

# Whole-genome analysis of papillary kidney cancer finds significant noncoding alterations

**Authors:** Shantao Li<sup>1</sup>, Brian M. Shuch<sup>2,\*</sup>, Mark B. Gerstein<sup>1,3,4,\*</sup>

<sup>1</sup>Program in Computational Biology and Bioinformatics, Yale University, New Haven, CT 06520, USA.

<sup>2</sup>Department of Urology, Yale School of Medicine, New Haven, CT, 06520, USA.

<sup>3</sup>Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, CT 06520, USA.

<sup>4</sup>Department of Computer Science, Yale University, New Haven, CT 06520, USA.

\*To whom correspondence should be addressed: [brian.shuch@yale.edu](mailto:brian.shuch@yale.edu), [pi@gersteinlab.org](mailto:pi@gersteinlab.org)

**Short title:** Whole-genome analysis of papillary kidney cancer

**Abstract:** To date, studies on papillary renal-cell carcinoma (pRCC) have largely focused on coding alterations in traditional drivers, particularly the tyrosine-kinases MET. However, for a significant fraction of tumors, researchers have been unable to determine clear molecular etiologies. To address this, we perform the first whole-genome analysis of pRCC. Elaborating on previous results on MET, we find a germline SNP in MET predicting prognosis (rs11762213). Interestingly, we detect no enrichment for small structural variants disrupting MET. Furthermore, we discover methylation dysregulation leads to cryptic promoter activation in MET, inducing alternate transcript expressing. Next, we scrutinize noncoding mutations, discovering potentially impactful ones in regulatory regions associated with MET and in a long noncoding RNA (NEATI). Moreover, NEATI mutations in pRCC are associated with increased expression and unfavorable outcome. Finally, we investigate genome-wide mutational patterns, finding they are governed mostly by methylation-associated C-to-T transitions. Also, we observe significantly more mutations in open chromatin and early replicating regions in tumors with chromatin-

Shantao 2/24/2017 2:44 PM

Deleted: non-coding

Shantao 2/24/2017 1:23 PM

Comment [1]: Somatic SNVs are covered in the Pan-RCC paper (cell reports); Therefore, I removed the somatic SNVs part in the abstract  
But it is still in the main text

Shantao 2/24/2017 1:27 PM

Deleted: in the coding regions of this gene

Shantao 2/23/2017 7:08 PM

Deleted: more somatic alternations and

Shantao 2/24/2017 1:24 PM

Deleted: find

Shantao 2/23/2017 7:07 PM

Deleted: associated with

Shantao 2/24/2017 2:44 PM

Deleted: non-coding

Shantao 2/24/2017 2:44 PM

Deleted: non-coding

Shantao 2/23/2017 7:09 PM

Deleted: is implicated in other cancers and its

33 modifier alterations. Last, we construct evolutionary trees and reveal various structures of tumor  
34 development. Our mutational processes study helps understand the origin of pRCC  
35 heterogeneity.

36 **Currently word count: 182, want <=150**

37

### 38 **Author Summary**

39 Renal cell carcinoma accounts for more than 90% of kidney cancers. Papillary renal cell  
40 carcinoma (pRCC) is the second most common subtype of renal cell carcinoma. Previous studies,  
41 focusing mostly on the protein-coding regions, have identified several key genomic alterations  
42 that are key to cancer initiation and development. However, researchers cannot find any key  
43 mutation in a significant portion of pRCC. Therefore, we carry out the first whole-genome study  
44 of pRCC to discover triggering DNA changes explaining these cases. By looking at the entire  
45 genome, we find additional potentially impactful alterations in and out of the protein-coding  
46 regions. These newly identified critical mutations from scrutinizing the entire genome help  
47 complete our understanding of pRCC genomes. Two alterations we found are associated with  
48 prognosis, which could aid clinical decisions. We are also able to recognize mutation patterns,  
49 signatures and tumor evolution structures, which reflect the mutagenesis processes and help  
50 understand how cancer develops. Our study provides valuable additional information to facilitate  
51 better tumor subtyping, risk stratification and potentially clinical management.

52

### 53 **Introduction**

Shantao 2/23/2017 7:14 PM

**Deleted:** give hints on

55 Renal cell carcinoma (RCC) makes up over 90% of kidney cancers and currently is the  
56 most lethal genitourinary malignancy (1). Papillary RCC (pRCC) accounts for 10%-15% of the  
57 total RCC cases (2). Unfortunately pRCC has been understudied and there are no current forms  
58 of effective systemic therapy for this disease. pRCC are further subtyped into two major groups:  
59 type 1 and type 2 based on histopathological features. For many years, the only prominent  
60 oncogene in pRCC (specifically, type 1) that physicians were able to identify was *MET*, a  
61 tyrosine kinase receptor for hepatic growth factor. An amino acid substitution that leads to  
62 constitutive activation and/or overexpression are two mechanisms of dysfunction of *MET* in  
63 tumorigenesis. Recently, the Cancer Genome Atlas (TCGA) published its first result on pRCC  
64 (3), which greatly improves our understanding of the genomic basis of this disease. Several more  
65 genes and specific sub-clusters were identified to be significantly mutated in pRCC.

66 Nevertheless, a significant portion of pRCC cases still remains without any known driver.

67 Therefore we think it is a good time to explore the rest 98% **noncoding** regions of the genome  
68 using whole genome sequencing (WGS). This is sensible because **noncoding** regions, previously  
69 overlooked in cancer, have been showed to be actively involved in tumorigenesis (4-6).

70 Mutations in **noncoding** regions may cause disruptive changes in both cis- and trans-regulatory  
71 elements, affecting gene expression. Understanding **noncoding** mutations helps fill the missing  
72 “dark matter” in cancer research.

73 Multiple endogenous and environmental mutation processes shape the somatic mutational  
74 landscape observed in cancers (7). Analyses of the genomic alterations associated with these  
75 processes give information on cancer development, shed light on mutational disparity between  
76 cancer subtypes and even indicate potential new treatment strategies (8). Additionally, genomic  
77 features such as replication time and chromatin environment govern mutation rate along the

Shantao 2/24/2017 2:44 PM

Deleted: non-coding

Shantao 2/24/2017 2:44 PM

Deleted: non-coding

Shantao 2/24/2017 2:44 PM

Deleted: non-coding

Shantao 2/24/2017 2:44 PM

Deleted: non-coding

82 genome, contributing to spatial mutational heterogeneity. While identifying mutation signatures  
83 is possible using data from whole exome sequencing (WXS), whole genome sequencing (WGS)  
84 gives richer information on mutation landscape and minimizes the potential confounding effects  
85 of exome capture process and driver selection.

86 In this study, we comprehensively analyzed 35 pRCC cases that were whole genome  
87 sequenced along with an extensive set of WXS data on multiple levels. We went from  
88 microscopic examination of driver genes to analyses of whole genome sequencing variants, and  
89 finally, to investigation of high-order mutational features. **We focused on two aims: exploring**

90 **potential noncoding drivers and better understanding the cancer heterogeneity.** First, we focused  
91 on *MET*, an oncogene which plays a central role in pRCC, especially in type 1. We found  
92 rs11762213, a germline exonic single nucleotide polymorphism inside *MET*, predicts cancer-  
93 specific survival (CSS) in type 2 pRCC. We also discovered several potentially impactful

94 **noncoding** mutation hotspots in *MET* promoter and its first two exons. The previous TCGA  
95 study identifies a *MET* **alternate transcript** as a driver, but without illustrating the etiology (3). We  
96 found that a cryptic promoter from a long interspersed nuclear element-1 (L1) triggers the

97 **alternate** isoform expression. Surprisingly, we did not find a significant amount of structural  
98 variations affecting *MET* besides polysomy 7. Then we went onto cases not as easily explained  
99 as those with *MET* alterations. We analyzed about 160,000 **noncoding** mutations throughout the  
100 entire genomes and found several potentially high-impact mutations in **noncoding** regions.

101 Further zooming out, we discovered pRCC exhibits mutational heterogeneity in both nucleotide  
102 context and genome location, indicating underlying vibrant mutational processes interplay. We  
103 found methylation is the leading factor influencing mutation landscape. Methylation status drives  
104 the intra-sample mutation variation by promoting more C-to-T mutations in the CpG context.

