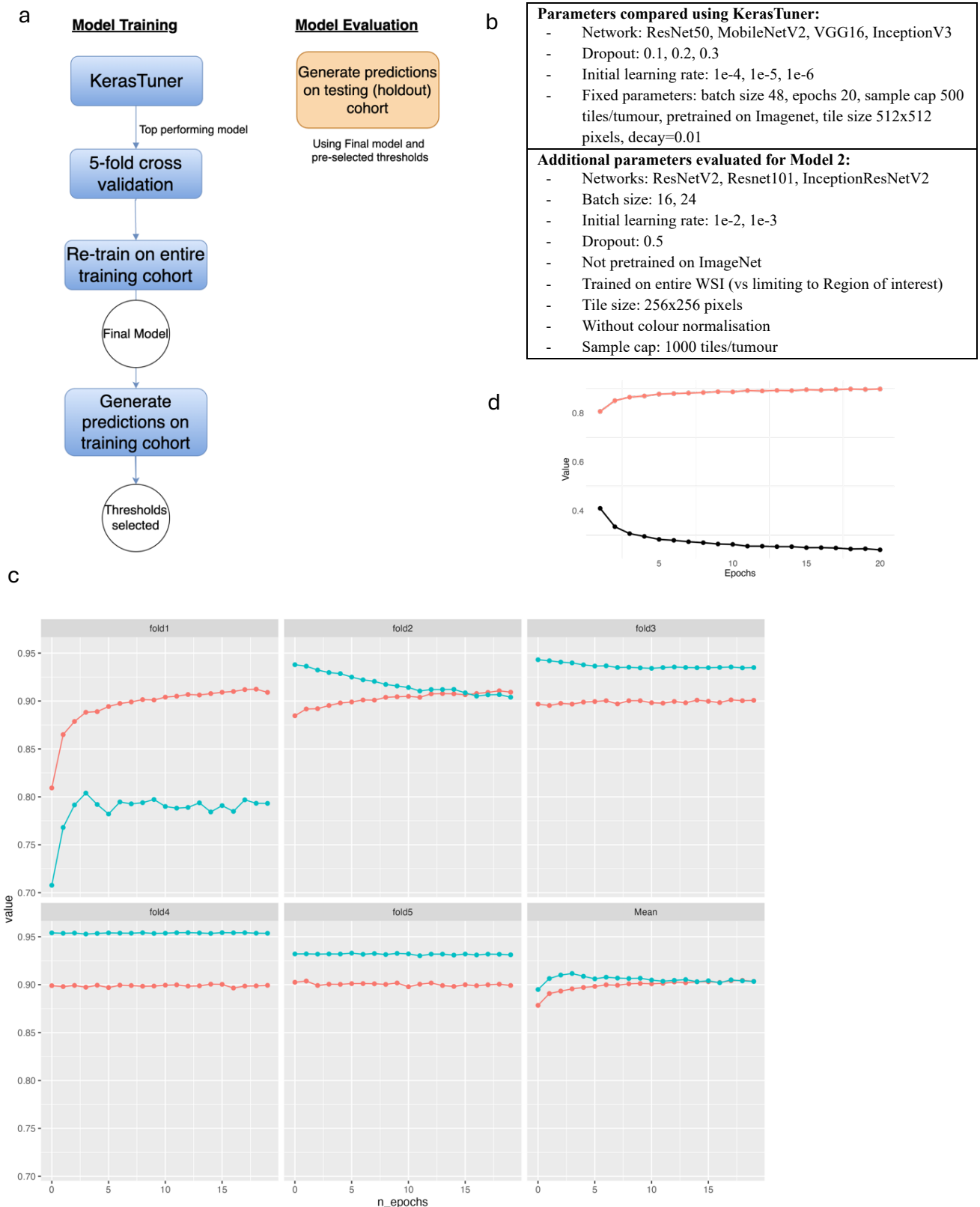


# Supplementary Material

