



24 **Author Summary**

25 Renal cell carcinoma accounts for more than 90% of kidney cancers. Papillary renal cell  
26 carcinoma (pRCC) is the second most common subtype of renal cell carcinoma. Previous studies,  
27 focusing mostly on the protein-coding regions, have identified several key genomic alterations  
28 that are key to cancer initiation and development. However, researchers cannot find any key  
29 mutation in a significant portion of pRCC. Therefore, we carry out the first whole-genome study  
30 of pRCC to discover triggering DNA changes explaining these cases. By looking at the entire  
31 genome, we find additional potentially impactful alterations in and out of the protein-coding  
32 regions. These newly identified critical mutations from scrutinizing the entire genome help  
33 complete our understanding of pRCC genomes. Two alterations we found are associated with  
34 prognosis, which could aid clinical decisions. We are also able to recognize mutation patterns  
35 and signatures, which reflect the mutagenesis processes and give hints on how cancer develops.  
36 Our study provides valuable additional information to facilitate better tumor subtyping, risk  
37 stratification and potentially clinical management.

38

39 **Introduction**

40 Renal cell carcinoma (RCC) makes up over 90% of kidney cancers and currently is the  
41 most lethal genitourinary malignancy (1). Papillary RCC (pRCC) accounts for 10%-15% of the  
42 total RCC cases (2). Unfortunately pRCC has been understudied and there are no current forms  
43 of effective systemic therapy for this disease. pRCC are further subtyped into two major groups:  
44 type 1 and type 2 based on histopathological features. For many years, the only prominent  
45 oncogene in pRCC (specifically, type 1) that physicians were able to identify was *MET*, a  
46 tyrosine kinase receptor for hepatic growth factor. An amino acid substitution that leads to

47 constitutive activation and/or overexpression are two mechanisms of dysfunction of *MET* in  
48 tumorigenesis. Recently, the Cancer Genome Atlas (TCGA) published its first result on pRCC  
49 (3), which greatly improves our understanding of the genomic basis of this disease. Several more  
50 genes and specific sub-clusters were identified to be significantly mutated in pRCC.  
51 Nevertheless, a significant portion of pRCC cases still remains without any known driver.  
52 Therefore we think it is time to explore the rest 98% non-coding regions of the genome using  
53 whole genome sequencing (WGS). This is sensible because non-coding regions, previously  
54 overlooked in cancer, have been showed to be actively involved in tumorigenesis (4-6).  
55 Mutations in non-coding regions may cause disruptive changes in both cis- and trans-regulatory  
56 elements, affecting gene expression. Understanding non-coding mutations helps fill the missing  
57 “dark matter” in cancer research.

58 Multiple endogenous and environmental mutation processes shape the somatic mutational  
59 landscape observed in cancers (7). Analyses of the genomic alterations associated with these  
60 processes give information on cancer development, shed light on mutational disparity between  
61 cancer subtypes and even indicate potential new treatment strategies (8). Additionally, genomic  
62 features such as replication time and chromatin environment govern mutation rate along the  
63 genome, contributing to spatial mutational heterogeneity. While identifying mutation signatures  
64 is possible using data from whole exome sequencing (WXS), whole genome sequencing (WGS)  
65 gives richer information on mutation landscape and minimizes the potential confounding effects  
66 of exome capture process and driver selection.

67 In this study, we comprehensively analyzed 35 pRCC cases that were whole genome  
68 sequenced along with an extensive set of WXS data on multiple levels. We went from  
69 microscopic examination of driver genes to analyses of whole genome sequencing variants, and

Shantao 1/24/2017 12:04 AM  
Deleted: 32

71 finally, to investigation of high-order mutational features. First, we focused on *MET*, an  
72 oncogene which plays a central role in pRCC, especially in type 1. We found rs11762213, a  
73 germline exonic single nucleotide polymorphism inside *MET*, predicts cancer-specific survival  
74 (CSS) in type 2 pRCC. We also discovered several potentially impactful non-coding mutation  
75 hotspots in *MET* promoter and its first two exons. The previous TCGA study identifies a *MET*  
76 alternative transcription event as a driver event but without illustrating the etiology (3). We  
77 found that a cryptic promoter from a long interspersed nuclear element-1 (L1) triggers the  
78 alternative isoform expression. Surprisingly, we did not find a significant amount of structural  
79 variations affecting *MET* besides polysomy <sup>7</sup>. Then we went onto cases not as easily explained  
80 as those with *MET* alterations. We analyzed nearly 150,000 non-coding mutations throughout the  
81 entire genomes and found several potentially high-impact mutations in non-coding regions.  
82 Further zooming out, we discovered pRCC exhibits mutational heterogeneity in both nucleotide  
83 context and genome location, indicating underlying vibrant mutational processes interplay. We  
84 found methylation is the leading factor influencing mutation landscape. Methylation status drives  
85 the intra-sample mutation variation by promoting more C-to-T mutations in the CpG context.  
86 APOBEC activity, although infrequently observed, leaves an unequivocal mutation signature in a  
87 pRCC genome but not in ccRCC. Last, we discovered samples with chromatin remodeler  
88 alternations accumulate more mutations in open chromatin regions.

