

1 **Pseudogene associated recurrent gene fusion in prostate cancer**

2 Balabhadrapatruni VSK Chakravarthi,¹ Pavithra Dedigama-Arachchige,² Shannon Carskadon,²
3 Shanker Kalyana Sundaram,³ Jia Li⁴, Kuan-Han Hank Wu⁴, Darshan Shimoga Chandrashekar,¹
4 James Peabody,² Hans Stricker,² Clara Hwang,⁷ Dhananjay Chitale,⁵ Sean Williamson,⁵ Nilesh
5 Gupta,⁵ Nora M Navone,⁸ Craig Rogers,² Mani Menon,² Sooryanarayana Varambally,^{1,6}
6 Nallasivam Palanisamy²

7 ¹Department of Pathology, University of Alabama at Birmingham, Birmingham, AL; ²Vattikuti
8 Urology Institute, Department of Urology, Henry Ford Health System, Detroit, MI; ³Michigan
9 Center for Translational Pathology, University of Michigan, Ann Arbor, MI; ⁴Department of
10 Public Health Sciences, Henry Ford Health System, Detroit, MI; ⁵Department of Pathology,
11 Henry Ford Health System, Detroit, MI; ⁶O'Neal Comprehensive Cancer Center, University of
12 Alabama at Birmingham, Birmingham, AL; ⁷Department of Hematology and Oncology, Henry
13 Ford Health System, Detroit, MI; ⁸Department of Genitourinary Medical Oncology-The
14 University of Texas, M D Anderson Cancer Center, TX.

15 **Running title:** Recurrent gene fusion in prostate cancer

16 **Keywords:** Prostate cancer, pseudogenes, gene fusion, non-coding RNA,

17 **Address for Correspondence:** Nallasivam Palanisamy, PhD, Department of Urology, Vattikuti
18 Urology Institute, Henry Ford Health System, Detroit, MI USA; Telephone: 313 874 6396; Fax:
19 313 874 4324; Email: npalani1@hfhs.org

20 **Funding:** National Cancer Institute, Grant number: R21CA176330; CMDRP W81XWH-16-1-
21 0544 to NP

22 **Conflicts of Interest:** None

23 **Word Count:** 4850

24 **Figures: 4**

25 **Tables: 6**

26 **Supplementary Tables and Figures: 1**

27 **ABSTRACT**

28 **Analysis of next generation transcriptome sequencing data of prostate cancer**
29 **identified a novel gene fusion formed by the fusion of a protein coding gene (*KLK4*) with a**
30 **non-coding pseudogene (*KLKP1*) and expression of its cognate protein. Screening of 659**
31 **prostate cancer TMA showed about 32% of positive cases predominantly expressed in**
32 **higher Gleason grade tumors. Concomitant expression with ERG but not with SPINK1 and**
33 **other ETS fusion positive tumors. Fusion gene expression potentially regulated by AR and**
34 **ERG. Antibody specific to the KLK4-KLKP1 fusion protein was validated by**
35 **immunohistochemistry and western blot methods. Oncogenic properties were validated by**
36 ***in vitro* and *in vivo* functional studies. Clinical data analysis shows significant association**
37 **with prostate cancer in young men and overall survival analysis indicate favorable prognosis.**
38 **Non-invasive detection in urine samples has been confirmed. Taken together, we present a**
39 **novel biomarker for routine screening of high Gleason grade prostate cancer at diagnosis.**

40

41 **SIGNIFICANCE**

42 We discovered and validated a novel prostate cancer (PCa) specific fusion gene involving
43 a protein coding (*KLK4*) and a pseudogene (*KLKP1*) and its cognate protein. The unique feature
44 of this fusion gene is the conversion of the noncoding pseudogene into a protein coding gene and
45 its unique expression only in about 30% of high Gleason grade PCa. Expression of this gene is
46 found to be concomitant in ERG fusion positive prostate cancer but mutually exclusive with
47 SPINK1, ETV1, ETV4 and ETV5 positive tumors. Like other ETS family gene fusions, KLK4-
48 KLKP1 can be detected in the urine samples of patients with prostate cancer enabling non-
49 invasive detection of high Gleason grade prostate cancer. Given the unique feature of this fusion

50 oncogenic potential, high Gleason grade specific expression and noninvasive detection, this
51 novel gene fusion has a potential to be used as a biomarker for early detection of high-grade
52 prostate cancer and a therapeutic target.

53

54

55

56

57

58

59

60

61

62

63

64 **INTRODUCTION**

65 Prostate cancer is the most common cancer among men in the United States. Advances in
66 diagnosis, treatment and management has resulted in increased survival rate, yet prostate cancer
67 still remains the second leading cause of cancer-related deaths among American men [1, 2]. One
68 of the major barriers to achieving successful prostate cancer control is the underlying molecular
69 complexity of the disease itself [3]. Morphologically, prostate cancer is well-known to be a
70 diverse disease with patients developing tumors with varying pathological characteristics [4, 5].
71 Many studies have also indicated that prostate cancer is highly heterogeneous with distinct
72 molecular aberrations observed in patient subgroups [6-8]. For example, roughly 50%-60% of

73 prostate cancer patients are known to carry E26 transformation-specific (ETS) family
74 rearrangements, where *ERG*, *ETV1*, *ETV4* or *ETV5* genes fused with androgen regulated 5'
75 partner genes [9]. Additionally, the overexpression of *SPINK1* has been observed in about 5%-
76 10% of prostate cancer patients [10]. Furthermore 1%-2% of the cases are known to carry *RAF*
77 kinase (*BRAF*, *RAF1*) gene fusions [11] while the genetic underpinnings in the remaining 30%-
78 40% of the prostate cancer cases are not known [6]. Importantly, distinct molecular changes have
79 been linked with unique disease outcomes [10, 12, 13], indicating complex heterogeneity among
80 patients with respect to disease progression. Therefore, discovery of new molecular markers for
81 further patient stratification and to classify indolent and aggressive prostate cancer is an urgent
82 unmet clinical need to facilitate targeted therapy and effective prostate cancer management.

83 Currently, prostate cancer diagnosis is primarily based on prostate-specific antigen (PSA)
84 levels and Gleason grade, a scoring system based on the morphology of the prostate tissue [14].
85 Following the detection of elevated PSA or pro-PSA levels, prostate cancer is identified by the
86 presence of Gleason graded cancer on needle biopsies. The decision to pursue immediate
87 treatment or continue active surveillance is mainly determined using the Gleason grade.
88 However, the rise in PSA is not prostate cancer specific and is multifactorial [15]. Therefore,
89 PSA has been an inadequate diagnostic marker, in some cases leading to overdiagnosis and
90 unnecessary treatment. Though high Gleason grade tumors are known to be clinically aggressive,
91 whether low Gleason grade tumors require treatment has been debated [16]. While intervention
92 in low Gleason grade cancers may result in overtreatment, watchful waiting may also pose an
93 unnecessary risk and additional burden of repeat biopsies. Given these limitations of the existing
94 markers and the recognition of prostate cancer as a heterogeneous disease, molecular markers

95 specific to distinct patient subgroups, are required as alternatives for both initial cancer diagnosis
96 and distinguishing aggressive cancer from indolent disease.

97 Although several recurrent molecular alterations have been identified in a subset of
98 prostate cancer cases, the genetic aberrations in prostate cancer patients negative for all the
99 known molecular makers remain to be studied. Moreover, most prostate cancer molecular studies
100 have been carried out on Caucasian American patients with little representation from the African
101 American (AA) population [17]. Given the unique ancestral background of AAs and the
102 aggressive nature of prostate cancer, the genetic underpinnings to understand the racial disparity
103 in the incidence of prostate cancer markers is not well studied. Therefore, the study of additional
104 molecular aberrations using large cohorts of a racially diverse population is a pressing need in
105 prostate cancer research. In addition to identifying subtype specific prostate cancer diagnostic
106 and prognostic markers, such studies may also facilitate the development of novel therapeutic
107 approaches by uncovering molecular alterations, which may be pharmacologically targeted in
108 distinct patient subgroups.

109 Given the need for identifying novel molecular markers in prostate cancer patients, we
110 investigated the expression patterns of pseudogenes in 89 prostate cancer patient samples using a
111 paired-end next generation sequencing approach [18]. Often considered as dysfunctional
112 relatives of known protein-coding genes, pseudogenes have recently been implicated in cancer
113 with roles in gene regulation [19]. While we observed distinct expression changes in several
114 pseudogenes in prostate cancer compared to normal prostate tissue, we also noted the rare
115 occurrence of a chimeric transcript formed through the fusion of the androgen regulated gene
116 *KLK4* (Kallikrein Related Peptidase 4) with the adjacent pseudogene *KLKP1* (Kallikrein
117 Pseudogene 1). Importantly, the fusion converts the *KLKP1* pseudogene to a protein-coding gene

118 with a predicted chimeric protein of 164 amino acids, of which 55 amino acids are derived from
119 the pseudogene part due to a shift in the open reading frame of the fusion formed by trans-
120 splicing mechanism rather than chromosomal rearrangement [18]. Although a few pseudogenes
121 have been previously reported to be expressed as proteins [20, 21], *KLK4-KLKPI* is a rare
122 example where gene fusion leading to conversion of a non-coding pseudogene to a protein-
123 coding gene. Further studies showed that *KLK4-KLKPI* fusion is both prostate tissue and cancer
124 specific, suggesting a possible role in prostate cancer formation [18]. Both the prostate cancer
125 specific expression and the intriguing nature of the *KLK4-KLKPI* fusion warrant further
126 functional studies to understand the role of *KLK4-KLKPI* in prostate cancer development.
127 Therefore, in this study, we explored the prevalence, the expression pattern, noninvasive
128 detection, and the oncogenic properties of *KLK4-KLKPI* to investigate the potential of *KLK4-*
129 *KLKPI* fusion gene as a novel molecular marker in prostate cancer.

130

131 **RESULTS**

132 Both *KLK4* and *KLKPI* belong to the kallikrein family of serine proteases, a cluster of
133 genes located on chromosome 19(q13.33-q13.41). The gene cluster contains 15 members
134 including *KLK3*, which is commonly known as PSA [22]. The *KLK4-KLKPI* fusion is formed by
135 a trans-splicing mechanism or an in-frame fusion due to a microdeletion of the region between
136 the adjacent genes, *KLK4* and *KLKPI*, leading to the fusion of the first 2 exons of *KLK4* with
137 exon 4 and 5 of *KLKPI* (**Fig. 1A and Supplementary Fig. S1**, GenBank ID 2227664). The
138 resulting chimeric sequence predicts a 164 amino acid protein, of which 55 amino acids are
139 derived from *KLKPI* (**Fig. 1B**). According to data on the GTEx portal, full length *KLKPI* is
140 exclusively expressed in normal prostate tissue (**Supplementary Fig. S2**). In contrast,

141 quantitative PCR (qRT-PCR) analysis of prostate cancer samples, prostate cell lines, benign
142 prostate tissues and other solid cancers revealed that *KLK4-KLKPI* fusion transcript is prostate
143 cancer specific and expressed in a subset of cases [18]. However, the study included only a
144 limited number of prostate cancer samples (n = 36) and the occurrence of *KLK4-KLKPI* in a
145 large, racially inclusive cohort must be explored to determine the prevalence of *KLK4-KLKPI* in
146 the prostate cancer patient population. Therefore, we studied the expression of *KLK4-KLKPI* on
147 a larger patient cohort using fusion transcript specific anti-sense oligonucleotide probe by RNA
148 *in situ* hybridization (RNA-ISH). Specifically, we constructed tissue microarrays (TMAs) using
149 prostate cancer tissues obtained from 659 radical prostatectomy (RP) specimens at the Henry
150 Ford Health Systems. The cohort was racially inclusive with 380 Caucasian, 250 AA and 29
151 patients belonging to other racial groups. Each TMA contained 3 cores obtained from different
152 regions of the RP prostate from each patient (**Supplementary Fig. S3**). The individual tissue
153 cores in each patient were reviewed and the highest tumor grade observed was assigned to each
154 case. Thus the TMAs included 612 patient cases with all cores carrying prostate cancer (Gleason
155 grade group 1 [3 + 3 = 6] - 110, Gleason grade group 2 [3 + 4 = 7] - 247, Gleason grade group 3
156 [4 + 3 = 7] - 119, Gleason grade group 4 [4 + 4 = 8] - 94, and Gleason grade group 5 [4 + 5 = 9;
157 5 + 4 = 9 and 5 + 5 = 10] - 42). The rest of the cases consisted of 23 cases with benign, 21 cases
158 with high grade prostate intraepithelial neoplasia, 2 cases with stroma and 1 case with atypical
159 cores. RNA-ISH was carried out using an antisense RNA probe specific to the *KLK4-KLKPI*
160 fusion. The TMA slides were then reviewed for the intensity of the RNA-ISH signal. A score of
161 expression ranging from 0 to 4+ was given according to the intensity of the RNA-ISH signal
162 where, 0 indicated no detectable RNA-ISH signal, while 4+ was assigned to the highest level of
163 RNA-ISH signal [23].

