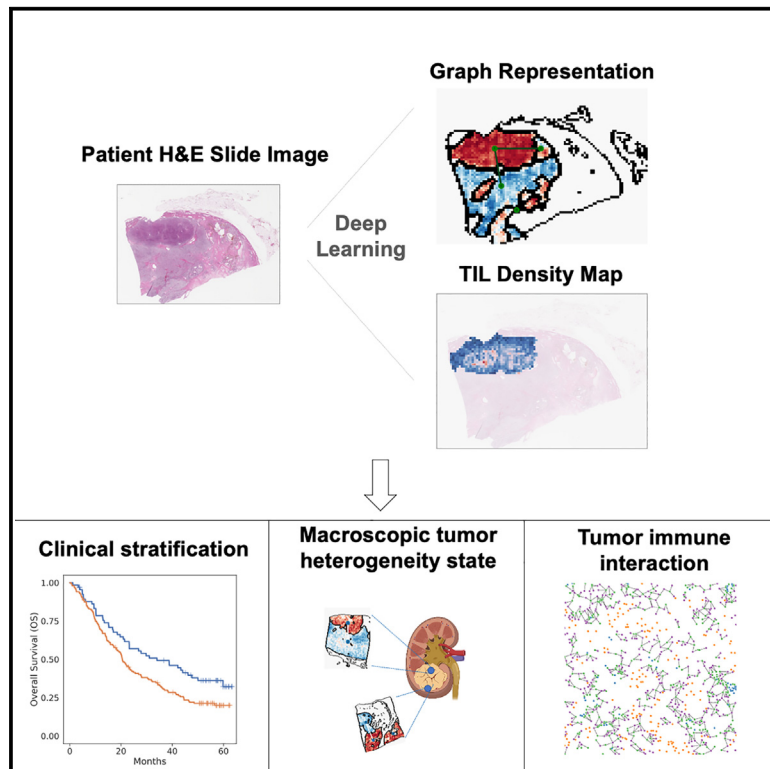


# Spatially aware deep learning reveals tumor heterogeneity patterns that encode distinct kidney cancer states

## Graphical abstract



## Authors

Jackson Nyman, Thomas Denize, Ziad Bakouny, ..., Toni K. Choueiri, Sabina Signoretti, Eliezer M. Van Allen

## Correspondence

eliezerm\_vanallen@dfci.harvard.edu

## In brief

Nyman et al. develop a deep-learning-derived representation of tumor-grade heterogeneity in clear cell renal cell carcinoma histology images and relate this feature to the immune microenvironment and distinct patient outcomes using multiple clinical cohorts.

## Highlights

- Deep learning on diagnostic kidney cancer images quantifies tumor and immune phenotypes
- Spatial microheterogeneity in tumor grade generalizes across macroscopic states
- Microheterogeneity associates with selective immunotherapy response



## Article

# Spatially aware deep learning reveals tumor heterogeneity patterns that encode distinct kidney cancer states

Jackson Nyman,<sup>1,2,3</sup> Thomas Denize,<sup>4</sup> Ziad Bakouny,<sup>1,3,5,6,18</sup> Chris Labaki,<sup>1,3,6,19</sup> Breanna M. Titchen,<sup>1,3,7</sup> Kevin Bi,<sup>1,3</sup> Surya Narayanan Hari,<sup>1,3</sup> Jacob Rosenthal,<sup>8,9</sup> Nicita Mehta,<sup>1,3,6</sup> Bowen Jiang,<sup>1,3,10</sup> Bijaya Sharma,<sup>17</sup> Kristen Felt,<sup>17</sup> Renato Umeton,<sup>8,9,11,12</sup> David A. Braun,<sup>13</sup> Scott Rodig,<sup>4</sup> Toni K. Choueiri,<sup>1,6,14</sup> Sabina Signoretti,<sup>3,4,6,15</sup> and Eliezer M. Van Allen<sup>1,3,6,16,20,\*</sup>

<sup>1</sup>Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA, USA

<sup>2</sup>Harvard Graduate Program in Systems Biology, Cambridge, MA, USA

<sup>3</sup>Broad Institute, Cambridge, MA, USA

<sup>4</sup>Department of Pathology, Brigham and Women's Hospital, Boston, MA, USA

<sup>5</sup>Department of Medicine, Brigham and Women's Hospital, Boston, MA, USA

<sup>6</sup>Harvard Medical School, Boston, MA, USA

<sup>7</sup>Harvard Graduate Program in Biological and Biomedical Sciences, Boston, MA, USA

<sup>8</sup>Department of Informatics & Analytics, Dana-Farber Cancer Institute, Boston, MA, USA

<sup>9</sup>Department of Pathology and Laboratory Medicine, Weill Cornell Medicine, New York, NY, USA

<sup>10</sup>Stanford University, Stanford, CA, USA

<sup>11</sup>Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge, MA, USA

<sup>12</sup>Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA, USA

<sup>13</sup>Center of Molecular and Cellular Oncology, Yale Cancer Center, Yale School of Medicine, New Haven, CT, USA

<sup>14</sup>Brigham and Women's Hospital, Boston, MA, USA

<sup>15</sup>Department of Oncologic Pathology, Dana-Farber Cancer Institute, Boston, MA, USA

<sup>16</sup>Department of Population Sciences, Dana-Farber Cancer Institute, Boston, MA, USA

<sup>17</sup>ImmunoProfile, Department of Pathology, Brigham & Women's Hospital and Dana-Farber Cancer Institute, Boston, MA, USA

<sup>18</sup>Department of Medicine, Memorial Sloan Kettering Cancer Center, New York, NY, USA

<sup>19</sup>Department of Medicine, Beth Israel Deaconess Medical Center, Boston, MA, USA

<sup>20</sup>Lead contact

\*Correspondence: [eliezerm\\_vanallen@dfci.harvard.edu](mailto:eliezerm_vanallen@dfci.harvard.edu)

<https://doi.org/10.1016/j.xcrm.2023.101189>

## SUMMARY

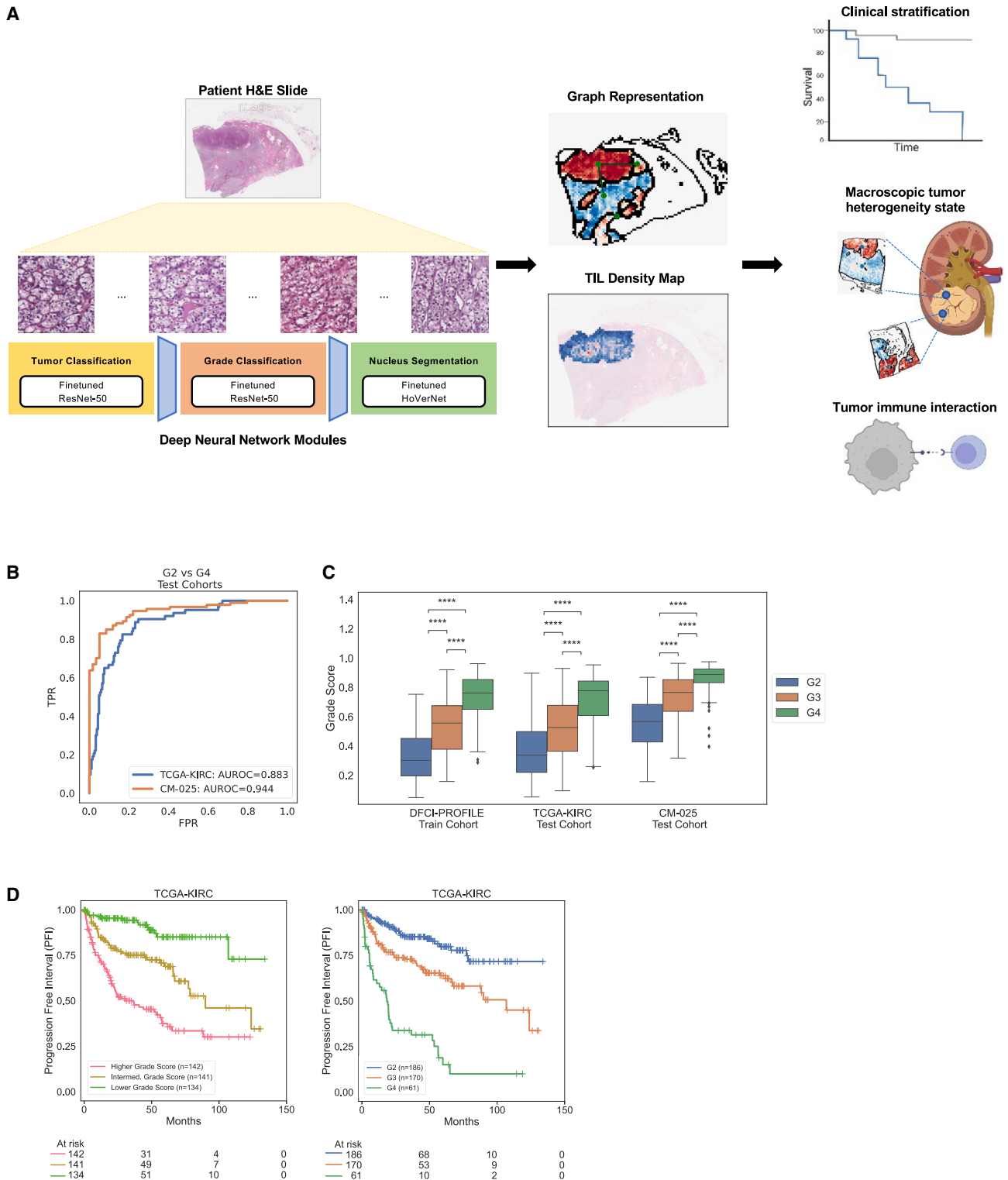
Clear cell renal cell carcinoma (ccRCC) is molecularly heterogeneous, immune infiltrated, and selectively sensitive to immune checkpoint inhibition (ICI). However, the joint tumor-immune states that mediate ICI response remain elusive. We develop spatially aware deep-learning models of tumor and immune features to learn representations of ccRCC tumors using diagnostic whole-slide images (WSIs) in untreated and treated contexts ( $n = 1,102$  patients). We identify patterns of grade heterogeneity in WSIs not achievable through human pathologist analysis, and these graph-based “microheterogeneity” structures associate with *PBRM1* loss of function and with patient outcomes. Joint analysis of tumor phenotypes and immune infiltration identifies a subpopulation of highly infiltrated, microheterogeneous tumors responsive to ICI. In paired multiplex immunofluorescence images of ccRCC, microheterogeneity associates with greater PD1 activation in CD8<sup>+</sup> lymphocytes and increased tumor-immune interactions. Our work reveals spatially interacting tumor-immune structures underlying ccRCC biology that may also inform selective response to ICI.

## INTRODUCTION

Renal cell carcinoma (RCC) is among the 10 most common cancers worldwide and is comprised of several histological subtypes.<sup>1</sup> The clear cell histological subtype (clear cell RCC [ccRCC]) is the most common form of RCC and accounts for the vast majority (75%–80%) of metastatic cases.<sup>1</sup> In addition to highly recurrent mutations in hypoxia (*VHL*) and chromatin regulator genes (e.g., *PBRM1*, *BAP1*, *SETD2*), ccRCC exhibits

extensive genomic intratumoral heterogeneity (ITH),<sup>2</sup> which was correlated with worse progression-free survival in both the TRACERx and TCGA-KIRC cohorts.<sup>3–5</sup> Nuclear grade, an established histopathologic score of tumor nuclei dedifferentiation, is a primary prognostic feature in ccRCC and can provide a histologic description of ITH<sup>6</sup> to pinpoint cell structures enriched for metastatic potential.<sup>7,8</sup> In addition, high nuclear grade has been associated with increased tumor-infiltrating lymphocytes (TILs) in ccRCC,<sup>9</sup> though whether molecular ITH or its relationship to





**Figure 1. A spatially aware deep-learning framework for studying ccRCC**

(A) Our approach builds a series of biologically relevant prediction models to provide both high-resolution and readily human-understandable representations of ccRCC slide images. The first two models identify tumor tissue and grade phenotype within predicted tumor regions, each using a fine-tuned ResNet-50 convolutional neural network (CNN). A third model identifies tumor-infiltrating lymphocytes (TILs) using a fine-tuned HoVerNet CNN. Local predictions are

(legend continued on next page)

histologic properties (e.g., grade, TILs) inform immunoresponsive tumor states in ccRCC is uncertain. Indeed, while immune checkpoint inhibitors (ICIs) are a standard therapy in ccRCC, this tumor type defies many conventions about molecular features that associate with selective ICI response identified in other solid tumors,<sup>10–13</sup> and both the underlying biology and clinical biomarkers to stratify patients for ICI in ccRCC remain elusive.

