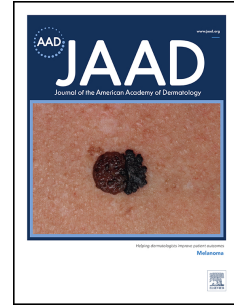# Journal Pre-proof

Transcriptomic analysis of cutaneous squamous cell carcinoma reveals a multi-gene prognostic signature associated with metastasis

Jun Wang, PhD, Catherine A. Harwood, MD PhD, Emma Bailey, PhD, Findlay Bewicke-Copley, PhD, Chinedu Anthony Anene, PhD, Jason Thomson, MD, Mah Jabeen Qamar, Rhiannon Laban, Craig Nourse, PhD, Christina Schoenherr, PhD, Mairi Treanor-Taylor, MD, Eugene Healy, MB PhD, Chester Lai, BM PhD, Paul Craig, Colin Moyes, William Rickaby, Joanne Martin, PhD, Charlotte Proby, MD, Gareth J. Inman, PhD, Irene M. Leigh, CBE, DSc

Please cite this article as: Wang J, Harwood CA, Bailey E, Bewicke-Copley F, Anene CA, Thomson J, Qamar MJ, Laban R, Nourse C, Schoenherr C, Treanor-Taylor M, Healy E, Lai C, Craig P, Moyes C, Rickaby W, Martin J, Proby C, Inman GJ, Leigh IM, Transcriptomic analysis of cutaneous squamous cell carcinoma reveals a multi-gene prognostic signature associated with metastasis, *Journal of the American Academy of Dermatology* (2023), doi: https://doi.org/10.1016/j.jaad.2023.08.012.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

1 **Article type:** Original article

2 **Transcriptomic analysis of cutaneous squamous cell carcinoma reveals a multi-gene**
3 **prognostic signature associated with metastasis**

4 Jun Wang, PhD,[a] Catherine A. Harwood, MD PhD,[a,b] Emma Bailey, PhD,[a] Findlay Bewicke-
5 Copley, PhD,[a] Chinedu Anthony Anene, PhD,[a,c] Jason Thomson, MD,[a,b] Mah Jabeen Qamar,[a,b]
6 Rhiannon Laban,[a,b] Craig Nourse, PhD,[d] Christina Schoenherr, PhD,[d] Mairi Treanor-Taylor,
7 MD,[d,e] Eugene Healy, MB PhD,[f,g] Chester Lai, BM PhD,[f,g] Paul Craig,[h] Colin Moyes,[i] William
8 Rickaby,[j] Joanne Martin, PhD,[a] Charlotte Proby, MD,[k] Gareth J. Inman, PhD,[d,e] Irene M. Leigh,
9 CBE, DSc[a]
10
11 a. Faculty of Medicine and Dentistry, Queen Mary University of London, London E1 1BB
12 b. Department of Dermatology, Royal London Hospital, Barts Health NHS Trust E1 IBB
13 c. Centre for Biomedical Science Research, School of Clinical and Applied Sciences, Leeds
14 Beckett University, Leeds, UK
15 d. Cancer Research UK Beatson Institute, Bearsden Glasgow, G61 1BD
16 e. School of Cancer Sciences, University of Glasgow, Bearsden G61 1QH
17 f. Dermatopharmacology, University of Southampton, Southampton General Hospital, SO16
18 6YD
19 g. Dermatology, University Hospital Southampton NHS Foundation Trust, SO16 6YD
20 h. Cellular Pathology, Gloucestershire Hospitals NHS Foundation Trust, Cheltenham General
21 Hospital, Cheltenham GL537AN
22 i. Queen Elizabeth University Hospital Glasgow G51 4TF
23 j. University College London NHS Trust, London W1T 4EU
24 k. Molecular and Clinical Medicine, School of Medicine, University of Dundee, DD 1 4HN

25 **Corresponding author**
26 Jun Wang PhD, Barts Cancer Institute, Faculty of Medicine and Dentistry, Queen Mary
27 University of London, Charterhouse Square, London EC1M 6BQ. E-mail:
28 j.a.wang@qmul.ac.uk

29 **Word count: 2,487**

30 **Number of references: 36**

31 **Number of Figures: 3; Number of Tables: 2**

32 **Supplemental Material:** https://data.mendeley.com/datasets/z77kdgddm9/1

35 **Disclosure:** Dr Harwood is honoraria for advisory boards from Sanofi/Regeneron, Almirall,
36 AmLo Biosciences, Incanthera, Leo Pharma, L'Oreal. Other authors have no disclosures.

37 **Patient consent** Not applicable

38 **ABSTRACT**
39 *Background:* Metastasis of cutaneous squamous cell carcinoma (cSCC) is uncommon. Current
40 staging methods are reported to have sub-optimal performances in metastasis prediction.
41 Accurate identification of patients with tumours at high risk of metastasis would have a
42 significant impact on management.
43 *Objective:* To develop a robust and validated gene expression profile (GEP) signature for
44 predicting primary cSCC metastatic risk using an unbiased whole transcriptome discovery-
45 driven approach.
46 *Methods:* Archival formalin-fixed paraffin-embedded primary cSCC with perilesional normal
47 tissue from 237 immunocompetent patients (151 non-metastasising and 86 metastasising)
48 were collected retrospectively from four centres. TempO-seq was used to probe the whole
49 transcriptome and machine learning algorithms were applied to derive predictive signatures,
50 with a 3:1 split for training and testing datasets.
51 *Results:* A 20-gene prognostic model was developed and validated, with an accuracy of 86.0%,
52 sensitivity of 85.7%, specificity of 86.1%, and positive predictive value of 78.3% in the testing
53 set, providing more stable, accurate prediction than pathological staging systems. A linear
54 predictor was also developed, significantly correlating with metastatic risk.
55 *Limitations:* This was a retrospective 4-centre study and larger prospective multicentre
56 studies are now required.
57 *Conclusion:* The 20-gene signature prediction is accurate, with the potential to be
58 incorporated into clinical workflows for cSCC.

59 *Key words:*

60 Cutaneous squamous cell carcinoma; Metastasis; Prognosis; Transcriptomics; Machine
61 learning; Risk stratification

63 **CAPSULE SUMMARY**

64 • A 20-gene expression profile signature derived from clinical archival tissue using an
65 unbiased whole-transcriptome approach showed superior performance for predicting
66 metastatic risks for primary cutaneous squamous cell carcinoma (cSCC).
67 • This prognostic signature could significantly improve risk stratification, identifying
68 patients with high-risk cSCC who may benefit from adjuvant treatment and reducing
69 overtreatment for patients with low-risk cSCC.