164 Of the 659 cases in the cohort, 209 (32%) were positive for *KLK4-KLKPI* fusion,
165 indicating the recurrent nature of *KLK4-KLKPI* among prostate cancer patients. Most of the
166 *KLK4-KLKPI* positive cases showed RNA-ISH signal intensity of 1+ (130 cases; **Fig. 1C**) while
167 more intense RNA-ISH signal 2+ was observed in 66 cases, 3+ in 12 cases and 4+ in 1 case (**Fig.**
168 **1C**) and remaining cases were “0” or negative, suggesting varying expression levels among
169 patients. To further confirm that *KLK4-KLKPI* is specific to prostate cancer, we then explored
170 the association of *KLK4-KLKPI* RNA-ISH signal with Gleason grade by using Pearson’s chi-
171 square test. The results showed that *KLK4-KLKPI* is exclusively expressed in prostate cancer
172 tissues compared to benign, high grade prostate intraepithelial neoplasia and atypical prostate
173 tissues (**Fig. 1D, Fig. 1E, and Table 1**), confirming that *KLK4-KLKPI* expression is prostate
174 cancer specific. Additionally, we also analyzed if *KLK4-KLKPI* expression is associated with
175 Gleason grade group and found no associations with distinct Gleason grade groups (**Table 2**).

176 Next, we investigated if *KLK4-KLKPI* fusion displays racial disparity in the incidence.
177 The 209 positive cases included 128 Caucasian Americans (34%), 69 AAs (28%) and 12 patients
178 from other races (41.4%). Pearson’s chi-square test analysis revealed that prevalence of *KLK4-*
179 *KLKPI* is high in Caucasian American compared with AA patients but not statistically
180 significant (**Table 3 and Supplementary Fig. S4**), demonstrating no racial bias in the incidence.
181 We also explored if *KLK4-KLKPI* expression is related with patient age. We categorized the
182 patients into 2 groups as young (age ranging from 40 to 50 years) and old (age ranging from 51
183 to 83 years). Pearson’s chi-square test showed significantly higher expression of *KLK4-KLKPI*
184 in young age group compared to the old age group (**Fig. 1F and Table 4**).

185 The other common prostate cancer specific mutations such as ETS gene fusions and
186 *SPINK1* overexpression are known to occur in a mutually exclusive manner. Therefore, we also

187 analyzed the association of *KLK4-KLKPI* fusion expression with ETS gene fusions and *SPINK1*
188 expression. We screened the same set of TMAs by using dual immunohistochemistry (IHC) for
189 *ERG* and *SPINK1* and dual RNA-ISH for *ETV1*, *ETV4* and *ETV5*. By using Pearson's chi-square
190 test, we observed that *KLK4-KLKPI* expression is associated with *ERG*⁺ cases (**Fig. 1G**,
191 **Supplementary Fig. S5, and Table 1**). However, no such association was observed with
192 *SPINK1*, *ETV1*, *ETV4* and *ETV5* (**Fig. 1G and Table 5**), suggesting concurrent expression of
193 *KLK4-KLKPI* with distinct ETS gene fusion positive cases. Next, we investigated if *KLK4-*
194 *KLKPI* is related with *PTEN* loss, another common prostate cancer mutation that is associated
195 *ERG*⁺ and aggressive disease [24-26]. We carried out IHC for *PTEN* on the same set of TMAs
196 and found that *PTEN* deletion was significantly lower in *KLK4-KLKPI* positive cases compared
197 to *KLK4-KLKPI* negative cases (**Fig. 1G and Table 1**). Given that *ERG* is known to co-occur
198 with *PTEN* loss [27], we further analyzed if there is any significant difference in *PTEN* loss in
199 cases showing both *ERG* fusion and *KLK4-KLKPI* compared to the rest of the cases and found
200 no significant difference in *PTEN* status in cases with *ERG* fusion and *KLK4-KLKPI* expression,
201 suggesting that *KLK4-KLKPI* may represent a distinct subtype of prostate cancer.

202 Having thus confirmed the recurrent and the prostate cancer specific occurrence of
203 *KLK4-KLKPI* fusion, we then studied the expression of *KLK4-KLKPI* fusion protein. Based on
204 the sequence, the *KLK4-KLKPI* fusion gene is predicted to generate a full-length protein of 164
205 amino acids of which 55 are derived from the *KLKPI* pseudogene (**Fig. 1B**). To validate the
206 *KLK4-KLKPI* expression as a full-length protein, we generated adenoviral constructs carrying
207 the N-FLAG-tagged *KLK4-KLKPI* fusion gene and transfected HEK293 cells. To stabilize the
208 protein levels of *KLK4-KLKPI*, the cells were treated with the proteasome inhibitor bortezomib.
209 As a control, bortezomib treated cells transfected with vector DNA alone were used. Expression

210 of the fusion transcript was confirmed by qRT-PCR using fusion specific primers (**Fig. 2A**). Cell
211 lysates were analyzed by western blotting using an anti-N-FLAG antibody. Importantly, we
212 observed a FLAG-specific protein band around 17kDA (**Fig. 2B**), confirming the expression of
213 *KLK4-KLKPI* as a full-length protein. For additional validation, we also checked the expression
214 of N-FLAG-tagged *KLK4-KLKPI* using the anti-FLAG antibody in the normal prostate cell line
215 RWPE-1, transfected with and without N-FLAG-tagged *KLK4-KLKPI* adenovirus construct.
216 Notably, we detected anti-FLAG specific protein band only in the transfected RWPE1 cells
217 (**Supplementary Fig. S6**). Furthermore, we also developed a *KLK4-KLKPI* specific polyclonal
218 antibody (Eurogentech, Seraing, Belgium) using the antigenic peptide “CTISATSSARTS” (**Fig.**
219 **1B**) derived from the *KLKPI* pseudogene region of the fusion protein. After cell lysis and SDS-
220 PAGE, we probed HEK293 lysates transfected with and without N-FLAG-tagged *KLK4-KLKPI*
221 adenovirus construct with the *KLK4-KLKPI* specific antibody using western blot. A protein band
222 around 17kDA was observed further confirming the expression of the chimeric *KLK4-KLKPI*
223 protein (**Fig. 2B, 2E**) and the specificity of the antibody to the fusion protein.

224 In order to assess the expression of *KLK4-KLKPI* in metastatic prostate cancer, we then
225 analyzed the expression of *KLK4-KLKPI* in prostate cancer patient derived xenografts
226 (PDX)[28]. We first screened the expression of *KLK4-KLKPI* using qRT-PCR and identified 17
227 out of 31 PDX models positive for endogenous expression of *KLK4-KLKPI* (**Supplementary**
228 **Fig. S7; Table 6**). Then, we selected one of the PDX tissues (MDA PCa 153-7) expressing high
229 levels of *KLK4-KLKPI* and one with no detectable levels of *KLK4-KLKPI* (MDA PCA 144-13).
230 After protein isolation and separation on SDS-PAGE, the lysates were probed with the *KLK4-*
231 *KLKPI* specific antibody using western blot. Importantly, we observed a protein band around
232 17kDA only in the *KLK4-KLKPI* positive PDX (**Fig. 2C, 2F**), indicating the endogenous

233 expression of *KLK4-KLKPI* fusion protein in metastatic prostate cancer patients. Additionally,
234 we also screened the expression of *KLK4-KLKPI* in xenograft tissues using IHC with the *KLK4-*
235 *KLKPI* specific antibody. While *KLK4-KLKPI* expression was observed in qRT-PCR positive
236 PDX tissues, minimal or no *KLK4-KLKPI* IHC signal was seen in qRT-PCR negative xenografts
237 (**Fig. 2D**), further suggesting the presence of *KLK4-KLKPI* protein in a subset of prostate cancer
238 patients. Comparison of IHC results with RNA-ISH positive and negative tissue showed
239 specificity of the antibody to *KLK4-KLKPI* RNA-ISH positive tissue only (**Fig.2G**).

240 Given the exclusive expression of *KLK4-KLKPI* in prostate cancer, next we explored the
241 functions of *KLK4-KLKPI* by studying the oncogenic properties of the fusion gene. Specifically,
242 we established RWPE-1 cells with stable expression of *KLK4-KLKPI* by transfection with
243 lentiviral constructs carrying FLAG-tagged *KLK4-KLKPI*. As controls, cells stably transfected
244 with a LACZ control (LACZ) and un-transfected RWPE-1 cells were used. We first confirmed
245 the expression of *KLK4-KLKPI* by qRT-PCR. The results showed significant expression of
246 *KLK4-KLKPI* in transfected cells compared to both the un-transfected cells and the LACZ
247 control (**Fig. 3A**). Then we investigated the effect of *KLK4-KLKPI* on cell proliferation by
248 measuring the number of cells using a Coulter particle counter. Compared to the un-transfected
249 cells and the LACZ control, a notable increase in the cell number was seen over time in *KLK4-*
250 *KLKPI* transfected cells (**Fig. 3B**), indicating a role of *KLK4-KLKPI* on cell proliferation. Next,
251 we studied the effect of *KLK4-KLKPI* in cell invasion using the Matrigel invasion assay.
252 Importantly, a significant increase in the number of invaded cells was observed with *KLK4-*
253 *KLKPI* transfected cells compared to both the un-transfected and the LACZ control (**Fig. 3C**).
254 For additional validation, we also transiently transfected PrEC, another normal prostate cell line
255 with *KLK4-KLKPI*. As controls, un-transfected cells and cells transfected with a LACZ control

256 were used. Additionally, we also used cells transfected with *EZH2*, which has been shown to
257 increase invasion of prostate cancer and other cancer cells [29, 30] as a positive control. The
258 invasion of cells was then examined by the Matrigel invasion assay. Like RWPE-1, PrEC cells
259 also showed a significant increase in the number of invaded cells compared to both the un-
260 transfected and the LACZ control (**Fig. 3D**). As expected, cells transfected with *EZH2* also
261 demonstrated increased invasion compared to the LACZ control and the un-transfected cells
262 (**Fig. 3D**). In all, our studies indicate that *KLK4-KLKPI* promote both cell proliferation and
263 invasion of prostate cells, suggesting an oncogenic role for *KLK4-KLKPI* fusion.

