

Classifying subtypes and predicting survival of renal cell carcinoma using histopathology image-based deep learning

Fei Wang

China Pharmaceutical University

Yuanzhe Geng

China Pharmaceutical University

Ting Wang

Nanjing University

Ke Zhao

The Affiliated Jiangyin Hospital of Southeast University Medical College

Bin Xu

Affiliated Zhongda Hospital of Southeast University

Dachuan Zhang

The Third Affiliated Hospital of Soochow University

Hongzhou Cai

The Affiliated Cancer Hospital of Nanjing Medical University & Jiangsu Cancer Hospital & Jiangsu Institute of Cancer Research

Qinbo Yuan

Huaiyin People's Hospital of Huaian City

Hang Lu

China Pharmaceutical University

Yue Zhang

China Pharmaceutical University

Xu Li

China Pharmaceutical University

Yangyang Sun

China Pharmaceutical University

Hang Gong

China Pharmaceutical University

Raphael N Alolga

China Pharmaceutical University

Xiangshan Fan

Nanjing Drum Tower Hospital

Gaoxiang Ma

China Pharmaceutical University <https://orcid.org/0000-0002-1774-6445>

Lian-Wen Qi (✉ Qilw@cpu.edu.cn)

China Pharmaceutical University <https://orcid.org/0000-0003-0728-4475>

Article

Keywords: Renal cell carcinoma, Deep learning, Histopathology, Diagnosis, Prognosis

Posted Date: May 19th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-533678/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Classifying subtypes and predicting survival of renal cell carcinoma using**
2 **histopathology image-based deep learning**

3 Fei Wang^{1a}, Yuanzhe Geng^{1a}, Ting Wang^{2a}, Ke Zhao^{3a}, Bin Xu^{4a}, Dachuan Zhang^{5a}, Hongzhou
4 Cai^{6a}, Qinbo Yuan⁷, Hang Lu¹, Yue Zhang¹, Xu Li¹, Yangyang Sun¹, Hang Gong¹, Raphael N.
5 Alolga¹, Xiangshan Fan^{2*}, Gaoxiang Ma^{8*}, Lian-Wen Qi^{1,8*}.

6 From ¹Clinical Metabolomics Center, China Pharmaceutical University, Nanjing, 211198,
7 China; ²Department of Pathology, The Affiliated Drum Tower Hospital, Nanjing University
8 Medical School, Nanjing, 210008, China; ³Department of Pathology, The Affiliated Jiangyin
9 Hospital of Southeast University Medical College, Wuxi, 214400, China; ⁴Department of
10 Urology, The Affiliated Zhongda Hospital of Southeast University, Nanjing, 210003, China;
11 ⁵Department of Pathology, The Third Affiliated Hospital of Soochow University, Changzhou,
12 213000, China; ⁶Department of Urology, The Affiliated Cancer Hospital of Nanjing Medical
13 University & Jiangsu Cancer Hospital & Jiangsu Institute of Cancer Research, Nanjing, 210009,
14 China; ⁷Department of Urology, Huaiyin People's Hospital of Huai'an City, Huai'an, 223300,
15 China; ⁸State Key Laboratory of Natural Medicines, School of Traditional Chinese Pharmacy,
16 China Pharmaceutical University, Nanjing, 211198, China.

17 ^a These authors contributed to this work equally.

18 **REPRINTS AND CORRESPONDENCE:** Address reprint requests to Dr. Lian-Wen Qi at the
19 Clinical Metabolomics Center of China Pharmaceutical University, E-mail address:
20 Qilw@cpu.edu.cn; Dr. Gaoxiang Ma at the School of Chinese Herbal Medicines of China
21 Pharmaceutical University, E-mail address: gaoxiang_ma@163.com; Dr. Xiangshan Fan at
22 Department of Pathology, The Affiliated Drum Tower Hospital, Nanjing University Medical
23 School, E-mail address: fanxiangshan@nju.edu.cn.

24 **Keyword:** Renal cell carcinoma; Deep learning; Histopathology; Diagnosis; Prognosis

25

26 **Abstract**

27 Classifying histopathological subtypes and predicting survival of renal cell carcinoma (RCC)
28 patients are critical steps towards treatment. In this work, we first proposed a deep learning
29 method involving patch-based segmentation, intelligent feature extraction and heatmap
30 visualization for classifying RCC into clear cell RCC, papillary RCC, chromophobe RCC, and
31 adjacent benign tissue. This algorithm was trained and validated using 2,374,446 patches, 6,340
32 whole-slide images, 2,399 patients from The Cancer Genome Atlas and 6 independent centers.
33 The classifiers provided areas under the curves of 0.979 to 0.996 in the internal phase, and 0.914
34 to 0.995 in the 6-center external phase. Furthermore, a modified deep learning approach
35 comprising automated detection of regions of interest, patch-level learning, and morphological
36 features-based risk scoring was developed for predicting survival of clear cell RCC patients.
37 The prognostication model provided a hazard ratio for poor versus good prognosis of 2.63 [95%
38 confidence interval (CI) 1.53–4.50, $P = 4.35e-4$] in the testing set, and 2.57 [95% CI 1.43–4.64,
39 $P = 1.68e-3$] in an independent validation set using multivariable analyses. In conclusion, the
40 developed histopathology image-based deep learning frameworks have the clinical potential to
41 assist pathologists in systematically evaluating histological information of RCC patients.

42

43 **Introduction**

44 Renal cell carcinoma (RCC) accounts for >90% of all kidney cancers. Epidemiological studies
45 show that RCC represents approximately 2.2% of all cancers, with ~400,000 new cases and
46 ~175,000 deaths yearly¹. It can be classified into three major subtypes: clear cell RCC (ccRCC),
47 the most common type accounting for 70% of all cases; papillary RCC (pRCC) which
48 represents 15% to 20% of all cases, and chromophobe RCC (chRCC), that accounts for 5% of
49 reported cases². The remaining subtypes are very rare with each accounting for $\leq 1\%$ of total
50 incidence. Each subtype of RCC has its specific histopathology, genetic characteristics, clinical
51 course, and response to therapy³.

52 Subtype classification and outcome prediction of RCC patients are critical steps towards
53 precise treatment. Histopathological slide is the gold standard of RCC subtype and stage^{2,4}.
54 Classification and prognostication based on human assessment remain time-consuming and
55 relatively subjective. In some cases, the distinction among RCC types is not clear as they may
56 share non-specific morphological patterns⁵. In addition, RCC is an extremely heterogeneous
57 disease, making prediction of prognosis a great challenge⁶.

58 Deep learning has achieved impressive successes in digital image analysis⁷. It enables
59 direct extraction of imaging features associated with classification and prognosis without
60 explicit programming. Deep learning has been applied to histopathology-based segmentation,
61 classification, and survival prediction of lung cancer⁸, colorectal cancer⁹, basal cell carcinoma¹⁰,
62 malignant mesothelioma¹¹, breast cancer¹², prostate cancer¹³, glioma¹⁴, and pan-cancer
63 analysis¹⁵. Although few models have been used in clinical practice¹⁶, deep learning algorithms
64 hold considerable promise and enormous potential in disease diagnosis and prognosis.

65 Development of deep learning models to study renal cancers have been very few.
66 Fenstermaker et al. trained a deep learning model to distinguish RCC subtypes based on biopsy
67 images of 42 patients from The Cancer Genome Atlas (TCGA)¹⁷. Tabibu et al. proposed a deep
68 learning method for the classification and survival prediction of RCC subtypes using images
69 from TCGA¹⁸. Chen et al. developed a strategy for integrating histology image and genomic
70 features to predict the outcomes of ccRCC patients from TCGA¹⁹. Marostica et al. diagnosed
71 RCC histological subtypes and predicted stage I ccRCC patients' survival outcomes²⁰.
72 Evidently, these studies have a limitation in terms of small sample size and do not include multi-
73 independent cohorts. Insufficient validation data hinder generalization of the underlying
74 hypotheses, leading to a possibility of overfitting²¹. To overcome the complexity and
75 heterogeneity of RCC, this work proposed histopathology imaged-based deep learning
76 frameworks for type classification and clinical outcome prediction.