70

71

72 **BACKGROUND**

73

74 Cutaneous squamous cell carcinoma (cSCC) is the commonest form of skin cancer with
75 metastatic potential and incidence and mortality are rising (1-4). Although the frequency of
76 metastasis arising from cSCC is relatively low at 2-5%, the sheer number of cases represents
77 a significant disease burden. Current management could be improved by more accurately
78 identifying tumours most likely to metastasise, targeting adjuvant therapy and intense clinical
79 supervision programmes to those at highest risk, whilst reducing unnecessary interventions
80 for people with low-risk tumours.

81

82 Multiple histopathological staging classifications for cSCC are available although reported to
83 be suboptimal in predicting poor outcomes (5,6). Recent studies suggest that genomic and
84 transcriptomic signatures may improve risk prediction for primary cSCC progression (7-10).
85 Using whole exome sequencing data, we previously identified 16 high-risk and 6 low-risk
86 specific significantly mutated genes (9). More recently, a 40-gene expression profiling (GEP)
87 signature based on candidate genes identified by a combination of literature review and
88 discovery efforts, was developed to predict metastatic risk (Castle Biosciences, Inc
89 Friendswood, Texas) (11,12). A positive predictive value (PPV) of 60% was achieved for the
90 highest-risk tumours, with overall sensitivity, specificity and PPV for differentiating Class 2
91 (high-) and Class 1 (low-risk) cSCC of 65.4%, 68.8%, and 28.8%, respectively (11). A completely
92 unbiased discovery-driven approach using information from the whole genome and
93 transcriptome to identify prognostic gene signatures is currently lacking. Such an approach
94 may also uncover key molecular mechanisms underpinning disease progression and
95 metastatic risk.

96

97 To develop a validated prognostic signature in an unbiased manner, we assembled a
98 multicentre cohort of primary cSCC archival tissue from 237 patients with known clinical
99 outcomes (no metastasis over 3 years, n=151; metastasis, n=86). Whole transcriptomic data
100 were generated from tumour and perilesional normal skin. A range of machine learning (ML)
101 techniques was applied and a 20-gene GEP model was developed which displayed a high level
102 of accuracy in differentiating metastasising and non-metastasising primary cSCC. A linear
103 predictor based on the 20-gene GEP was then developed to further aid the implementation
104 of the GEP signature for risk stratification in clinical practice. Ultimately, use of this GEP to
105 guide management decisions may significantly improve patient management for this
106 common cancer.

107

108

109

## METHODS

**Ethical approval and sample identification**

This study was approved as IRAS project 266559 (Diagnostic marker panel development for progression in skin cancer, PERMEDID). Four collaborating pathology centres identified consecutive patients with primary cSCC which had metastasised, or primary cSCC which had not metastasised within 3 years (**Table I).** Immunosuppressed patients were excluded. Formalin fixed paraffin embedded (FFPE) sections were reviewed by an expert dermatopathologist and tumour and perilesional normal skin marked for subsequent analysis (see Supplemental Materials).

**Pathology review and pathological tumour staging**

Haematoxylin and eosin (H&E) stained sections were digitally scanned by Leica scanner and Aperio software. Images were reviewed centrally by two expert dermatopathologists and primary tumours typed, graded and histologically staged using Union for International Cancer Control (UICC)-8 and Brigham and Women's Hospital (BWH) classifications.

**Transcriptomics investigation**

Transcriptomic analysis was performed using the TempO-Seq whole-protein coding transcriptome platform with a proprietary processing pipeline (Bioclavis Ltd, Glasgow, UK) (13). Data pre-processing and normalisation were performed using limma R package (14). Batch effect was removed using the ComBat package (15). Differential expression (DE) analysis using limma was performed between clinical groups, followed by gene set over-representation and gene set enrichment analysis (GSEA) using DAVID (16) and clusterProfiler (17).

**Gene signature analysis using machine learning**

To derive a set of genes that could distinguish two groups (i.e., metastasising versus non-metastasising cSCC), the caret R package (18) was used for machine learning (ML) analysis. A range of ML techniques were used and compared (Supplemental Materials). We randomly split the samples into training (75%) and testing (25%) sets. Starting with an initial set of genes in the training set (i.e. all DE genes from the DE analysis comparing metastasising and non-metastatic cSCC), the best performing set of genes for each ML algorithm (i.e., feature selection) was determined using the Recursive Feature Elimination procedure, with 10-fold repeated cross validation of five repeats. A final model for each ML algorithm was then trained using the final selected number of genes with 10-fold repeated cross validation of ten repeats, and used to predict the two classes in the testing set. The performances of predictions were measured using accuracy, precision, along with sensitivity and specificity, positive predive value (PPV) and negative predictive value (NPV).

149    A weighted linear predictor was generated for each sample based on the expression of the

150    final set of genes in the model and their fold changes in the DE analysis (see Supplemental

151    Materials) Linear predictors were compared between clinical groups and correlated with

152    classes. The area under the ROC curve (AUC) was calculated using the pROC package (19).

153

**RESULTS**

154

**Clinicopathologic characteristics**

155

156    Demographic details of patients and histologic features of primary cSCC are presented in

157    **Table I**.

158

159    **Transcriptomic analysis between primary cSCC groups**

160    Gene expression profiles (GEP) of 19,072 genes across a total of 433 samples were sufficiently

161    profiled for analysis. Four sample groups were compared; cSCC tumour from metastasising

162    (n=84) and non-metastasising (n=146) cSCC, and matched perilesional normal skin from

163    metastasising (n=71) and non-metastasising (n=132) cSCC (Supplemental Table I). Principal

164    component analysis based on genes across all samples showed a clear separation between

165    cSCC and perilesional normal skin samples from metastasising and non-metastasising cSCC

166    (Supplemental Fig 1). Differential gene expression analysis revealed that 1,038 genes were

167    upregulated and 236 genes downregulated in metastasising cSCC compared to non-

168    metastasising cSCC (absolute $\log_2$ fold change >1 and adjusted *p*-value <0.05). The gene set

169    over-representation test showed keratinisation, B-cell receptor (BCR), innate immune

170    response, cell cycle, DNA replication and DNA repair were highly over-represented in the DE

171    genes (hypergeometric test *q*<0.05, Supplemental Fig 2A). Over-representation analysis

172    against cellular signatures showed that signatures associated with neural progenitor,

173    endothelial and cancer stem cells were highly enriched within the DE genes (Supplemental

174    Fig 2B), suggesting that cell differentiation is a key factor distinguishing the two cSCC groups.

175    GSEA against MSigDB canonical pathways further suggested that cell cycle related, DNA

176    replication and repair, and immune pathways (BCR regulation, interferon and interleukin-12

177    signalling), were all significantly upregulated in metastasising cSCC, while formation of the

178    cornified envelope, keratinisation, and many metabolism pathways (sphingolipid,

179    triglyceride, creatine and fatty acid metabolism) were significantly downregulated (**Figure 1**).