Current approaches to simultaneously quantify tumor-intrinsic heterogeneity and its potential relationship to immune microenvironmental interactions in patients are hamstrung by (1) a lack of spatial resolution in molecular sequencing, (2) difficulty with simultaneous multi-regional measurements of tumor and immune molecular properties in sufficient cohort sizes, and (3) practical limitations related to pathologists being incapable of manually performing such measurements from histopathology data at scale. However, by leveraging biologically guided deep learning applied to whole-slide images (WSIs), highly detailed evaluations of both established pathology features (e.g., nuclear grade) and spatial structures that arise from these features are possible at a scale otherwise intractable via manual pathologist review.<sup>14,15</sup> Thus, we hypothesized that spatially aware deep-learning models of ccRCC WSIs could provide a unified understanding of distinct tissue structures that dictate biological and clinical states in ccRCC, and we examined this hypothesis in multiple clinical ccRCC cohorts.

## RESULTS

### Development of a deep-learning framework for ccRCC diagnostic images

We developed prediction models that provide high-resolution, quantitative, and human-understandable representations of ccRCC hematoxylin and eosin (H&E) WSIs to identify established pathology features like tumor tissue and nuclear grade at scale<sup>16,17</sup> (Figure 1A; STAR Methods; Figures S1, and S2). In total, after quality control analysis, our study included WSIs from 1,102 patients with ccRCC (n = 208 Dana-Farber Cancer Institute [DFCI]-PROFILE, 421 The Cancer Genome Atlas Clear Cell Renal Cell Carcinoma [TCGA-KIRC], 439 Checkmate-025 [CM-025], 21 multiblock nephrectomy cases, 13 paired multiplex immunofluorescence [mIF] ccRCC cases; STAR Methods). We first built a latent representation of the nuclear-grade topology of ccRCC WSIs by combining multiple layers of inferred histological features produced by neural network models. In the first layer, we used a ResNet50 convolutional neural network (CNN) fine-tuned with pixel-level expert pathologist annotations from a retrospective cohort of patients collected at the DFCI (DFCI-PROFILE; n = 208 total patients; STAR Methods; Figure S2) to produce a classifier

that distinguishes regions of tumor tissue from adjacent non-tumor tissue. We next used putative tumor tissue to train a second CNN classifier to distinguish low- (G2) from high-grade (G4) cases in the DFCI-PROFILE cohort. We validated the performance of these models in two unseen test cohorts: the TCGA-KIRC (n = 421) and CM-025 (n = 439) datasets using an ensemble of four models (area under the receiver operating characteristic curve [AUROC] = 0.88 [TCGA] and 0.944 [CM-025]; Figures 1B, 1C, and S3).

We subsequently averaged model grade predictions across an entire slide image to produce a continuous grade score and compared this score with the categorical pathologist-assigned nuclear grade (STAR Methods). We discretized continuous grade scores into tercile bins to mirror G2/G3/G4 categories, which produced significant patient stratifications for both progression-free interval (PFI) and overall survival (OS) (Figures 1D and S1A;  $p < 1e-5$  [PFI],  $p < 1e-5$  [OS], multivariate log-rank test). Thus, in these cohorts, a deep-learning computer vision model could both mimic and refine clinically standard categorical nuclear-grade assignments in ccRCC.

Finally, to represent each patient slide compactly, we formed region adjacency graphs (RAGs) that describe where regions of distinct tumor and grade prediction phenotypes occur in a slide, as well as whether these regions directly or indirectly contact one another (STAR Methods). In aggregate, this framework produces a multi-layered, information-rich latent representation of ccRCC patient tumor images. Moreover, by condensing the local predictions made by each model, we also represent spatial patterns that arise between these image-derived features.

### Spatial microheterogeneity in ccRCC

Upon inspection of the model representations, we observed a distinct heterogeneity phenomenon in continuous nuclear-grade prediction graphs: some WSIs demonstrated co-occurrence of different grade phenotypes within the same slide, while others were markedly homogeneous. This co-occurrence, which we termed “microheterogeneity,” can be described in two primary (but not mutually exclusive) forms: (1) “proximal,” wherein heterogeneity occurred between tumor tissues that directly contacted one another (Figure 2A), and (2) “distal,” wherein stromal barriers or separation in the slide image interrupted the differing tumor tissues (Figure 2B). We identified microheterogeneity (any proximal or distal occurrences) in 40.6% of TCGA-KIRC cases and in 34.7% of CM-025 (Figures 2C, 2D, and S4). WSI microheterogeneity was present in varying frequencies within pathologist-assigned grade labels in each cohort, without any consistent pattern between pathologist grade label groups (Figure 2D; frequency of microheterogeneity = 0.36/0.494/0.317 [G2/G3/G4, TCGA-KIRC], 0.524/0.333/0.230 [G2/G3/G4, CM-025]). To produce a continuous measurement of the amount of microheterogeneity

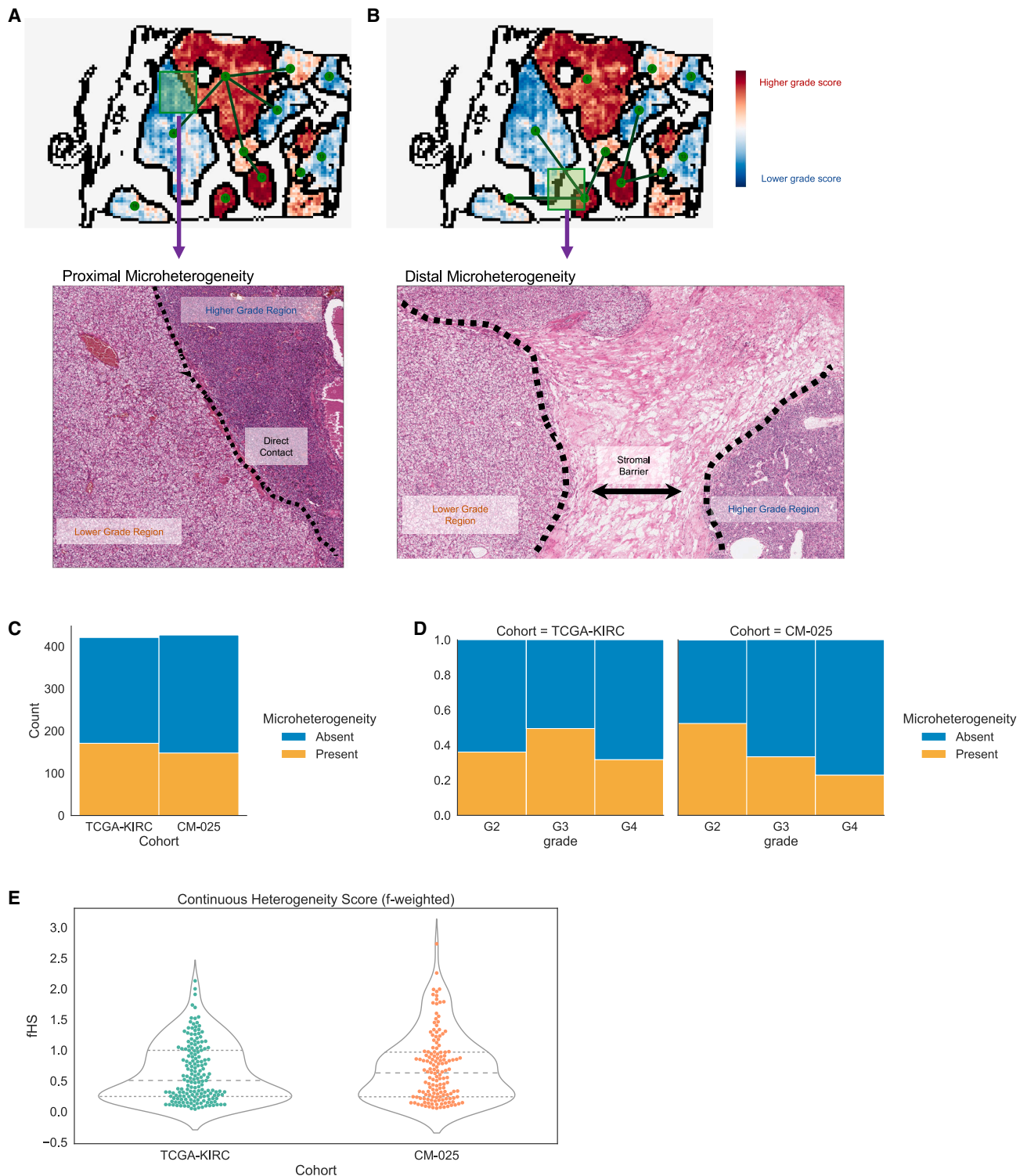
grouped via watershed segmentation and assembled into graph representations for slide-level description of patients. Computationally inferred patient representations capture both clinically relevant and biologically informative characteristics of ccRCC.

(B) Area under ROC curves for the task of distinguishing G2 from G4 in held-out test cohorts. TPR, true positive rate; FPR, false positive rate.

(C) Comparison of assigned pathologist grade and grade score on held-out test cohorts (TCGA-KIRC, CM-025) in-house training set used for tumor- and grade-classifier development (DFCI-PROFILE).

(D) Kaplan-Meier curves for progression-free interval (PFI) in TCGA-KIRC based on tercile bins of computationally inferred continuous grade score (left) and assigned pathologist grade (right).

\*:  $p \leq 0.05$ , \*\*:  $p \leq 0.001$ , \*\*\*:  $p \leq 0.0001$ , \*\*\*\*:  $p \leq 0.00001$ .



**Figure 2. Computationally inferred phenotypic variation in ccRCC**

(A) Representative example of proximally occurring grade microheterogeneity (dashed line indicates interface of region contact).

(B) Representative example of distally occurring grade microheterogeneity.

(C–E) Summary statistics surrounding microheterogeneity in the TCGA-KIRC and CM-025 cohorts.

(C) Number of patients with/without microheterogeneity.

(D) Frequency of microheterogeneity by assigned pathologist grade (where available).

(E) Distribution of continuous heterogeneity score (f-weighted) in non-homogeneous cases, which describes the extent of microheterogeneity in a given slide.

in a single WSI among slides that had microheterogeneity, we calculated the weighted sum of the number of heterogeneous contacts (RAG edges) per WSI, where larger weights are given to contacts with similar tumor region areas (STAR Methods). In two independent ccRCC cohorts, tumors exhibited a wide distribution of microheterogeneity abundance per WSI (Figures 2E, S4, and S5). Thus, in ccRCC WSIs, distinct nuclear-grade patterns create microheterogeneity structures that can be quantitatively represented as graphs for further investigation.

### Establishing the linkage between micro- and macrolevel heterogeneity

Given the distribution of microheterogeneity abundance per WSI, we then examined how this local, slide-level microheterogeneity related to variation throughout a whole tumor (“macroheterogeneity”). We evaluated a cohort of multiple spatially separated tumor blocks from the same nephrectomy specimen (Figures 3A, 3B, and S6). For a given patient’s tumor, the maximum microheterogeneity abundance in any single WSI correlated with the presence of microheterogeneity across all WSIs from that tumor, and this correlation was not driven by patient sample size (Figures 3C, S7A, and S7B; STAR Methods). In contrast, variation in image-derived grade scores did not correlate with sample size or frequency of microheterogeneity (Figures 3D, S7C, and S7D). Moreover, subsequent predictive modeling demonstrated that a single WSI could predict microheterogeneity for the remaining WSIs from the same patient (minimum [min.]  $\log_{10}(\text{Bayes factor}) = 3.04$ ; Figure S8; STAR Methods). These findings indicate that observing a single reference slide is predictive of macrolevel tumor phenotypes, particularly when that reference slide contains higher-grade phenotypes, and that a single ccRCC WSI encodes latent information regarding spatial structures present throughout the tumor.

### Molecular correlates of microheterogeneity in ccRCC

Since certain somatic mutations have been associated with macrolevel tumor heterogeneity, we subsequently evaluated whether computationally derived microheterogeneity structures from a single WSI were associated with recurrent somatic driver mutations in ccRCC, even though direct prediction of mutations from ccRCC images without multi-layered analysis has thus far been limited.<sup>18,19</sup> WSIs from tumors with somatic *PBRM1* loss of function (LOF), previously associated with molecular ITH, were also associated with a higher frequency of microheterogeneity compared with WSIs from non-LOF tumors (Figure S9). We also examined other common driver mutations in ccRCC and found a similar trend of higher microheterogeneity frequency in *SETD2* LOF mutants but inconclusive trends for *BAP1* and *PTEN* (Figure S9). Regarding somatic copy-number alterations, tumors with 9p21.3 deletions, a molecular feature previously implicated in ccRCC oncogenesis,<sup>20–22</sup> were enriched for microheterogeneity patterns (Figure S9). Whereas specific molecular events associated with global molecular ITH were also linked to histologic microheterogeneity, global molecular ITH itself was decoupled from histologic microheterogeneity (STAR Methods; Figure S10). However, grade score itself appears to correlate with weighted genome integrity index (wGII),<sup>23</sup> a measure of copy-number alteration burden, and this does not appear to depend on microheterogeneity state (Figure S10). Thus, mi-

croheterogeneity patterns also encoded features related to recurrent somatic alterations in ccRCC, even though these patterns are distinct from molecular ITH.