Shantao 2/24/2017 7:29 PM  
Formatted: Highlight

Shantao 2/24/2017 2:44 PM  
Deleted: non-coding

Shantao 2/24/2017 2:37 PM  
Deleted: alternative

Shantao 2/24/2017 2:37 PM  
Deleted: splicing

Shantao 2/24/2017 2:37 PM  
Deleted: event

Shantao 2/24/2017 2:37 PM  
Deleted: event

Shantao 2/24/2017 2:37 PM  
Deleted: alternative

Shantao 2/24/2017 2:44 PM  
Deleted: non-coding

Shantao 2/24/2017 2:44 PM  
Deleted: non-coding

113 APOBEC activity, although infrequently observed, leaves an unequivocal mutation signature in a  
114 pRCC genome but not in ccRCC. Also, we discovered samples with chromatin remodeler  
115 alternations accumulate more mutations in open chromatin and early-replicating regions. Last,  
116 we inferred evolution tree for each individual samples. Tree structures vary, reflect tumor  
117 heterogeneity and correlate with tumor subtypes.

Shantao 2/24/2017 2:14 PM

Deleted: and found

Shantao 2/24/2017 2:15 PM

Deleted: t

Shantao 2/24/2017 2:16 PM

Deleted: .

118

## 119 Results

### 120 1. An exonic SNP in *MET*, rs11762213, predicts prognosis in type 2 pRCC.

121 We begin with coding variants in the long known driver *MET*. The TCGA study of 161  
122 pRCC patients found 15 samples carrying somatic, nonsynonymous single nucleotide variant  
123 (SNV) in *MET*. By analyzing 117 extra WXS samples (see Methods), we found six more  
124 nonsynonymous somatic mutations in six samples (Table S1). V1110I and M1268T were two  
125 recurrent mutations in this extra set. Both of them were observed in the TCGA study as well.  
126 Additionally, we found two samples carrying H112Y and Y1248C respectively. H112Y has  
127 been observed in two patients the original TCGA study cohort and H1118R is a long-known  
128 germline mutation associated with hereditary papillary renal carcinoma (HPRC, 13). Y1248C  
129 has been observed in type 1 pRCC before and the TCGA cohort has a case carrying Y1248H. All  
130 mutations occur in the hypermutated tyrosine kinase catalytic domain of *MET*. Two out of these  
131 six samples were identified as type 1 pRCC while the subtypes of the rest four were unknown.

Shantao 2/24/2017 2:20 PM

Deleted: (rs121913246)

132 Although many *MET* somatic mutations are believed to play a central role in pRCC,  
133 some germline *MET* mutations have also been associated with the disease. In particular, a  
134 germline SNP, rs11762213, has been discovered to predict recurrence and survival in a mixed

139 | RCC cohort (14, [Figure 1A](#)). ccRCC predominated the initial discovery RCC cohort. This  
 140 | conclusion was later validated in a ccRCC cohort but never in pRCC (9). We wondered whether  
 141 | this SNP has a prognostic effect in pRCC. Using an extensive WXS set of 277 patients (see  
 142 | Methods; Figure S1 and Table S1;), we found 14 patients carry one risk allele of rs11762213  
 143 | (G/A, Table 1, minor allele frequency (MAF) = 2.53%). No homozygous A/A was observed.  
 144 | Cancer specific deceases are concentrated in type 2 pRCC. Among 96 type 2 pRCC cases, seven  
 145 | patients carry the minor A allele (MAF = 3.65%, Table 1). Survival is significantly worse in type  
 146 | 2 patients carrying the risk allele of rs11762213 ( $p = 0.034$ , Figure 1B). But we did not find  
 147 | significant association of this germline SNP with survival in type 1 patients. We did not find  
 148 | statistically significant association of rs11762213 with *MET* RNA expression in either tumor  
 149 | samples or normal controls ( $p > 0.1$ , two-sided rank-sum test). *Met* pY1235 levels in tumor  
 150 | samples, as measured by Reverse phase protein array (RPPA), were not significantly different in  
 151 | patients carrying the minor G allele compared to patients with A/A genotype ( $p > 0.1$ , two-sided  
 152 | rank-sum test).

Shantao 2/24/2017 2:21 PM  
 Formatted: Font:Not Italic

Characteristic	G/A (n = 7)	A/A (n = 89)
<b>Sex, No. (%)</b>		
Male (%)	4 (57)	25 (28)
Female (%)	3 (43)	64 (72)
<b>Age, median (IQR), year</b>	54 (47-61)	65 (57-73)
<b>Race, No. (%)</b>		
White	6 (86)	65 (73)
Black	1 (14)	16 (18)
Asian	0	4 (4)
NA	0	4 (4)
<b>T stage, No. (%)</b>		
T1	4 (57)	47 (53)
T2	1 (14)	10 (11)
T3	2 (29)	31 (35)
T4	0	1 (1)
<b>N stage, No. (%)</b>		

N0	3 (43)	20 (22)
N1	0	15 (17)
N2	1 (14)	2 (2)
NX	3 (43)	52 (58)
<b>M stage, No. (%)</b>		
M0	3 (43)	54 (61)
M1	1 (14)	4 (4)
MX/NA	3 (43)	31 (35)
<b>AJCC stage, No. (%)</b>		
I	4 (57)	43 (48)
II	0	7 (8)
III	1 (14)	29 (33)
IV	2 (29)	6 (7)
NA	0	4 (4)
<b>Median follow-up for surviving patients, days (IQR)</b>	243 (132-354)	579 (219-1247)

153

154 **Table 1. Patient clinical profiles of the type 2 pRCC cohort in rs11762213 survival analysis.** AJCC: American  
 155 Joint Committee on Cancer; IQR: interquartile range; NA: not available. Percentages may not add up to 100%  
 156 because of rounding.

157

158 **2. Epigenetic alterations and mutation hotspots in noncoding regions**

159 The TCGA study has identified a *MET* alternate transcript as a driver event (3). However,  
 160 the etiology of this new isoform is unknown. We identified this alternate transcript results from  
 161 the usage of a cryptic promoter from an L1 element (Figure 1A), likely due to a local loss of  
 162 methylation (REF). This event was reported in several other cancer types (REF). To test its  
 163 relationship with methylation, we found a closet probe (cg06985664, ~3kb downstream) on the  
 164 methylation array shows marginally statistically significant ( $p=0.055$ , one-sided rank-sum test).  
 165 Additionally, this event is associated with methylation group 1 (odds ration (OR)= 4.54,  
 166  $p<0.041$ ), indicating genome-wide methylation dysfunction. This association is stronger in type  
 167 2 pRCC and it shows a significant association with the C2b cluster (OR= 17.5,  $p<0.007$ ).

Shantao 2/24/2017 2:26 PM

Deleted: tive

Shantao 2/24/2017 2:26 PM

Deleted: lation isoform

Shantao 2/24/2017 2:26 PM

Deleted: isoform

Shantao 2/24/2017 5:45 PM

Deleted: one-side

Shantao 2/24/2017 2:38 PM

Deleted: as expected,

173 Despite the fact *MET* is the most common driver alteration, about 20% presumably *MET*-  
174 driven yet *MET* wild-type pRCC samples were still left unexplained (3). Therefore, we scanned  
175 the *MET* **noncoding** regions. We observed one mutation in *MET* promoter region in a type 1  
176 pRCC sample (Figure 1A and Table S2). This sample shows no evidence of a nonsynonymous  
177 mutation in *MET* gene but it has copy number gain of *MET*. Additionally, we observed 6/35  
178 (17.1%) samples carry mutations in the intronic regions between exon 1-3 of *MET* (Figure 1A  
179 and Table S2). Previously it is been established that ~~an alternate transcript involving these~~ exons  
180 is a driver event (3). Therefore we speculated that these **noncoding** variants might correlate with  
181 the alternative splicing. However, likely being hindered by a small size, we were not able to find  
182 statistically significant association between the alternative splicing event and these intronic  
183 mutations.

184 We further expanded our scope and ran FunSeq (4-5) to identify potentially high-impact,  
185 **noncoding** variants in pRCC. First, we identified a high-impact mutation hotspot on chromosome  
186 1. 6/35 (17.1%) samples have mutations within this 6.5kb region (Figure 2A and Table S2). This  
187 hotspot locates at the upstream of *ERRFI1* (ERBB Receptor Feedback Inhibitor 1) and overlaps  
188 with the predicted promoter region. *ERRFI1* is a negative regulator of EGFR family members,  
189 including EGFR, HER2 and HER3, all have been implicated in cancer. Due to a limited sample  
190 size here, our test power was inevitably low. We didn't observe statistically significant changes  
191 among mutated samples in mRNA expression level, protein level and phosphorylation level of  
192 EGFR, HER2 and HER3.