77

78 **Results**

79 **Patient information.**

80 The samples for classification covered 6,340 whole-slide images from 2,399 RCC patients, of
81 which 3,260 slides of 941 patients were from TCGA and 3,080 slides of 1,458 patients were
82 from 6 independent cohorts ([Supplementary Table 1](#)). Images from TCGA included 1,550
83 ccRCC, 678 pRCC, 304 chRCC, and 728 adjacent benign tissues. Images from the 6
84 independent cohorts included 1,872 ccRCC, 454 pRCC, 586 chRCC, and 168 adjacent benign
85 tissues. Patient characteristics are shown in [Supplementary Tables 2–8](#). Patients from TCGA
86 were divided into 55% training set, 15% tuning set, and 30% internal testing set, and patients

87 from the 6 independent cohorts were used as external validation sets.

88 The samples for prognostication included 493 whole-slide images of 488 ccRCC patients
89 from TCGA and 343 slides of 316 ccRCC patients from an independent TADTH cohort. The
90 patient characteristics and follow-up information are shown in [Supplementary Table 9](#).
91 Inclusion criteria were age 18 years or older, ccRCC histologically proven to be R0 stage I–IV
92 and grade 1–4. The median follow-up period of ccRCC patients was 39.0 months (0.1–151.2)
93 for TCGA and 60.9 months (1.0–100.5) for TADTH. For TADTH cohort, follow-up consisted
94 of telephone calls following the initial diagnosis. Of these, 72 (72/388, 18.6%) were excluded
95 because of incomplete follow-up data. A total of 343 FFPE slides were obtained from the 316
96 patients (316/388, 81.4%; 87 women and 229 men). Clinicopathological findings were based
97 on tumor–node–metastasis (TNM) classification. There were 219 patients with stage I ccRCC,
98 31 with stage II, 59 with stage III, and 7 patients with distant metastases classified as stage IV
99 ccRCC.

100 HistoQC analysis was employed to check the image quality in terms of a number of metrics
101 and features ([Supplementary Fig. 2a](#)). These include microns per pixel ([Supplementary Fig. 2b](#)),
102 brightness ([Supplementary Fig. 2c](#)) and contrast ([Supplementary Fig. 2d](#)), etc. Of the 6,384
103 whole-slide images from TCGA and the independent testing set, 44 slides were excluded for
104 being unsuitable for subsequent computational analysis.

105 **Development of classifier models.**

106 The construction of classifier models involved steps of patch-based segmentation, intelligent
107 feature extraction, and heatmap visualization ([Fig. 1](#)). Slides and patches information for each
108 set are shown in [Supplementary Table 10](#). Inception V3 achieved a bit better performance than

109 the other four state-of-the-art CNN architectures, including VGG19, DenseNet, ResNet, and
110 MobileNet (Supplementary Table 11). Models trained at 5× provided higher performance
111 compared with 10× and 20× (Supplementary Table 12). Models trained with Inception V3 at
112 5× were then employed for subsequent analysis.

113 In the internal testing set, a binary classifier achieved an AUC of 0.996 [0.990–1.000] at the
114 patient level using FFPE slides, and 0.993 [0.984–0.998] for frozen slides (Supplementary
115 Table 12) in differentiating tumor from normal tissues. In the classification of RCC subtypes, a
116 three-way classification provided a macro-average AUC of 0.988 [0.979–0.995], a micro-
117 average AUC of 0.987 [0.979–0.995] for the FFPE samples, and a macro-average AUC of 0.979
118 [0.960–0.996], a micro-average AUC of 0.983 [0.970–0.995] for the frozen samples. The
119 performances of classifier models were also satisfactorily evaluated by Youden index
120 (Supplementary Table 13), PRC (Supplementary Fig. 3a–d), confusion matrices
121 (Supplementary Fig. 4a–d), and PDI (Supplementary Table 14).

122 **Validation of the classifier models in 6 independent cohorts.**

123 The classifier models were validated with 3,080 slides from 1,458 patients in the 6 multi-center
124 external phase. In identifying malignancy from adjacent benign tissues, the binary classifier
125 provided an AUC of 0.995 [0.987–0.999] at the patient level for TADTH cohort (Fig. 2a). In
126 differentiating between the 3 RCC types, the trained three-way classifier generated a macro-
127 average AUC of 0.968 [0.957–0.978] and a micro-average AUC of 0.967 [0.957–0.975] for
128 FFPE slides from TADTH (Fig. 2b). For frozen slides from TADTH, it yielded a macro-average
129 AUC of 0.916 [0.817–0.989] and a micro-average AUC of 0.914 [0.835–0.977] (Fig. 2c). The
130 three-way classifier offered a macro-average AUC of 0.964 [0.937–0.987] and a micro-average

131 AUC of 0.947 [0.919–0.971] for FFPE slides from TAJH (Fig. 2d), 0.960 [0.905–0.994] and
132 0.965 [0.934–0.989] from TAZH (Fig. 2e), 0.958 [0.914–0.986] and 0.951 [0.912–0.984] for
133 FFPE slides from TTAH (Fig. 2f), 0.964 [0.907–0.999] and 0.975 [0.945–0.995] for TACH
134 (Fig. 2g), and 0.958 [0.900–0.991] and 0.937 [0.872–0.980] for HPH (Fig. 2h). The
135 performances of classifier models were also satisfactorily evaluated by Youden index
136 (Supplementary Table 13), PRC (Supplementary Fig. 3e–l), confusion matrices (Supplementary
137 Fig. 4e–l), and PDI (Supplementary Table 14).

138 Heatmaps of the prediction probability over slides were generated to discern the tumor
139 regions associated with histological patterns. We have shown representations of images of the
140 four types (Fig. 3a). Heatmap visualizations were produced for which color is proportional to
141 the predicted probability (0~1) of the patch. The generated heatmap clearly differentiated tumor
142 from normal region at the patch level (Fig. 3b) by a binary classifier, and then classified the three
143 RCC subtypes of tumor patches by a three-way classifier (Fig. 3c). Most of the predicted
144 patches had a strong true positive probability for a certain type of RCC (Fig. 3d).

145 **Comparison of the classifier with pathologists.**

146 A subset of 221 FFPE and frozen slides from the internal testing set was used for algorithm-
147 pathologist comparison²². These included 24 out of 32 misclassified cases and 197 out of 492
148 correctly classified by the classifier. Pathologist 1 with 3 years' experience and pathologist 2
149 with 5 years' experience from TADTH were instructed to assess the digital whole-slide images
150 alone, independently of the clinical data provided by TCGA. The pathologists were free to use
151 a slide analytical software (ASAP; version 1.8) at varying zoom levels of up to 40×
152 magnification. For FFPE slides, the agreement of the classifier with TCGA using Cohen's

153 Kappa statistic was comparable with the pathologists (0.839 versus 0.758 for pathologist 1 and
154 0.813 for pathologist 2) ([Supplementary Fig. 5](#), [Supplementary Table 15](#)). For frozen slides, the
155 performance of the classifier was also comparable with the pathologists (0.820 versus 0.709 for
156 pathologist 1 and 0.767 for pathologist 2) ([Supplementary Fig. 5](#), [Supplementary Table 15](#)).

157 Of the images misclassified by the classifier, 83% (20/24) were also incorrectly diagnosed
158 by at least one pathologist, while 60% (30/50) of those misclassified by at least one pathologist
159 were classified successfully by the classifier. A few slides showed mixed histologic features,
160 poor differentiation or massive necrosis ([Supplementary Fig. 6](#)), leading to possible
161 misclassification by the model and the two pathologists.