180

181    Normal perilesional samples from metastasising and non-metastasising primary cSCC were

182    also compared. GSEA indicated many immune pathways (such as BCR and T cell receptor

183    signalling, Fc gamma receptor activation, and chemokine receptor binding) and cell cycle

184    related pathways (synthesis, replication and repair of DNA) were significantly upregulated in

185    perilesional skin samples from metastasising tumours (Supplementary Table II).

186

187    **Development of the 20-GEP prognostic signature**

188  To identify a smaller set of genes that were predictive for primary cSCC metastasis, a range of
189  ML classification algorithms were applied after splitting the primary cSCC samples into
190  training and validation sets. A 20-gene model derived from K-nearest neighbours (KNN) was
191  identified (Supplemental Table III) which provided the best performance in differentiating the
192  two cSCC groups in the validation set (n=57: 36 non-metastasising; 21 metastasising), with an
193  accuracy of 86.0% (95% confidence interval 74.2-93.7%), a sensitivity of 85.7% and a
194  specificity of 86.1% (**Table II**). Patients predicted as high-risk of metastasis by the 20-GEP
195  signature (n=23) had significantly worse metastasis-free survival (MFS) rates than those
196  predicted as low-risk (n=34) (3-year MFS, 91.7% for low-risk versus 21.7% for high-risk) (**Figure
197  2**). In this 20-gene GEP model, 18 genes were upregulated in non-metastasising cSCC and 2
198  genes (*MDK* and *STMN1)* were upregulated in metastasising cSCC (Supplemental Table III,
199  Supplemental Fig 3). Functional annotation of the 20 genes suggested the significant
200  enrichment in the signatures from keratinisation, GnRH, oxytocin, Ras and MAPK signalling
201  pathways (hypergeometric test, *p*<0.01). Using the same ML procedure based on perilesional
202  normal skin samples, a 22-gene KNN model was also developed with an accuracy of 64.0%
203  (95% CI: 49.2-77.1%), sensitivity of 41.2%, and specificity of 75.8% (**Table II**).

204

205  **Prognostic accuracy of the 20-GEP test compared to pathological staging classifications**
206  Using the Royal College of Pathologists dataset for histopathological reporting of primary
207  invasive cSCC, tumours were staged by both UICC-8 TNM and BWH T-staging classifications
208  after central consensus histopathological review. Prognostic metrics for UICC-8 (low T1/T2 vs.
209  high T3/T4) and BWH (low T1/T2a vs high T2b/T3) staging showed performance with an
210  accuracy of 85.4% for both systems in the validation set, compared to 86.0% for the 20-GEP
211  signature (**Table II**). Performance of BWH T-staging based on original pathology reports
212  without central consensus review (BWH v1), was marginally inferior in predicting metastasis,
213  with an accuracy of 81.8%. This was largely due to differences between the scoring of poor
214  differentiation after central review compared to the original report (**Table I**, Supplemental
215  Table IV).

216  The 20-GEP signature showed strong correlations with staging for risk prediction in the
217  validation set. Of 23 metastasising cases predicted by the 20-GEP test, 21/23 (91.3%) were
218  T2b/T3 by BWH staging versus 15/23 (65.2%) UICC-8 T3/4. Of 32 non-metastasising cSCC
219  predicted by the 20-GEP, 26/32 were T1/T2a by BWH and 26/32 were UICC-8 T1/T2 (81.3%).
220  Accuracy of the histology staging systems dropped to 81.1% and 76.5% for BWH and UICC8,
221  respectively, when the whole cohort (n=237) was considered (**Table II**).

222  **Generation of a linear predictor for metastatic prediction**
223  To further enhance the potential clinical application of the 20-GEP signature, a linear
224  predictor for metastasis combining the expression values and fold-changes of these 20 genes
225  in the DE analysis was generated: the higher the linear predictor value, the higher the risk of
226  developing metastasis. The previously reported 40-GEP (11) stratifies tumours into 3 classes

227    of risk (low, high, highest), whereas a linear predictor allows a more detailed assessment of
228    risk that can be used alongside pathological risk factors to influence clinical management. The
229    linear predictor had a very high correlation with metastatic risk, with an area under the ROC
230    curve (AUC) of 0.85 (95% CI, 0.80-0.91) and 0.88 (95% CI, 0.78-0.99) for the training and
231    validation (testing) sets, respectively (**Figure 3**). In comparison, the KNN binary classification
232    model (i.e., yes or no for metastasis prediction) had an AUC of 0.86 (0.76-0.96). As expected,
233    the linear predictor was significantly higher in metastasising versus non-metastasising cSCC
234    in both training and testing sets (Wilcoxon rank sum test, *p*<0.0001, Supplemental Fig 4).

235    Finally, the linear predictors across both tumour and perilesional skin for both metastasising
236    and non-metastasising cSCC were compared (Supplemental Fig 5). There was no difference in
237    linear predictors between non-metastasising cSCC and both normal adjacent groups.
238    However, linear predictors increased significantly for metastasising cSCC compared to other
239    groups (*p*<0.0001), suggesting that our linear predictor was only associated with
240    metastasising primary tumours.