193 Another potentially impactful mutation hotspot is in *NEATI*. We saw mutations inside  
194 this nuclear long **noncoding** RNA in 6/35(17.1%) samples (Figure 2B and Table S2). Several  
195 studies indicated *NEATI* is associated in many other cancers (15-16). It promotes cell

Shantao 2/24/2017 2:44 PM

Deleted: non-coding

Shantao 2/24/2017 5:17 PM

Deleted: alternative

Shantao 2/24/2017 5:17 PM

Deleted: splicing

Shantao 2/24/2017 5:17 PM

Deleted: of these

Shantao 2/24/2017 5:17 PM

Deleted:

Shantao 2/24/2017 2:44 PM

Deleted: non-coding

Shantao 2/24/2017 2:44 PM

Deleted: non-coding

Shantao 2/24/2017 2:39 PM

Deleted: very

Shantao 2/24/2017 2:44 PM

Deleted: non-coding



205 proliferation in hypoxia (17) and alters the epigenetic landscape, increasing transcription of  
206 target genes (18).

207 All the mutations we found fell into a putative promoter region of *NEATI*. We noticed  
208 *NEATI* mutations were associated with higher *NEATI* expression (Figure 2C,  $p < 0.032$ , two-  
209 sided rank sum test). We also found *NEATI* mutations were associated with worse prognosis  
210 (Figure 2D,  $p < 0.041$ , log-rank test). To further investigate the role of *NEATI* in RCCs, we  
211 found *NEATI* overexpression is significantly associated with shorted overall survival in TCGA  
212 ccRCC cohort ( $p=0.0132$ , Fig SXX). Moreover, *MALAT1*, another noticeable lncRNA in cancer,  
213 is tightly co-expressed with *NEATI* in both pRCC and ccRCC (Spearman's correlation: 0.79 and  
214 0.87 respectively). Catalogue of Somatic Mutations in Cancer (COSMIC) (REF) annotates  
215 *MALAT1* as cancer consensus gene, associating it with pediatric RCCs and lung cancer.  
216 Overexpression of *MALAT1* is reported to be associated with cancer progression (REF).

### 218 3. Structural variations in pRCC

219 We used DELLY (10) to perform structural variants (SVs) discovery from WGS reads  
220 information (see Methods and Table S3). The SV discovery approach has higher sensitivity and  
221 resolution than array-based methods, which were employed in the TCGA analysis. In the end we  
222 found 424 somatic SV events, includes 170 deletions, 53 duplications, 105 inversions and 96  
223 translocations (Figure SXX). Samples clearly split into two types based on SV counts (0-88):  
224 genome unstable (6 samples, >40 events/per samples) and genome stable (29 samples, <10  
225 events/per sample). The unstable type is significantly associated with type 2 versus type 1  
226 ( $p < 0.015$ , two-tailed Fisher exact test) and enriched in C2b cluster ( $p < 0.002$ , two-tailed Fisher  
227 exact test).

Shantao 2/24/2017 2:46 PM

Formatted: Font:Italic

Shantao 2/24/2017 2:45 PM

Formatted: Indent: First line: 0"

Shantao 2/24/2017 3:18 PM

Formatted: Font:Bold

228 First, by overlapping SVs with curated cancer genes from COSMIC (REF), we found two  
229 cases with deletion in *SDHB*. The median *SDHB* expression is ~50% compared to cases without  
230 alteration (Figure SXX). We validated the deletions affecting *SDHB* with another SV caller,  
231 Lumpy (SV). We confirmed three cases carrying deletions affecting *CDKN2A* called by TCGA  
232 array-based methods but not the other two cases. Notably, three confirmed cases have  
233 significantly lower *CDKN2A* expression but not in the unconfirmed two cases (Figure SXX).  
234 This suggests SV calling from WGS is accurate and predicts expression better. One sample,  
235 TCGA-B9-4116, which has extensive amplification of *MET*, showed multiple SVs of various  
236 classes hitting *MET* regions. However, surprisingly, we did not find SVs affecting *MET* except  
237 this one example. We postulate trisomy/polysomy 7 is the main mechanism of *MET* structural  
238 alteration rather than duplication in a smaller scale. Besides duplication, we did not expect to  
239 find deletion, inversion or translocation disrupting oncogene *MET*. These SVs are likely to cause  
240 loss-of-function rather than gain-of-function mutations. This is consistent with the putative role  
241 of *MET* as an oncogene, rather than a tumor suppressor.

242 Last, we observed several high impact sporadic events, including duplication in *EGFR*  
243 and *HIF1A* duplication and deletions in *DNMT3A* and *STAG2* (see SXX).

#### 245 4. Mutation spectra and mutation processes of pRCC

246 To further get a high-order overview of the mutation landscape, we summarized the  
247 mutation spectra of 35 whole genome sequenced pRCC samples (Figure 3A). C-to-T in CpGs  
248 showed the highest mutation rates, which were roughly three to six-fold higher than mutation  
249 rates in other nucleotide contexts.

Shantao 2/24/2017 6:49 PM

Formatted: Font:Italic

Shantao 2/24/2017 6:49 PM

Formatted: Font:Italic

Shantao 2/24/2017 6:49 PM

Formatted: Font:Italic

Shantao 2/24/2017 2:45 PM

Deleted: 3

251 We used principle components analysis (PCA) to reveal factors that explain the most  
252 inter-sample variation. The loadings on the first principle component (which explained 12.5% of  
253 the variation) demonstrated C-to-T in CpGs contributed the most to inter-sample variation  
254 (Figure 3B). C-to-T in CpGs is highly associated with methylation. It reflects the spontaneous  
255 deamination of cytosines in CpGs, which is much more frequent in 5-methyl-cytosines (REF).  
256 So we further explored the association between C-to-T in CpGs and tumor methylation status.  
257 First, we validated the TCGA identified methylation cluster 1 showed higher methylation level  
258 than cluster 2 in all annotated regions (Figure S2, see Methods), prominently in CpG Islands  
259 (Odds ratio of sites being differentially hypermethylated: 1.29, 95%CI: 1.20-1.39,  $p < 0.0001$ ).  
260 We confirmed this association by showing samples from methylation cluster 1 had higher PC1  
261 scores as well as higher C-to-T mutation counts and mutation percentages in CpGs (Figure 3C).  
262 This trend was further validated using a larger WXS dataset as well. Especially, the most  
263 hypermethylated group, CpG island methylation phenotype (CIMP), showed the greatest C-to-T  
264 in CpGs (Figure S2). Therefore, methylation status is the most prominent factor shaping the  
265 mutation spectra across patients.

266 We further explored the functional impact of the excessive mutations driven by  
267 methylation. C-to-T mutations in CpGs were more likely to be in the coding region (OR=1.54,  
268 95%CI: 1.27-1.85,  $p < 0.0001$ ) and nonsynonymous (OR=1.47, 95%CI: 1.17-1.84,  $p < 0.001$ ). This  
269 indicates hypermethylation tends to cause high impact mutations. However, C-to-T mutations in  
270 CpGs did not show functional bias between two methylation groups in noncoding regions (based  
271 on FunSeq score distribution).

272 Recently, 30 somatic mutation signatures were identified. Many have putative etiology,  
273 revealing the underlying mutation processes and help understand tumor development (7). We

Shantao 2/24/2017 4:05 PM

Deleted: several

275 used a LASSO-based approach (see Methods) to decompose mutations into a linear combination  
276 of these canonical mutation signatures in both WGS and WXS samples (Figure S3). The leading  
277 signature was signature 5, which is consistent with previous studies (7). Interestingly, we found  
278 one type 2 pRCC case out of 155 somatic WXS sequenced samples exhibited APOBEC-  
279 associated mutation signature 2 and 13. APOBEC mutation pattern enrichment analysis (see  
280 Method) further confirmed the presence of APOBEC activity (Figure 3D). This sample was  
281 statistically enriched of APOBEC mutations (adjusted p-value < 0.0003).

282 Prominent APOBEC activities were also incidentally detected in three upper track  
283 urothelial cancer (UC) samples sequenced and processed in the same pipeline with pRCC  
284 samples. UC often carries APOBEC mutation signatures and our result is consistent with TCGA  
285 bladder urothelial cancer study (19).

286 The APOBEC-signature carrying pRCC case was centrally reviewed by six pathologists  
287 in the original study and confirmed to be type 2 pRCC (3). Thus this tumor is likely a special  
288 case of type 2 with genomic alterations share some similarities with UC. It has non-silent  
289 mutations in *ARID1A* and *MLL2* and a synonymous mutation in *RXRA*, all are identified as  
290 significantly mutated genes in UC but not in pRCC. Potential [type 2](#) pRCC driver events, for  
291 example low expression of *CDKN2A* and nonsynonymous alternations in significantly mutated  
292 genes of pRCC, are absent in this sample.