264 In order to further understand the oncogenic properties of *KLK4-KLKPI*, we also studied
265 the effects of *KLK4-KLKPI* fusion on intravasation and tumor formation using the chicken
266 chorioallantoic membrane (CAM) *in vivo* assay [31, 32]. We implanted eggs with RWPE-1 cells
267 stably expressing *KLK4-KLKPI* and then checked for the presence of intravasated cells in the
268 lower CAM by using quantitative human Alu-specific PCR. As controls, eggs implanted with
269 either un-transfected cells or cells stably transfected with a LACZ control were used. Notably,
270 we observed a marked intravasation by *KLK4-KLKPI* transfected cells in the lower CAM
271 compared to both un-transfected cells and LACZ control (**Fig. 3E**). Additionally, we also
272 isolated and weighed the extraembryonic tumors from eggs implanted with either *KLK4-KLKPI*
273 transfected cells or controls. The tumors isolated from eggs implanted with cells expressing
274 *KLK4-KLKPI* showed significantly higher weight than the tumors isolated from eggs treated
275 with the un-transfected cells and the LACZ control (**Fig. 3F**). Overall, the results establish that
276 *KLK4-KLKPI* drives intravasation and tumor formation in prostate cells, indicating a potential
277 role in prostate cancer development.

278 Further, we investigated the molecular mechanisms underlying the oncogenic functions
279 of *KLK4-KLKPI* fusion. We conducted a gene expression microarray analysis using RWPE-1
280 cells stably transfected with *KLK4-KLKPI*. As the control, cells transfected with LACZ control
281 were used. After RNA isolation, and microarray analysis, we observed a significant number of
282 genes expressed differently between the RWPE-1 cells transfected with *KLK4-KLKPI* and the
283 LACZ control. We selected the genes showing a fold change value of more than 1 in 2
284 independent replicates and generated a heat map with the top 100 genes differentially expressed
285 (**Fig. 4A**). We noted genes both upregulated and downregulated in cells expressing the *KLK4-*
286 *KLKPI* fusion, suggesting a possible function for *KLK4-KLKPI* in gene expression regulation.
287 Further, we also carried out a gene set enrichment analysis [33] to explore any overlap between
288 the differentially expressed genes observed with *KLK4-KLKPI* transfection and other curated
289 gene sets. Importantly, we noted enrichment of 2 curated gene sets, one involving genes
290 upregulated in endometroid endometrial metastatic tumor and the other containing genes
291 overexpressed in melanoma metastatic cancer (**Fig. 4B**), indicating that the genes affected by
292 *KLK4-KLKPI* are associated with metastatic cancer. As a further step, we also carried out a
293 KEGG pathway analysis using the DAVID tool [34]. The genes differentially affected by *KLK4-*
294 *KLKPI* were shown to be associated with several cancer-related pathways (**Fig. 4C**), further
295 implying that *KLK4-KLKPI* may regulate the expression of genes involved in cancer and
296 metastasis.

297 Given the well-established role of androgen receptor (*AR*) in gene expression in prostate
298 cancer [35], we also explored if *AR* is driving the expression of *KLK4-KLKPI* in prostate cancer.
299 Additionally, since we observed concurrent expression of *ERG* with *KLK4-KLKPI* (**Fig. 1G**),
300 we also studied if *ERG* is involved in the expression of *KLK4-KLKPI*. Therefore, to identify any

301 AR or ERG binding sequences on *KLK4* or *KLKPI*, we examined data from a previous study
302 where a chromatin immunoprecipitation assay was carried out using antibodies specific to *AR*
303 and *ERG* [36]. Notably, we observed both *AR* and *ERG* binding sites at the fusion junction of
304 *KLKPI* (**Supplementary Fig. S8**), suggesting that both *AR* and *ERG* may modulate the
305 expression of *KLK4-KLKPI* during prostate cancer formation.

306 For further characterization of the functional role of *KLK4-KLKPI*, we also studied the
307 cellular localization of *KLK4-KLKPI*. We carried out immunofluorescence studies of RWPE-1
308 cells transfected with adeno-FLAG tagged-*KLK4-KLKPI* using fluorescent anti-FLAG antibody.
309 As a control, cells transfected with adeno-LacZ were used. While cells transfected with adeno-
310 Lacz showed minimal immunofluorescence as expected, notably, we observed colocalization of
311 *KLK4-KLKPI* immunofluorescence signal with 4, 6-diamidino-2-phenylindole (**Supplementary**
312 **Fig. S9**), indicating that *KLK4-KLKPI* is localized in the nucleus of the cells.

313 The prostate cancer exclusive expression of *KLK4-KLKPI* in a considerable subset of
314 patients indicates the possible use of *KLK4-KLKPI* as a biomarker for prostate cancer.
315 Therefore, to further explore the potential utility of *KLK4-KLKPI* as a prostate cancer marker,
316 we investigated the association between *KLK4-KLKPI* expression and preoperative PSA of the
317 659 patients in our cohort. Specifically, we performed a t-test to evaluate difference in log-
318 transformed preoperative PSA between cases with and without *KLK4-KLKPI* expression.
319 Interestingly, patients with *KLK4-KLKPI* expression showed slightly lower preoperative PSA
320 values compared to patients without *KLK4-KLKPI* expression (**Supplementary Fig. S10**). As a
321 further step, we also analyzed the association between *KLK4-KLKPI* and the time to biochemical
322 recurrence, using multivariable Cox regression model. Patients with *KLK4-KLKPI* showed a
323 lower risk of biochemical recurrence (HR = 0.58; **Supplementary Fig. S11**) after adjusting for

324 age, Gleason grade, and tumor stage. However, the difference in recurrence was not statistically
325 significant ($p = 0.12$), possibly due to small power as the number of patients showing recurrence
326 was small ($n = 49$). Additionally, we also analyzed the association of *KLK4-KLKPI* with other
327 clinical and pathological parameters such as family history, tumor stage, tumor volume,
328 metastasis to lymph nodes, perineural invasion and presence of lymph vascular invasion using
329 Pearson's chi-square test. No statistically significant association was observed between *KLK4-*
330 *KLKPI* and the clinicopathological variables. Lastly, like TMPRSS2-ERG gene fusions in
331 prostate cancer, we explored the feasibility of detecting *KLK4-KLKPI* in urine samples of
332 prostate cancer patients for noninvasive detection of this marker. We collected urine samples
333 from 90 unselected prostate cancer patients. All patients had confirmed prostate cancer, with
334 most having metastatic or biochemically recurrent disease. Then we screened for *KLK4-KLKPI*
335 transcript using qRT-PCR. As a positive control, RWPE-1 cells stably expressing *KLK4-KLKPI*
336 was used. Importantly, *KLK4-KLKPI* expression was detected in 15 out of 90 (17%) patient
337 samples (**Supplementary Fig. S12**), suggesting the potential for noninvasive detection in patient
338 urine samples. Overall, our study establishes *KLK4-KLKPI* as a recurrent chimeric transcript
339 exclusively expressed in prostate cancer tissues with implications on disease progression and
340 feasibility of being noninvasively detected in patient urine samples.

341

342 **DISCUSSION**

343 Given the complex heterogeneous nature of prostate cancer, the identification of distinct
344 patient subgroups based on molecular markers is a necessary step towards targeted disease
345 management. Therefore, in this study we further explored and characterized a pseudogene
346 associated gene fusion *KLK4-KLKPI*. We established that *KLK4-KLKPI* is a recurrent, prostate

347 cancer exclusive fusion transcript that occurs at a significant incidence rate (32%) among
348 prostate cancer patients. Similar to other distinct molecular aberrations such as ETS
349 rearrangements [9] and *SPINK1* mutation [10], *KLK4-KLKPI* was observed only in a subset of
350 prostate cancer patients. However, unlike the mutually exclusive pattern of expression of ETS
351 rearrangements and *SPINK1*, *KLK4-KLKPI* showed concomitant expression with *ERG*,
352 indicating possible cross-talk with ERG. Notably, *KLK4-KLKPI* expression was associated with
353 intact *PTEN* status, suggesting these fusion positive tumors are distinct molecular subtypes from
354 ERG+/PTEN- tumors. Interestingly, full-length normal *KLKPI* transcript showed normal
355 prostate specific expression (GTEx portal) and not in prostate cancer. Furthermore, despite
356 *KLKPI* being categorized as a pseudogene, we showed that *KLK4-KLKPI* is expressed as a full-
357 length protein in a rare phenomenon where gene fusion leads to the inclusion of a pseudogene
358 segment in an expressed protein. Importantly, *KLK4-KLKPI* promoted proliferation, invasion,
359 intravasation and tumor formation, suggesting functional implications on prostate cancer
360 development. Moreover, gene expression studies revealed considerable transcriptional changes
361 in cancer-related genes in cells transfected with *KLK4-KLKPI*, which may indicate that *KLK4-*
362 *KLKPI* may play a role in transcription during prostate cancer formation. In agreement with a
363 role in transcriptional regulation, *KLK4-KLKPI* was also seen to be localized in the nucleus.
364 Furthermore, both *ERG* and *AR* were found to have binding sites on *KLKPI*, indicating that
365 *KLK4-KLKPI* expression may be *ERG* and *AR* modulated. Finally, we showed that *KLK4-*
366 *KLKPI* can be easily detected in patient urine samples, suggesting the feasibility for possible
367 future use as a biomarker for early detection of high Gleason grade prostate cancer. Altogether,
368 our study establishes *KLK4-KLKPI* as a novel player in a subset of prostate cancer cases with
369 likely roles in tumor formation.

370 Long thought to be junk or nonfunctional units of the human genome, pseudogenes have
371 been recently acknowledged to have key cellular roles, particularly in diseases such as cancer
372 [37]. While some pseudogenes are known to be transcribed into non-coding RNA [37], a few
373 pseudogenes have been shown to be even express proteins [20]. Studies have revealed that
374 several different variants of *KLKP1* pseudogene are transcribed exclusively in prostate tissues
375 (**Supplementary Fig. S1**) in an androgen regulated manner [21, 38]. Of the different variants, at
376 least one *KLKP1* variant has been shown to be expressed as a protein in a transfected cell,
377 although not *in vivo* [21]. Even though the variant chimeric transcripts of *KLK4-KLKP1* has
378 been previously described [39, 40], it has not been reported to be expressed as a protein and the
379 functional characteristics have not been validated. Importantly, we verified that *KLK4-KLKP1* is
380 expressed as a full-length protein in both transfected cells and endogenously in castration
381 resistant prostate cancer (PDX), suggesting the occurrence in prostate cancer tissues. In contrast
382 to *KLK4*, which is overexpressed in prostate cancer with roles in cell proliferation, migration and
383 cancer metastasis [41-43], all *KLKP1* variants are known to be expressed more in normal
384 prostate tissues compared to prostate cancer [21, 38]. However, *KLK4-KLKP1* is exclusively
385 expressed in prostate cancer with co-occurrence with ERG+ tumors. Thus, our results indicate
386 novel complexity in the *KLK4* and *KLKP1* locus and hint at differential expression of the loci in
387 prostate cancer cells compared to normal prostate cells. Given the presence of *AR* and *ERG*
388 binding sites on *KLKP1* and the previous reports demonstrating *AR* regulation of *KLKP1*
389 expression [21, 38], it is likely that prostate cancer specific expression of *KLK4-KLKP1* is
390 modulated by *AR* and *ERG*. Furthermore, additional variants of *KLK4-KLKP1*, which are
391 different from the *KLK4-KLKP1* transcript observed in prostate cancer, have also been reported
392 in renal cell cancer [40]. While the alternative *KLK4-KLKP1* transcripts were found to occur in a

393 considerable subset of renal cell cancer cases (27%), none of the variants were shown to be
394 expressed as proteins. Thus *KLK4-KLKPI* may be spliced and expressed differently in a tissue
395 specific manner in distinct cancers. Taken together, our results suggest that *KLK4* and *KLKPI*
396 may be a diverse locus that undergo differential splicing and transcription with functional
397 implications in cancer. Consequently, our work highlights unprecedented roles of pseudogenes
398 and complex molecular events involved in cancer.