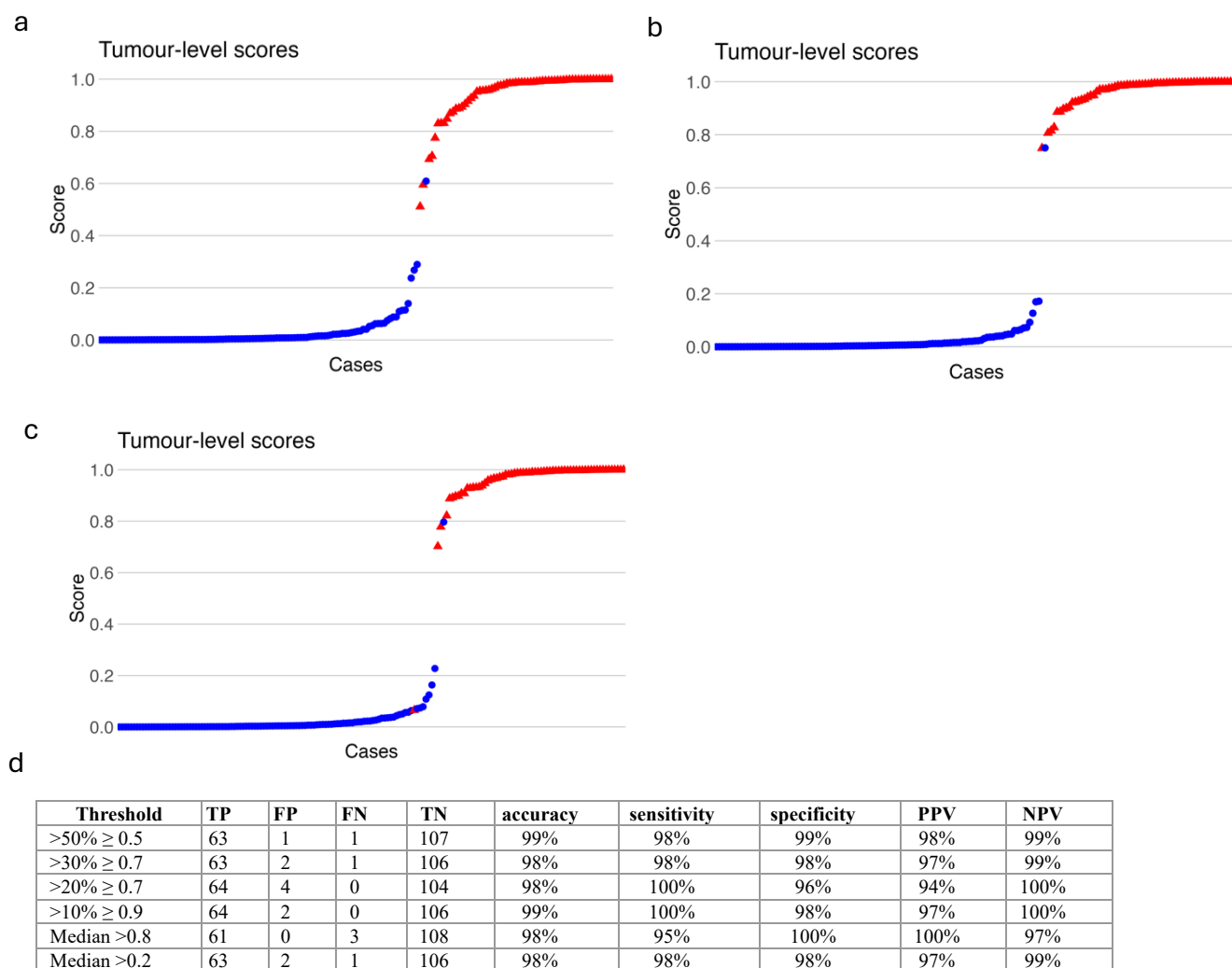
Supplementary Figure 1. Model training and evaluation



