

1 **Supplementary Methods**

2

3 **Patients and samples**

4 At the discovery stage, we retrospectively collected paired cancer and adjacent normal tissues
5 from 119 ESCC patients with at least 5 years of follow up. The end of follow-up time for the
6 patients was Dec. 2012 and the median follow-up time was 32.2 months. The clinical endpoints
7 were defined as: died from ESCC before the time of follow-up or being alive at the time of
8 follow-up (at least 5 years after surgery). The exclusion criteria of patients included: received
9 radiotherapy or chemotherapy before surgery; died from reasons rather than ESCC; tumor or
10 normal tissues too small for the assay (<50ug); RNA quality inadequate for microarray assay;
11 tumor cell percentage of tumor tissue less than 60% or normal tissue contaminated by tumor
12 tissues (Supplementary Figure 1). The tissue samples were snap frozen in liquid nitrogen shortly
13 after resection and stored at -80°C till RNA extraction. A small part from every sample was
14 paraffin embedded, sectioned, and H&E stained using routine methods for pathological
15 examination. The tumor histology was independently confirmed by two pathologists (F-XL and
16 S-SS). The percentage of cancer cells was 60 or more in all tumor samples and no cancer cells
17 were found in the normal tissues. To validate the prognostic signature, we enrolled an
18 independent cohort of 60 ESCC patients who underwent surgery at the Cancer Institute and
19 Hospital, CAMS from Jan. 2008 to Dec. 2008. The paired cancer and normal tissues of these
20 patients were tested by the same microarray as the 119 patients. The inclusion and exclusion
21 criterion of these patients was also same as the 119 patients (Supplementary Figure 1). For all
22 samples, clinical and pathologic information (age, gender, pathology, differentiation, TNM stage,
23 co-morbidities, post-operative complications and survival time after surgery) was available
24 (Supplementary Table 1). The TNM stage was based on the American Joint Committee on Cancer
25 staging manual (seventh edition). The surgical procedures the patients received were chosen
26 according to the tumor location, mediastinal lymph nodes, and the general performance status of
27 them. All patients enrolled in the study received R0 resection of the tumor with microscopically
28 negative resection margins. The study was approved by the medical ethics committee of Cancer
29 Institute and Hospital, CAMS.

30

31 **RNA extraction**

32 Total RNA was extracted using the Trizol reagent (Invitrogen) and purified with mirVana miRNA
33 Isolation Kit (Ambion, Austin, TX, USA) according to the manufacturer's protocol. The purity and
34 concentration of RNA were determined by OD260/280 using spectrophotometer (NanoDrop
35 ND-1000). RNA integrity was determined by 1% formaldehyde denaturing gel electrophoresis.

36

37 **RNA amplification and labeling.**

38 cDNA labeled with a fluorescent dye (Cy5 and Cy3-dCTP) was produced by Eberwine's linear RNA
39 amplification method and subsequent enzymatic reaction.[1] In detail, double-stranded cDNAs

1 (containing a T7 RNA polymerase promoter sequence) were synthesized from 1 ug total RNA
2 using the CbcScript reverse transcriptase with cDNA synthesis system according to the
3 manufacturer's protocol (Capitalbio) with the T7 Oligo (dT) and T7 Oligo (dN). After completion of
4 double-stranded cDNA (dsDNA) synthesis using DNA polymerase and RNase H, the dsDNA
5 products were purified using PCR NucleoSpin Extract II Kit (MN) and eluted with 30 ul elution
6 buffer. The eluted double-stranded cDNA products were vacuum evaporated to 16 ul and
7 subjected to 40 ul in vitro transcription reactions at 37°C for 14 hr using a T7 Enzyme Mix. The
8 amplified cRNA was purified using the RNA Clean-up Kit (MN). Klenow enzyme labeling strategy
9 was adopted after reverse transcription using CbcScript II reverse transcriptase. Briefly, 2 ug
10 amplified RNA was mixed with 4 ug random nanomer, denatured at 65°C for 5min, and cooled
11 on ice. Then, 5 ul of 4×first-strand buffer, 2 ul of 0.1M DTT, and 1.5 ul CbcScript II reverse
12 transcriptase were added. The mixtures were incubated at 25°C for 10 min, then at 37°C for 90
13 min. The cDNA products were purified using a PCR NucleoSpin Extract II Kit (MN) and vacuum
14 evaporated to 14 ul. The cDNA was mixed with 4 ug random nanomer, heated to 95°C for 3 min,
15 and snap cooled on ice for 5 min. Then, 5 ul Klenow buffer, dNTP, and Cy5-dCTP or Cy3-dCTP (GE
16 Healthcare) were added to final concentrations of 240 M dATP, 240 M dGTP, 240 M dTTP, 120 M
17 dCTP, and 40 M Cy-dCTP. 1.2 ul Klenow enzyme was then added, and the reaction was performed
18 at 37°C for 90 min. Labeled cDNA was purified with a PCR NucleoSpin Extract II Kit (MN) and
19 resuspended in elution buffer. Labeled controls and test samples labeled with Cy5-dCTP and
20 Cy3-dCTP were dissolved in 80 ul hybridization solution containing 3×SSC, 0.2% SDS,
21 5×Denhardt's solution and 25% formamide.

