1    Transcriptional landscape of DNA repair genes underpins a pan-cancer prognostic

2    signature associated with cell cycle dysregulation and tumor hypoxia

3

4

5

6    Wai Hoong Chang and Alvina G. Lai

7

8

9    Nuffield Department of Medicine, University of Oxford,

10    Old Road Campus, Oxford, OX3 7FZ, United Kingdom

11

12

13    For correspondence: alvina.lai@ndm.ox.ac.uk

14    <u>Abstract</u>

15

16    Overactive DNA repair contributes to therapeutic resistance in cancer. However, pan-cancer

17    comparative studies investigating the contribution of *all* DNA repair genes in cancer

18    progression employing an integrated approach have remained limited. We performed a multi-

19    cohort retrospective analysis to determine the prognostic significance of 138 DNA repair genes

20    in 16 cancer types (n=16,225). Cox proportional hazards analyses revealed a significant

21    variation in the number of prognostic genes between cancers; 81 genes were prognostic in

22    clear cell renal cell carcinoma while only two genes were prognostic in glioblastoma. We

23    reasoned that genes that were commonly prognostic in highly correlated cancers revealed by

24    Spearman's correlation analysis could be harnessed as a molecular signature for risk

25    assessment. A 10-gene signature, uniting prognostic genes that were common in highly

26    correlated cancers, was significantly associated with overall survival in patients with clear cell

27    renal cell (P<0.0001), papillary renal cell (P=0.0007), liver (P=0.002), lung (P=0.028), pancreas

28    (P=0.00013) or endometrial (P=0.00063) cancers. Receiver operating characteristic analyses

29    revealed that a combined model of the 10-gene signature and tumor staging outperformed

30    either classifiers when considered alone. Multivariate Cox regression models incorporating

31    additional clinicopathological features revealed that the signature was an independent

32    predictor of overall survival. Tumor hypoxia is associated with adverse outcomes. Consistent

33    across all six cancers, patients with high 10-gene and high hypoxia scores had significantly

34    higher mortality rates compared to those with low 10-gene and low hypoxia scores. Functional

35    enrichment analyses revealed that high mortality rates in patients with high 10-gene scores

36    were attributable to an overproliferation phenotype. Death risk in these patients was further

37    exacerbated by concurrent mutations of a cell cycle checkpoint protein, *TP53*. The 10-gene

38      signature identified tumors with heightened DNA repair ability. This information has the

39      potential to radically change prognosis through the use of adjuvant DNA repair inhibitors with

40      chemotherapeutic drugs.

41      [298 words]

42

43      <u>Keywords:</u> DNA repair, pan-cancer, cell cycle, hypoxia, tumor microenvironment

44

45      <u>List of abbreviations:</u>

|  |  |
|---|---|
| DDR | DNA damage response |
| BER | Base excision repair |
| NER | Nucleotide excision repair |
| MR | Mismatch repair |
| HDR | Homology-directed repair |
| NHEJ | Non-homologous end joining |
| FA | Fanconi anemia |
| TCGA | The Cancer Genome Atlas |
| GO | Gene Ontology |
| KEGG | Kyoto Encyclopedia of Genes and Genomes |
| HR | Hazard ratio |
| ROC | Receiver operating characteristic |
| AUC | Area under the curve |
| TNM | Tumor, node and metastasis |
| CDK | Cyclin-dependent kinase |
| DEG | Differentially expressed genes |

46

## Introduction

Genetic material must be transmitted in its original, unaltered form during cell division. However, DNA faces continuous assaults from both endogenous and environmental agents contributing to the formation of permanent lesions and cell death. To overcome DNA damage threats, living systems have evolved highly coordinated cellular machineries to detect and repair damages as they occur. However, DNA repair mechanisms and consequently DNA damage responses (DDR) are often deregulated in cancer cells and such aberrations may contribute to cancer progression and influence prognosis. Overexpression of DNA repair genes allow tumor cells to overcome the cytotoxic effects of radiotherapy and chemotherapy. As such, inhibitors of DNA repair can increase the vulnerability of tumor cells to chemotherapeutic drugs by preventing the repair of deleterious lesions[1].

There are six main DNA repair pathways in mammalian cells. Single-strand DNA damage are repaired by the base excision repair (BER), nucleotide excision repair (NER) and mismatch repair (MR) pathways. The poly(ADP-ribose) polymerase (PARP) gene family encodes key players of the BER pathway involved in repairing damages induced by ionizing radiation and alkylating agents[2,3]. Replication errors are corrected by the MR pathway while the NER pathway is responsible for removing bulky intercalating agents[4,5]. Tumor cells with deficiencies in the NER pathway have increased sensitivity to platinum-based chemotherapeutic drugs (cisplatin, oxaliplatin etc.)[6,7]. Double-strand breaks induced by ionizing radiation are more difficult to repair and thus are highly cytotoxic. Dysregulation of genes involved in the homology-directed repair (HDR), non-homologous end joining (NHEJ) and Fanconi anemia (FA) pathways are associated with altered repair of double-strand breaks.

71

72     Aberrations in DNA repair genes are widespread in most cancers; hence they represent

73     attractive candidates for pharmacological targeting to improve radiosensitivity and

74     chemosensitivity[8]. In a process known as 'synthetic lethality', faults in two or more DNA repair

75     genes or pathways together would promote cell death, while defects in a single pathway may

76     be tolerated[1]. Functional redundancies in repair pathways allow tumor cells to rely on a second

77     pathway for repair in the event that the first pathway is defective. Based on the principles of

78     synthetic lethality, inhibition of the second pathway will confer hypersensitivity to cytotoxic

79     drugs in cells with another malfunctioning pathway. This promotes cell death because DNA

80     lesions can no longer be repaired by either pathway. For instance, PARP inhibitors (targeting

81     the BER pathway) could selectively kill tumor cells that have *BRCA1* or *BRCA2* mutations

82     (defective HDR pathway) while not having any toxic effects on normal cells[9,10].

83

84     Since one DDR pathway could compensate for another, there is a need for a pan-cancer, large-

85     scale, systematic study on *all* DNA repair genes to reveal similarities and differences in DDR

86     signaling between cancer types, which is limited at present. In this study, we explored pan-

87     genomic expression patterns of 138 DNA repair genes in 16 cancer types. We developed and

88     validated the prognostic significance of a 10-gene signature that can be used for rapid risk

89     assessment and patient stratification. There are considerable variations in the success of

90     chemotherapy and radiotherapy regimes between cancer types. Such differences may be

91     explained by the complex cancer-specific nature of DDR defects. Prognostic biomarkers of DNA

92     repair genes are needed to allow the use of repair inhibitors in a stratified, non-universal

93     approach to expose the selective vulnerabilities of tumors to therapeutic agents.