399 In agreement with previous reports indicating significant molecular heterogeneity among
400 prostate cancer cases [7], *KLK4-KLKPI* was expressed only in a subset of prostate cancer
401 patients (32%). Additionally, *KLK4-KLKPI* expression was significantly higher in younger
402 patients compared to older prostate cancer patients. Given the oncogenic properties and the
403 transcriptional changes observed with *KLK4-KLKPI*, our results suggest that distinct molecular
404 changes may dictate unique prostate cancer clinical outcomes among patients. Thus, our study
405 further emphasizes the need for subtype specific molecular markers in prostate cancer control.
406 In addition to enhancing cell proliferation, invasion and tumor formation, *KLK4-KLKPI* also
407 caused marked changes in gene expression. Notably, genes affected by *KLK4-KLKPI* were
408 cancer-related and were involved in metastasis of other cancers, implicating a functional role for
409 *KLK4-KLKPI* in prostate cancer. Additionally, ERG was found to have a binding site on *KLK4-*
410 *KLKPI*. Given that ERG expression was associated with *KLK4-KLKPI*, ERG may bind to the
411 *KLKPI* locus and may promote the expression of *KLK4-KLKPI* in a subset of prostate cancer
412 patients.

413 Even though *KLK4-KLKPI* was implicated in metastatic prostate cancer, the association
414 of *KLK4-KLKPI* with intact *PTEN* status and lower preoperative PSA values suggests indolent
415 disease in prostate cancer patients with *KLK4-KLKPI* expression. However, larger studies

416 exploring the association between *KLK4-KLKPI* expression and prostate cancer clinical
417 outcomes are necessary to establish *KLK4-KLKPI* as a biomarker for prostate cancer.
418 Furthermore, detailed studies are also necessary to fully understand the molecular mechanisms
419 through which *KLK4-KLKPI* promotes prostate cancer formation. Consequently, such studies
420 will explore the potential of *KLK4-KLKPI* as a biomarker and a therapeutic target in prostate
421 cancer, eventually making significant contributions towards achieving effective prostate cancer
422 control.

423

424

425

426

427 **MATERIALS AND METHODS**

428 **Tissue microarray construction**

429 Prostatectomy samples collected from 659 patients who underwent radical prostatectomy
430 at Henry Ford Health Systems (HFHS), were reviewed and tissue cores from different regions of
431 the tumor were isolated to construct paraffin embedded tissue microarrays. In most cases, a total
432 of three tissue cores were obtained from each prostatectomy sample. In all cases, appropriate
433 informed consent and Institutional Review Board approval were obtained. The Gleason grade of
434 each tissue core and the race of the patients were reviewed by the study pathologists (NG and SW).
435 Clinical and pathological information of patients such as age, race, family history of prostate
436 cancer, pre-operative PSA, prostatectomy date, Gleason Grade group, tumor stage, cancer status
437 of the lymph nodes, tumor volume, perineural invasion, presence of lymph vascular invasion, last

438 PSA, last PSA date, presence of biochemical recurrence, date of biochemical recurrence were also
439 recorded.

440 **KLK4-KLKP1 RNA *in situ* Hybridization (RNA-ISH)**

441 RNA-ISH was performed as described previously using RNAscope 2.5 HD Reagent Kit
442 (ACDBio, catalog #322350) according to the manufacturer's instructions (1). Briefly, after baking,
443 deparaffinization, and target retrieval per manufacturer's instructions, TMA slides were incubated
444 with target probes for KLK4-KLKP1 (ACDBio, catalog #405501, NM_001136154, region 2933–
445 3913) for 2 hours at 40°C in a humidity chamber. After detection and color development, slides
446 were washed twice in deionized water and then counterstained in hematoxylin (Agilent DAKO,
447 catalog #K800821-2) for 5 minutes. Slides were washed several times in tap water, then dried,
448 dipped in xylene, and mounted in EcoMount (Fisher, catalog #50–828-32). Next the slides were
449 scanned using a digital imaging system (Aperio Scanner, Leica). The images were reviewed and
450 the RNA-ISH signal on the TMAs was scored. A staining pattern of distinct punctuate cytoplasmic
451 dots was considered as a positive RNA-ISH signal for KLK4-KLKP1 expression. Depending on
452 the intensity of the RNA-ISH staining, a score ranging from +1 to +4 was given to tissue cores
453 with positive RNA-ISH signal, with +1 assigned to the weakest RNA-ISH staining, and +4 given
454 to the cores showing the most intense RNA-ISH staining. A score of 0 was assigned to tissue cores
455 with no visible RNA-ISH staining. The highest score observed among the tissue cores was then
456 assigned to each patient case. If all tissue cores of a patient was 0, the case was recorded as
457 negative.

458 **Cell culture**

459 HEK-293 cells and prostate benign epithelial cells (RWPE-1, #CRL-11609) were
460 purchased from American Type Culture Collection (Manassas, VA). Primary prostate epithelial

461 cells (PrEC) were purchased from Lonza (Walkersville, MD). HEK-293 cells were cultured in
462 MEM media (Thermo Fisher Scientific, catalog #11095080,) supplemented with 10% FBS (fetal
463 bovine serum, Thermo Fisher Scientific, catalog number #10082147). RWPE-1 cells were cultured
464 in Keratinocyte serum free medium (K-SFM, Gibco™, Thermo Fisher Scientific, catalog #17005-
465 042, Carlsbad, CA) supplemented with Bovine Pituitary Extract (BPE, 0.05 mg/ml, Thermo Fisher
466 Scientific, catalog #17005-042), human recombinant Epidermal Growth Factor 1-53 (EGF 1-53,
467 5 ng/ml, Thermo Fisher Scientific, catalog #17005-042) and 1% penicillin/streptomycin. PrEC
468 cells were cultured in Prostate Epithelial Cell Basal Medium (PrEGM) supplemented with Prostate
469 Epithelial Cell Growth Kit (Clonetics™ PrEGM™, BulletKit™, Lonza). All cell cultures were
470 maintained at 37°C in an incubator with a controlled humidified atmosphere composed of 95% air
471 and 5% CO₂.

472 ***In vitro* overexpression of KLK4-KLKP1**

473 KLK4-KLKP1 cDNA was PCR amplified using a forward primer with DDK tag and a
474 reverse primer from KLK4-KLKP1 template and was cloned into Gateway expression system
475 (Life Technologies). To generate lentiviral and adenoviral constructs, PCR8-KLK4-KLKP1 (DDK
476 tagged) was recombined with pLenti6/V5-Dest™ (Life Technologies) or pAD/CMV/V5-Dest™
477 (Life Technologies), respectively using LR Clonase II (Life Technologies). For transient
478 overexpression in HEK-293, RWPE-1 and PrEC cells, adenoviruses carrying KLK4-KLKP1,
479 EZH2 or lacZ were added to the culture media after cells reached 50-70% confluency. At the same
480 time, cells were treated with or without bortezomib (100nM in ethanol, 10 µL, Cayman Chemical,
481 catalog #10008822). After incubation for 48 hours at 37°C, cells were harvested by scraping. For
482 stable overexpression, RWPE-1 cells were infected with lentiviruses expressing KLK4-KLKP1 or
483 lacZ, and stable clones were selected with blasticidin (3.5 µg/ml, Sigma-aldrich, MO, USA). Lenti

484 and adeno viruses were generated by the University of Michigan Vector Core (Ann Arbor, MI,
485 USA).

486 **Western blotting**

487 Harvested cells were spun down (1000 rpm, 5 min, 4 °C). For HEK-293 cells, the cell pellet
488 was re-suspended in RIPA lysis buffer (Thermo Fisher Scientific, catalog #PI89900) supplemented
489 with protease inhibitor (1X, genDEPOT, catalog #50-101-5488). For RWPE-1 cells, NP-40 lysis
490 buffer (Boston BioProducts, Ashland, MA) with protease inhibitor was used to lyse the cells. With
491 xenograft tissues, frozen tissues were cut into small pieces and then sonicated on ice in RIPA lysis
492 buffer. The debris from cells or tissues were removed by centrifugation (13.2 rpm, 10 minutes, 4
493 0C). Protein concentration of the supernatant was determined using Micro BCA protein assay kit
494 (Thermo Fisher Scientific, catalog #23235). The lysates were separated on a 12% SDS-PAGE or
495 a NuPAGE™ 4-12% Bis-Tris protein gel. After separation, proteins were transferred onto a PVDF
496 membrane (Milipore Immobilon-P, Fisher, catalog #IPVH00010). Then the membranes were
497 probed with specific antibodies: Flag (Sigma, catalog #F1804), KLK4/KLKP1 (Eurogentec
498 custom synthesized antibody) and β -actin (Sigma, catalog #A2228). The membranes were
499 visualized on an imaging system (ChemiDoc, BIO-RAD) using a chemiluminescence developing
500 kit (Clarity™ Western ECL Blotting Substrates, BIO-RAD, catalog #1705060).

501 **Measurement of cell proliferation**

502 Cell proliferation was measured by cell counting. For this, stable RWPE-1 cells
503 overexpressing KLK4-KLKP1 (DDK-tagged) or lacZ were used. The cells were seeded at a
504 density of 10 000 cells per well in 24-well plates (n=3). Next, the cells were trypsinized and
505 counted at specified time points by Z2 Coulter particle counter (Beckman Coulter, Brea, CA,

506 USA). LacZ cells were served as controls. Each experiment has been performed with three
507 replicates per sample.

508 **Matrigel invasion assay**

509 Matrigel invasion assays were performed using BD BioCoat Matrigel matrix (Corning Life
510 Sciences, Tewksbury, MA, USA). The parental and transfected clones of RWPE-1 and PrEC cells
511 were seeded at 1×10^5 cells in serum-free medium in the upper chamber of a 24-well culture plate.
512 The lower chamber containing respective medium was supplemented with 10% serum as a
513 chemoattractant. After 48 h, the non-invading cells and Matrigel matrix from the upper side of the
514 chamber were gently wiped with a cotton swab. Invasive cells located on the lower side of the
515 chamber were stained with 0.2% crystal violet in methanol, air-dried and photographed using an
516 inverted microscope (4x). Invasion was quantified by colorimetric assay or by counting the number
517 of cells. For colorimetric assays, the inserts were treated with 150 μ l of 10% acetic acid and the
518 absorbance measured at 560 nm.