### Prognostic relevance of microheterogeneity

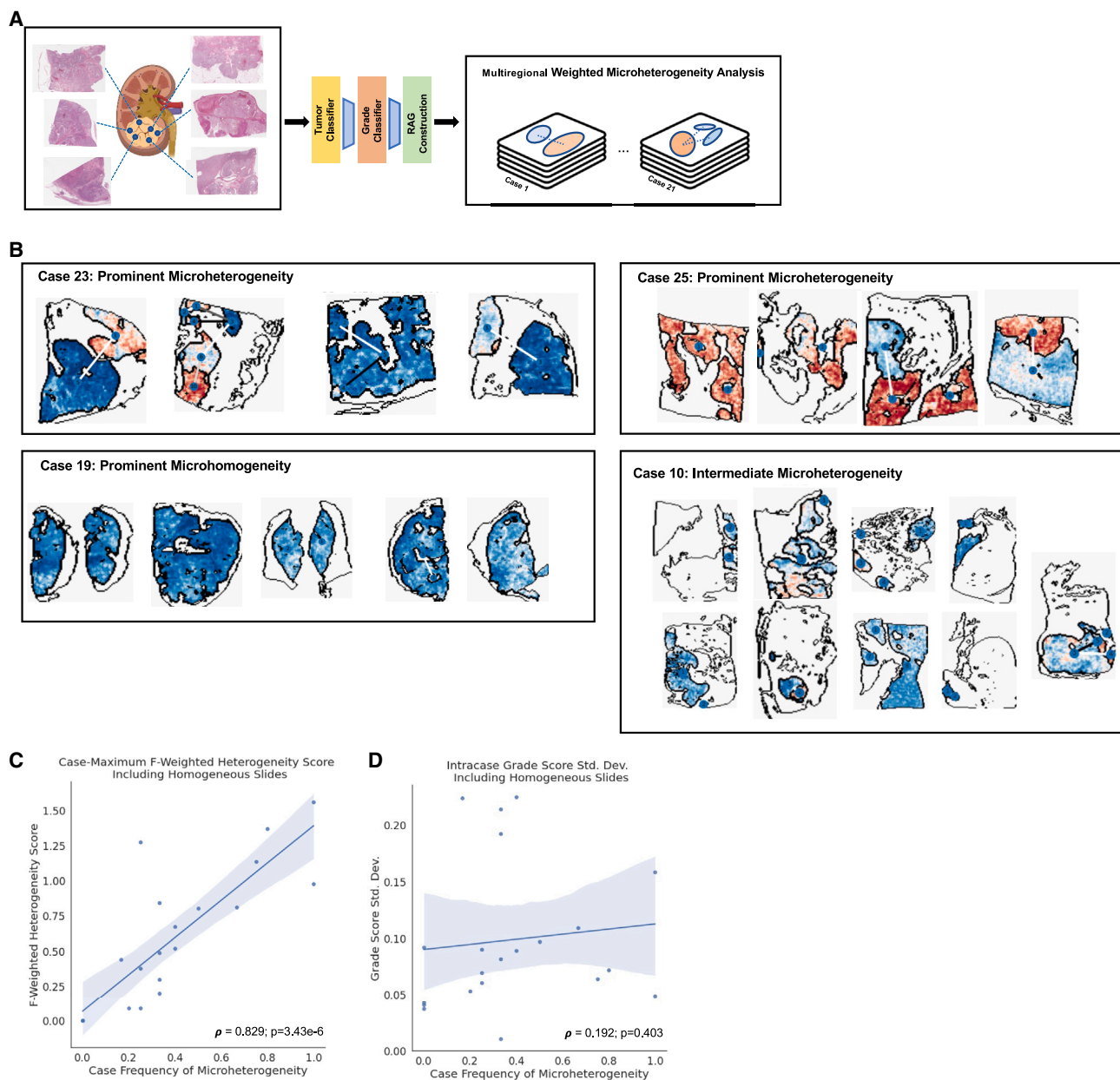
Since certain somatic mutations have prognostic value in ccRCC, we assessed whether computationally derived microheterogeneity from WSIs contained additional prognostic information beyond pathologist-derived nuclear grade. We compared univariate Cox proportional hazards models, using either pathologist-assigned grade or computationally inferred continuous grade, with bivariate models that introduced a binary indicator of whether microheterogeneity was observed. In both univariate and bivariate models in TCGA-KIRC, continuous grade had a stronger concordance index (C-Index) for PFI but not for OS (Figure S11). Within bivariate models for both survival contexts and grading types, the presence of microheterogeneity was negatively correlated with survival (hazard ratios all above 1), most notably in the continuous grade model (Figure S12). Thus, the presence of microheterogeneity in a single localized, untreated ccRCC WSI can identify tumors with poor prognosis and greater metastatic potential, consistent with phenomena previously described using multi-region molecular profiling.<sup>7</sup>

### Microheterogeneity and response to ICI

In addition to prognostic clinical value, we assessed whether this computer-vision-derived feature may be predictive for certain ccRCC therapeutics. We assessed spatial microheterogeneity patterns within both treatment arms of CM-025, a phase 3 randomized clinical trial cohort that compared anti-PD1 blockade (nivolumab) to mTOR inhibition (everolimus) in patients with anti-angiogenic refractory metastatic ccRCC (STAR Methods).<sup>13,24</sup> The presence of microheterogeneity was associated with improved OS and PFS in the ICI arm but not in the mTOR inhibitor arm (Figures 4A and S13). Given that continuous grade score correlated with OS for each trial arm, we also examined whether microheterogeneous cases correlated with changes in survival due to having lower overall grade scores. However, within microheterogeneous cases in the ICI arm, grade score did not contribute a statistically significant predictive signal for PFS or OS, though it trended toward significance for OS (Figures 4B and S14). Thus, in CM-025, microheterogeneity was selectively associated with improved response to ICI even though it was a poor prognostic marker in the primary, untreated setting.

### High immune infiltration combined with grade microheterogeneity identifies a further population of ICI responders

Immune infiltration as measured by CD8 immunofluorescence was not associated with response to ICI,<sup>13,25</sup> despite its predictive value in other immune-responsive cancers. We hypothesized that TIL patterns may still be relevant for predicting response to ICI in ccRCC but that joint inference of tumor spatial heterogeneity with TIL patterns is required for adequate context. Thus, we inferred TILs in the CM-025 WSIs (Figures S15–S18), confirmed the fidelity of this approach in high-grade tumors (Figures S19 and S20), and related these features to microheterogeneity (Figure 5A; STAR Methods). In WSIs



**Figure 3. Linkage between microheterogeneity and whole-tumor variation**

(A) Schematic for creating a multi-regional weighted microheterogeneity analysis using computer vision models.

(B) Example data collections from four patients, showing RAG plots for the scanned slide of each tissue block.

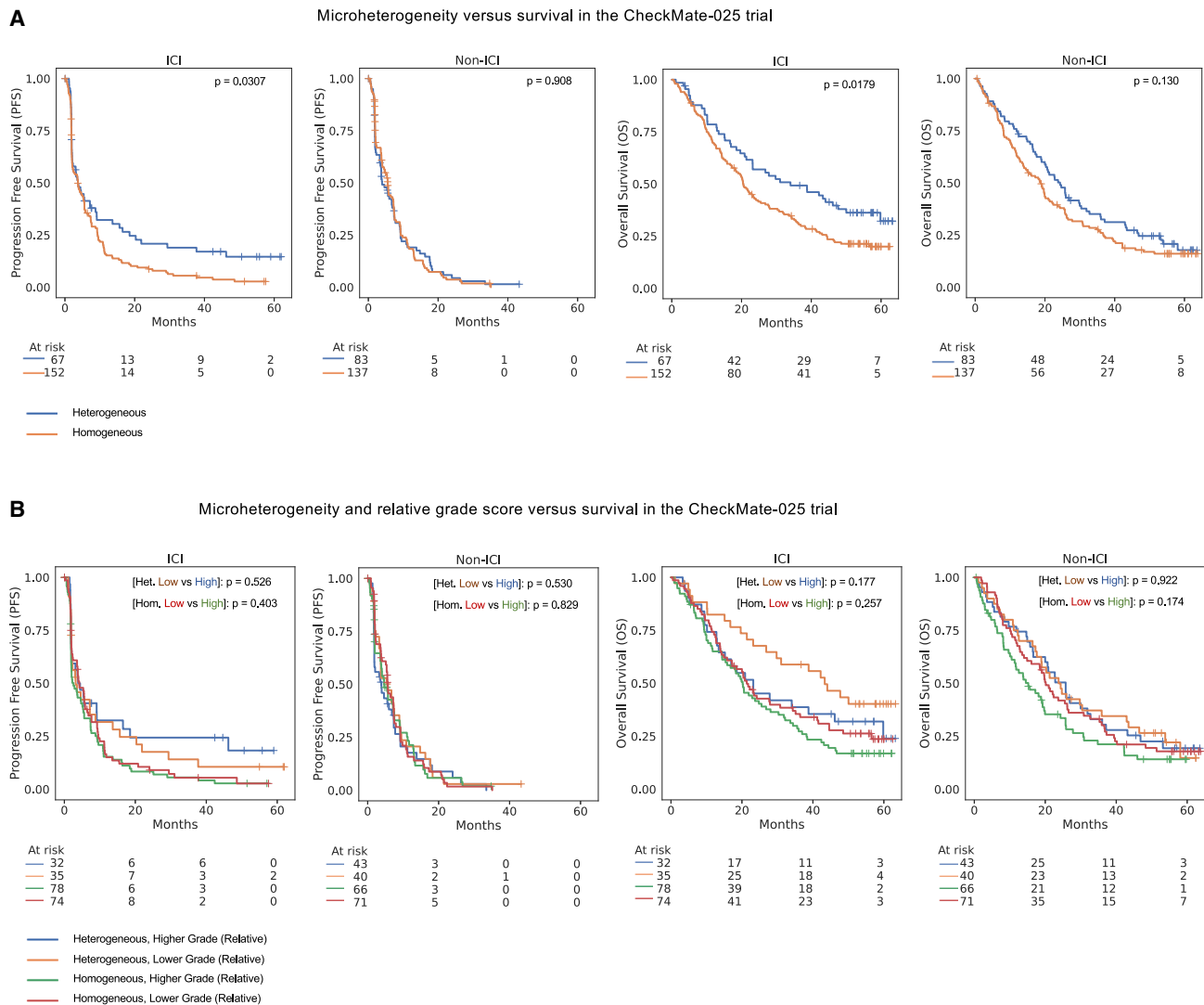
(C) Case-wise frequency of microheterogeneity vs. the maximum observed f-weighted heterogeneity score, which describes the largest extent of heterogeneity observed in a patient. Statistics are aggregated within a given patient's set of scanned tissue blocks (1 slide per block).

(D) Case-wise frequency of microheterogeneity vs. standard deviation of grade score predictions within the same case. Statistics are aggregated within a given patient's set of scanned tissue blocks (1 slide per block). Pearson's Rho p values were calculated via exact distribution.

with microheterogeneity, highly infiltrated cases associated with improved OS only in the ICI arm (Figure 5C;  $p = 0.0220$ , log-rank test). This subset of ICI-treated patients also demonstrated a consistent trend in improved PFS, but it did not reach statistical significance (Figure 5B;  $p = 0.0662$ , log-rank test).

We also compared the performance of predictive models that exclusively use image-derived or previously nominated molecular

features.<sup>13</sup> For OS in the ICI arm of CM-025, models using computer vision features had similar performance to those only using genomic features (*PBRM1* LOF, 9p21.3 deletion) (Figures S21–S30; STAR Methods). Moreover, combining these features resulted in net improvements while retaining consistent parameter associations (i.e., *PBRM1* LOF and microheterogeneity each retained positive coefficient weights). We lastly introduced clinical



**Figure 4. Inferred patterns of grade microheterogeneity associate with improved survival in the CM-025 cohort but only for ICI-treated patients**

(A) Kaplan-Meier curves for overall survival (OS) and progression-free survival (PFS) in the CM-025 cohort based on the presence of microheterogeneity. Significance values were calculated via log-rank test.

