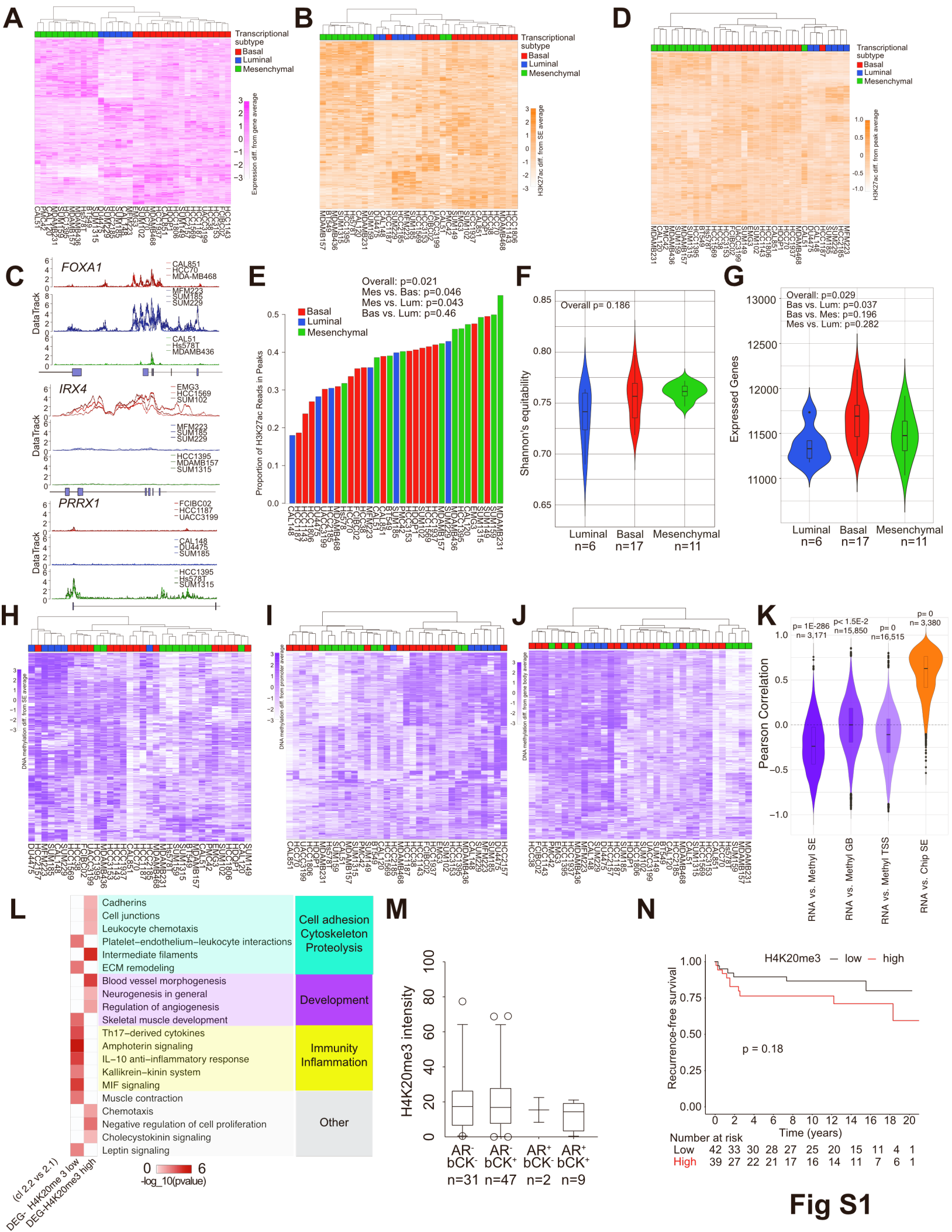


## Supplemental information

### Heterogeneity and transcriptional drivers of triple-negative breast cancer

**Bojana Jovanović, Daniel Temko, Laura E. Stevens, Marco Seehawer, Anne Fassl, Katherine Murphy, Jayati Anand, Kodie Garza, Anushree Gulvady, Xintao Qiu, Nicholas W. Harper, Veerle W. Daniels, Huang Xiao-Yun, Jennifer Y. Ge, Maša Alečković, Jason Pyrdol, Kunihiko Hinohara, Shawn B. Egri, Malvina Papanastasiou, Raga Vadhi, Alba Font-Tello, Robert Witwicki, Guillermo Peluffo, Anne Trinh, Shaokun Shu, Benedetto Diciaccio, Muhammad B. Ekram, Ashim Subedee, Zachary T. Herbert, Kai W. Wucherpfennig, Anthony G. Letai, Jacob D. Jaffe, Piotr Sicinski, Myles Brown, Deborah Dillon, Henry W. Long, Franziska Michor, and Kornelia Polyak**



**Figure S1. TNBC transcriptional and epigenetic subtypes. Related to Figure 1.**

(A) Heatmap showing clustering of 34 TNBC cell lines based on the expression of the top 20% most variable genes. Values shown are expression differences from gene average. Values are capped at +/- 3 for the purpose of visualization.

(B) Heatmap showing clustering of 33 TNBC cell lines based on H3K27ac signal in the top 20% most variable super-enhancers. Values are capped at +/-3 for the purpose of visualization.

(C) Examples of subtype-specific super-enhancer regions.

(D) Heatmap showing clustering of 33 TNBC cell lines based on H3K27ac signal in the top 10% most variable peaks.

(E) Proportion of H3K27ac reads in peaks in individual TNBC cell lines. Overall p value ( $p=0.021$ ) from Kruskal-Wallis test, pairwise p values ( $p=0.046$  mes vs bas,  $p=0.043$  mes vs lum,  $p=0.46$  bas vs lum) from Dunn's test, adjusted using Holm's method.

(F) Violin plot of Shannon's equitability across cell lines of each TNBC type. Bottom and top hinges of inset box plots show the 25<sup>th</sup> and 75<sup>th</sup> percentiles. Upper whisker extends from the upper hinge to the highest value that is no further than 1.5 times the interquartile range (IQR) from the hinge. Lower whisker extends from the lower hinge to the lowest value no further than 1.5 times the IQR from the hinge. Any data beyond the end of the whiskers are plotted as individual points. Overall p value ( $p=0.186$ ) from Kruskal-Wallis test.