519 **Chicken Chorioallantoic Membrane Assay (CAM) assay**

520 CAM assay was performed as described earlier [31]. Briefly, fertilized eggs were incubated
521 in a rotary humidified incubator at 38°C for 10 days. CAM was dropped by making two holes, one
522 through the eggshell into the air sac and a second hole near the allantoic vein that penetrates the
523 eggshell membrane but not the CAM. Subsequently a cutoff wheel (Dremel) was used to cut a 1
524 cm² window to expose the underlying CAM near the allantoic vein. After 3 days of implanting
525 the 2×10^6 cells in 50 μ l medium on the top of each egg, lower CAM was harvested and analyzed
526 for the presence of tumor cells by quantitative human Alu-specific PCR. Genomic DNA from
527 lower CAM and livers were prepared using Puregene DNA purification system (Qiagen USA) and
528 quantification of human-Alu was performed as described (Ref). After 7 days of implantation,

529 extraembryonic tumors were isolated and weighed. An average of 8 eggs per group was used in ea

530 **Gene expression microarray analysis**

531 Two-channel microarray experiment was performed with two replicates using the Agilent
532 Whole Human Genome Oligo Microarray (Agilent, catalog #G4851C Whole Human Genome
533 Microarray 8x60K). Raw data from each replicates were independently processed using
534 Bioconductor packages. “agilp” Bioconductor package (1) was used to apply loess normalization
535 on raw expression values. Fold change for each probe was obtained by taking difference of loess-
536 normalized, log₂-transformed signal intensity between sample with KLK4-KLKP1 gene fusion
537 and control sample. Probes showing differential expression in both two-channel experiments were
538 considered for functional analysis. In total, 1956 probes were up-regulated (with Log₂FC >=1)
539 and 1918 probes were down-regulated (with log₂FC <= -1) in KLK4-KLKP1 gene fusion sample.
540 Heatmap of differentially expressed genes was created using heatmap.2 of “gplots” R package.

541 **Gene set enrichment analysis (GSEA)**

542 Gene set enrichment analysis (GSEA) was performed using the curated gene sets [C2]
543 (n=1267) from Molecular Signature Database (MSigDB v5.0) provided by Broad institute (2)
544 Differentially expressed genes were ranked by average log₂FC from two arrays and submitted to
545 GSEAPreranked module in GSEA software.

546 **KEGG pathway analysis**

547 DAVID (Database for Annotation, Visualization and Integrated Discovery) v6.8 (3) was
548 used to identify enriched KEGG pathways in these differentially expressed genes. With default
549 parameters (gene count of 2 and EASE of 0.1), functional annotation chart was obtained and
550 KEGG pathways with p-value <0.05 were considered to be enriched.

551 **Screening of KLK4-KLKP1 in the urine samples of prostate cancer patients**

552 Random urine samples were collected with informed consent and Institutional Review
553 Board approval from PCa patients visiting the Hematology Oncology clinic at Henry Ford hospital
554 in Detroit, MI. RNA was isolated using ZR urine RNA isolation kit™ (Zymo Research, catalog #
555 R1038 & R1039) according to manufacturer's instructions. cDNA synthesis and qRT-PCR were
556 performed as described earlier.

557 **Statistical analysis**

558 Pearson's chi-square test was used to evaluate the association of KLK4-KLKP1 fusion with
559 race, age, Gleason score and other molecular markers. For association between KLK4-KLKP1 and
560 pre-operative PSA, two-sample t-test was performed to evaluate difference in log-transformed pre-
561 operative PSA between KLK4-KLKP1 positive and negative cases. Multivariable Cox regression
562 was used to estimate the association between KLK4-KLKP1 and the risk of biochemical
563 recurrence. Cox regression model was adjusted for patients' age group (<50; ≥50), Gleason score
564 (6 or 3+4; 4+3 or 8+), and tumor stage (pT2; pT3 or pT4). For all analyses, a p-value of <.05 was
565 considered statistically significant. All analyses were performed using the Statistical Analysis
566 System (SAS) statistical software package, version 9.1.3. For the rest of the experiments, Student's
567 two-sample t-test was used to determine significant differences between two groups. P-values
568 <0.05 were considered significant.

569

570 **ACKNOWLEDGEMENTS**

571 We would like to thank Mireya Diaz-Insua for helping to identify the prostate patient
572 cohort. Natalia Draga and Jingli Yang for making the prostate tissue microarray. We thank the
573 University of Michigan Sequencing Core and Vector Core facilities for their assistance in
574 sequencing of the clones and construction of the adenoviral and lentiviral constructs used in this

575 study. Arul Chinnaiyan for his support at the Michigan Center for Translational Pathology,
576 University of Michigan during the early phase of the project.

577

578 **AUTHOR CONTRIBUTIONS**

579 Conception and design: NP

580 Development and methodology: SV, NP

581 Acquisition of data: BVSKC, PDA, SC

582 Analysis and interpretation of data: SK, JL, KHHW, DSC, NP, SV, PDA, NG, SW, DC, NP

583 Writing, review and/or revision of the manuscript: NP, PDA

584 Administrative, technical, or material support: NN, JP, HS, CR, MM,

585 Study supervision: SV, NP

586 Other:

587

588 **REFERENCES**

- 589 [1] Cronin KA, Lake AJ, Scott S, Sherman RL, Noone AM, Howlader N, Henley SJ,
590 Anderson RN, Firth AU, Ma J, et al. (2018). Annual Report to the Nation on the Status of
591 Cancer, part I: National cancer statistics *Cancer* **124**, 2785-2800.
- 592 [2] Bray F, Ferlay J, Soerjomataram I, Siegel RL, Torre LA, Jemal A (2018). Global cancer
593 statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36
594 cancers in 185 countries *CA Cancer J Clin* **68**, 394-424.
- 595 [3] Abate-Shen C, Shen MM (2000). Molecular genetics of prostate cancer *Genes &*
596 *development* **14**, 2410-2434.

- 597 [4] Arora R, Koch MO, Eble JN, Ulbright TM, Li L, Cheng L (2004). Heterogeneity of
598 Gleason grade in multifocal adenocarcinoma of the prostate *Cancer* **100**, 2362-2366.
- 599 [5] Cheng L, MacLennan GT, Lopez-Beltran A, Montironi R (2012). Anatomic,
600 morphologic and genetic heterogeneity of prostate cancer: implications for clinical
601 practice *Expert Rev Anticancer Ther* **12**, 1371-1374.
- 602 [6] Cancer Genome Atlas Research N (2015). The molecular taxonomy of primary prostate
603 cancer *Cell* **163**, 1011-1025.
- 604 [7] Tomlins SA, Alshalalfa M, Davicioni E, Erho N, Yousefi K, Zhao S, Haddad Z, Den RB,
605 Dicker AP, Trock BJ, et al. (2015). Characterization of 1577 primary prostate cancers
606 reveals novel biological and clinicopathologic insights into molecular subtypes *Eur Urol*
607 **68**, 555-567.
- 608 [8] Yang L, Wang S, Zhou M, Chen X, Jiang W, Zuo Y, Lv Y (2017). Molecular
609 classification of prostate adenocarcinoma by the integrated somatic mutation profiles and
610 molecular network *Sci Rep* **7**, 738.
- 611 [9] Tomlins SA, Bjartell A, Chinnaiyan AM, Jenster G, Nam RK, Rubin MA, Schalken JA
612 (2009). ETS gene fusions in prostate cancer: from discovery to daily clinical practice *Eur*
613 *Urol* **56**, 275-286.
- 614 [10] Tomlins SA, Rhodes DR, Yu J, Varambally S, Mehra R, Perner S, Demichelis F,
615 Helgeson BE, Laxman B, Morris DS, et al. (2008). The role of SPINK1 in ETS
616 rearrangement-negative prostate cancers *Cancer cell* **13**, 519-528.
- 617 [11] Palanisamy N, Ateeq B, Kalyana-Sundaram S, Pflueger D, Ramnarayanan K, Shankar S,
618 Han B, Cao Q, Cao X, Suleman K, et al. (2010). Rearrangements of the RAF kinase
619 pathway in prostate cancer, gastric cancer and melanoma *Nat Med* **16**, 793-798.

- 620 [12] Baena E, Shao Z, Linn DE, Glass K, Hamblen MJ, Fujiwara Y, Kim J, Nguyen M, Zhang
621 X, Godinho FJ, et al. (2013). ETV1 directs androgen metabolism and confers aggressive
622 prostate cancer in targeted mice and patients *Genes Dev* **27**, 683-698.
- 623 [13] Nam RK, Sugar L, Yang W, Srivastava S, Klotz LH, Yang LY, Stanimirovic A, Encioiu
624 E, Neill M, Loblaw DA, et al. (2007). Expression of the TMPRSS2:ERG fusion gene
625 predicts cancer recurrence after surgery for localised prostate cancer *Br J Cancer* **97**,
626 1690-1695.
- 627 [14] Gleason DF (1992). Histologic grading of prostate cancer: a perspective *Hum Pathol* **23**,
628 273-279.
- 629 [15] Adhyam M, Gupta AK (2012). A Review on the Clinical Utility of PSA in Cancer
630 Prostate *Indian J Surg Oncol* **3**, 120-129.
- 631 [16] Carter HB, Partin AW, Walsh PC, Trock BJ, Veltri RW, Nelson WG, Coffey DS, Singer
632 EA, Epstein JI (2012). Gleason score 6 adenocarcinoma: should it be labeled as cancer? *J*
633 *Clin Oncol* **30**, 4294-4296.
- 634 [17] Tan SH, Petrovics G, Srivastava S (2018). Prostate cancer genomics: recent advances and
635 the prevailing underrepresentation from racial and ethnic minorities *Int J Mol Sci* **19**,
636 E1255.
- 637 [18] Kalyana-Sundaram S, Kumar-Sinha C, Shankar S, Robinson DR, Wu YM, Cao X,
638 Asangani IA, Kothari V, Prensner JR, Lonigro RJ, et al. (2012). Expressed pseudogenes
639 in the transcriptional landscape of human cancers *Cell* **149**, 1622-1634.
- 640 [19] Poliseno L (2012). Pseudogenes: newly discovered players in human cancer *Sci Signal* **5**,
641 re5.

- 642 [20] Brosch M, Saunders GI, Frankish A, Collins MO, Yu L, Wright J, Verstraten R, Adams
643 DJ, Harrow J, Choudhary JS, et al. (2011). Shotgun proteomics aids discovery of novel
644 protein-coding genes, alternative splicing, and "resurrected" pseudogenes in the mouse
645 genome *Genome Res* **21**, 756-767.
- 646 [21] Kaushal A, Myers SA, Dong Y, Lai J, Tan OL, Bui LT, Hunt ML, Digby MR,
647 Samaratunga H, Gardiner RA, et al. (2008). A novel transcript from the KLKP1 gene is
648 androgen regulated, down-regulated during prostate cancer progression and encodes the
649 first non-serine protease identified from the human kallikrein gene locus *Prostate* **68**,
650 381-399.
- 651 [22] Clements J, Hooper J, Dong Y, Harvey T (2001). The expanded human kallikrein (KLK)
652 gene family: genomic organisation, tissue-specific expression and potential functions *Biol*
653 *Chem* **382**, 5-14.
- 654 [23] Warrick JI, Tomlins SA, Carskadon SL, Young AM, Siddiqui J, Wei JT, Chinnaiyan
655 AM, Kunju LP, Palanisamy N (2014). Evaluation of tissue PCA3 expression in prostate
656 cancer by RNA in situ hybridization--a correlative study with urine PCA3 and
657 TMPRSS2-ERG *Mod Pathol* **27**, 609-620.
- 658 [24] Leinonen KA, Saramaki OR, Furusato B, Kimura T, Takahashi H, Egawa S, Suzuki H,
659 Keiger K, Ho Hahm S, Isaacs WB, et al. (2013). Loss of PTEN is associated with
660 aggressive behavior in ERG-positive prostate cancer *Cancer Epidemiol Biomarkers Prev*
661 **22**, 2333-2344.
- 662 [25] Ahearn TU, Pettersson A, Ebot EM, Gerke T, Graff RE, Morais CL, Hicks JL, Wilson
663 KM, Rider JR, Sesso HD, et al. (2016). A prospective investigation of PTEN loss and
664 ERG expression in lethal prostate cancer *J Natl Cancer Inst* **108**, djv346.