94    Materials and methods

95    A list of 138 DNA repair genes is available in Table S1.

96    Study cohorts

97    We obtained RNA-sequencing datasets for the 16 cancers from The Cancer Genome Atlas

98    (TCGA)[11] (n=16,225) (Table S2). TCGA Illumina HiSeq rnaseqv2 Level 3 RSEM normalized data

99    were retrieved from the Broad Institute GDAC Firehose website. Gene expression profiles for

100   each cancer types were separated into tumor and non-tumor categories based on TCGA

101   barcodes and converted to $\log_2(x + 1)$ scale. To compare the gene-by-gene expression

102   distribution in tumor and non-tumor samples, violin plots were generated using R. The

103   nonparametric Mann-Whitney-Wilcoxon test was used for statistical analysis.

104

105   Calculation of 10-gene scores and hypoxia scores

106   The 10-gene scores for each patient were determined from the mean $\log_2$ expression values

107   of 10 genes: *PRKDC, NEIL3, FANCD2, BRCA2, EXO1, XRCC2, RFC4, USP1, UBE2T* and *FAAP24*).

108   Hypoxia scores were calculated from the mean $\log_2$ expression values of 52 hypoxia signature

109   genes[12]. For analyses in Figure 5, patients were delineated into four categories using median

110   10-gene scores and hypoxia scores as thresholds. The nonparametric Spearman's rank-order

111   correlation test was used to determine the relationship between 10-gene scores and hypoxia

112   scores.

113

114   Differential expression analyses comparing expression profiles of high-score and low-score

115   patients

116    Patients were median dichotomized into low- and high-score groups based on their 10-gene

117    scores in each cancer type. Differential expression analyses were performed using the linear

118    model and Bayes method executed by the limma package in R. P values were adjusted using

119    the Benjamini-Hochberg false discovery rate procedure. We considered genes with $\log_2$ fold

120    change of > 1 or < -1 and adjusted P-values < 0.05 as significantly differentially expressed

121    between the two patient groups.

122

123

124    <u>Functional enrichment and pathway analyses</u>

125    To determine which biological pathways were significantly enriched, differentially expressed

126    genes were mapped against the Gene Ontology (GO) and Kyoto Encyclopedia of Genes and

127    Genomes (KEGG) databases using GeneCodis[13]. The Enrichr tool was used to investigate

128    transcription factor protein-protein interactions that were associated with the differentially

129    expressed genes[14,15].

130

131

132    <u>Survival analysis</u>

133    Univariate Cox proportional hazards regression analyses were performed using the R survival

134    and survminer packages to determine if expression levels of individual DNA repair genes as

135    well as those of the 10-gene scores were significantly associated with overall survival.

136    Multivariate Cox regression was employed to determine the influence of additional clinical

137    variables on the 10-gene signature. Hazard ratios (HR) and confidence intervals were

138    determined from the Cox models. HR greater than one indicated that a covariate was positively

139    associated with even probability or increased hazard and negatively associated with survival

140    duration. Non-significant relationship between scaled Schoenfeld residuals supported the

141    proportional hazards assumption in the Cox model. Both survival and survminer packages were

142    also used for Kaplan-Meier analyses and log-rank tests. For Kaplan-Meier analyses, patients

143    were median dichotomized into high- and low-score groups using the 10-gene signature. To

144    determine the predictive performance (specificity and sensitivity) of the signature in relation

145    to tumor staging parameters, we employed the receiver operating characteristic (ROC) analysis

146    implemented by the R survcomp package, which also calculates area under the curve (AUC)

147    values. AUC values can fall between 1 (perfect marker) and 0.5 (uninformative marker).

148

149    *TP53* mutation analysis

150    TCGA mutation datasets (Level 3) were retrieved from GDAC Firehose to annotate patients

151    with mutant *TP53*. To ascertain the association of *TP53* mutation with the 10-gene signature

152    on overall survival, we employed the Kaplan-Meier analysis and log-rank tests implemented in

153    R.

154

155    All plots were generated using R pheatmap and ggplot2 packages[16]. Venn diagram was

156    generated using the InteractiVenn tool[17].

157    Results

158

159    **Prognosis of DNA repair genes in 16 cancer types and the development of a 10-gene signature**

160    A total of 187 genes associated with six DDR pathways found in mammalian cells were curated:

161    BER (33 genes), MR (23 genes), NER (39 genes), HDR (26 genes), NHEJ (13 genes) and FA (53

162    genes)[18] (Fig. 1, Table S1). Of the 187 genes, 49 were represented in two or more pathways,

163    yielding 138 non-redundant candidates. To determine which of the 138 DNA repair genes

164    conferred prognostic information, we employed Cox proportional hazards regression on all

165    genes individually on 16 cancer types to collectively include 16,225 patients[11] (Table S2). In

166    clear cell renal cell carcinoma, 81 genes were found to be significantly associated with overall

167    survival; this cancer had the highest number of prognostic DNA repair genes (Table S3). This is

168    followed by 54, 53, 46, 44 and 33 prognostic genes in cancers of the pancreas, papillary renal

169    cell, liver, lung and endometrium respectively (Table S3). In contrast, cancers of the brain

170    (glioblastoma: 2 genes), breast (5 genes), cervix (6 genes) and esophagus (7 genes) had some

171    of the lowest number of prognostic DNA repair genes (Table S3), suggesting that there is a

172    significant degree of variation in the contribution of DNA repair genes in predicting survival

173    outcomes. Spearman's rank-order correlation analysis revealed a hub of five highly correlated

174    cancers (lung, papillary renal cell, pancreas, liver and endometrium), indicating that a good

175    number of prognostic DNA repair genes were shared between these cancers (Spearman's

176    rho=0.21 to 0.44) (Fig. S1). We rationalized that prognostic genes that are common in these

177    highly correlated cancers could form a new multigenic risk assessment classifier. Ten genes

178    were prognostic in the five highly correlated cancers: *PRKDC* (NHEJ)*, NEIL3* (BER)*, FANCD2* (FA)*,*

179    *BRCA2* (HDR and FA)*, EXO1* (MR)*, XRCC2* (HDR)*, RFC4* (MR and NER)*, USP1* (FA)*, UBE2T* (FA) and

180    *FAAP24* (FA), which, interestingly, represent members from all six DDR pathways.

181

## A 10-gene signature predictive of DDR signaling is an independent prognostic classifier in 6 cancer types

184     The aforementioned ten genes were employed as a new prognostic model to evaluate whether

185     they were significantly associated with overall survival in all 16 cancer types. A 10-gene score

186     for each patient was calculated by taking the mean expression of all ten genes. Patients were

187     median dichotomized based on their 10-gene scores into a low- and high-score groups. The

188     10-gene signature could predict patients at significantly higher risk of death in the five cancers

189     that were originally highly correlated (Fig. S1), and in one additional cancer (clear cell renal cell

190     carcinoma) (Fig. 2). Kaplan-Meier analyses demonstrated that patients categorized within

191     high-score groups had significantly poorer survival rates: clear cell renal cell (log-rank

192     $P<0.0001$), papillary renal cell ($P=0.0007$), liver ($P=0.002$), lung ($P=0.028$), pancreas

193     ($P=0.00013$) and endometrium ($P=0.00063$) (Fig. 2). Expression profiles of the 10 genes in

194     tumor and non-tumor samples showed a general distribution that were comparable among

195     the six cancer types. Mann-Whitney-Wilcoxon tests revealed that a vast majority of genes were

196     significantly upregulated in tumor samples with a few minor exceptions (Fig. S2). *USP1* was

197     significantly downregulated in tumors of papillary renal cell and endometrium (Fig. S2). Only

198     four non-tumor samples were available in the pancreatic cancer cohort, precluding robust

199     statistical analyses. Due to limitations in sample size, only *UBE2T* was observed to be

200     significantly upregulated in pancreatic tumors (Fig. S2).