(G) Violin plot of number of expressed genes (FPKM > 1) across cell lines of each TNBC type. Bottom and top hinges of inset box plots show the 25<sup>th</sup> and 75<sup>th</sup> percentiles. Upper whisker extends from the upper hinge to the highest value that is no further than 1.5 times the interquartile range (IQR) from the hinge. Lower whisker extends from the lower hinge to the lowest value no further than 1.5 times the IQR from the hinge. Any data beyond the end of the whiskers are plotted as individual points. Overall p value ( $p=0.029$ ) from Kruskal-Wallis test, pairwise p values ( $p=0.037$  bas vs lum,  $p=0.196$  bas vs mes,  $p=0.282$  mes vs lum,) from Dunn's test, adjusted using Holm's method.

(H) Heatmap showing clustering of 34 TNBC cell lines based on DNA methylation levels in the top 20% most variable super-enhancers. Values shown are methylation differences from super-enhancer average. Values are capped at +/-3 for the purpose of visualization.

(I) Heatmap showing clustering of 34 TNBC cell lines based on DNA methylation levels in top 20% most variable promoters. Values shown are methylation differences from promoter average. Values are capped at +/- 3 for the purposes of visualization.

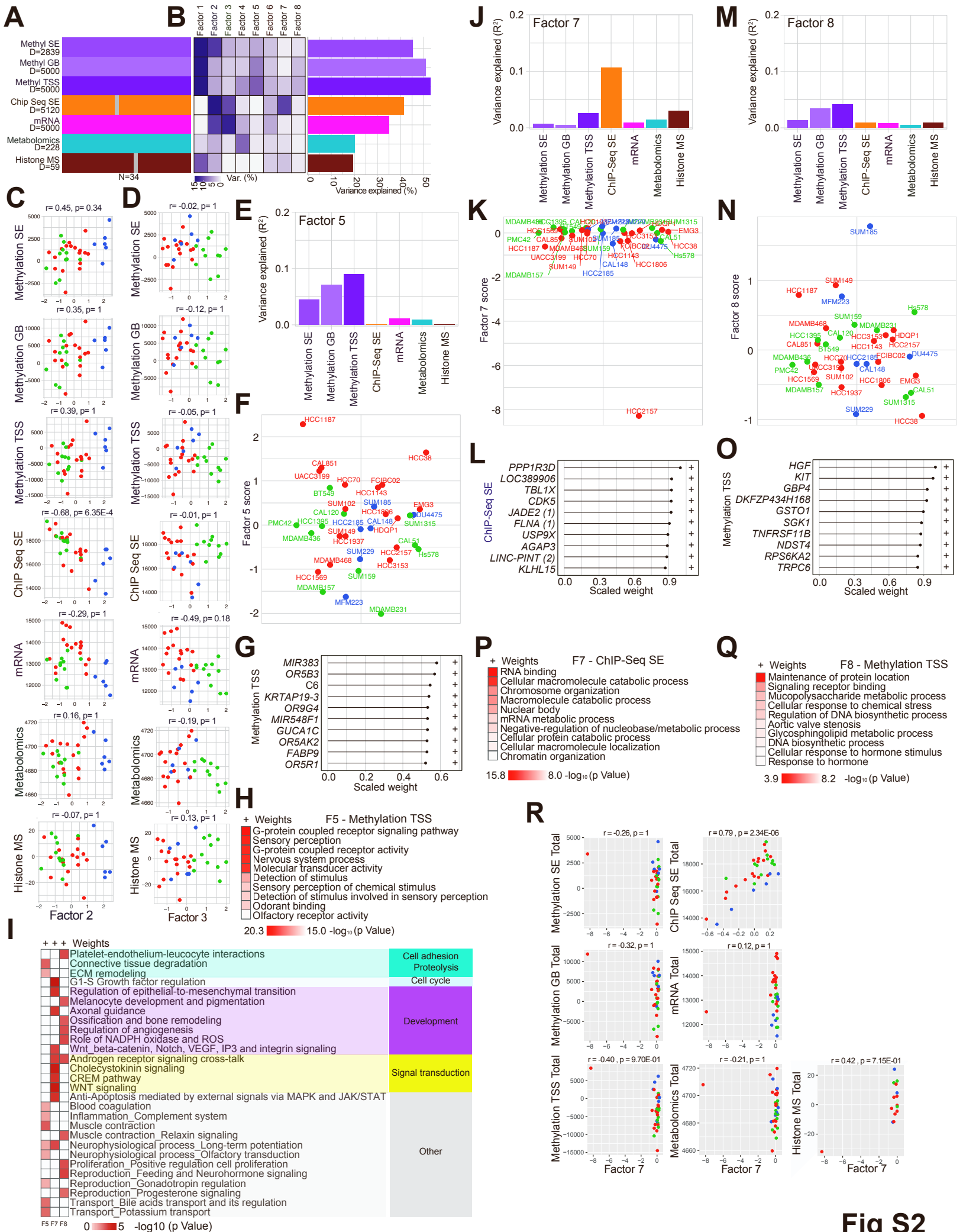
(J) Heatmap showing clustering of 34 TNBC cell lines based on DNA methylation levels in the top 20% most variable gene bodies. Values shown are methylation differences from gene body average. Values are capped at +/- 3 for the purposes of visualization.

(K) Violin plot of Pearson correlations between gene expression level, and DNA methylation and H3K27ac levels in regulatory elements for individual genes. GB: Gene body, SE: Super-enhancer, TSS: Gene promoter. Mann-Whitney U test p values are shown. DNA methylation in SEs median correlation across genes = -0.239. H3K27ac in SEs median correlation across genes = 0.626.

(L) Heatmap depicting Metacore network enrichment in differentially expressed genes (DEG) between cell lines in histone mass spectrometry cluster 2.2 versus cell lines in histone mass spectrometry cluster 2.1 from Figure 1I. DGE low/high refers to H4K20me3 levels in the two clusters.