241

242    **DISCUSSION**

243    This study reports a 20-GEP signature that predicts metastatic risk of primary cSCC. It was
244    developed and validated in a UK cohort of 237 primary cSCC from immunocompetent
245    individuals using archival FFPE tissue in which whole-transcriptome analysis with an unbiased
246    discovery approach was performed. The 20-GEP signature achieved an accuracy of 86.0%, a
247    negative predictive value of 91.2% and a positive predictive value of 78.3% for predicting
248    metastasis in the validation set (n=57). A linear predictor to facilitate potential clinical use of
249    the 20-GEP was created based on the expression and fold changes of signature genes and had
250    an AUC of 0.88. UICC-8 TNM and BWH pathological staging systems performed unexpectedly
251    well in risk prediction compared with previous reports. Nonetheless, the 20-GEP remained
252    overall the most stable and accurate predictor of metastatic risk, and in contrast to histology,
253    the GEP signature is unbiased and not dependent on human evaluation and interpretation.
254
255    There appeared to be a strong association between the 20-GEP and keratinisation. Key
256    keratinisation genes, such as *LCE1C*, *LCE2B/C*, *LCE3C* and *CDSN*, were all significantly
257    downregulated in metastasising primary cSCC as were two genes involved in alpha-Linolenic
258    acid and ether lipid metabolism (*PLA2G4E/F*), consistent with our GSEA results. Only two
259    genes, *STMN1* and *MDK*, were significantly upregulated in metastasising samples. STMN1 a
260    microtubule-destabilising protein, regulates the dynamics of microtubules and cell cycle
261    progress (20). Its high expression is associated with poor prognosis in oesophageal (ESCC),
262    lung (LUSC) and oral SCC (21-23). In ESCC and LUSC, it was reported to promote cell
263    proliferation, migration, chemoradiation resistance (21,22,24), and is strongly associated with
264    lymph node metastasis in ESCC (25,26). Midkine (MDK), a heparin-binding growth factor, is

7

265 also associated with cancer progression, drug resistance and a tolerogenic and immune-
266 resistant state (27-30). A recent study showed that MDK was highly expressed by stem-like
267 tumour cells and led to mTOR inhibition persistence and an immune-suppressive
268 microenvironment (31). MDK represents an interesting therapeutic target for advanced cSCC.
269

270 Currently, clinical pathways determining treatment plans for patients with cSCC use
271 clinicopathological staging systems. In practice, the predictive accuracy of staging systems for
272 primary cSCC can vary significantly across reported studies (11, 32-35). Factors possibly
273 accounting for the variability in pathology staging include non-standardised reporting of high-
274 risk features (particularly poor differentiation and perineural invasion); problems defining the
275 state of differentiation of an individual tumour; and variable practice in the use of Mohs'
276 surgery which may affect detection of high-risk features and lead to understaging (11). In our
277 study, careful central review by two highly experienced dermatopathologists adhering to the
278 Royal College of Pathologists dataset led to a much higher performance of pathology staging
279 systems than previously published. This highlights the need for a more objective grading
280 system such as that used worldwide in breast carcinoma (36).
281

282 Additional strengths of our study include an unbiased discovery-driven approach using the
283 whole transcriptome of FFPE clinical samples to develop a prognostic signature suitable for
284 routine clinical use. We also excluded immunosuppressed patients as iatrogenic and disease-
285 associated immunosuppression is an important risk factor for poor outcomes in cSCC and
286 variations in immune status and effects of immunosuppressive drugs are likely to impact on
287 the transcriptome. Excluding confounding factors due to immunosuppression may have
288 permitted generation of a more metastasis-specific gene signature of greater use for risk
289 prediction. More work is needed to test our 20-gene signature in other patient populations,
290 such as those with darker skin and in immunosuppressive populations.
291

292 The retrospective nature of this study was a limitation and, although consecutive eligible
293 primary cSCC were enrolled at each centre, the possibility of some bias relating to patient and
294 sample selection cannot be excluded. The study size for the validation set was also a limitation
295 and further validation will require larger, prospective studies (5).
296

297 In conclusion, we have used an unbiased discovery-driven approach to generate a promising
298 candidate 20-GEP prognostic signature for cSCC metastasis. The GEP not only represents a
299 novel and potentially clinically applicable prognostic tool but has also provided biological
300 insights into the process of metastasis and potential therapeutic targets. In addition, there
301 are biological and genomic mechanisms common to cSCC across different tissue types and
302 this signature may provide further insights into common differentiation and stem-like
303 pathways underpinning these SCCs. Further prospective evaluation is now underway to
304 confirm clinical utility of this GEP in management of primary cSCC.
305

306 ***Abbreviations used:***

307 AJCC: American Joint Committee on Cancer

308 UICC: Union for International Cancer Control

309 BWH: Brigham and Women's Hospital

310 cSCC: cutaneous squamous cell carcinoma

311 GEP: gene-expression profile

312 DE: differential expression

313 GSEA: Gene set enrichment analysis

314 HR: hazard ratio

315 LR: likelihood ratio

316 NPV: negative predictive value

317 PPV: positive predictive value

318 KNN: K-nearest neighbourhood

319 BCR: B-cell receptor

320

321

322 **Author contributions**

323

324 IML, JW and GJI conceived and designed the study, and acquired the funding. CAH, CP, EH,

325 CL, CM, JT and MJQ recruited the patient cohort and collected the clinical data. PC, WR and

326 JM performed the tumour grading and histology analysis. EB, JW, FBC and CAA performed the

327 bioinformatics and machine learning data analysis. CN, CS and MTT performed the TempO-

328 seq experiment. JW, IML, GJI and CAH supervised the study, analysed and interpreted the

329 data, and drafted the manuscript. All authors critically revised the manuscript and approved

330 the final version to be submitted.

331

332

343

344

### References

1. Venables ZC, Nijsten T, Wong KF, Autier P, Broggio J, Deas A, Harwood CA, Hollestein LM, Langan SM, Morgan E, Proby CM, Rashbass J, Leigh IM. Epidemiology of basal and cutaneous squamous cell carcinoma in the U.K. 2013-15: a cohort study. Br J Dermatol. 2019 Sep;181(3):474-482.

2. Venables ZC, Autier P, Nijsten T, Wong KF, Langan SM, Rous B, Broggio J, Harwood C, Henson K, Proby CM, Rashbass J, Leigh IM. Nationwide Incidence of Metastatic Cutaneous Squamous Cell Carcinoma in England. JAMA Dermatol. 2019 Mar 1;155(3):298-306.