293 Noticeably, all four samples with APOBEC activities showed significantly higher  
294 *APOBEC3A* and *APOBEC3B* mRNA expression level ( $p < 0.0022$  and  $p < 0.0039$  respectively,  
295 [one-sided](#) rank sum test, Figure S4). This is in concordance with previous studies of APOBEC  
296 mutagenesis in various types of cancer (12).

Shantao 2/24/2017 5:45 PM

Deleted: one-side

298 Consistent with previous studies (12), we failed to detect statistically significant  
299 APOBEC activities in an extensive WXS dataset consisting of 418 clear cell RCC (ccRCC)  
300 samples, even after resampling to avoid p-value adjustment eroding the power. Very low levels  
301 of APOBEC signatures (<15%) was found in less than 1%(4/418) samples. With a much larger  
302 sample size, this result was unlikely to be confounded by detecting power.

303

304

### 305 5 Defects in chromatin remodeling affects mutation landscape

306 Chromatin remodeling genes are frequently mutated in pRCC and many other cancers  
307 including ccRCC (20, [REF CR paper](#)). Defects in chromatin remodeling cause dysregulation of  
308 chromatin environment. Open chromatin regions show lower mutation rate, presumably due to  
309 more effective DNA repair (21). Thus chromatin remodeler alternations could possibly alter the  
310 mutation landscape, specifically increase mutation rate in previously open chromatin regions. To  
311 test this hypothesis, we tallied the number of mutations inside DNase I hypersensitive sites  
312 (DHS) [inferred from](#) eleven normal fetal kidney cortex samples (The NIH Roadmap  
313 Epigenomics Mapping Consortium, REF), which represent the normal, physiological condition.  
314 9/35 samples with disruptive mutations in ten chromatin remodeling, cancer associated genes  
315 show higher genome-wide mutation counts ( $p < 0.021$ , [one-sided](#) rank-sum test), partially driven  
316 by higher mutation counts in DHS region ( $p < 0.0023$ , [one-sided](#) rank-sum test). The median  
317 number of mutations in DHS region considerably increases by 60% (67.5 versus 108) in samples  
318 carrying chromatin remodeling defects. The effect is significant after normalizing against the  
319 total mutation counts ( $p < 0.019$ , [one-sided](#) rank-sum test, Figure 3E), [demonstrating a true shift](#)  
320 [in mutation landscape](#).

Shantao 2/24/2017 2:45 PM

Deleted: 4

Shantao 2/24/2017 4:09 PM

Deleted: in

Shantao 2/24/2017 5:46 PM

Deleted: one-side

Shantao 2/24/2017 5:46 PM

Deleted: one-side

Shantao 2/24/2017 5:46 PM

Deleted: one-side

Shantao 2/24/2017 4:09 PM

Deleted: .

327 Replication time is known to correlate greatly with mutation rate. Early replicating  
328 regions have lower mutation rate compared to late replicating ones. Researchers reason  
329 replication errors are more likely to be corrected by DNA repair system in early replicating  
330 regions. With defects in mutated chromatin remodeling, we observed this trend became less  
331 pronounced ( $p < 0.031$ , one-sided rank-sum test, Figure S5). This is likely because dysregulation  
332 of the chromatin environment hinders replication error repair by changing the accessibility of  
333 newly synthesized DNA chains.

Shantao 2/24/2017 5:46 PM  
Deleted: one-side

334

### 335 6. Evolutionary tree reveals the heterogeneity of tumor evolution profile,

336 With the richness of SNVs in WGS samples, we further tackle the mutational process  
337 heterogeneity of pRCC by constructing individual evolutionary trees for 35 tumors (Figure  
338 SXX). Three trees have largest population fraction  $< 0.5$  (likely due to low mutation number, high  
339 sequence error and/or high heterogeneity) and thus are excluded from downstream analysis. We  
340 could further classify the trees into four types based on topology (Figure 4A, 4B): (1) no branch,  
341 fewer subclones (10, 32.3%), (2) short branches (12, 37.5%), (3) no branch, more subclones (5,  
342 15.6%) and (4) long branches (5, 15.6%). Both (3) and (4) show significant clonal evolution,  
343 indicated by long mutations distances between populations.

Shantao 2/24/2017 2:45 PM  
Deleted: 5

Shantao 2/24/2017 10:47 AM  
Deleted: analysis

Shantao 2/24/2017 4:11 PM  
Deleted: , we inferred

Shantao 2/24/2017 4:11 PM  
Deleted: 35

Shantao 2/24/2017 4:11 PM  
Deleted: a

344 Short branch type is significantly enriched in Type I pRCC ( $p < 0.011$ , two-tailed fisher exact test,  
345 Figure 4B) while the more heterogeneous types: long branches and no branch, more subclones  
346 type are significantly depleted in Type I ( $p < 0.0034$ , two-tailed fisher exact test). This indicates  
347 type I tumors are more homogenous and show less complex evolution features compared to type  
348 II and unclassified samples.

Shantao 2/24/2017 10:39 AM  
Deleted: 2

356

357 **Discussion**

358 Our study is the first one that comprehensively looks into the noncoding regions of  
359 pRCC. Doing so allow us to tackle an open question in the field of cancer genomics, whether  
360 whole genome sequencing adds additional value over whole exome sequencing. We  
361 comprehensively analyzed both WGS and an extensive set of WXS of pRCC, scrutinizing local  
362 high-impact events as well as giving a macro overlook of the mutation landscape and evolution.  
363 Our work further completed the genomic alteration landscape of pRCC (Figure 4B). Beyond  
364 traditionally driver events, we suggested several novel noncoding alterations potentially drive  
365 tumorigenesis. We also provide valuable insights to tumor heterogeneity though investigating the  
366 mutation patterns, landscape and evolution profiles.

367 First, we elaborated on previous results of the long known driver *MET*. In an extended  
368 117 WXS dataset, we found six additional nonsynonymous somatic mutations in the  
369 hypermutated tyrosine kinase catalytic domain. These somatic mutations are highly recurrent,  
370 concentrated on a few critical amino acids. This is in line with *MET* being an oncogene and  
371 supports the central role of *MET* in pRCC. Then we found an exonic SNP in *MET*, rs11762213,  
372 to be a prognostic germline variance in type 2 pRCC. Previously, rs11762213 was found to  
373 predict outcome in a mixed RCC samples, predominated by ccRCC (14). Later, the result is  
374 confirmed in a large TCGA ccRCC cohort (9). However, it is never clear whether rs11762213  
375 only predicts the outcome in ccRCC or other histological types as well. In this study, we  
376 concluded that the minor alternative allele of rs11762213 also forecasts unfavorable outcome in  
377 type 2 pRCC patients. The mechanism of this exonic germline SNP remains unsettled. A  
378 previous study proposes it disrupts a putative enhancer of *MET*. However, researchers cannot

379 find significant difference in *MET* expression in either tumor or normal tissues. We noticed there  
380 is no other gene within 100 kb of this SNP. Given the significant role of *MET* in pRCC, we also  
381 think rs11762213 is affecting survival through *MET*, although the mechanism unknown.

382 Remarkably, similar to ccRCC, type 2 pRCC is not primarily driven by *MET*. Not as  
383 significantly mutated in ccRCC and type 2 pRCC, *rs11762213 correlating with survival shows*  
384 *MET* nonetheless seems to play a role in cancer development. Our finding on rs11762213 is  
385 potentially meaningful in clinical management of patients with the more aggressive type 2  
386 pRCC. rs11762213 genotyping could become a reliable, low-cost risk stratification tool for these  
387 patients. Also, rs11762213 might become a biomarker for predicting response to *Met* inhibitors.

388 Interestingly, rs11762213 is prevalent mostly in European and American populations but  
389 not in African populations and rare in Asian populations. MAF of rs11762213 among African  
390 American patients in our cohort is 2.73%, higher than MAFs in general African populations  
391 observed in 1000 Genome phase 3 dataset (0.2%, with 0% in Americans with African ancestry  
392 (ASW))) and the ExAC dataset (1.1%, excluding TCGA cohorts). This implies a possible effect  
393 of rs11762213 on pRCC incidence among African Americans that is worth further investigation.

394 Besides, in *MET* noncoding regions, we first find a cryptic promoter from a  
395 retrotransposon in the second intron initiates the alternate transcript, which is classified as a  
396 driver by the TCGA study (3). Methylation is a major source of silencing retrotransposon  
397 activities in human genome (REF). Indeed, we observed evidence for a local loss of methylation  
398 and global methylation dysregulation in samples expressing the alternate isoforms. Our finding  
399 indicates methylation change might directly drive pRCC growth through *MET*.