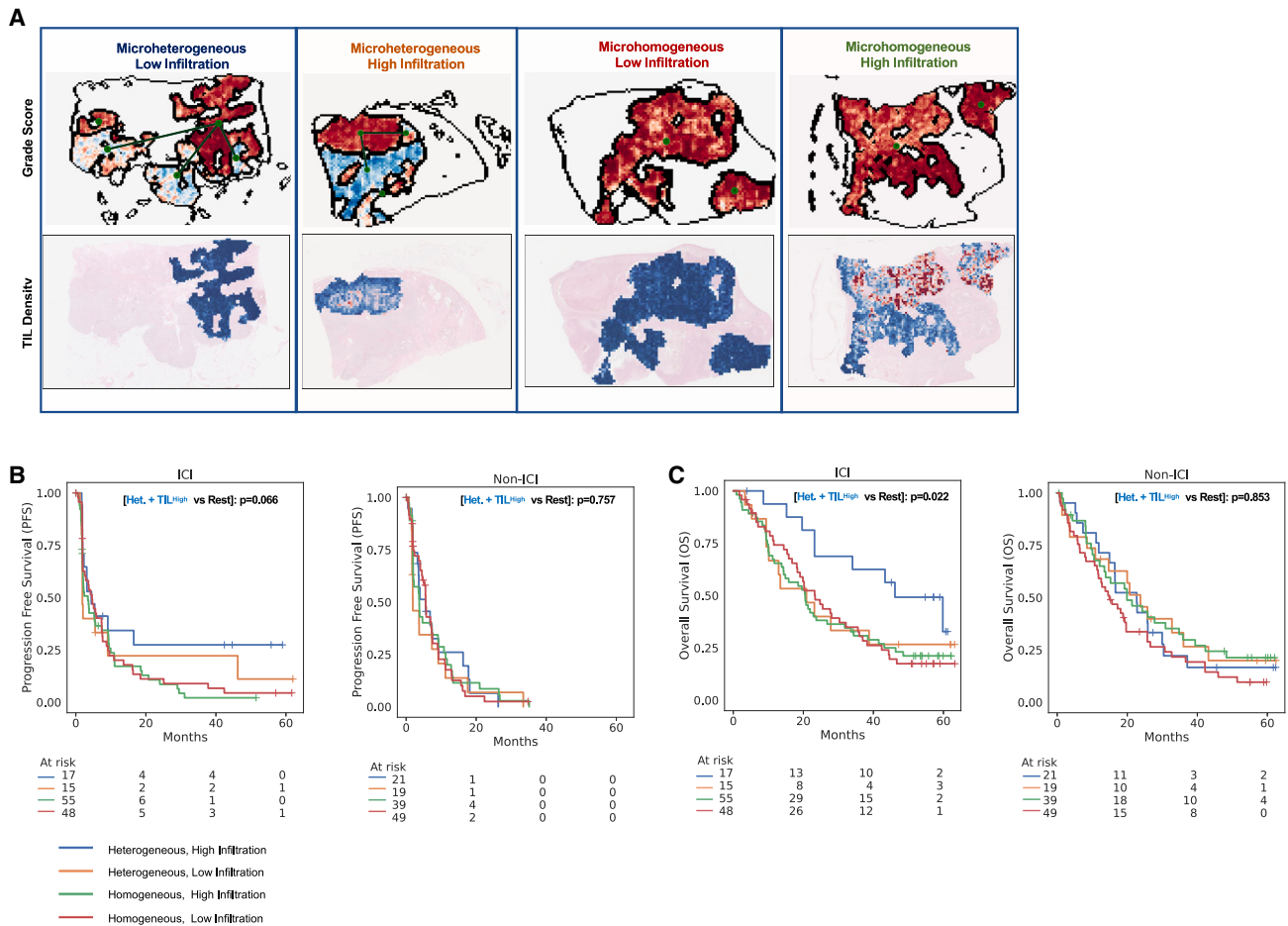
(B) Kaplan-Meier curves for OS and PFS in the CM-025 cohort based on the presence of microheterogeneity, stratifying further based on relative grade score within a grade microheterogeneity category. Significance values were calculated via (pairwise) log-rank test.

risk covariates into a full parameter model, which produced further improvements to C-Index metrics (Figures S24 and S28; STAR Methods). Taken together, tandem consideration of tumor-intrinsic spatial microheterogeneity and TIL features in WSIs learned by the computer vision models captured meaningful representations of selective ICI response.

### Tumor-immune interactions are more extensive and involve greater CD8<sup>+</sup> PD-1 activation in advanced ccRCC

To more precisely understand the tumor-immune spatial interactions identified from WSIs and linked to selective ICI response, we evaluated advanced ccRCC tumors with paired H&E and

mIF images derived from the same tissue (markers = {PAX8, CD8, DAPI, PD1, PDL1, FOXP3}).<sup>26,27</sup> To describe spatial phenotypes, we built a nearest-neighbor graph of CD8<sup>+</sup> and tumor cells and classified cells as “tumor-immune interacting” if they were adjacent to a distinct cell type in the graph (Figures S31 and S32; STAR Methods). Through analysis of regions with high tumor-immune interaction density from each patient (STAR Methods), we observed that microheterogeneous tumors had higher CD8<sup>+</sup> cell density, while tumor cell density was similar between heterogeneous and homogeneous cases (Figures S33–S35). The frequency of tumor cells adjacent to CD8<sup>+</sup> cells was higher in heterogeneous cases, suggesting a greater presence of “desert”-like regions of non-infiltrated tumor tissue in



**Figure 5. Combining computationally inferred tumor and immune states identifies a further subset of responders to ICI in the CM-025 trial**

(A) Representative examples of four classes of patients identified using computational inference, based on the presence of grade microheterogeneity, and the relative abundance of TILs in high-grade tumor regions (TIL density). Top row: representative RAG plots based on grade score (blue: lower score, red: higher score). Bottom row: representative inferred TIL densities (blue: lower infiltration, red: higher infiltration; uncolored: lower-grade regions not considered for TIL density evaluation).

(B and C) Kaplan-Meier curves for PFS and OS in both arms of the CM-025 trial based on the groups demonstrated in (A). Significance values were calculated via log-rank test between “microheterogeneous and high infiltration” patients and all remaining patients within a trial arm.

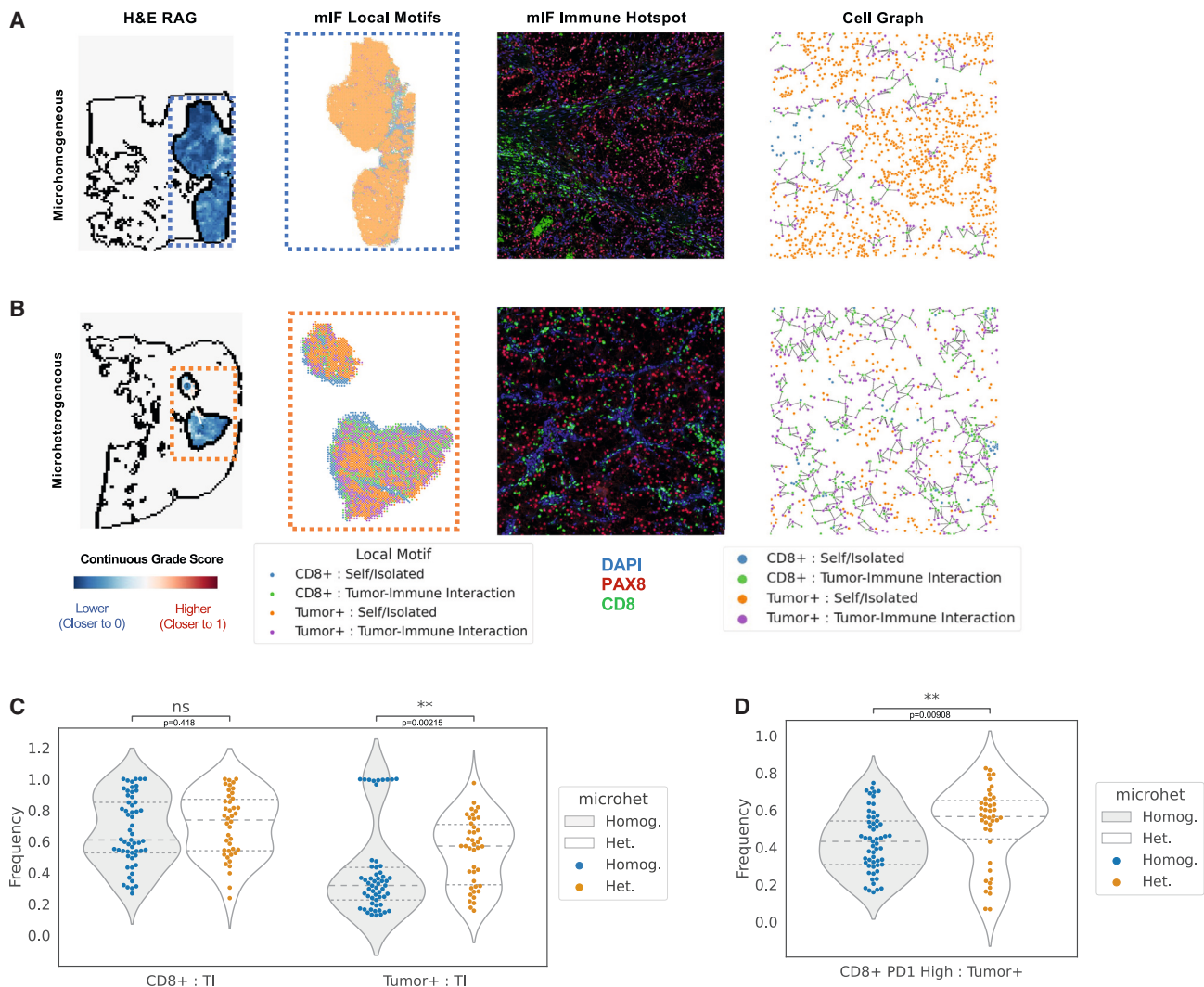
homogeneous cases (Figure 6C;  $p = 0.00215$  [tumor+, tumor-immune], Wilcoxon rank-sum test). In contrast, the frequency of CD8<sup>+</sup> cells adjacent to tumor cells was similar between heterogeneous and homogeneous cases (Figure 6C;  $p = 0.418$  [CD8<sup>+</sup> tumor-immune]). Thus, the observed increase in tumor-immune interaction frequency in microheterogeneous tumors resulted from increased infiltration deeper within tumor-dense regions rather than a uniform increase across the tumor microenvironment.

We lastly asked whether any of these observed differences related to tumor-immune cell subtypes, specifically PD1-low/high CD8<sup>+</sup> and PDL1-low/-high tumor cells. In general, PD1-high CD8<sup>+</sup> cells were common, and PDL1-high tumor cells were sparser (PD1 median frequency [freq.] = 0.480, PDL1 median freq. = 0.150; Figure S35). Within CD8<sup>+</sup> cells engaged in a tumor interaction, microheterogeneous cases had a higher freq. of PD1 high cells compared with homogeneous cases (Fig-

ure 6D;  $p = 0.00908$ , Wilcoxon rank-sum test). Thus, spatial microheterogeneity structures in ccRCC, which exhibited enrichment for ICI response, may foster an immune compartment that is both more tumor experienced and abundant.

## DISCUSSION

Simultaneous quantitative measurements of key tumor and microenvironmental properties that represent distinct modes of oncogenesis, evolution, and immune evasion may unlock insights in ccRCC biology and potential modes of patient stratification. To this end, we developed a series of biologically informed neural network models to perform spatially aware computer vision analysis on multiple independent ccRCC cohorts. In doing so, we produced a continuous, quantitative, and automated grading approach that reproduces existing manual histological assessments of nuclear grade and provides comparable



**Figure 6. Exploration of microheterogeneity implications in paired multiplexed imaging data**

(A and B) Representative examples of a microhomogeneous case (A) and a microheterogeneous case (B). H&E Rag: slide-level representation of H&E-inferred tumor and grade properties. mIF local motifs indicates the primary cell type and context present within overlapping 200 pixel windows. mIF immune hotspot: representative example of an area of high tumor-immune interaction density. Cell graph: visualization of tumor and CD8<sup>+</sup> cells and their interaction context; edges are drawn between interacting tumor and CD8<sup>+</sup> cells (nearest neighbors in a Delaunay triangulation).

(C) Comparison of the frequency of tumor-immune interaction within CD8<sup>+</sup> cells (left) and tumor cells (right) vs. H&E-inferred microheterogeneity status (TI, “tumor-immune” interaction context).

(D) Comparison of the frequency of PD1-high within CD8<sup>+</sup> cells that interact with tumor cells vs. H&E-inferred microheterogeneity status. Significance was calculated via Wilcoxon rank-sum test.

