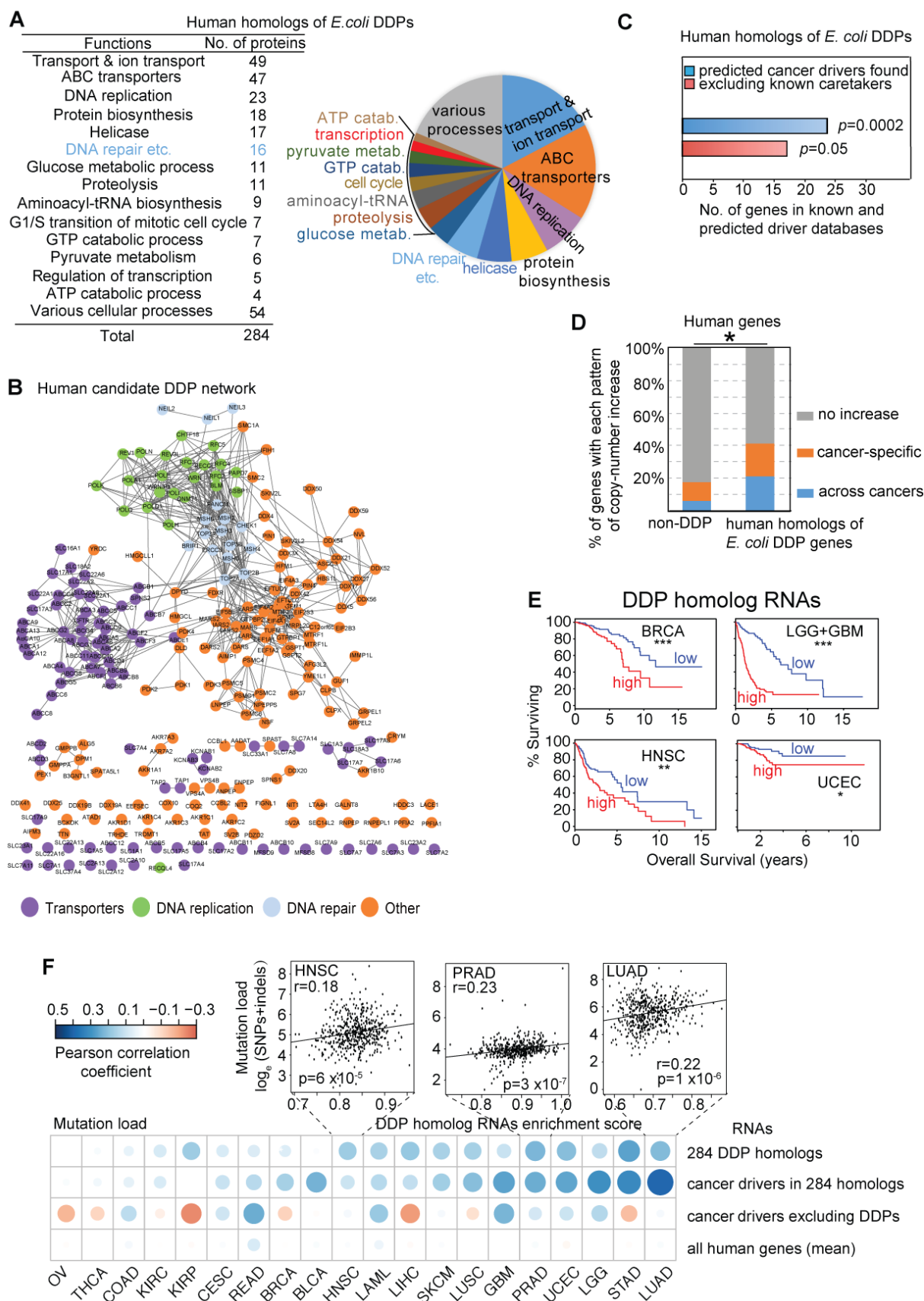1044 Figure S2 (discussed, text).
1045 (H) LexA-dependence of increased fluorescence of representative DDPs shows that fluorescence
1046 results from activation of the SOS DNA-damage response.
1047 (I) Correlation of DDP (SOS-positive) phenotype with RecA*GFP foci, indicating persistent
1048 single-stranded (ss)DNA. A representative sample of 67 DDPs overproduced showed 32 (48%)
1049 with significantly more RecA*GFP foci than the vector control ($p < 0.05$, unpaired two-tail *t*-test),
1050 a less sensitive assay than the stringent flow-cytometric assay for SOS-positive cells. RecA*GFP
1051 foci are correlated positively with the SOS-response assay: r = 0.7, $p = 1.3 \times 10^{-10}$, Pearson's
1052 correlation. **Left**, DDP clones assayed. Each dot, one DDP clone, assayed twice (mean). Blue,
1053 negative control; orange, low SOS activity (DNA damage); green moderate DNA damage; purple,
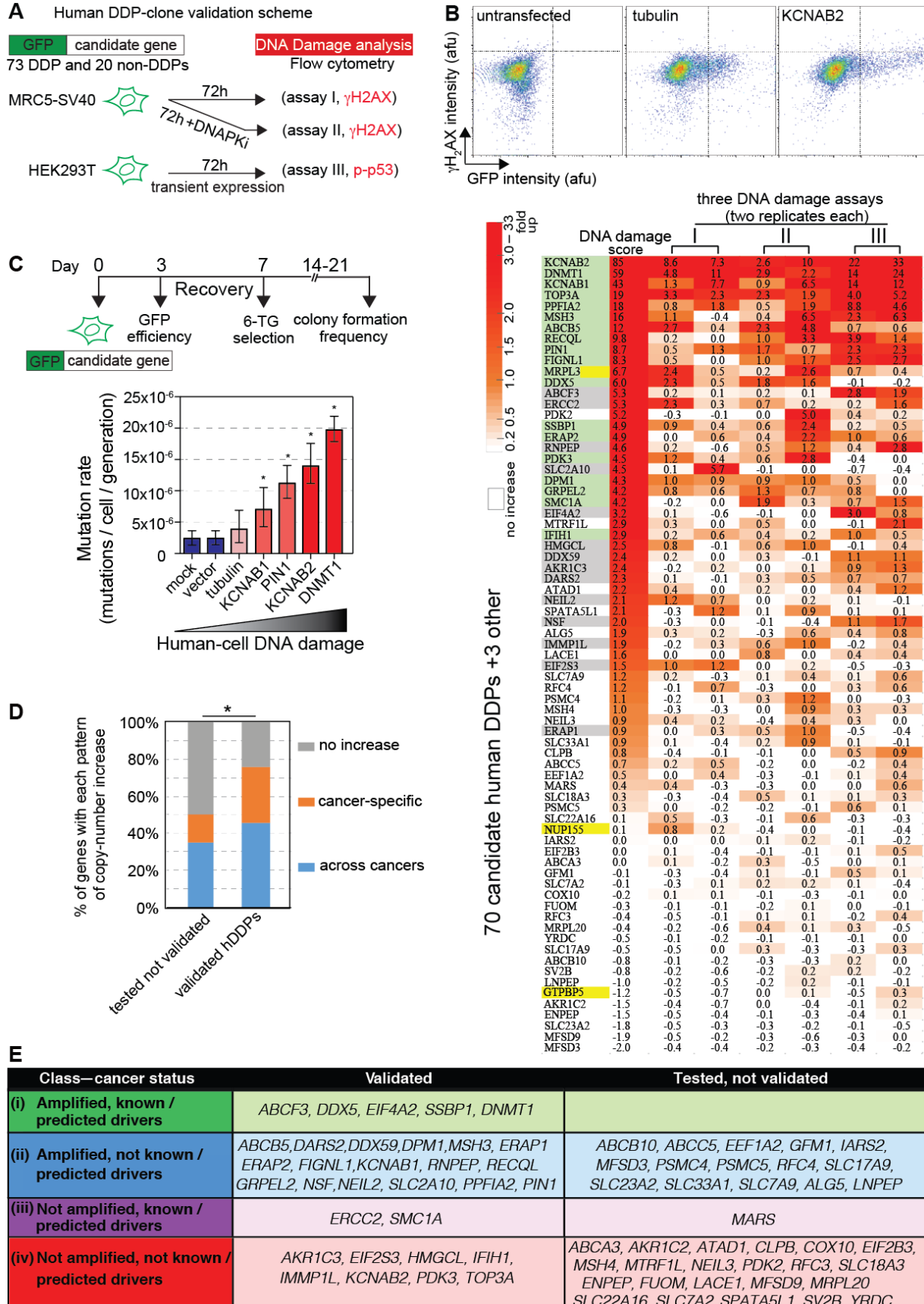1054 high DNA damage, strain-by-strain data, Table S1. **Right**, representative foci.
1055 (J) Mutation-rate increase in representative DDP-overproducing clones. Above, *c*I mutation assay
1056 design (Gutierrez et al., 2013). Loss-of-function mutations in a chromosomal *c*I gene, encoding
1057 phage lambda transcription repressor, allow transcription of *tetA*, which confers tetracycline
1058 resistance (TetR) (STAR Methods). Below, increased mutation rates are associated with increased
1059 DNA-damage levels (SOS induction, flow-cytometric assay) in *E. coli* DDP-overproducing
1060 clones. Each bar, the mean mutation rate (± SEM) of each strain, 3 experiments (fluctuation tests,
1061 STAR Methods). The DDPs overproduced (Supplemental Discussion 3) represent various classes
1062 (Table S1). Table S1 for mutation rates. *P*-values, one-way Fisher's exact test of the number of

28

1063    clones with mutation rate significantly higher than the vector-only control (unpaired 2-tailed
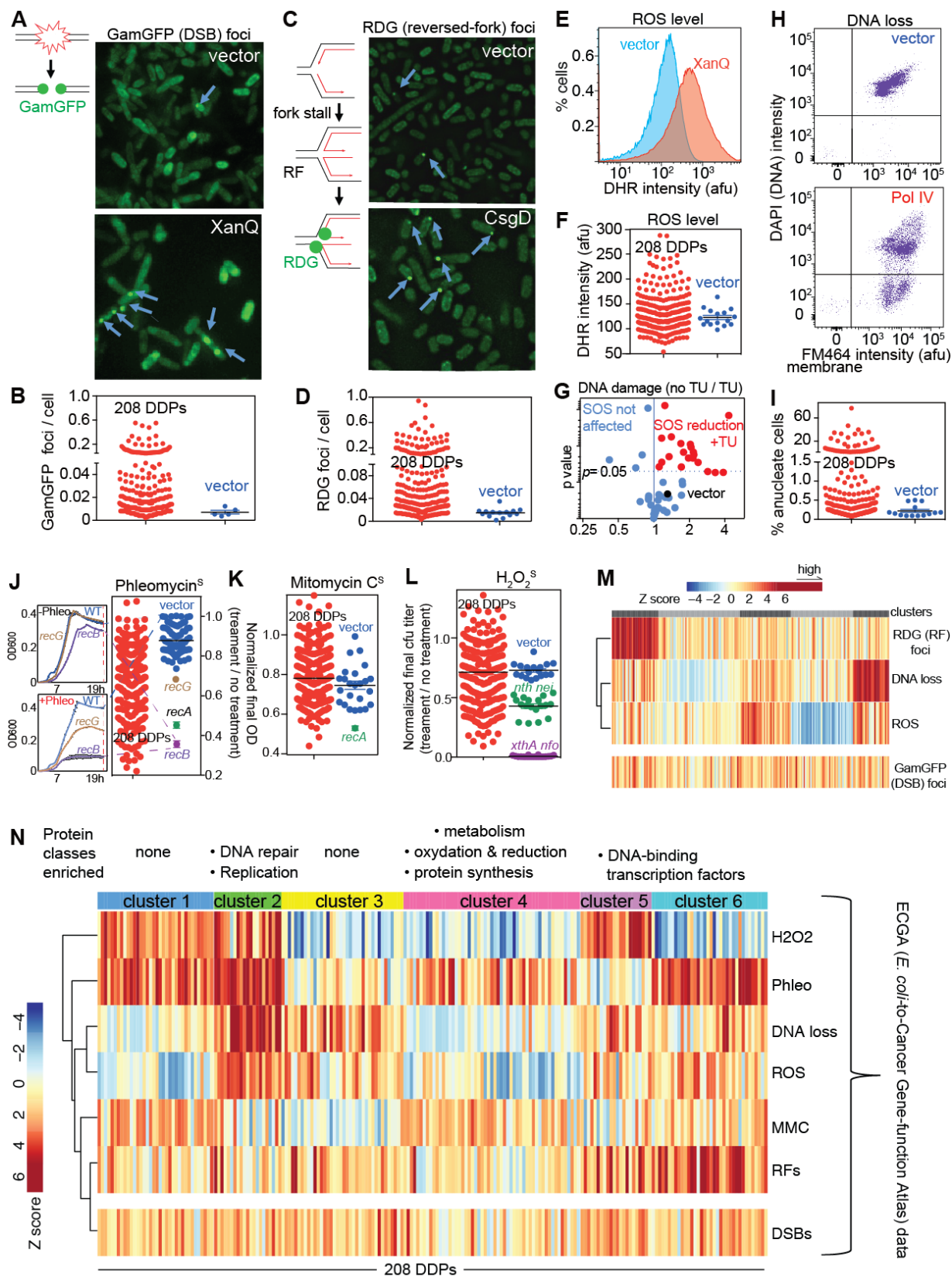1064    Student's *t*-test).
1065

1066



**A** Human homologs of *E.coli* DDPs

| Functions | No. of proteins |
| --- | --- |
| Transport & ion transport | 49 |
| ABC transporters | 47 |
| DNA replication | 23 |
| Protein biosynthesis | 18 |
| Helicase | 17 |
| DNA repair etc. | 16 |
| Glucose metabolic process | 11 |
| Proteolysis | 11 |
| Aminoacyl-tRNA biosynthesis | 9 |
| G1/S transition of mitotic cell cycle | 7 |
| GTP catabolic process | 7 |
| Pyruvate metabolism | 6 |
| Regulation of transcription | 5 |
| ATP catabolic process | 4 |
| Various cellular processes | 54 |
| Total | 284 |

**B** Human candidate DDP network

Transporters · DNA replication · DNA repair · Other

**C** Human homologs of *E. coli* DDPs

predicted cancer drivers found
excluding known caretakers

$p=0.0002$
$p=0.05$

No. of genes in known and predicted driver databases

**D** Human genes

% of genes with each pattern of copy-number increase

no increase
cancer-specific
across cancers

non-DDP    human homologs of *E. coli* DDP genes

**E** DDP homolog RNAs

BRCA ***    LGG+GBM ***
low / high

HNSC **    UCEC *
low / high

% Surviving

Overall Survival (years)

**F**

Pearson correlation coefficient

HNSC r=0.18    p=6 ×10⁻⁵
PRAD r=0.23    p=3 ×10⁻⁷
LUAD r=0.22    p=1 ×10⁻⁶

Mutation load $\log_e$ (SNPs+indels)

DDP homolog RNAs enrichment score

Mutation load

RNAs

284 DDP homologs
cancer drivers in 284 homologs
cancer drivers excluding DDPs
all human genes (mean)

OV THCA COAD KIRC KIRP CESC READ BRCA BLCA HNSC LAML LIHC SKCM LUSC GBM PRAD UCEC LGG STAD LUAD

30

**Figure 2**. **Human Homologs of *E. coli* DDPs a Network Associated with Cancers**

(A) Summary of 284 human proteins identified as homologs of *E. coli* DDPs (STAR Methods). Only 5.6% are known DNA-repair genes (blue, Table S2).

(B) Protein-protein association network of human homologs of *E. coli* DDPs (Figure S2 for specific interactions). Network displayed per Figure 1G. Other, defined Table S2.

(C) Human homologs of *E. coli* DDPs (hDDP candidates) are significantly overrepresented among known (Forbes et al., 2015) and predicted (D'Antonio and Ciccarelli, 2013) cancer drivers (blue bar). After subtracting known DNA-repair ("caretaker") genes, the remaining hDDP candidate genes are still enriched for known and predicted cancer drivers (red bar).

(D) hDDP candidate genes are enriched among genes with cancer-associated copy-number increases, indicating selection for overexpression in cancers. Pan-Cancer copy-number-increase analysis (GISTIC threshold copy-number gain ≥ 1) of the 284 hDDP candidates in 26 cancer types (per Figure S3). Blue, human genes with increased copy numbers across cancers ($p < 0.05$, FDR $< 0.10$, Wilcoxon test); orange, genes with cancer-specific copy-number increase; grey, not particularly cancer associated. Complete data for the network Table S4.

(E) Decreased cancer survival is associated with high DDP-homolog RNA levels in cancers [our analyses of data from TCGA (Gao et al., 2013), STAR Methods]. BRCA, breast invasive carcinoma; LGG+GBM, gliomas (low-grade glioma + glioblastoma multiforme); HNSC, head and neck squamous cell carcinoma; UCEC, uterine corpus endometrial carcinoma. *, **, ***, survival of the cancers with high and low levels of the 284 RNAs differ at $p \leq 0.05$; $\leq 0.01$, and $\leq 0.001$ respectively, log-rank test.

(F) High RNA levels of the 284 hDDP candidates are associated with tumor mutation burden in data from TCGA (Gao et al., 2013). Each dot represents a Pearson correlation coefficient between the RNA-enrichment score, relative to total RNAs, of a gene set and mutation burden in the tumor. The average correlation strength of 284 hDDP-candidate RNA levels with mutation loads across 20 TCGA cancers was in the top 0.5% of correlations for randomly selected groups of genes across all human genes. Blown-up: correlation of increased mutation loads (y axis) with increased hDDP-candidate RNA enrichment scores (x axis, STAR Methods) in three represented cancers. Cancer-types: OV ovarian serous cystadenocarcinoma; THCA thyroid carcinoma; COAD colon adenocarcinoma; KIRC kidney renal clear cell carcinoma; KIRP kidney renal papillary cell carcinoma; CESC cervical squamous cell carcinoma and endocervical adenocarcinoma; READ rectum adenocarcinoma; BLCA bladder urothelial carcinoma; LAML acute myeloid leukemia; LIHC liver hepatocellular carcinoma; SKCM skin cutaneous melanoma; LUSC lung squamous cell carcinoma; PRAD prostate adenocarcinoma; STAD stomach adenocarcinoma; LUAD lung adenocarcinoma.

31

**Figure 3. Overproduced Human Homologs Promote DNA Damage in Human Cells**

**Figure 3**. **Overproduced Human Homologs Promote DNA Damage in Human Cells**

(A) Scheme for validating hDDPs. 70 full-length sequence-verified human homolog candidate-hDDP-GFP N-terminal fusions (and 3 damage-down-, plus 20 non-DDP-GFP fusion controls) were transiently overproduced in MRC5-SV40 or HEK293T cell lines and green cells screened for high DNA damage by flow cytometry.