201

202     To evaluate the independent predictive value of the signature over the current tumor, node

203     and metastasis (TNM) staging system, we applied the signature on patients separated by TNM

204     stage: early (stages 1 and/or 2), intermediate (stages 2 and/or 3) and late (stages 3 and/or 4)

205    disease stages. Remarkably, the signature successfully identified high risk patients in early

206    (liver, lung, pancreas, endometrium), intermediate (papillary renal cell, liver, pancreas,

207    endometrium) and late (clear cell renal cell, papillary renal cell, liver, endometrium) TNM

208    stages (Fig. 3). Collectively, this implied that the signature offered an additional resolution of

209    prognosis within similarly staged tumors and that the signature retained excellent prognostic

210    ability in individual tumor groups when considered separately.

211

212    To evaluate the predictive performance of the 10-gene signature on 5-year overall survival, we

213    employed receiver operating characteristic (ROC) analyses on all six cancers. Comparing the

214    sensitivity and specificity of the signature in relation to TNM staging revealed that the signature

215    outperformed TNM staging in cancers of the papillary renal cell (AUC=0.832 vs. AUC=0.640),

216    pancreas (AUC=0.697 vs. AUC=0.593) and endometrium (AUC=0.700 vs. AUC=0.674) (Fig. 4).

217    Importantly, when the signature was used in conjunction with TNM staging as a combined

218    model, its performance was superior to either classifiers when they were considered

219    individually: clear cell renal cell (AUC=0.792), papillary renal cell (AUC=0.868), liver

220    (AUC=0.751), lung (AUC=0.693), pancreas (AUC=0.698) and endometrium (AUC=0.764) (Fig.

221    4).

222

223    We next employed multivariate Cox regression models to examine whether the association

224    between high 10-gene scores and increased mortality was not due to underlying clinical

225    characteristics of the tumors. Univariate analysis revealed that TNM staging is not prognostic

226    in pancreatic cancer (hazard ratio [HR]=1.339, P=0.153), hence this cancer was excluded from

227    the multivariate model involving TNM (Table 1). For the five remaining cancer types, even

228    when TNM staging was considered, the signature significantly distinguished survival outcomes

229    in high- versus low-score patients, confirming that it is an independent prognostic classifier:

230    clear cell renal cell (HR=1.555, P=0.0058), papillary renal cell (HR=1.677, P=0.032), liver

231    (HR=1.650, P=0.029), lung (HR=1.301, P=0.032) and endometrium (HR=2.113, P=0.013) (Table

232    1).

233

234

235    <u>Crosstalk between DDR signaling and tumor hypoxia</u>

236    Tumor hypoxia is a well-known barrier to curative treatment. It is often associated with poor

237    prognosis[19,20], which may be a result of tumor resistance to chemotherapy and

238    radiotherapy[21,22]. Since both the upregulation of DNA repair genes and hypoxia are linked to

239    therapeutic resistance, we rationalized that incorporating hypoxia information in the 10-gene

240    signature would allow further delineation of patient risk groups. Patients with high 10-gene

241    scores had significantly poorer survival outcomes and we predict that these patients have

242    tumors that are more hypoxic, and that oxygen deprivation could influence DDR signaling to

243    enhance tumor resistance to apoptotic stimuli leading to more aggressive disease states. We

244    calculated hypoxia scores for each patient using a mathematically derived hypoxia gene

245    signature consisting of 52 genes[12]. Hypoxia scores were defined as the mean expression of the

246    52 genes. Patients for each of the six cancer types were divided into four categories using the

247    median 10-gene and hypoxia scores: 1) high scores for both 10-gene and hypoxia, 2) high 10-

248    gene and low hypoxia scores, 3) low 10-gene and high hypoxia scores and 4) low scores for

249    both 10-gene and hypoxia (Fig. 5A). Remarkably, significant positive correlations were

250    observed between 10-gene scores and hypoxia scores consistent across all six cancer types:

251    clear cell renal cell (rho=0.363, P<0.0001), papillary renal cell (rho=0.518, P<0.0001), liver

252    (rho=0.615, P<0.0001), lung (rho=0.753, P<0.0001), pancreas (rho=0.582, P<0.0001) and

253    endometrium (rho=0.527, P<0.0001) (Fig. 5A). This suggests that tumor hypoxia may influence

254    DDR signaling and potentially, patient outcomes.

255

256    We generated Kaplan-Meier curves and employed the log-rank test to determine whether

257    there were differences in overall survival outcomes among the four patient groups. Combined

258    relation of hypoxia and 10-gene scores revealed significant associations with overall survival in

259    all six cancers (Fig. 5B). Patients classified within the 'high 10-gene and high hypoxia' category

260    had significantly poorer survival rates compared to those with low 10-gene and low hypoxia

261    scores: clear cell renal cell (HR=2.316, P<0.0001), papillary renal cell (HR=7.635, P=0.0011),

262    liver (HR=2.615, P=0.00013), lung (HR=1.832, P=0.0021), pancreas (HR=2.680, P=0.00079) and

263    endometrium (HR=2.707, P=0.0075) (Table 2; Fig. 5B). Our results suggest that the combined

264    effects of hypoxia and heightened expression of DNA damage repair genes may be linked to

265    tumor progression and increased mortality risks. It remains unknown in this context whether

266    the basis for differential sensitivity to chemotherapy would be explained, in part, by DNA repair

267    ability of tumor cells exposed to chronic hypoxia environments.

268

269

270    **Patients with high 10-gene scores had an overproliferation phenotype due to cell cycle**

271    **dysregulation**

272    The cell cycle represents a cellular gatekeeper that controls how cells grow and proliferate.

273    Cyclins and cyclin-dependent kinases (CDKs) allow cells to progress from one cell cycle stage

274    to the next; a process that is antagonized by CDK inhibitors. Many tumors overexpress cyclins

275    or inactivate CDK inhibitors, hence resulting in uncontrolled cell cycle entry, loss of checkpoint

276    and uninhibited proliferation[23–25]. Targeting proteins responsible for cell cycle progression

277 would thus be an attractive measure to limit tumor cell proliferation. This has led to the

278 development of numerous CDK inhibitors as anticancer agents[26,27]. DNA repair is tightly

279 coordinated with cell cycle progression. Certain DNA repair mechanisms are dampened in non-

280 proliferating cells, while repair pathways are often perturbed during tumor development.

281 Perturbation can take the form of defective DNA repair or over-compensation of a pathway

282 arising from defects in another pathway[28]. As a result, DNA repair inhibitors could prevent the

283 repair of lesions induced by chemotherapeutic drugs to trigger apoptosis and to enhance the

284 elimination of tumor cells.