162 **Development of the prognostic model.**

163 The ccRCC is the most common subtype and accounts for the majority of deaths. To obtain a
164 prognostic biomarker for ccRCC patients, a modified deep learning architecture using CNN
165 combined with Cox regression was developed. Patients in the testing set were divided into poor
166 and good prognostic groups based on the median risk score obtained from the training set. The
167 prognostication framework consisted of automated detection of regions of interest, patch-level
168 learning and morphological features-based risk scoring ([Fig. 4](#)). VGG19 generated a bit better
169 performance than Inception V3, DenseNet, ResNet, and MobileNet ([Supplementary Table 16](#)).
170 We compared the performances of 10, 20 and 30 representative image patches with the largest
171 average nuclei size for each slide. We observed that the c-index value did not increase with
172 more patches but slightly decreased using 20 and 30 patches compared to 10 patches
173 ([Supplementary Table 17](#)). As a result, we used 10 image patches of each slide as inputs to
174 VGG19 for subsequent prognostication.

175 Our survival deep-learning model showed strong prognostic power with a c-index of 0.779,
176 outperforming manual histologic-grade system of 0.748. Kaplan–Meier analysis showed that
177 the prognostic model had better prediction capability than histologic grade (log rank $P = 3.49\text{e-}$
178 5 versus $P = 2.22\text{e-}3$) (Fig. 5a,b). The risk index (low risk, high risk) was a prognostic factor
179 for overall survival in a univariate Cox analysis (HR 2.93 [95% CI 1.72–5.00], $P = 7.84\text{e-}5$)
180 (Supplementary Table 18). In multivariate analysis, the risk index was independently prognostic
181 of survival in all-stage (HR = 2.63 [95% CI 1.53–4.50], $P = 4.35\text{e-}4$), early-stage (HR = 2.41
182 [95% CI 1.03–5.61], $P = 0.04$), and late-stage (HR = 2.81 [95% CI 1.38–5.72], $P = 4.30\text{e-}3$)
183 tumors adjusting for stage and age (Fig. 5c). The histologic grade (G1+G2, G3+G4) was only
184 prognostic in all-stage but not in early-stage and late-stage tumors. The predicted risk score was
185 highly consistent with patient outcomes, especially for 3-year follow-up (Fig. 5d). The model’s
186 performance remained satisfactory as measured by accelerated failure time analyses
187 (Supplementary Table 19). The predicted risk scores were significantly associated with tumor
188 grading and staging (Spearman’s $\rho = 0.34$ and 0.22 , respectively, both $P < 0.01$) (Supplementary
189 Table 20). The established scoring algorithm allowed patches to be scored for each patient.
190 Comparison between high-risk patches and low-risk patches was beneficial in obtaining
191 morphological characteristics associated with prognosis (Supplementary Fig. 7).

192 **Validation of the prognostic model in an independent cohort.**

193 We validated the model in an independent cohort of 316 RCC patients from TADTH. The model
194 showed a c-index of 0.769, better than manual histologic grade of 0.753. This prognostic model
195 provides significant survival differences between high- and low-risk groups (log rank $P = 5.88\text{e-}$
196 4), better than the histologic grade (log rank $P = 1.02\text{e-}3$) (Fig. 5e,f). The risk index was a

197 predictor of patient outcomes in a univariate Cox analysis (HR 2.70 [95% CI 1.50–4.87], $P =$
198 $9.65e-4$) (Supplementary Table 18). In multivariate analysis, the risk index was independently
199 prognostic of survival in all-stage (HR = 2.57 [95% CI 1.43–4.64], $P = 1.68e-3$), early-stage
200 (HR = 2.30 [95% CI 1.07–4.94], $P = 0.03$), and late-stage tumors (HR = 2.65 [95% CI 1.03–
201 6.83], $P = 0.04$) (Fig. 5g). The histologic grade showed prognostication in all-stage but not in
202 early-stage and late-stage tumors. We have presented the distribution plot of risk scores and
203 survival time of ccRCC patients in Fig. 5h. The performance in independent cohort was also
204 evaluated by accelerated failure time analyses (Supplementary Table 19). The risk index was
205 also positively correlated with tumor grade and stage (Spearman's $\rho = 0.45$ and 0.16 ,
206 respectively, $P < 0.01$). (Supplementary Table 20).

207

208 Discussion

209 In this study, we developed histopathology image-based deep learning frameworks: a classifier
210 model with Inception V3 to classify RCC subtypes, and a prognostication model with VGG19
211 to predict outcomes of ccRCC. The classifier model involves patch-based segmentation,
212 intelligent feature extraction, and heatmap visualization using whole-slide images. It was
213 trained and validated using 3,260 slides of 941 patients from TCGA and 3,080 real-world slides
214 of 1,458 patients from 6 independent centers. The algorithm could unambiguously distinguish
215 normal from tumor tissues with AUC >0.99 and effectively distinguished between the three
216 RCC subtypes with AUC >0.91 , with accuracy comparable to or even better than two clinical
217 pathologists. To obtain a prognostic biomarker, the developed classifier model first
218 automatically outlined cancerous regions of ccRCC, and then the prognostication model was

219 built by combining deep learning with Cox regression. Independent validation confirmed that
220 the predicted risk index is significantly associated with clinical outcome of ccRCC.

221 In general, most cases of RCC can be easily classified based on histological criteria. We
222 observed that the misclassified cases displayed a combination of morphological features.
223 Indeed, the presence of clear cells is not unique to ccRCC but can also be observed in some
224 cases of pRCC and chRCC. Similarly, papillary structures characteristic of pRCC, can also be
225 present in other subtypes²³.

226 When applied to external validation, there was a slight drop of the classifier' performance
227 in differentiating tumor from normal (AUC 0.1%) and subtypes (AUC 1%–7%). It should be
228 noted that these slides were collected without data curation. The interlaboratory differences in
229 slide preparation might explain the decrease. The difference between the scanners in brightness
230 and contrast could also somewhat affect the prediction accuracy.

231 ccRCC is the most common variant and has wide ranging clinical outcomes attributed to its
232 tumor heterogeneity⁶. Historically, the most widely used histological grading system for ccRCC
233 is the Fuhrman grading system²³, which relies on the experience of pathologists. In recent years,
234 automatic models depending on extracting engineered imaging features have been proposed^{18,24}.
235 More recently, survival analysis has been approached with a classification task by predicting
236 several periods of survival time divided by specific time points²⁰. These classifiers, however,
237 cannot model the risk values with certain survival times and lack independent follow-up cohort.
238 This work predicted the patient outcomes by learning prognostic features directly from slides.
239 It employed Cox regression, a time-to-event model that can better fit the prediction of survival
240 and model all patients' risks for a range of survival time.

241 Our prognostic model highlights the importance of dedifferentiation, stroma component,
242 cellular diversity, inflammation, and growth pattern for prognosis. We observed that a
243 considerable proportion of high-risk patches were located in the sarcomatoid and rhabdoid
244 regions, which are well known as malignant and dedifferentiated components^{24,26}. Some
245 contained pleomorphic or giant cells and coagulative necrosis, consistent with previous
246 observations²⁷. Some had unique characteristics of infiltrative growth patterns²⁸. Others
247 presented stromal components in the tumor microenvironment²⁹. In contrast, good prognostic
248 patches were located in compact small/large nests and cystic regions. Some showed
249 intratumoral inflammatory reaction, tubular, clear cell papillary or chromophobe cell-like
250 patterns³⁰.

251 This study has some limitations. The algorithms focused on classifying the three common
252 subtypes of renal cancer, and did not identify rare subtypes such as Mit family translocation
253 RCC and clear cell papillary RCC due to unavailability of samples. For clinical application of
254 the model, more samples are encouraged to fine-tune the algorithm to increase its
255 generalizability. Another limitation is that our diagnostic algorithms were only trained on H&E-
256 stained images. Integrating features from immunochemical-stained images remains to be
257 explored.