(B) DNA-damage assays with candidate hDDPs identify 33 validated hDDPs. **Upper:** representative flow cytometric DNA-damage assay data (STAR Methods, Figure S3G-I; Table S6). **Lower:** heatmap of the flow-cytometric data normalized to GFP-tubulin. Data from each of the three DNA-damage assays, with 2 replicates each, ranked by a cumulative DNA-damage score that sums the fold changes of each DNA-damage assay for each candidate protein. Highlighting colors: green, significantly damage-up in $\geq 2$ assays (two-tailed unpaired $t$-test with FDR correction); gray, significantly damage-up in one assay; yellow, homologs of *E. coli* damage-down proteins; white, not damage-up. 45% validated, significantly more than among 20 random human genes ($p$ <0.0001, two-tailed unpaired $t$-test with FDR correction, Figure S3G-I), indicating that homologs of *E. coli* DDPs are enriched for hDDPs.

(C) Increased mutation rates in human cells overproducing validated hDDPs, assayed by human-cell *HPRT* forward-mutation assays in fluctuation tests (STAR Methods). **Upper**: *HPRT* assay scheme. *HPRT* loss-of-function mutants are selected as 6-thioguaine (6-TG)-resistant clones. **Lower**: mutation rates of selected hDDP overproducers shown with their DNA-damage levels; error bars, 95% confidence intervals.

(D) Validated hDDP genes are enriched among genes with cancer-associated copy-number increases compared with the candidates that were tested but not validated ($p$ = 0.02, one-way Fisher's exact test).

(E) New and known potential cancer-promoters predicted among 33 validated hDDPs. The 33 validated hDDP genes comprise genes that are—(i) both amplified in TCGA cancers and known (Forbes et al., 2015) or predicted (D'Antonio and Ciccarelli, 2013) cancer drivers (16%); (ii) amplified in TCGA cancers and were not known or predicted cancer-driving genes (53%); (iii) known or predicted cancer drivers that are not found to be amplified in cancers (6%); and (iv) not found to be amplified in cancers in TCGA, and not known or predicted cancer drivers (25%). The data suggest potential overexpression cancer-promoting roles for the genes in all classes.

33

1137 **Figure 4. Kinds, Causes and Consequences of DNA Damage from *E. coli* DDPs Reveal**
1138 **Function Clusters**

1139 **Figure 4. Kinds, Causes and Consequences of DNA Damage from *E. coli* DDPs Reveal**
1140 **Function Clusters**
1141 (A) Identification of DDPs that increase DNA double-strand breaks (DSBs), detected as GamGFP
1142 foci, per (Shee et al., 2013). Representative image in cells overproducing DDP XanQ, a membrane-
1143 spanning transporter. Diagram: lines, DNA strands; green balls, GamGFP.
1144 (B) 87 of the 208 *E. coli* DDPs promote DSBs. Quantification of GamGFP foci in 208 DDP-
1145 overproducing clones (means of two experiments, >1000 cells each; data by clone Table S1).
1146 (C) Detection of stalled, reversed replication forks (RFs) as RDG foci in *recA⁻* cells. Engineered
1147 protein RuvCDefGFP (RDG) traps 4-way DNA-junctions, and in *recA* cells detects only RFs (Xia
1148 et al., 2016). Representative image in cells overproducing DDP CsgD. Diagram: lines, DNA
1149 strands; red lines, newly synthesized strands; green balls, RDG.
1150 (D) 106 of the 208 *E. coli* DDPs cause fork stalling and reversal. Quantification of RDG foci in
1151 208 DDP-overproducing clones (data by clone, Table S1).
1152 (E-I) Flow-cytometric assays for—
1153 (E) elevated ROS measured as peroxide by dihydrorhodamine (DHR) fluorescence. Representative
1154 flow-cytometric histogram, XanQ overproducer, and
1155 (F) quantified in all 208 DDP-overproducing clones, mean fluorescence intensity; afu, arbitrary
1156 fluorescence units (2 experiments, mean) shows 56 (27%) of DDP clones (data by clone, Table
1157 S1).
1158 (G) Increased DNA damage in 16 of 56 high-ROS DDP clones is reduced by ROS-quenching
1159 agent thiourea (TU) indicating that the high ROS underlie the DNA damage. *P* values, unpaired
1160 two-tailed *t* test. Data by clone, Table S1.
1161 (H) DNA loss: the fraction of cells with no DNA (anucleate cells), representative example, Pol IV
1162 overproducer (DAPI indicates DNA, membrane dye FM464 indicates cells); events below the
1163 horizontal line scored as anucleate, DNA loss,
1164 (I) quantified in all 208 DDP-overproducing clones (2 experiments, mean), showing 67 (32%) of
1165 DDP clones (data by clone, Table S1).
1166 (J-L) Sensitivity to DNA-damaging agents in DDP-overproducing clones implies DNA-repair-
1167 pathway reduction (possible saturation) as a potential consequence of elevated DNA damage.
1168 Positive controls: relevant DNA-repair-defective mutants indicated. Each measured as slowed
1169 growth curves per J left (STAR Methods). Data by clone, Table S1.
1170 (J) Phleomycin sensitivity (reduced homology-directed DSB repair) seen in 106 of the DDP clones
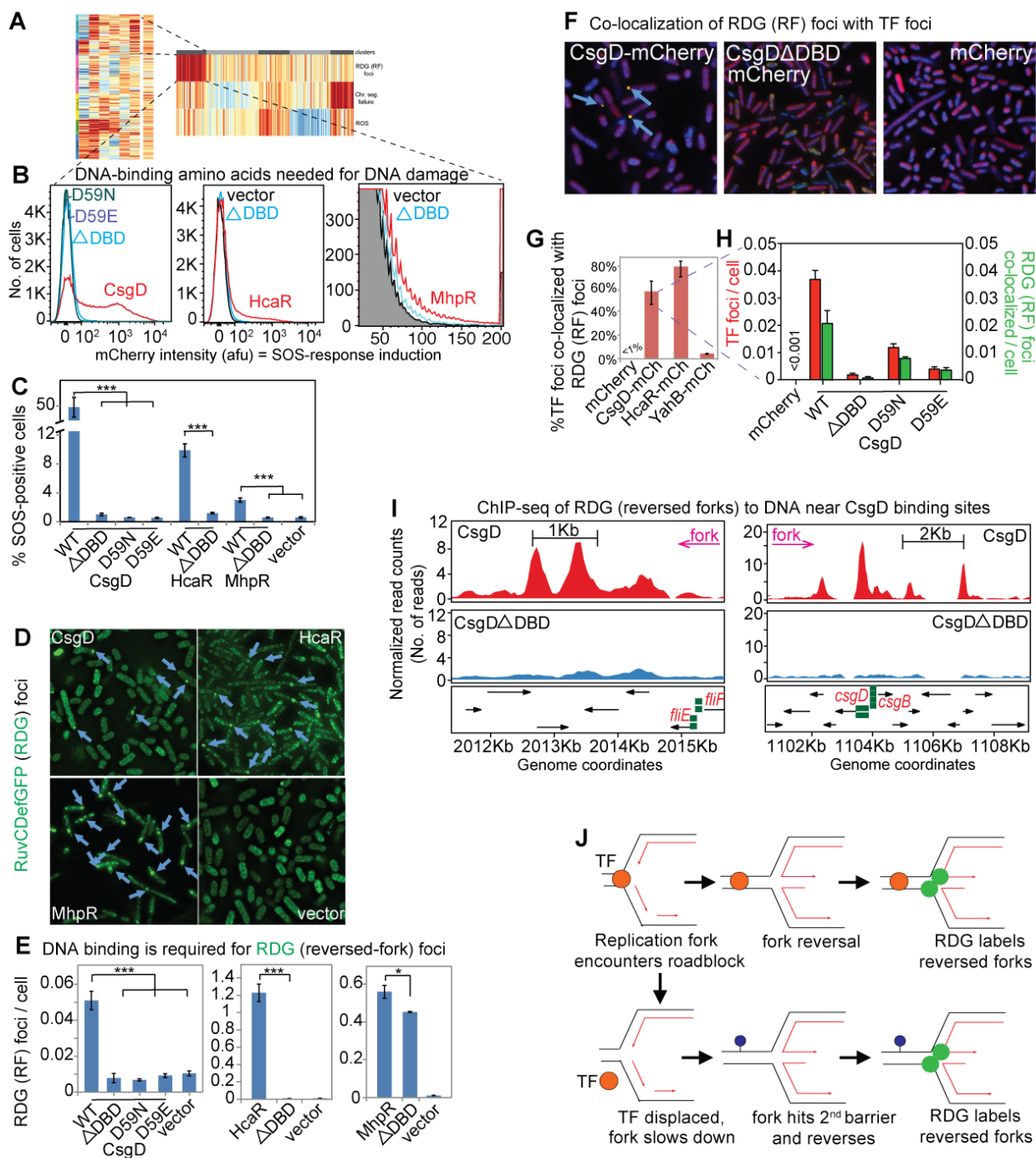1171 (51%).
1172 (K) Mitomycin C (MMC) sensitivity (reduced nucleotide-excision repair and/or homology-
1173 directed repair) seen 10 of the DDP clones (5%).
1174 (L) $H_2O_2$ sensitivity (reduced base-excision repair) seen in 75 of the DDP clones (36%).
1175 (M) Stalled replication (RFs) is clustered among particular DDP overproducers; DNA breakage is
1176 not. Progeny clustering (STAR Methods).
1177 (N) Cluster analysis of Z (significance) scores of assays: $H_2O_2$, hydrogen-peroxide sensitivity;
1178 Phleo, phleomycin sensitivity; DNA loss (anucleate cells); ROS (ROS levels); MMC, mitomycin-
1179 C sensitivity; RFs (reversed forks), RDG (RF) foci; DSBs, GamGFP (DSB) foci. Vertical bars
1180 along the x axis: the phenotype of each DDP clone. The 6 clusters indicate 6 DNA-damage
1181 signatures and suggest at least 6 different mechanisms of DNA-damage generation in the DDP
1182 network. Protein categories significantly increased in each cluster shown above (one-way Fisher's
1183 exact test), cluster 2 at $p = 0.01$, cluster 4 at $p = 0.01$, and clusters 5 and 6 at $p = 0.03$.
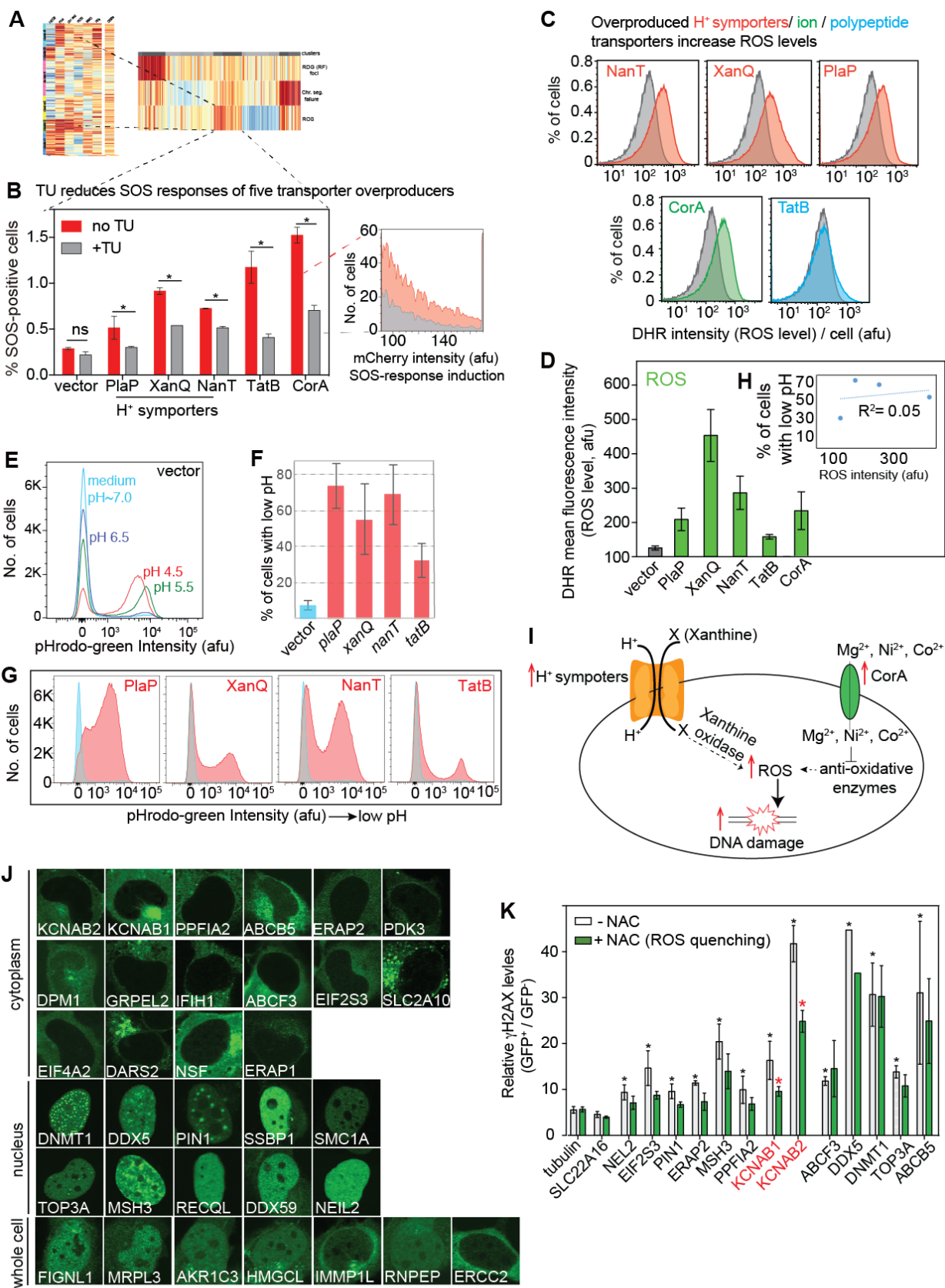1184

**Figure 5. *E. coli* Transcription Factors Promote Replication-fork Stalling and Reversal DNA-binding-domain Dependently**

(A) DNA-binding transcription factors (TFs) are enriched among DDP clones with high RFs ($p$ =0.002, one-way Fisher's exact test).

(B) DNA-binding ability is required for DNA damage/SOS activity caused by overproduced DNA-binding TFs. Representative flow cytometry histograms of SOS induction, three TFs and their corresponding mutants: ΔDBD, DNA-binding domain in-frame deletion; D59N, D59E, single amino-acid changes that reduce CsgD DNA-binding (Ogasawara et al., 2011).

36

1193   (C) Mean ± SEM of ≥3 experiments.
1194   (D) DNA-binding ability of overproduced TFs is required for increased RDG (RF) foci (blue
1195   arrows). Representative images. Figure S6A, all genotypes.
1196   (E) Mean ± SEM of ≥3 experiments.
1197   (F and G) mCherry-protein fusions of DNA-binding TFs form foci that co-localize with RGD (RF)
1198   foci, placing TF binding and RFs in relative proximity (within 50kb, text) in the 4.6MB *E. coli*
1199   genome.
1200   (F) Representative data. CsgD-mCherry foci co-localize with RDG foci dependently on the CsgD
1201   DNA-binding domain (DBD). Blue arrows, co-localized red and green (CsgD-mCherry and RDG)
1202   foci. Figure S6B, all genotypes.
1203   (G) Mean ± SEM of ≥3 experiments with CsgD-, HcaR-, and YahB-mCherry co-localization with
1204   RDG.
1205   (H) Co-localization of CsgD-mCherry foci with RDG foci requires CsgD DNA-binding ability;
1206   quantification.
1207   (I) RDG ChIP-Seq peaks in HR-defective Δ*recA* cells (RFs) are enriched near CsgD-binding sites
1208   (green squares; *p* =0.01, two-tailed z-test compared with simulated data, see Supplemental
1209   Discussion 12). Representative peaks shown; Figure S7 for the complete set of RF peaks.
1210   (J) Model: overproduced TFs (orange circles) binding to DNA (parallel lines, basepaired strands)
1211   cause RFs by replication roadblock.  Green circles, RDG bound to RF.
1212

**A**

**B** TU reduces SOS responses of five transporter overproducers

**C** Overproduced H+ symporters/ ion / polypeptide transporters increase ROS levels

**D** ROS

**E**

**F**

**G**

**H**

**I**

**J**

**K**

1213

**Figure 6. *E. coli* and Human Transmembrane Transporters Promote DNA Damage via Increased ROS**

(A) *E. coli* high ROS cluster is enriched for membrane-spanning transporters ($p = 0.004$ one-way Fisher's exact test).

(B) DNA damage (SOS activity) from five overproduced *E. coli* transporters is partially reversed by ROS-scavenger thiourea (TU), implying ROS-dependent DNA damage. Quantification of flow cytometry per blow up. Mean ± range, 2 experiments. Blow up: representative flow cytometry for DNA damage: cells with chromosomal SOS-promoter-mCherry fusion.

(C) ROS levels increase upon overproduction of various *E. coli* membrane-spanning transporters. ROS measured as $H_2O_2$ shown by DHR stain and flow cytometry. Representative data (Table S1 for all). Gray, vector only; red, $H^+$ symporters; green, ion transporter; cyan, polypeptide transporter.

(D) Means ± range of 2 experiments.

(E-H) Increased *E. coli* $H^+$ symporter activity is caused by overproduction, shown by reduced cellular pH. (E) Detection of *E. coli* intracellular pH by pHrodo-green dye staining followed by flow cytometry. Cells with the vector were exposed to buffers with varied pH levels and pHrodo-green dye was used to stain the cells. Control vector-bearing cells exposed to low pH cells had subpopulations with increased pHrodo-green intensity.