285

286 We rationalize that patients with high 10-gene scores would have heightened ability for DNA

287 repair thus allowing tumor cells to progress through the cell cycle and continue to proliferate.

288 Using Spearman's rank-order correlation, we observed that the expression of each of the 10

289 signature genes were positively correlated with the expression of genes involved in cell cycle

290 progression (cyclins and CDKs) and negatively correlated with genes involved in cell cycle arrest

291 (CDK inhibitors) (Fig. 6A). Interestingly, the patterns of correlation were remarkably similar

292 across all six cancer types, implying that elevated expression of DNA repair genes is associated

293 with a hyper-proliferative phenotype. We next asked whether patients within the high 10-gene

294 score category had an overrepresentation of processes associated with cell cycle dysregulation

295 as this could provide an explanation on the elevated mortality risks in these patients. To answer

296 this, we divided patients from each of the six cancer types into two groups (high score and low

297 score) based on the mean expression of the 10 signature genes using the 50th percentile cut-

298 off. Differential expression analyses between the high- and low-score groups revealed that

299 394, 425, 1259, 1279, 714 and 977 genes were differentially expressed (-1 > $\log_2$ fold-change

300    > 1, P<0.05) in clear cell renal cell, papillary renal cell, liver, lung, pancreas and endometrial

301    cancers respectively (Table S4).

302

303    Analyses of biological functions of these genes revealed functional enrichment of ontologies

304    associated with cell division, mitosis, cell cycle, cell proliferation, DNA replication and

305    homologous recombination consistent in all six cancer types (Fig. 6B). This suggests that the

306    significantly higher mortality rates in patients with high 10-gene scores were due to enhanced

307    tumor cell proliferation exacerbated by the ability of these cells to repair DNA lesions as they

308    arise. Additional ontologies related to tumorigenesis such as *PPAR* and *TP53* signaling were

309    also associated with poor prognosis (Fig. 6B). A total of 87 differentially expressed genes (DEGs)

310    were found to be in common in all six cancer types (Fig. S3) (Table S5). To dissect the underlying

311    biological roles of the 87 DEGs at the protein level, we evaluated the enrichment of

312    transcription factor protein-protein interactions using the Enrichr platform[14].*TP53* represents

313    the most enriched transcription factor involved in the regulation of the DEGs as evidenced by

314    the highest combined score, which takes into account both Z score and P value (Table S6). This

315    indirectly corroborated our results on enriched *TP53* signaling obtained from the KEGG

316    pathway analysis (Fig. 6B). Taken together, these results highlight the interplay between DDR

317    signaling, cell cycle regulation and *TP53* function in determining prognosis.

318

319

320    <u>Prognostic relevance of a combined model involving the 10-gene signature and *TP53* mutation</u>

321    <u>status</u>

322    An important role of *TP53* is its tumor suppressive function through *TP53*-mediated cell cycle

323    arrest and apoptosis[29]. Hence, somatic mutations in *TP53* can confer tumor cells with growth

324  advantage and indeed, this is a well-known phenomenon in many cancers[30–32]. We rationalized

325  that *TP53* deficiency resulting in defective checkpoint may synergize with the overexpression

326  of DNA repair genes to prevent growth arrest and promote tumor proliferation. To test this

327  hypothesis, we examined *TP53* mutation status in all six cancer types and observed that *TP53*

328  mutation frequency was the highest in pancreatic cancer patients (58%) followed by lung

329  cancer (57%), endometrial cancer (21%), liver cancer (16%), papillary renal cell (1.8%) and clear

330  cell renal cell (1.2%) (Table S7). Cancers with *TP53* mutation frequency of at least 10% were

331  selected for survival analyses. Univariate Cox regression analyses revealed that *TP53* mutation

332  status only conferred prognostic information in pancreatic (HR=1.657, P=0.044), endometrial

333  (HR=1.780, P=0.041) and liver (HR=2.603, P<0.0001) cancers but not in lung cancer (HR=1.428,

334  P=0.056) (Table 1). Cancers where *TP53* mutation offered predictive value were taken forward

335  for analyses in relation to the 10-gene signature. Cox regression analyses revealed that a

336  combination of *TP53* mutation and high 10-gene score resulted in significantly higher risk of

337  death (Table 3; Fig. 6C). Survival rates were significantly diminished in patients harboring high

338  10-gene scores and the mutant variant of *TP53* compared to those with low 10-gene scores

339  and wild-type *TP53*: liver (HR=3.876, P<0.0001), pancreas (HR=4.881, P=0.0002) and

340  endometrium (HR=3.719, P=0.00028) (Table 3; Fig. 6C). Moreover, in multivariate Cox models

341  involving TNM staging and *TP53* mutation status, the 10-gene signature remained a significant

342  prognostic factor (Table 1). This suggests that although the 10-gene signature provided

343  additional resolution in risk assessment when used in combination with *TP53* mutation status,

344  its function is independent. However, in the multivariate model *TP53* was significant only in

345  liver cancer (HR=2.085, P=0.0044), suggesting that *TP53* mutation was not independent of the

346  signature or TNM staging in pancreatic and endometrial cancers (Table 1). Overall, the results

347  suggest that defects in cell cycle checkpoint combined with augmented DNA repair ability were

348     adverse risk factors contributing to poor prognosis. Both *TP53* mutation status and 10-gene

349     scores could offer additional predictive value in risk assessment by further delineation of

350     patients into additional risk groups.

351

352 <u>Discussion and Conclusion</u>

353

354 We systematically examined the associations between the expression patterns of 138 DNA

355 repair genes in 16 cancer types and prognosis. Our pan-cancer multigenic approach revealed

356 genes that work synergistically across cancers to inform patient prognosis that would

357 otherwise remain undetected in analysis involving a single gene or a single cancer type. We

358 developed a 10-gene signature that incorporates the expression profiles of 10 highly correlated

359 DNA repair genes for use as risk predictors in six cancer types (n=2,257). This signature offers

360 a more precise discrimination of patient risk groups in these six cancers where high expression

361 of signature genes is associated with poor survival outcomes. Importantly, we demonstrated

362 that the signature can improve the prognostic discrimination of TNM when used as a combined

363 model, which is particularly useful to allow further stratification of patients within similar TNM

364 stage groups (Fig. 4).

365

366 Intrinsic differences in DNA repair machineries in cancer cells may pose a significant challenge

367 to successful therapy. Mutations in DNA repair genes allow the generation of persistent DNA

368 lesions that would otherwise be repaired. Germline mutations of DNA repair genes are linked

369 to increased genome instability and cancer risks[33] and abrogation of genes in one DNA repair

370 pathway can be compensated by another pathway[1]. *BRCA1* and *BRCA2* mutations sensitize

371 cells to PARP1 inhibition, a protein involved in the BER pathway[10]. Since *BRCA1* and *BRCA2* are

372 important for homology-directed repair, PARP1 inhibition in *BRCA1/2*-defective cells would

373 result in dysfunctional HDR and BER pathways preventing lesion repair and thus leading to

374 apoptosis[10].