**a** Diagram summarising the model training and evaluation pipeline followed for both models. Firstly, hypertuning using the KerasTuner was used for selection of the ideal model architecture and hyperparameters, through systematic comparison of widely-used convolutional neural network backbones and parameters. The top performing model was then evaluated using 5-fold cross-validation. A final model was then re-trained on the entire

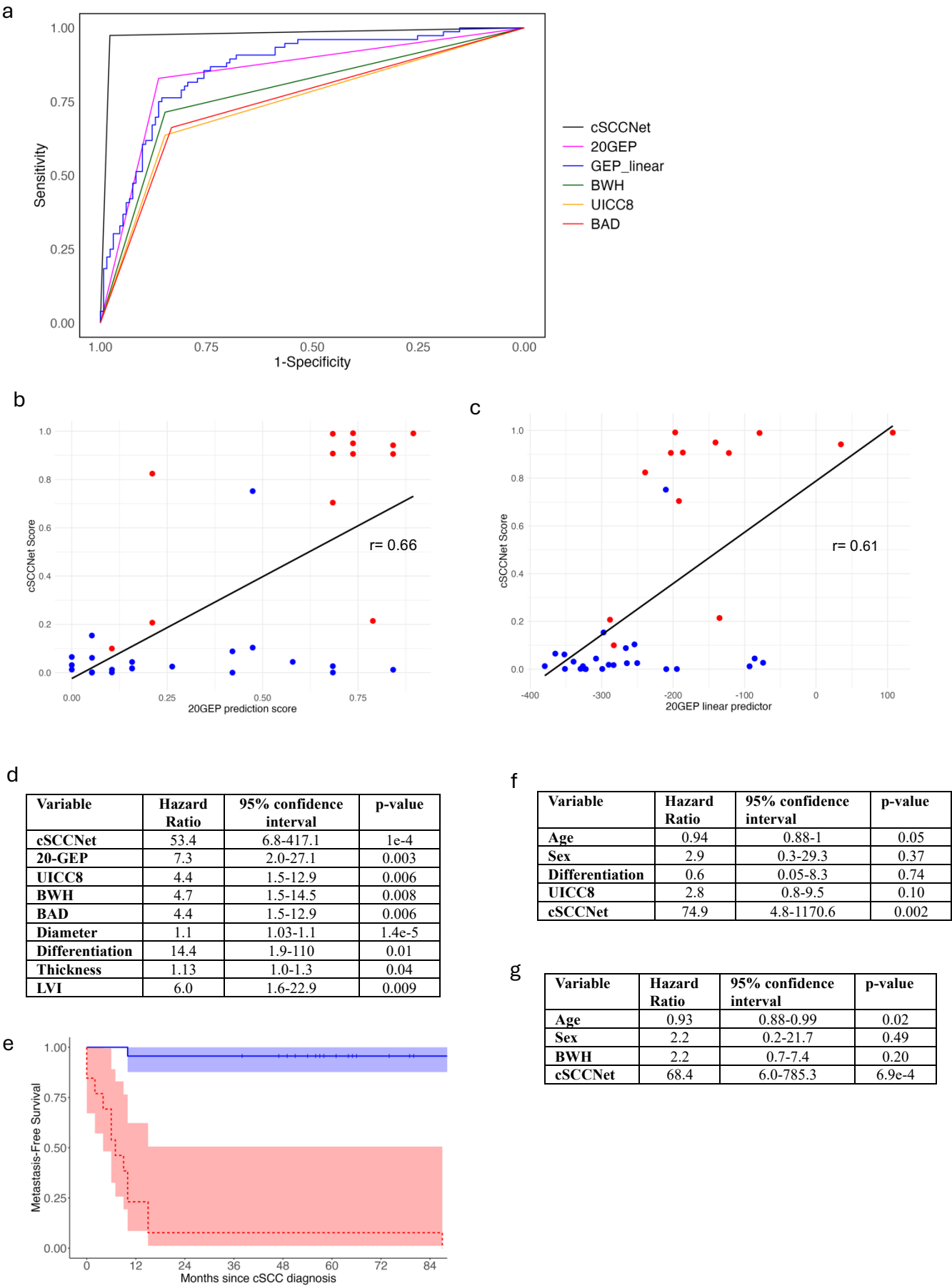
training cohort. To select a threshold for binary classification, the final model was used to generate predictions on the training cohort. Model evaluation was performed by generating predictions on the hold-out testing cohort and comparing model predictions to ground truth. **b** List of parameters compared using the KerasTuner and additional important parameters evaluated individually. **c** Five-fold cross validation curves for Model 2, with training accuracy in red and validation accuracy in blue. **d** Model 2 final training curve, with training accuracy in red and training loss in black. After 20 epochs, the model reached accuracy 0.90 and loss 0.24.