- 665 [26] Fontugne J, Lee D, Cantaloni C, Barbieri CE, Caffo O, Hanspeter E, Mazzoleni G, Dalla
666 Palma P, Rubin MA, Fellin G, et al. (2014). Recurrent prostate cancer genomic
667 alterations predict response to brachytherapy treatment *Cancer Epidemiol Biomarkers*
668 *Prev* **23**, 594-600.
- 669 [27] Bismar TA, Hegazy S, Feng Z, Yu D, Donnelly B, Palanisamy N, Trock BJ (2018).
670 Clinical utility of assessing PTEN and ERG protein expression in prostate cancer
671 patients: a proposed method for risk stratification *J Cancer Res Clin Oncol* **144**, 2117-
672 2125.
- 673 [28] Navone NM, van Weerden WM, Vessella RL, Williams ED, Wang Y, Isaacs JT, Nguyen
674 HM, Culig Z, van der Pluijm G, Rentsch CA, et al. (2018). Movember GAP1 PDX
675 project: An international collection of serially transplantable prostate cancer patient-
676 derived xenograft (PDX) models *Prostate* **78**, 1262-1282.
- 677 [29] Ren G, Baritaki S, Marathe H, Feng J, Park S, Beach S, Bazeley PS, Beshir AB, Fenteany
678 G, Mehra R, et al. (2012). Polycomb protein EZH2 regulates tumor invasion via the
679 transcriptional repression of the metastasis suppressor RKIP in breast and prostate cancer
680 *Cancer Res* **72**, 3091-3104.
- 681 [30] Bryant RJ, Cross NA, Eaton CL, Hamdy FC, Cunliffe VT (2007). EZH2 promotes
682 proliferation and invasiveness of prostate cancer cells *Prostate* **67**, 547-556.
- 683 [31] Chakravarthi BV, Goswami MT, Pathi SS, Robinson AD, Cieslik M, Chandrashekar DS,
684 Agarwal S, Siddiqui J, Daignault S, Carskadon SL, et al. (2016). MicroRNA-101
685 regulated transcriptional modulator SUB1 plays a role in prostate cancer *Oncogene* **35**,
686 6330-6340.

- 687 [32] Chakravarthi BV, Pathi SS, Goswami MT, Cieslik M, Zheng H, Nallasivam S, Arekapudi
688 SR, Jing X, Siddiqui J, Athanikar J, et al. (2014). The miR-124-prolyl hydroxylase
689 P4HA1-MMP1 axis plays a critical role in prostate cancer progression *Oncotarget* **5**,
690 6654-6669.
- 691 [33] Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich
692 A, Pomeroy SL, Golub TR, Lander ES, et al. (2005). Gene set enrichment analysis: a
693 knowledge-based approach for interpreting genome-wide expression profiles *Proc Natl*
694 *Acad Sci U S A* **102**, 15545-15550.
- 695 [34] Huang DW, Sherman BT, Tan Q, Kir J, Liu D, Bryant D, Guo Y, Stephens R, Baseler
696 MW, Lane HC, et al. (2007). DAVID Bioinformatics Resources: expanded annotation
697 database and novel algorithms to better extract biology from large gene lists *Nucleic*
698 *Acids Res* **35**, W169-175.
- 699 [35] Heinlein CA, Chang C (2004). Androgen receptor in prostate cancer *Endocr Rev* **25**, 276-
700 308.
- 701 [36] Yu J, Yu J, Mani RS, Cao Q, Brenner CJ, Cao X, Wang X, Wu L, Li J, Hu M, et al.
702 (2010). An integrated network of androgen receptor, polycomb, and TMPRSS2-ERG
703 gene fusions in prostate cancer progression *Cancer Cell* **17**, 443-454.
- 704 [37] Pink RC, Wicks K, Caley DP, Punch EK, Jacobs L, Carter DR (2011). Pseudogenes:
705 pseudo-functional or key regulators in health and disease? *RNA* **17**, 792-798.
- 706 [38] Lu W, Zhou D, Glusman G, Utleg AG, White JT, Nelson PS, Vasicek TJ, Hood L, Lin B
707 (2006). KLK31P is a novel androgen regulated and transcribed pseudogene of kallikreins
708 that is expressed at lower levels in prostate cancer cells than in normal prostate cells
709 *Prostate* **66**, 936-944.

- 710 [39] Lai J, Lehman ML, Dinger ME, Hendy SC, Mercer TR, Seim I, Lawrence MG, Mattick
711 JS, Clements JA, Nelson CC (2010). A variant of the KLK4 gene is expressed as a cis
712 sense-antisense chimeric transcript in prostate cancer cells *RNA* **16**, 1156-1166.
- 713 [40] Pflueger D, Mittmann C, Dehler S, Rubin MA, Moch H, Schraml P (2015). Functional
714 characterization of BC039389-GATM and KLK4-KRSP1 chimeric read-through
715 transcripts which are up-regulated in renal cell cancer *BMC Genomics* **16**, 247.
- 716 [41] Veveris-Lowe TL, Lawrence MG, Collard RL, Bui L, Herington AC, Nicol DL,
717 Clements JA (2005). Kallikrein 4 (hK4) and prostate-specific antigen (PSA) are
718 associated with the loss of E-cadherin and an epithelial-mesenchymal transition (EMT)-
719 like effect in prostate cancer cells *Endocr Relat Cancer* **12**, 631-643.
- 720 [42] Klock TI, Kilander A, Xi Z, Waehre H, Risberg B, Danielsen HE, Saatcioglu F (2007).
721 Kallikrein 4 is a proliferative factor that is overexpressed in prostate cancer *Cancer Res*
722 **67**, 5221-5230.
- 723 [43] Gao J, Collard RL, Bui L, Herington AC, Nicol DL, Clements JA (2007). Kallikrein 4 is
724 a potential mediator of cellular interactions between cancer cells and osteoblasts in
725 metastatic prostate cancer *Prostate* **67**, 348-360.

726

727

728 **FIGURE LEGENDS**

729 **Figure 1:** The structure of *KLK4-KLKPI* fusion and the RNA-ISH screening of *KLK4-KLKPI* in
730 tissue micro arrays. (A) Schematic diagram of the structure of *KLK4-KLKPI* fusion. *KLK4-*
731 *KLKPI* is formed through the fusion of exon 1 and 2 of *KLK4* gene with exon 4 and 5 of *KLKPI*.
732 (B) The predicted sequence of *KLK4-KLKPI* fusion protein. The sequence in purple is derived
733 from *KLK4* while the sequence in red is originating from *KLKPI*. (C) The expression of *KLK4-*
734 *KLKPI* in prostate tissue cores detected by RNA-ISH. The bottom set of images show an
735 enlarged section of the corresponding tissue core in the top set of images. 1+ to 4+ indicate the
736 intensity of *KLK4-KLKPI* RNA-ISH staining. (D) Prostate cancer specific expression of *KLK4-*
737 *KLKPI*. *KLK4-KLKPI* RNA-ISH staining in benign, HGPIN and prostate cancer tumor cores are
738 shown. The bottom set of images contains a magnified area of the images on the top. 1+ to 4+
739 refer to the intensity of the *KLK4-KLKPI* RNA-ISH staining. (E) *KLK4-KLKPI* is expressed
740 more in the prostate cancer patients (GG1-5) compared to non-cancer (benign, HGPIN, atypical
741 and stroma) cases. The percentage of cases showing a positive *KLK4-KLKPI* RNA-ISH signal
742 among non-cancer and GG1-5 groups is shown. P-value was calculated based on Pearson's chi-
743 square test. (F) *KLK4-KLKPI* is expressed more in young patients. The percentages of cases with
744 positive *KLK4-KLKPI* RNA-ISH signal in the young patient (age lower than 50 years) and old
745 patient groups (age equal to or higher than 50 years) are shown. P-value was calculated based on
746 Pearson's chi-square test. (G) *KLK4-KLKPI* expression is associated with *ERG* overexpression.
747 *SPINK1*, *ETV1*, *ETV4* and *ETV5* overexpression is mutual from *KLK4-KLKPI* expression. *PTEN*
748 loss is significantly lower in cases with *KLK4-KLKPI* expression. The percentages of cases
749 showing positive signal for *ERG*, *SPINK1*, *ETV1*, *ETV4*, *ETV5* or *PTEN* loss among *KLK4-*
750 *KLKPI* RNA-ISH positive cases (dark grey bars) and *KLK4-KLKPI* RNA-ISH negative cases

751 (light grey bars) are shown. P-value was calculated based on Pearson's chi-square test.

752 Abbreviations: GG, Gleason grade; HGPIN, high grade prostate intraepithelial neoplasia; ISH, in
753 situ hybridization.

754

755 **Figure 2:** Validation of the expression of *KLK4-KLKP1* protein in HEK-293 cells and PDX
756 tissues. (A) The qRT-PCR analysis HEK-293 cells transfected with and without FLAG tagged-
757 *KLK4-KLKP1*. HEK-293 cells were transfected with adenoviral vectors carrying FLAG tagged-
758 *KLK4-KLKP1* (adeno-FLAG-*KLK4-KLKP1*). As a control untransfected cells treated with
759 bortezomib were used. The expression of *KLK4-KLKP1* was confirmed by qRT-PCR. (B)
760 Western blot analysis of HEK-293 cells transfected with FLAG-tagged *KLK4-KLKP1* using anti-
761 FLAG, anti-*KLK4-KLKP1* and anti- β -actin antibody. (C) Western blot analysis of *KLK4-KLKP1*
762 qRT-PCR negative (MDA PCa144-13) and qRT-PCR positive (MDA PCa 153-7) PDX tissues
763 using anti-*KLK4-KLKP1* and anti- β -actin antibody. (D) IHC staining of *KLK4-KLKP1* qRT-
764 PCR positive and qRT--PCR negative PDX models. Abbreviations: PDX, patient derived
765 xenografts; qRT-PCR, quantitative PCR. Images of original western blots show anti-N-FLAG
766 antibody (2E-left), and anti-*KLK4-KLKP1* antibody (2E-right). Images of original western blots
767 show, lysates from the prostate xenografts each positive (MDA PCa 153-7) and negative (MDA
768 PCa 144-13) for endogenous expression of *KLK4-KLKP1* transcript were probed with anti-
769 *KLK4-KLKP1* antibody (2F). Validation of *KLK4-KLKP1* specific antibody in comparison with
770 RNA-ISH (2G). PCa tissue confirmed to be positive by RNA-ISH (left) is positive for the
771 antibody whereas the tumor negative for *KLK4-KLKP1* by RNA-ISH also negative for the
772 antibody by IHC, thus confirming the specificity of the new antibody to the *KLK4-KLKP1*
773 fusion protein.

774

775 **Figure 3:** Functional characterization of *KLK4-KLKPI*. (A) qRT-PCR validation of *KLK4-*

776 *KLKPI* expression in RWPE-1 cells after stable transfection with FLAG tagged-*KLK4-KLKPI*.

777 As controls untransfected cells (control) and cells transfected with LacZ were used. (B) Analysis

778 of cellular proliferation in RWPE-1 cells stably expressing FLAG tagged *KLK4-KLKPI*. Cells

779 were plated in 96-well plates. The number of cells was measured on days 2, 4, 6 and 8 using a

780 Coulter particle counter. Cells untransfected and transfected with LACZ were used as controls.

781 (C) Analysis of cell invasion in RWPE-1 cells. The invasion of RWPE-1 cells stably transfected

782 with either FLAG tagged-*KLK4-KLKPI* or LacZ was studied using the Boyden chamber assay.