(M) Box plot demonstrating H4K20me3 signal intensity across TNBC patient tumors positive or negative for AR and/or bCK. No significant differences were observed among the groups (ANOVA test  $p=0.5111$ ). Bottom and top hinges of inset box plots show the 25<sup>th</sup> and 75<sup>th</sup> percentiles. Upper whisker extends from the upper hinge to the highest value that is no further than 1.5 times the interquartile range (IQR) from the hinge. Lower whisker extends from the lower hinge to the lowest value no further than 1.5 times the IQR from the hinge. Any data beyond the end of the whiskers are plotted as individual points.

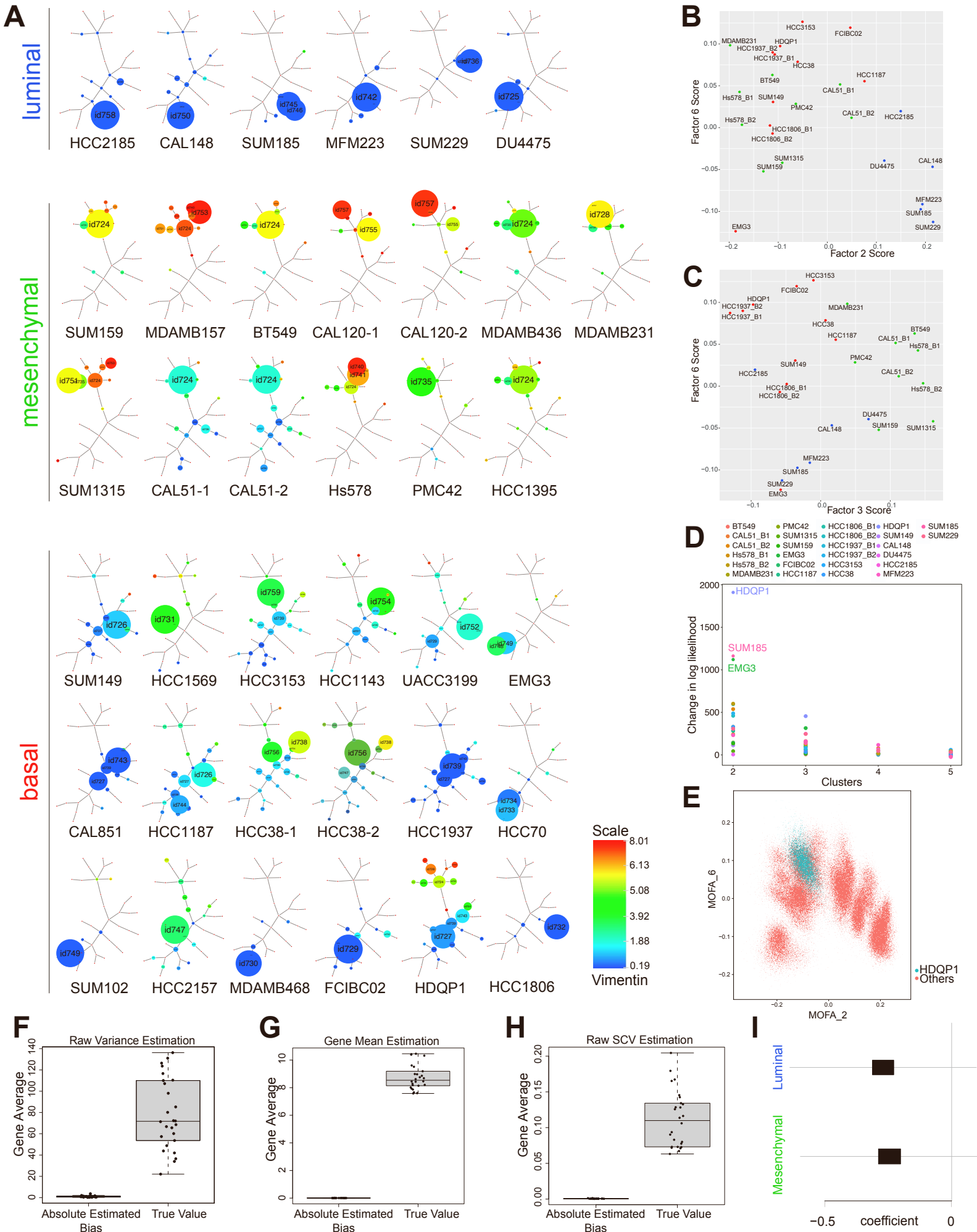
(N) Kaplan-Meier plot depicting differences in recurrence-free survival (RSF) of TNBC patients with low or high H4K20me3 levels based on median-split of H4K20me3 immunofluorescence-based intensity. Number of patients at risk over time are depicted under graph ( $p=0.18$ , Cox regression).