\*:  $p \leq 0.05$ , \*\*:  $p \leq 0.001$ , \*\*\*:  $p \leq 0.0001$ , \*\*\*\*:  $p \leq 0.00001$ .

prognostic value without interobserver variability. More importantly, by formalizing tumor phenotype predictions into spatial maps and subsequent RAGs using a single WSI per patient, we discovered histological ITH properties not feasibly measurable by manual review that were informative for multiple phenomena and that represented patterns present throughout a patient’s tumor. Namely, the graph-based microheterogeneity feature contained additional prognostic value beyond established pathology scores, as well as predictive value specifically for response to ICI in CM-025. Furthermore, this feature corre-

lated with a series of molecular characteristics, such as *PBRM1* LOF, and thus may provide a unified histological representation for connecting clinically relevant molecular features. Upon simultaneously integrating tumor and immune microenvironmental features, we identified a subset of ICI responders enriched for microheterogeneity and a higher degree of TILs. Moreover, microheterogeneity in advanced ccRCC associated with greater PD1 activation in CD8<sup>+</sup> lymphocytes and a greater extent of tumor-immune interaction, suggesting a more active tumor-immune interaction landscape that is

more likely to respond to ICI. Taken together, these findings suggest that tumor and immune features of ccRCC can be jointly considered in a spatially aware manner to guide biological and clinical investigations using widely available H&E WSIs. Indeed, future study of other cancers that exhibit similar phenomena of phenotypic heterogeneity (e.g., Gleason phenotypes in prostate cancer) may benefit from a similar series of approaches as those followed herein.

### Limitations of the study

There are several challenges and limitations to this analysis. The histological data analyzed from CM-025 consisted of pre-treatment primary tumor samples and thus may differ from the tumor state at the time of trial accrual due to ongoing tumor evolution. As such, the specimens we analyzed may be uncoupled from eventual metastatic progression. Similarly, larger sample sizes in additional clinical cohorts are necessary to generalize these findings regarding both the prognostic value of our continuous grading approach in untreated disease as well as our observations relating microheterogeneity and selective response to immunotherapy, and additional histologic features could be added to our model framework and RAG construction (e.g., necrosis, refined TIL subtypes, proximal vs. distal patterns, stromal heterogeneity) to more completely describe the complexities of ccRCC biology. Relatedly, our approach to validating spatial features is limited by a lack of exhaustive, pixel-level annotation of all images. Future work could address this via larger-scale pathologist evaluation to reach consensus annotations. Additionally, occurrence of microheterogeneity cannot be evaluated like typical histological features (e.g., tumor-stroma interfaces). Nuclear grade itself is a composite of several imprecise tumor microenvironment phenotypes and as such does not have a tractable definition at the resolution required for consistent pathologist annotation generation.

Similarly, while we were able to produce a model of TILs that aligned with complementary CD8 immunofluorescence results, our estimation of TILs is limited by its inability to distinguish morphologically similar cells (e.g., CD8<sup>+</sup> vs. CD4<sup>+</sup>, B vs. T cells). Although we derived candidate associations between specific molecular events and complex spatial patterns derived in this study, further work is necessary to understand which combinations of molecular features result in these microheterogeneity patterns using both model systems and additional patient cohorts and the relationship between molecular ITH and this WSI property. While we were able to provide an orthogonal glimpse at the specific cell populations that might underlie the tumor-immune phenotypes associated with microheterogeneity, our analysis of paired mIF and H&E data also had limitations. In particular, the sample size was small and composed of varying biopsy sites, and larger paired cohorts representing diverse biopsy sites using emerging highly multiplex spatial imaging technologies will guide extensions of these tumor-immune interactions and their relevance to selective immunotherapy responsiveness in ccRCC. Finally, as multiple-instance learning and saliency mapping methods continue to evolve in digital pathology domains, strategies that incorporate these approaches (particularly in the context of phenotypic variation similar to

grade microheterogeneity) to complement our human-interpret-able approach may further enhance the ability to understand spatial structures that determine tumor-immune interactions in ccRCC or, potentially, in other cancers.<sup>28–31</sup>

### Conclusion

Taken together, we propose spatially aware deep-learning models that build upon inference of known histological features (e.g., nuclear grade) to learn interacting distinct features (graph-based microheterogeneity) and that reveal distinct oncogenic paths and ICI response phenotypes in ccRCC. The occurrence of microheterogeneity and its predictive capacity for PFS and OS in ICI warrants further study, including via model systems, to unravel how this phenomenon influences tumor evolution and anti-tumor immunity. Broadly, the use of biologically guided computer vision strategies for cancer histopathology to automatically infer tumor and microenvironmental features, their respective higher-order interactions, and their relationship to molecular and clinical states may have general utility across tumor types and therapeutic modalities.

### STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- **KEY RESOURCES TABLE**
- **RESOURCE AVAILABILITY**
  - Lead contact
  - Materials availability
  - Data and code availability
- **EXPERIMENTAL MODEL AND STUDY PARTICIPANT DETAILS**
  - Clinical cohorts
  - Data acquisition
- **METHOD DETAILS**
  - Image quality control
  - Cross validation
  - Tumor classification
  - Tumor grade classification
  - Inference post-processing
  - Phenotype segmentation
  - Adjacency descriptions
  - Heterogeneity description
  - Reference slide modeling
  - Tumor infiltrating lymphocyte inference
  - Tumor infiltration classification
  - Truncating mutation categorization
  - Molecular ITH metrics
  - Survival analysis
  - Image registration
  - Multiplexed immunofluorescence image preprocessing
  - Cell phenotype calling
  - Cell graph construction
  - Immune hotspot analysis
- **QUANTIFICATION AND STATISTICAL ANALYSIS**

**SUPPLEMENTAL INFORMATION**

Supplemental information can be found online at <https://doi.org/10.1016/j.xcrm.2023.101189>.

**ACKNOWLEDGMENTS**

This work was supported in part by Bristol Myers Squibb through its International Immuno-Oncology Network. The results shown here are also in part based on data generated by the TCGA Research Network: <https://www.cancer.gov/tcga>. This work was supported by Dana-Farber/Harvard Cancer Center Kidney SPORE (P50CA101942-15). This work was also supported in part by the National Institutes of Health (R37 CA222574, R01CA227388, and P50CA191842 [E.M.V.A.]; T32GM007753 [N.M.]; P30CA016359 [D.A.B.]; and F31 CA250136 [J.N.]); a Dunkin' Donuts Breakthrough Grant (J.N. and E.M.V.A.); and the DOD Academy of Kidney Cancer Investigators (KC190128) (D.A.B.). T.K.C. is supported in part by the Dana-Farber/Harvard Cancer Center Kidney SPORE (2P50CA101942-16) and Program 5P30CA006516-56, the Kohlberg Chair at Harvard Medical School and the Trust Family, Michael Brigham, Pan Mass Challenge, the Hinda and Arthur Marcus Fund, and Loker Pinard Funds for Kidney Cancer Research at DFCI.

**AUTHOR CONTRIBUTIONS**

Conceptualization, J.N., E.M.V.A., S.S., D.A.B., T.K.C.; formal analysis, J.N.; interpretation of the data, S.N.H., B.J., N.M., C.L., Z.B., T.D., D.A.B., K.B., and B.M.T.; resources, R.U. and J.R.; data curation and annotation, T.D., Z.B., C.L., S.S., S.R., K.F., and B.S.; writing – original draft, J.N. and E.M.V.A.; writing – review & editing, S.S., D.A.B., S.S., T.K.C., B.J., N.M., S.N.H., K.B., and B.M.T.; visualization, J.N.; supervision, E.M.V.A., S.S., and T.K.C.; funding acquisition, E.M.V.A., S.S., and T.K.C.

**DECLARATION OF INTERESTS**

T.K.C. reports institutional and personal paid and unpaid support for research, advisory boards, consultancy, and honoraria from Alkermes, AstraZeneca, Aravive, Aveo, Bayer, Bristol Myers-Squibb, Calithera, Circle Pharma, Deciphera Pharmaceuticals, Eisai, EMD Serono, Exelixis, GlaxoSmithKline, Gilead, IQVA, Infinity, Ipsen, Jansen, Kanaph, Lilly, Merck, Nikang, Nuscan, Novartis, Oncohost, Pfizer, Roche, Sanofi/Aventis, Scholar Rock, Surface Oncology, Takeda, Tempest, UpToDate, and CME events (Peerview, Onclive, MJH, CCO, and others), outside the submitted work; institutional patents filed on molecular alterations and immunotherapy response/toxicity and ctDNA; and equity in Tempest, Pionyr, Osel, Precede Bio, CureResponse, and InnDura. T.K.C. serves on the committees of NCCN, GU Steering Committee, and ASCO/ESMO. Medical writing and editorial assistance support may have been funded by communications companies in part. T.K.C. does not report any speaker's bureau. T.K.C. has mentored several non-US citizens on research projects with potential funding (in part) from non-US sources/foreign components. The institution (Dana-Farber Cancer Institute) may have received additional independent funding of drug companies and/or royalties potentially involved in research around the subject matter. E.M.V.A. reports advisory/consulting with Tango Therapeutics, Genome Medical, Invitae, Monte Rosa, Enara Bio, Manifold Bio, Riva Therapeutics, Serinus Bio, and Janssen; research support from Novartis and BMS; equity in Tango Therapeutics, Genome Medical, Syapse, Manifold Bio, Monte Rosa, Enara Bio, Riva Therapeutics, and Serinus Bio; institutional patents filed on chromatin mutations and immunotherapy response and on methods for clinical interpretation; and intermittent legal consulting on patents for Foaley & Hoag. D.A.B. reports personal fees from LM Education and Exchange, Advovate Strategies, MDedge, Cancer Network, Cancer Expert Now, Onclive, Catenion, and AVEO and grants and personal fees from Exelixis outside the submitted work. C.L. reports research funding from Genentech/imCORE. Z.B. reports research funding from Bristol-Myers Squibb and Genentech/imCORE and honoraria from UpToDate. S.S. reports grants from Exelixis, grants from Bristol-Myers Squibb, personal fees from Merck, grants and personal fees from AstraZeneca, personal fees from CRISPR Therapeutics, personal fees from

NCI, and personal fees from AACR as well as a patent for Biogenex with royalties paid. K.B. has consulted for Related Sciences (RS) outside of the scope of this work. S.R. receives research funding from Bristol-Myers Squibb and KITE/Gilead and is a member of the SAB for Immunitas Therapeutics.

Received: January 19, 2023

Revised: June 20, 2023

Accepted: August 16, 2023

Published: September 19, 2023

**REFERENCES**

- Hsieh, J.J., Purdue, M.P., Signoretti, S., Swanton, C., Albiges, L., Schmieder, M., Heng, D.Y., Larkin, J., and Ficarra, V. (2017). Renal cell carcinoma. *Nat. Rev. Dis. Prim.* 3, 17009. <https://doi.org/10.1038/nrdp.2017.9>.
- Gerlinger, M., Rowan, A.J., Horswell, S., Math, M., Larkin, J., Endesfelder, D., Gronroos, E., Martinez, P., Matthews, N., Stewart, A., et al. (2012). Intratumor heterogeneity and branched evolution revealed by multiregion sequencing. *N. Engl. J. Med.* 366, 883–892. <https://doi.org/10.1056/NEJMoa1113205>.
- Turajlic, S., Xu, H., Litchfield, K., Rowan, A., Chambers, T., Lopez, J.I., Nicol, D., O'Brien, T., Larkin, J., Horswell, S., et al. (2018). Tracking Cancer Evolution Reveals Constrained Routes to Metastases: TRACERx Renal. *Cell* 173, 581–594.e12. <https://doi.org/10.1016/j.cell.2018.03.057>.
- Mitchell, T.J., Turajlic, S., Rowan, A., Nicol, D., Farmery, J.H.R., O'Brien, T., Martincorena, I., Tarpey, P., Angelopoulos, N., Yates, L.R., et al. (2018). Timing the Landmark Events in the Evolution of Clear Cell Renal Cell Cancer: TRACERx Renal. *Cell* 173, 611–623.e17. <https://doi.org/10.1016/j.cell.2018.02.020>.
- Turajlic, S., Xu, H., Litchfield, K., Rowan, A., Horswell, S., Chambers, T., O'Brien, T., Lopez, J.I., Watkins, T.B.K., Nicol, D., et al. (2018). Deterministic Evolutionary Trajectories Influence Primary Tumor Growth: TRACERx Renal. *Cell* 173, 595–610.e11. <https://doi.org/10.1016/j.cell.2018.03.043>.
- Delahunt, B., Egevad, L., Montironi, R., and Srigley, J.R. (2013). International Society of Urological Pathology (ISUP) consensus conference on renal neoplasia: rationale and organization. *Am. J. Surg. Pathol.* 37, 1463–1468. <https://doi.org/10.1097/PAS.0b013e318299f14a>.
- Zhao, Y. et al. Selection of Metastasis Competent Subclones in the Tumour Interior: TRACERx Renal. <https://doi.org/10.21203/rs.3.rs-61979/v1>.
- Sirohi, D., Chipman, J., Barry, M., Albertson, D., Mahlow, J., Liu, T., Raps, E., Haaland, B., Sayegh, N., Li, H., et al. (2022). Histologic Growth Patterns in Clear Cell Renal Cell Carcinoma Stratify Patients into Survival Risk Groups. *Clin. Genitourin. Cancer* 20, e233–e243. <https://doi.org/10.1016/j.clgc.2022.01.005>.
- Ghatalia, P., Gordetsky, J., Kuo, F., Dulaimi, E., Cai, K.Q., Devarajan, K., Bae, S., Naik, G., Chan, T.A., Uzzo, R., et al. (2019). Prognostic impact of immune gene expression signature and tumor infiltrating immune cells in localized clear cell renal cell carcinoma. *J. Immunother. Cancer* 7, 139. <https://doi.org/10.1186/s40425-019-0621-1>.
- Choueiri, T.K., and Motzer, R.J. (2017). Systemic Therapy for Metastatic Renal-Cell Carcinoma. *N. Engl. J. Med.* 376, 354–366. <https://doi.org/10.1056/NEJMra1601333>.
- Miao, D., Margolis, C.A., Gao, W., Voss, M.H., Li, W., Martini, D.J., Norton, C., Bossé, D., Wankowicz, S.M., Cullen, D., et al. (2018). Genomic correlates of response to immune checkpoint therapies in clear cell renal cell carcinoma. *Science* 359, 801–806. <https://doi.org/10.1126/science.aan5951>.
- Braun, D.A., Ishii, Y., Walsh, A.M., Van Allen, E.M., Wu, C.J., Shukla, S.A., and Choueiri, T.K. (2019). Clinical Validation of PBRM1 Alterations as a Marker of Immune Checkpoint Inhibitor Response in Renal Cell Carcinoma. *JAMA Oncol.* 5, 1631–1633. <https://doi.org/10.1001/jamaoncol.2019.3158>.