(F, G) Reduced intracellular pH in clones overproducing the PlaP, XanQ, or NanT $H^+$ symporters, or the TatB polypeptide transporter, indicates that overproduction causes gain-of-function, overall increased activity per cell of these transporters. Cyan, vector only. (F) Quantification: mean ± range of two experiments. (G) Representative flow-cytometry histograms.
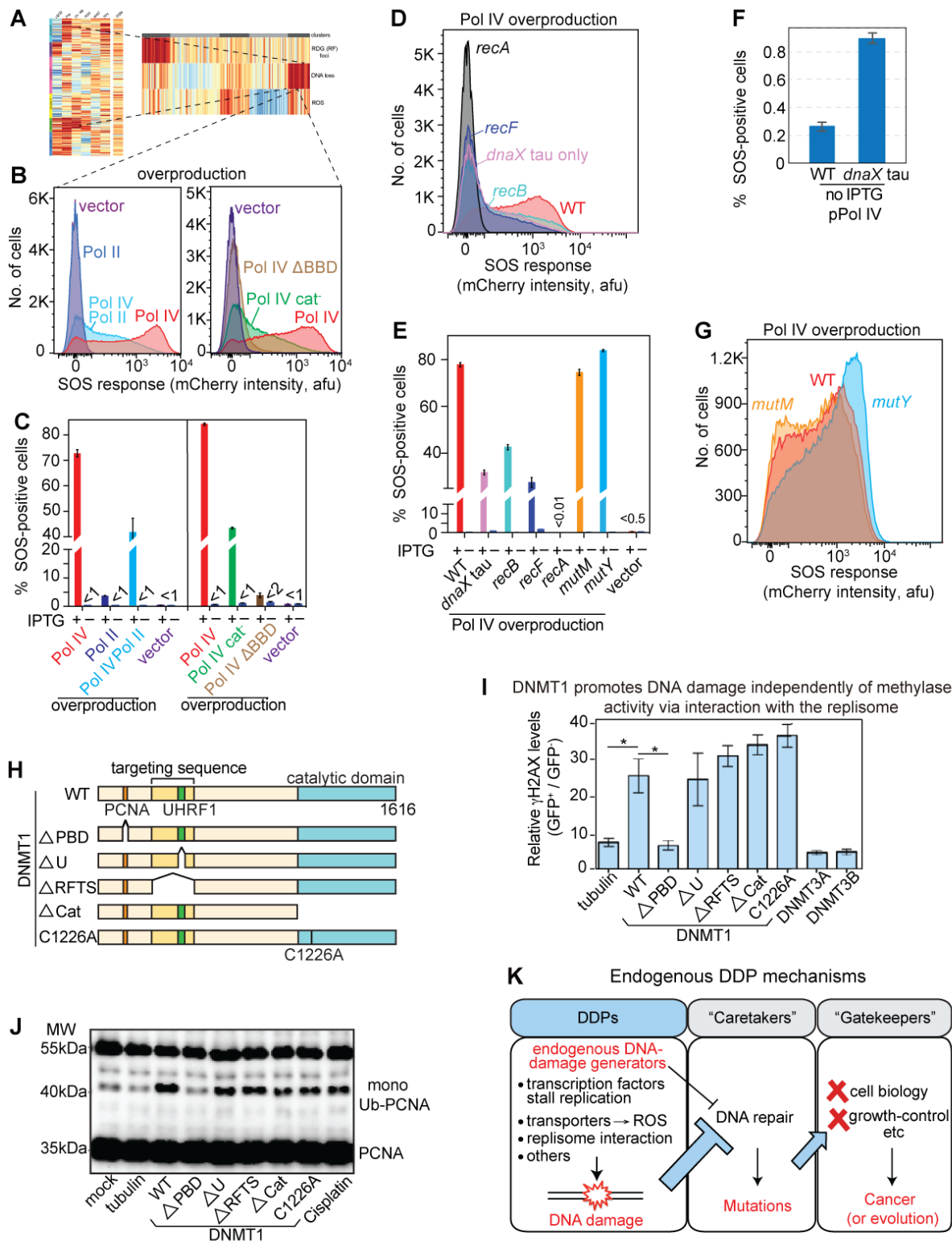
(H) Reduced pH in *E. coli* transporter-overproducing clones—the three $H^+$ symporters and TatB polypeptide transporter—is not correlated quantitatively with increased ROS ($R^2=0.05$, Pearson's correction analysis), suggesting that the specific cargoes, not low pH, promote DNA damage.

(I) Models for ROS-dependent DNA-damage promotion by overproduction of the *E. coli* XanQ and CorA transporters. Left: Overproduced $H^+$ symporter XanQ might cause ROS by increased import of xanthine which is oxidized by the ROS-generating xanthine oxidase (Kelley et al., 2010). Right: CorA overproduction might cause DNA damage via increased import of $Ni^{2+}$, which inhibits anti-oxidative enzymes (Schmidt et al., 2009) and causes DNA damage ROS-dependently (Cameron et al., 2011).

(J) Overproduced validated human (h)DDPs are localized in various subcellular compartments, implying various DDP mechanisms. Of the 33 overproduced validated hDDPs, 16 were detected only in the cytoplasm, 10 only in the nucleus, and 7 throughout the cell. All show qualitatively repeatable subcellular localization.

(K) ROS underlie at least some of the DNA damage caused by human KCNAB1/2 transporter overproduction. NAC: N-acetyl-cysteine, an ROS quencher (STAR Methods). * $p < 0.05$ relative to the NAC-untreated GFP-tubulin control; * $p < 0.05$ relative to the corresponding NAC-untreated control.

1254
**Figure 7. *E. coli* DNA Pol IV and Human DNMT1 Promote DNA Damage via Interaction**
1256 **with the Replisome Clamp**
1257 (A) *E. coli* function cluster with high DNA loss includes DNA Pol IV.

1258      (B) *E. coli* Pol IV promotion of DNA damage is reduced by co-overproduction of its competitor
1259      at the replisome, DNA Pol II (left). Right: Pol IV promotion of DNA damage requires its
1260      interaction with beta replisome sliding clamp, seen by dependence on the Pol IV beta clamp-
1261      binding domain (BBD), and is partly independent of Pol IV catalytic activity (cat$^-$ mutant).
1262      Representative data.
1263      (C) Mean ± SEM of ≥3 experiments.
1264      (D) Reduction of *E. coli* Pol IV-induced DNA damage in a clamp loader tau-only mutant, which
1265      interacts with replicative DNA Pol III in preference to Pol IV (Dohrmann et al., 2016). RecB- and
1266      RecF-dependence of the SOS response induced by overproduced Pol IV implicates DSBs and
1267      single-strand gaps, respectively, as DNA-damage types produced. Representative data.
1268      (E) Mean ± SEM of ≥3 experiments. Pol IV is induced by IPTG.
1269      (F) The clamp-loader *dnaX* tau-only mutant is not generally deficient in SOS-response induction,
1270      but merely reduces the DNA damage (SOS activity) caused by Pol IV upregulation (D, E).
1271      (G) Neither the *E. coli* MutM 8-oxo-dG glycosylase nor the MutY adenine glycosylase are required
1272      for DNA damage instigated by Pol IV, indicating DNA damage produced independently of
1273      incorporation of 8-oxo-dG into DNA. Representative data, quantified (E).
1274      (H) Constructs for production of human wild-type and truncated DNA methyltransferase DNMT1
1275      in human cells. PBD, PCNA-binding domain; U, UHRF1, ubiquitin-like containing PHD and
1276      RING-finger domains 1 binding domain; RFTS, replication-focus-targeting sequence, recruits
1277      DNMT1 to DNA-methylation sites; Cat, catalytic domain for methyltransferase activity; C1226A,
1278      mutation of the catalytic active site.
1279      (I) Human DNMT1 overproduced in human cells promotes γH2AX (DNA-break indicator)
1280      accumulation methylase-independently and replisome-clamp-interaction dependently.
1281      Overproduction of two DNMT1 catalytically dead mutants increased DNA damage similarly to
1282      overproduced WT DNMT1. Overproduction of two other de novo DNA methyltransferases
1283      (DNMT3A, DNMT3B) did not elevate DNA damage. DNA-damage promotion by DNMT1
1284      requires the DNMT1 PBD domain, required for DNMT1 binding to the human replisome clamp:
1285      PCNA.
1286      (J) Human DNMT1 overproduction promotes PCNA monoubiquitination, a replication-stress
1287      indicator, in a replisome-interaction-dependent manner. Monoubiquitination determined by
1288      western blot with anti-PCNA antibody (STAR Methods).
1289      (K) The endogenous DDP model for cancer promotion and mechanisms: DDPs a cancer-protein
1290      functional class upstream of DNA repair. Excessive endogenous DNA damage is proposed to
1291      titrate (thick blue -|) or inhibit (thin black -|) DNA repair causing DNA-repair ("caretaker")-protein
1292      deficiency in cells without a DNA-repair-gene mutation. Repair deficiency increases mutation rate,
1293      leading to cancer- (or evolution-) driving mutations in the cell-biology altering "gatekeeper" genes
1294      that cause the cancer cell-biological phenotypes.
1295
1296
1297

1298 **STAR★METHODS**

1299 **KEY RESOURCES TABLE**

| REAGENT or RESOURCE | SOURCE | IDENTIFIER |
|---|---|---|
| **Antibodies** | | |
| anti-PCNA | Santa Cruz | sc-56 |
| anti-γH2AX | Millipore | 05-636 |
| anti-phospho-P53 | Cell Signaling | 9286 |
| anti-RAD18 | Cell Signaling | 9040S |
| anti-Tubulin | Abcam | ab6046 |
| anti-GFP | ThermoFisher Scientific | A11122 |
| Anti-rabbit IgG, HRP-linked Antibody | Cell Signaling | 7074 |
| Anti-mouse IgG, HRP-linked Antibody | Cell Signaling | 7076 |
| Goat anti-Mouse IgG (H+L) Highly Cross-Adsorbed Secondary Antibody, Alexa Fluor 647 | Invitrogen | A21236 |
| anti-Pol IV polyclonal | Kim et al., 2001 | N/A |
| anti-RuvC (5G9/3) monoclonal | Santa Cruz | sc-53437 |
| anti-Mouse IgG | Bethyl Laboratories | A90-116D5 |
| anti-Goat IgG | Bethyl Laboratories | A50-100D5 |
| **Bacterial and Virus Strains** | | |
| *E. coli* mobile plasmid overexpression library | (Saka et al., 2005) | N/A |
| See Table S7 for all bacterial strains used | | |
| **Human cell lines** | | |
| MRC5-SV40 | Stephen P. Jackson Lab | N/A |
| HEK293T | ATCC | CRL3216 |
| **Plasmid** | | |
| See Tables S5 and 7 for all plasmids used | | |
| **Chemicals** | | |
| GenJet™ In Vitro DNA Transfection Reagent | SignaGen Laboratories | SL100489 |
| Polyethylenimine | Polysciences | 23966 |
| Lipofectamine RNAiMAX Transfection Reagent | Invitrogen | 13778030 |
| N-acetyl cysteine | Sigma | A7250 |
| DNA-PK inhibitor | Tocris Bioscience | NU7441 |
| 6-thioguanine | Sigma | A4882 |
| **Critical Commerial assays** | | |
| FM® 4-64FX | Thermal Fisher | F34653 |

| DHR123 | Thermal Fisher | D23806 |
|---|---|---|
| pHrodo® Green AM Intracellular pH indicator | Thermal Fisher | P35373 |
| RNeasy Mini Kit | Qiagen | 74104 |
| RiboZero | Illumina | MRZB12424 |
| TruSeq Stranded mRNA Library Preparation Kit | Illumina | RS-122-2001 |
| qPCR-based Illumina Library Quantification Kit | KAPA Biosystems | KK4828 |
| Dneasy Blood & Tissue kits | Qiagen | 69506 |
| QIAprep Spin Miniprep Kit | Qiagen | 27106 |
| Gateway™ LR Clonase™ II Enzyme mix | Invitrogen | 11791100 |
| SuperScript™ III Reverse Transcriptase | Invitrogen | 18080-093 |
| Q5 High-Fidelity DNA Polymerase | New England Biolabs | M0491S |
| **Oligonucleotides** | | |
| ON-TARGETplus Non-targeting Pool | Dharmacon | D-001810-10-05 |
| siRAD18 ACUCAGUGUCCAACUUGCU | Sigma | N/A |
| cI forward primer ACCGCGGCGTGGGTAGTAAAGT | (Gutierrez et al., 2013) | N/A |
| cI reverse primer GCCAATCCCCATGGCATCGAGTAAC | (Gutierrez et al., 2013) | N/A |
| **Deposited Data** | | |
| RNA-Seq data | This paper | ENA: PRJEB21034 |
| ChIP-Seq data | This paper | ENA: PRJEB21035 |
| **Software and Algorithms** | | |
| cBioportal | Cerami et al., 2012; Gao et al., 2013 | http://www.cbioportal.org/ |
| R programming language | R Development Core Team, 2015. | https://www.R-project.org/ |
| *E. coli*-to-Cancer Gene-function Atlas (ECGA) | This paper | https://microbialphenotypes.org/wiki/index.php/Special:ECGA |
| R package for progeny clustering: *ProgenyClust* | CRAN | https://cran.r-project.org/ |
| Trimmomatic | Bolger et al., 2014 | N/A |
| BWA-MEM | Li, 2013 | N/A |
| deepTools | Ramirez et al., 2014 | N/A |
| MOSAiCS | Sun et al., 2013 | N/A |
| Rockhopper | McClure et al., 2013 | N/A |
| Prism | GraphPad | https://www.graphpad.com/scientific-software/prism/ |

| | | |
|---|---|---|
| R programming language | R Development Core Team, 2015. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. | https://www.R-project.org/ |
| FlowJo 10.2 | FLOWJO | https://www.flowjo.com/ |
| STRING 10.0 | (Szklarczyk et al., 2015) | https://string-db.org/ |
| FACSDivaTM | BD Biosciences | http://www.bdbiosciences.com |
| Advanced imaging collection in Pipeline pilot 8.5 or 9.2 | Biovia-Dassault Systems | N/A |
| Softworx | GE | N/A |
| BLASTp and delta BLAST | NCBI | https://blast.ncbi.nlm.nih.gov/Blast.cgi |
| RNA-sequencing data | The Cancer Genome Atlas (Gao et al., 2013) | N/A |
| Gene-Set Enrichment Analysis (ssGSEA) using GSVA package | (Hanzelmann et al., 2013) | N/A |

1300
1301
1302 **CONTACT FOR REAGENT AND RESOURCE SHARING**
1303 Corresponding authors, S. M. Rosenberg (smr@bcm.edu) and K. M. Miller
1304 (kyle.miller@austin.utexas.edu) are the contacts for reagents and resource sharing.

1305

1306 **EXPERIMENTAL MODEL AND SUBJECT DETAILS**
1307 *Escherichia coli* K-12 (strains MG1655 and W3110) and isogenic derivatives were used for all
1308 bacterial experiments. Human MRC5-SV40 and HEK293T cells were used for all human cell line
1309 experiments.
1310

1311 **METHOD DETAILS**
1312 *Escherichia coli* **strains and media**
1313 *E. coli* K12 strains and plasmids used in this work are shown in Table S7. Strains were grown in
1314 Luria Bertani Herskowitz (LBH) (Torkelson et al., 1997) rich medium or M9 minimal medium
1315 (Miller, 1993) supplemented with thiamine (10μg/ml) and 0.1% glucose or glycerol as a carbon
1316 source. Other additives were used at the following concentrations: ampicillin (100μg/ml),
1317 carbenecillin (20μg/ml), chloramphenicol (25μg/ml), kanamycin (30μg/ml), and sodium citrate
1318 (20 mM). P1 transductions were performed as described (Thomason et al., 2007). Genotypes were
1319 verified by antibiotic resistance, polymerase chain reaction (PCR) followed by sequencing, and,
1320 when relevant, UV sensitivity.
1321

44

**Synthesis and generation of *E. coli* mutant and fusion genes**

Mutant or truncated genes were synthesized to introduce site-specific mutations or small deletions in (GenScript) pUC57 backbone plasmids, and subsequently cloned into plasmid pNT3-SD to allow *E. coli* conjugation. Genes that encode wild-type and mutant DNA-binding transcription factors were fused with *mCherry* of (Shee et al., 2013) with a 4-6 alanine linker as described (Heckman and Pease, 2007). The plasmids mentioned above are shown in Table S7.

*E. coli* **mobile plasmid overexpression library**

The mobile-plasmid collection is an ordered library of all 4229 *E. coli* protein-coding genes in a conjugation transferrable plasmid (Saka et al., 2005).  Of these genes, 1017 (or 24%) encode the native *E. coli* protein, whereas 3212 (or 76%) encode the *E. coli* protein with an additional three N-terminal amino acids (Met-Arg-Ala) and an additional two C-terminal amino acids (Gly-Leu), with genes randomly distributed to one or the other kind. We found that native proteins were over-represented significantly as positive for DNA damage in our screens (**RESULTS, A Larger Network Predicted**), implying that the 5 additional amino acids in some of the clones are more likely to confer false-negative results for the proteins that carry them than false-positive results. Table S1 shows the 208 validated DDP-gene clones and indicates the clones that produce native proteins with * next to the clone-ID number, and those with the five additional amino acids with an unmarked clone-ID number.

**Whole-genome primary DDP screen of ordered *E. coli* overexpression library**

The ordered mobile-plasmid collection of 4229 *E. coli* genes in a conjugation transferrable plasmid (Saka et al., 2005) was mobilized into SOS-response-reporter strain SMR17962 (Nehring et al., 2016), to generate a DNA-damage screenable *E. coli* overproduction library.  Protein overproduction is controlled by the IPTG-inducible $P_{tac}$ promoter in cells that fluoresce red when they experience SOS-inducing DNA damage (single-stranded DNA) (Nehring et al., 2016). We adapted a high-throughput 96-well plate reader and robotics to screen for potential DDP-positive strains with increased mCherry fluorescence. Fluorescence intensity per unit of $OD_{600}$ was compared in each well.  Primary screens were performed on cells grown in M9 glucose or M9 glycerol medium (to survey two different conditions), each in duplicate.  Ordered *E. coli* overproduction strains were grown to saturation overnight with shaking at 37°C in clear 96-well plates containing 150µl medium per well, then each well diluted 1:100 into 150µl IPTG-containing medium in 96-well plates (µ-clear, black, Greiner Bio-One, Monroe, NC, USA). The plates were shaken at 37°C for another 24h and analyzed in a Synergy 2 fluorescence plate reader (BioTek, Winooski, VT). We set thresholds of 20% (for glucose with IPTG induction) or 30% (for glycerol with IPTG induction) compared with the median fluorescence intensity per unit of $OD_{600}$ of each individual 96-well plate to identify primary hits. Primary hits were called when two replicates done in the same medium were both above the threshold. Altogether, 414 candidate proteins were identified in this high-throughput plate-reader screen, then tested by flow cytometry for increased endogenous DNA-damage levels to eliminate false positives from the lower resolution/noisier plate-reader assay.