258 In conclusion, we have presented deep learning frameworks for systematic assessment of
259 histological information on RCC. The established models have been evaluated in multi-center,
260 independent cohorts. Our models have the potential of accelerating the pathologist's evaluation
261 for kidney tissue slides, and making it more objective and replicable in clinic.

262

263 **Methods**

264 **Data sources.**

265 Our biospecimen were collected from TCGA Resource Network³¹ and 6 independent cohorts
266 in China. Based on a previous research³², 15 initial samples originally submitted as ccRCC in
267 TCGA were reclassified as chRCC. The 6 independent cohorts were: The Affiliated Drum
268 Tower Hospital (TADTH), The Affiliated Jiangyin Hospital of Southeast University Medical
269 College (TAJH), The Affiliated Zhongda Hospital of Southeast University (TAZH), The Third
270 Affiliated Hospital of Soochow University (TTAH), The Affiliated Cancer Hospital of Nanjing
271 Medical University & Jiangsu Cancer Hospital & Jiangsu Institute of Cancer Research (TACH),
272 and Huaiyin People's Hospital of Huai'an City (HPH). Both formalin-fixed and paraffin-
273 embedded (FFPE) sections and frozen sections were included. There were no exclusion criteria
274 on age, gender or race.

275 Cases were reviewed by two pathologists with more than 3 years' experience. The
276 diagnoses were made based on morphology following the World Health Organization
277 recommendation. All slides were finally checked by Prof. Fan Xiangshan, a pathologist with
278 many years of experience and the head of pathology department at TADTH. In case of
279 inconsistent decisions, immunostaining was further performed for classification using a panel
280 of antibodies.

281 Slides from TTAH were scanned by a KF-PRO-120 scanner (Konfoong Biotechnology) at
282 20× magnification, and samples from the other independent cohorts were scanned by a
283 NanoZoomer 2.0-RS scanner (Hamamatsu). Quality checking was performed for the whole-
284 slide images by an open-source tool *HistoQC* to exclude unsuitable samples³³.

285 **Construction of classification models.**

286 Several commonly used convolutional neural networks (CNNs) were compared including
287 Inception V3, 19-layer Visual Geometry Group (VGG19), DenseNet, ResNet, and MobileNet.
288 We optimized the parameters of fully connected layers as well as the convolution layers by
289 unfreezing the whole network. The parameters of each network were initialized with pretrained
290 parameters from the ImageNet classification challenge as a starting point. For every image in
291 the training set, rotational invariance was achieved through data augmentation with random
292 horizontal and vertical flips before being fed into the network. For the training process, the
293 batch size was set to 32. The cross-entropy loss was set as the loss function on the training and
294 tuning sets. All parameters of the networks were trained jointly using stochastic gradient
295 descent³⁴ as a backpropagation method with learning rate of 0.05, weight decay of 0.9 and
296 momentum of 0.9. The number of training epochs was 50, and the network with the lowest loss
297 on the tuning set was selected.

298 Patients of each RCC type from TCGA dataset were divided into 55% training set, 15%
299 tuning set, and 30% internal testing set. For the training and tuning sets, tumor lesions in each
300 malignant slide were manually annotated by pathologists ([Supplementary Fig. 1](#)). Each whole-
301 slide image was automatically split into hundreds of nonoverlapping 299×299-pixel patches
302 using the Openslide (version 3.4.1). Benign patches for training and tuning sets were sampled
303 from both adjacent benign tissue slides and non-tumor regions in malignant slides.
304 Segmentation at magnifications of 5×, 10× and 20× were compared. Patches from patients in
305 the training set were input into the architectures to intelligently extract morphological features
306 and make decisions with parameters constantly optimized. The tuning set was then employed

307 to evaluate the trained algorithm and avoid overfitting. After tuning, the algorithm was then
308 kept locked in subsequent validations.

309 For internal and external validation, the slides were automatically segmented into patches
310 without manual annotation. Each patch was predicted as tumor or benign by a binary classifier
311 with a possibility threshold of 0.5. Tumor patches were then classified as ccRCC, pRCC or
312 chRCC based on the possibility predicted by a three-way subtype classifier. Predictions were
313 made with probability and indicated using colors, producing a heatmap for visualization of each
314 slide. The prediction results were then aggregated by calculating percentage of normal/tumor
315 patches over all patches, or percentage of one-RCC type patches over predicted tumor patches
316 in slides from each patient.

317 **Establishment of a prognostic model.**

318 ccRCC patients from TCGA were randomly divided into 55% training set, 15% tuning set and
319 30% internal testing set. An independent cohort from TADTH was used for external validation.
320 Slides were segmented into nonoverlapping 1196×1196-pixel patches at 20×. The developed
321 classifier model first automatically outlined cancerous patches from benign ones in each ccRCC
322 slide. The nuclei size of cancerous patches was calculated by a Python package
323 (<https://github.com/DigitalSlideArchive/HistomicsTK>). Top 10, 20 or 30 tumor patches with
324 the largest average nuclei size for each slide were selected and compared.

325 Subpatches (224×224-pixel) were further randomly sampled from the representative
326 patches as inputs. The prognostic deep learning model was trained end-to-end, directly from
327 slide to survival time. We adapted the predictive algorithm by combining CNN architecture and
328 Cox regression model. The performances of CNNs were compared, including VGG19,

329 Inception V3, DenseNet, ResNet, and MobileNet. The CNN made use of various convolutions
330 from original architecture to extract image features. The last layer was changed to a fully
331 connected layer containing one node, which predicted a risk score for the input subpatch. The
332 risk scores in each subpatch were presented to a Cox proportional hazards layer allowing for
333 the use of censored data to calculate the Cox loss function. The loss function was optimized
334 using stochastic gradient descent with a learning rate of $5.0e-5$, weight decay of 0.9 and
335 momentum of 0.9. The training process was run for 100 epochs, and the model with the lowest
336 loss on the tuning set was saved.

337 For internal and external validation sets, each representative tumor patch was split into
338 nonoverlapping 16 subpatches. The median risk of these subpatches was calculated as the value
339 for each patch, and then the highest value among all representative patches was selected as a
340 patient risk. These processing procedures were designed to address tumor heterogeneity by
341 emulating histological assessment by pathologists. In routine clinical practice, the
342 pathologically assigned grade is based on the most malignant region when ccRCC shows
343 morphologic grade variation³⁵. The patients were then separated into equivalent groups of low-
344 and high-risk on the basis of their predicted risk scores with the median risk in the training set
345 as the cut-off point.

346 **Statistics.**

347 To evaluate the performance of classifiers, areas under the receiver operating characteristic
348 (ROC) and precision-recall curves (PRC), Youden index, and polytomous discrimination index
349 (PDI)³⁶ were computed using the Python library sklearn. Confidence intervals (CIs) 95% were
350 measured using 1,000 iterations of the bootstrap method. The agreements between the classifier

351 and pathologists were measured by the Cohen's Kappa statistic³⁷. For survival prediction,
352 overall survival time was defined as the time from nephrectomy to death by any cause or the
353 date of last follow-up. Harrell's c-index was used to quantify the concordance between
354 predicted risks and true survival time. The survival curves of the predicted poor and good
355 survival subgroups were plotted by the Kaplan-Meier method. The Cox regression model was
356 used to obtain hazard ratios (HRs) and 95% CIs. Associations between the models and other
357 risk factors were evaluated by Spearman's correlation coefficients. A two-sided *P* value of less
358 than 0.05 was considered significant. All statistical analyses were done in R unless otherwise
359 noted (R version 3.6.1). The open-source deep learning framework Keras (<https://keras.io>) was
360 used to train and evaluate the algorithms.