**Fig S2**

**Figure S2. Multi-omics factor analysis (MOFA) of TNBC. Related to Figure 2.**

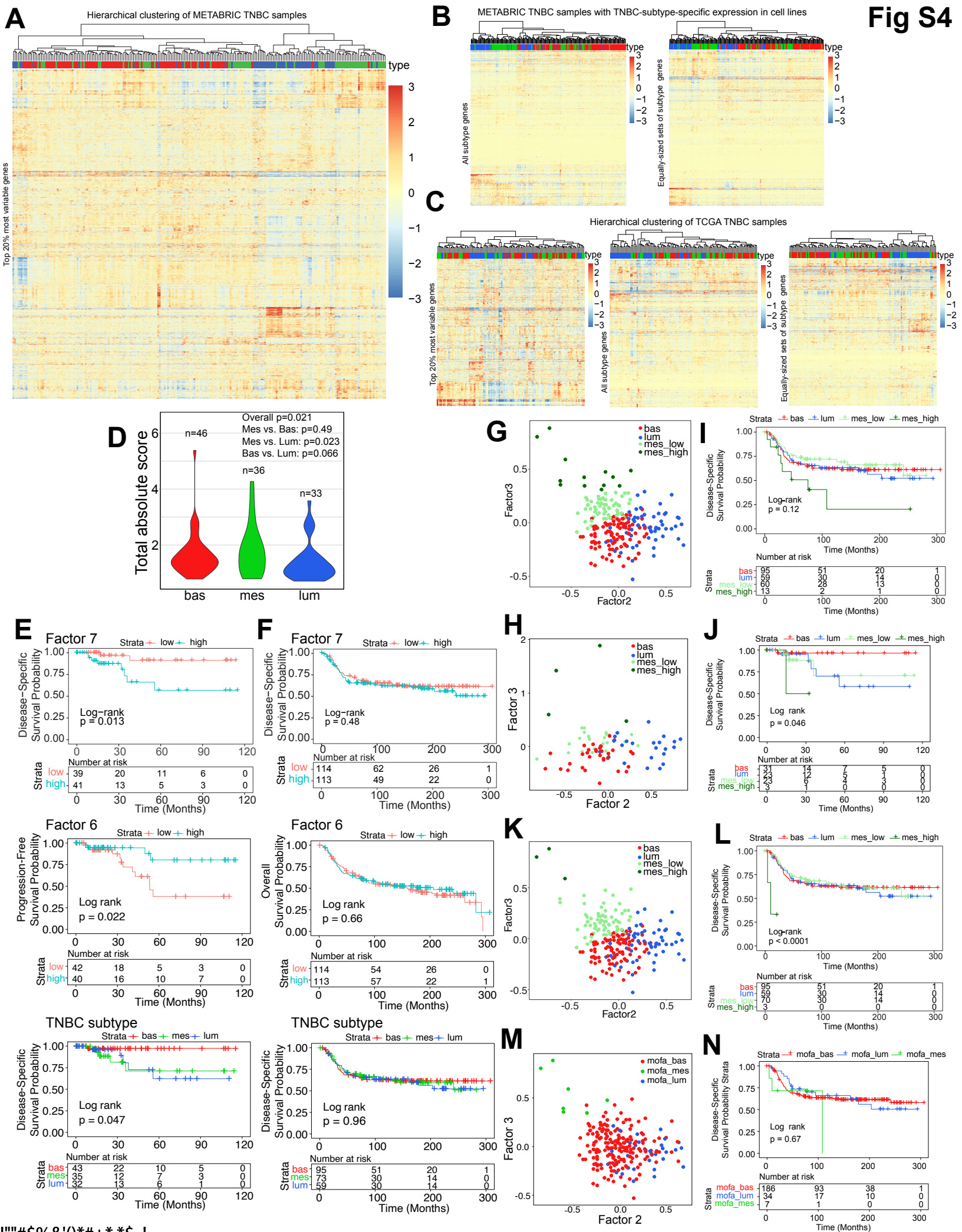
- (A) Samples and features across each dataset used for MOFA analysis.
- (B) Heatmap of variance explained in each dataset (rows) by each MOFA factor (columns). Bar-graph shows total variance explained in each dataset considering all MOFA factors in the full model.
- (C) Scatter plots of total signal in each dataset against Factor 2 scores. Holm-adjusted Pearson correlation test p values are shown.
- (D) Scatter plots of total signal in each dataset against Factor 3 scores. Holm-adjusted Pearson correlation test p values are shown.
- (E) Bar graph showing variance explained by Factor 5 in each dataset.
- (F) Bee swarm plot showing Factor 5 scores in TNBC cell lines. Points are jittered along the horizontal axis for the purpose of visualization.
- (G) Scaled Factor 5 weights for the promoter methylation features with the largest positive weights for this factor. Scaled weights for each factor in each dataset are derived from the weights for that factor in that dataset by linearly rescaling the values to lie between -1 and 1 (See Methods).
- (H) Heatmap for ontology gene set feature set enrichment analysis based on features in the promoter DNA methylation dataset with positive Factor 5 (F5) weights.
- (I) Heatmap depicting Metacore network enrichment for promoter DNA methylation features with positive Factor 5 weights (F5), H3K27ac SE features with positive Factor 7 weights (F7) and promoter methylation features with positive Factor 8 weights (F8).
- (J) Bar graph showing variance explained by Factor 7 in each dataset.
- (K) Bee swarm plot showing Factor 7 scores in TNBC cell lines. Points are jittered along the horizontal axis for the purpose of visualization.
- (L) Scaled Factor 7 weights for the H3K27ac SE features with the largest positive weights for this factor. Scaled weights for each factor in each dataset are derived from the weights for that factor in that dataset by linearly rescaling the values to lie between -1 and 1 (See Methods).
- (M) Bar graph showing variance explained by Factor 8 in each dataset.
- (N) Bee swarm plot showing Factor 8 scores in TNBC cell lines. Points are jittered along the horizontal axis for the purpose of visualization.
- (O) Scaled Factor 8 weights for the promoter methylation features with the largest positive weights for this factor. Scaled weights for each factor in each dataset are derived from the weights for that factor in that dataset by linearly rescaling the values to lie between -1 and 1 (See Methods).
- (P) Heatmap for ontology gene set feature enrichment analysis based on features in the H3K27ac SE dataset with positive Factor 7 (F7) weights.
- (Q) Heatmap for ontology gene set feature enrichment analysis based on features in the promoter methylation dataset with positive Factor 8 (F8) weights.
- (R) Scatter plots of total signal in each dataset against Factor 7 scores. Holm-adjusted Pearson correlation test p values are shown.