*E. coli* **flow-cytometry secondary screen for increased endogenous DNA damage**

We screened candidate-protein hits from the primary (plate-reader) screen with our more sensitive flow-cytometric assays for SOS-induction/DNA damage (Pennington and Rosenberg, 2007). Each strain identified as positive from the primary plate-reader screen was grown at 37°C to saturation

45

1368  overnight in M9 glucose medium, diluted 1:100 into M9 glycerol medium, and grown for 9 hrs to
1369  early exponential phase at which time IPTG was added to 100μM to induce plasmid-protein
1370  overproduction. After 8 hours of induction, the cultures were diluted 1:100 in filtered M9 glycerol
1371  medium. Samples were analyzed in a LSR Fortessa flow cytometer (BD Biosciences) and
1372  analyzed with BD FACSDivaTM and FlowJo software. For these analyses, $10^5$ events were
1373  collected per strain, per experiment, with each strain assayed in three independent repeats.
1374  Student's *t*-test (*p* value ≤ 0.05) and False Discovery Rate (FDR) q < 0.1 were calculated and
1375  applied based on Benjamini multiple comparison (Benjamini and Hochberg, 1995) to determine
1376  whether overproduction strains had significantly increased levels of endogenous DNA damage.
1377  Two-hundred and eight of the original 414 *E. coli* candidate proteins were validated as genuine
1378  DNA-damage-inducing DDPs when overproduced (shown Table S1).
1379
1380  ### *E. coli* assay for RecA*GFP foci indicating single-stranded DNA
1381  *E. coli* containing the chromosomal *recA4155gfp* allele, encoding RecA*GFP (Renzette et al.,
1382  2005), and the flow-cytometrically validated mobile-plasmid carriers were grown to saturation in
1383  M9 glucose medium at 37°C, then diluted 1:100 into M9 0.1% glycerol and grown for 9h to early
1384  log phase. IPTG was added to 100μM to induce protein overproduction for 8h as described above,
1385  then images taken and analyzed.
1386
1387  ### *E. coli* assays for GamGFP and RDG foci
1388  Saturated cultures of *E. coli* strains (GamGFP: SMR14334; RDG [RuvCDefGFP]: SMR19406)
1389  containing each of the 208 validated DDP-encoding mobile plasmids were grown and induced as
1390  described in flow-cytometric assays for DNA damage. 100ng/ml of doxycycline were added to
1391  induce GamGFP for 3h and RDG for 2h prior to harvesting. Cells were fixed with 1%
1392  paraformaldehyde for 15 min. and washed with PBS buffer three times before being concentrated
1393  for microscopy.
1394
1395  ### *E. coli* microscopy and image analysis for RecA*GFP, GamGFP and RDG foci
1396  Images were acquired using a 100x /NA = 1.4 immersion oil objective (Olympus) on a DeltaVision
1397  Elite deconvolution microscope (Applied Precision, GE). A z-series was acquired sampling every
1398  0.2 microns for a total of 15-25 sections. The z-series was then deconvolved, and a maximum
1399  projection image rendered using Softworx (GE). Image analysis was performed using the
1400  Advanced Imaging collection in Pipeline Pilot 8.5 or 9.2 (Biovia-Dassault Systemes, San Diego).
1401  Projected images from the DeltaVision were read into Pipeline Pilot and metadata data parsed from
1402  the file name and path. A rolling ball background subtraction was applied to improve the signal-
1403  to-background ratio, and to facilitate further segmentation. Individual bacterial cells were then
1404  identified and segmented by applying a global threshold on images of the fluorescently labeled
1405  protein. Morphological manipulations (smoothing, opening and closing) were applied to refine the
1406  segmentation edges and a watershed was then performed to separate neighboring objects. Filtering
1407  was then applied to remove bacteria that fell outside a certain area threshold and that did not
1408  contain DNA. Foci were then identified using a more aggressive per-object background subtraction
1409  and peak identification method. Objects tentatively identified by this method were subsequently
1410  filtered by circularity, signal-to-background ratio, and size. Focus-positive bacteria were then
1411  determined using the co-localized objects component in the Advanced-imaging library in Pipeline
1412  Pilot. A binary metric, whether the cells were focus-positive or not, was calculated in addition to
1413  recording the total area and count of foci for those bacteria that were positive.

1414

1415 **STRING/network analyses**
1416 Known protein-protein interactions were displayed using CytoScape V3.4.0 software. Protein-
1417 protein interaction linkage scores were taken from the STRING 10.0 database (Szklarczyk et al.,
1418 2015) to identity interaction pairs. We used STRING, all parameters, with an interaction score cut-
1419 off of ≥0.6 (medium-to-high confidence). Random controls were produced by examining equal-
1420 size groups of random *E. coli* genes. *P* values were calculated with a hypergeometric test
1421 (Berkopec, 2007). The *E. coli* DDP network has network properties that are defined as scale-free
1422 and "small-world", and it has significantly more edges (connectivity) compared with a random
1423 network (Figure 1G, S2A, Results). The human candidate-DDP network was generated similarly,
1424 and also has more connectivity than a random human-gene network or random human genes with
1425 *E. coli* homologs (Figures 2B, S2B, Results).

1426

1427 *E. coli* **forward-mutation assay**
1428 We used the forward-mutation assay of Matic and colleagues (Gutierrez et al., 2013) in which *E.*
1429 *coli* wild-type strain MG1655 harbors a chromosomal phage lambda *c*I transcriptional repressor
1430 gene, and a CI-repressible *tetA* gene, such that mutations that inactivate *c*I are scored as
1431 tetracycline-resistant (Tet$^R$) mutant cfu. Into this strain, we conjugated 32 validated *E. coli* DDP
1432 genes in their mobile-plasmid-library vector (genes tested Table S1; Supplemental Discussion 3,
1433 and Table S1 for their mobile-library clone names). We developed a modified higher-throughput
1434 fluctuation-test assay for determining numbers of cultures with Tet$^R$ mutants from which to
1435 calculate Tet$^R$ mutation rates. Each DDP overproducer was grown overnight to saturation in M9
1436 glucose with 20μg/ml carbenicillin at 37°C shaking, then diluted 1:10,000 into M9 glycerol
1437 carbenicillin and each culture split into 24 or 32 wells in 96-well-plates at 100μl per well. The
1438 plates were shaken at 37°C for 15h (early log phase), and IPTG added to attain 100μM in each
1439 well to induce protein overproduction for 8h, as described in flow-cytometric validation. From the
1440 end cultures, 5-10 μl were moved into LBH medium containing 10μg/ml tetracycline to determine
1441 the fraction of cultures that contained no Tet$^R$ cells after incubation and scoring of the wells in the
1442 plate reader for OD (Tet$^R$ cells) versus failure to grow (no Tet$^R$ cells). The viable cell counts were
1443 estimated by sampling three wells chosen randomly. The $P_0$ method was used to estimate mutation
1444 rates for each genotype as described with correction for the fraction sampled (Foster, 2006). The
1445 data reported (Figure 1J; Table S1) are the mean mutation rates (± SEM) of three experiments of
1446 at least 24 cultures per strain for each of the 32 strains assayed.

1447

1448 *E. coli* **Tet$^R$ mutation verification by sequencing**
1449 We selected strains that overproduce the following 10 different DDPs with strong DNA-damage-
1450 up phenotypes: CsgD, TopB, CheA, YegL, MdtA, GrpE, HslU, YicR, UvrA, and Mrr. We selected
1451 3-10 independent Tet$^R$ mutant colonies, each from a separate culture from each strain, from which
1452 to sequence *c*I mutations. For the vector-only negative control, 19 independent Tet$^R$ colonies were
1453 isolated. We amplified and sequenced a 1122nt region encompassing the *c*I gene as described
1454 (Gutierrez et al., 2013) to identify the mutations. For those Tet$^R$ mutants that failed to yield PCR
1455 products, implying deletion of the *c*I gene, further outside primers (forward:
1456 ACCGCGGCGTGGGTAGTAAAGT, and reverse: GCCAATCCCCATGGCATCGAGTAAC)
1457 were used for PCR, and the products sequenced. In two cases (both TopB overproducers), whole-
1458 genome sequencing (WGS) was performed to determine the end-points for deletions that could not
1459 be determined via PCR and sequencing.

47

1460

1461 ### *E. coli* whole-genome sequencing and analysis

1462 Tet$^R$ mutants were grown at 37°C to saturation overnight in LBH with 10μg/ml tetracycline, and
1463 genomic DNA was extracted and purified using DNeasy Blood & Tissue kits (Qiagen). Libraries
1464 were prepared using Nextera XT kits (Illumina); sequencing was performed on an Illumina Mi-
1465 Seq, and sequencing data analyzed as described (Xia et al., 2016). Sequencing reads were mapped
1466 to the MG1655 genome (NCBI RefSeq Accession: NC_000913.3). Low-quality reads and
1467 duplicates were removed. WGS files were visualized and deletion endpoints were analyzed using
1468 IGV software (Broad Institute, MA).

1469

1470 ### Flow-cytometric assays for DNA loss

1471 Quantification of anucleate cells by flow cytometry was adapted from (Joshi et al., 2013).
1472 Saturated cultures of *E. coli* strains derived from SMR21384 containing each of the 208 validated
1473 DDP-producing mobile plasmids were grown and induced as described in flow-cytometric assays
1474 for DNA damage. Cells were resuspended in 100 μl PBS, and stained with membrane dye
1475 FM® 4-64FX (Thermal Fisher) with a final concentration of 10 μg/ml. The mix was kept on
1476 ice for 10 min. and then washed three times with PBS. A final concentration of 70% ethanol
1477 (pre-chilled) was used to fix the cells at -20°C for 1h, after which cells were washed with
1478 twice with PBS and resuspended in 100 μl PBS. 100 μl DAPI (5μg/μl) were used to stain
1479 DNA at room temperature (RT) for another 10 min. Samples were filtered and analyzed as
1480 described above.

1481

1482 ### Flow-cytometric assay for intracellular ROS levels

1483 Saturated cultures of *E. coli* strains derived from SMR21384 containing each of the 208 validated
1484 DDP-producing mobile plasmids were grown and induced as described in flow-cytometric assays
1485 for DNA damage. The ROS measurement protocol was modified from Gutierrez et al.
1486 (Gutierrez et al., 2013). In brief, cells were incubated with ROS-staining dye DHR123
1487 (Invitrogen), which measures $H_2O_2$, for 30 min. at 4°C in M9 buffer. After washing twice with
1488 M9 buffer, flow cytometry analyses were performed immediately as described above.

1489

1490 ### Flow-cytometric assay for intracellular pH

1491 pHrodo® Green AM Intracellular pH Indicator (Thermal Fisher) was used to measure
1492 intracellular pH in live *E. coli*. Protocols were adapted from (Loiselle and Casey, 2010). Cells
1493 were first washed with live-cell imaging solution (LCIS) and then 10 μl of pHrodo™ Green AM
1494 with 100 μl of PowerLoad™ concentrate were added to 10 ml of LCIS. The pHrodo™
1495 AM/PowerLoad™/LCIS was mixed with cells and incubated at 37°C for 30 minutes. Cells were
1496 then washed twice with PBS to remove excess dye before flow-cytometric analysis. Intracellular
1497 pH calibration buffers (Thermal Fisher) were used as standards.

1498

1499 ### Assays for sensitivity to DNA-damaging agents

1500 Cultures of *E. coli* strain (SMR21384) containing each of the 208 validated DDP-producing mobile
1501 plasmids were grown as described in flow-cytometric assays for DNA damage with the following
1502 modifications: For hydrogen-peroxide ($H_2O_2$) treatment, 100 μM IPTG was used to induce
1503 overproduction of each DDP. Each culture was split into two tubes, prior to addition of 5mM $H_2O_2$
1504 into one of the tubes for 15min. The cells with and without $H_2O_2$ were immediately diluted and
1505 plated onto LBH plates for assay of viable cells as cfu after incubation for a day at 37°C. For
1506 phleomycin or mitomycin C (MMC) treatment, saturated M9 glucose cultures were diluted into

1507  M9 glycerol medium with 100 µM IPTG to induce overproduction in 96-well plates. The plates
1508  were grown with shaking for 8 hours at 37°C to early log phase prior to addition of 1 µg/ml
1509  phleomycin or 0.05µg/ml MMC to each well. After 20 hours of continuous shaking, the OD600
1510  was read using a BioTek microplate reader Synergy 2 (BioTek). DNA-damaging-agent
1511  sensitivities of the DDP-producing clones are normalized to sensitivity of vector-only controls:
1512  (treated/untreated DDP overproducer) / (treated/untreated vector-only) so that values < 1 indicate
1513  sensitivity. For all three assays for sensitivity to DNA-damaging agents, Student's *t*-test (*p* value
1514  ≤ 0.05) with FDR adjustment (q ≤ 0.1) was used to determine whether DDP-overproducing strains
1515  were significantly more sensitive to DNA-damaging agents than the vector-only control.
1516
1517  **Clustering methods**
1518  For each DDP and DNA-damage outcome measure, raw data for each functional assay
1519  (overproduction versus vector) were converted into z scores and were used to delineate groupings
1520  of proteins with similar properties and patterns of response. Unsupervised discovery methods K-
1521  means in combination with Progeny Clustering (Hu et al., 2015) were performed using the R
1522  package *ProgenyClust* (Hu and Qutub, 2016) to determine the optimal number of protein clusters
1523  for the 208 DDPs. Seven functional tests were clustered by hierarchical clustering to assess the
1524  association of kinds, causes, and consequences of DNA damage.
1525
1526  **RNA-seq library preparation and sequencing**
1527  *E. coli* cultures were grown as described for flow-cytometric assays for DNA damage, and RNA
1528  was isolated from 1 ml of culture (~$10^8$ cells) for each of two biological replicates. Total RNA was
1529  isolated using the RNeasy Mini Kit (Qiagen), according to the manufacturer's protocol.
1530  RNAprotect Bacterial Reagent (Qiagen) was used to stabilize RNA during harvest and enzymatic
1531  cell lysis. After elution, total RNA was treated with RNase-free DNase I (NEB), according to the
1532  manufacturer's protocol. RNA was recovered by phenol-chloroform extraction and ethanol
1533  precipitation. Ribosomal RNA was depleted using RiboZero (Epicentre/Illumina), according to
1534  the manufacturer's protocol. Remaining RNA was concentrated by ethanol precipitation and
1535  approximately 100 ng of rRNA-depleted RNA was used to construct libraries using the TruSeq
1536  Stranded mRNA Library Preparation Kit (Illumina). Libraries were prepared according to the
1537  manufacturer's protocol, using recommended modifications for previously isolated mRNA
1538  (McClure et al., 2013) (poly-A RNA enrichment steps excluded). Final RNA-seq libraries were
1539  run on a BioAnalyzer (Agilent) to estimate the average fragment size (~800 bp) and the
1540  concentration of adapter-ligated library fragments was determined using the qPCR-based Illumina
1541  Library Quantification Kit (KAPA Biosystems). Libraries were pooled and sequenced on an
1542  Illumina NextSeq 500 using a High Output v2 Kit (2 x 75 bp paired-end reads).
1543
1544  **Analysis and deposition of RNA-seq data**
1545  Read mapping, transcript assembly, and differential expression analysis were performed using
1546  Rockhopper (McClure et al., 2013), a bacteria-specific RNA-seq analysis pipeline, using MG1655
1547  (NC_000913.3) as the reference genome. Genes were considered as differentially expressed if the
1548  fold change was greater than or equal to 2 and q-value was less than 0.01. Sequencing data are
1549  available in the European Nucleotide Archive (ENA) under study accession no. PRJEB21034.
1550
1551  **RDG ChIP-seq library preparation, sequencing, and data analysis**
1552  Cells were grown as for focus quantification, then crosslinked, lysed and sonicated as described

1553    (Xia et al., 2016). Immunoprecipitation and library preparation methods are based on those of
1554    (Bonocora and Wade, 2015) with small modifications as follows. RuvC antibody (Santa Cruz) was
1555    first pre-incubated with Dynabead protein A, then the RuvC-antibody-coated Dynabeads were
1556    incubated with cell lysates at 4°C overnight. Library preparation was performed while DNA
1557    fragments were still on Dynabeads. Samples were barcoded using NEBNext Multiplex Oligos for
1558    Illumina. Size selection of adaptor ligated DNA was performed on AMPure XP Beads as described
1559    in NEBNext ChIP-Seq Library Prep guidelines. Because the concentrations of eluted ChIP DNA
1560    are low, samples were amplified briefly prior to size selection, and a second amplification was
1561    performed after size selection. Sequencing was performed on an Illumina MiSeq. The pipeline for
1562    data analysis consists of the following steps: (i) reads were trimmed by Trimmomatic (Bolger et
1563    al., 2014) removing sequencing adaptors and low quality bases; (ii) reads were aligned by BWA-
1564    MEM (Li, 2013) to the W3110 genome [National Center for Biotechnology Information (NCBI)
1565    Reference Sequence (RefSeq) Database accession: NC_007779.1] and the plasmid pNT3 (Saka et
1566    al., 2005); (iii) Secondary alignment and multiple-mapped reads were discarded, this results in
1567    zero coverage in repetitive regions and regions present in both the genome and the plasmid,
1568    including the *csgD* gene; (iv) potential PCR duplicates were removed by Picard Tools
1569    MarkDuplicates; (v) bedGraph files were generated with deepTools (Ramirez et al., 2014) and
1570    imported to R for plotting; and (vi) peak calling was performed with MOSAiCS (Sun et al., 2013).
1571    Sequencing data are available in the European Nucleotide Archive (ENA) under study accession
1572    no. PRJEB21035.
1573
1574    **Western analyses of Pol IV protein levels**
1575    M9 glycerol cultures inducing wild-type, catalytically inactive, and β-binding-defective Pol IV
1576    were normalized to $OD_{600}$ of 1.0, and 1 ml of each was pelleted, resuspended and boiled as
1577    described (Kim et al., 2001). Proteins were separated by 10% SDS-PAGE and transferred to PVDF
1578    membrane according to the manufacturer's instructions (Amersham, GE Healthcare). The
1579    membranes were blocked with ECL Prime blocking agent (GE Healthcare) and probed with
1580    primary anti-Pol IV polyclonal antibody (Kim et al., 2001) (1:2000). The membrane was further
1581    probed with secondary polyclonal goat anti-rabbit IgG-Cy5 antibody (Bethyl Laboratories) and
1582    visualized by scanning in multicolor imager Typhoon detection system (GE Healthcare).
1583
1584    **Identification of human homologs using BLASTp and delta BLAST**
1585    "Homologs" are defined here as proteins with amino-acid similarity that could result from possible
1586    evolutionary relatedness. We used two basic local alignment search tools: the BLASTp and Delta-
1587    BLAST algorithms, searching protein sequences obtained from GenBank and other NCBI database
1588    resources. For both we used e-value < 0.01 (≤1 gene is identified by random chance in 100 queries)
1589    and sequence identity of ≥20%. Note that ≥20% sequence identity between *E. coli* and human is
1590    considerable. For example, known orthologs *E. coli* RecA and human RAD51 have 25% amino-
1591    acid identity. Given a protein query, BLASTp returns the most similar protein sequences from the
1592    protein database with e-value < 0.01 and identity ≥ 20%. Delta-BLAST uses multiple sequence
1593    alignment with conserved domains found in the CDD (Conserved Domains database from NCBI)
1594    and computes a Position Specific Score Matrix (PSSM) (Boratyn et al., 2012) with e-value < 0.01.
1595    Both methods were compared against the human protein database of NCBI. Proteins identified
1596    from either algorithm were identified as human homologs of the *E. coli* DDPs.
1597
1598    **Analyses of cancer survival and mutation loads**