## Supplementary Figure 2. Model 2 threshold selection



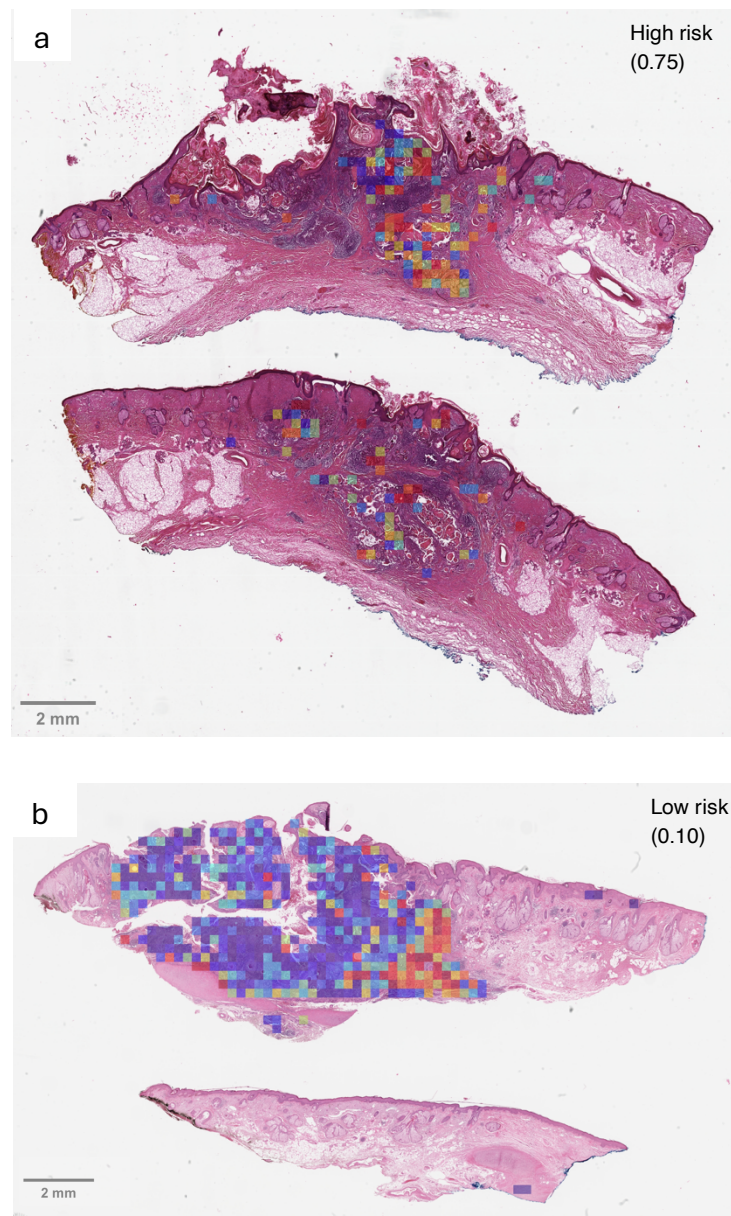
Results given for 172 samples in the training cohort. **a** Scatterplot showing Model 2 median tile scores for metastasising cases as red triangles and non-metastasising cases as blue circles. All tiles within the pathologist-annotated regions of interest (ROI) were included. **b** To improve the separation between low- and high-risk tumours, borderline tiles (with scores 0.3-0.7) were excluded. Scatterplot showing Model 2 median tile scores, for tiles within pathologist ROI and after excluding tiles with borderline scores. A median score  $>0.2$  was selected as the definition for 'high-risk' cSCC, based on graphical representations and accuracy statistics. **c** Both models were then used in series. Scatterplot showing Model 2 median tile scores, for tiles selected by Model 1 and after excluding borderline tiles. **d** Comparison of various Model 2 thresholds in predicting risk of cSCC metastasis when both models were used in series and after excluding borderline tiles. Thresholds: '>50%  $\geq 0.5$ ' (>50% of tiles have score  $\geq 0.5$ ), '>30%  $\geq 0.7$ ' (>30% of tiles have score  $\geq 0.7$ ), '>20%  $\geq 0.7$ ' (>20% of tiles have score  $\geq 0.7$ ), '>10%  $\geq 0.9$ ' (>10% of tiles have score  $\geq 0.9$ ), median tile score  $>0.8$  and median tile score  $>0.2$ . FN: false negatives; FP: false positives; NPV: negative predictive value; PPV: positive predictive value; TN: true negatives; TP: true positives.