**Fig S3**

**Figure S3. Single cell analyses of TNBC lines. Related to Figure 4.**

- (A) CyTOF-derived x-shift clustering plots for cell line samples from each TNBC type. Plots are colored by Vimentin expression in single cells from each cluster. Identifiers within plots indicate unique clusters. Three cell lines (CAL51, CAL120, and HCC38) were run in duplicates across the two batches; sample names ending in “-1” and “-2” indicate CY1 and CY2 replicates of these three cell lines, respectively.
- (B) Average MOFA Factor 2 and Factor 6 scores of single cells from each sample (see Methods)
- (C) Average MOFA Factor 3 and Factor 6 scores of single cells from each sample (see Methods)
- (D) For each sample, the change in log likelihood compared to a model with one fewer cluster is shown for models with 2, 3, 4, or 5 clusters.
- (E) Inferred MOFA Factor 2 and MOFA Factor 6 scores for HDQP1 single cells (blue) and all other single cells (red).
- (F) Average absolute estimated bias for raw variance estimation across genes for parameter settings corresponding to each single-cell sample (left) and average ground truth raw variance values for each single-cell sample (right). Bias was estimated as the average difference between the ground truth value and the estimated value across 1,000 simulated datasets.
- (G) Average absolute estimated bias for gene mean expression estimation across genes for parameter settings corresponding to each single-cell sample (left) and average ground truth gene mean expression value across genes for each single-cell sample (right). Bias was estimated as the average difference between the ground truth value and the estimated value across 1,000 simulated datasets.
- (H) Average absolute estimated bias for raw SCV estimation across genes for parameter settings corresponding to each single-cell sample (left) and average ground truth raw SCV value across genes for each single-cell sample (right). Bias was estimated as the average difference between the ground truth value and the estimated value across 1,000 simulated datasets.
- (I) Confidence interval for estimated luminal and mesenchymal coefficients in a linear mixed effects model of the logarithm of raw transcriptional variance (see Methods). The basal subtype is the base level.



!""#\$%&'()\*#+\*,,\$ !



**Figure S4. Clinical relevance of TNBC subtype heterogeneity. Related to Figure 1.**

(A) Results of unbiased hierarchical clustering of METABRIC TNBC samples based on expression data, using the top 20% most variable genes. Annotation track represents assigned subtype using TNBC subtype signatures defined from cell lines. Values are capped at +/-3 for the purpose of visualization.

(B) Hierarchical clustering of METABRIC TNBC samples using all genes with TNBC-subtype-specific expression in cell lines (See Table S1) (left), and equally sized sets of genes with subtype-specific expression in cell lines (right). Values are capped at +/-3 for the purpose of visualization.

(C) Heatmaps showing clustering of TCGA TNBC sample using top 20% most variable genes (left), all genes with TNBC-subtype-specific expression in cell lines (See Table S1) (middle), and equally sized sets with subtype-specific expression in cell lines (right).

(D) Total absolute CIBERSORTx scores for TCGA samples assigned to each of the three TNBC types. Overall  $p=0.021$ , Kruskal-Wallis test; pairwise  $p$  values ( $p=0.49$  mes vs bas,  $p=0.023$  mes vs lum,  $p=0.066$  bas vs lum from Dunn's test, adjusted using Holm's method).

(E) Kaplan-Meier survival curves of patients in the TCGA cohort split by median Factor 7 score (top), median Factor 6 score (middle), and assigned TNBC subtype (bottom).

(F) Kaplan-Meier survival curves of patients in the METABRIC cohort split by median Factor 7 score (top), median Factor 6 score (middle), and assigned TNBC subtype (bottom).

(G) Scatter plot of inferred Factor 2 and 3 scores for METABRIC TNBC samples, colored by refined TNBC subtype.

(H) Scatter plot of inferred Factor 2 and 3 scores for TCGA TNBC samples, colored by refined TNBC subtype.

(I) METABRIC-based Kaplan-Meier survival curves for samples split by refined TNBC subtype.

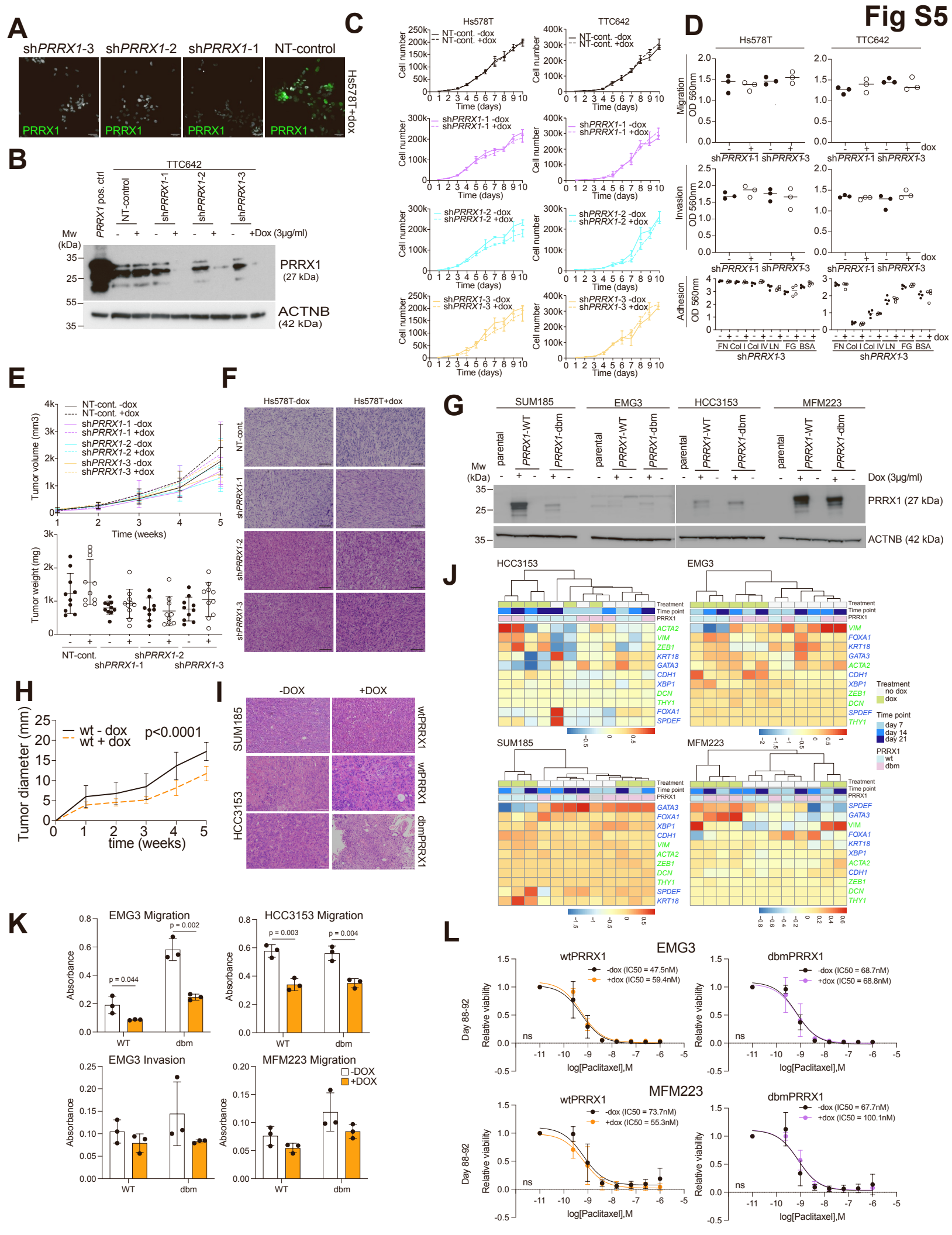
(J) TCGA-based Kaplan-Meier survival curves for samples split by refined TNBC subtype.