Shantao 2/24/2017 4:16 PM

Deleted: MET

Shantao 2/24/2017 2:44 PM

Deleted: non-coding

Shantao 2/24/2017 4:20 PM

Formatted: Font:Italic



402 We also discovered mutations associated with *MET* promoter and first two introns.  
403 Although the implication is unknown, our analysis suggests there is a mutation hotspot in *MET*  
404 that calls for further research.

405 Expanding our scope from coding to noncoding and use FunSeq to group SNVs by  
406 functional elements, we found several potentially significant noncoding mutation hotspots  
407 relevant to tumorigenesis throughout the entire genome. A mutation hotspot was found upstream  
408 of *ERRF1*, an important regulator of the EGFR pathway, which may serve as a potential tumor  
409 suppressor. EGFR inhibitors have been used in papillary kidney cancer with an 11% response  
410 rate observed (22). These mutations potentially disrupt regulatory elements of *ERRF1* and thus  
411 play a role in tumorigenesis. However, likely limited by a small sample size, we were not able to  
412 detect statistically significant functional changes in *ERRF1* and related pathways. Another

413 noncoding hotspot is in *NEAT1*, a long noncoding RNA that has been speculated to involve in  
414 cancer. All mutations locate in a putative regulatory region of the gene. Patients carrying  
415 mutations in *NEAT1* have significantly higher *NEAT1* expression and worse prognosis. High  
416 expression of *NEAT1* predicts significantly worse survival in ccRCC as well. *NEAT1* has been  
417 shown to be hypermutated in other cancers and some studies also linked high *NEAT1* association  
418 with unfavorable prognosis in several other tumors (23-24). Last, a downstream lncRNA,  
419 *MALAT1*, shows tight co-expression pattern with *NEAT1* in both pRCC and ccRCC. *MALAT1* is  
420 in COSMIC consensus cancer gene list and annotated as related with pediatric RCCs (REF).

421 With abundant reads from WGS, we generated a high confident SV dataset for 35 pRCC  
422 samples. Our method shows high accuracy and predicts CDKN2A expression level compared to  
423 the array-based approach by TCGA (3).

Shantao 2/24/2017 2:44 PM

Deleted: non-coding

Shantao 2/24/2017 2:44 PM

Deleted: non-coding

Shantao 2/24/2017 2:44 PM

Deleted: non-coding

Shantao 2/24/2017 2:44 PM

Deleted: non-coding

428 We found the pRCC clearly split into two groups: one stable group has less than 10  
429 events per sample while the unstable group all has above 40. Moreover, the unstable type is  
430 tightly associated with C2b group, which has inferior outcomes (3). Our SV study also finds  
431 recurrent cases of *SDHB* deletion and expression data supports our finding. *SDHB* is a subunit of  
432 succinate dehydrogenase. Previous studies indicated the loss of *SDHB* being a driver event as it  
433 disturbed tumor metabolic environment (REF BMS paper, AAH paper). Besides *SDHB*, we also  
434 found some other sporadic events involving known tumor drivers. Surprisingly, despite extensive  
435 *MET* copy number gain in pRCC, we did not detect an enrichment of smaller SV events of *MET*.  
436 We postulate polysomy 7 might be the major mechanism of *MET* gain and lack of smaller SVs  
437 disrupting *MET* further support the oncogene role of *MET*.

Shantao 2/24/2017 4:46 PM  
Formatted: Font:Italic

Shantao 2/24/2017 4:46 PM  
Formatted: Font:Italic

Shantao 2/24/2017 4:46 PM  
Formatted: Font:Italic

Shantao 2/24/2017 4:48 PM  
Formatted: Font:Italic

Shantao 2/24/2017 4:49 PM  
Formatted: Font:Italic

Shantao 2/24/2017 4:32 PM  
Deleted: .

Shantao 2/24/2017 4:49 PM  
Formatted: Font:Italic

Shantao 2/24/2017 4:49 PM  
Formatted: Font:Italic

438 WGS provides many times more SNVs compared to WXS, and noncoding SNVs are less  
439 constrained by selection pressure. Thus it gives us an opportunity to look into the high-level  
440 landscape of mutations in pRCC. Several recent landmark pan-cancer studies lead to the wide  
441 recognition of significance and great research interests in cancer mutational processes (REF).  
442 DNA mutation is one of the driving forces of cancer development. Understanding the underlying  
443 processes and affecting factors that generate the mutations is vital in cancer studies. In particular  
444 we focus on revealing the underlying sources that fuel tumor heterogeneity, which is a key  
445 feature in pRCC.

Shantao 2/24/2017 7:43 PM  
Formatted: Highlight

446 We identified mutation rate dispersion of C-to-T in CpG motif contributes the most to the  
447 inter-sample mutation spectra variations. We further pinned down the cause of dispersion by  
448 showing the hypermethylated cluster, identified in the previous TCGA study (3), has higher C-  
449 to-T rate in CpGs. This hypermethylated cluster is associated with later stage, type 2 pRCC,  
450 *SETD2* mutation and worse prognosis (3). Although increased C-to-T in CpG is likely the result

Shantao 2/18/2017 8:20 PM  
Deleted: Last, focusing on the high-level landscape of mutations in pRCC, w

454 of hypermethylation, we cannot rule out the possibility the change of mutation landscape plays a  
455 role in cancer development. For example, C-to-T in methylated CpG causes loss of methylation,  
456 which could have effects on local chromatin environment, trans-elements recruitment and gene  
457 expression regulation. In our study, we observed C-to-Ts in CpG are enriched in coding regions,  
458 which indicates they have higher functional impacts in cancer genome.

459 Significant APOBEC activities and consequential mutation signatures were observed in  
460 one type 2 pRCC case. APOBEC activities were known to be prevalent in UCs (12, 19). We also  
461 successfully detected prominent APOBEC signatures in all three UC samples processed in the  
462 same pipeline as pRCCs. Intriguingly, despite being considered to have the same cellular origin  
463 with pRCC, we were not able to detect significant APOBEC activities in ccRCC. This is in  
464 agreement with previous studies (12). APOBEC mutation signature was also found in a small  
465 percentage of chromophobe renal cell carcinoma (25), although they are believed to have a  
466 different cellular origin. APOBEC activities have been linked with genetic predisposition and  
467 viral infection (26). Given a statistically robust signal in our conservative algorithm, it is  
468 plausible that a small fraction of otherwise driver mutation absent type 2 pRCCs might share  
469 some etiologically and gnomically similarity with UC. Standard treatment for UC differs  
470 significantly from the one for pRCC. Pending further research, this finding might lead to  
471 actionably clinical implications.

472 Chromatin remodeling pathway is highly mutated in pRCC (3). Several chromatin  
473 remodelers, for example *SETD2* and *PBRM1*, have been identified as cancer drivers in pRCC.  
474 We investigate the relationship between samples with mutated chromatin remodelers and those  
475 without such mutations in terms of overall mutational spectrum. We demonstrated pRCC with  
476 defects in chromatin remodeling genes shows higher mutation rate in general, driven by an even

- Shantao 1/24/2017 2:18 AM  
**Deleted:** b
- Shantao 1/24/2017 2:18 AM  
**Deleted:** e
- Shantao 1/24/2017 2:18 AM  
**Deleted:** genomically
- Shantao 1/24/2017 2:19 AM  
**Deleted:** to
- Shantao 1/24/2017 2:19 AM  
**Deleted:** ince s
- Shantao 2/24/2017 4:27 PM  
**Deleted:**
- Shantao 2/24/2017 4:27 PM  
**Deleted:** involves cytotoxic chemotherapy and radiation
- Shantao 1/24/2017 2:22 AM  
**Formatted:** Highlight
- Shantao 1/24/2017 2:19 AM  
**Deleted:** !
- Shantao 1/24/2017 2:19 AM  
**Deleted:** this
- Shantao 1/24/2017 2:20 AM  
**Deleted:** could have
- Shantao 1/24/2017 2:20 AM  
**Deleted:** a very
- Shantao 1/24/2017 2:20 AM  
**Deleted:** meaningful clinical impact.
- Shantao 2/18/2017 8:05 PM  
**Deleted:** , *BAP1*

491 stronger mutation rate increase in putative open chromatin regions in normal kidney tissues. This  
492 is likely because chromatin remodeling defects affect normal open chromatin environment and  
493 impede DNA repairing in these regions.

