

1 **Supplementary materials and methods**

2 **Single cell preparation**

3 Fresh tissue of each region was washed with 10 mL cold PBS for three times.
4 Necrotic tissues were excluded artificially under microscope. Then the tissues were
5 cut into small pieces by surgical scissors. 2 mg/mL Collagenase/Dispase (Roche, cat.
6 # 10269638001) was used to digest the tissues into single cells at 37°C for 30-40 min.
7 Then the single cell suspension was centrifuged at 800g for 5min at 4°C. The cell
8 pellets were washed with cold PBS three times and finally re-suspended in 1% HAS
9 (Human Serum Albumin). Single cells were picked into 2µl scRNA-seq lysis buffer
10 (0.0475% Triton X-100 (Sigma-Aldrich, X100), 0.1 U/µL RNase Inhibitor (Takara),
11 2.5 µM dNTP mixture (Thermo) and 2.5 µM Transcription reverse (RT) primers)
12 gently with mouth pipette. And the remaining cell pellets were used to extract
13 genomic DNA according to the manufacturer protocol of DNeasy Blood & Tissue Kit
14 (QIAGEN).

15

16 **Single cell RNA-seq library construction**

17 Single cell RNA-seq libraries were carried out using STRT-seq protocol (reference).
18 We modified the protocol to make it compatible to multiplexed single cell sequencing
19 as mentioned in the published papers by our group[1,2]. Briefly, we add 8bp barcode
20 sequence to the RT primer
21 (5'-TCAGACGTGTGCTCTTCCGATCT-XXXXXXXXX-NNNNNNNN-T25-3', X
22 represents the 8bp predesigned barcodes, 'N' stands for the UMI (unique molecular
23 identifier). We lysed the single cells at 72°C for 3min and then add 2.85µl RT Mix (5
24 U RNase Inhibitor; 40 U SuperScript II reverse transcriptase (Invitrogen); 5×
25 Superscript II first-strand buffer (Invitrogen); 25 mM dithiothreitol; 5 M betaine
26 (Sigma-Aldrich); 30 mM MgCl₂ (Sigma-Aldrich); 1.75 µM TSO primer (same
27 sequence as described in STRT protocol)). Reverse transcription was performed in
28 thermocycler using program as follows: at 25 °C for 5 min; at 42 °C for 60 min; at
29 50 °C for 30 min; at 70 °C for 10 min. Then 7.5µl PCR Mix (6.25 µL 2× KAPA HiFi

30 HotStart ReadyMix (KAPA); 300 nM ISPCR oligo primer
31 (5'-AAGCAGTGGTATCAACGCAGAGT-3'); 1 μ M 3' Anchored oligo primer
32 (5'-Oligo-GTGACTGGAGTTCA GACGTGTGCTCTTCCGATC-3')) to amplify the
33 cDNA. The cDNA was amplified at 98°C for 20 s; 65°C for 30 s; 72°C for 5 min; 14
34 cycle of 98°C for 20 s, 67°C for 15 s, 72 °C for 5 min; 72 °C for 5 min. The cDNA
35 of cells with different barcodes can be pooled together for next steps. The pooled
36 DNA was purified once with DNA Clean & Concentrator-5 (DC2005; Vistech) and
37 twice with 0.8 \times Ampure XP beads (Beckman, A63882). Next, we use bio- tynlated
38 index primer (5'-/Biotin/CAAGCAGAAGACGG
39 CATACGAGAT-index-GTGACTGGAGTTCAGACGTGTGCTCTTCCGATC-3')
40 and ISPCR primer to amplify the pooled DNA. Then the biotin-labeled DNA was
41 fragmented into 300-bp using ultra-sonication. The fragmented DNA with biotin was
42 enriched by the affinity interaction between biotin and Dynabeads® MyOne™
43 Streptavidin C1 beads (Invitrogen, cat. #65602). Then the fragmented DNA inserts
44 were library constructed according to the manufacturer procedures of KAPA Hyper
45 Prep Kits with PCR Library Amplification/Illumina series (KAPA, cat. KK8054) and
46 sequencing on Illumina HiSeq 4000 for obtaining 150 bp paired end reads.