13. Braun, D.A., Hou, Y., Bakouny, Z., Ficial, M., Sant' Angelo, M., Forman, J., Ross-Macdonald, P., Berger, A.C., Jegede, O.A., Elagina, L., et al. (2020). Interplay of somatic alterations and immune infiltration modulates response to PD-1 blockade in advanced clear cell renal cell carcinoma. *Nat. Med.* 26, 909–918. <https://doi.org/10.1038/s41591-020-0839-y>.
14. Diao, J.A., Wang, J.K., Chui, W.F., Mountain, V., Gullapally, S.C., Srinivasan, R., Mitchell, R.N., Glass, B., Hoffman, S., Rao, S.K., et al. (2021). Human-interpretable image features derived from densely mapped cancer pathology slides predict diverse molecular phenotypes. *Nat. Commun.* 12, 1613. <https://doi.org/10.1038/s41467-021-21896-9>.
15. Tellez, D., Litjens, G., van der Laak, J., and Ciompi, F. (2018). Neural Image Compression for Gigapixel Histopathology Image Analysis. Preprint at arXiv. <https://doi.org/10.10109/TPAMI.2019.2936841>.
16. He, K., Zhang, X., Ren, S., and Sun, J. (2015). Preprint at Deep Residual Learning for Image Recognition. <https://doi.org/10.1048550/arXiv.1512.03385>.
17. Graham, S., Vu, Q.D., Raza, S.E.A., Azam, A., Tsang, Y.W., Kwak, J.T., and Rajpoot, N. (2019). Hover-Net: Simultaneous segmentation and classification of nuclei in multi-tissue histology images. *Med. Image Anal.* 58, 101563. <https://doi.org/10.1016/j.media.2019.101563>.
18. Acosta, P.H., Panwar, V., Jarmale, V., Christie, A., Jasti, J., Margulis, V., Rakheja, D., Chevillat, J., Leibovich, B.C., Parker, A., et al. (2022). Intratumoral Resolution of Driver Gene Mutation Heterogeneity in Renal Cancer Using Deep Learning. *Cancer Res.* 82, 2792–2806. <https://doi.org/10.1158/0008-5472.CAN-21-2318>.
19. Fu, Y., Jung, A.W., Torne, R.V., Gonzalez, S., Vöhringer, H., Shmatko, A., Yates, L., Jimenez-Linan, M., Moore, L., and Gerstung, M. (2019). Pan-cancer computational histopathology reveals mutations, tumor composition and prognosis. Preprint at bioRxiv, 813543. <https://doi.org/10.101101/813543>.
20. Baietti, M.F., Zhao, P., Crowther, J., Sewduth, R.N., De Troyer, L., Debiec-Rychter, M., and Sablina, A.A. (2021). Loss of 9p21 Regulatory Hub Promotes Kidney Cancer Progression by Upregulating HOXB13. *Mol. Cancer Res.* 19, 979–990. <https://doi.org/10.1158/1541-7786.MCR-20-0705>.
21. El-Mokadem, I., Fitzpatrick, J., Rai, B., Cunningham, J., Pratt, N., Fleming, S., and Nabi, G. (2014). Significance of chromosome 9p status in renal cell carcinoma: a systematic review and quality of the reported studies. *Bio-Med Res. Int.* 2014, 521380. <https://doi.org/10.1155/2014/521380>.
22. Bakouny, Z., Braun, D.A., Shukla, S.A., Pan, W., Gao, X., Hou, Y., Flaifel, A., Tang, S., Bosma-Moody, A., He, M.X., et al. (2021). Integrative molecular characterization of sarcomatoid and rhabdoid renal cell carcinoma. *Nat. Commun.* 12, 808. <https://doi.org/10.1038/s41467-021-21068-9>.
23. Endesfelder, D., Burrell, R., Kanu, N., McGranahan, N., Howell, M., Parker, P.J., Downward, J., Swanton, C., and Kschischo, M. (2014). Chromosomal instability selects gene copy-number variants encoding core regulators of proliferation in ER+ breast cancer. *Cancer Res.* 74, 4853–4863. <https://doi.org/10.1158/0008-5472.CAN-13-2664>.
24. Motzer, R.J., Escudier, B., McDermott, D.F., George, S., Hammers, H.J., Srinivas, S., Tykodi, S.S., Sosman, J.A., Procopio, G., Plimack, E.R., et al. (2015). Nivolumab versus Everolimus in Advanced Renal-Cell Carcinoma. *N. Engl. J. Med.* 373, 1803–1813. <https://doi.org/10.1056/NEJMoa1510665>.
25. Fridman, W.H., Zitvogel, L., Sautès-Fridman, C., and Kroemer, G. (2017). The immune contexture in cancer prognosis and treatment. *Nat. Rev. Clin. Oncol.* 14, 717–734. <https://doi.org/10.1038/nrclinonc.2017.101>.
26. Carey, C.D., Gusenleitner, D., Lipschitz, M., Roemer, M.G.M., Stack, E.C., Gjini, E., Hu, X., Redd, R., Freeman, G.J., Neuberg, D., et al. (2017). Topological analysis reveals a PD-L1-associated microenvironmental niche for Reed-Sternberg cells in Hodgkin lymphoma. *Blood* 130, 2420–2430. <https://doi.org/10.1182/blood-2017-03-770719>.
27. Griffin, G.K., Weirather, J.L., Roemer, M.G.M., Lipschitz, M., Kelley, A., Chen, P.H., Gusenleitner, D., Jeter, E., Pak, C., Gjini, E., et al. (2021). Spatial signatures identify immune escape via PD-1 as a defining feature of T-cell/histiocyte-rich large B-cell lymphoma. *Blood* 137, 1353–1364. <https://doi.org/10.1182/blood.2020006464>.
28. Oakden-Rayner, L., Dunnmon, J., Carneiro, G., and Ré, C. (2020). Hidden Stratification Causes Clinically Meaningful Failures in Machine Learning for Medical Imaging. *Proc. ACM Conf. Health Inference Learn.* 2020, 151–159. <https://doi.org/10.1145/3368555.3384468>.
29. Ghassemi, M., Oakden-Rayner, L., and Beam, A.L. (2021). The false hope of current approaches to explainable artificial intelligence in health care. *Lancet. Digit. Health* 3, e745–e750. [https://doi.org/10.1016/S2589-7500\(21\)00208-9](https://doi.org/10.1016/S2589-7500(21)00208-9).
30. Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi, A., Acharya, U.R., et al. (2021). A review of uncertainty quantification in deep learning: Techniques, applications and challenges. *Inf. Fusion* 76, 243–297. <https://doi.org/10.1016/j.infus.2021.05.008>.
31. Saporta, A., Gui, X., Agrawal, A., Pareek, A., Truong, S.Q.H., Nguyen, C.D.T., Ngo, V.D., Seekins, J., Blankenberg, F.G., Ng, A.Y., et al. (2022). Benchmarking saliency methods for chest X-ray interpretation. *Nat. Mach. Intell.* 4, 867–878. <https://doi.org/10.1038/s42256-022-00536-x>.
32. Paszke, A., et al. (2019). PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems*, 32, H. Wallach, et al., eds.
33. Janowczyk, A., Zuo, R., Gilmore, H., Feldman, M., and Madabhushi, A. (2019). HistoQC: An Open-Source Quality Control Tool for Digital Pathology Slides. *JCO Clin. Cancer Inform.* 3, 1–7. <https://doi.org/10.1200/CCI.18.00157>.
34. Greenwald, N.F., Miller, G., Moen, E., Kong, A., Kagel, A., Dougherty, T., Fullaway, C.C., McIntosh, B.J., Leow, K.X., Schwartz, M.S., et al. (2022). Whole-cell segmentation of tissue images with human-level performance using large-scale data annotation and deep learning. *Nat. Biotechnol.* 40, 555–565. <https://doi.org/10.1038/s41587-021-01094-0>.
35. Davidson-Pilon, C. (2021). Lifelines, Survival Analysis in Python. <https://doi.org/10.5281/zenodo.5745573>.
36. van der Walt, S., Schönberger, J.L., Nunez-Iglesias, J., Boulogne, F., Warner, J.D., Yager, N., Gouillart, E., and Yu, T.; scikit-image contributors (2014). scikit-image: image processing in Python. *PeerJ* 2, e453. <https://doi.org/10.7717/peerj.453>.
37. Falcon, W. (2019). The PyTorch Lightning team. Pytorch lightning. <https://doi.org/10.5281/zenodo.3828935>.
38. Dagher, J., Delahunt, B., Rioux-Leclercq, N., Egevad, L., Srigley, J.R., Coughlin, G., Dungleon, N., Gianduzzo, T., Kua, B., Malone, G., et al. (2017). Clear cell renal cell carcinoma: validation of World Health Organization/International Society of Urological Pathology grading. *Histopathology* 71, 918–925. <https://doi.org/10.1111/his.13311>.
39. Trémeau, A., and Colantoni, P. (2000). Regions adjacency graph applied to color image segmentation. *IEEE Trans. Image Process.* 9, 735–744. <https://doi.org/10.10109/83.841950>.
40. Gamper, J., Alemi Koohbanani, N., Benet, K., Khuram, A., and Rajpoot, N. (2019). PanNuke: An Open Pan-Cancer Histology Dataset for Nuclei Instance Segmentation and Classification. *Digital Pathology*, 11–19. [https://doi.org/10.1007/978-3-030-23937-4\\_2](https://doi.org/10.1007/978-3-030-23937-4_2).
41. Rosenthal, J., Carelli, R., Omar, M., Brundage, D., Halbert, E., Nyman, J., Hari, S.N., Van Allen, E.M., Marchionni, L., Umetsu, R., and Loda, M. (2022). Building Tools for Machine Learning and Artificial Intelligence in Cancer Research: Best Practices and a Case Study with the PathML Toolkit for Computational Pathology. *Mol. Cancer Res.* 20, 202–206. <https://doi.org/10.1158/1541-7786.MCR-21-0665>.
42. Liu, C.C., et al. (2022). Robust phenotyping of highly multiplexed tissue imaging data using pixel-level clustering. *bioRxiv*. <https://doi.org/10.1101/2022.08.16.504171>.