494 It is known that replication time strongly governs local mutation rate. Early replication  
495 regions have fewer mutations. But the difference dissipates when DNA mismatch repair becomes  
496 defective (21). In our study, we found this correlation weakened in chromatin remodeling genes  
497 mutated samples, presumably caused by failure of replication error repair in an abnormal  
498 chromatin environment. By adapting defects in chromatin remodeling genes, tumor alters its  
499 mutation rate and landscape, which might further provide advantage in cancer evolution. Yet,  
500 high mutation burden in functional important open chromatin regions also raises the chance that  
501 tumor antigens activate host immune system. Researchers found tumors with DNA mismatch  
502 repair deficiency response better to PD-1 blockage (27). These tumors also accumulate more  
503 mutations in early replicated regions (21). Thus chromatin remodeler alterations might as well  
504 correlate with higher response rate of immunotherapy,

505 Last, we constructed individual evolutionary trees for all 35 samples. This is the first  
506 study inferring tumor evolutionary trees using large number of SNVs from WGS in pRCC.  
507 Benefited from a large number of SNVs, the tree construction becomes more accurate and  
508 reveals more details. Evolution trees reveal the history of tumor evolution and how mutations  
509 accumulate. We discovered the trees show four major types of topologies and reflected tumor  
510 heterogeneity. Type 2 pRCCs show a distinct evolutionary profile, indicating they are more  
511 heterogeneous. Evolutionary trees give us the opportunity to observe how pRCC heterogeneity  
512 develops over time.

Shantao 1/24/2017 2:24 AM

Deleted: due to

514 In this first whole genome study of pRCC, we found several novel noncoding alterations  
515 that might drive tumor development and we explored mutation landscape and evolution profiles  
516 to better understand tumor heterogeneity†. However, due to a limited sample size, some of our  
517 statistical tests were underpowered. As the cost of sequencing keeps dropping, we expect to have  
518 more pRCC whole genome sequenced in the near future (28). With a larger cohort, we hope to  
519 gain enough power to test the hypotheses we formed as well as further explore the noncoding  
520 regions of pRCC.

Shantao 2/24/2017 2:44 PM

Deleted: non-coding

Shantao 2/24/2017 7:43 PM

Deleted: have meaningful clinical impacts

521

## 522 **Materials and Methods**

### 523 **Data acquisition**

524 We downloaded pRCC and ccRCC WXS and pRCC WGS variation calls from TCGA  
525 Data Portal (<https://tcga-data.nci.nih.gov/tcga/tcgaDownload.jsp>) and TCGA Jamboree  
526 respectively. pRCC RNAseq, RPPA and methylation data were downloaded from TCGA Data  
527 Portal as well. Repli-seq data was obtained from ENCODE (<https://www.encodeproject.org/>).

528 DHS data was obtained from Roadmap Epigenomics Project

529 (<http://www.roadmapepigenomics.org>)

Shantao 2/24/2017 4:51 PM

Deleted: and DHS data

Shantao 2/24/2017 4:51 PM

Deleted: were

530

### 531 **Testing rs11762213 on prognosis and exploring somatic mutations in *MET***

532 We downloaded pRCC clinical outcomes from TCGA Data Portal ([https://tcga-](https://tcga-data.nci.nih.gov/tcga/tcgaDownload.jsp)  
533 [data.nci.nih.gov/tcga/tcgaDownload.jsp](https://tcga-data.nci.nih.gov/tcga/tcgaDownload.jsp)). pRCC samples that failed the histopathological review  
534 were excluded (3). In total, we included 277 patients in our analyses (Figure S1, Table S1). For  
535 germline calls, the majority of samples, 163 out of 277, were supported by germline SNV

540 callings from two centers (BCM and BI), 100% genotype concordance rate was observed. Also,  
541 162 curated rs11762213 genotypes were in agreement with automated callsets. All calls has  
542 alternative allelic fraction of 0.42 to 0.68, supporting a heterozygous genotype (9). Calls from BI  
543 all have genotype quality scores >125 and all calls in BCM pass the filter. With proved high  
544 confidence in accuracy of genotyping rs11762213 in germline, we recruited additional 114  
545 samples from single-center (BCM), automated calls to form an extensive patients set (Figure S1).  
546 For somatic SNVs in *MET*, after excluding cases that were recruited in the TCGA study, we  
547 formed an additional set encompassing 117 patients. Five callings were supported by two  
548 centers. The rest were supported by single-center (BCM) automated calls.

549 Cancer-specific survival was defined using the same criteria as described in a ccRCC  
550 study (9). Deaths were considered as cancer-specific if the “Personal Neoplasm Cancer Status” is  
551 “With Tumor”. If “Tumor Status” is not available, then the deceased patients were classified as  
552 cancer-specific death if they had metastasis (M1) or lymph node involvement ( $\geq$  N1) or died  
553 within two years of diagnosis. An R package, “survival”, was used for the survival analysis.

554

#### 555 **SV calling procedure**

556 We remapped the reads using bwa 0.7.12, which supports split read mapping. Then  
557 we used DELLY (10) with default parameters for somatic SV calling. To avoid sample  
558 contamination or germline SVs, we filtered our callsets against the entire TCGA pRCC WGS  
559 dataset, regardless of sample match or pathological reviews. We discharge all callings that were  
560 marked “LowQual” (PE/SR support below 3 or mapping quality below 20). Last, to further  
561 eliminate germline contamination, we filtered out SVs that show at least 0.8 reciprocally  
562 overlapping with 1000 Genome Phase 3 SV callset (only 1/425 filtered out).

Shantao 2/21/2017 6:33 PM

Deleted: at least

Shantao 2/21/2017 6:33 PM

Deleted: (102 from three centers).

Shantao 2/17/2017 7:05 PM

Deleted: procedure

Shantao 2/17/2017 7:07 PM

Formatted: Indent: First line: 0"

Shantao 2/17/2017 7:07 PM

Deleted: -

... [1]

Shantao 2/17/2017 7:07 PM

Deleted: 2

Shantao 2/17/2017 7:08 PM

Deleted: Lastly, we

570 For Lumpy(REF), we ran it with default parameters. We also filtered the results using the  
571 1000 Genome Phase 3 callset and required the SV have both paired-end and split reads support.

Shantao 2/17/2017 7:07 PM

Deleted: .

572

### 573 **Mutation spectra study**

574 WGS Mutations were extracted from flanking 5' and 3' nucleotide context. The raw  
575 mutation counts were normalized by trinucleotide frequencies in the whole genome.

576 To identify signatures in the mutation spectra, we used a robust, objective LASSO-based  
577 method. First, 30 known signatures were downloaded from COSMIC  
578 (<http://cancer.sanger.ac.uk/cosmic/signatures>). Then we solve a positive, zero-intercept linear  
579 regression problem with L1 regularizer to obtain signatures and corresponding weights for each  
580 genome. Specifically, we solve the problem:

$$\min_W (\|SW - M\|_2 + \lambda \|W\|)$$

581 Where M is the mutation matrix, containing the mutations of each sample in 96  
582 nucleotide contexts. S is the 96×30 signature matrix, representing the mutation probability in 96  
583 nucleotide contexts of the 30 signatures. W is the weighting matrix, representing the contribution  
584 of 30 signatures to each sample.

585 The penalty parameter lambda ( $\lambda$ ) was determined empirically using 10-fold cross-  
586 validation individually for every sample.  $\lambda$  was chosen to maximize sparsity and constrained to  
587 keep mean-square error (MSE) within one standard error of its minimum. Last, we discharged  
588 signatures that composite less than 5% of the total detectable signatures.

589

591 **Methylation association analysis**

592 In total, we collected HumanMethylation450 BeadChip array data for 139 samples that  
593 are either methylation cluster 1 or 2. We used an R package “IMA” to facilitate analysis (11).  
594 After discharging sites with missing values or on sex chromosomes, we obtained beta-values on  
595 366,158 CpG sites in total. Then we test beta-values of each site by Wilcoxon rank sum test  
596 between two methylation clusters. After adjusting p-value using Benjamini-Hochberg procedure,  
597 we called 9,324(2.55%) hypermethylation sites. These sites have an adjusted p-value of less than  
598 0.05 and mean beta-values in methylation cluster 1 are 0.2 or higher than the ones in methylation  
599 cluster 2.