22

23 **Array hybridization**

24 DNA in hybridization solution was denatured at 95°C for 3 min prior to loading onto microarray.
25 Array hybridization was performed in a CapitalBio BioMixer™ II Hybridization Station overnight
26 at a rotation speed of 8 rpm and a temperature of 42°C and washed with two consecutive
27 solutions (0.2% SDS, 2× SSC at 42°C for 5 min, and 0.2× SSC for 5 min at room temperature).

28

29 **Microarray fabrication**

30 The Agilent human lncRNA+mRNA Array v2.0 was designed with four identical arrays per slide (4 x
31 180K format). Each array contained probes interrogating about 39,000 human lncRNAs and about
32 32,000 human mRNAs. Each RNA was detected by two probe repeats. The array also contained 4974
33 Agilent control probes.

34

35 **Filtering procedure for lncRNAs in the microarray**

36 The probes with same sequence were merged into one, thereafter, resulted in 35,025 unique probes.
37 To obtain the map from the probe to annotated lncRNAs, the UCSC data base, GENCODE(V13) data
38 base and lincRNAs from Cabili et al[2] were taken as the reference annotation (totally 13812 long
39 intergenic and 6528 antisense non-coding RNAs). Then, we employed the blast program to map the
40 probes uniquely to the annotated lncRNA sequences, and 8900 lncRNAs with at least one unique

1 probe were retrieved. For each of the 8900 lncRNAs, the median of the expression values of the
2 probes mapped to it was used as its expression value. If part of the probes mapped to a lncRNA have
3 missing values, the rest of the probes mapped to it was taken to calculate the median expression
4 value. The expression value of a lncRNA was defined missing value when all the probes mapped to had
5 missing value.

6 7 **Missing value imputation using random forest unsupervised learning.**

8 Random Jungle[3] was used to impute missing values by unsupervised learning. It began by filling a
9 rough value of the missing data. Then, a forest including 10,000 trees ran and the proximities were
10 computed. The missing values were estimated based on the proximities between the sample and
11 non-missing value samples. The forest was constructed iteratively and the missing values were
12 re-estimated iteratively. The number of iterations was set to 5.

13 14 **Random forest supervised classification algorithm**

15 In the random forest supervised classification algorithm, an iteration procedure was
16 implemented to narrow down the gene set in which the 1/3 least important lncRNAs were
17 discarded at each iteration step. Ten thousand trees were grown at each step, and the square
18 root of the number of input lncRNAs at each step was set to the size of randomly chosen lncRNAs
19 at each node of single classification tree. Because the number of good-prognostic and
20 poor-prognostic patients were not equal, the class weights were adjusted accordingly. The
21 generalization error was estimated on the out-of-bag samples. Finally, 9 lncRNAs were selected
22 (Figure 1C).

23 24 **Quantitative RT-PCR**

25 In qRT-PCR, the reverse transcription (RT) reactions were carried out with Reverse Transcriptase
26 (SuperScript III, Invitrogen) according to the manufacture's instruction. Around 3ug total RNAs
27 were added to each reaction. Quantitative PCR reactions were then performed on ABI 7900 in a
28 10ul system. The reactions were incubated at 95°C for 5 min, followed by 40 cycles of 95°C for
29 15s, and 60°C for 40s. All quantitative PCR reactions were performed in triplicate. The Ct value of
30 each candidate lncRNA was then normalized to the expression value of GAPDH. Relative
31 expression levels of the lncRNAs were calculated using $2^{-\Delta Ct}$. The sequences of primers used in
32 qRT-PCR of the lncRNAs are listed in Supplementary Table 2.