783 Untransfected cells were also used as a control. After invasion of cells into the invasion chamber,

784 cells were fixed and visualized using crystal violet. Additionally, the invasion chamber

785 membranes carrying the fixed cells were dipped in glacial acetic acid and the absorbance at 560

786 nm was also measured. Representative images of the crystal violet stained cells that underwent

787 invasion in each case and the absorbance at 560 nm are shown. (D) Analysis of cell invasion in

788 PrEC cells. The cellular invasion in PrEC cells transfected with FLAG tagged-*KLK4-KLKPI*

789 was performed as described in Figure 2C. The number of invaded cells were counted and plotted.

790 In addition to LACZ and untransfected cells, PrEC cells transfected with *EZH2* were also used a

791 control. (E) Intravasation of RWPE-1 cells measured using CAM assay. RWPE-1 cells stably

792 transfected with FLAG tagged-*KLK4-KLKPI*, were implanted on eggs. The presence of

793 intravasated cells in the lower CAM was assessed by quantitative human Alu-specific PCR.

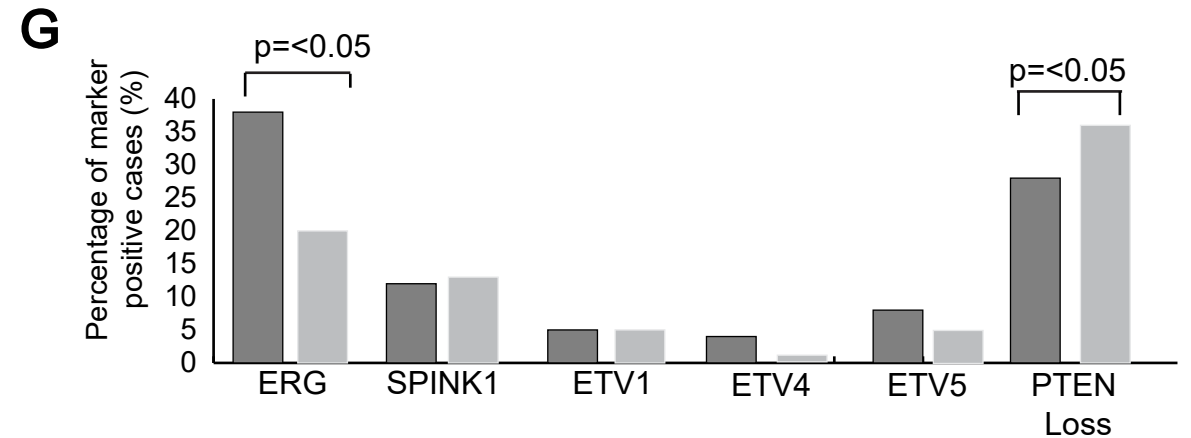
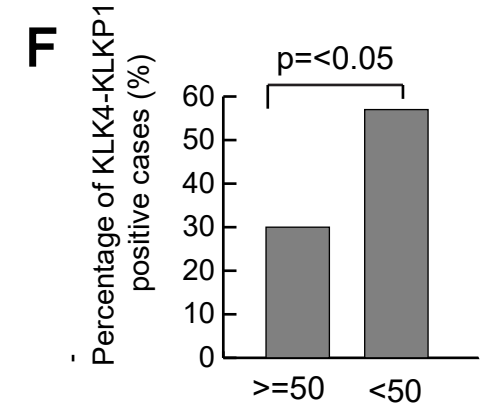
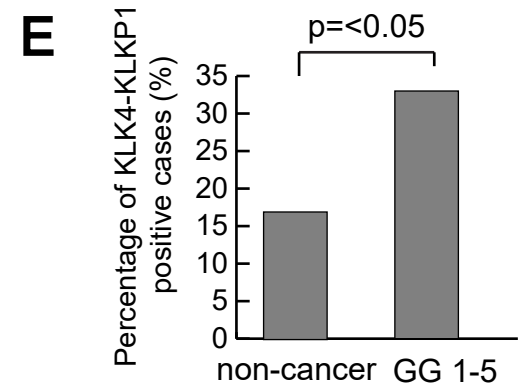
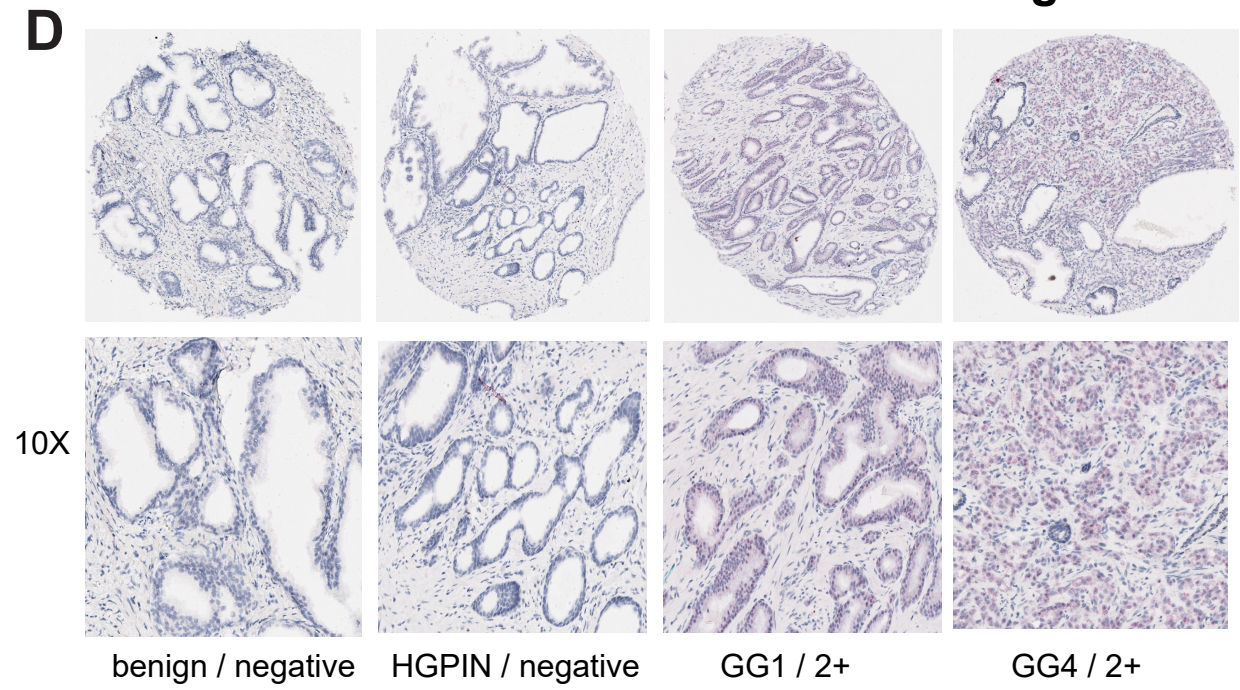
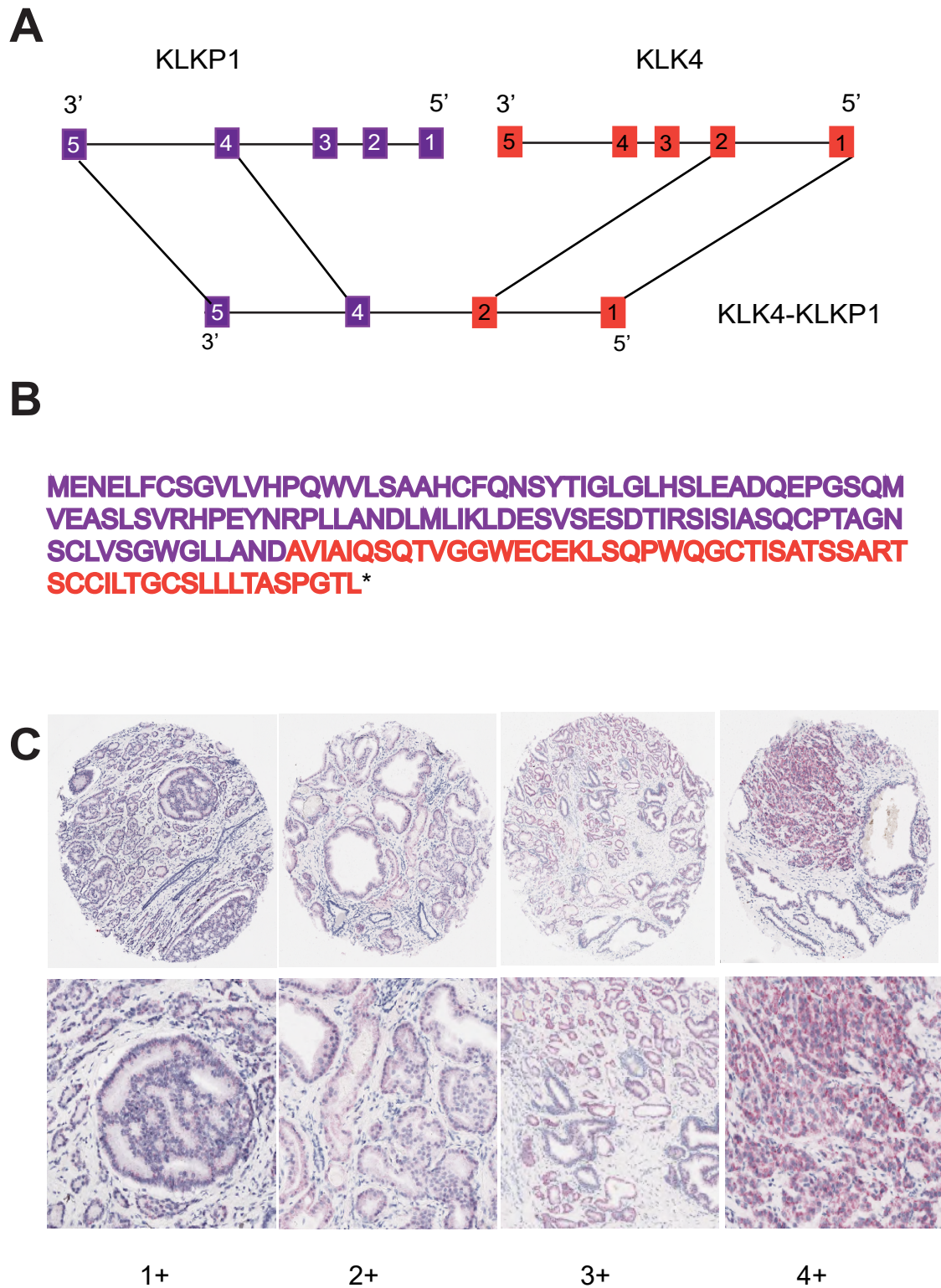
794 Untransfected cells and cells transfected with LACZ were used as controls. (F) Analysis of

795 weight of extraembryonic tumors isolated from eggs implanted with RWPE-1 cells stably

796 expressing FLAG-tagged *KLK4-KLKPI*. Cells transfected with LACZ and untransfected cells
797 were used as controls. Abbreviations: CAM, chicken chorioallantoic membrane.