Shantao 1/23/2017 5:32 PM

Formatted: Font:Italic

Shantao 1/24/2017 2:29 AM

Deleted: 3

89

## 90 **Results**

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

92 We begin with coding variants in the long known driver *MET*. The TCGA study of 161  
93 pRCC patients found 15 samples carrying somatic, nonsynonymous single nucleotide variant

95 (SNV) in *MET*. By analyzing 117 extra WXS samples (see Methods), we found six more  
96 nonsynonymous somatic mutations in six samples (Table S1). V1110I and M1268T were two  
97 recurrent mutations in this extra set. Both of them were observed in the TCGA study as well.  
98 Additionally, we found two samples carrying H112Y and Y1248C respectively. H112Y has  
99 been observed in two patients the original TCGA study cohort and H1118R is a long-known  
100 germline mutation associated with hereditary papillary renal carcinoma (HPRC, 13). Y1248C  
101 has been observed in type 1 pRCC before (rs121913246) and the TCGA cohort has a case  
102 carrying Y1248H. All mutations occur in the hypermutated tyrosine kinase catalytic domain of  
103 *MET*. Two out of these six samples were identified as type 1 pRCC while the subtypes of the rest  
104 four were unknown.

105         Although many *MET* somatic mutations are believed to play a central role in pRCC,  
106 some germline *MET* mutations have also been associated with the disease. In particular, a  
107 germline SNP, rs11762213, has been discovered to predict recurrence and survival in a mixed  
108 RCC cohort (14). ccRCC predominated the initial discovery RCC cohort. This conclusion was  
109 later validated in a ccRCC cohort but never in pRCC (9). We wondered whether this SNP has a  
110 prognostic effect in pRCC. Using an extensive WXS set of 277 patients (see Methods; Figure S1  
111 and Table S1;), we found 14 patients carry one risk allele of rs11762213 (G/A, Table 1, minor  
112 allele frequency (MAF) = 2.53%). No homozygous A/A was observed. Cancer specific deceases  
113 are concentrated in type 2 pRCC. Among 96 type 2 pRCC cases, seven patients carry the minor  
114 A allele (MAF = 3.65%, Table 1). Survival is significantly worse in type 2 patients carrying the  
115 risk allele of rs11762213 ( $p = 0.034$ , Figure 1B). But we did not find significant association of  
116 this germline SNP with survival in type 1 patients. We did not find statistically significant  
117 association of rs11762213 with *MET* RNA expression in either tumor samples or normal controls

118 ( $p > 0.1$ , two-sided rank-sum test). *Met* pY1235 levels in tumor samples, as measured by Reverse  
 119 phase protein array (RPPA), were not significantly different in patients carrying the minor G  
 120 allele compared to patients with A/A genotype ( $p > 0.1$ , two-sided rank-sum test).

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), y</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)

121  
 122 **Table 1. Patient clinical profiles of the type 2 pRCC cohort in rs11762213 survival analysis.** AJCC: American  
 123 Joint Committee on Cancer; IQR: interquartile range; NA: not available. Percentages may not add up to 100%  
 124 because of rounding.

125

126 2. **Epigenetic alterations and mutation hotspots in non-coding regions**

127 The TCGA study has identified a *MET* alternative translation isoform as a driver event  
128 (3). However, the etiology of this new isoform is unknown. We identified this isoform results  
129 from the usage of a cryptic promoter from an L1 element, likely due to a local loss of  
130 methylation (REF). This event was reported in several other cancer types (REF). To test its  
131 relationship with methylation, we found a closet probe (cg06985664, ~3kb downstream) on the  
132 Methylation array show marginally statistically significant (p=0.055, one-side rank-sum test).  
133 Additionally, as expected, this event is associated with methylation group 1 (odds ration (OR)=  
134 4.54, 95%CI: 1.07-19.34, p<0.041), indicating genome-wide methylation dysfunction. This  
135 association is stronger in type 2 pRCC and it shows a significant association with the C2b cluster  
136 (OR= 17.5, 95%CI: 1.72-32.6, p<0.007).

137 Despite the fact *MET* is the most common driver alteration, about 20% presumably *MET*-  
138 driven yet *MET* wild-type pRCC samples were still left unexplained (3). Therefore, we scanned  
139 the *MET* non-coding regions. We observed one mutation in *MET* promoter region in a type 1  
140 pRCC sample (Figure 2A and Table S2). This sample shows no evidence of a nonsynonymous  
141 mutation in *MET* gene but it has copy number gain of *MET*. Additionally, we observed 6/35  
142 (17.1%) samples carry mutations in the intronic regions between exon 1-3 of *MET* (Figure 2A  
143 and Table S2). Previously it is been established that alternative splicing of these exons is a driver  
144 event (3). Therefore we speculated that these non-coding variants might correlate with the  
145 alternative splicing. However, likely being hindered by a small size, we were not able to find  
146 statistically significant association between the alternative splicing event and these intronic  
147 mutations.