375

376    In addition to genetic polymorphism, upregulation of DNA repair genes in tumors could

377    promote resistance to radiotherapy and chemotherapy as the cells would have enhanced

378    ability to repair cytotoxic lesions induced by these therapies. Overexpression of *ERCC1* involved

379    in the NER pathway in non-small-cell lung cancer is linked to poor survival in cisplatin-treated

380    patients[7]. The 1,2-d(GpG) cross-link lesion generated by cisplatin treatment is readily repaired

381    by the NER pathway, hence *ERCC1* overexpression would promote cisplatin resistance. Low

382    *MGMT* expression in astrocytoma is associated with longer survival outcomes in patients

383    treated with temozolomide[34]; an observation that is consistent with the role of *MGMT* in

384    repairing lesions caused by temozolomide thus allowing *MGMT* deficient tumor cells to

385    accumulate enough unrepairable damage. *TP53* plays essential roles in cell-cycle arrest and

386    apoptosis through the activation of checkpoint genes[29]. We show that patients with high 10-

387    gene scores that concurrently have mutant *TP53* exhibited significantly higher mortality rates

388    (Fig. 6C), suggesting that defects in cell cycle checkpoint coupled with an increase propensity

389    for DNA repair may lead to dramatically poorer outcomes.

390

391    Multiple studies have reported the associations between dysfunctional DNA repair pathways

392    and cancer, but most of these studies are restricted to investigations on a limited number of

393    genes and on one cancer at a time. One of the key advantages of our study is that it is an

394    unbiased exploration transcending the candidate-gene approach that takes into account the

395    multifaceted interplay of DNA repair genes in diverse cancer types. We rationalize that since

396    ionizing radiation and chemotherapy are the main treatment options currently available for

397    cancer patients, a molecular signature capable of discriminating patients with increased

398    expression of DNA repair genes that would benefit from adjuvant therapy through

399    pharmacological inhibition of DNA repair to overall improve therapeutic outcomes.

400

401     Tumor hypoxia is also a well-known cause of therapy resistance. A notable finding of our study

402     is that patients having both high 10-gene and hypoxia scores had significantly poorer survival

403     rates compared to those with low 10-gene and hypoxia scores (Fig. 5). Previous reports suggest

404     that low oxygen conditions may interfere with DNA damage repair. For example, hypoxia could

405     compromise HR function through decreased *RAD51* expression[35]. However, results concerning

406     the effects of hypoxia on DDR signaling have remained inconclusive. Genes associated with

407     NHEJ were reported to be downregulated under hypoxia in prostate cancer cell lines[36], while

408     hypoxia drove the upregulation of NHEJ-associated genes, *PRKDC* and *XRCC6,* in hepatoma cell

409     lines[37]. The authors proposed an interaction between *PRKDC* and the hypoxia-responsive

410     transcriptional activator, HIF-1α, hence suggesting that tumor hypoxia may lead to increase in

411     NHEJ. Tumor cells within their 3D space are subjected to differential levels of oxygen over time

412     and chronic exposures to these fluctuating conditions could result in very different biological

413     outcomes. *In vitro* studies retain a significant caveat as many hypoxia assays are carried out

414     short term using constant, predefined oxygen tensions. Although further work is needed to

415     ascertain the clinical relevance of these findings, our results clearly demonstrate that the

416     integration of hypoxia assessment in molecular stratification using the 10-gene signature

417     revealed a subset of high-risk individuals accounting for approximately 31% to 38% in each

418     cohort (Fig. 5B). Whether hypoxia could directly promote DNA damage repair *in vivo* remains

419     an open question.

420

421     We reasoned that the expression patterns of DNA repair genes would positively correlate with

422     genes involved in cell cycle progression since lesions could be repaired more effectively to

423     prevent cell cycle arrest (Fig. 6A). Enhanced DNA repair ability may also confer tumor cells with

424     growth advantage. Consistent with this hypothesis, differential expression analyses between

425     patients with high versus low 10-gene scores revealed an enrichment of ontologies involved in

426     growth stimulation as a consequence of increased DNA repair gene expression (Fig. 6B).

427     Enrichment of biological pathways involved in cell cycle, mitosis, cell division and DNA

428     replication implied that the shorter life expectancy in patients with high 10-gene scores could

429     in part be explained by an overproliferation phenotype commonly present in more aggressive

430     tumors.

431

432     In summary, we developed a prognostic signature involving DNA repair genes and confirmed

433     its utility as a powerful predictive marker for six cancer types. Although not currently afforded

434     by this work due to its retrospective nature, it will be useful to determine if the signature can

435     predict response to radiotherapy and chemotherapy in future research. While prospective

436     validation is warranted, we would expect, based on our encouraging retrospective data, that

437     the signature can guide decision making and treatment pathways. The confirmation of this

438     hypothesis by a clinical trial using the 10-gene signature to select patients that would benefit

439     from treatment with adjuvant DNA repair inhibitors could have a substantial impact on

440     treatment outcomes.

449    References

450

451    1.    Curtin NJ. DNA repair dysregulation from cancer driver to therapeutic target. *Nat Rev*

452          *Cancer*. 2012;12(12):801.

453    2.    Krishnakumar R, Kraus WL. The PARP side of the nucleus: molecular actions,

454          physiological outcomes, and clinical targets. *Mol Cell*. 2010;39(1):8-24.

455    3.    Sweasy JB, Lang T, DiMaio D. Is base excision repair a tumor suppressor mechanism?

456          *Cell cycle*. 2006;5(3):250-259.

457    4.    Martin SA, Lord CJ, Ashworth A. Therapeutic targeting of the DNA mismatch repair

458          pathway. *Clin cancer Res*. 2010:432-1078.

459    5.    Marteijn JA, Lans H, Vermeulen W, Hoeijmakers JHJ. Understanding nucleotide

460          excision repair and its roles in cancer and ageing. *Nat Rev Mol cell Biol*.

461          2014;15(7):465.

462    6.    Usanova S, Piée-Staffa A, Sied U, et al. Cisplatin sensitivity of testis tumour cells is due

463          to deficiency in interstrand-crosslink repair and low ERCC1-XPF expression. *Mol*

464          *Cancer*. 2010;9(1):248.

465   7.   Rosell R, Taron M, Barnadas A, Scagliotti G, Sarries C, Roig B. Nucleotide excision repair

466        pathways involved in Cisplatin resistance in non-small-cell lung cancer. *Cancer Control*.

467        2003;10(4):297-305.

468   8.   Martin LP, Hamilton TC, Schilder RJ. Platinum resistance: The role of DNA repair

469        pathways. *Clin Cancer Res*. 2008;14(5):1291-1295. doi:10.1158/1078-0432.CCR-07-

470        2238.

471   9.   Bryant HE, Schultz N, Thomas HD, et al. Specific killing of BRCA2-deficient tumours

472        with inhibitors of poly (ADP-ribose) polymerase. *Nature*. 2005;434(7035):913.

473   10.  Farmer H, McCabe N, Lord CJ, et al. Targeting the DNA repair defect in BRCA mutant

474        cells as a therapeutic strategy. *Nature*. 2005;434(7035):917.