3. Kwiatkowska M, Ahmed S, Ardern-Jones MR, Bhatti LA, Bleiker TO, Gavin A, Hussain S, Huws DW, Irvine L, Langan SM, Millington GWM, Mitchell H, Murphy R, Paley L, Proby CM, Thomson CS, Thomas R, Turner C, Vernon S, Venables ZC. A summary of the updated report on the incidence and epidemiological trends of keratinocyte cancers in the UK 2013-2018. Br J Dermatol. 2022 Feb;186(2):367-369.

4. Detailed statistics from the Get Data Out (>Skin Tumours) programme, National Disease Registration Service, UK. https://www.cancerdata.nhs.uk/getdataout/skin

5. Venables ZC, Tokez S, Hollestein LM, Mooyaart AL, van den Bos RR, Rous B, Leigh IM, Nijsten T, Wakkee M. Validation of four cutaneous squamous cell carcinoma staging systems using nationwide data. Br J Dermatol. 2022 May;186(5):835-842.

6. Tokez S, Venables ZC, Hollestein LM, Qi H, Bramer EM, Rentroia-Pacheco B, van den Bos RR, Rous B, Leigh IM, Nijsten T, Mooyaart AL, Wakkee M. Risk factors for metastatic cutaneous squamous cell carcinoma: Refinement and replication based on 2 nationwide nested case-control studies. J Am Acad Dermatol. 2022 Jul;87(1):64-71.

7. South AP, Purdie KJ, Watt SA, Haldenby S, denBreem N, Dimon M, Arron ST, McHugh A, Xue DJ, Jasbani HS, Dayal HS, Proby CM, Harwood CA, Leigh IM, Prevalent NOTCH1 mutations in squamous cell carcinogenesis are an early event J Invest Dermatol 2014 :134; 2630-8

8. Cammameri P, Rose AM, Vincent DF, Wang J, Nagano A, Libertini S, Ridgeqay RA, Athineos D, Coates PJ, McHugh A, Pourreyron C, Dayal JH, Larsson J, Weidlich S, Spender LC, Sapkota GP, Purdie KJ, Proby CM, Harwood CA, Leigh IM, Clevers H, Barker N, Karlsson S, Pritchard C, Marais R, Chelala C, South AP, Sansom OJ, Inman GJ. Inactivation of TGFb receptors in stem cell drives cutaneous squamous cell carcinoma. Nature Commun 2016 25: 12493

9. Inman GJ, Wang AJ, Ai N, Alexandreev L, Chelala C, Stratton M, Harwood CA, Sherwood V, Proby CM, Leigh IM. The genomic landscape of cutaneous squamous cell carcinoma from immunosuppressed and immunocompetent patients reveals common drivers and a novel mutational signature associated with chronic azathioprine exposure. Nat Commun 2018 Sept 10 9(1) 3667.

10. Thomson J, Bewicke-Copley F, Anene CA, Gulati A, Nagano A, Purdie K, Inman GJ, Proby CM, Leigh IM, Harwood CA,Wang J. The genomic landscape of actinic keratosis. J Invest Dermatol 2021 Jul;141(7):1664-1674.e7.

390

391 11. Wysong A, Newman JG, Covington KR, Kurley SJ, Ibrahim SF, Farberg AS, Bar A, Cleaver NJ,
392 Somani AK, Panther D, Brodland DG, Zitelli J, Toyohara J, Maher IA, Xia Y, Bibee K, Griego R,
393 Rigel DS, Meldi Plasseraud K, Estrada S, Sholl LM, Johnson C, Cook RW, Schmults CD, Arron ST.
394 Validation of a 40-gene expression profile test to predict metastatic risk in localized high-risk
395 cutaneous squamous cell carcinoma. J Am Acad Dermatol. 2021 Feb;84(2):361-369.

396

397 12. Ibrahim SF, Kasprzak JM, Hall MA, Fitzgerald AL, Siegel JJ, Kurley SJ, Covington KR, Goldberg
398 MS, Farberg AS, Trotter SC, Reed K, Brodland DG, Koyfman SA, Somani AK, Arron ST, Wysong
399 A. Enhanced metastatic risk assessment in cutaneous squamous cell carcinoma with the 40-
400 gene expression profile test. Future Oncol. 2022 Mar;18(7):833-847. doi: 10.2217/fon-2021-
401 1277.

402

403 13. Yeakley JM, Shepard PJ, Goyena DE, VanSteenhouse HC, McComb JD, Seligmann BE. A
404 trichostatin A expression signature identified by TempO-Seq targeted whole transcriptome
405 profiling. PloS One. 2017;12(5):e0178302.

406

407 14. Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, Smyth GK. Limma powers differential
408 expression analyses for RNA-sequencing and microarray studies. Nucleic Acids Res. 2015 Apr
409 20;43(7):e47.

410

411 15. Leek JT, Johnson WE, Parker HS, Fertig EJ, Jaffe AE, Zhang Y, Storey JD, Torres LC. Sva:
412 Surrogate Variable Analysis. 2022 R package version 3.44.0.

413

414 16. Huang da W, Sherman BT, Lempicki RA. Systematic and integrative analysis of large gene
415 lists using DAVID bioinformatics resources. Nat Protoc. 2009;4(1):44-57.

416

417 17. Wu T, Hu E, Xu S, Chen M, Guo P, Dai Z, Feng T, Zhou L, Tang W, Zhan L, Fu X, Liu S, Bo X,
418 Yu G. clusterProfiler 4.0: A universal enrichment tool for interpreting omics data. Innovation
419 (Camb). 2021 Jul 1;2(3):100141.

420

421 18. Kuhn M. Building predictive models in R using the caret package. J Stat Software.
422 2008;28:1–26.

423

424 19. Robin X, Turck N, Hainard A, Tiberti N, Lisacek F, Sanchez JC, Müller M. pROC: an open-
425 source package for R and S+ to analyze and compare ROC curves. BMC Bioinformatics. 2011
426 Mar 17;12:77.

427

428 20. Rubin CI, Atweh GF. The role of stathmin in the regulation of the cell cycle. J Cell Biochem.
429 2004 Oct 1;93(2):242-50.