33 34 **Multiple imputation of Markov chain Monte Carlo (MCMC) for missing value of adjuvant in Cox** 35 **regression analysis**

36 For our data, the probability of missing adjuvant therapy information could be dependent on fully
37 observed clinical factors like age, N stage and TNM stage, and is independent of the unobserved
38 covariable (adjuvant therapy). Thus the missing at random mechanism should be suitable for our
39 data.[4, 5] Here we used the Multiple Imputation procedure in SAS, which is popularly used for
40 missing at random data.[4, 6, 7]

41
42 For the training set, adjuvant therapy information of ten out of 60 patients was missing, and we
43 created ten imputations by multiple Markov chains. At first, we did univariable Cox regression
44 analysis for adjuvant therapy. The results of seven out of ten imputations showed that adjuvant

1 therapy was not significantly associated with survival ($p>0.05$), which were consistent with the
2 combining inference result ($p=0.1694$) (Supplementary Table 3). Then, we did multivariable Cox
3 regression analysis. As that of the univariable analysis, the results of seven in ten imputations
4 showed that adjuvant therapy was not an independent prognostic factor ($p>0.05$ by stepwise
5 regression, the p -value cutoffs of entry and stay were both set to 0.1), and the combining
6 inference result was similar ($p=0.3694$ by full model) (Supplementary Table 3).

8 For the combined test and independent cohort, adjuvant therapy information for 20 of the 119
9 patients was missing and we performed the same procedure as in the training set. In univariable
10 Cox regression analysis, adjuvant therapy was significantly associated with survival in all of the
11 ten imputations ($p<0.05$), and also in the combining inference result ($p=0.0077$) (Supplementary
12 Table 3). In multivariable Cox regression analysis, the results of seven out of ten imputations
13 showed that adjuvant therapy was an independent prognostic factor ($p<0.05$ by stepwise
14 regression, the p -value cutoffs of entry and stay were both set to 0.1), and so did the combining
15 inference result ($p=0.0406$ by full model) (Supplementary Table 3).

17 In Table 2, we reported the result of one from the ten imputations that has similar result with the
18 combining inference of the ten imputations.

20 **ROC analysis**

21 The patients could be classified to “good” or “poor” prognostic groups according to the survival
22 time being longer than 5 years or not. We compared the survival prediction abilities for training
23 and test set among 3 factors: TNM stage (I-II vs III), the 3-lncRNA signature (low-risk vs high risk),
24 and the combination of the two factors.

25 Next, we constructed prognostic score models for the two factors and the combined model by
26 following the method of Liu N et al.[8] In the prognostic score models, the coefficients of low-risk
27 in the signature and I-II stages in TNM were set to 1, and the coefficients of low-risk in the
28 signature and I-II stages in TNM were set to the hazard ratio in univariable Cox regression. A
29 cumulative risk score was calculated for each patient in training and test set and was used to
30 perform receiver operating characteristic (ROC) analysis. In the comparison of area under the
31 ROC (AUROC) among the 3 models, the bootstrap test was used with 10,000 trials.

34 **References**

- 35 1 Patterson TA, Lobenhofer EK, Fulmer-Smentek SB, Collins PJ, Chu TM, Bao W, *et al.* Performance
36 comparison of one-color and two-color platforms within the MicroArray Quality Control (MAQC)
37 project. *Nat Biotechnol* 2006;**24**:1140-50.
- 38 2 Cabili MN, Trapnell C, Goff L, Koziol M, Tazon-Vega B, Regev A, *et al.* Integrative annotation of
39 human large intergenic noncoding RNAs reveals global properties and specific subclasses. *Genes Dev*
40 2011;**25**:1915-27.
- 41 3 Schwarz DF, Konig IR, Ziegler A. On safari to Random Jungle: a fast implementation of Random
42 Forests for high-dimensional data. *Bioinformatics* 2010;**26**:1752-8.
- 43 4 Ibrahim JG, Chu H, Chen MH. Missing data in clinical studies: issues and methods. *J Clin Oncol*
44 2012;**30**:3297-303.

1 5 Ibrahim JG, Molenberghs G. Missing data methods in longitudinal studies: a review. *Test (Madr)*
2 2009;**18**:1-43.

3 6 Yuan YC. Multiple imputation for missing data: Concepts and new developments. *Proceedings*
4 *of the Twenty-Fifth Annual SAS Users Group International Conference, 2000*:267.

5 7 Harel O, Zhou XH. Multiple imputation: review of theory, implementation and software. *Stat Med*
6 2007;**26**:3057-77.

7 8 Liu N, Chen NY, Cui RX, Li WF, Li Y, Wei RR, *et al.* Prognostic value of a microRNA signature in
8 nasopharyngeal carcinoma: a microRNA expression analysis. *Lancet Oncol* 2012;**13**:633-41.

9
10
11