Shantao 1/23/2017 5:41 PM  
Formatted: Font:Bold, Not Italic

Shantao 1/23/2017 5:40 PM  
Deleted: M

Shantao 1/24/2017 12:09 AM  
Deleted: 32

Shantao 1/24/2017 12:09 AM  
Deleted: 18.8

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

Shantao 1/24/2017 12:05 AM  
Deleted: 2  
Shantao 1/24/2017 12:05 AM  
Deleted: 8.8

160 Another potentially impactful mutation hotspot is in *NEAT1*. We saw mutations inside  
161 this nuclear long non-coding RNA in 6/35 (17.1%) samples (Figure 2C and Table S2). Several  
162 studies indicated *NEAT1* is associated in many other cancers (15-16). It promotes cell  
163 proliferation in hypoxia (17) and alters the epigenetic landscape, increasing transcription of  
164 target genes (18).

Shantao 1/24/2017 12:04 AM  
Deleted: 5  
Shantao 1/24/2017 12:05 AM  
Deleted: 2  
Shantao 1/24/2017 12:05 AM  
Deleted: 15.6

165 All the mutations we found fell into a putative promoter region of *NEAT1*. We noticed  
166 *NEAT1* mutations were associated with higher *NEAT1* expression (Figure 2D,  $p < 0.044$ , two-  
167 sided rank sum test). We also found *NEAT1* mutations were associated with worse prognosis  
168 (Figure 2E,  $p < 0.022$ , log-rank test). However, without mutation status, *NEAT1* expression level  
169 is not significantly linked with pRCC survival. Nonetheless, *NEAT1* is overexpressed in about  
170 6% ccRCC samples from the TCGA cohort. *NEAT1* overexpression is significantly associated  
171 with shorted overall survival (Fig SXX). *MALAT1*, another noticeable lncRNA in cancer, is  
172 tightly co-expressed with *NEAT1* in both pRCC and ccRCC. Overexpression of *MALAT1* is  
173 reported to be associated with cancer progression (REF).

Shantao 1/24/2017 12:15 AM  
Formatted: Font:Italic  
Shantao 1/24/2017 12:17 AM  
Formatted: Font:Italic  
Shantao 1/24/2017 12:18 AM  
Formatted: Font:Italic  
Shantao 1/24/2017 12:19 AM  
Formatted: Font:Not Italic  
Shantao 1/24/2017 12:18 AM  
Formatted: Font:Italic  
Shantao 1/24/2017 12:20 AM  
Formatted: Font:Italic



179 We used DELLY (10) to perform structural variants (SVs) discovery from WGS reads  
180 information (see Methods and Table S3). The SV discovery approach has higher sensitivity and  
181 resolution than array-based methods, which were employed in the TCGA analysis. In the end we  
182 found 343 somatic SV events, includes deletions, duplications, inversions and translocations. We  
183 confirmed three cases carrying deletions affecting *CDKN2A* called by TCGA array-based  
184 methods but not the other two cases, possibly due to large-scale events (aneuploidy). One  
185 sample, TCGA-B9-4116, which has extensive amplification of *MET*, showed multiple SVs of  
186 various classes hitting *MET* regions. However, surprisingly, we did not find SVs affecting *MET*  
187 except this one example. We postulate trisomy/polysomy 7 is the main mechanism of *MET*  
188 structural alteration rather than duplication in a smaller scale. Besides duplication, we did not  
189 expect to find deletion, inversion or translocation disrupting oncogene *MET*. These SVs are  
190 likely to cause loss-of-function rather than gain-of-function mutations. This is consistent with the  
191 putative role of *MET* as an oncogene, rather than a tumor suppressor. (Will work on this after  
192 the SV results come back)

Shantao 1/24/2017 2:29 AM  
Formatted: Highlight

### 194 3. Mutation spectra and mutation processes of pRCC

195 To further get a high-order overview of the mutation landscape, we summarized the  
196 mutation spectra of 35 whole genome sequenced pRCC samples (Figure 3A). C-to-T in CpGs  
197 showed the highest mutation rates, which were roughly ten to twenty-fold higher than mutation  
198 rates in other nucleotide contexts.

Shantao 1/24/2017 2:27 AM  
Deleted: 32

199 We used principle components analysis (PCA) to reveal factors that explain the most  
200 inter-sample variation. The loadings on the first principle component (which explained 12.5% of  
201 the variation) demonstrated C-to-T in CpGs contributed the most to inter-sample variation

203 (Figure 3B). C-to-T in CpGs is highly associated with methylation. It reflects the spontaneous  
204 deamination of cytosines in CpGs, which is much more frequent in 5-methyl-cytosines (REF).  
205 So we further explored the association between C-to-T in CpGs and tumor methylation status.  
206 First we validated the TCGA identified methylation cluster 1 showed higher methylation level  
207 than cluster 2 in all annotation regions (Figure S2, see Methods), prominently in CpG Islands  
208 (OR of sites being differentially hypermethylated: 1.29, 95%CI: 1.20-1.39, p<0.0001). We  
209 confirmed this association by showing samples from methylation cluster 1 had higher PC1 scores  
210 as well as higher C-to-T mutation counts and mutation percentages in CpGs (Figure 3C). This  
211 trend was further validated using a larger WXS dataset as well. Especially, the most  
212 hypermethylated group, CpG island methylation phenotype (CIMP), showed the greatest C-to-T  
213 in CpGs (Figure S2). As expected, C-to-T mutations in CpGs in group 1 showed higher but not  
214 statistically significant percentage overlapping with CpG islands compared with group 2 (1.8%  
215 versus 1.4%, p=0.14). Therefore, methylation status is the most prominent factor shaping the  
216 mutation spectra across patients. We further tried to explore the functional impact of the  
217 excessive mutations driven by methylation. C-to-T mutations in CpGs were more likely to be in  
218 the coding region (OR=1.54, 95%CI: 1.27-1.85, p<0.0001) and nonsynonymous (OR=1.47,  
219 95%CI: 1.17-1.84, p<0.001). Yet, C-to-T mutations in CpGs did not show functional bias  
220 between two methylation groups nor in non-coding regions (Figure SXX).