430

431 21. Suzuki S, Yokobori T, Altan B, Hara K, Ozawa D, Tanaka N, Sakai M, Sano A, Sohda M, Bao
432 H, Fukuchi M, Miyazaki T, Kaira K, Asao T, Kuwano H. High stathmin 1 expression is associated

433 with poor prognosis and chemoradiation resistance in esophageal squamous cell carcinoma.
434 Int J Oncol. 2017 Apr;50(4):1184-1190.
435

436 22. Bao P, Yokobori T, Altan B, Iijima M, Azuma Y, Onozato R, Yajima T, Watanabe A, Mogi A,
437 Shimizu K, Nagashima T, Ohtaki Y, Obayashi K, Nakazawa S, Bai T, Kawabata-Iwakawa R, Asao
438 T, Kaira K, Nishiyama M, Kuwano H. High STMN1 Expression is Associated with Cancer
439 Progression and Chemo-Resistance in Lung Squamous Cell Carcinoma. Ann Surg Oncol. 2017
440 Dec;24(13):4017-4024.
441

442 23. Ma HL, Jin SF, Tao WJ, Zhang ML, Zhang ZY. Overexpression of stathmin/oncoprotein 18
443 correlates with poorer prognosis and interacts with p53 in oral squamous cell carcinoma. J
444 Craniomaxillofac Surg. 2016 Oct;44(10):1725-1732.
445

446 24. Ni PZ, He JZ, Wu ZY, Ji X, Chen LQ, Xu XE, Liao LD, Wu JY, Li EM, Xu LY. Overexpression of
447 Stathmin 1 correlates with poor prognosis and promotes cell migration and proliferation in
448 oesophageal squamous cell carcinoma. Oncol Rep. 2017 Dec;38(6):3608-3618.
449

450 25. Li J, Qi Z, Hu YP, Wang YX. Possible biomarkers for predicting lymph node metastasis of
451 esophageal squamous cell carcinoma: a review. J Int Med Res. 2019 Feb;47(2):544-556.
452

453 26. Jiang W, Huang S, Song L, Wang Z. STMN1, a prognostic predictor of esophageal squamous
454 cell carcinoma, is a marker of the activation of the PI3K pathway. Oncol Rep. 2018
455 Feb;39(2):834-842.
456

457 27. Filippou PS, Karagiannis GS, Constantinidou A. Midkine (MDK) growth factor: a key player
458 in cancer progression and a promising therapeutic target. Oncogene. 2020 Mar;39(10):2040-
459 2054.
460

461 28. Yu X, Zhou Z, Tang S, Zhang K, Peng X, Zhou P, Zhang M, Shen L, Yang L. MDK induces
462 temozolomide resistance in glioblastoma by promoting cancer stem-like properties. Am J
463 Cancer Res. 2022 Oct 15;12(10):4825-4839.
464

465 29. Olmeda D, Cerezo-Wallis D, Riveiro-Falkenbach E, Pennacchi PC, Contreras-Alcalde M, Ibarz
466 N, Cifdaloz M, Catena X, Calvo TG, Cañón E, Alonso-Curbelo D, Suarez J, Osterloh L, Graña O,
467 Mulero F, Megías D, Cañamero M, Martínez-Torrecuadrada JL, Mondal C, Di Martino J, Lora D,
468 Martinez-Corral I, Bravo-Cordero JJ, Muñoz J, Puig S, Ortiz-Romero P, Rodriguez-Peralto JL,
469 Ortega S, Soengas MS. Whole-body imaging of lymphovascular niches identifies pre-metastatic
470 roles of midkine. Nature. 2017 Jun 28;546(7660):676-680.
471

472 30. Cerezo-Wallis D, Contreras-Alcalde M, Troulé K, Catena X, Mucientes C, Calvo TG, Cañón E,
473 Tejedo C, Pennacchi PC, Hogan S, Kölblinger P, Tejero H, Chen AX, Ibarz N, Graña-Castro O,

474 Martinez L, Muñoz J, Ortiz-Romero P, Rodriguez-Peralto JL, Gómez-López G, Al-Shahrour F,
475 Rabadán R, Levesque MP, Olmeda D, Soengas MS. Midkine rewires the melanoma
476 microenvironment toward a tolerogenic and immune-resistant state. Nat Med. 2020
477 Dec;26(12):1865-1877.
478

479 31. Tang Y, Kwiatkowski DJ, Henske EP. Midkine expression by stem-like tumor cells drives
480 persistence to mTOR inhibition and an immune-suppressive microenvironment. Nat Commun.
481 2022 Aug 26;13(1):5018.
482

483 32. Jambusaria-Pahlajani A, Kanetsky PA, Karia PS, Hwang WT, Gelfand JM, Whalen FM,
484 Elenitsas R, Xu X, Schmults CD. Evaluation of AJCC tumor staging for cutaneous squamous cell
485 carcinoma and a proposed alternative tumor staging system. JAMA Dermatol. 2013
486 Apr;149(4):402-10. doi: 10.1001/jamadermatol.2013.2456. PMID: 23325457.
487

488 33. Karia PS, Jambusaria-Pahlajani A, Harrington DP, Murphy GF, Qureshi AA, Schmults CD.
489 Evaluation of American Joint Committee on Cancer, International Union Against Cancer, and
490 Brigham and Women's Hospital tumor staging for cutaneous squamous cell carcinoma. J Clin
491 Oncol. 2014;32(4):327-334. doi: 10.1200/JCO.2012.48.5326
492

493 34. Karia PS, Morgan FC, Califano JA, Schmults CD. Comparison of tumor classifications for
494 cutaneous squamous cell carcinoma of the head and neck in the 7th vs 8th edition of the AJCC
495 Cancer Staging Manual. JAMA Dermatol. 2018; 154:175-181.
496

497 35. Ruiz ES, Karia PS, Besaw R, Schmults CD. Performance of the American Joint Committee on
498 Cancer Staging Manual, 8th Edition vs the Brigham and Women's Hospital Tumor Classification
499 System for Cutaneous Squamous Cell Carcinoma. *JAMA Dermatol.* 2019;155(7):819–825.
500 doi:10.1001/jamadermatol.2019.0032
501

502 36. Blamey RW, Ellis IO, Pinder SE, Lee AH, Macmillan RD, Morgan DA et al. Survival of invasive
503 breast cancer according to the Nottingham Prognostic Index in cases diagnosed in 1990–1999.
504 Eur J Cancer 2007;43:1548–1555.
505

506

507

508

509

510

511 **Figures Legend**

512 **Fig 1.** Normalised enrichment scores (NES) of the top dysregulated canonical pathways
513 between metastasising and non-metastasising cSCC. Pathways with positive NES (in red) were
514 upregulated while pathways with negative NES (in blue) were downregulated in metastasising
515 compared to non-metastasising primary cSCC.