43. Wolf, F.A., Angerer, P., and Theis, F.J. (2018). SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol.* *19*, 15. <https://doi.org/10.1186/s13059-017-1382-0>.
44. Blondel, V.D., Guillaume, J.-L., Lambiotte, R., and Lefebvre, E. (2008). Fast unfolding of communities in large networks. *J. Stat. Mech.* *2008*, P10008. <https://doi.org/10.1088/1742-5468/2008/10/p10008>.
45. Virtanen, P., Gommers, R., Oliphant, T.E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., et al. (2020). SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat. Methods* *17*, 352–272. <https://doi.org/10.1038/s41592-020-0772-5>.
46. Charlier, F. Statannotations: add statistical significance annotations on seaborn plots. Further development of statannot, with bugfixes, new features, and a different API. (Github). <https://doi.org/10.5281/zenodo.7213391>
47. Waskom, M. (2021). seaborn: statistical data visualization. *J. Open Source Softw.* *6*, 3021. <https://doi.org/10.21105/joss.03021>.

## STAR★METHODS

### KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
<b>Biological samples</b>		
CheckMate-025 Clear cell renal cell carcinoma H&E diagnostic whole-slide images	Bristol-Myers Squibb	N/A
Clear cell renal cell carcinoma multiplexed immunofluorescence whole-slide images	This paper	N/A
Clear cell renal cell carcinoma H&E diagnostic whole-slide images	This paper	N/A
<b>Deposited data</b>		
TCGA-KIRC H&E diagnostic whole-slide images,	ISB-CGC Mirror	<a href="https://isb-cancer-genomics-cloud.readthedocs.io/en/latest/sections/data/TCGA-images.html?">https://isb-cancer-genomics-cloud.readthedocs.io/en/latest/sections/data/TCGA-images.html?</a>
TCGA-KIRC WES-seq, clinical, and mutational data	TCGA GDC Portal	<a href="https://gdc.cancer.gov/about-data/publications/pancanatlas">https://gdc.cancer.gov/about-data/publications/pancanatlas</a>
TCGA-KIRC revised clinical data	Ricketts et al., 2018	N/A
CheckMate-025 normalized bulk RNA-seq, clinical, and mutational data	Braun et al., <sup>13</sup>	EGA: EGAS00001004290, EGAS00001004291, EGAS00001004292
<b>Software and algorithms</b>		
python 3	<a href="https://www.python.org/downloads/">https://www.python.org/downloads/</a>	<a href="https://www.python.org/downloads/">https://www.python.org/downloads/</a>
PyTorch	Paszke et al., <sup>32</sup>	<a href="https://pytorch.org/get-started/locally/">https://pytorch.org/get-started/locally/</a>
HistoQC	Janowczyk et al., <sup>33</sup>	<a href="https://github.com/choosehappy/HistoQC">https://github.com/choosehappy/HistoQC</a>
PathFlow-MixMatch	Levy et al., 2020	<a href="https://github.com/jlevy44/PathFlow-MixMatch">https://github.com/jlevy44/PathFlow-MixMatch</a>
Mesmer	Greenwald et al., <sup>34</sup>	<a href="https://github.com/vanvalenlab/deepcell-tf">https://github.com/vanvalenlab/deepcell-tf</a>
Lifelines	Davidson-Pilon, <sup>35</sup>	<a href="https://lifelines.readthedocs.io/en/latest/">https://lifelines.readthedocs.io/en/latest/</a>
scikit-image	Van Der Walt et al., <sup>36</sup>	<a href="https://scikit-image.org/">https://scikit-image.org/</a>
Custom training and analysis code	This paper	<a href="https://doi.org/10.5281/zenodo.8244800">https://doi.org/10.5281/zenodo.8244800</a>

### RESOURCE AVAILABILITY

#### Lead contact

Further information and requests for resources and reagents should be directed to and will be fulfilled by the Lead Contact, Eliezer M. Van Allen ([eliezerm\\_vanallen@dfci.harvard.edu](mailto:eliezerm_vanallen@dfci.harvard.edu)).

#### Materials availability

This study did not generate new unique reagents.

#### Data and code availability

Code used to perform the analyses described in this study are available in a public github repository (<https://github.com/vanallenlab/microhet-paper>). For the TCGA-KIRC cohort, we obtained clinical data, and normalized bulk RNA and genomic sequencing from the GDC PanCanAtlas (<https://gdc.cancer.gov/about-data/publications/pancanatlas>), and downloaded whole-slide H&E stained diagnostic images from the ISB-CGC mirror of TCGA data. For the CheckMate 025 cohort, we directly obtained H&E stained diagnostic images via a Bristol Myers Squibb IION agreement, and used clinical and molecular data previously generated by Braun et al., 2020 (European Genome-Phenome Archive: EGAS00001004290, EGAS00001004291, EGAS00001004292). DFCI-PROFILE images were obtained via the Dana-Farber/Brigham and Women's PROFILE project. Multiplexed Immunofluorescence (and accompanying H&E images) were obtained via the Dana-Farber/Brigham and Women's ImmunoProfile project. Multi-block nephrectomy images were obtained from the Signoretti Lab. These images are available upon request with provision of IRB and adherence to institutional policies regarding storage security and other parameters. Restrictions apply to the availability of the raw in-house and external data (including patient whole-slide images and molecular sequencing data), which were used with institutional permission through IRB

approval for the current study. Please e-mail all requests for academic use of raw and processed data to the lead author. All requests will be evaluated based on institutional and departmental policies to determine whether the data requested is subject to intellectual property or patient privacy obligations. Data can only be shared for non-commercial academic purposes and will require a formal Data Use Agreement. Any additional information required to reanalyze the data reported in this work paper is available from the [lead contact](#) upon request.

### EXPERIMENTAL MODEL AND STUDY PARTICIPANT DETAILS

#### Clinical cohorts

Three distinct patient cohorts were used in the analysis: TCGA-KIRC, the CheckMate 025 (CM-025) phase III clinical trial (NCT01668784), and the DFCI-PROFILE of ccRCC patients from the Dana-Farber Cancer Institute (also under DFCI IRB #20–293 and 20–376). Additional datasets include a set of multi-block nephrectomy samples from Dana-Farber/Brigham (under DFCI IRB #20–293 and 20–376), and multiplexed immunofluorescence data from Dana-Farber/Brigham the ImmunoProfile project (under DFCI IRB #20–293 and 20–376).

#### Data acquisition

For the TCGA-KIRC cohort, we obtained clinical data, and normalized bulk RNA and genomic sequencing from the GDC PanCanAtlas (<https://gdc.cancer.gov/about-data/publications/pancanatlas>), and downloaded whole-slide H&E stained diagnostic images from the ISB-CGC mirror of TCGA data. For the CheckMate 025 cohort, we directly obtained H&E stained diagnostic images from the Signoretti Lab via an established Bristol Myers Squibb IION agreement, and used clinical and molecular data previously generated by Braun et al., 2020 (European Genome-Phenome Archive: EGAS00001004290, EGAS00001004291, EGAS00001004292). DFCI-PROFILE images were obtained via the Dana-Farber/Brigham and Women's PROFILE project. Multiplexed Immunofluorescence (and accompanying H&E images) were obtained via the Dana-Farber/Brigham and Women's ImmunoProfile project. Multi-block nephrectomy images were obtained from the Signoretti Lab. These images are available upon request with provision of IRB and adherence to institutional policies regarding storage security and other parameters.

### METHOD DETAILS

#### Image quality control

Quality control of H&E whole-slide images was performed using the HistoQC<sup>33</sup> toolkit. A full set of modules used is available as a '.ini' file. Custom examples used to train pen detection modules are available in a forked repository (<https://github.com/jmnyman/HistoQC>). Following HistoQC filtration, we further removed small slide images with fewer than 500 tiles (512px at 20X resolution).

#### Cross validation

Training and validation were performed in the DFCI-PROFILE cohort, with additional testing done in the unseen TCGA-KIRC and CM-025 cohorts. In subsequent finetuning experiments for tumor and grade classification, we utilized 4-fold cross-validation. Each dataset was split into folds at the patient level to ensure no patient-level information bled between training and validation contexts. Datasets were then subsampled to ensure a balanced label composition. Following balanced patient-level fold creation, we then sample a fixed number (per specified hyperparameters) of 512 pixel tiles (20X) from each patient slide to create folds that were label-balanced.

#### Tumor classification

A model to classify general renal cell carcinoma tumor tissue versus adjacent normal stromal tissue was trained using pixel-level pathologist annotations from the DFCI-PROFILE cohort ( $n = 36$  slides). Following quality control, training data slides were split into 512 pixel tiles (20X), and assigned labels according to pathologist annotations (i.e., whether a region contained tumorous tissue or not) (Figure S2). For the tumor classifier, direct slide annotation masks were obtained from the reviewing pathologist using the DFCI-PROFILE cohort, producing "strong" labels for the tumor vs. non-tumor content in a given region (pathologist drawing of borders in QuPath over slide image to produce a pixel map of tumor vs. non-tumor areas). A pretrained ResNet-50 neural network model was then finetuned using color jitter, and the highest performing model in a series of 4-fold cross validation was selected for subsequent inference. All neural network architecture and training code used the PyTorch and PyTorch-Lightning libraries, and 1–2 NVIDIA Tesla V100 GPU units on Google Cloud VM instances.<sup>32,37</sup> Hyperparameters used for training are available in the project repository.

#### Tumor grade classification

A second finetuned ResNet-50 neural network model was trained to distinguish low (G2) from high (G4) ccRCC cases from the DFCI-PROFILE cohort ( $n = 190$  slides) that was scored via manual pathologist review. This cohort contained samples collected prior to and following the adoption of the WHO/ISUP grading changes, and as such contains both Fuhrman and WHO/ISUP grades.<sup>6</sup> These conventions share significant overlap and are generally highly concordant.<sup>38</sup> Additionally, manual pathologist review for sarcomatoid and rhabdoid (S/R) tumor content was performed, and cases with S/R content were upgraded to G4 if previously assigned a lower grade

to ensure greater concordance with whom/ISUP guidelines. Following quality control, training data slides were again split into 512 pixel tiles (20X), and tiles were assigned labels according to pathologist annotation of the source slide. Tiles were only considered for model training if their predicted probability of containing tumor tissue was  $\geq 0.5$ , and slides were restricted to those with at least 500 putative tumor-containing tiles. We chose this cutoff to reflect an agnostic biological prior, intending to focus model training on tissue examples likely to contain mostly tumor content and thus learn appropriate signal relevant to nuclear grade. A pretrained ResNet-50 neural network model was again finetuned using heavy color jitter. An ensemble composed of each model trained in 4-fold cross validation was then used for subsequent inference, taking the average across all model softmax outputs to make predictions. All neural network architecture and training code used the PyTorch and PyTorch-Lightning libraries, and 1–2 NVIDIA Tesla V100 GPU units on Google Cloud VM instances.<sup>32,37</sup> Hyperparameters used for training are available in the project repository.

### Inference post-processing

Following tumor and grade inference, tile-level model scores were smoothed using uniform nearest neighbor averaging ( $n = 4$  nearest tiles). For grade score smoothing, we again considered tiles if their predicted probability of containing tumor tissue was  $\geq 0.5$  for the same reasons described previously. Smoothed tumor and grade scores were used for all downstream analysis. When training both the tumor and grade classifier models, we used a cross-entropy loss in PyTorch (CrossEntropyLoss) with default parameters. In both models, we formulated training as a binary classification problem, with tiles being assigned labels based on either pathologist-region annotation (tumor classifier) or slide-level pathologist annotation (grade classifier).

### Phenotype segmentation

Regions of tumor tissue were identified using watershed segmentation in scikit-image.<sup>36</sup> We first performed segmentation on smoothed tumor prediction scores, and classified regions as “tumor” if their average segment score was  $\geq 0.7$ . This more stringent cutoff was necessary to avoid a failure mode in which some regions with tumor probability scores between 0.5 and 0.7 did reliably contain tumor tissue. We support the performance of this approach in our evaluation of candidate tumor borders by pathologist review (see [Figure S2](#)). Following an initial watershed segmentation, regions were merged if they were similar (region score difference  $< 0.2$ ). These putative tumor regions were then considered for secondary segmentation using smoothed grade scores to identify regions of distinct grade. Specifically, the outputs of the first series of segmentations served as inputs for a second set of segmentations, such that a putative tumor region would become segmented into one or more subset regions based on the grade phenotypes it contains. Furthermore, a slide-average grade score was obtained using the average grade score across all putative tumor area. All thresholding and cutoffs were chosen prior to molecular or clinical association analyses.

### Adjacency descriptions

Following watershed segmentation of tumor and grade scores, we represented each slide as a region adjacency graph<sup>39</sup> (RAG) to describe the connectivity of each region produced, wherein directly contacting regions are assigned an edge in the graph. Small area nodes were removed ( $n < 50$  total tiles) following RAG construction. Subsequently, we performed a series of segmentation expansions to recover missing connectivity locally (e.g., regions that are visually in contact, but are separated by a thin layer of non-tumor predicted area), and also to describe long-range differences (e.g., regions that are distinct in grade score, but separated by 10+ tiles). RAG edges forming either directly or at an expansion distance of 1 tile were classified as “proximal”, and those forming at an expansion distance of up to 25 tiles were classified as “distal”. In analyses using TIL predictions, we only considered edges containing at least one node with an average grade score above 0.8 (see *Tumor Infiltration Classification*).