Supplementary Figure 3. Model 2 evaluation



**a** Receiver operating characteristic (ROC) curves are shown for cSCCNet (black), the 20-gene expression profile (20-GEP) prediction score based on k-nearest neighbours analysis (20GEP, purple) and linear predictor (GEP\_linear, blue), Brigham and Women's Hospital classification (BWH, green), 8th edition Union for International Cancer Control classification (UICC8, yellow), and British Association of Dermatologists' cSCC guidelines (BAD, red) on the entire cohort (n = 212). Of note, this includes the training cases for the cSCCNet and 20-GEP models, which might bias their performance. The area under the ROC curves (AUC) and their 95% confidence intervals were: 0.98 (0.95-1) for cSCCNet, 0.85 (0.79-0.90) for 20-GEP prediction score, 0.86 (0.81-0.91) for the 20-GEP linear predictor, 0.74 (0.68-0.80) for UICC8, 0.78 (0.72-0.84) BWH, and 0.74 (0.69-0.81) for BAD. For 37 cases in the testing cohort where complete data were available, the Pearson correlation between the cSCCNet score and 20-GEP test showed a moderate positive correlation: **b** 0.66 ( $p = 9e-6$ ) for the 20-GEP prediction score and **c** 0.61 ( $p = 7e-5$ ) for the 20-GEP linear predictor score. **d** On univariate analysis, features predictive of metastasis (Wald test,  $p < 0.05$ ) in the testing cohort included the cSCCNet classification, 20-GEP, UICC8, BWH, BAD Very High risk grade, tumour diameter, differentiation, thickness, and presence of lymphovascular invasion (LVI). Age, sex, site of primary cSCC and presence of perineural invasion were not statistically significant in the testing cohort; however, all were significant ( $p < 0.05$ ) when assessed in the entire cohort (n=212), suggesting an impact of sample size. Margin status was not a significant predictor of outcome. **e** Kaplan-Meier curve showing metastasis-free survival after cSCC diagnosis, stratified by cSCCNet prediction, with high-risk cases in red and low-risk cases in blue. **f** On multivariate analysis for predicting the risk of metastasis in the testing cohort (n=35), cSCCNet was a significant predictor, independently of age, sex, tumour differentiation, or UICC8. **g** cSCCNet was also a significant predictor when multivariate analysis was repeated with BWH. As differentiation is already included within BWH, it was not included as a separate variable.