475   11.  Weinstein JN, Collisson EA, Mills GB, et al. The cancer genome atlas pan-cancer

476        analysis project. *Nat Genet*. 2013;45(10):1113.

477   12.  Buffa FM, Harris AL, West CM, Miller CJ. Large meta-analysis of multiple cancers

478        reveals a common, compact and highly prognostic hypoxia metagene. *Br J Cancer*.

479        2010;102(2):428-435. doi:10.1038/sj.bjc.6605450.

480   13.  Tabas-Madrid D, Nogales-Cadenas R, Pascual-Montano A. GeneCodis3: a non-

481        redundant and modular enrichment analysis tool for functional genomics. *Nucleic

482        Acids Res*. 2012;40(W1):W478--W483.

483   14.  Kuleshov M V, Jones MR, Rouillard AD, et al. Enrichr: a comprehensive gene set

484        enrichment analysis web server 2016 update. *Nucleic Acids Res*. 2016;44(W1):W90--

485        W97.

486   15.  Chen EY, Tan CM, Kou Y, et al. Enrichr: interactive and collaborative HTML5 gene list

487        enrichment analysis tool. *BMC Bioinformatics*. 2013;14(1):128.

488   16.  Wickham H. *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York;

489        2016. http://ggplot2.org.

490    17.    Heberle H, Meirelles GV, da Silva FR, Telles GP, Minghim R. InteractiVenn: a web-based

491        tool for the analysis of sets through Venn diagrams. *BMC Bioinformatics*.

492        2015;16(1):169.

493    18.    Wood RD, Mitchell M, Sgouros J, Lindahl T. Human DNA repair genes. *Science (80- )*.

494        2001;291(5507):1284-1289.

495    19.    Chang WH, Forde D, Lai AG. Dual prognostic role for 2-oxoglutarate oxygenases in ten

496        diverse cancer types: Implications for cell cycle regulation and cell adhesion

497        maintenance. *bioRxiv*. 2018. doi:10.1101/442947.

498    20.    Chang WH, Forde D, Lai AG. A novel signature derived from immunoregulatory and

499        hypoxia genes predicts prognosis in liver and five other cancers. *J Transl Med*.

500        2019;17(1):14. doi:10.1186/s12967-019-1775-9.

501    21.    Samanta D, Gilkes DM, Chaturvedi P, Xiang L, Semenza GL. Hypoxia-inducible factors

502        are required for chemotherapy resistance of breast cancer stem cells. *Proc Natl Acad*

503        *Sci*. 2014;111(50):E5429--E5438.

504    22.    Semenza GL. Hypoxia-inducible factors: mediators of cancer progression and targets

505        for cancer therapy. *Trends Pharmacol Sci*. 2012;33(4):207-214.

506    23.    Sutherland RL, Musgrove EA. Cyclin D1 and mammary carcinoma: new insights from

507        transgenic mouse models. *Breast cancer Res*. 2001;4(1):14.

508    24.    Buckley MF, Sweeney KJ, Hamilton JA, et al. Expression and amplification of cyclin

509        genes in human breast cancer. *Oncogene*. 1993;8(8):2127-2133.

510    25.    Musgrove EA, Caldon CE, Barraclough J, Stone A, Sutherland RL. Cyclin D as a

511        therapeutic target in cancer. *Nat Rev Cancer*. 2011;11(8):558.

512    26.    Schwartz GK, Shah MA. Targeting the cell cycle: a new approach to cancer therapy. *J*

513      *Clin Oncol*. 2005;23(36):9408-9421.

514    27.    Finn RS, Crown JP, Lang I, et al. The cyclin-dependent kinase 4/6 inhibitor palbociclib in

515      combination with letrozole versus letrozole alone as first-line treatment of oestrogen

516      receptor-positive, HER2-negative, advanced breast cancer (PALOMA-1/TRIO-18): a

517      randomised phase 2 study. *Lancet Oncol*. 2015;16(1):25-35.

518    28.    Helleday T, Petermann E, Lundin C, Hodgson B, Sharma RA. DNA repair pathways as

519      targets for cancer therapy. *Nat Rev Cancer*. 2008;8(3):193.

520    29.    Sengupta S, Harris CC. p53: traffic cop at the crossroads of DNA repair and

521      recombination. *Nat Rev Mol cell Biol*. 2005;6(1):44.

522    30.    Petitjean A, Achatz MIW, Borresen-Dale AL, Hainaut P, Olivier M. TP53 mutations in

523      human cancers: functional selection and impact on cancer prognosis and outcomes.

524      *Oncogene*. 2007;26(15):2157.

525    31.    Olivier M, Langer A, Carrieri P, et al. The clinical value of somatic TP53 gene mutations

526      in 1,794 patients with breast cancer. *Clin cancer Res*. 2006;12(4):1157-1167.

527    32.    Skinner HD, Sandulache VC, Ow TJ, et al. TP53 disruptive mutations lead to head and

528      neck cancer treatment failure through inhibition of radiation-induced senescence. *Clin*

529      *cancer Res*. 2012;18(1):290-300.

530    33.    Hoeijmakers JHJ. Genome maintenance mechanisms for preventing cancer. *Nature*.

531      2001;411(6835):366.

532    34.    Hegi ME, Diserens A-C, Godard S, et al. Clinical trial substantiates the predictive value

533      of O-6-methylguanine-DNA methyltransferase promoter methylation in glioblastoma

534      patients treated with temozolomide. *Clin cancer Res*. 2004;10(6):1871-1874.

535    35.    Bindra RS, Schaffer PJ, Meng A, et al. Down-regulation of Rad51 and decreased

536      homologous recombination in hypoxic cancer cells. *Mol Cell Biol*. 2004;24(19):8504-

537        8518.

538    36.    Meng AX, Jalali F, Cuddihy A, et al. Hypoxia down-regulates DNA double strand break

539        repair gene expression in prostate cancer cells. *Radiother Oncol*. 2005;76(2):168-176.

540    37.    Um JH, Kang CD, Bae JH, et al. Association of DNA-dependent protein kinase with

541        hypoxia inducible factor-1 and its implication in resistance to anticancer drugs in

542        hypoxic tumor cells. *Exp Mol Med*. 2004;36(3):233.

543

544    <u>Figure legends</u>

545

546    **Figure 1.** Schematic representation of the study design and development of the 10-gene

547    signature. DNA repair genes from six major pathways were manually curated to generate a

548    non-redundant list containing 138 genes. Cox proportional hazards regression was employed

549    to determine the significance of each individual genes in predicting overall survival in 16 cancer

550    types. Spearman's correlation analyses revealed that five cancer types exhibited a high degree

551    of correlation in terms of their prognostic genes. Ten genes were found to be prognostic in all

552    five cancers; these genes subsequently formed the 10-gene signature. The ability of the

553    signature in predicting survival outcomes was tested using Kaplan-Meier, Cox regression and

554    receiver operating characteristic methods. The signature could predict high-risk patients in six

555    cancer types (n=2,257). Associations of the signature with tumor hypoxia, cell cycle

556    deregulation and *TP53* mutation were investigated. Potential clinical applications of the

557    signature were proposed.