1599 RNA-sequencing data from The Cancer Genome Atlas (TCGA) (Gao et al., 2013) were processed
1600 in the form of transcripts per million (TPM) as described (Rahman et al., 2015) and obtained via
1601 Gene Expression Omnibus (accession number GSE62944). Only the TCGA cancer types that had
1602 over 100 patients with RNA- and DNA-sequencing data were analyzed. Upon defining our gene
1603 sets of interest, RNA data were subjected to single sample Gene-Set Enrichment Analysis
1604 (ssGSEA) using GSVA package (Hanzelmann et al., 2013) in R. The resulting gene-set enrichment
1605 score for each sample was used as a representation of gene-set RNA level in each sample. Somatic
1606 mutation data for TCGA cancers were obtained in the form of mutation annotation files (raw or
1607 final) from the Broad Institute Genome Data Analysis Center (GDAC). For each sample, the sum
1608 of base-substitution and indel mutations was taken as the total mutation count, and $\log_e$ of this
1609 value was referred to as "mutation load." Correlation analysis for "all human genes" was
1610 performed via bootstrapping. Briefly, we computed the mean correlation coefficient of mutation
1611 load with gene-set enrichment scores for 1000 randomly sampled gene sets, each consisting of a
1612 random number, between 10 to 1000, of genes out of over all human genes for which expression
1613 data were available. Kaplan Meier survival analysis was performed using "survival" package in R
1614 comparing the top and bottom tertiles of samples based on their gene-set enrichment score.
1615 Correlation analyses with mutation loads was performed in base R and correlation coefficients
1616 were plotted using the "corrplot" package in R.
1617
1618 **Cloning of human genes for DNA-damage analyses in human cells**
1619 Fifty-eight human DDP and 19 non-DDP cDNA clones (Table S5) in the Gateway entry vectors
1620 pDONR221 and pDONR223 (Invitrogen) were subcloned from an augmented library of ~32,000
1621 Orfeome V8.1 (Yang et al., 2011) stated to contain sequenced human full-length cDNA clones,
1622 and additional full-length and commonest splice-variant length clones obtained from others
1623 including CCsBroad gene libraries. The size of cDNA from each gene was confirmed by restriction
1624 enzyme digestions. We also cloned, de novo, 15 candidate hDDP genes, one non-hDDP gene,
1625 tubulin and two *de novo* methylase genes (Table S5) that were not present as full-length clones in
1626 the Orfeome V8.1 (Yang et al., 2011) or CCsBroad gene libraries. These candidate genes were
1627 amplified from cDNAs generated from mRNAs extracted from the human cancer-cell lines U2OS
1628 or MRC5-SV40. PCR products of the correct size were cloned into the Gateway entry vector
1629 pENTR11 at restriction enzyme cut sites or into pDONR201 using *attB* site-specific recombination
1630 sites. Five DNMT1 truncated constructs were modified by using site-directed mutagenesis (Table
1631 S5). Clones were sequenced and verified as the correct gene sequence based on the Reference
1632 Sequence (RefSeq) database from NCBI. We subcloned each gene into a mammalian expression
1633 vector containing a GFP epitope tag (pcDNA6.2/N-EmGFP-DEST, Invitrogen), which allows us
1634 to analyze transfection efficiency and visualize protein localization in transfected cells. All human-
1635 cell overexpression plasmids used in this study are listed in Table S5.
1636
1637 **Human cell lines, plasmids, and reagents**
1638 MRC5-SV40 and HEK293T cells were maintained in Dulbecco's modified Eagle's medium
1639 (DMEM) (Invitrogen) supplemented with 10% fetal bovine serum (FBS), 2 mM L-glutamine, 100
1640 μg/mL penicillin, 100 μg/mL and streptomycin. Transient transfections into human cells were
1641 performed using GenJet (SignaGen Laboratories) for MRC5-SV40 and PEI (polyethylenimine,
1642 Sigma) for HEK293T. Transfections for siRNA were carried out with lipofectamine RNAiMax
1643 (Invitrogen) following the manufacturer's instructions. The siRNAs were siNT: non-targeting pool
1644 (Dharmacon) and siRAD18: ACUCAGUGUCCAACUUGCU (Sigma). DNA-PK inhibitor

1645 (NU7441, Tocris Bioscience) was used at 2.5 μM 6 h prior to harvesting cells for flow cytometry.
1646 NAC (N-acetyl-cysteine, Sigma) treatment was performed twice, with a final concentration of 5
1647 mM, post-24hr and -48hr transfection. To create inducible stable clones to verify DNMT1 and
1648 PCNA interaction, GFP-tubulin, GFP-DNMT1 and GFP-DNMT1-ΔPBD cDNAs were cloned into
1649 pcDNA5/FRT/TO/Intron vector (Invitrogen, CA). Inducible HEK293T FlpIn Trex GFP-tubulin,
1650 GFP-DNMT1 and GFP-DNMT1-ΔPBD cells were generated followed by manufacturer's protocol
1651 and were cultured in the same normal medium with 15μg/ml Blasticidin and 80μg/ml hygromycin.
1652 Doxycycline (Sigma) was added to medium to trigger the production of GFP fusions.
1653
1654 **Human-cell DNA-damage screens by flow cytometry**
1655 We screened for increased DNA damage by flow-cytometric quantification of γH2AX- and
1656 phospho-P53-antibody signals among GFP-positive transfectants. Immunostaining was performed
1657 according to a standard procedure with minor modifications. Seventy-two hours post-transfection,
1658 cells were collected and approximately $1 \times 10^6$ cells taken for staining. For staining, cells were fixed
1659 with 2% (v/v) formalin for 15 min on ice, washed twice in cold-PBS and permeabilized with 0.05%
1660 (v/v) Triton-X for 15 min on ice followed by two washes with PBS.  The fixed cells were then
1661 blocked with 5% BSA-PBS for 1 hr, and stained with either γH2AX (Millipore) or phosphorylated
1662 p53 primary antibodies (Cell Signaling) overnight at 4°C. Cells were washed three times in 1%
1663 BSA-PBS followed by an incubation of Alexa Fluor 647 goat anti-mouse IgG in 5% BSA-PBS
1664 (Invitrogen) for 1 hr at room temperature in the dark, then washed three times with 1% BSA-PBS.
1665 Stained samples were measured by a BD LSRFortessa flow cytometer and analyzed using FlowJo
1666 software. Cells without transfection were used to set the threshold gating to determine the
1667 percentage of GFP- and γH2AX- or phosphorylated p53-positive cells, with 0.5% of control cells
1668 gated as the damage threshold. The DNA-damage ratio caused by protein overproduction is
1669 defined by (Q2/Q3)/(Q1/Q4), where Q2 is the number of transfected damage-positive cells; Q3 is
1670 the number of transfected damage-negative cells; Q1 is the number of untransfected damage-
1671 positive cells; and Q4 is the number of untransfected damage-negative cells. Results were obtained
1672 from at least two independent experiments. Statistical significance (*p* value) was determined using
1673 two-tailed unpaired Student's *t*-test followed by false discovery rate (q value) correction.  Both the
1674 γH2AX and phosphorylated-p53 assays show linear responses to exogenous DNA damage caused
1675 by ionizing radiation (Figures S3J and S3K), indicating their quantitative validity.
1676
1677 **HPRT mutagenesis assay**
1678 MRC5-SV40 cells were transfected with the plasmids indicated, and harvested 72 hours post-
1679 transfection. The percentage of GFP-positive cells of each transfectant was scored as transfection
1680 efficiency using a BD Accuri flow cytometer. The remaining cells were re-grown in 15 cm dishes
1681 for an additional 4 days. After a week of transfection, $3 \times 10^6$ cells were plated in 15 cm dishes
1682 containing medium with 20 mM 6-thioguanine (Sigma), with five 15 cm dishes for each gene. In
1683 addition, 600 cells were plated in triplicate, per well, in a 6-well plate without 6-TG to determine
1684 plating efficiency. The plates were incubated at 37°C in a humidified incubator until colonies
1685 formed. The colonies were stained with 0.005% crystal violet. These colonies were counted, and
1686 mutation rates determined using the MSS-maximum likelihood estimator method with correction
1687 for transfection efficiency. We verified that 6-TG resistant clones result from *HPRT* mutations by
1688 sequencing the cloned *HPRT* cDNAs from four independent mutants (Supplemental Discussion
1689 8). The *HPRT* cDNA is 657bp long, whereas *HPRT* including introns is 42kb, making sequencing
1690 the cDNAs more practical.

52

1691
1692 **Human-cell immunoprecipitation and western blot analysis**
1693 After induction of protein production using doxycycline in FlpIn-inducible HEK293T cells
1694 producing GFP-tubulin, GFP-DNMT1-WT or GFP-DNMT1-ΔPBD, cells were lysed with NETN
1695 buffer (150 mM NaCl, 1mM EDTA, 10 mM Tris-HCl, pH 8.0, and 0.5% NP-40) containing
1696 TurboNuclease (Accelagen) and 1 mM $MgCl_2$ for 1 h at 4°C. Cell lysates were then centrifuged
1697 for 30 min at 4°C. GFP-tagged proteins were immunoprecipitated with 20 μl of GFP-Trap_A
1698 (Chromotek) for 1 h at 4°C. Beads were then washed three times with NETN buffer. Protein
1699 mixtures were eluted by boiling at 95°C with Laemmli buffer (4% (v/v) SDS, 20% (v/v) glycerol
1700 and 120 mM Tris-HCl, pH 6.8). For whole cell extracts, cells were collected with Laemmli buffer,
1701 and heated for 5 min at 95°C before loading. Samples were resolved by SDS-PAGE followed by
1702 western blot analysis. Primary antibodies were used as follows: anti-GFP (Invitrogen), anti-PCNA
1703 (Santa Cruz), anti-beta tubulin (Abcam), anti-RAD18 (Cell Signaling). Blots were analyzed by
1704 standard chemiluminescence (GE Healthcare, Amersham ECL Prime system) using a Bio-Rad
1705 molecular imager ChemiDoc XRS+ system.
1706
1707 **Statistics**
1708 All *E. coli* wet-bench experiments were performed at least three times independently, and a two-
1709 tailed unpaired *t*-test was used to determine significant differences, unless otherwise specified.
1710 Error bars represent 1 SEM except where otherwise indicated. Pearson's correlation coefficient
1711 was computed to assess the relationship between two parameters. STRING enrichment analysis
1712 was performed using hypergeometric tests with the correction for multiple comparisons. False
1713 discovery rate (FDR) adjustments are used to limit the overall type I errors in both *E. coli* and
1714 human DNA-damage flow-cytometry assays. The FDR (Benjamini Hochberg) method (Benjamini
1715 and Hochberg, 1995) is the default *p*-value adjustment method in this paper. Fisher exact test is
1716 used to determine whether two proportions are different. Wilcoxon rank-sum test was used to
1717 determine whether each gene has cancer-associated copy-number increases.
1718
1719 **QUANTIFICATION AND STATISTICAL ANALYSIS**
1720 Statistical details can be found in the main text, figure legends, or in the Method Details section.

53

## SUPPLEMENTAL INFORMATION

Supplemental information includes a discussion file, seven figures and seven tables that can be found with article online at ***

### Supplemental Discussion 1
### Significant protein-protein interactions of random human homologs of *E. coli* proteins

Hypergeometric test analyses show that the 284 human homologs of *E. coli* DDPs have far more significant association ($p = 1.2$ x $10^{-327}$) than 284 random human proteins ($p = 0.80$), and 284 random human homologs of *E. coli* proteins ($p = 1.8$ x $10^{-49}$). Although not associated nearly as strongly as the human homologs of DDPs ($1.2$ x $10^{-327}$), the significant association of random human homologs of *E. coli* proteins ($1.8$ x $10^{-49}$) compared with random human proteins ($p = 0.80$) could potentially result if highly conserved proteins generally have more interactions with each other than random proteins. This might be because the most highly conserved, fundamental aspects of biology, and proteins that participate in them, are enriched for conserved protein machines (ribosomes, replisomes, transcription complexes, etc.), and/or fundamental pathways the actors in which have remained associated. Alternatively, it might be that proteins that function as part of interacting protein groups evolve more slowly, and so are overrepresented among conserved proteins.

### Supplemental Discussion 2
### A larger network predicted and estimate of additional *E. coli* DDPs not discoverable in the mobile plasmid library

The 208 proteins are a large network, and occupy 5% of *E. coli* genes, but are likely to represent just over half of overproduction DDPs encoded in the *E. coli* genome. Per Figure S1E, we found that 1 of 99 random proteins not identified in the primary screen was positive in the sensitive flow-cytometry secondary assay, predicting an additional undiscovered 38 DDPs in the overproduction library used (Figure S1E). Further, although it is the most complete and least adulterated *E. coli* overexpression library, the mobile plasmid library (Saka et al., 2005) contains some genes that encode five additional amino acids, which our data indicate were biased against in our screens. Twenty-four percent of clones in the library (STAR Methods) produce native *E. coli* proteins, and the rest produce proteins with three extra N-terminal (Met-Arg-Ala) and two extra C-terminal (Gly-Leu) amino acids, with the composition of genes in each class being random (Saka et al., 2005) (STAR Methods). We found that both the initial DDP candidates identified in the plate-reader primary screen and the 208 flow-cytometry-validated DDPs carried significantly higher fractions of native proteins than the library; there were 158 native proteins in the initial 414 candidates identified in the primary plate-reader screen (38%, differs from the library at $p = 1.7$ x $10^{-11}$, Fisher's exact test), and 85 native proteins in the 208 validated DDPs, or 41% (shown in Table S1, differs from the library at $p = 4.1$ x $10^{-8}$, Fisher's exact test). The data imply that some of the non-native proteins may have lost full function, and, because of that, gave false-negative readings in the screens. We found that the native genes in the library were "hit" in the primary screen at 16% (158 discovered out of 1015 native genes in the library), whereas the non-native genes were identified at 8% efficiency (256 discovered out of 3214 non-native genes in the library). If there are an additional 7.6% of the non-native proteins that would score as DNA-damage-promoting in our primary screen, if they did not carry the extra amino acids, then among the 3214 non-native-protein-encoding genes in the library, we predict that there would be an additional 244 overproduction DDP candidates found in the primary screen (7.6% of 3214). We found that

1767 candidates from the primary screen were validated in the secondary screen at 208 validated out of
1768 414 candidates (Figure 1F; Tables S1 and S2), or just over 50%, which predicts 123 additional
1769 genuine DDPs among the predicted additional candidates.
1770

1771 **Supplemental Discussion 3**
1772 *E. coli* **clones assayed for mutation rate**
1773 In Figure 1J, we assayed mutation rates in mutation-assay strains overproducing the following
1774 DDPs. DDPs that cause < 5-fold increase in DNA damage: DsbG, YijF, CadA, FolD, YddG,
1775 LeuO, UvrB, YajR, YbgQ, ORF 6106.1. DDPs that cause ≥ 5-fold increase in DNA damage:
1776 HypF, ZipA, YedA, CueO, YefU, MacB, HcaR, MdtB, SetB, DinD, RusA, YdcR, CsgD, HslU,
1777 SfsA, TopB, CorA, YegI, GrpE, PgrR, Mrr, MhpR. Non-DDPs: AceF, HprT, AceE, YaeG, YadF,
1778 PdhR, HrpB, MrcB, FhuD, YadG, Dgt, FhuA, HtrE, EcpD, FhuC, YacH, YadK.
1779

1780 **Supplemental Discussion 4**
1781 **Cancer association of human proteins from a select DNA-damage screen**
1782 We analyzed published data from the limited human overexpression DNA-damage-up screen of
1783 (Lovejoy et al., 2009) by Fisher exact test against known (Forbes et al., 2015) and predicted
1784 (D'Antonio and Ciccarelli, 2013) cancer-driving genes. This overexpression screen of a set of
1785 nucleus/DNA-associated proteins discovered 96 human proteins (Lovejoy et al., 2009), which we
1786 found are overrepresented among known and predicted cancer drivers at $p = 0.0001$ and $p = 0.0002$,
1787 with DNA-repair proteins excluded (Fisher exact test, identities, Table S3). Only one protein was
1788 identified in common between the *E. coli* DDP homologs and the human overproduction screen
1789 (FIGNL1), indicating that the *E. coli* screen identified many new hDDP candidates, then validated
1790 hDDPs. Overall, the candidate hDDPs identified from human screens and the *E. coli* screen are
1791 highly significantly overrepresented among known (Forbes et al., 2015) and predicted (D'Antonio
1792 and Ciccarelli, 2013) cancer driver genes, independently of DNA-repair proteins, supporting the
1793 importance of DDPs to human cancer. We note that an unbiased screen of all human proteins for
1794 DNA damage on overproduction is not possible because the best human overexpression libraries
1795 contain a fraction of all human protein-coding genes, and many clones that are not full length. See
1796 STAR Methods, **Cloning of human genes for DNA-damage analyses in human cells.**
1797

1798 **Supplemental Discussion 5**
1799 **Choice of candidate hDDP and control proteins for validation in DNA-damage assays**
1800 Of the 284 human homologs, we identified 121 candidates of particular interest according to the
1801 following criteria: (i) Many are encoded by genes amplified at high frequencies in cancer genomes
1802 from TCGA (Gao et al., 2013) (Table S4). (ii) For a minority, the genes are mutated or deleted at
1803 impressive, high frequencies in TCGA (Gao et al., 2013). (iii) Full-length clones that encode 90
1804 of these appeared to be available in the Orfeome V8.1 or CCsBroad cDNA-clone collections (Yang
1805 et al., 2011). We determined by restriction mapping that many of the human genes in those libraries
1806 are not full length (Table S5), and cloned 18 genes including 15 candidate hDDP genes and three
1807 controls de novo as full-length cDNA clones that we sequence-verified (STAR Methods). We
1808 ultimately created 70 full-length overexpression GFP-fusion clones of human homologs of *E. coli*
1809 DDP genes, and overexpression GFP-fusion clones of 3 human homologs of *E. coli* damage-down
1810 genes, as possible negative controls (STAR Methods, Tables S5), 9 random human genes, and 11
1811 random human homologs of *E. coli* non-DDP genes (Tables S5).
1812