361

362 **Data availability**

363 The publicly shared RCC histopathology images in TCGA dataset to train and test the models
364 is available at <https://portal.gdc.cancer.gov/>. The dataset consists of 3,260 whole-slide images
365 from 941 patients, including 1,550 ccRCC, 678 pRCC, 304 chRCC, and 728 adjacent benign
366 tissues. The independent datasets are not publicly available due to hospital regulations and
367 patient privacy. Source data are provided with this paper. The remaining data are available
368 within the Article, Supplementary information, or available from the authors upon request.

369

370 **Code availability**

371 The code is available at <https://github.com/cilcmc/dlrc>.

372

373 **Author contributions**

374 L.-W.Q., G.M., and X.F. conceived and directed the project. F.W. and Y.G., implemented and
375 trained the deep learning model. H.L., Y.Z., X.L., Y.S., and H.G. collected and scanned slides
376 from the independent datasets. F.W. and G.M. contributed to the analysis of the data. T.W., K.Z.,
377 B.X., D.Z., H.C., and Q.Y. helped identify and label the slides. Y.Z. and H.G. performed the
378 follow up and collected the prognostic information of patients. F.W., L.-W.Q., G.M. and R.N.A.
379 wrote the manuscript. All authors discussed the results and reviewed the manuscript.

380

381 **Declaration of interests**

382 The authors declare no competing interests.

383

384 **Source of funding**

385 This work was supported by the National Natural Science Fund of China for Distinguished
386 Young Scholars (No. 81825023).

387

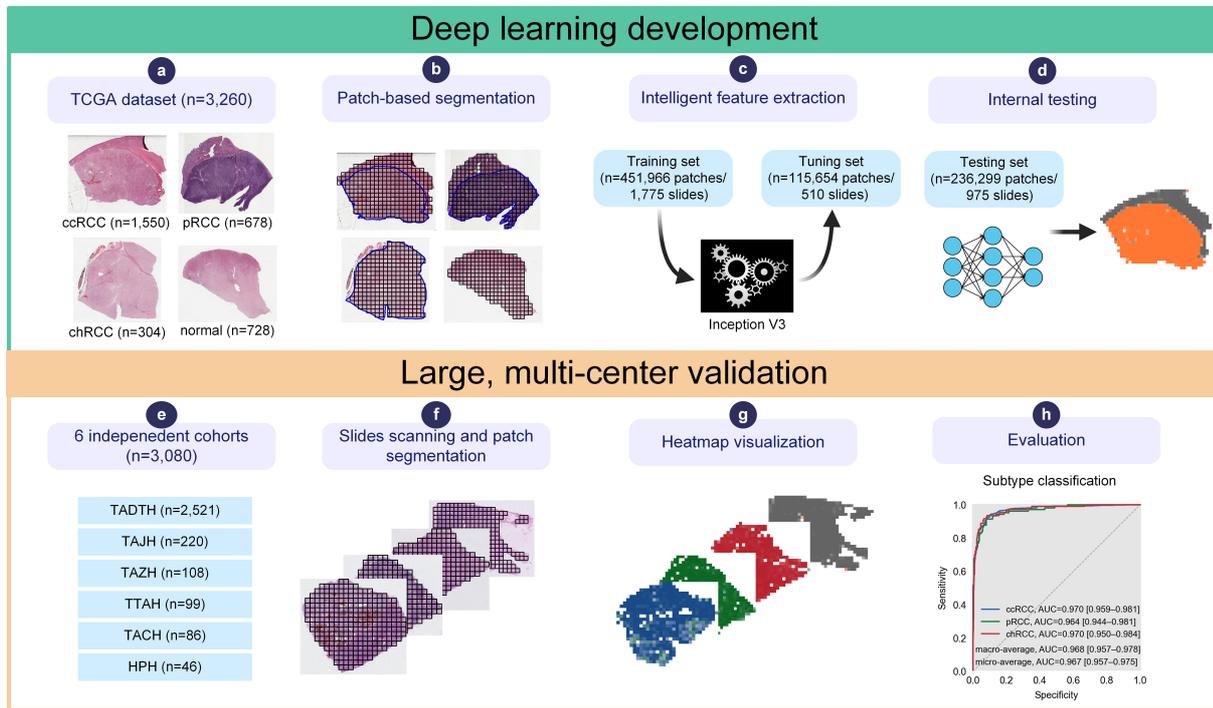
388 **References**

- 389 1. Bray, F. et al. Global cancer statistics 2018: GLOBOCAN estimates of incidence and
390 mortality worldwide for 36 cancers in 185 countries. *CA Cancer J. Clin.* **68**, 394–424
391 (2018).
- 392 2. Moch, H., Cubilla, A.L., Humphrey, P.A., Reuter, V.E. & Ulbright, T.M. The 2016 WHO
393 classification of tumours of the urinary system and male genital organs—part A: renal,
394 penile, and testicular tumours. *Eur. Urol.* **70**, 93–105 (2016).
- 395 3. Escudier, B. et al. Renal cell carcinoma: ESMO clinical practice guidelines for diagnosis,
396 treatment and follow-up. *Ann. Oncol.* **27**, v58–v68 (2016).

- 397 4. Leibovich, B.C. et al. Predicting oncologic outcomes in renal cell carcinoma after surgery.
398 *Eur. Urol.* **73**, 772–780 (2018).
- 399 5. Hsieh, J.J. et al. Renal cell carcinoma. *Nat. Rev. Dis. Primers* **3**, 17009 (2017).
- 400 6. Gulati, S. et al. Systematic evaluation of the prognostic impact and intratumour
401 heterogeneity of clear cell renal cell carcinoma biomarkers. *Eur. Urol.* **66**, 936–948 (2014).
- 402 7. Bera, K., Schalper, K.A., Rimm, D.L., Velcheti, V. & Madabhushi, A. Artificial
403 intelligence in digital pathology—new tools for diagnosis and precision oncology. *Nat.*
404 *Rev. Clin. Oncol.* **16**, 703–715 (2019).
- 405 8. Coudray, N. et al. Classification and mutation prediction from non–small cell lung cancer
406 histopathology images using deep learning. *Nat. Med.* **24**, 1559–1567 (2018).
- 407 9. Skrede, O.-J. et al. Deep learning for prediction of colorectal cancer outcome: a discovery
408 and validation study. *Lancet* **395**, 350–360 (2020).
- 409 10. Campanella, G. et al. Clinical-grade computational pathology using weakly supervised
410 deep learning on whole slide images. *Nat. Med.* **25**, 1301–1309 (2019).
- 411 11. Courtiol, P. et al. Deep learning-based classification of mesothelioma improves prediction
412 of patient outcome. *Nat. Med.* **25**, 1519–1525 (2019).
- 413 12. Bejnordi, B.E. et al. Diagnostic assessment of deep learning algorithms for detection of
414 lymph node metastases in women with breast cancer. *JAMA* **318**, 2199–2210 (2017).
- 415 13. Bulten, W. et al. Automated deep-learning system for Gleason grading of prostate cancer
416 using biopsies: a diagnostic study. *Lancet Oncol.* **21**, 233–241 (2020).
- 417 14. Mobadersany, P. et al. Predicting cancer outcomes from histology and genomics using
418 convolutional networks. *Proc. Natl. Acad. Sci. USA.* **115**, E2970–E2979 (2018).
- 419 15. Kalra, S. et al. Pan-cancer diagnostic consensus through searching archival histopathology
420 images using artificial intelligence. *NPJ. Digit. Med.* **3**, 31 (2020).
- 421 16. Liu, X. et al. A comparison of deep learning performance against health-care professionals
422 in detecting diseases from medical imaging: a systematic review and meta-analysis. *Lancet*
423 *Digit. Health* **1**, e271–e297 (2019).
- 424 17. Fenstermaker, M., Tomlins, S.A., Singh, K., Wiens, J. & Morgan, T.M. Development and
425 validation of a deep-learning model to assist with renal cell carcinoma histopathologic
426 interpretation. *Urology* **144**, 152–157 (2020).