221 Recently, several somatic mutation signatures were identified. Many have putative  
222 etiology, revealing the underlying mutation processes (7). We used a LASSO-based approach  
223 (see Methods) to decompose mutations into a linear combination of these canonical mutation  
224 signatures in both WGS and WXS samples (Figure S3). The leading signature was signature 5,  
225 which is consistent with previous studies (7). Interestingly, we found one type 2 pRCC case out

Shantao 1/23/2017 6:07 PM

Deleted: (hypermethylated group, Figure S2)

Shantao 1/23/2017 9:09 PM

Deleted: .

228 of 155 somatic WXS sequenced samples exhibited APOBEC-associated mutation signature 2  
229 and 13. APOBEC mutation pattern enrichment analysis (see Method) further confirmed the  
230 presence of APOBEC activity (Figure 3D). This sample was statistically enriched of APOBEC  
231 mutations (adjusted p-value < 0.0003).

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

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

243 Noticeably, all four samples with APOBEC activities showed significantly higher  
244 *APOBEC3A* and *APOBEC3B* mRNA expression level ( $p < 0.0022$  and  $p < 0.0039$  respectively,  
245 one-side rank sum test, Figure S4). This is in concordance with previous studies of APOBEC  
246 mutagenesis in various types of cancer (12).

247 Consistent with previous studies (12), we failed to detect statistically significant  
248 APOBEC activities in an extensive WXS dataset consisting of 418 clear cell RCC (ccRCC)  
249 samples, even after resampling to avoid p-value adjustment eroding the power. Very low levels

250 of APOBEC signatures (<15%) was found in less than 1%(4/418) samples. With a much larger  
251 sample size, this result was unlikely to be confounded by detecting power.

252

253

#### 254 4. Defects in chromatin remodeling affects mutation landscape

255 Chromatin remodeling genes are frequently mutated in pRCC and many other cancers  
256 including ccRCC (20). Defects in chromatin remodeling cause dysregulation of chromatin  
257 environment. Open chromatin regions show lower mutation rate, presumably due to more  
258 effective DNA repair (21). Thus chromatin remodeler alternations could possibly alter the  
259 mutation landscape, specifically increase mutation rate in previously open chromatin regions. To  
260 test this hypothesis, we tallied the number of mutations inside DNase I hypersensitive sites  
261 (DHS) in HEK293, a cell line derived from human embryonic kidney cells, the closest match we  
262 could find in the ENCODE DHS database. 12/32 samples with non-silent mutations in eleven  
263 chromatin remodeling, cancer associated genes show higher genome-wide mutation counts ( $p <$   
264  $0.032$ , one-side rank-sum test), partially driven by higher mutation counts in DHS region ( $p <$   
265  $0.003$ , one-side rank-sum test). The median number of mutations in DHS region considerably  
266 increases by about 50% (75.5 versus 112) in samples carrying chromatin remodeling defects.  
267 The effect is significant after normalizing against the total mutation counts ( $p < 0.015$ , one-side  
268 rank-sum test, Figure 3E).

269 Replication time is known to correlate greatly with mutation rate. Early replicating  
270 regions have lower mutation rate compared to late replicating ones. Researchers reason  
271 replication errors are more likely to be corrected by DNA repair system in early replicating

272 regions. With defects in mutated chromatin remodeling, we observed this trend became less  
273 pronounced (Figure S5). This is likely because dysregulation of the chromatin environment  
274 hinders replication error repair by changing the accessibility of newly synthesized DNA chains.  
275 However, a non-parametric permutation Kolmogorov–Smirnov test (see Methods) failed to  
276 detect a statistical significance ( $p > 0.05$ ), likely because of the small number of samples and the  
277 prudence of our conserved test.

278

## 279 **Discussion**

280 We comprehensively analyzed both WGS and an extensive set of WXS of pRCC,  
281 scrutinizing local high-impact events as well as giving a macro overlook of the mutation  
282 landscape. Our work further completed the genomic alteration landscape of pRCC (Figure 4).  
283 Beyond traditionally driver events, we suggested several novel noncoding alterations potentially  
284 | drive tumorigenesis.