600

601 **APOBEC enrichment analysis**

602 We used the method described by Roberts et al. (12). For every  $C \in \{T, G\}$  and  $G \in \{A, C\}$   
603 mutation we obtained 20bp sequence both upstream and downstream. Then enrichment fold was  
604 defined as:

$$Enrichment\ Fold = \frac{Mutation_{TCW/WGA} \times Context_{C/G}}{Mutation_{C/G} \times Context_{TCW/WGA}}$$

605 Here TCW/WGA stands for  $T[C \in \{T, G\}]W$  and  $W[G \in \{A, C\}]A$ . W stands for A or T. p-  
606 value for enrichment were calculated using one-sided Fisher-exact test. To adjust for multiple  
607 hypothesis testing, p-values were corrected using Benjamini-Hochberg procedure.

608 WXS data for APOBEC enrichment and signature analysis was obtained from a high  
609 quality somatic callset: hgsc.bcm.edu\_KIRP.IlluminaGA\_DNASeq.1.protected.maf. This dataset

Shantao 2/24/2017 5:46 PM

Deleted: one-side



611 includes 155 pRCC samples and three UC samples. We use  
612 hgsc.bcm.edu\_KIRC.Mixed\_DNASeq.1.protected.maf for ccRCC analyses.

613

#### 614 Chromatin remodeling genes and replication time association

615 We identified chromatin remodeling genes based on its significance in pRCC and  
616 function. Our gene list is the intersection of gene lists in the original TCGA pRCC study  
617 molecular feature table (supplementary table 3) with the chromatin remodeling and SNI/SWF  
618 pathway gene lists (supplementary table 4). Our gene set include ten genes: *SETD2, KDM6A,*  
619 *PBRM1, SMARCB1, ARID1A, ARID2, MLL2 (KMT2D), MLL3(KMT2C), MLL4(KMT2B),*  
620 *EP300.* We noticed *BAP1* is not in the gene list. However, adding BAP1 into the list does not  
621 change the significance of our key tests (p<0.0115 for mutation counts in DHS and p<0.020 for  
622 mutation percentage in DHS). We defined chromatin remodeling defect as nonsynonymous  
623 mutations in these genes. For missense mutations, we additionally filtered out mutations with  
624 polyphen score less then 0.8 (benign).

625 For replication time, in order to avoid cell type redundancy, we only kept GM12878 as  
626 the representative of all lymphoblastoid cell lines. Eleven cell types were included in our  
627 analysis: BG02ES, BJ, GM12878, HeLaS3, HEPG2, HUVEC, IMR90, K562, MCF7, NHEK,  
628 SK-NSH. Wave smoothed replication time signal was averaged in a  $\pm 10$ kb region from every  
629 mutation. To avoid potential selection effects, we removed mutations in exome and flanking 2bp.  
630 Regions overlap with reference genome gaps and DAC blacklist (<https://genome.ucsc.edu/>) were  
631 removed as well. Last, we picked the median number from 11 cell types at each mutation  
632 position for further analysis.

Shantao 2/17/2017 7:03 PM  
Deleted: .  
Shantao 2/17/2017 7:02 PM  
Formatted: Indent: First line: 0.5"

Shantao 2/17/2017 7:05 PM  
Formatted: Font:Italic

Shantao 2/24/2017 4:55 PM  
Formatted: Font:Italic

Shantao 2/17/2017 7:03 PM  
Formatted: Highlight

Shantao 2/17/2017 7:03 PM  
Formatted: Highlight

Shantao 2/17/2017 7:03 PM  
Formatted: Highlight

Shantao 2/17/2017 7:03 PM  
Formatted: Highlight

Shantao 2/17/2017 7:01 PM  
Deleted: included eleven genes. They are  
*ARID1A, ARID2, BAP1, DNMT3A, KDM6A,*  
*MLL2, MLL3, MLL4, PBRM1, SETD2,*  
*SMARCB1.*

Shantao 2/17/2017 7:03 PM  
Formatted: Highlight

Shantao 2/24/2017 5:02 PM  
Deleted: I

639 To test the significance of replication time of noncoding mutations between two groups,  
640 we assigned all the mutation with its local replication time and then defined the ones stand above  
641 90 percentile in all pooled mutations as “mutations in early replicating regions”. Then we  
642 calculate the percentage of “mutations in early replicating regions” in total mutations for each  
643 sample and compare between two groups using rank-sum test.  
644  
645 **Evolution tree inference:**  
646 We use PhyloWGS (REF) to infer the evolution trees for each individual tumor. To mitigate the  
647 effects on copy number change, we removed all the SNVs inside the copy number change  
648 regions as defined by assay-based method in the original TCGA study (REF). To be prudent, we  
649 defined any region with an absolute log tumor copy number to normal ratio larger than 0.3.  
650 Additionally, we then removed all SNVs with allele frequency higher than 0.6 as they are likely  
651 affected by copy number loss.

Shantao 2/24/2017 2:44 PM

**Deleted:** non-coding

Shantao 2/17/2017 6:51 PM

**Deleted:** we adapted a conservative non-parametric Kolmogorov–Smirnov test (K-S test) using empirical p-value.

Shantao 2/17/2017 6:51 PM

**Deleted:** W

Shantao 2/17/2017 6:52 PM

**Deleted:** its percentile among all mutations replication time shifted  $\pm 100\text{kb}$  from the origin (represents the background replication time)

Shantao 2/17/2017 6:52 PM

**Deleted:** K-S test

Shantao 2/17/2017 6:54 PM

**Deleted:** statistics

Shantao 2/17/2017 6:54 PM

**Deleted:** in

Shantao 2/17/2017 6:54 PM

**Deleted:** and compare

Shantao 2/17/2017 6:52 PM

**Deleted:** To obtain the empirical p-value, we randomly permuted the chromatin remodeling genes mutation labels for 1,000 times to estimate the test statistics distribution under null hypothesis.

Shantao 2/18/2017 8:10 PM

**Formatted:** Font:Bold

652

671 **Author contributions:** SL, BMS and MG conceived and designed the study. SL carried out the  
672 computation and data analysis, SL, BMS and MG interpreted the results. SL wrote the  
673 manuscript. BMS and MG co-directed this work. All authors have read and approved the final  
674 manuscript. **Competing interests:** The authors declare no competing interests.

675 **Acknowledgments:** This work was supported by the National Institutes of Health, AL Williams  
676 Professorship, and in part by the facilities and staffs of the Yale University Faculty of Arts and  
677 Sciences High Performance Computing Center. We thank Patrick McGillivray for his help in  
678 manuscript preparation.

679

## 680 **References**

681

- 682 1. Siegel, R, Naishadham, D, Jemal, A. Cancer statistics, 2015. CA: a cancer journal for  
683 clinicians. 2015; 65(1), 5-29.
- 684 2. Shuch B, Amin A, Armstrong AJ, Eble JN, Ficarra V, Lopez-Beltran A, et al.  
685 Understanding pathologic variants of renal cell carcinoma: distilling therapeutic  
686 opportunities from biologic complexity. European urology. 2015;67(1):85-97.
- 687 3. Cancer Genome Atlas Research Network. Comprehensive molecular characterization of  
688 papillary renal-cell carcinoma. N Engl J Med. 2016;2016(374):135-45.
- 689 4. Khurana E, Fu Y, Colonna V, Mu XJ, Kang HM, Lappalainen T, et al. Integrative  
690 annotation of variants from 1092 humans: application to cancer genomics. Science.  
691 2013;342(6154):1235587.
- 692 5. Fu Y, Liu Z, Lou S, Bedford J, Mu XJ, Yip KY, et al. FunSeq2: a framework for  
693 prioritizing noncoding regulatory variants in cancer. Genome biology. 2014;15(10):1.

- 694 6. Huang FW, Hodis E, Xu MJ, Kryukov GV, Chin L, Garraway LA. Highly recurrent  
695 TERT promoter mutations in human melanoma. *Science*. 2013;339(6122):957-9.
- 696 7. Alexandrov LB, Nik-Zainal S, Wedge DC, Campbell PJ, Stratton MR. Deciphering  
697 signatures of mutational processes operative in human cancer. *Cell reports*.  
698 2013;3(1):246-59.
- 699 8. Alexandrov LB, Nik-Zainal S, Siu HC, Leung SY, Stratton MR. A mutational signature  
700 in gastric cancer suggests therapeutic strategies. *Nature communications*. 2015;6.
- 701 9. Hakimi AA, Ostrovnaya I, Jacobsen A, Susztak K, Coleman JA, Russo P, et al.  
702 Validation and genomic interrogation of the MET variant rs11762213 as a predictor of  
703 adverse outcomes in clear cell renal cell carcinoma. *Cancer*. 2016;122(3):402-10.
- 704 10. Rausch T, Zichner T, Schlattl A, Stütz AM, Benes V, Korbel JO. DELLY: structural  
705 variant discovery by integrated paired-end and split-read analysis. *Bioinformatics*.  
706 2012;28(18):i333-9.
- 707 11. Wang D, Yan L, Hu Q, Sucheston LE, Higgins MJ, Ambrosone CB, et al. IMA: an R  
708 package for high-throughput analysis of Illumina's 450K Infinium methylation data.  
709 *Bioinformatics*. 2012;28(5):729-30.
- 710 12. Roberts SA, Lawrence MS, Klimczak LJ, Grimm SA, Fargo D, Stojanov P, et al. An  
711 APOBEC cytidine deaminase mutagenesis pattern is widespread in human cancers.  
712 *Nature genetics*. 2013;45(9):970-6.
- 713 13. Schmidt L, Junker K, Weirich G, Glenn G, Choyke P, Lubensky I, et al. Two North  
714 American families with hereditary papillary renal carcinoma and identical novel  
715 mutations in the MET proto-oncogene. *Cancer research*. 1998;58(8):1719-22.