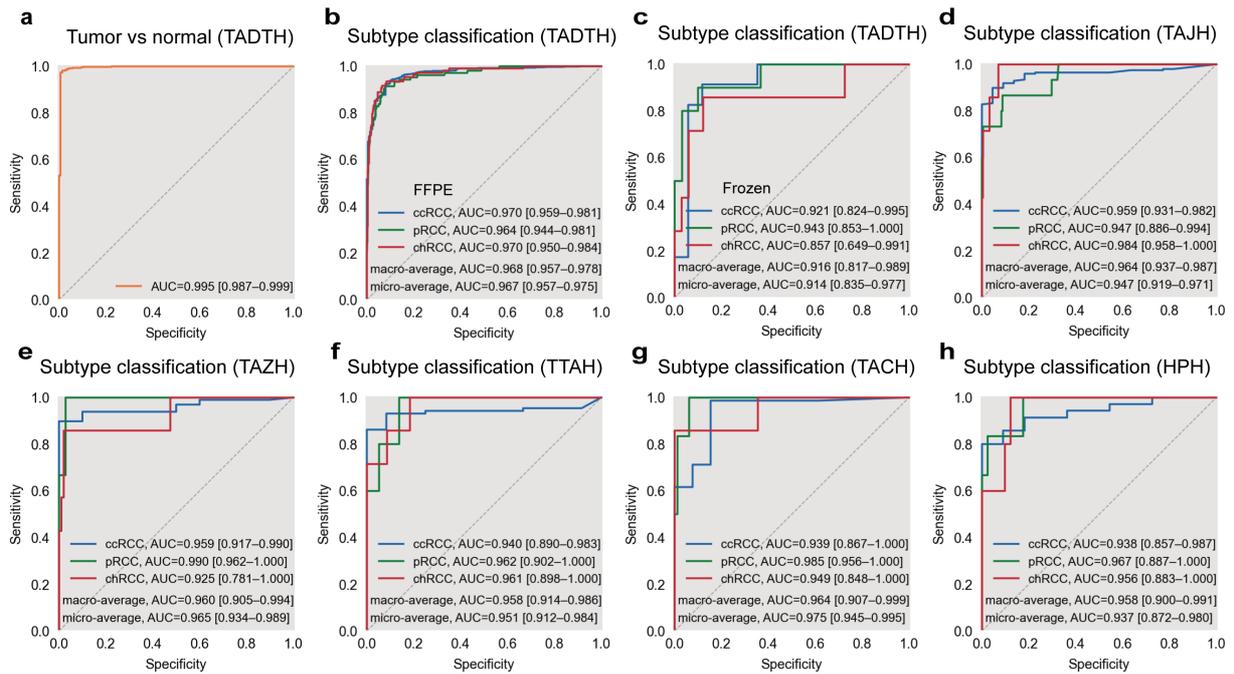
- 427 18. Tabibu, S., Vinod, P. & Jawahar, C. Pan-renal cell carcinoma classification and survival
428 prediction from histopathology images using deep learning. *Sci. Rep.* **9**, 1–9 (2019).
- 429 19. Chen, R.J. et al. Pathomic fusion: An integrated framework for fusing histopathology and
430 genomic features for cancer diagnosis and prognosis. *IEEE Trans. Med. Imaging* (2020).
- 431 20. Marostica, E. et al. Development of a histopathology informatics pipeline for classification
432 and prediction of clinical outcomes in subtypes of renal cell carcinoma. *Clin. Cancer Res.*
433 (2021).
- 434 21. Topol, E.J. High-performance medicine: the convergence of human and artificial
435 intelligence. *Nat. Med.* **25**, 44–56 (2019).
- 436 22. van Smeden, M., Van Calster, B. & Groenwold, R.H.H. Machine learning compared with
437 pathologist assessment. *JAMA* **319**, 1725–1726 (2018).
- 438 23. Fuhrman, S.A., Lasky, L.C. & Limas, C. Prognostic significance of morphologic
439 parameters in renal cell carcinoma. *Am. J. Surg. Pathol.* **6**, 655–664 (1982).
- 440 24. Cheng, J. et al. Integrative analysis of histopathological images and genomic data predicts
441 clear cell renal cell carcinoma prognosis. *Cancer Res.* **77**, e91–e100 (2017).
- 442 25. Cheville, J.C. et al. Sarcomatoid renal cell carcinoma: an examination of underlying
443 histologic subtype and an analysis of associations with patient outcome. *Am. J. Surg.*
444 *Pathol.* **28**, 435–441 (2004).
- 445 26. Singh, R.R. et al. Intratumoral morphologic and molecular heterogeneity of rhabdoid renal
446 cell carcinoma: challenges for personalized therapy. *Mod. Pathol.* **28**, 1225–1235 (2015).
- 447 27. Delahunt, B. et al. Grading of clear cell renal cell carcinoma should be based on nucleolar
448 prominence. *Am. J. Surg. Pathol.* **35**, 1134–1139 (2011).
- 449 28. Cai, Q. et al. Ontological analyses reveal clinically-significant clear cell renal cell
450 carcinoma subtypes with convergent evolutionary trajectories into an aggressive type.
451 *EBioMedicine* **51**, 102526 (2020).
- 452 29. Junttila, M.R. & de Sauvage, F.J. Influence of tumour micro-environment heterogeneity
453 on therapeutic response. *Nature* **501**, 346–354 (2013).
- 454 30. Verine, J. et al. Architectural patterns are a relevant morphologic grading system for clear
455 cell renal cell carcinoma prognosis assessment: Comparisons with WHO/ISUP grade and
456 integrated staging systems. *Am. J. Surg. Pathol.* **42**, 423–441 (2018).

- 457 31. Grossman, R.L. et al. Toward a shared vision for cancer genomic data. *N. Engl. J. Med.*
458 **375**, 1109–1112 (2016).
- 459 32. Ricketts, C.J. et al. The Cancer Genome Atlas comprehensive molecular characterization
460 of renal cell carcinoma. *Cell Rep.* **23**, 313–326.e315 (2018).
- 461 33. Janowczyk, A., Zuo, R., Gilmore, H., Feldman, M. & Madabhushi, A. HistoQC: An open-
462 source quality control tool for digital pathology slides. *JCO Clin. Cancer Inform* **3**, 1–7
463 (2019).
- 464 34. Bottou, L. Stochastic gradient descent tricks. in *Neural networks: Tricks of the trade* 421–
465 436 (Springer, 2012).
- 466 35. Delahunt, B. et al. A novel grading system for clear cell renal cell carcinoma incorporating
467 tumor necrosis. *Am. J. Surg. Pathol.* **37**, 311–322 (2013).
- 468 36. Van Calster, B. et al. Extending the c-statistic to nominal polytomous outcomes: the
469 Polytomous Discrimination Index. *Stat. Med.* **31**, 2610–2626 (2012).
- 470 37. McHugh, M.L. Interrater reliability: the kappa statistic. *Biochem. Med.* **22**, 276–282
471 (2012).
472