285 First, we elaborated on previous results of the long known driver *MET*. In an extended  
286 117 WXS dataset, we found six additional nonsynonymous somatic mutations in the  
287 hypermutated tyrosine kinase catalytic domain. These somatic mutations are highly recurrent,  
288 concentrated on a few critical amino acids. This is in line with *MET* being an oncogene and  
289 supports the central role of *MET* in pRCC. Then we found an exonic SNP in *MET*, rs11762213,  
290 to be a prognostic germline variance in type 2 pRCC. Previously, rs11762213 was found to  
291 predict outcome in a mixed RCC samples, predominated by ccRCC (14). Later, the result is  
292 confirmed in a large ccRCC cohort (9). However, it is never clear whether rs11762213 only  
293 predicts the outcome in ccRCC or other histological types as well. In this study, we concluded  
294 that the minor alternative allele of rs11762213 also forecasts unfavorable outcome in type 2

295 | pRCC patients. The mechanism of this exonic germline SNP remains unsettled. Remarkably,  
296 | similar to ccRCC, type 2 pRCC is not primarily driven by *MET*. Not significantly mutated in  
297 | ccRCC and type 2 pRCC, *MET* nonetheless seems to play a role in cancer development. This  
298 | finding is potentially meaningful in clinical management of patients with the more aggressive  
299 | type 2 pRCC. rs11762213 genotyping could become a reliable, low-cost risk stratification tool  
300 | for these patients. Theoretically, the subgroup of patients with rs11762213 might benefit from  
301 | *MET* inhibitors.

Shantao 1/24/2017 1:54 AM

**Deleted:** central

Shantao 1/24/2017 1:55 AM

**Deleted:** of

Shantao 1/24/2017 1:55 AM

**Deleted:** Potentially

Shantao 1/24/2017 2:32 AM

**Deleted:** also

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

Shantao 1/24/2017 2:11 AM

**Deleted:** Perhaps this variant could play a role in the significant racial disparities are known to exist in the overall incidence, histologic distribution, and survival of African Americans with kidney cancer.

308 | Besides, in *MET* non-coding regions, we also discovered mutations associated with *MET*  
309 | promoter and first two introns. Although the implication is unknown, our analysis suggests there  
310 | is a mutation hotspot in *MET* that calls for further research.

311 | Expanding our scope from coding to non-coding, we found several potentially significant  
312 | non-coding mutation hotspots relevant to tumorigenesis throughout the entire genome. A  
313 | mutation hotspot was found upstream of *ERRF1*, an important regulator of the EGFR pathway,  
314 | which may serve as a potential tumor suppressor. EGFR inhibitors have been used in papillary  
315 | kidney cancer with an 11% response rate observed (22). These mutations potentially disrupt  
316 | regulatory elements of *ERRF1* and thus play a role in tumorigenesis. However, likely limited by  
317 | a small sample size, we were not able to detect statistically significant functional changes in

327 ERRFI1 and related pathways. Another non-coding hotspot is in *NEATI*, a long non-coding RNA  
328 that has been speculated to involve in cancer. All mutations locate in a putative regulatory region  
329 of the gene. Patients carrying mutations in *NEATI* have significantly higher *NEATI* expression  
330 and worse prognosis. *NEATI* has been shown to be hypermutated in other cancers and some  
331 studies also linked high *NEATI* association with unfavorable prognosis in several other tumors  
332 (23-24).

333 Last, focusing on the high-level landscape of mutations in pRCC, we identified mutation  
334 rate dispersion of C-to-T in CpG motif contributes the most to the inter-sample mutation spectra  
335 variations. We further pinned down the cause of dispersion by showing the hypermethylated  
336 cluster, identified in the previous TCGA study (3), has higher C-to-T rate in CpGs. This  
337 hypermethylated cluster is associated with later stage, type 2 pRCC, *SETD2* mutation and worse  
338 prognosis (3). Although increased C-to-T in CpG is likely the result of hypermethylation, we  
339 cannot rule out the possibility the change of mutation landscape plays a role in cancer  
340 development. For example, C-to-T in methylated CpG causes loss of methylation, which could  
341 have effects on local chromatin environment, trans-elements recruitment and gene expression  
342 regulation. In our study, we observed C-to-Ts in CpG are enriched in coding regions, which  
343 supports their roles in cancer development.

344 Significant APOBEC activities and consequential mutation signatures were observed in  
345 one type 2 pRCC case. APOBEC activities were known to be prevalent in UCs (12, 19). We also  
346 successfully detected prominent APOBEC signatures in all three UC samples processed in the  
347 same pipeline as pRCCs. Intriguingly, despite being considered to have the same cellular origin  
348 with pRCC, we were not able to detect significant APOBEC activities in ccRCC. This is in  
349 agreement with previous studies (12). APOBEC mutation signature was also found in a small

350 percentage of chromophobe renal cell carcinoma (25), although they are believed to have a  
351 different cellular origin. APOBEC activities have been linked with genetic predisposition and  
352 viral infection (26). Given a statistically robust signal in our conservative algorithm, it is  
353 plausible that a small fraction of otherwise driver mutation absent type 2 pRCCs might share  
354 some etiologically and genomically similarity with UC. Standard treatment for UC involves  
355 cytotoxic chemotherapy and radiation while RCC shows low response rate to cytotoxic therapy.  
356 Pending further research, this finding might lead to actionably clinical implications (still too  
357 strong?).