- 716 14. Schutz FA, Pomerantz MM, Gray KP, Atkins MB, Rosenberg JE, Hirsch MS, et al.  
717 Single nucleotide polymorphisms and risk of recurrence of renal-cell carcinoma: a cohort  
718 study. *The lancet oncology*. 2013;14(1):81-7.
- 719 15. Guo S, Chen W, Luo Y, Ren F, Zhong T, Rong M, et al. Clinical implication of long  
720 noncoding RNA NEAT1 expression in hepatocellular carcinoma patients. *International*  
721 *journal of clinical and experimental pathology*. 2015;8(5):5395.
- 722 16. Nik-Zainal S, Davies H, Staaf J, Ramakrishna M, Glodzik D, Zou X, et al. Landscape of  
723 somatic mutations in 560 breast cancer whole-genome sequences. *Nature*.  
724 2016;534(7605):47-54.
- 725 17. Choudhry H, Albukhari A, Morotti M, Haider S, Moralli D, Smythies J, et al. Tumor  
726 hypoxia induces nuclear paraspeckle formation through HIF-2 $\alpha$  dependent transcriptional  
727 activation of NEAT1 leading to cancer cell survival. *Oncogene*. 2015;34(34):4482-90.
- 728 18. Chakravarty D, Sboner A, Nair SS, Giannopoulou E, Li R, Hennig S, et al. The oestrogen  
729 receptor alpha-regulated lncRNA NEAT1 is a critical modulator of prostate cancer.  
730 *Nature communications*. 2014;5.
- 731 19. Cancer Genome Atlas Research Network. Comprehensive molecular characterization of  
732 urothelial bladder carcinoma. *Nature*. 2014;507(7492):315-22.
- 733 20. Sato Y, Yoshizato T, Shiraishi Y, Maekawa S, Okuno Y, Kamura T, et al. Integrated  
734 molecular analysis of clear-cell renal cell carcinoma. *Nature genetics*. 2013;45(8):860-7.
- 735 21. Supek F, Lehner B. Differential DNA mismatch repair underlies mutation rate variation  
736 across the human genome. *Nature*. 2015;521(7550):81-4.

Shantao 2/24/2017 2:44 PM

**Deleted:** non-coding

- 738 22. Gordon MS, Hussey M, Nagle RB, Lara PN, Mack PC, Dutcher J, et al. Phase II study of  
739 erlotinib in patients with locally advanced or metastatic papillary histology renal cell  
740 cancer: SWOG S0317. *Journal of Clinical Oncology*. 2009;27(34):5788-93.
- 741 23. Li Y, Li Y, Chen W, He F, Tan Z, Zheng J, et al. NEAT expression is associated with  
742 tumor recurrence and unfavorable prognosis in colorectal cancer. *Oncotarget*.  
743 2015;6(29):27641.
- 744 24. He C, Jiang B, Ma J, Li Q. Aberrant NEAT1 expression is associated with clinical  
745 outcome in high grade glioma patients. *Apmis*. 2016;124(3):169-74.
- 746 25. Davis CF, Ricketts CJ, Wang M, Yang L, Cherniack AD, Shen H, et al. The somatic  
747 genomic landscape of chromophobe renal cell carcinoma. *Cancer cell*. 2014;26(3):319-  
748 30.
- 749 26. Henderson S, Chakravarthy A, Su X, Boshoff C, Fenton TR. APOBEC-mediated  
750 cytosine deamination links PIK3CA helical domain mutations to human papillomavirus-  
751 driven tumor development. *Cell reports*. 2014;7(6):1833-41.
- 752 27. Le DT, Uram JN, Wang H, Bartlett BR, Kemberling H, Eyring AD, et al. PD-1 blockade  
753 in tumors with mismatch-repair deficiency. *New England Journal of Medicine*.  
754 2015;372(26):2509-20.
- 755 28. Muir P, Li S, Lou S, Wang D, Spakowicz DJ, Salichos L, et al. The real cost of  
756 sequencing: scaling computation to keep pace with data generation. *Genome biology*.  
757 2016;17(1):1.

758

759

760

761

Shantao 2/24/2017 5:04 PM

**Formatted: Font:Bold**

Shantao 2/24/2017 5:04 PM

**Formatted: Normal**

762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789

**Figure 1. MET noncoding alterations and Survival analysis of rs11762213 in pRCC patients.**

(A) A schematics diagram of noncoding mutations on MET. The germline exonic SNP, rs11762213, is also shown. Thin black lines indicate the alternate transcript initiated by retrotransposon.  
(B) Genotypes of rs11762213 are shown in the legend. Peto & Peto modification of the Gehan-Wilcoxon test.

**Figure 2. Noncoding alterations in pRCC.**

(A) A schematics diagram of noncoding mutations in ERFF1. (B) A schematics diagram of noncoding mutations in NEAT1. One tumor carries two mutations on NEAT1. (C) Tumors with mutations on NEAT1 show higher NEAT1 expression. (D) Survival analysis shows mutations in NEAT1 are associated with worse prognosis. To avoid potential confounding effects, we removed one subject who carries rs11762213 but not NEAT1 mutation. Log-rank test.

**Figure 3. Mutation spectra and mutation processes in pRCC.**

(A) The mutation spectrum of all pRCC WGS samples. Mutations are ordered in alphabetical order of the reference trinucleotides (with the mutated nucleotide in the middle, from A[C>A]A to T[T>G]T) from left to right. Then we use PCA to maximize inter-sample variation. The loadings on the first principle component is strongly dominated by C>T in CpGs. (B) PC1, along with C>T in CpGs mutation counts and the fractions of such mutations among total mutations are significantly different between two methylation groups. (C) APOBEC mutation signatures are shown for both pRCC (along with three UC sampels, which have blue outer circles) and ccRCC TCGA cohorts. Red dashed line represents the median APOBEC enrichment. (D) Comparison of total mutation counts, mutations counts in open chromatin regions and percentages of mutations in open chromatin regions of total mutations between tumors with chromatin remodeling genes alterations and the ones without.

**Figure 4. Evolution trees and genomic alteration landscape of 35 whole genome sequenced pRCC samples.**

(A) Two individual evolutions trees. Mutations in cancer related gene are shown in colors corresponding to where it first appear. (B) Summary table of alterations in pRCC WGS. Index: patient index, see Table S2

Shantao 2/18/2017 9:59 PM  
**Deleted:** (A) A schematics diagram of non-coding mutations on MET. The germline SNP, rs11762213, is also shown.

Shantao 2/18/2017 10:00 PM  
**Deleted:** B

Shantao 2/24/2017 2:44 PM  
**Deleted:** non-coding

Shantao 2/24/2017 5:18 PM  
**Deleted:** o

Shantao 2/18/2017 10:00 PM  
**Deleted:** C

Shantao 2/24/2017 2:44 PM  
**Deleted:** non-coding

Shantao 2/24/2017 5:18 PM  
**Deleted:** o

Shantao 2/18/2017 10:00 PM  
**Deleted:** D

Shantao 2/18/2017 10:00 PM  
**Deleted:** E

Shantao 2/18/2017 10:00 PM  
**Deleted:** (B)

Shantao 2/18/2017 10:00 PM  
**Deleted:** W

Shantao 2/18/2017 10:00 PM  
**Deleted:** C

Shantao 2/18/2017 10:00 PM  
**Deleted:** D

Shantao 2/18/2017 10:00 PM  
**Deleted:** E

Shantao 2/18/2017 10:01 PM  
**Deleted:** The

Shantao 2/18/2017 10:01 PM  
**Deleted:** 32

Shantao 2/18/2017 10:01 PM  
**Deleted:** Grey cells represent genomic alterations. CN: copy number.