(K) Scatter plot of inferred Factor 2 and 3 values for METABRIC TNBC samples, colored by refined TNBC subtype, using an alternative clustering-based approach to define the refined subtypes.

(L) METABRIC-based Kaplan-Meier survival curves for samples split by refined TNBC subtype, using an alternative clustering-based approach to define the refined subtypes.

(M) Scatter plot of inferred Factor 2 and 3 values for METABRIC TNBC samples, colored by TNBC subtype using an alternative clustering-based approach to define the TNBC subtypes.

(N) METABRIC-based Kaplan-Meier survival curves for samples split by TNBC subtype, using an alternative clustering-based approach to define the TNBC subtypes.



**Figure S5. Characterization of PRRX1. Related to Figure 5.**

(A) Immunofluorescence analyses of PRRX1 protein expression in Hs578T cells, respectively expressing three independent TET-doxycycline inducible *PRRX1*-targeting shRNAs as well as non-targeting control. Scale bar corresponds to 50 $\mu$ m.

(B) Immunoblot analyses of PRRX1 protein expression in TTC642 cells expressing three independent TET-doxycycline inducible *PRRX1*-targeting shRNAs as well as non-targeting control.

(C) Proliferation of Hs578T and TTC642 cell lines after downregulation of *PRRX1* by TET-inducible shRNAs. Error bars represent SEM; n=3 replicates.

(D) Migration, invasion, and cell adhesion assays in Hs578T and TTC642 cell lines after downregulation of *PRRX1* by TET-inducible shRNAs. Error bars represent SD, n=3 replicates.

(E) Tumor weights and volumes of xenografts derived from Hs578T cell line expressing TET-inducible shRNAs with and without doxycycline in diet. Error bars represent SEM, n=10/group.

(F) Hematoxylin and eosin-stained images of xenografts of Hs578T TNBC line expressing TET-inducible sh*PRRX1* from mice with and without doxycycline in diet. Scale bar corresponds to 50 $\mu$ m.

(G) Immunoblot analysis for PRRX1 protein expression in basal (HCC3153, EMG3) and luminal (SUM185, MFM223) lines expressing doxycycline inducible wt and dbm PRRX1 protein.

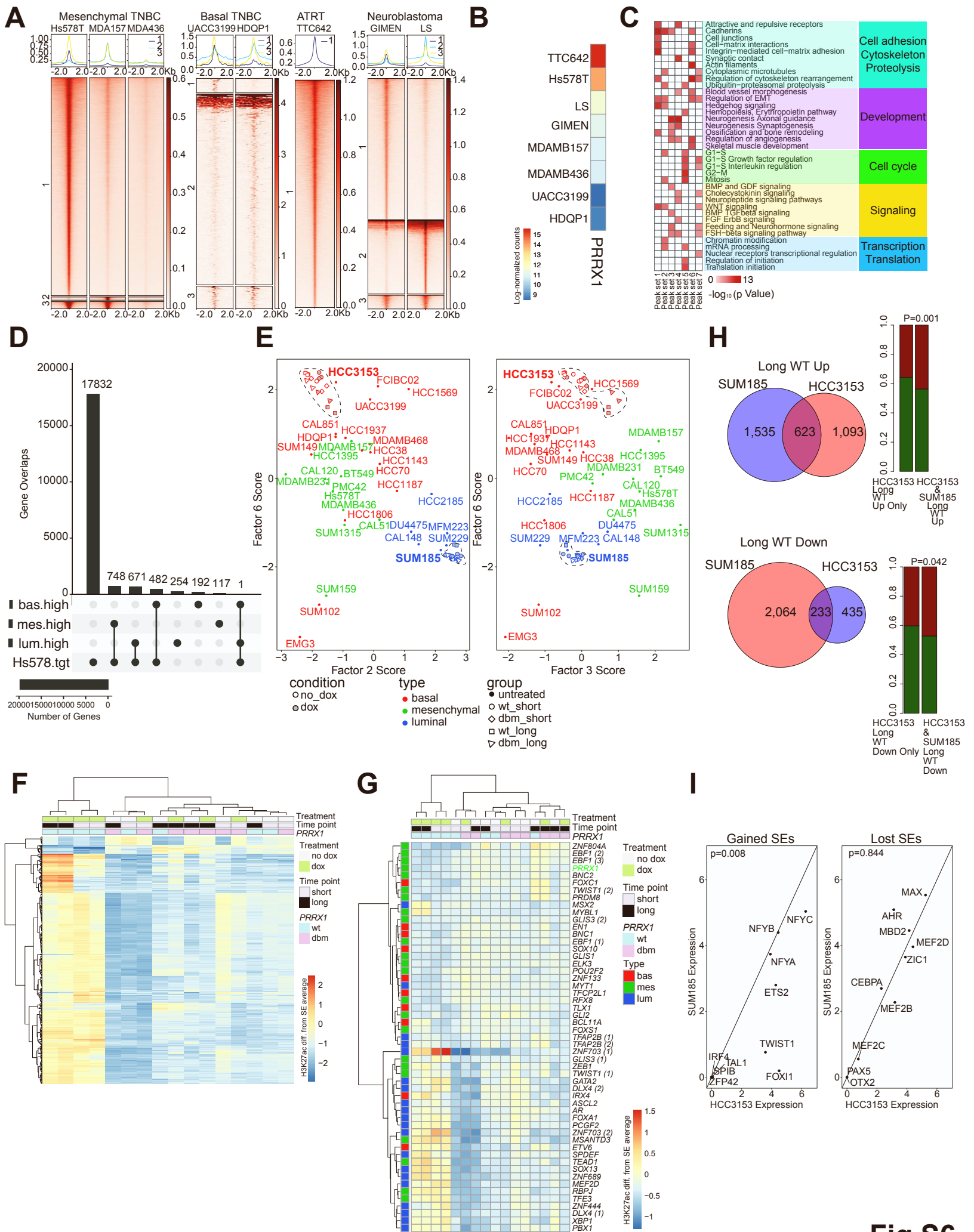
(H) Plot depicting change in tumor diameter over time of xenografts derived from SUM185 cell line expressing wt PRRX1. Each data bar represents the mean tumor diameter of 10 tumors; error bar represents SEM (P <0.0001 for a two-tailed Student's t-test).