47

48 **Whole genome sequencing**

49 The extracted genomic DNA for each region was fragmented by sonication into
50 300-bp and then was used to construct the sequencing library according to the
51 manufacturer procedures of KAPA Hyper Prep Kits with PCR Library
52 Amplification/Illumina series (KAPA, cat. KK8054) and sequenced on Illumina
53 HiSeq 4000 for obtaining 150-bp paired end reads.

54

55 **Hematoxylin-eosin staining and Immunohistochemistry**

56 The sampled tissues were fixed in neutral buffered formalin and embedded in paraffin.
57 H&E staining and immunohistochemistry staining were performed on 5-um thick
58 slides according to standard procedures. For immunohistochemistry staining,

secondary horseradish peroxidase (HRP)-conjugated antibodies were used to reveal the color in combination with DAP Peroxidase substrate. The primary antibodies used included Anti-Ki67 antibody (1:3000) (Abcam, ab15580), Anti-TGFBI antibody (1:250) (Abcam, ab170874) and Anti-CD8 antibody (1:200) (Abcam, ab17147). We captured the images with a Nikon Eclipse 90i and the pathological grades of tumors were assessed by two independent pathologists.

Quality control, mapping and quantification of single cell RNA-Seq data

Paired-end sequence data was firstly split according to the cell barcode. Low quality reads, TSO and poly A were removed and then the cleaned reads were mapped to human genome (hg19) by TopHat[3]. Only unique mapping reads were kept for further analysis. By using HTSeq[4], transcript abundance of each gene was calculated and reads with unique UMIs of each genes were used to measure the expression levels of genes. The expression levels of genes were normalized to the transcripts per million (TPM). Cells that expressed less than 1000 genes or showed similar expression pattern with less than 2 samples (pearson correlation < 0.6) were filtered. Genes expressed by less than 5 cells were removed for subsequent analysis.

Quality control and mapping for whole genome sequence

Low quality and adapter contaminate reads were firstly removed by applying quality control pipeline. The cleaned reads were mapped to human genome (hg19) with BWA (v0.7.12) and duplicated reads were marked by picard “MarkDuplicates”. As for copy number variants, mapped reads of each 10M window were calculated and normalized by the total mapped read depths of each samples. For each 10M window, reads were further normalized by dividing the average reads of diploid samples and the results were visualized by R ggplot2 package.

Whole exome sequence and SNV calling

First low quality and adapter contaminate reads were removed by applying quality control pipeline. Clean reads were mapped to human genome (hg19) through BWA.

89 GATK (Genome Analysis Toolkit, Version 3.8) was performed on duplicates marked
90 reads. Haplotype[5] and Mutect2[6] were used to call germline and somatic mutations
91 respectively according to the online manual. “PASS” somatic mutations were filtered
92 for further analysis. For each germline mutation, the reads sequenced more than 30X
93 were kept for further analysis. For somatic mutations, WES data of blood from each
94 patient were used as control. cnLOH includes the various loss of heterogeneity (LOH)
95 events including, copy gain, copy loss and copy-neutral, and it was inferred from
96 whole exome sequencing (WES) by using ‘Sequenza’ tool according to the online
97 manual (<https://cran.r-project.org/web/packages/sequenza/vignettes/sequenza.html>)[7].
98 Sanger sequencing was used to verify the germline and somatic mutation sites in APC,
99 as well as the germline mutations in *MUTYH*. The primers used for Sanger
100 sequencing were summarized in Supplementary figure 3.

101

102 **Prediction of potential pathogenic mutations**

103 For the somatic mutations called by Mutect2, further filtrations were used for
104 prediction of potential pathogenic mutations. Firstly, synonymous mutations and
105 mutations in noncoding regions were excluded. Then, non-synonymous mutations that
106 would cause amino acid changes or have “High effect” on the protein were kept.
107 Finally, the filtered mutation sites on recurrently mutated genes in CRC were
108 supposed to be potential pathogenic mutations (Supplementary table 4).

109

110 **Tumor phylogenetic tree reconstruction**

111 The somatic mutations of each patient were used for the construction of tumor
112 phylogenetic tree. The ‘dist.gene’ function of ‘ape’ package was used to calculate the
113 distance between different lesions of each patient. Then the “nj” function was used to
114 construct the phylogenetic tree, which was finally displayed by running “plot.phylo”
115 function.