558

559    **Figure 2.** Patient stratification using the 10-gene signature in six cancer types. Kaplan-Meier

560    analyses of overall survival on patients stratified into high- and low-score groups using the 10-

561    gene signature. P values were determined from the log-rank test.

562

563    **Figure 3.** Independence of the 10-gene signature over TNM staging. Kaplan-Meier analyses

564    were performed on patients categorized according to tumor TNM stages that were further

565    stratified using the 10-gene signature. The signature successfully identified patients at higher

566    risk of death in all TNM stages. P values were determined from the log-rank test. TNM: tumor,

567    node, metastasis.

568

569    **Figure 4.** Predictive performance of the 10-gene signature. Receiver operating characteristic

570    (ROC) was employed to determine the specificity and sensitivity of the signature in predicting

571    5-year overall survival in all six cancer types. ROC curves generated based on the 10-gene

572    signature, TNM staging and a combination of 10-gene signature and TNM staging were

573    depicted. AUC: area under the curve.  TNM: tumor, node, metastasis. AUCs for TNM staging

574    were in accordance with previous publications employing TCGA datasets[19,20].

575

576    **Figure 5.** Association between the 10-gene signature and tumor hypoxia. **(A)** Scatter plots

577    depict significant positive correlation between 10-gene scores and hypoxia scores in all six

578    cancers. Patients were color-coded and separated into four categories based on their 10-gene

579    and hypoxia scores. **(B)** Kaplan-Meier analyses were performed on the four patient categories

580    to assess the effects of combined relationship of hypoxia and the signature on overall survival.

581

582    **Figure 6.** Elevated DNA repair gene expression is associated with an overproliferation

583    phenotype. **(A)** Significant positive correlations between individual signature gene expression

584    and genes involved in cell cycle progression, while negative correlations were observed with

585    genes involved in cell cycle arrest. Heatmaps were generated using the R pheatmap package.

586    Cell cycle genes were depicted on the y-axis and the 10 signature genes on the x-axis. **(B)**

587    Patients were median-stratified into low- and high-score groups using the 10-gene signature

588    for differential expression analyses. Enrichment of GO and KEGG pathways associated with

589    differentially expressed genes were depicted for all six cancers. **(C)** Investigation of the

590    relationship between a gene involved in cell cycle checkpoint regulation, *TP53*, and the

591     signature. Patients were categorized into four groups based on their *TP53* mutation status and

592     10-gene scores for Kaplan-Meier analyses. P values were determined from the log-rank test.

593

594     **Table 1.** Univariate and multivariate Cox proportional hazards analyses of the 10-gene

595     signature and additional clinical risk factors associated with overall survival in six cancers.

596

597     **Table 2.** Univariate Cox proportional hazards analysis of the relation between the 10-gene

598     signature and hypoxia score.

599

600     **Table 3.** Univariate Cox proportional hazards analysis of the relation between the 10-gene

601     signature and *TP53* mutation status.

602     <u>Supplementary information</u>

603

604     **Figure S1.** Correlation analyses of 138 prognostic DNA repair genes. Spearman's correlation

605     coefficients were determined from pairwise comparisons prognostic genes from 16 cancer

606     types. Five cancers were highly correlated as shown in the blue area of the heatmap. Numbers

607     represent correlation coefficient values. Refer to Table S2 for cancer abbreviations.

608

609     **Figure S2.** Expression distribution of the ten signature genes in tumor and non-tumor samples.

610     Boxplots overlaying violin plots were used to illustrate tumor and non-tumor distribution in six

611     cancers: **(A)** clear cell renal cell, **(B)** papillary renal cell, **(C)** liver, **(D)** lung, **(E)** pancreas and **(F)**

612     endometrium. Nonparametric Mann-Whitney-Wilcoxon tests were employed to determine

613     whether there were significant differences in expression distributions. Asterisks represent

614     significant P values: * < 0.05, *** < 0.0001.

615

616     **Figure S3**. Venn diagram depicts a six-way comparison of the differentially expressed genes (-

617     $1 > \log_2$ fold-change > 1, P<0.05) identified from high-score versus low-score patients in all six

618     cancers. Numbers in parentheses represent the number of differentially expressed genes in

619     each cancer. The Venn intersection of all cancers indicated that 87 genes were common.

620

621     **Table S1.** List of 138 DNA repair genes and associated pathways.

622

623     **Table S2.** Description of TCGA cancer cohorts.

624

625     **Table S3.** Univariate Cox proportional hazards analysis of the 138 genes in 16 cancers.

626

627     **Table S4.** Differentially expressed genes between high- and low-score patient groups in six

628     cancers.

629

630     **Table S5.** List of 87 differentially expressed genes that are common in all six cancers.

631

632     **Table S6.** Enrichr transcription factor protein-protein interaction analysis of the 87

633     differentially expressed genes.