(I) Hematoxylin and eosin-stained images of xenografts of HCC3153 and SUM185 cell line exogenously expressing WT or MUT *PRRX1* from mice with and without doxycycline in diet. Scale bar 100 $\mu$ m.

(J) Heat maps depicting clustering of the samples based on the expression of known luminal (blue) and mesenchymal (green) genes.

(K) Migration and invasion assays in the indicated cell lines after induction of wt or dbm *PRRX1* by doxycycline (DOX). Error bars represent SEM (EMG3: p=0.044 (wt) and p=0.002 (dbm); HCC3153: p=0.003 (wt) and p=0.004 (dmb) values for a two-tailed Student's t-test), n=3 replicates.

(L) Plots depicting viable cell numbers of EMG3 and MFM223 cells following paclitaxel treatment and induction of wt or dbm PRRX1 expression by doxycycline for the indicated days. Error bars represent SEM (all p values not significant for a nonlinear fit test), n=3 replicates.



**Fig S6**

**Figure S6. PRRX1-mediated expression and chromatin changes.**

(A) PRRX1 ChIP-seq peaks across mesenchymal and basal TNBC, rhabdoid (ATRT), and neuroblastoma cell lines. All combinations represented in 1% of peaks in each type of cell line are shown (Mesenchymal plot: 1: Hs578T (80,757 peaks), 2: MDAMB157 (1,188 peaks), 3: Hs578T and MDAMB157 (2,617 peaks), Basal plot: 1: UACC3199 (29 peaks), 2: HDQP1 (391 peaks), 3: UACC3199 and HDQP1 (46 peaks), Rhabdoid plot: 1: TTC642 (46,296 peaks), neuroblastoma plot: 1: GIMEN (16,187 peaks), 2: LS (8,437 peaks), 3: GIMEN and LS (1,566 peaks)). The color scheme represents PRRX1 ChIP-seq signal score +/- 2,000 bp from the peak center in each cell line. For each cell line data is shown from the ChIP-seq experiment which yielded a larger number of peaks.

(B) PRRX1 mRNA levels in the cell lines used for ChIP-seq.

(C) Metacore network enrichment analysis for shared and unique PRRX1 peaks across TNBC, neuroblastoma and rhabdoid tumor cell line shown in Fig. 5C.

(D) Number of genes specifically expressed in each TNBC sub-type overlapping the PRRX1 Hs578T target set (see Methods). Top right inset bar plot shows the total size of each gene set.

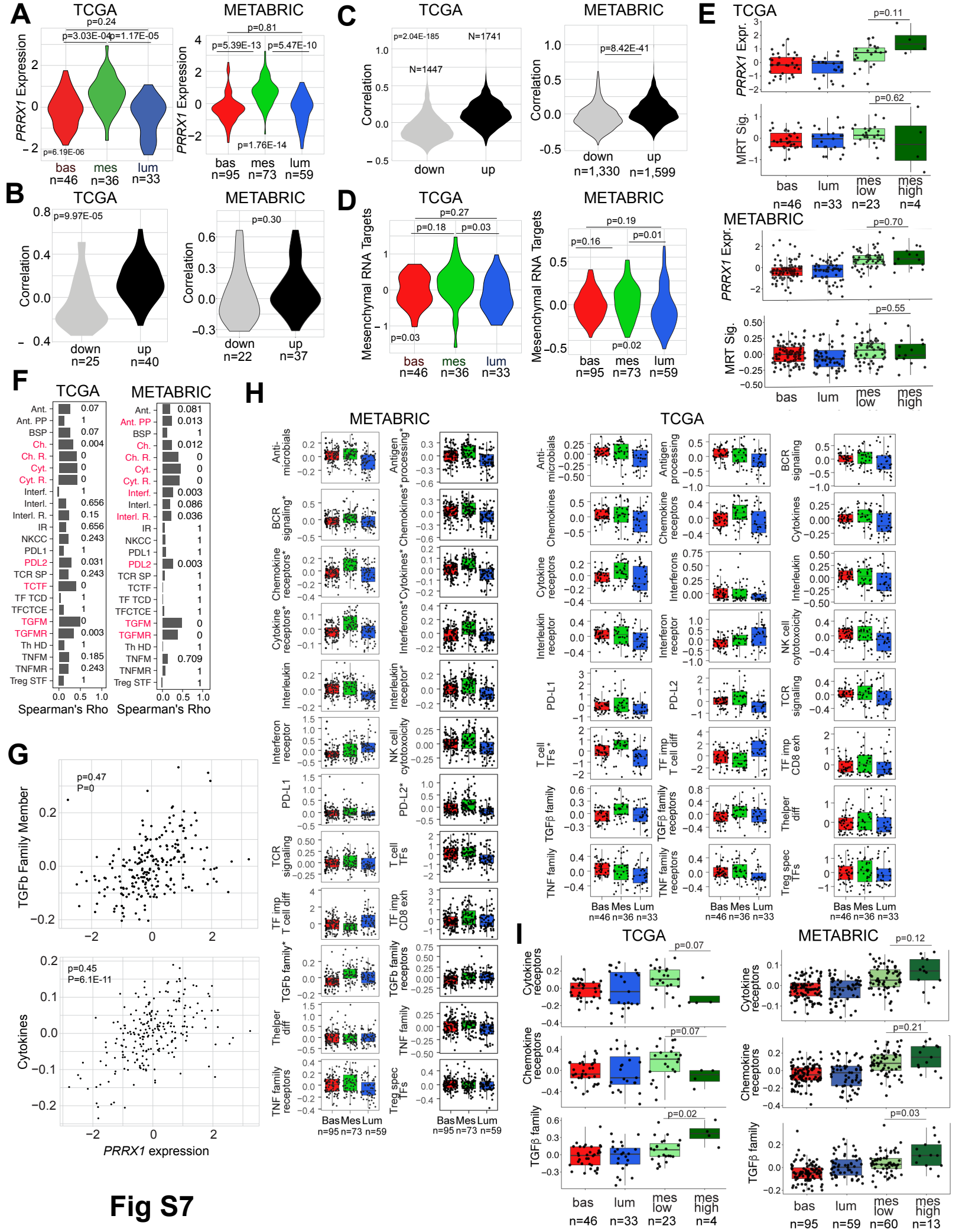
(E) MOFA Factor 2 and Factor 6 scores (left) and MOFA Factor 3 and Factor 6 scores (right) for each sample, calculated based on super-enhancer H3K27ac signal (see Methods). Red and blue outlined shapes within dotted lines represent HCC3153 and SUM185 samples from the PRRX1 over-expression H3K27ac experiment respectively.