116

117 **tSNE clustering based on the transcription factor regulation network**

118 The python package pySCENIC[8] was used to construct the transcription factor

119 regulation network following the online manual
120 (<https://github.com/aertslab/pySCENIC>)[8]. The cellular regulon enrichment matrix
121 was used to calculate the distance between samples. Non-linear dimensional reduction
122 was analyzed by performing Rtsne package based on the distance matrix. Cell clusters
123 were identified based on the first nine PCA dimensions and annotated by known cell
124 type markers.

125

126 **PCA and tSNE clustering based on the transcriptome data**

127 The FactorMinR package was used to perform the PCA analysis based on the top 500
128 high variants genes and the “RunTSNE” function in Seurat package was used to
129 perform the tSNE clustering with parameters `dims.use=1:10, do.fast=T`.

130

131 **Identification of differentially expressed genes**

132 The “FindAllMarkers” function of Seurat package[9] were used to identified DEGs
133 among different groups with setting the parameters as “`min.pct = 0.25, only.pos = T,`
134 `logfc.threshold = 1.5`”.

135

136 **Pseudotime analysis**

137 Firstly, DEGs between adjacent normal tissue and tumor were identified by Seurat
138 package. Then Monocle2[10] was used to order cells according to the progression of
139 tumorigenesis based on the DEGs between normal and tumor tissues. We first reduced
140 the dimension of our dataset by running “reduceDimension” function and then
141 “orderCells” function was used to order cells along the pseudotime trajectory. The
142 DEGs between epithelial cells from adjacent normal tissue and multi-grade lesions, as
143 well as the TCA related genes (from KEGG
144 <https://www.genome.jp/kegg/pathway.html>) were plotted with
145 ‘plot_pseudotime_heatmap’ function in monocle package with setting
146 ‘show_rownames = T, num_clusters=5’

147

148

149 **References**

- 150 1 Dong J, Hu Y, Fan X, *et al.* Single-cell RNA-seq analysis unveils a prevalent
 151 epithelial/mesenchymal hybrid state during mouse organogenesis. *Genome Biol*
 152 2018;**19**:1–20. doi:10.1186/s13059-018-1416-2
- 153 2 Islam S, Kjällquist U, Moliner A, *et al.* Characterization of the single-cell
 154 transcriptional landscape by highly multiplex RNA-seq. *Genome Res*
 155 2011;**21**:1160–7. doi:10.1101/gr.110882.110
- 156 3 Trapnell C, Pachter L, Salzberg SL. TopHat: Discovering splice junctions with
 157 RNA-Seq. *Bioinformatics* 2009;**25**:1105–11.
 158 doi:10.1093/bioinformatics/btp120
- 159 4 Anders S, Pyl PT, Huber W. HTSeq-A Python framework to work with
 160 high-throughput sequencing data. *Bioinformatics* 2015;**31**:166–9.
 161 doi:10.1093/bioinformatics/btu638
- 162 5 Depristo MA, Banks E, Poplin R, *et al.* A framework for variation discovery and
 163 genotyping using next-generation DNA sequencing data. *Nat Genet*
 164 2011;**43**:491–501. doi:10.1038/ng.806
- 165 6 Cibulskis K, Lawrence MS, Carter SL, *et al.* Sensitive detection of somatic point
 166 mutations in impure and heterogeneous cancer samples. *Nat Biotechnol*
 167 2013;**31**:213–9. doi:10.1038/nbt.2514
- 168 7 Favero F, Eklund AC, Joshi T, *et al.* Sequenza: allele-specific copy number and
 169 mutation profiles from tumor sequencing data. *Ann Oncol* Published Online First:
 170 2014. doi:10.1093/annonc/mdl479
- 171 8 Aibar S, Bravo González-Blas C, Moerman T, *et al.* SCENIC: Single-Cell
 172 Regulatory Network Inference And Clustering. *bioRxiv* 2017;;1–41.
 173 doi:10.1101/144501
- 174 9 Satija R, Farrell JA, Gennert D, *et al.* Spatial reconstruction of single-cell gene
 175 expression data. *Nat Biotechnol* 2015;**33**:495–502. doi:10.1038/nbt.3192
- 176 10 Trapnell C, Cacchiarelli D, Grimsby J, *et al.* The dynamics and regulators of cell
 177 fate decisions are revealed by pseudotemporal ordering of single cells. *Nat*
 178 *Biotechnol* 2014;**32**:381–6. doi:10.1038/nbt.2859