358 Chromatin remodeling pathway is highly mutated in pRCC (3). Several chromatin  
359 remodelers, for example *SETD2*, *BAP1* and *PBRM1*, have been identified as cancer drivers in  
360 pRCC. We investigate the relationship between samples with mutated chromatin remodelers and  
361 those without such mutations in terms of overall mutational spectrum. We demonstrated pRCC  
362 with defects in chromatin remodeling genes shows higher mutation rate in general, driven by an  
363 even stronger mutation rate increase in putative open chromatin regions. This is likely because  
364 chromatin remodeling defects affect open chromatin environment and impede DNA repairing in  
365 these regions.

366 It is known that replication time strongly governs local mutation rate. Early replication  
367 regions have fewer mutations. But the difference dissipates when DNA mismatch repair becomes  
368 defective (21). In our study, we found this correlation weakened in chromatin remodeling genes  
369 mutated samples, presumably caused by failure of replication error repair in an abnormal  
370 chromatin environment. By adapting defects in chromatin remodeling genes, tumor alters its  
371 mutation rate and landscape, which might further provide advantage in cancer evolution. Yet,  
372 high mutation burden in functional important open chromatin regions also raises the chance that

- 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 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 1/24/2017 2:24 AM  
Deleted: due to



384 tumor antigens activate host immune system. Researchers found tumors with DNA mismatch  
385 repair deficiency response better to PD-1 blockage (27). Thus chromatin remodeler alterations  
386 might as well correlate with higher response rate of immunotherapy,

387         In this first whole genome study of pRCC, we found several novel non-coding alterations  
388 that might have meaningful clinical impacts. However, due to a limited sample size, our  
389 statistical tests were underpowered. As the cost of sequencing keeps dropping, we expect to have  
390 more pRCC whole genome sequenced in the near future (28). With a larger cohort, we hope to  
391 gain enough power to test the hypotheses we formed as well as further explore the noncoding  
392 regions of pRCC.

393

## 394 **Materials and Methods**

### 395 **Data acquisition**

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

401

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

403         We downloaded pRCC clinical outcomes from TCGA Data Portal ([https://tcga-](https://tcga-data.nci.nih.gov/tcga/tcgaDownload.jsp)  
404 [data.nci.nih.gov/tcga/tcgaDownload.jsp](https://tcga-data.nci.nih.gov/tcga/tcgaDownload.jsp)). pRCC samples that failed the histopathological review  
405 were excluded (3). In total, we included 277 patients in our analyses (Figure S1, Table S1). For

406 germline calls, the majority of samples, 163 out of 277, were supported by SNV callings from at  
407 least two centers (102 from three centers). 100% genotype concordance rate was observed. Also,  
408 162 curated rs11762213 genotypes were in agreement with automated callsets. With proved high  
409 confidence in accuracy of genotyping rs11762213 in germline, we recruited additional 114  
410 samples from single-center (BCM), automated calls to form an extensive patients set (Figure S1).  
411 For somatic SNVs in *MET*, after excluding cases that were recruited in the TCGA study, we  
412 formed an additional set encompassing 117 patients. Five callings were supported by two  
413 centers. The rest were supported by single-center (BCM) automated calls.

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

419

#### 420 **SV calling procedure**

421 We use DELLY2 (10) with default parameters for somatic SV calling. To avoid sample  
422 contamination or germline SVs, we filtered our callsets against the entire TCGA pRCC WGS  
423 dataset, regardless of sample match or pathological reviews. Lastly, we discharge all callings that  
424 were marked “LowQual” (PE/SR support below 3 or mapping quality below 20).

425

#### 426 **Mutation spectra study**

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

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

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

434 Where M is the mutation matrix, containing the mutations of each sample in 96  
435 nucleotide contexts. S is the 96×30 signature matrix, representing the mutation probability in 96  
436 nucleotide contexts of the 30 signatures. W is the weighting matrix, representing the contribution  
437 of 30 signatures to each sample.

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

442

#### 443 **Methylation association analysis**

444 In total, we collected HumanMethylation450 BeadChip array data for 139 samples that  
445 are either methylation cluster 1 or 2. We used an R package “IMA” to facilitate analysis (11).  
446 After discharging sites with missing values or on sex chromosomes, we obtained beta-values on  
447 366,158 CpG sites in total. Then we test beta-values of each site by Wilcoxon rank sum test

448 between two methylation clusters. After adjusting p-value using Benjamini-Hochberg procedure,  
449 we called 9,324(2.55%) hypermethylation sites. These sites have an adjusted p-value of less than  
450 0.05 and mean beta-values in methylation cluster 1 are 0.2 or higher than the ones in methylation  
451 cluster 2.

452

### 453 **APOBEC enrichment analysis**

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

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

457 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-  
458 value for enrichment were calculated using one-side Fisher-exact test. To adjust for multiple  
459 hypothesis testing, p-values were corrected using Benjamini-Hochberg procedure.

460 WXS data for APOBEC enrichment and signature analysis was obtained from a high  
461 quality somatic callset: hgsc.bcm.edu\_KIRP.IlluminaGA\_DNASeq.1.protected.maf. This dataset  
462 includes 155 pRCC samples and three UC samples. We use  
463 hgsc.bcm.edu\_KIRC.Mixed\_DNASeq.1.protected.maf for ccRCC analyses.