474

475 **Fig. 1. Schematic presentation of the diagnostic system. a–d, Deep learning development. a,**
 476 **Whole-slide images of renal cell carcinoma (RCC) from The Cancer Genome Atlas (TCGA).**
 477 **A total of 3,260 images were collected, including 1,550 clear cell RCC (ccRCC), 678 papillary**
 478 **RCC (pRCC), 304 chromophobe RCC (chRCC), and 728 normal tissues. b, Manual annotation**
 479 **of tumor lesions for each slide and automatic segmentation by nonoverlapping 299 × 299-pixel**
 480 **patches at 5× magnification with background removed. c, Training Inception V3 architecture to**
 481 **intelligently extract diagnostic features. The training set was input into the algorithm to make**
 482 **prediction, and the tuning set was used to avoid overfitting. d, Patch-based classification for the**
 483 **internal testing set. The results were aggregated per slide to generate the heatmaps. e–h, Large,**
 484 **multi-center validation. e, Assembling real-world slides from 6 independent cohorts. The**
 485 **framework was validated with 3,080 slides from The Affiliated Drum Tower Hospital (TADTH),**
 486 **The Affiliated Jiangyin Hospital of Southeast University Medical College (TAJH), The**
 487 **Affiliated Zhongda Hospital of Southeast University (TAZH), The Third Affiliated Hospital of**
 488 **Soochow University (TTAH), The Affiliated Cancer Hospital of Nanjing Medical University &**
 489 **Jiangsu Cancer Hospital & Jiangsu Institute of Cancer Research (TACH), Huaiyin People's**
 490 **Hospital of Huai'an City (HPH). f, Slides scanning and patch segmentation. The slides were**
 491 **scanned and automatically segmented without manual annotation. g, Heatmap visualization.**
 492 **Patch-level diagnosis were made with predicted probability and are indicated using colors,**
 493 **producing a heatmap for visualization of each slide. h, Receiver operating characteristic and**
 494 **precision-recall curves, Youden index, and polytomous discrimination index statistics for**
 495 **classifications of normal tissue versus tumor and between RCC subtypes.**



496

497

498

499

500

501

502

503

504

505

506

507

508

509

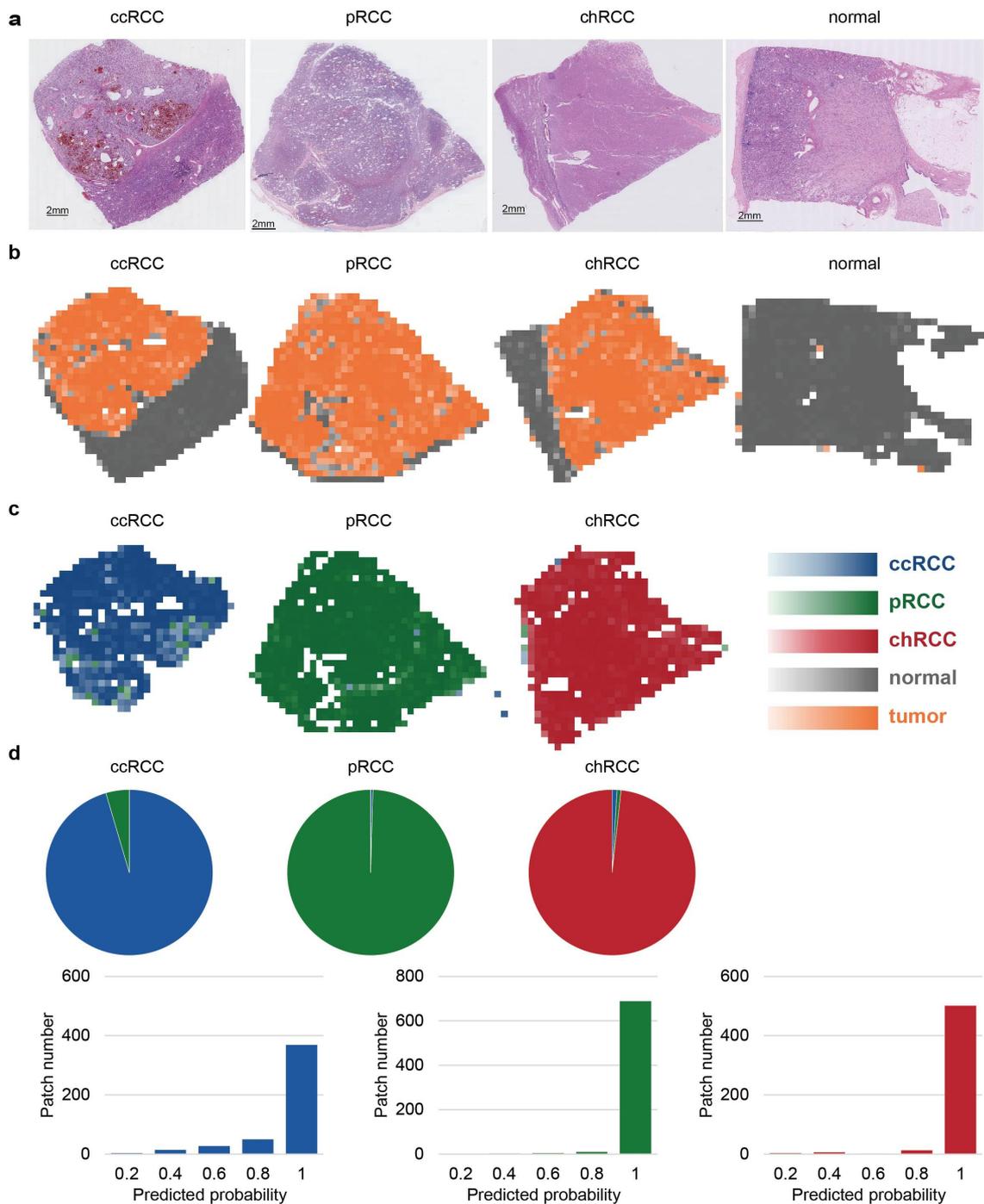
510

511

512

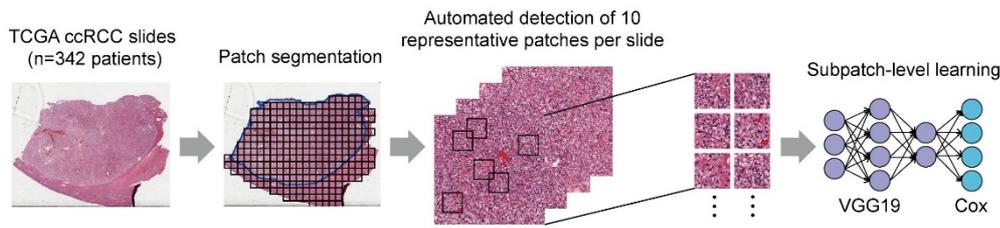
513

Fig. 2. Differential diagnosis of normal tissue from tumor and between renal cell carcinoma (RCC) types at the patient level. a. Receiver operating characteristic curves (ROC) for identifying tumor versus normal using FFPE slides from The Affiliated Drum Tower Hospital (TADTH, n=2,474 slides from 859 patients). **b.** ROC for classifying RCC types using FFPE slides from TADTH (n=2,306 slides from 859 patients). **c.** ROC for classifying RCC subtypes using frozen slides from TADTH (n=47 slides from 40 patients). **d.** ROC for classifying RCC types using FFPE slides from The Affiliated Jiangyin Hospital of Southeast University Medical College (TAJH, n=220 slides from 220 patients). **e.** ROC for classifying RCC subtypes using FFPE slides from The Affiliated Zhongda Hospital of Southeast University (TAZH, n=108 slides from 108 patients). **f.** ROC for classifying RCC subtypes using FFPE slides from The Third Affiliated Hospital of Soochow University (TTAH, n=99 slides from 99 patients). **g.** ROC for classifying RCC subtypes using FFPE slides from The Affiliated Cancer Hospital of Nanjing Medical University & Jiangsu Cancer Hospital & Jiangsu Institute of Cancer Research (TACH, n=86 slides from 86 patients). **h.** ROC for classifying RCC subtypes using FFPE slides from Huaiyin People's Hospital of Huai'an City (HPH, n=46 slides from 46 patients).