1813 **Supplemental Discussion 6**
1814 **Superiority of transient transfection to stable integration of genes encoding hDDP candidates**
1815 We found transient transfection to be superior to creation of stable clones because of apparent
1816 selection for mutations in the inducible hDDP candidate genes upon integration. Mutations in the
1817 hDDP candidates or other DNA damage-response pathways are selected probably because the
1818 gene products are toxic when overproduced and the genes are difficult to keep tightly "off". The
1819 GFP-hDDP-gene fusions allow transient transfection assays to identify immediate effects of the
1820 DNA damage and to analyze only the minority population of cells that have been transfected
1821 successfully and produce the protein of interest. This cell subpopulation is GFP-positive, and easily
1822 identified in the flow-cytometric assays (e.g., Figure 3B).
1823
1824 **Supplemental Discussion 7**
1825 **Estimation of additional human DDPs demonstrable in assays used here**
1826 We evaluated the validation efficiencies of four classes of human homologs of *E. coli* DDP genes
1827 (shown Figure 3D):  genes that are—(i) both known (Forbes et al., 2015) or predicted (D'Antonio
1828 and Ciccarelli, 2013) cancer drivers and amplified in TCGA cancers; (ii) amplified in cancers and
1829 not known or predicted drivers; (iii) known/predicted cancer drivers that are not known to be
1830 amplified in cancers; and (iv) neither amplified in cancers nor previously known/predicted cancer
1831 drivers. Based on the number of candidates that we tested in each class among the 70 DDP
1832 homologs tested, these data correspond to the following validation rates as DNA-damage-
1833 promoting for each class: (i) 100%; (ii) 53%; (iii) 67%; and (iv) 27%. Based on the numbers of
1834 homologs not yet tested in each of these classes, our data predict that the following numbers of
1835 proteins among the remaining (284 - 70=214) human-homolog candidate hDDPs would be likely
1836 to be validated in these particular DNA-damage assays: (i) 6; (ii) 38; (iii) 34; and (iv) 7, for a total
1837 of at 85 more demonstrable hDDPs predicted among the 284-protein candidate hDDP network.
1838 We note, however, that the human-cell DNA-damage assays used favor detection of DNA double-
1839 strand breaks, not all DNA-damage types comprehensively. Thus, many more of the human
1840 homologs may be DNA-damage promoting for other kinds of DNA damage than is estimated here.
1841
1842 **Supplemental Discussion 8**
1843 **Verification of *HPRT* Mutations in 6-Thioguanine Resistant Human-Cell Clones**
1844 Four independent 6-thioguanine-resistant clones were shown to result from *HPRT* mutations by
1845 sequencing the cloned *HPRT* cDNAs. The mutations are: a single-basepair insertion between the
1846 206-207nt of *HPRT* gene, and three identical deletions (from 403nt to 485nt). Two of the
1847 sequenced clones were independent DNMT1-overproducing transfectants, and two were from
1848 independent vector-only control transfected cells. The *HPRT* cDNA is 657bp long, whereas
1849 *HPRT* including introns is 42kb, making sequencing the cDNAs more practical.
1850
1851 **Supplemental Discussion 9**
1852 **Controls**
1853 While analyzing RNA-Seq data, we identified a 2177bp deletion including the *lacI*$^q$ region on the
1854 pNT3 empty vector. We have determined that this deletion does not alter results in any of our
1855 assays or any of our conclusions. The phenotypes of the truncated empty vector were compared
1856 with 10 non-DDP overproducers, and then with the full-length empty vector, in all 7 functional
1857 assays, and, by one-way ANOVA analysis, there were no significant differences between the
1858 means of all 11 strains in any of the 7 assays ($p$=0.19 GamGFP foci; $p$=0.28 RDG (reversed-fork)

1859 foci; $p=0.99$ ROS; $p=0.99$ anucleate cells/DNA loss; $p=0.26$ phleomycin sensitivity; $p=0.08$ H$_2$O$_2$
1860 sensitivity; $p=0.21$ mitomycin C sensitivity), and no difference between it and the full-length
1861 vector ($p=0.31$; $p=0.44$; $p=0.62$; $p=0.32$; $p=0.62$; $p=0.28$; and $p=0.78$, respectively, two-tailed
1862 unpaired $t$-test).

1863

1864 **Supplemental Discussion 10**
1865 **DNA-damage sensitivity not from mutations or *E. coli* DDP overproduction**
1866 We show that DNA-damage sensitivity does not result from heritable mutations in a sample of
1867 DDP-producing clones tested. Even highly sensitive DDP-producing strains were sensitive during
1868 overproduction, but not afterward, when colonies recovered after exposure were cultured and re-
1869 exposed without DDP-gene induction (Figure S5H). The data imply that heritable mutations did
1870 not confer the DNA-damage sensitivities. Further, we used RNA-seq to quantify mRNAs of a
1871 panel of 32 *E. coli* DNA-repair genes, representing five DNA-repair mechanisms, in seven DDP-
1872 overproducing clones that display DNA-damage sensitivity (Figure S5F; Table S1), and that
1873 represent the six major biological DDP clusters (Figure 4N). The repair pathways represented were
1874 nucleotide excision repair (*uvrA*, *uvrB*, *uvrC*, *uvrD*), base-excision repair (BER: *mutT*, *mutM*, *xthA*,
1875 *nfo*, *ung*, *mug*, *nth*, *tag*, *alkA*, *nei*), mismatch repair (*mutS*, *mutL*, *uvrD*, *mutH*), homology-directed
1876 repair (*recA*, *radA*, *ruvA*, *ruvB*, *ruvC*, *recT*, *recF*, *recO*, *recR*, *recG*, *recN*), and homology-directed
1877 DSB repair (as for homology-directed repair with the addition of *recB*, *recC*, *recD*, and without
1878 *recFOR*). Thirty-one of the DNA-repair gene mRNAs were not downregulated with DDP
1879 overproduction, and 19 (*bamC*), 20 (*cusR*), 17 (*topA*), 4 (*aroP*), 8 (*hemX*), 14 (*yicR*) and 1 (*yqjD*)
1880 were significantly upregulated, presumably resulting from SOS or other DNA-damage or
1881 oxidative-stress responses (Figure S5F, Table S1). The sole DNA-repair-gene mRNA
1882 downregulated was that of *mutM*, which encodes the BER protein MutM, a DNA glycosylase that
1883 removes oxidized deoxy-guanine from DNA (Michaels et al., 1991), and which was decreased
1884 with DNA Topoisomerase (Topo I) overproduction (Figure S5F, Table S1). The data from all of
1885 the other genes show that reduced DNA-repair activities, inferred from DNA-damage sensitivities
1886 during overproduction of the seven representative DDPs (Figure S5F), did not result from
1887 transcriptional down-regulation of these DNA-repair genes. In the case of *mutM* mRNA reduction
1888 during Topo I overproduction, this is unlikely to cause the H$_2$O$_2$ sensitivity of Topo I
1889 overproducing cells because *mutM* null mutants are resistant to the low H$_2$O$_2$ levels we used (Asad
1890 et al., 2004). The data exclude the hypothesis that most or all of the DNA-damage sensitivities
1891 caused by overproduction of DDPs (Figures 4J-L and S5A-E; Table S1) result from transcriptional
1892 downregulation of DNA-repair genes. Potential post-transcriptional regulation mechanisms are
1893 not excluded. The data support the hypothesis that the DNA-damage sensitivities result from
1894 excess DNA damage overwhelming DNA-repair capacity, potentially titrating DNA-repair
1895 enzymes, or direct inhibition of DNA repair by the overproduced protein. Reduction of DNA-
1896 repair capacity could produce phenotypes like those of DNA-repair mutants, which would drive
1897 genome instability, in many DDP-gene dysregulated cells that do not possess DNA-repair-gene
1898 mutations.

1899

1900 **Supplemental Discussion 11**
1901 **DNA damage reductions in some function clusters**
1902 Reduced ROS levels are apparent in a cluster in Figure 4M, as are reduced anucleate cells (DNA
1903 loss) in that figure. Similarly, clusters 3, 4, and 6 (Figure 4N) show reduced sensitivity (greater
1904 resistance than wild-type cells) to H$_2$O$_2$, and 2 and 3 (Figure 4N) show greater resistance to

1905 mitomycin C. A probable explanation for these "better-than-wild-type" phenotypes is that some
1906 imbalance in these cells may induce a stress response, one of the consequences of which may be
1907 improvement of the fidelity of, e.g., chromosome segregation (preventing anucleate cells),
1908 reduction of ROS levels, or resistance to DNA-damaging agents, as is documented, for example,
1909 for DNA-damage resistance induced by mild induction of the SOS response (Friedberg et al.,
1910 2005), among other stress responses.
1911
1912 **Supplemental Discussion 12**
1913 **RDG ChIP-seq signals are enriched near known CsgD-binding sites, directionally**
1914 The RDG ChIP-seq experiments (Figures 5I and S7A-C) were performed twice, with peaks called
1915 against the matched control DNA-binding domain mutant, Csg*D*ΔDBD. We identified 155 total
1916 reproducible RDG ChIP-seq peaks with at least 2.5-fold increase in both repeats compared with
1917 CsgDΔDBD done in parallel. The following control simulations show that RDG ChIP-seq signals
1918 are enriched near known CsgD binding sites.  We simulated CsgD-overproduction ChIP-Seq data
1919 by randomly distributing the called RDG peaks across the genome and analyzed, first, the number
1920 of binding sites with at least one RDG peak within 10kb, and, second, the distance between each
1921 known CsgD-binding site and the nearest RDG peak (median), for each of the 10 experimentally
1922 validated CsgD-binding sites (Brombacher et al., 2003; Dudin et al., 2014; Keseler et al., 2017;
1923 Ogasawara et al., 2011). Strikingly, there is at least one observed (actual) RDG peak within 10kb
1924 of 9 out of the 10 known CsgD-binding sites (Figure S7A-C), whereas there are only 5 known
1925 binding sites with at least one simulated RDG peak within 10kb of the 10 known binding sites in
1926 our simulations, which is significantly less than the observed ($p = 0.01$, two-tailed z-test). Also,
1927 the median distance between a given known CsgD-binding site and the nearest RDG ChIP-seq
1928 peak is 2.8kb for the real peaks, which is significantly closer than the 10.3kb median distance in
1929 our simulations ($p = 0.009$).
1930     **CsgD-DBD-dependent stalled-fork RDG peaks not at known CsgD-binding sites.**  For
1931 the 142 CsgD-DBD-dependent RDG peaks that are not significantly near to a known CsgD binding
1932 site, these could result from—(i) binding of CsgD to sites not yet identified in the literature, in
1933 which no comprehensive high-resolution (ChIP-seq) genome-wide binding study has been done;
1934 (ii) relaxation of the binding specificity of CsgD when overproduced such that sites not normally
1935 bound are bound on overproduction; and (ii) downstream (indirect) effects of regulation of CsgD-
1936 regulated genes by CsgD binding to its sites.  For example, CsgD upregulation of other DNA-
1937 binding transcription factors could cause peaks at their binding sites, among other indirect but
1938 biologically real (CsgD-DBD-dependent) possibilities. The 142 of 155 RDG ChIP-seq peaks that
1939 are not near known CsgD binding sites do not overlap significantly with palindromic REP
1940 sequences, which are prone to form DNA cruciform structures (Stern et al., 1984) (four RDG peaks
1941 overlap with REPs*, p* =0.43, two-tailed z-test compared with median random distribution), or *terA-*
1942 *terC* regions, which have higher frequencies of fork convergence (11 RDG peaks overlap with *ter*
1943 sequence*s, p* =0.47, two-tailed z-test, compared with compared with median random distribution).
1944     **Upstream bias of RDG stalled-fork peaks at CsgD binding sites.** The 9 known CsgD-
1945 binding sites with one or more observed CsgD-DBD-dependent RDG peaks within 10kb are of
1946 two types: six show RDG peaks at or upstream of the binding sites in the replication path and five
1947 show RDG peaks downstream. The observed number of CsgD binding sites with upstream or co-
1948 localized RDG peaks is significantly higher than the median in the simulation ($p = 0.04$, two-tailed
1949 z-test), whereas the number of CsgD binding sites with RDG peaks downstream is not quite
1950 significant, at $p = 0.13$, two-tailed z-test. The analysis above suggests that direct binding of CsgD

1951    to its known binding sites underlies the RDG upstream and co-localized peaks, whereas the
1952    downstream peaks may reflect a component of direct CsgD-DNA-induced fork reversal, and a
1953    component of indirect effects, or binding to as yet unknown CsgD binding sites. One way that an
1954    RDG stalled-fork peak might result *downstream* in the replichore, but still via direct interaction of
1955    replication with bound CsgD, could be if some of the downstream stalled forks are caused by two
1956    events (illustrated Figure 5J, lower): first, slowing of replication forks by CsgD binding at its site,
1957    which might make the replisome susceptible to an otherwise surmountable second barrier of any
1958    kind downstream in the replication path.

1959

1960    **Supplemental Discussion 13**
1961    **Model for DNA-damage promotion by the overproduced XanQ xanthine-proton symporter**
1962    Membrane transporters including proton ($H^+$) symporters are overrepresented among DDPs that
1963    cause increased ROS when overproduced (Figures 6A-D), and cause DNA damage ROS-
1964    dependently (Figure 6B). The strongest generator of ROS-dependent DNA damage and ROS is
1965    the XanQ xanthine symporter. Although overproduction of each of three proton symporters
1966    decreases pH (increases $H^+$, Figures 6F,G), their ROS promotion is not well correlated with
1967    decreased pH (Figure 6H), suggesting that the molecules they transport with $H^+$ might underlie
1968    generation of ROS and DNA damage. For XanQ, overproduction of which causes high ROS but
1969    a moderate drop in pH (Figures 6F,G), one possibility (Figure 6I) is that the excess xanthine
1970    imported may be oxidized by ROS-generator xanthine oxidase (Kelley et al., 2010), causing the
1971    increased ROS, which damages DNA (Figure 6I). Other explanations are possible. Regardless of
1972    specific hypotheses, overproduction of membrane transporters generally may promote DNA
1973    damage, and for several, increased ROS, by any of many mechanisms that result from
1974    compromising compartmentalization of the cell from its environment.
1975

1976    **Figure S1.**

1977 **Figure S1. Genuine SOS Response, Detection of DNA Damage by the SOS-reporter Gene,**
1978 **and Evaluation of False Negatives in the Primary Screen**
1979 The *E. coli* DNA-damage assay detects fluorescence caused by upregulation of the SOS DNA-
1980 damage-response-activated promoter P$_{sulA}$ fused to mCherry in a non-genic chromosomal site
1981 (Nehring et al., 2016; Pennington and Rosenberg, 2007). This assay was shown to report on DNA
1982 damage, not spurious promoter firing with the demonstration that fluorescence induction requires
1983 the DNA-damage-sensing protein RecA, and is inhibited by a mutant SOS-response transcriptional
1984 repressor, LexAInd⁻, which does not de-repress the SOS genes during DNA damage (Pennington
1985 and Rosenberg, 2007), both also shown here (A). Further, of the spontaneous SOS-inducing DNA
1986 damage, previously 60% was shown to reflect DNA double-strand breaks (DSBs) and 40% single
1987 stranded DNA not at DSBs (Pennington and Rosenberg, 2007), and the spontaneous DSB
1988 frequency and rates per cell division and chromosome replication were confirmed in a direct,
1989 independent assay for DSB ends (Shee et al., 2013). Thus, this reporter reports on DNA damage.
1990 (A)  Fluorescence requires the ability to activate the SOS response, and so is absent in SOS-
1991 induction-deficient Δ*recA* or *lexA3*(Ind⁻) mutant cells, demonstrating that DNA damage, not
1992 spurious promoter firing, underlies fluorescence increases, per (Pennington and Rosenberg, 2007).
1993 Using these negative controls, a flow-cytometry gate is set in each experiment, per (Pennington
1994 and Rosenberg, 2007), at fluorescence at which $10^{-4}$ of the *lexA3*(Ind⁻) negative- control cells are
1995 positive, shown as a vertical dashed line. The SOS-positive cells in the wild-type (WT) strain
1996 represent spontaneous DNA damage, per (Pennington and Rosenberg, 2007). Representative flow-
1997 cytometry histograms generated under the growth conditions used in the screens reported here.
1998 (B) Positive control, mCherry⁺ *lexA51*(Def) cells with constitutively activated SOS response, and
1999 negative control, *lexA3*(Ind⁻) SOS-off cells.
2000 (C) RecA-dependence of increased fluorescence in representative DDP clones.
2001 (D) Representative DDPs do not enhance mCherry fluorescence generally, when *mCherry* is
2002 controlled by non-SOS promoter P$_{BAD}$.
2003 (E) Screen of 99 random *E. coli* proteins *not* identified in the primary plate-reader screen finds 1
2004 SOS-positive clone. Because of its inherent noise level, the plate-reader primary screen is likely
2005 to generate some false-negative results; *i.e.*, it might miss clones with subtle but real DNA damage-
2006 up phenotypes. We therefore tested a sample of 99 random *E. coli* proteins that were not identified
2007 in the plate-reader primary screen for increased DNA-damage fluorescence in the sensitive flow-
2008 cytometry secondary assay, to estimate the frequency of possible false negatives. We found that 1
2009 of the 99, or ~1%, was positive in the flow-cytometry assay. Thus, if the 3815 *E. coli* proteins that
2010 were not identified in the primary screen (4229 protein-coding genes in the library - 414 proteins
2011 found in the primary screen) also harbor 1% real but undiscovered DDPs, then an additional 38
2012 overproduction DDPs are predicted to reside in the *E. coli* overproduction library, for an estimated
2013 total network size of 246 (38 + 208) DDP genes in the overexpression library used.  In
2014 Supplementary Discussion 2, we estimate an additional 123 *E. coli* overexpression DDP genes that
2015 are likely to be in the *E. coli* genome, but would not be discoverable using the mobile plasmid
2016 overexpression library.
2017 (F) *E. coli* TetR clones from *cI* mutation assay harbor genuine mutations. Sequences of *cI*
2018 mutations from 10 different DDP-overproducing strains with strong damage-up phenotypes:
2019 CsgD, TopB, YegL, GrpE, HslU, Mrr, CheA, MdtA, YicR, and UvrA.  The compiled sequences
2020 from the 10 DDP clones are shown in red and blue (n = 3-10 independent TetR mutants per clone).
2021 The 69 *cI* mutations sequenced include 3 IS (mobile)-element insertions, and 2 clones without a
2022 PCR product (both isolated from TopB overproducing cultures; PCR failure is due to large
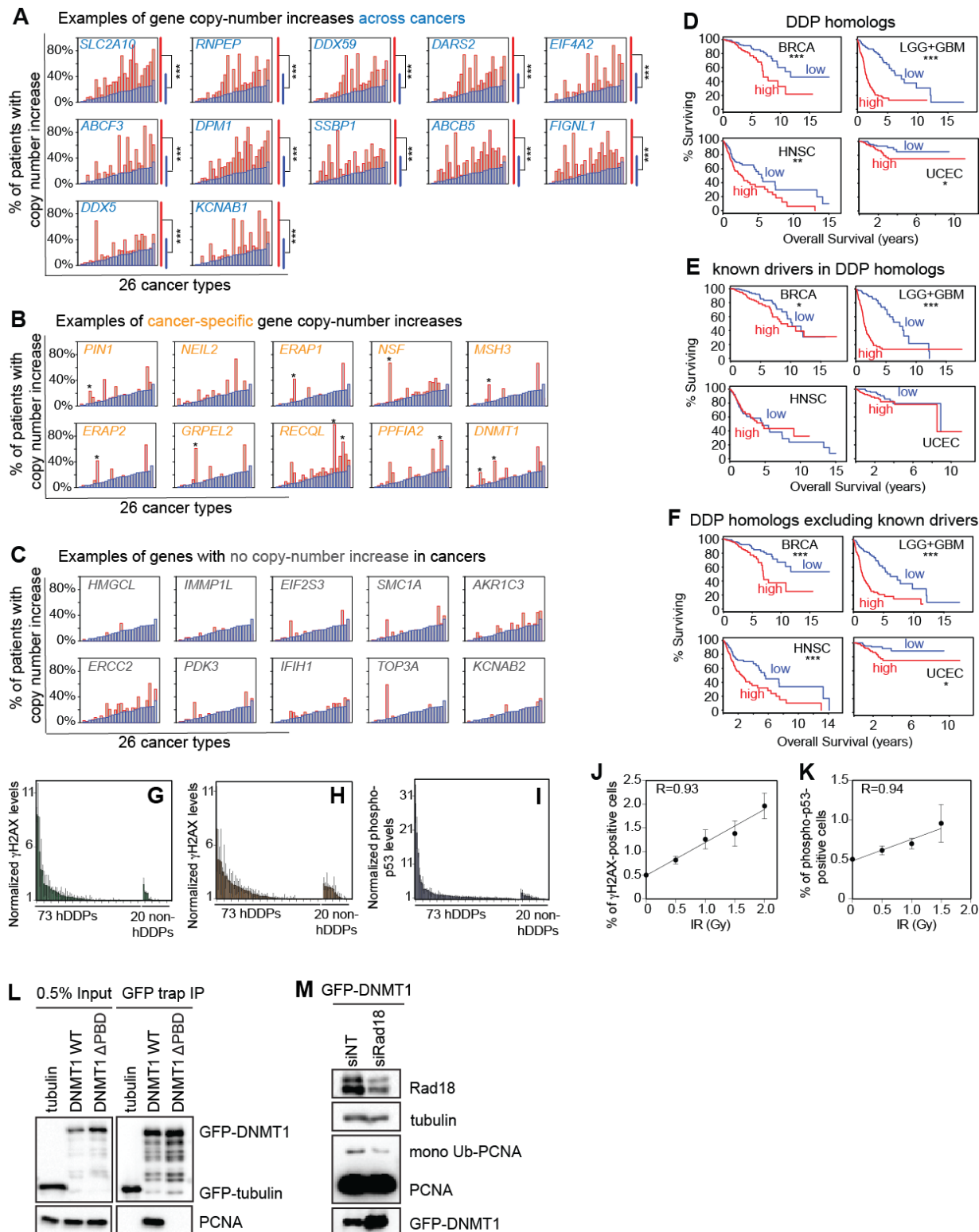
61