464

### 465 **Chromatin remodeling genes and replication time association**

466 We identified chromatin remodeling genes based on its significance in pRCC and  
467 function. Our gene list included eleven genes. They are *ARID1A*, *ARID2*, *BAP1*, *DNMT3A*,  
468 *KDM6A*, *MLL2*, *MLL3*, *MLL4*, *PBRM1*, *SETD2*, *SMARCB1*.

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

477 To test the significance of replication time of non-coding mutations between two groups,  
478 we adapted a conservative non-parametric Kolmogorov–Smirnov test (K-S test) using empirical  
479 p-value. We assigned all the mutation with its percentile among all mutations replication time  
480 shifted  $\pm 100$ kb from the origin (represents the background replication time). Then we calculate  
481 the K-S test statistics in two groups and compare. To obtain the empirical p-value, we randomly  
482 permuted the chromatin remodeling genes mutation labels for 1,000 times to estimate the test  
483 statistics distribution under null hypothesis.

484

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

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

493

#### 494 **References**

495

- 496 1. Siegel, R, Naishadham, D, Jemal, A. Cancer statistics, 2015. CA: a cancer journal for  
497 clinicians. 2015; 65(1), 5-29.
- 498 2. Shuch B, Amin A, Armstrong AJ, Eble JN, Ficarra V, Lopez-Beltran A, et al.  
499 Understanding pathologic variants of renal cell carcinoma: distilling therapeutic  
500 opportunities from biologic complexity. European urology. 2015;67(1):85-97.
- 501 3. Cancer Genome Atlas Research Network. Comprehensive molecular characterization of  
502 papillary renal-cell carcinoma. N Engl J Med. 2016;2016(374):135-45.
- 503 4. Khurana E, Fu Y, Colonna V, Mu XJ, Kang HM, Lappalainen T, et al. Integrative  
504 annotation of variants from 1092 humans: application to cancer genomics. Science.  
505 2013;342(6154):1235587.
- 506 5. Fu Y, Liu Z, Lou S, Bedford J, Mu XJ, Yip KY, et al. FunSeq2: a framework for  
507 prioritizing noncoding regulatory variants in cancer. Genome biology. 2014;15(10):1.

- 508 6. Huang FW, Hodis E, Xu MJ, Kryukov GV, Chin L, Garraway LA. Highly recurrent  
509 TERT promoter mutations in human melanoma. *Science*. 2013;339(6122):957-9.
- 510 7. Alexandrov LB, Nik-Zainal S, Wedge DC, Campbell PJ, Stratton MR. Deciphering  
511 signatures of mutational processes operative in human cancer. *Cell reports*.  
512 2013;3(1):246-59.
- 513 8. Alexandrov LB, Nik-Zainal S, Siu HC, Leung SY, Stratton MR. A mutational signature  
514 in gastric cancer suggests therapeutic strategies. *Nature communications*. 2015;6.
- 515 9. Hakimi AA, Ostrovnaya I, Jacobsen A, Susztak K, Coleman JA, Russo P, et al.  
516 Validation and genomic interrogation of the MET variant rs11762213 as a predictor of  
517 adverse outcomes in clear cell renal cell carcinoma. *Cancer*. 2016;122(3):402-10.
- 518 10. Rausch T, Zichner T, Schlattl A, Stütz AM, Benes V, Korbel JO. DELLY: structural  
519 variant discovery by integrated paired-end and split-read analysis. *Bioinformatics*.  
520 2012;28(18):i333-9.
- 521 11. Wang D, Yan L, Hu Q, Sucheston LE, Higgins MJ, Ambrosone CB, et al. IMA: an R  
522 package for high-throughput analysis of Illumina's 450K Infinium methylation data.  
523 *Bioinformatics*. 2012;28(5):729-30.
- 524 12. Roberts SA, Lawrence MS, Klimczak LJ, Grimm SA, Fargo D, Stojanov P, et al. An  
525 APOBEC cytidine deaminase mutagenesis pattern is widespread in human cancers.  
526 *Nature genetics*. 2013;45(9):970-6.
- 527 13. Schmidt L, Junker K, Weirich G, Glenn G, Choyke P, Lubensky I, et al. Two North  
528 American families with hereditary papillary renal carcinoma and identical novel  
529 mutations in the MET proto-oncogene. *Cancer research*. 1998;58(8):1719-22.