516 **Fig 2.** Kaplan-Meier analysis of the 20-GEP prognostic test and outcomes in terms of
517 metastasis free survival in the validation dataset. No. at risk in the follow-up was shown in
518 the table below.

519 **Fig 3.** Area under the receiver operating characteristic curve (AUC) of the performance of
520 linear predictors correlating with the metastatic incidences. Linear predictors were produced
521 based on the 20-GEP signature, and both training and testing data sets were included in the
522 calculation. AUC and 95% confidence interval were shown.

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537    **Table I.** Clinicopathologic details of patients and primary cSCC samples

| Feature | All (n=237) | No metastasis (n=151)* | Metastasis (n=86)** | P value |
|---|---|---|---|---|
| Age, y, median (range) | 80 (39-100) | 78 (39-100) | 80 (64-93) | .57 |
| Male, n (%) | 142 (60) | 90 (60) | 52 (60) | 1 |
| Located on head and neck, n (%)*** | 155 (65) | 91 (61) | 64 (74) | .033 |
| Tumour diameter, cm, mean (range) # | 1.85 (0.18-9) | 1.31 (0.18-4.1) | 2.82 (1.6-9) | <.0001 |
| Tumour thickness, mm, mean (range) ## | 3.94 (0.2-26.7) | 2.96 (0.2-13) | 5.65 (0.3-26.7) | <.0001 |
| Poorly differentiated, n (%) | 115 (48.3) | 47 (30.9) | 68 (79.1) | <.0001 |
| Clark level > V (beyond fat), n (%)§ | 43 (18.6) | 10 (6.7) | 33 (40.2) | <.0001 |
| PNI, n (%)¶ | | | | .0004 |
| Present (≥ 0.1mm) | 20 (8.6) | 8 (5.3) | 12 (14.6) | |
| Present (<0.1mm or unknown) | 11 (4.7) | 3 (1.99) | 8 (9.8) | |
| Not present | 202 (86.7) | 140 (92.7) | 62 (75.6) | |
| Lymphovascular invasion∞ | 15 (6.5) | 1 (0.66) | 14 (17.5) | <.0001 |
| UICC T stage, n (%)§§ | | | | <.0001 |
| T1 | 134 (59.3) | 115 (78.8) | 19 (23.75) | |
| T2 | 25 (11.1) | 11 (7.5) | 14 (17.5) | |
| T3 | 67 (29.6) | 20 (13.7) | 47 (58.75) | |
| T4 | - | - | - | |
| BWH T stage, n (%)¶¶ | | | | <.0001 |
| T1 | 86 (37.7) | 84 (56.75) | 2 (2.5) | |
| T2a | 65 (28.5) | 44 (29.7) | 21 (26.25) | |
| T2b | 71 (31.1) | 20 (13.5) | 51 (63.75) | |
| T3 | 6 (2.6) | - | 6 (7.5) | |
| | | | | |

538

539    *Total number of primary cSCC which did not metastasise =152 (one patient had 2

540    separate primary cSCCs); median follow-up was 76 months

541    ** median time from primary cSCC to metastasis was 9.9 months

542    *** Location not recorded for 2 cSCCs (both non-metastasising)

543    # not available for 10 cSCC (5 non-metastasising and 5 metastasising)

544    ## not available for 15 cSCC (10 non-metastasising and 5 metastasising)

545    § Invasion through or beyond subcutaneous fat: not available for 7 cSCC (3 non-

546    metastasising and 4 metastasising)

547    ¶ not available for 5 cSCC (1 non-metastasising and 4 metastasising cSCC)

548    ∞ Lymphovascular invasion not available for 6 cSCC (all metastasising)

549    §§ not available for 12 cSCC (6 non- metastasising and 6 metastasising)

550    ¶¶ not available for 10 cSCC (4 non-metastasising and 6 metastasising)
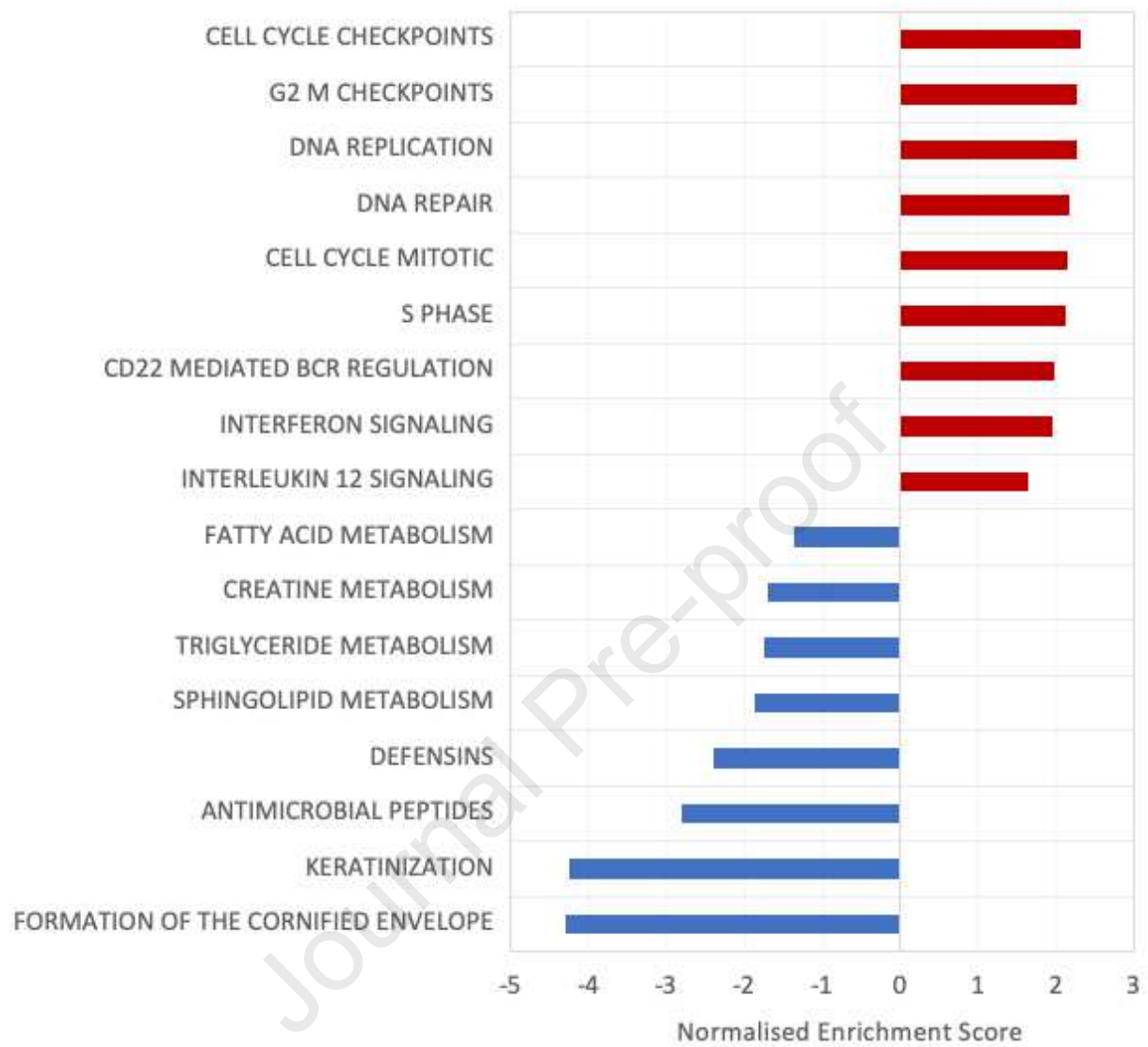
551