2023      deletions, see STAR Methods). Mutations in the vector-only control strain are shown in black (19
2024      independent mutants). Red font, mutations attributable to common errors of the SOS-upregulated
2025      error-prone DNA polymerases V and IV (Kobayashi et al., 2002; Maor-Shoshani et al., 2000;
2026      Wagner and Nohmi, 2000); blue font, mutations not attributable to common Pol V or Pol IV errors.
2027      Capital letters, basepairs that were changed or deleted; carrots, insertion points of bases indicated;
2028      Δs, deletions of the bases indicated. Upper right, diagram of gross chromosomal rearrangements
2029      (GCRs) found in TopB-overproducing cells, among which, 18% (7/40) of the TetR mutations are
2030      GCRs, including large deletions (from 200bp-6200bp), an inversion, and a transposon insertion.
2031

**Figure S2. Specific Protein-Protein Interactions in *E. coli* and Human DDP Networks**

2033 **Figure S2. Specific Protein-Protein Interactions in *E. coli* and Human DDP Networks**

2034 (A) Specific protein-protein interactions in *E. coli* DDP network. Diagram details are as indicated

2035 in Figure 1G. Image enlarged to illustrate the specific proteins with interactions. Connectivity in

2036 random protein networks is discussed in Supplemental Discussion 1.

2037 (B) Specific protein-protein interactions of human candidate DDP network. Diagram details are

2038 as indicated in Figures 1G and 2B. Image enlarged to illustrate the specific proteins with

2039 interactions, similarly to the human candidate DDP network in Figure 2B. Although, the human

2040 homolog (candidate hDDP) network has DNA-repair and -replication genes at its center, removal

2041 of these proteins leaves the remainder with still significant connectivity: $p = 4.1 \times 10^{-215}$

2042 (hypergeometric test), though less than the whole network ($p = 1.2 \times 10^{-327}$, hypergeometric test).

2043 Random human proteins do not form robust protein-protein association networks and the

2044 significant protein-protein interactions of well conserved proteins—the human homologs of

2045 random *E. coli* proteins—are discussed Supplemental Discussion 1.

2046

2047
2048 **Figure S3. Association of Human Homolog Network with Cancers: Copy-number Gain,**
2049 **Survival, and Mutation Load, and Controls for Human DNA-Damage Assays and DNMT1**
2050 (A to C) Twenty-six cancer types in TCGA data (Gao et al., 2013) are displayed along the x axes.
2051 Blue, median % of patient cancers with increased copy number of any gene in their genome; red,

65

% of patient cancers with increased copy number of the gene indicated. Copy-number-increases in human homologs of *E. coli* DDP genes are higher than those of human homologs of random *E. coli* genes ($p < 0.05$, FDR $< 0.10$, Wilcoxon test). Examples of the Pan-Cancer copy-number-increase analysis (GISTIC threshold copy-number gain $\geq 1$) of the 284 human homologs of *E coli* DDP genes in 26 cancer types are shown here (complete analysis Table S4). The human genes fell into three categories:

(A) genes with increased copy numbers across cancers (fold change $> 1.5$, $p < 0.05$, FDR $< 0.10$, Wilcoxon test);

(B) genes with cancer-specific copy-number-increases ($p < 0.05$); and

(C) not particularly cancer associated.

(D) Decreased cancer survival is associated with high DDP-homolog RNA levels in cancers [our analyses of data from TCGA (Gao et al., 2013), STAR Methods]. BRCA, breast invasive carcinoma; LGG+GBM, gliomas (low-grade glioma + glioblastoma multiforme); HNSC, head and neck squamous cell carcinoma; UCEC, uterine corpus endometrial carcinoma. *, **, *** indicate that survival of the cancers with high and low levels of the 284 RNAs differ at $p \leq 0.05$; $\leq 0.01$, and $\leq 0.001$ respectively, log-rank test.

(E) Decreased cancer survival is not confined to the previously known cancer drivers in the network. RNAs of the known (Forbes et al., 2015) and predicted (D'Antonio and Ciccarelli, 2013) cancer-driving genes among the homologs show less association with poor survival than the whole homolog network.

(F) RNAs of the human DDP homologs excluding known (Forbes et al., 2015) and predicted (D'Antonio and Ciccarelli, 2013) drivers are still associated with decreased patient survival, indicating that genes in the homolog network other than and in addition to the previously known drivers are also associated with decreased survival.

(G-I) Human MRC5-SV40 or HEK293T were transfected with human candidate DDP and random human genes and random human non-DDP homologs of *E. coli* genes, and DNA-damage markers were analyzed by flow cytometry.

(G) Assay I: γH2AX in MRC5-SV40 cells.

(H) Assay II: γH2AX in MRC5-SV40 cells treated with DNA-PK inhibitor, which inhibits DNA break repair by non-homologous end joining and so is a sensitizing DNA-damage screen.

(I) Assay III: phosphorylated p53 in HEK293T cells. Data represent mean ± range, n ≥ 2. The candidate hDDP genes differ from the 20 random human genes, $p < 0.0001$, Fisher exact test.

(J and K) Linear response of both human-cell DNA-damage-detection assays with exogenous ionizing radiation treatment indicates quantitative validity of these assays. Percent of MRC5-SV40 cells that are positive for (J) γH2AX, and (K) phosphorylated p53 (p53-S15p), in flow-cytometric assays. MRC5-SV40 cells were treated with IR with the indicated dose and analyzed by flow cytometry.

(L) DNMT1 wild-type (WT) but not ΔPBD interacts with replisome sliding clamp, PCNA. The tubulin negative control also does not interact with PCNA. GFP-trap immunoprecipitation was performed in flipIn stable inducible HEK293T cells expressing GFP-tubulin, GFP-DNMT1-WT or GFP-DNMT1-ΔPBD. Interactions were determined in immunoprecipitation samples by western blotting with anti-GFP and -PCNA antibodies.

(M) Overproduction of wild-type DNMT1, but not the DNMT1-PBD-defective mutant protein, enhances RAD18-mediated ubiquitylation of replisome clamp, PCNA. Knockdown of *RAD18* reduced the level of ubiquitylated PCNA in DNMT-WT overexpressing cells. Western analyses

2097    of PCNA ubiquitylation in MRC5-SV40 cells transfected with the plasmids indicated in
2098    combination with non-targeting (NT) or *RAD18* siRNA.
2099

**Figure S4. DSBs, Reversed Forks, Reactive Oxygen, and DNA Loss in DDP Clones**

**Figure S4. DSBs, Reversed Forks, Reactive Oxygen, and DNA Loss in DDP Clones**

(A) Increased GamGFP foci indicate DSBs caused by overproduction of 87 of the 208 *E. coli* DDPs. DSBs in all 208 *E. coli* DDP-overproducing clones were visualized and quantified as GamGFP foci (Shee et al., 2013) using automated microscopy. 87 of the 208 clones, were significantly different from the vector-only control at $p < 0.05$, q <0.10 (unpaired two-tailed *t*-test with FDR adjustment). Each bar represents a DDP clone (mean ± range, n=2 experiments of >1000 cells per strain, STAR Methods). DDP-overproducing clones are grouped by the fold change of GamGFP focus levels compared with the vector-only strain. Blue, vector only; grey, < 2-fold increase; brown, 2-5-fold increase; green, 5-10-fold increase; magenta, > 10-fold increase, per Table S1 (data summary). Representative images are shown above and to the right of the bar graphs. None of 25 non-DDP overproducers had increased GamGFP foci, showing that clones with high levels of GamGFP foci are enriched in the DDP-overproducing clones ($p = 3.6 \times 10^{-6}$, one-way Fisher's exact test).

(B) Stalled reversed replication forks in most DDP-overproducing strains. Significantly increased reversed forks (RFs), visualized as RDG foci in Δ*recA* cells (Xia et al., 2016) via automated microscopy, occur in 106 of the 208 (51% of) *E. coli* DDP-overproducing clones. The 106 were significantly different from the vector-only control, $p < 0.05$, q <0.10 unpaired two-tailed *t*-test with FDR adjustment. Each bar represents a DDP clone (mean ± range, n=2 experiments of >1000 cells per strain). Blue, vector only; grey, < 2-fold increase; brown, 2-5-fold increase; green, 5-10-fold increase; magenta, >10-fold increase, per Table S1 (complete data summary). Representative images shown above and on the right side of the bar graphs. None of the 30 non-DDP overproducers had increased RDG foci; so, DDP overproducers with high RDG-focus (RF) loads are enriched in the DDP-overproducing strains ($p = 4.0 \times 10^{-9}$, one-way Fisher's exact test).
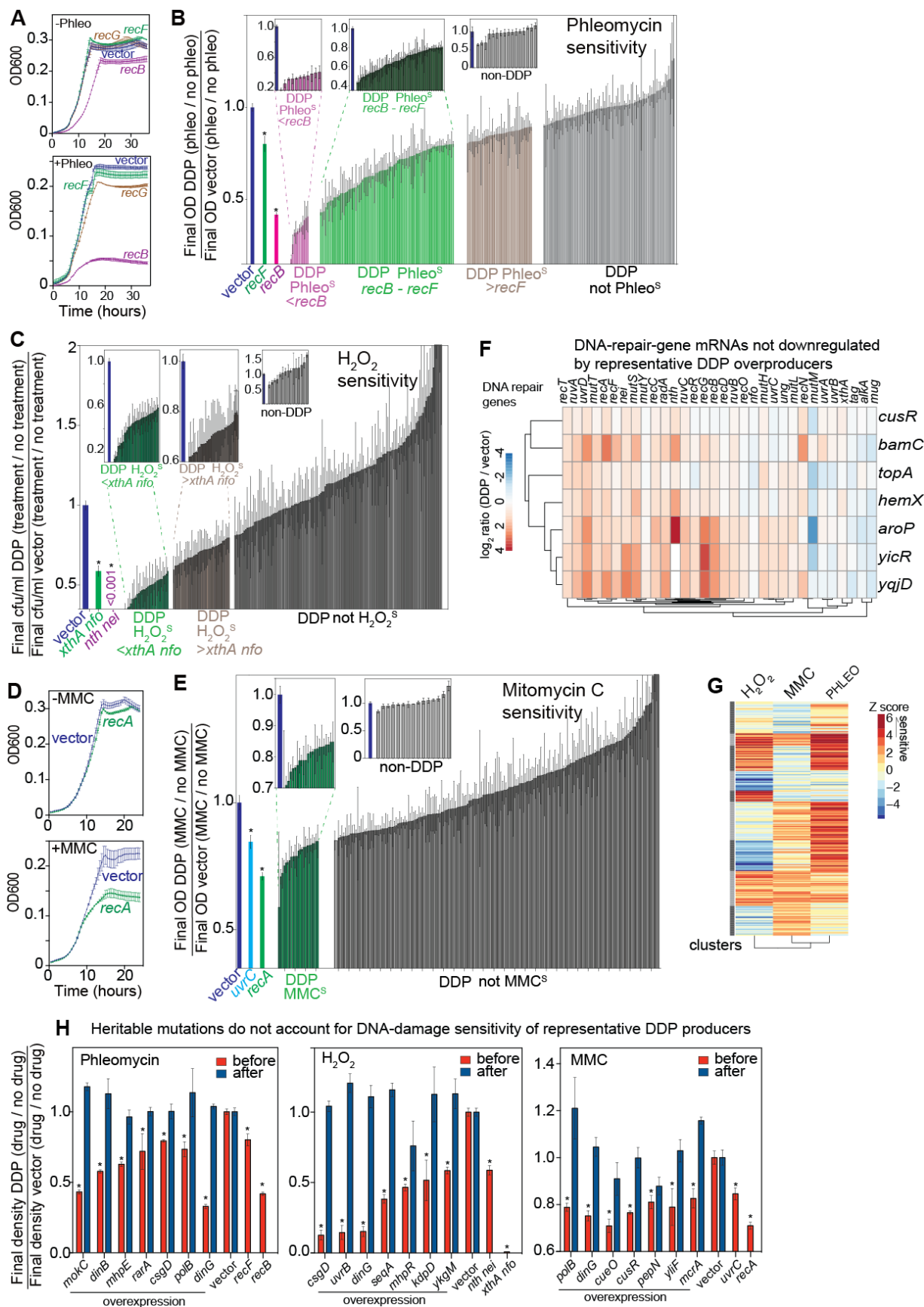
(C-D) Increased intracellular ROS levels in 56 *E. coli* DDP-overproducing clones identified by a flow-cytometric assay after di-hydrorhodamine (DHR) staining. Intracellular ROS levels in all 208 *E. coli* DDP-overproducing clones were measured with the peroxide-specific dye DHR (Gutierrez et al., 2013) and fluorescence was measured by flow cytometry. We found that 56 of the 208 (27%) were significantly different from the vector-only control, $p < 0.05$, q <0.10 unpaired two-tail *t*-test with FDR adjustment. Each bar represents a DDP clone (mean ± range, of two experiments). Blue, vector only; green, DDP overproducers with increased ROS levels, per Table S1. Representative flow cytometry histograms are shown above the bar graphs. None of 17 non-DDP overproducers had increased ROS, such that high-ROS clones are enriched in the DDP-overproducing clones at $p = 0.006$ (one-way Fisher's exact test).

(E-F) DNA loss in 67 *E. coli* DDP-overproducing clones identified by flow-cytometric quantification of anucleate cells.

(E) Increased anucleate cells can be detected in a *parC*TS mutant, which has a severe chromosome-segregation defect at non-permissive temperature. Live cells were stained with FM464 membrane-staining dye, and then fixed and stained with DAPI for DNA staining. Cells with positive membrane but negative DAPI staining are anucleate cells (lower right quadrant).

(F) We found that 67 (32%) were significantly different from the vector-only control, $p < 0.05$, q <0.10 unpaired two-tail *t*-test with FDR adjustment. Each bar represents a DDP-overproducing clone (mean ± range, n=2). DDP overproducers are grouped by the fold change of anucleate cells levels compared with the vector only strain. Blue, vector only; grey, < 2-fold increase; brown, 2-5-fold increase; green, 5-10-fold increase; magenta, >10-fold increase in anucleate cells, per Table

69

2145    S1. Representative flow cytometric histograms are shown above and to the right of the bar graphs.
2146    Only 1 of 16 non-DDP-overproducing clones had increased anucleate cells, such that clones with
2147    high-levels of anucleate cells are enriched in the DDP-overproducing strains at $p = 0.02$ (one-way
2148    Fisher's exact test).
2149

2150 **Figure S5. Sensitivities to DNA-Damaging Agents**

2151 **Figure S5. Sensitivities to DNA-Damaging Agents**

2152 (A) Phleomycin sensitivity detected by growth inhibition in known mutants with reduced
2153 homologous recombinational (HR) repair efficiency. *recF*: defective in single-strand gap HR-
2154 repair; *recG*: reduced ability to branch migrate Holliday junctions; *recB*: defective in DSB repair
2155 by HR (Kuzminov, 2011).