- 530 14. Schutz FA, Pomerantz MM, Gray KP, Atkins MB, Rosenberg JE, Hirsch MS, et al.  
531 Single nucleotide polymorphisms and risk of recurrence of renal-cell carcinoma: a cohort  
532 study. *The lancet oncology*. 2013;14(1):81-7.
- 533 15. Guo S, Chen W, Luo Y, Ren F, Zhong T, Rong M, et al. Clinical implication of long non-  
534 coding RNA NEAT1 expression in hepatocellular carcinoma patients. *International*  
535 *journal of clinical and experimental pathology*. 2015;8(5):5395.
- 536 16. Nik-Zainal S, Davies H, Staaf J, Ramakrishna M, Glodzik D, Zou X, et al. Landscape of  
537 somatic mutations in 560 breast cancer whole-genome sequences. *Nature*.  
538 2016;534(7605):47-54.
- 539 17. Choudhry H, Albukhari A, Morotti M, Haider S, Moralli D, Smythies J, et al. Tumor  
540 hypoxia induces nuclear paraspeckle formation through HIF-2 $\alpha$  dependent transcriptional  
541 activation of NEAT1 leading to cancer cell survival. *Oncogene*. 2015;34(34):4482-90.
- 542 18. Chakravarty D, Sboner A, Nair SS, Giannopoulou E, Li R, Hennig S, et al. The oestrogen  
543 receptor alpha-regulated lncRNA NEAT1 is a critical modulator of prostate cancer.  
544 *Nature communications*. 2014;5.
- 545 19. Cancer Genome Atlas Research Network. Comprehensive molecular characterization of  
546 urothelial bladder carcinoma. *Nature*. 2014;507(7492):315-22.
- 547 20. Sato Y, Yoshizato T, Shiraishi Y, Maekawa S, Okuno Y, Kamura T, et al. Integrated  
548 molecular analysis of clear-cell renal cell carcinoma. *Nature genetics*. 2013;45(8):860-7.
- 549 21. Supek F, Lehner B. Differential DNA mismatch repair underlies mutation rate variation  
550 across the human genome. *Nature*. 2015;521(7550):81-4.



- 551 22. Gordon MS, Hussey M, Nagle RB, Lara PN, Mack PC, Dutcher J, et al. Phase II study of  
552 erlotinib in patients with locally advanced or metastatic papillary histology renal cell  
553 cancer: SWOG S0317. *Journal of Clinical Oncology*. 2009;27(34):5788-93.
- 554 23. Li Y, Li Y, Chen W, He F, Tan Z, Zheng J, et al. NEAT expression is associated with  
555 tumor recurrence and unfavorable prognosis in colorectal cancer. *Oncotarget*.  
556 2015;6(29):27641.
- 557 24. He C, Jiang B, Ma J, Li Q. Aberrant NEAT1 expression is associated with clinical  
558 outcome in high grade glioma patients. *Apmis*. 2016;124(3):169-74.
- 559 25. Davis CF, Ricketts CJ, Wang M, Yang L, Cherniack AD, Shen H, et al. The somatic  
560 genomic landscape of chromophobe renal cell carcinoma. *Cancer cell*. 2014;26(3):319-  
561 30.
- 562 26. Henderson S, Chakravarthy A, Su X, Boshoff C, Fenton TR. APOBEC-mediated  
563 cytosine deamination links PIK3CA helical domain mutations to human papillomavirus-  
564 driven tumor development. *Cell reports*. 2014;7(6):1833-41.
- 565 27. Le DT, Uram JN, Wang H, Bartlett BR, Kemberling H, Eyring AD, et al. PD-1 blockade  
566 in tumors with mismatch-repair deficiency. *New England Journal of Medicine*.  
567 2015;372(26):2509-20.
- 568 28. Muir P, Li S, Lou S, Wang D, Spakowicz DJ, Salichos L, et al. The real cost of  
569 sequencing: scaling computation to keep pace with data generation. *Genome biology*.  
570 2016;17(1):1.
- 571
- 572
- 573
- 574 29.

575 **Figure 1. Survival analysis of rs11762213 in pRCC patients.**

576 Genotypes are shown in the legend. Peto & Peto modification of the Gehan-Wilcoxon test.

577

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

579 (A) A schematics diagram of non-coding mutations on *MET*. The germline SNP, rs11762213, is also shown. (B) A

580 schematics diagram of non-coding mutations on *ERFF11*. (C) A schematics diagram of non-coding mutations on

581 *NEAT1*. One tumor carries two mutations on *NEAT1*. (D) Tumors with mutations on *NEAT1* show higher *NEAT1*

582 expression. (E) Survival analysis shows mutations in *NEAT1* are associated with worse prognosis. To avoid potential

583 confounding effects, we removed one subject who carries rs11762213 but not *NEAT1* mutation. Log-rank test.

584

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

586 (A) The mutation spectrum of all pRCC WGS samples. Mutations are ordered in alphabetical order of the reference

587 trinucleotides (with the mutated nucleotide in the middle, from A[C>A]A to T[T>G]T) from left to right. (B) We

588 use PCA to maximize inter-sample variation. The loadings on the first principle component is strongly dominated by

589 C>T in CpGs. (C) PC1, along with C>T in CpGs mutation counts and the fractions of such mutations among total

590 mutations are significantly different between two methylation groups. (D) APOBEC mutation signatures are shown

591 for both pRCC (along with three UC sampels, which have blue outer circles) and ccRCC TCGA cohorts. Red

592 dashed line represents the median APOBEC enrichment. (E) Comparison of total mutation counts, mutations counts

593 in open chromatin regions and percentages of mutations in open chromatin regions of total mutations between

594 tumors with chromatin remodeling genes alterations and the ones without.

595

596 **Figure 4. The genomic alteration landscape of 32 whole genome sequenced pRCC samples.**

597 Grey cells represent genomic alterations. CN: copy number. Index: patient index, see Table S2

598