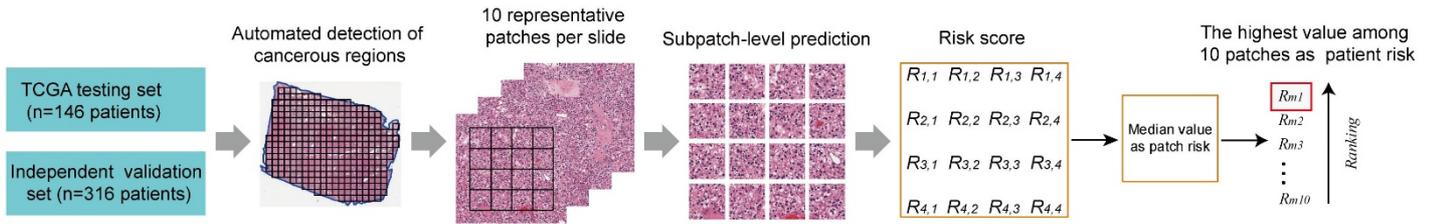


514
 515 **Fig. 3. Heatmap visualization of renal cell carcinoma (RCC) slides.** a, Representative
 516 whole-slide images of the three types and adjacent benign tissues from The Affiliated Drum
 517 Tower Hospital (TADTH) cohort. b, Binary classifier for heatmap visualization of tumor versus
 518 normal. c, Three-way classifier for heatmap visualization of ccRCC versus pRCC versus
 519 chRCC. d, Patch-based prediction probability for a slide. The pie chart shows the predicted
 520 patch number for each category, and the bar chart shows the predicted probability of each patch.
 521 In b–d, the predicted heatmaps with probabilities are assigned to each patch, where grey is for
 522 patches classified as normal, orange for tumor, blue for ccRCC, green for pRCC, and red for
 523 chRCC.

a Deep learning development



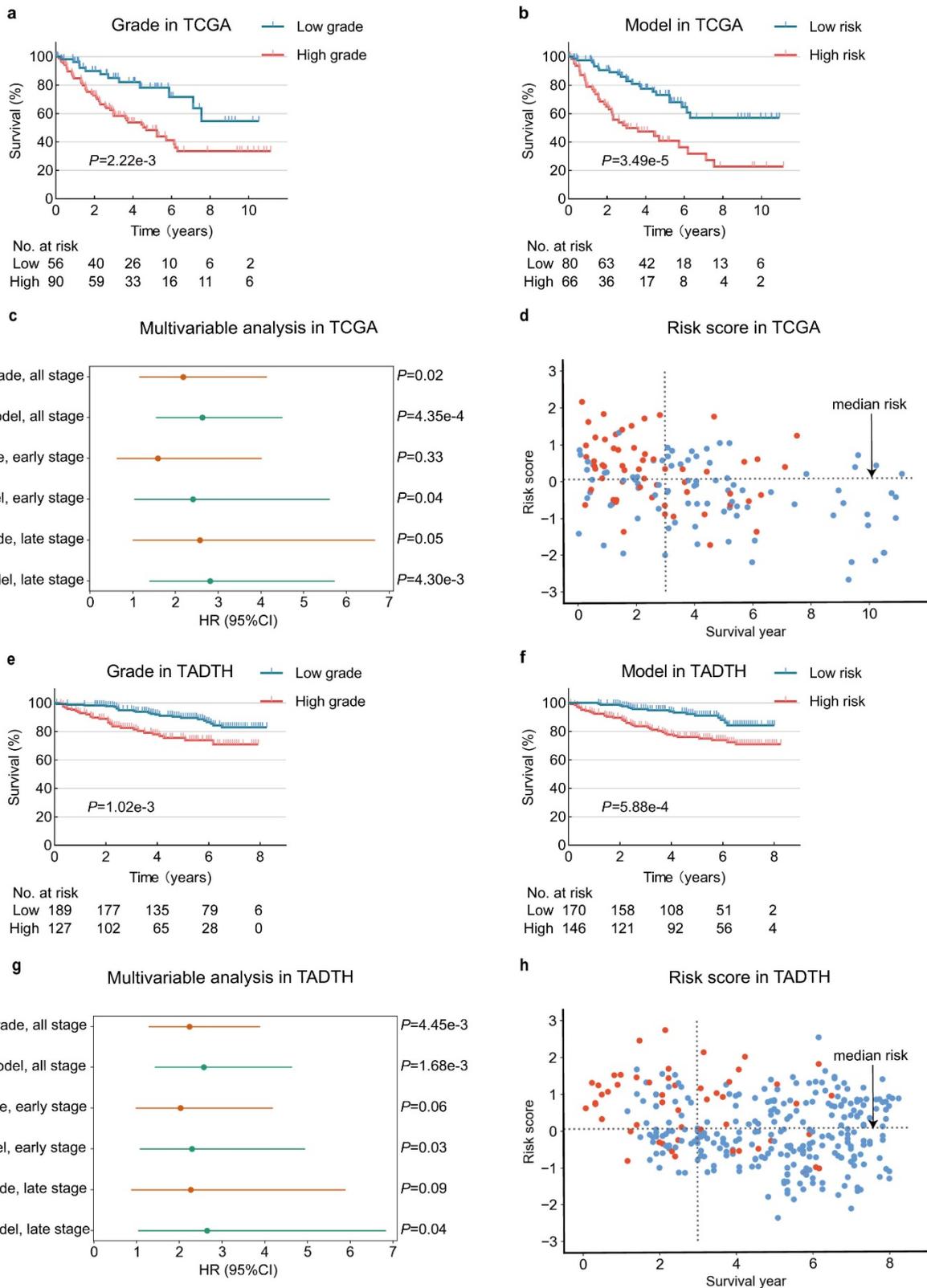
b Prognosis prediction



5_ .

525 **Fig. 4. Pipeline of the prognostic framework.** **a**, Deep learning development. Whole-slide
 526 images of clear cell renal cell carcinoma (ccRCC) patients from The Cancer Genome Atlas
 527 (TCGA) dataset were randomly divided into training (55%), tuning (15%) and testing (30%)
 528 sets. After patch-based segmentation, the classifier model first automatically outlined cancerous
 529 regions of ccRCC. Ten representative patches (1196×1196-pixel at 20× magnification) with the
 530 largest average nuclei in tumor regions were selected from each slide. Subpatches (224×224-
 531 pixel) were sampled from the representative patches. The 19-layer Visual Geometry Group
 532 (VGG19) architecture combined with Cox regression was employed to predict patient outcome
 533 end-to-end. **b**, Prognosis prediction. For slides in the TCGA testing and independent validation
 534 sets, cancerous regions were automatically detected by the diagnostic model, and 10
 535 representative patches were selected from each slide. Each representative patch was split into
 536 16 nonoverlapping subpatches; the median subpatch risk was calculated for each patch, and
 537 then the highest value among 10 patches was selected as patient risk.

538



539
540
541
542
543
544

Fig. 5. Performances of prognostic model and manual histologic grade for clear cell renal cell carcinoma (ccRCC) patients. **a**, Kaplan-Meier analysis of histologic grade in The Cancer Genome Atlas (TCGA) testing set (n=146 patients). **b**, Kaplan-Meier analysis of the prognostic model in TCGA testing set. **c**, Hazard ratios (HRs) with 95% confidence interval (CI) of grade (G1+G2, G3+G4) and prognostic model (low risk, high risk) in multivariable analyses adjusting

545 for stage and age across different stages in TCGA testing set. **d**, The distribution of risk scores
546 and outcomes of ccRCC patients in TCGA testing set. **e**, Kaplan-Meier analysis of histologic
547 grade in The Affiliated Drum Tower Hospital (TADTH) independent validation set (n=316
548 patients). **f**, Kaplan-Meier analysis of the prognostic model in TADTH. **g**, HRs with 95% CI of
549 histologic grade (G1+G2, G3+G4) and prognostic model (low risk, high risk) in multivariable
550 analyses adjusting for stage and age across different stages in TADTH. **h**, The distribution of
551 risk scores and outcomes of ccRCC patients in TADTH. The blue dots represent censored data
552 and the red dots represent the dead patients. Prognostic model was highly associated with
553 patient survival time, especially for a 3-year follow-up (**d**, **h**). The median risk score obtained
554 in the TCGA training set was used as threshold to divide patients into low- and high-risk groups
555 (**a-h**).

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Supplementarydata.pdf](#)