634

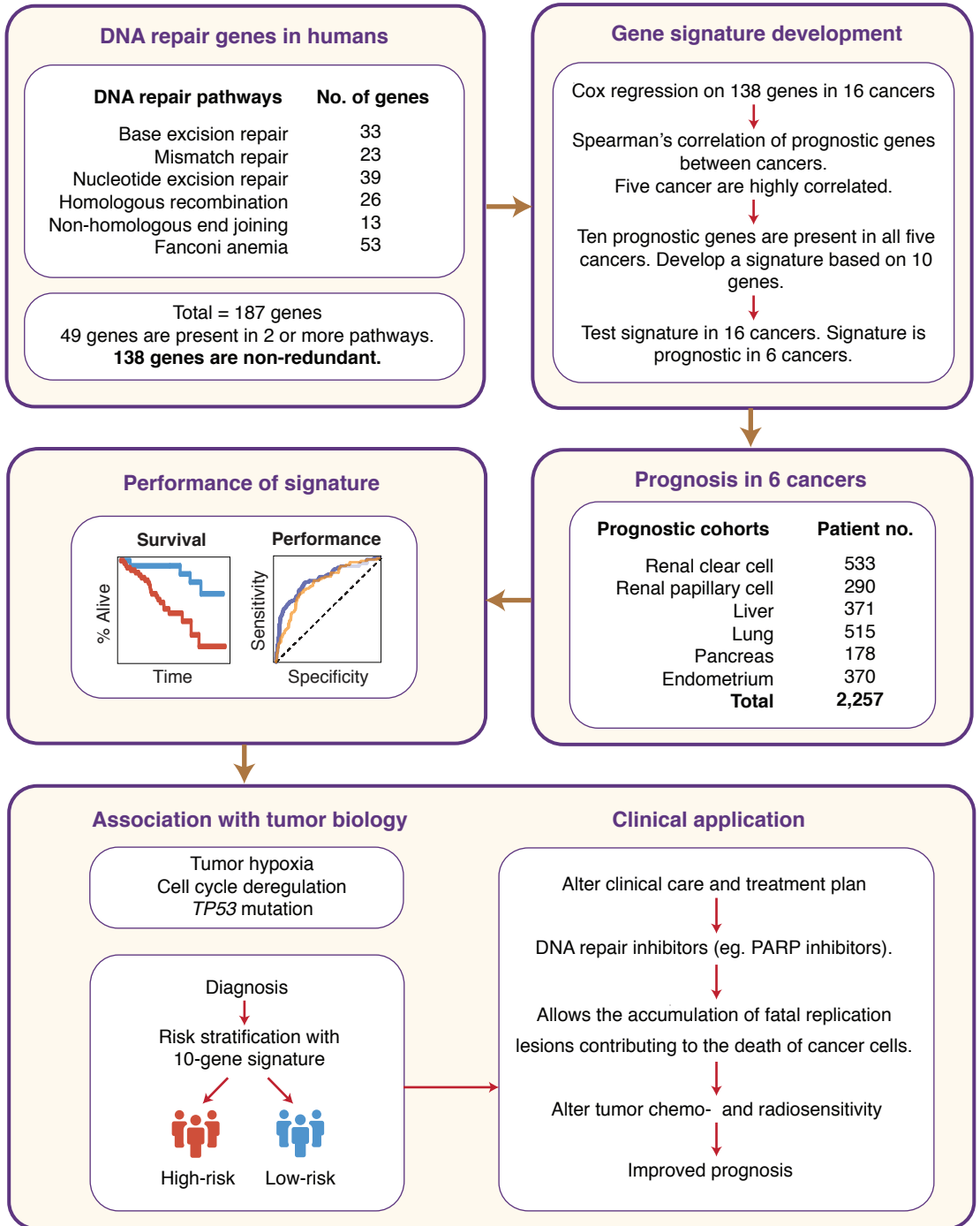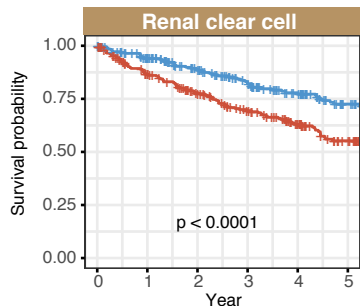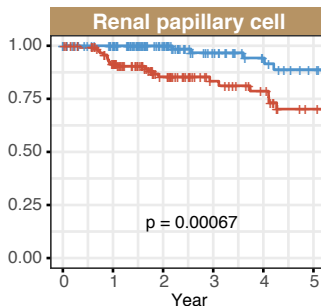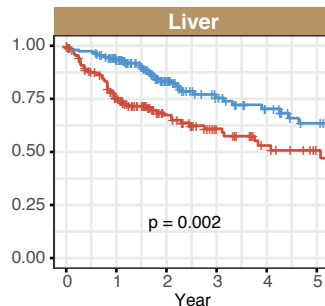635     **Table S7.** *TP53* mutation analysis in liver, pancreatic, endometrial and lung cancers.

# Figure 1

## DNA repair genes in humans

| DNA repair pathways | No. of genes |
| --- | --- |
| Base excision repair | 33 |
| Mismatch repair | 23 |
| Nucleotide excision repair | 39 |
| Homologous recombination | 26 |
| Non-homologous end joining | 13 |
| Fanconi anemia | 53 |

Total = 187 genes
49 genes are present in 2 or more pathways.
**138 genes are non-redundant.**

## Gene signature development

Cox regression on 138 genes in 16 cancers

Spearman's correlation of prognostic genes between cancers.
Five cancer are highly correlated.

Ten prognostic genes are present in all five cancers. Develop a signature based on 10 genes.

Test signature in 16 cancers. Signature is prognostic in 6 cancers.

## Performance of signature



## Prognosis in 6 cancers

| Prognostic cohorts | Patient no. |
| --- | --- |
| Renal clear cell | 533 |
| Renal papillary cell | 290 |
| Liver | 371 |
| Lung | 515 |
| Pancreas | 178 |
| Endometrium | 370 |
| **Total** | **2,257** |

## Association with tumor biology

Tumor hypoxia
Cell cycle deregulation
*TP53* mutation

Diagnosis

Risk stratification with 10-gene signature



High-risk          Low-risk

## Clinical application

Alter clinical care and treatment plan

DNA repair inhibitors (eg. PARP inhibitors).

Allows the accumulation of fatal replication lesions contributing to the death of cancer cells.

Alter tumor chemo- and radiosensitivity

Improved prognosis

**Figure 2**



| Renal clear cell | | | | | |
|---|---|---|---|---|---|
| **Number at risk** | | | | | |
| Low score | 261 | 225 | 182 | 150 | 117 | 80 |
| High score | 263 | 210 | 173 | 139 | 104 | 72 |

p < 0.0001

| Renal papillary cell | | | | | |
|---|---|---|---|---|---|
| **Number at risk** | | | | | |
| Low score | 129 | 114 | 76 | 48 | 35 | 27 |
| High score | 127 | 102 | 62 | 39 | 30 | 19 |

p = 0.00067

| Liver | | | | | |
|---|---|---|---|---|---|
| **Number at risk** | | | | | |
| Low score | 157 | 133 | 74 | 48 | 36 | 25 |
| High score | 156 | 101 | 57 | 36 | 23 | 14 |

p = 0.002

| Lung | | | | | |
|---|---|---|---|---|---|
| **Number at risk** | | | | | |
| Low score | 210 | 181 | 103 | 64 | 38 | 28 |
| High score | 209 | 163 | 86 | 48 | 29 | 21 |

p = 0.028

| Pancreas | | | | | |
|---|---|---|---|---|---|
| **Number at risk** | | | | | |
| Low score | 74 | 54 | 23 | 15 | 9 | 6 |
| High score | 75 | 49 | 11 | 4 | 1 | 1 |

p = 0.00013

| Endometrium | | | | | |
|---|---|---|---|---|---|
| **Number at risk** | | | | | |
| Low score | 184 | 164 | 122 | 88 | 70 | 55 |
| High score | 184 | 157 | 119 | 82 | 63 | 41 |

p = 0.00063

# Figure 3

# Figure 4



**Renal clear cell**

| Classifier | AUC |
|---|---|
| TNM | 0.716 |
| Signature | 0.711 |
| Signature + TNM | 0.792 |

**Renal papillary cell**

| Classifier | AUC |
|---|---|
| TNM | 0.640 |
| Signature | 0.832 |
| Signature + TNM | 0.868 |

**Liver**

| Classifier | AUC |
|---|---|
| TNM | 0.697 |
| Signature | 0.689 |
| Signature + TNM | 0.751 |

**Lung**

| Classifier | AUC |
|---|---|
| TNM | 0.663 |
| Signature | 0.639 |
| Signature + TNM | 0.693 |

**Pancreas**

| Classifier | AUC |
|---|---|
| TNM | 0.593 |
| Signature | 0.697 |
| Signature + TNM | 0.698 |

**Endometrium**

| Classifier | AUC |
|---|---|
| TNM | 0.674 |
| Signature | 0.700 |
| Signature + TNM | 0.764 |

# Figure 5

# Figure 6

**Figure S1**

# Figure S2

# Figure S3