2156 (B) 106 *E. coli* DDP overproducers are sensitive to phleomycin. Phleomycin sensitivities of the
2157 DDP-overproducing clones are shown as normalized to sensitivity of vector-only controls:
2158 (treated/untreated DDP overproducer) / (treated/untreated vector-only) so that values < 1 indicate
2159 sensitivity. Among the 106 sensitive DDP clones, 11 are more sensitive than a *recB* mutant; 72 are
2160 within the range of *recB* and *recF* mutants; 23 are more sensitive than the vector-only control but
2161 less than *recF* mutant. One of the 16 non-DDP-overproducing strains has phleomycin sensitivity,
2162 such that overproducing clones with phleomycin sensitivity are enriched in the DDP-
2163 overproducing clones at $p = 0.0003$ (one-way Fisher's exact test).

2164 (C) Sensitivity of 75 *E. coli* DDP-overproducing clones to oxidative-damage-inducing agent $H_2O_2$.
2165 $H_2O_2$ sensitivity was detected by reduced viability, measured by colony forming units in known
2166 mutants with reduced base excision repair (BER) efficiency. *xthA* (inset) encodes exonuclease III;
2167 *nfo* encodes endonuclease IV. *xthA nfo* (green) double mutants are reported and confirmed by us
2168 to be more sensitive to $H_2O_2$ than each single mutant (Galhardo et al., 2000). *nth* encodes
2169 endonuclease III; *nei* encodes endonuclease VIII. The *nth nei* (purple) double mutant has almost
2170 no AP lyase activity, and so has extreme sensitivity to $H_2O_2$ (Saito et al., 1997), consistent with
2171 our observation. In our assay, $H_2O_2$ sensitivities of the DDP-overproducing clones are shown as
2172 normalized to sensitivity of vector-only controls: (treated/untreated DDP overproducer) /
2173 (treated/untreated vector-only) so that values < 1 indicate sensitivity. Among the 75 $H_2O_2$-sensitive
2174 DDP overproducers, 36 were more sensitive than the *xthA nfo* mutant; 39 were more sensitive than
2175 the vector-only control but less sensitive than the *xthA nfo* mutant. None of the 15 nonDDP-
2176 overproducers had $H_2O_2$ sensitivity, such that overproducing clones with $H_2O_2$ sensitivity are
2177 enriched in the DDP-overproducing clones at $p = 0.002$ (one-way Fisher's exact test).

2178 (D and E) Sensitivity of 10 *E. coli* DDP-overproducing clones to interstrand-crosslinking agent
2179 Mitomycin C (MMC). (D) MMC sensitivity was detected by growth inhibition in known mutants
2180 with defective HR repair or nucleotide excision repair (NER): *recA*, defective in HR-repair and
2181 SOS response; and *uvrC*, defective in NER (shown in E, cyan). (E) Ten *E. coli* DDP-overproducing
2182 clones were sensitive to MMC. Sensitivities of the DDP-overproducing clones are normalized to
2183 sensitivity of vector-only controls: (treated/untreated DDP overproducer) / (treated/untreated
2184 vector-only) so that values < 1 indicate sensitivity. None of the 16 non-DDP overproducers was
2185 MMC sensitive. Overproduction clones with MMC sensitivity are not enriched among the DDP-
2186 overproducing clones at $p = 0.47$ (one-way Fisher's exact test).

2187 (F) Representative DDPs do not downregulate RNAs of DNA-repair genes upon overproduction.
2188 We chose 7 DDPs that confer DNA-damage sensitivity on overproduction representative of the
2189 six DDP function clusters (Figure 4N; Table S1). We assayed RNA levels of a panel of 32 DNA-
2190 repair genes by RNA-seq with and without overproduction of each of 7 DDPs that confer
2191 sensitivity to the following agents: TopA and BamC ($H_2O_2$); CusR, HemX, AroP, YicR, and YqjD
2192 (phleomycin). $Log_2$ ratio (DDP overproducer / vector) for each DNA-repair gene in each of the
2193 seven DDP-overproducing strains was calculated and clustered by hierarchy clustering (Johnson,
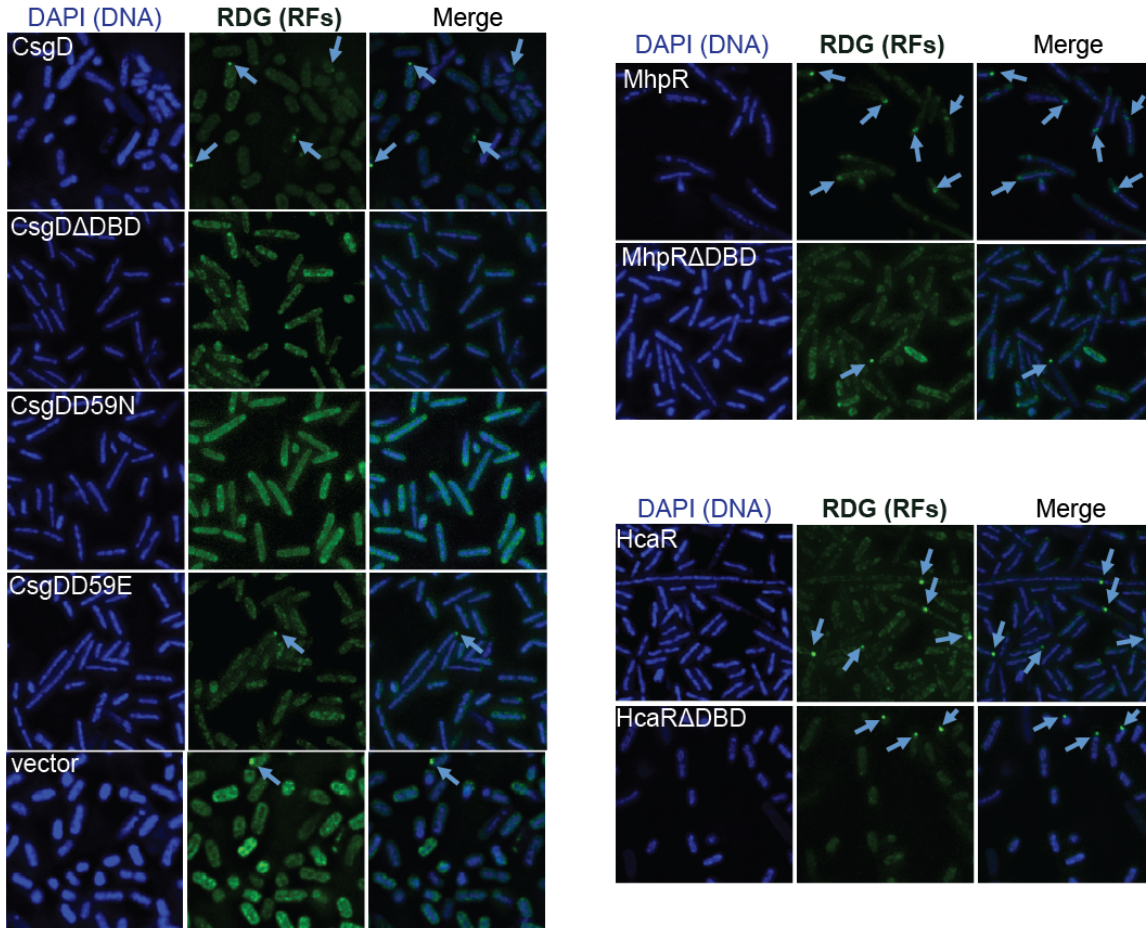2194 1967). Most of the DNA-repair genes are up-regulated during DDP overproduction, with the

2195 exception of *mutM* RNA, which is downregulated in the TopA-overproducing strain. However,
2196 *mutM* downregulation is unlikely to account for the $H_2O_2$ sensitivity of the TopA-overproducing
2197 strain because even *mutM* null mutants are not sensitive to $H_2O_2$ (Asad et al., 2004).

2198 (G) Cluster analysis of quantitative data on DNA-damaging-agent sensitivities in the *E. coli* DDP
2199 network. Each bar represents a quantitative phenotype of each strain producing each of the 208 *E.*
2200 *coli* DDPs, arrayed along the x axis. Assays indicated to right of the heatmap: $H_2O_2$, hydrogen-
2201 peroxide sensitivity (reduced base excision repair); PHLEO, phleomycin sensitivity (reduced
2202 DSB-repair); MMC, mitomycin-C sensitivity (reduced NER and/or HR repair). Red bar: high Z
2203 score, increased DNA-damage sensitivity of DDP-overproducing clones compared with the
2204 vector-only negative control.

2205 (H) Heritable mutations do not account for DNA-damage sensitivities of seven highly DNA-
2206 damage-sensitive DDP-overproducing strains. DDP-overproducing strains with robust sensitivity
2207 to each of the three DNA-damaging agents were tested for sensitivity with overproduction of the
2208 DDP (red, "before"), then again after the drug treatment, with the DDP-gene repressed (blue,
2209 "after") to determine whether their sensitivity resulted from induction of mutations in genes needed
2210 for resistance. After drug treatment, three colonies were recovered on plates with no DDP-gene
2211 inducing IPTG, and then tested for DNA-damaging-agent sensitivity again. In all cases, the
2212 recovered colonies not induced for DDP overproduction showed no sensitivity to DNA-damaging
2213 agents (blue, after). Thus, their DNA-damage sensitivities did not result from heritable mutations.
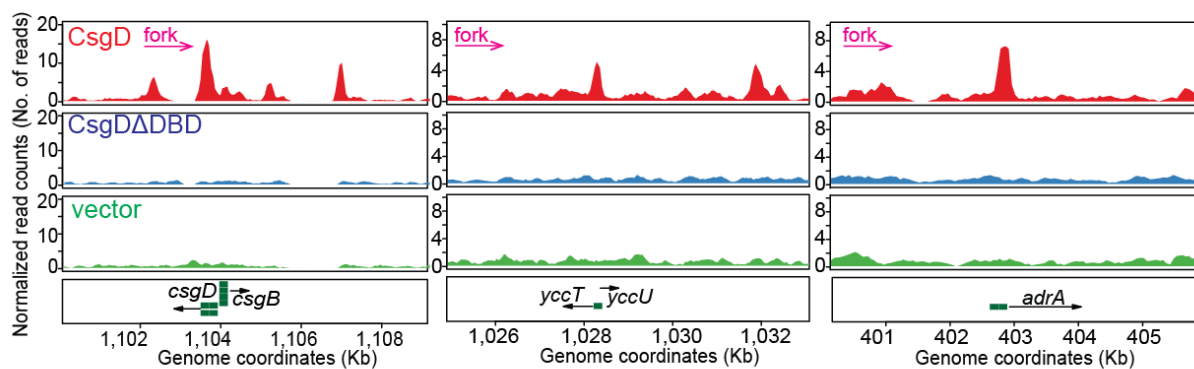2214

2215

2216 **Figure S6. Examples of DNA-binding Transcription Factor Induction of and Co-localization**
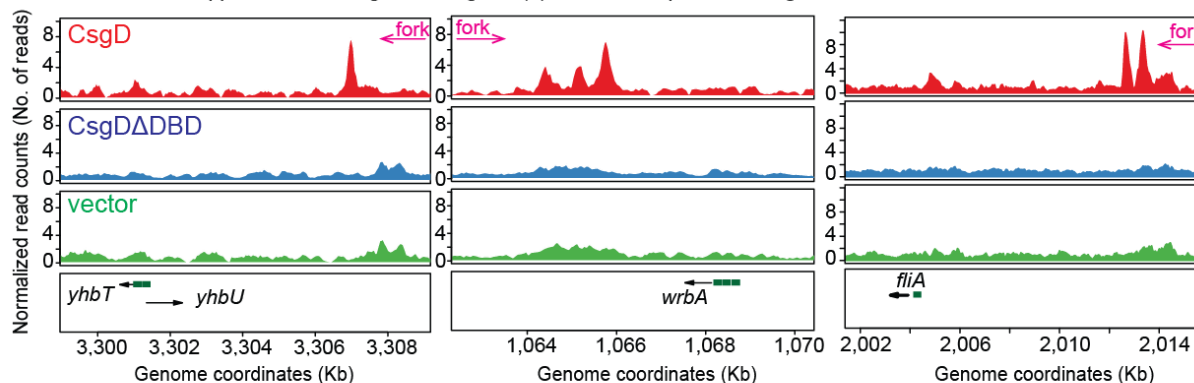2217 **with Replication-Stall (RDG) Foci**

2218 (A) DNA-binding ability of DNA-binding transcription factors is required for their promotion of
2219 increased RDG (reversed-fork) foci upon overproduction. ΔDBD, in-frame deletion of the DNA-
2220 binding domain; CsgDD59N, D59E: single amino-acid changes that reduce CsgD DNA binding
2221 (Ogasawara et al., 2011). Three wild-type transcription factors and the corresponding mutants with
2222 reduced DNA-binding ability were overproduced in Δ*recA* cells, RDG (reversed-fork) foci
2223 quantified (Figure 4C; Table S1), and representative images are shown here.

2224 (B) Foci of transcription factors CsgD-mCherry and HcaR-mCherry co-localize with RDG
2225 (reversed-fork) foci. Representative examples. Most of the CsgD-mCherry and HcaR-mCherry
2226 foci were co-localized with RDG (reversed-fork) foci. Showing weak co-localization, about 5% of
2227 the YahB-mCherry transcription-factor foci were co-localized with RDG (reversed-fork) foci. By
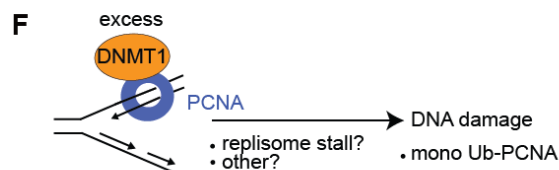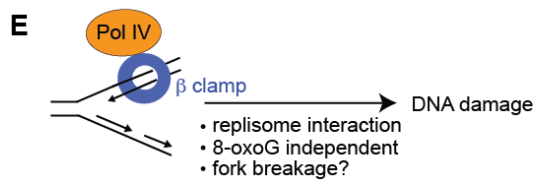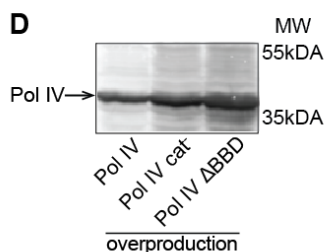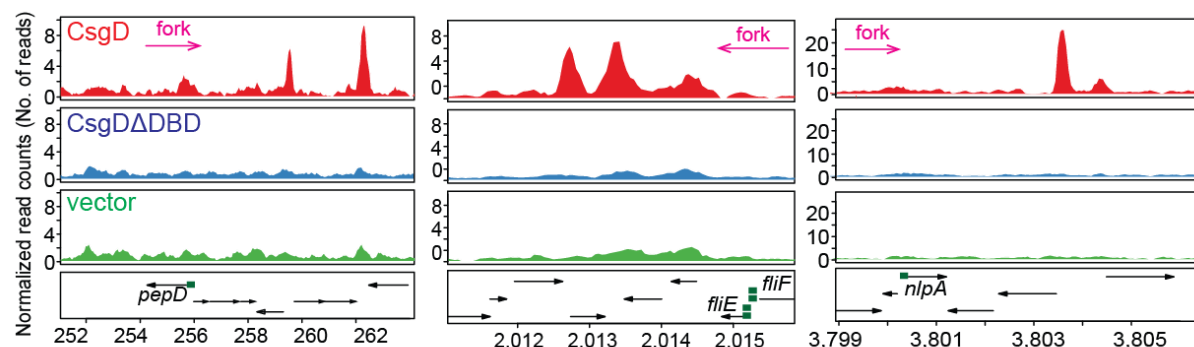2228 contrast, mCherry alone rarely forms foci.
2229

**Figure S7.**

**Figure S7. RDG ChIP-seq Detects Stalled-Fork Enrichment Near CsgD-Binding Sites and Controls for Pol IV**

Nine of the 10 known, validated CsgD-binding sites (Brombacher et al., 2003; Dudin et al., 2014; Keseler et al., 2017; Ogasawara et al., 2011) showed a CsgD-DNA-binding-domain (DBD)-dependent RDG peak nearby: within 10kb (median 2.8kb), which differs from random simulations of the genomic distribution of the total number of CsgD-DBD-dependent RDG peaks ($p = 0.01$ two-tailed z test, see Supplemental Discussion 12). Isogenic cells overproducing—CsgD, red; CsgD ΔDBD (deletion of the DNA-binding domain), blue; vector only green. Green boxes, CsgD binding site(s).

(A) CsgD-DBD-dependent RDG ChIP-seq peaks overlap with three known CsgD-binding sites.

(B) CsgD-DBD-dependent RDG ChIP-seq peaks near three known CsgD-binding sites, upstream in the replication path. Increased negative supercoiling between the oncoming fork and the CsgD-bound site may stall replication and causes fork reversal, illustrated Figure 5J.

(C) CsgD-DBD-dependent RDG ChIP-seq peaks located near three known CsgD binding sites, downstream in the replication path. Because the upstream peaks (B) are more significantly associated with the known CsgD-binding sites relative to simulation of random genomic distributions (Supplemental Discussion 12), it is possible that the downstream CsgD-DBD-dependent RDG peaks could result from either CsgD binding to sites not yet known, or from indirect consequences of CsgD DNA binding, such as effects of CsgD-regulated gene products interacting with other sites in DNA (Supplemental Discussion 12). Alternatively, some of the downstream CsgD-DBD-dependent RDG peaks could be a direct result of CsgD binding its known site, slowing replication, then encountering an otherwise surmountable obstacle downstream, per Figure 5J.

(D) The Pol IV R49F catalytic-mutant (Pol IV cat⁻) and ΔBBD-mutant proteins do not display reduced protein levels in western blots, in agreement with previous studies (Uchida et al., 2008; Wagner et al., 1999).

(E) Model: Pol IV overproduction induces DNA damage by binding the beta replisome sliding clamp. Excess Pol IV interaction with the replisome could potentially slow the replisome causing fork breakage or collapse, or displace DNA-repair proteins that interact with the beta clamp, or otherwise promote DNA damage.
(F) Model and hypotheses: overproduced DNMT1 provokes DNA damage independently of its DNA-methylase activity but dependently on binding PCNA, the mammalian replisome sliding clamp and structural homolog of *E. coli* beta. Excess DNMT1 might promote DNA damage by stalling DNA replication, interfering with PCNA-coordinated DNA-repair or translesion-synthesis processes, or by other PCNA-dependent means.