552

553

554 **Table II.** Accuracy of the prediction of metastatic risks of the 20-GEP signature and other risk
555 assessment methods (n=57).

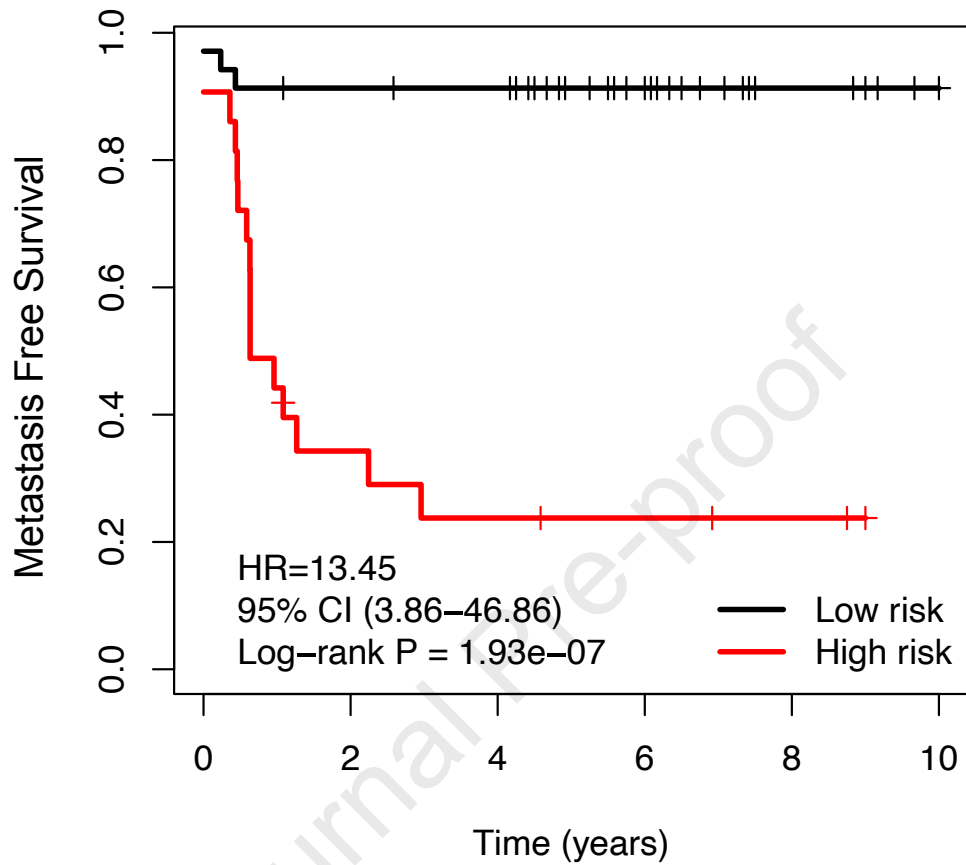| Classifier | Accuracy% | Sensitivity% | Specificity% | PPV% | NPV% | +LR | -LR |
|------------|-----------|--------------|--------------|------|------|------|------|
| 20-GEP | 86.0 | 85.7 | 86.1 | 78.3 | 91.2 | 6.17 | 0.17 |
| UICC-8 | 85.4 | 81.0 | 88.2 | 81.0 | 88.2 | 6.88 | 0.22 |
| BWH | 85.4 | 95.2 | 79.4 | 74.1 | 96.4 | 4.63 | 0.06 |
| BWH v1 | 81.8 | 76.2 | 85.3 | 76.2 | 85.3 | 5.18 | 0.28 |
| 22-GEP* | 64.0 | 41.2 | 75.8 | 46.7 | 71.4 | 1.70 | 0.78 |
| UICC-8** | 76.5 | 58.8 | 86.3 | 70.1 | 79.2 | 4.29 | 0.48 |
| BWH** | 81.1 | 71.2 | 86.5 | 74.0 | 84.8 | 5.27 | 0.33 |

556 UICC: Union for International Cancer Control; BWH, Brigham and Women's Hospital Staging
557 System after the central review; BWH v1: derived from original pathology reports before
558 central pathology review; GEP, gene expression profile; PPV: Positive Predictive Value; NPV:
559 Negative Predictive Value; +LR: Positive Likelihood Ratio; -LR: Negative Likelihood Ratio.

560 22-GEP* was derived from normal adjacent samples only.

561 ** Statistics were derived from the whole cohort (n=237)

# UK−cSCC cohort



HR=13.45
95% CI (3.86−46.86)
Log−rank P = 1.93e−07

Low risk
High risk

| # at risk | 0 | 2 | 4 | 6 | 8 | 10 |
|---|---|---|---|---|---|---|
| Low-risk | 34 | 30 | 29 | 17 | 6 | 2 |
| High-risk | 21 | 6 | 4 | 3 | 2 | 0 |

# ROC curves of linear predictor