### Heterogeneity description

Patients with at least one proximal or distal RAG edge were classified as “heterogeneous”, and those lacking any edges as “homogeneous”. In analyses using TIL predictions, we only considered proximal/distal edges containing at least one node with an average grade score above 0.8 (see *Tumor Infiltration Classification*). To describe microheterogeneity continuously, we derived two related metrics. First was a “total-weighted” heterogeneity score (t-HS), we calculated a weighted sum of the number of RAG edges, wherein each edge was weighted by the total fractional tumor area occupied by the node pair involved in that edge (e.g., if two nodes comprise nearly all of the tumor area, that edge is highly weighted, while smaller regions contribute less to the sum). The second score was “f-weighted” (f-HS), and instead used the harmonic mean of the area fractions of each node in an edge to produce a weight, which describes both the contribution scale and balance of the nodes involved in that edge.

### Reference slide modeling

Null models, which follow the scenario of making a spurious predictions, were configured for each case based on the number of blocks available for a given patient (1 H&E stained slide per block,  $n = 21$  patients, minimum 3 blocks, average = 4.57 blocks, median = 4 blocks, max = 10 blocks), and a Beta prior was set according to the empirical observations of microheterogeneity in the CM-025 cohort (prior parameters:  $a = 148$ ,  $b = 279$ ). The null model is set up in which a given collection of slides from one patient’s tumor follows draws from a (Beta-)Binomial distribution whose frequency parameters are based on empirical counts from the CM-025 cohort. Thus, this null model expects random “miscalculation” or spurious detection. Alternative models for each patient were configured by first setting a uniform prior ( $a = 1$ ,  $b = 1$ ), and then updating based on the microheterogeneity status of a reference slide

(e.g.,  $\{a = 2, b = 1\}$  when observing microheterogeneous reference). We selected the reference slide based on grade score in 4 different ways: highest slide-average, highest by segment, highest by segment (with 10% area minimum), and highest by segment (with 25% area minimum). We also considered near-ties, considering a tie if two candidate grade scores were within 0.01 of each other, taking the average log likelihood of the competing alternative models for a given patient when comparing to the null in downstream testing. Comparison testing was done with a log likelihood ratio test. Bayes factors were calculated analytically under beta-binomial distributions.

### Tumor infiltrating lymphocyte inference

A HoVerNet model trained on the PanNuke dataset was used for nuclei segmentation.<sup>17,40</sup> We leveraged the implementation and pretrained model from the PathML toolkit<sup>41</sup> for computational pathology (<https://github.com/Dana-Farber-AIOS/pathml>). While this pretrained model produces accurate nucleus segmentation, its subtype classification notably fails on clear cell renal carcinoma, likely due to a near-absence of this histology within PanNuke. Consequently, we finetuned the classification head of the model to predict tumor-context vs. stromal-context nuclei using a pseudo-labeling scheme; nuclei in a tile were randomly assigned “tumor nuclei” labels proportionally to the predicted probability produced by the tumor classifier (ex., 90% of nuclei randomly assigned “tumor nuclei” if tumor score = 0.9). Following inference, we further stratified nuclei predicted to be “tumor-context nuclei” to distinguish tumor cells from infiltrating lymphocytes (TILs) using heuristic cutoffs chosen via manual pathologist review; lymphocytes were selected by a combined criteria of increased circularity, smaller area, and darker pixels.

### Tumor infiltration classification

Pathologist review was performed on tile examples to create a ground-truth thresholding set for TIL prediction, wherein the extent of TIL infiltration in a given ROI was used to tune model cutoffs to establish an optimal choice for predicting TILs. Following nuclei inference, we aggregated nuclei calls at the tile-level, and utilized the described pathologist annotations to establish effective cutoffs to determine whether a tile was infiltrated or not. We classified a tile as “infiltrated” if it contained 14 or more TILs, a cutoff selected by maximizing concordance with pathologist annotations for the presence of lymphocytes in a given tile (Figure S18; bootstrapped AUROC comparing “infiltrated” vs. “non-infiltrated” tile-level labels). We then considered region-level descriptions of infiltration, describing the proportion of tiles above the “infiltrated” cutoff as the “area infiltration fraction”. Next, we binarized samples into “low” versus “high” infiltration by splitting at the median area infiltration fraction value (cutoff = 15.16%). This was restricted to high-grade regions (grade score  $\geq 0.8$ ) to avoid excessive false positive infiltration calls, as lower-grade ccRCC nuclei can be visually ambiguous from TILs, even to expert pathologists ( $n = 256$  patient slides post filtration) (See Figures S15–S17). We extended evaluation of TIL classification to the key task of tissue-level infiltration description. We compared our H&E based measurement of lymphocyte infiltrated tumor area to CD8<sup>+</sup> immunofluorescence quantification performed on the same tumor independently (recognizing limitations given different sections of the tumor performed for these assays).<sup>11–13</sup> We evaluated whether there was a correlation between H&E derived TIL infiltration fraction measure and orthogonal CD8<sup>+</sup> immunofluorescence “tumor center” density using Pearson rho (Figure S19B).

### Truncating mutation categorization

When considering somatic alterations, we consider mutations only if they are truncating (likely loss of function). Within MAF annotation data, this comprised the following variant categories: {'Nonsense\_Mutation', 'Frame\_Shift\_Ins', 'Frame\_Shift\_Del', 'Splice\_Site'}.

### Molecular ITH metrics

We calculated the intratumoral heterogeneity index score (ITH) as described in Turajlic 2018, where ITH index = # subclonal drivers/# clonal drivers, calling indel mutational events as “drivers”, and setting a CCF cutoff of 0.5 to determine clonality. We calculated weighted genome integrity (wGI) index as originally described by Endesfelder et al. 2014.<sup>23</sup>

### Survival analysis

Survival analysis in the TCGA-KIRC And CM-025 cohorts was performed using the python package *Lifelines*.<sup>35</sup> Kaplan-Meier regression and plotting was performed using the *KaplanMeierFitter* function with default parameters, and multivariate Cox Proportional Hazards regression was performed using the *CoxPHFitter* function with moderate regularization (L1 ratio = 0.1 [multivariate models], L1 ratio = 0 [univariate models], penalizer scale = 0.1 [multivariate models], penalizer = 0.0 [univariate models]). We also further excluded slide images with fewer than 200 tiles predicted to contain tumor tissue. We considered only slides obtained from primary biopsy sites, excluding metastatic biopsies to remain consistent between cohorts. When annotations were available, we excluded Grade 1 (G1) cases due to their rarity. We only considered cases where watershed segmentation successfully produced at least one segment containing 50 or more tiles with an average tumor score  $\geq 0.7$ . When describing TIL infiltration content in CM-025 in Kaplan-Meier curves, we used binary (lower/higher) groups as described above, and continuous area infiltration fraction for Cox modeling.

### Image registration

We adapted PathFlow MixMatch, displacing an input H&E image against a fixed mIF image at 1.25X, and using GPU acceleration (Tesla V100) when learning each case’s alignment/displacement tensor. Learned displacement tensors were then used to shift

H&E-based grade segmentation maps into the same coordinate space as mIF data. These aligned maps were then used in K-nearest neighbors regression to assign cell predictions in the mIF data to a grade segmentation label. Since the image pairs are not from the same exact tissue section, alignments were assessed visually via overlay to assess quality, resulting in 13 total passing cases. Within successful alignments, putative tumor regions were manually reviewed, resulting in omission of two false-positive regions (adjacent metastatic tissue misclassified as “tumor”).

### Multiplexed immunofluorescence image preprocessing

To first predict cellular locations, we used a pretrained Mesmer model,<sup>34</sup> which produced a candidate mask of cell segmentations for each mIF WSI. We used DAPI as the “nuclear channel”, and PAX8 with CD8 as the “membrane channels”. Full resolution (20X) mIF images were broken into 10,000 pixel bands as batch inputs to Mesmer using GPU acceleration (Tesla V100). A subset of images (3) that failed this batch procedure were re-run with 2500 pixel square tiles as batched inputs. To quantify area-normalized cellular expression, we used the Ark analysis toolkit’s ‘create\_marker\_count\_matrix’ function,<sup>34,42</sup> which produces a description of each predicted cell segmentation that contains both morphological and arcsinh-transformed, area-normalized expression values.

### Cell phenotype calling

Following expression quantification, we inspected the histograms of each case’s channel values, as well as the ratio of CD8<sup>+</sup>: Autofluorescence, and determined manual cutoffs to gate each primary cell population (CD8<sup>+</sup> vs. Tumor+ vs. ungated). These cutoffs were then used to make a coarse-grained estimate of each cell subpopulation. We then performed an orthogonal clustering analysis, using cell expression and morphology features produced by the Ark toolkit (*{centroid\_dif, num\_concavities, convex\_hull\_resid, major\_axis\_equiv\_diam\_ratio, perim\_square\_over\_area, arcsinh(Cell Area)}*) and the Louvain method for community detection in scanpy (number of principal components = 5, nearest neighbors = 15, cluster resolution = 10).<sup>43,44</sup> Cells with low DAPI (<7 arcsinh units) were also excluded. Clusters with outlier morphology (>50% of its cells having 3+ features outside of 5th/95th percentile values), or high autofluorescence-to-CD8<sup>+</sup> signal (>35% below case-specific CD8<sup>+</sup>: Autofluorescence ratio cutoffs) were excluded. Remaining clusters were assigned to “CD8+” or “Tumor+” identities if at least 60% of a cluster’s cells were assigned that label when using purely manual cutoffs. Remaining cells were labeled “ungated” and excluded from downstream analysis. To determine cell subpopulations, we subsetted each primary cell population (CD8<sup>+</sup>, Tumor+), and fit a linear regression model of autofluorescence vs. submarker expression, using the resulting residuals as a noise-corrected expression value. Resulting submarker distributions (PDL1 for Tumor+ cells and PD1 for CD8<sup>+</sup> cells) were inspected, and binary cutoffs were chosen for each case individually. We lastly performed a filtering of false-positive tumor cell predictions which exhibited high PAX8 and high DAPI (likely to be B cell lineage), and again used manual histogram inspection to remove the DAPI-high subpopulation.

### Cell graph construction

To construct a graph of cell-interactions, we first removed ungated cells, and then performed Delaunay triangulation with a maximum radius of 100px to form a parsimonious nearest neighbor graph. We then defined “self” interactions as edges between cells of the same type (e.g., CD8<sup>+</sup>), and “tumor-immune” interactions as those occurring between CD8<sup>+</sup> and Tumor+ cells. Cells disconnected from the graph were deemed “isolated”, and grouped with “self” interactions for downstream analysis.

### Immune hotspot analysis

To select regions of interest with high tumor-immune activity, we split full resolution mIF WSI data into tiles 2000px (approx 1mm) wide, and further selected for regions with at least 50% area overlap with H&E-inferred tumor region predictions. We then filtered for regions with at least 50 CD8<sup>+</sup> and 50 Tumor+ cells that were engaged in tumor-immune interactions. From each patient slide, we sampled up to 10 hotspots, selecting those with the most CD8<sup>+</sup> density involved in tumor-immune interactions (min = 2 samples, mean = 9.0 samples; 42 total microheterogeneous samples [n = 5 slides], 57 total microhomogeneous samples [n = 6 slides]) (See [Figures S31](#) and [S32](#)). Two patients lacked any hotspots and were excluded from this analysis.

## QUANTIFICATION AND STATISTICAL ANALYSIS

All statistical analysis was performed using python 3. For comparison of group counts, Fisher’s Exact test was used via the scipy function *fisher\_exact*.<sup>45</sup> Other continuous, score-based comparisons were performed using a two-sided Wilcoxon rank-sum (Mann-Whitney U) test using the *statannotations* package to directly annotate Seaborn plots with p value results.<sup>46,47</sup> For survival analyses, cohort subgroup survival distributions were compared using the log rank test using the *multivariate\_logrank\_test* and *pairwise\_logrank\_test* functions in *Lifelines*. For Cox models, the concordance index (C-Index) and (Log) Likelihood Ratio Test (LLRT) were used to evaluate goodness of fit. When comparing continuous to categorical grade, the relative likelihood was estimated using the partial AIC produced for each Cox model, and interpreted as the probability one model minimizes the AIC of the other. Barplot error bars indicated standard error. Boxplot elements are as follows: center line, median; box limits, upper and lower quartiles; whiskers, 1.5 interquartile range (IQR); points, outliers past 1.5 IQR. Violinplot dotted interior lines indicate median, and upper and lower quartiles.