798 **Figure 4:** Gene expression analysis of *KLK4-KLKPI*. (A) Heat map showing the top 100 genes
799 differentially expressed in RWPE-1 cells stably transfected with *KLK4-KLKPI* compared to cells
800 transfected with LACZ. The results from 2 independent trials are shown. (B) Gene set
801 enrichment analysis of differentially expressed genes. The genes were enriched in 2 curated gene
802 sets, one involving genes upregulated in endometrioid endometrial metastatic tumor
803 “BIDUS_METASTASIS_UP” (top image) and the other including genes overexpressed in
804 melanoma metastatic cancer “WINNEPENNINCKX_METASTASIS_UP” (bottom image). (C)
805 Top 10 KEGG pathways enriched in differentially expressed genes obtained using DAVID tool.

Figure 1



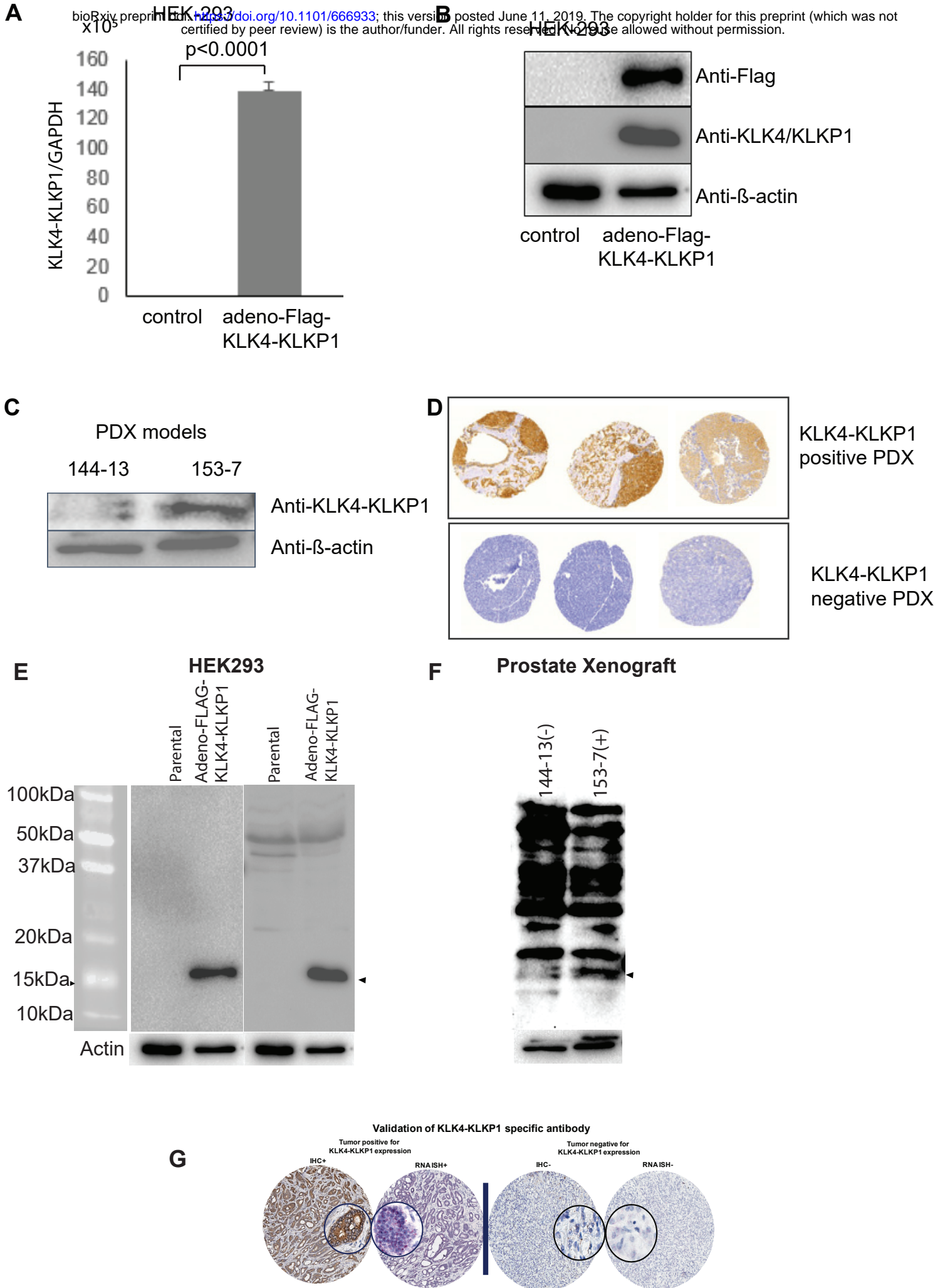
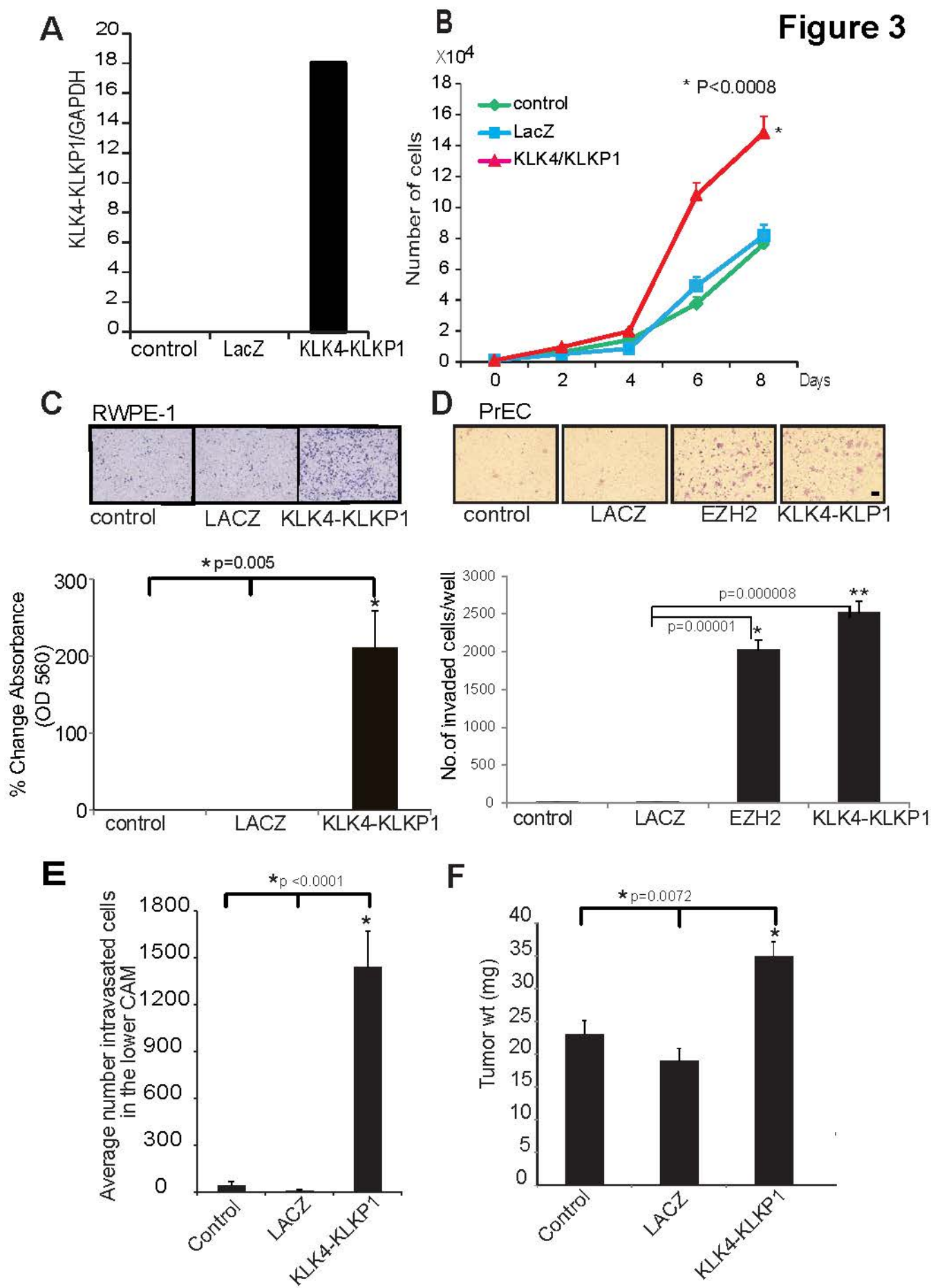
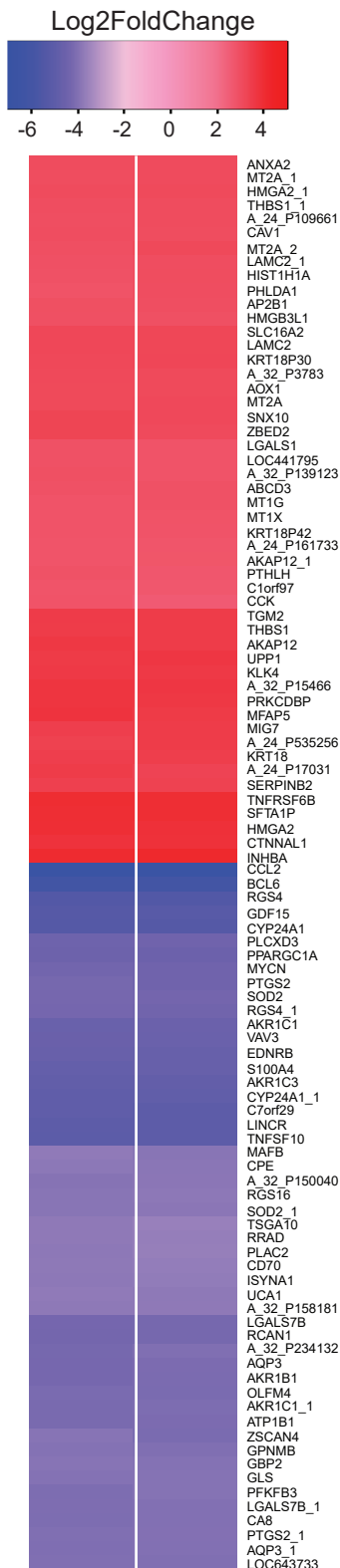
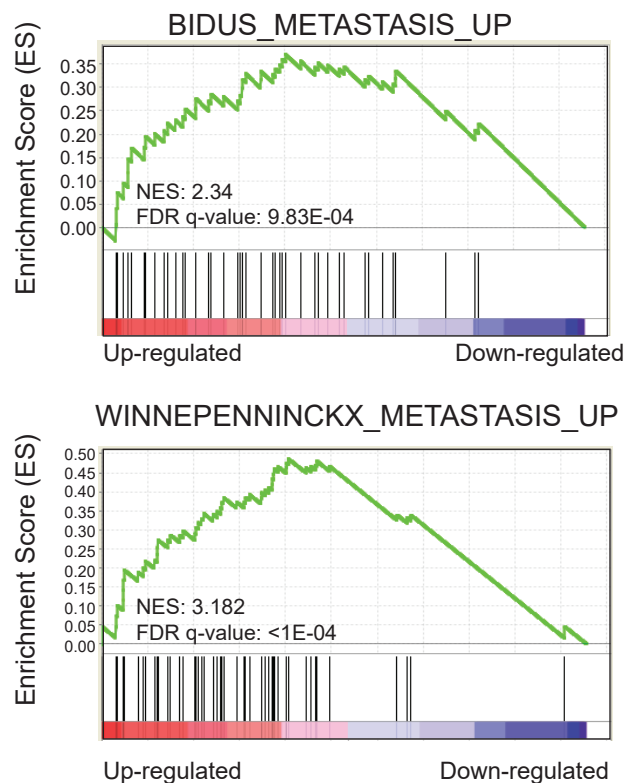


Figure 3

A



B



C