#### Supplementary Figure 4. Review of incorrect cases



Heatmaps for the testing cohort cases misclassified by cSCCNet ( $n=2/40$ ), with Model 2 tile scores converted to colour using a blue to red scale, for scores 0-1 (low to high-risk). The tumour-level aggregate scores are on the top right corner of each case, with scores  $>0.20$  representing 'high-risk' tumours. **a** Non-metastasising scalp cSCC classified as high-risk by the cSCCNet model and by staging criteria (8th edition Union for International Cancer Control classification, UICC8, and Brigham and Women's Hospital classification, BWH). It is poorly differentiated and invades beyond subcutis. Heatmaps show that Model 1 had failed to select  $>60\%$  of the ROI; the small number of tiles passed to Model 1 contained poorly differentiated carcinoma. **b** A metastasising pinna cSCC with incomplete excision margins was classified as low-risk by the cSCCNet model. The majority of the tumour was moderately-differentiated with good keratinisation; however, there was extension beyond cartilage. It was high grade on staging criteria (AJCC8 T3/BWH T2b). A small area of poorly differentiated carcinoma was present and was classified as 'high-risk' by the model.

**Supplementary Table 1. Baseline clinicopathological characteristics**

	<b>Training cohort n = 172<sup>a</sup></b>	<b>Testing cohort n=40</b>
Metastasising cases	64 (37)	14 (35)
Non-metastasising cases	108 (63)	26 (65)
Age, years	80 (71-84)	82 (75-86)
<b>Sex</b>		
Male	114 (67)	29 (73)
Female	57 (33)	11 (28)
<b>Site</b>		
Head and Neck	112 (65)	24 (60)
<b>Tumour diameter</b>		
Median, mm	15 (10-23)	15 (10-25)
≥20 mm	58/167 (35)	15 (38)
<b>Differentiation</b>		
Poorly differentiated	87 (51)	21 (53)
Moderately differentiated	63 (37)	13 (33)
Well differentiated	22 (13)	6 (15)
<b>Thickness, mm</b>		
Median	3	3
>6mm	24/170 (14)	3/36 (8)
<b>Invasion to</b>		
Dermis	89 (53)	24 (60)
Subcutis	45 (27)	9 (23)
Beyond subcutaneous fat	35 (21)	7 (18)
Perineural invasion	22/170 (13)	3/39 (8)
Lymphovascular invasion	11/167 (7)	3/38 (8)
<b>Follow-up, months</b>	64 (44-80)	65 (54-80)
<b>UICC8</b>		
pT1	96 (57)	22 (56)
pT2	16 (9)	5 (13)
pT3	57 (34)	12 (31)
<b>BWH</b>		
T1	61 (36)	12 (31)
T2a	45 (27)	15 (38)
T2b	58 (34)	10 (26)
T3	5 (3)	2 (5)
<b>BAD</b>		
LR	45 (27)	11 (28)
HR	64 (38)	15 (38)
VHR	60 (36)	13 (33)

This table does not include the 15 cases used to train Model 1 only. Numbers in brackets are percentages or interquartile range, as appropriate. There were no statistically significant differences ( $p > 0.05$ ) between the training and testing cohorts, using the Mann-Whitney U test for continuous variables and Fisher's exact test for categorical variables. Although all cSCC were treated with wide local excision, 30/210 tumours (21/172 training and 9/40 testing) had incomplete excision margins, and all had either re-excision or adjuvant radiotherapy, except one case in the training cohort which had already metastasized at the time of presentation of the primary lesion. <sup>a</sup>172 tumours from 171 patients. BAD: British Association of Dermatologists' cSCC guidelines; BWH: Brigham and Women's Hospital classification; UICC8: the 8th edition Union for International Cancer Control classification.