(F) Heatmap showing clustering of SUM185 PRRX1 over-expression samples and corresponding controls based on H3K27ac signal in the top 20% most variable super-enhancers.

(G) Heatmap of super-enhancer H3K27ac signal for TNBC subtype-specific transcription factors in SUM185 PRRX1 over-expression samples and corresponding controls.

(H) Venn diagrams: Overlaps between positively and negative differentially acetylated super-enhancer regions under long-term wt PRRX1 over-expression in HCC3153 and SUM185. Bar plots: Proportion of HCC3153-unique and shared (between HCC3153 and SUM185) long-term changed superenhancers in HCC3153 that were classified as indirect (red) or direct changes (green) in HCC3153.

(I) Expression of putative PRRX1 co-binding TFs in SUM185 and HCC3153. Left and right plots show transcription factors identified from long-term HCC3153 gained and lost regions, respectively. Black diagonal lines are shown for reference and indicate equal expression in the two cell lines.  $P=0.008$  (gained SEs) and  $p=0.844$  (lost SEs) for Mann-Whitney  $U$  test.



### Figure S7. Clinical relevance of PRRX1. Related to Figure 5.

(A) Mean-centered *PRRX1* expression in TCGA and METABRIC samples assigned to each of the three TNBC types. Holm-adjusted Dunn test p values are shown. *PRRX1* expression differs among subtypes ( $p=6.2E-6$  (TCGA);  $p=1.76E-14$  (METABRIC) for Kruskal-Wallis test).

(B) Pearson correlation with *PRRX1* expression for putatively positively and negatively regulated genes in the 'Mesenchymal RNA target' (MRT) gene set for TCGA and METABRIC TNBC samples. Up, positive targets; down, negative targets. Mann-Whitney *U* test p values are shown.

(C) Pearson correlation with *PRRX1* expression for putatively positively and negatively regulated genes in the 'Hs578T RNA target' (HsRT) gene set for TCGA and METABRIC TNBC samples. Up, positive targets; down, negative targets. Mann-Whitney *U* test p values are shown.

(D) Mesenchymal RNA target signature scores in TCGA and METABRIC samples assigned to each TNBC type. Overall p value from Kruskal-Wallis test; pairwise p values from Dunn's test, adjusted using Holm's method.

(E) Mean-centered *PRRX1* expression (top panels) and Mesenchymal RNA target signature score (bottom panels) in TCGA and METABRIC samples assigned to each of the four refined TNBC types. Mann-Whitney *U* p values are shown.

(F) Correlations between immune signature scores and *PRRX1* expression in the TCGA ( $n=115$ ) and METABRIC ( $n=214$ ) cohorts. Samples with tied immune signature scores or *PRRX1* expression values were removed for the purposes of the correlation test. Holm-adjusted p value from Spearman correlation tests are shown. Signatures with significant adjusted p values ( $<0.05$ ) are colored red. Ant.: Anitmicrobials, Ant. PP: Antigen\_Processing\_and\_Presentation, BSP: BCRSignalingPathway, Ch.: Chemokines, Ch. R.: Chemokine\_Receptors, Cyt.: Cytokines, Cyt. R.: Cytokine\_Receptors, Interf.: Interferons, Interl.: Interleukins, Interl. R.: Interleukins\_Receptor, IR: Interferon\_Receptor, NKCC: NaturalKiller\_Cell\_Cytotoxicity, TCR SP: TCRsignalingPathway, TCTF: T\_cell\_TF, TF TCD: TF\_important\_for\_T\_cellll\_differentiation, TFCTCE: TF\_important\_for\_CD8\_T\_cell\_exhaustion, TGFM: TGFb\_Family\_Member, TGFMR: TGFb\_Family\_Member\_Receptor, Th HD: Th\_helper\_diff, TNFM: TNF\_Family\_Members, TNFMR: TNF\_Family\_Members\_Receptors, Treg STF: Treg\_specific\_TF.

(G) Scatter plots of TGF- $\beta$  Family Member signature score (top) and Cytokines signature score (bottom) and mean-centered *PRRX1* expression in the METABRIC cohort.  $\rho$ : Spearman's rho, P: Holm-adjusted Spearman correlation test p value. Samples with tied immune signature expression values or *PRRX1* expression values were removed for the purposes of the correlation test. The trend line is omitted since the relationship between *PRRX1* expression and signature expression is not assumed to be linear.  $n=227$  (scatter plot),  $n=214$  (correlation test).

(H) Immune signature scores in TCGA and METABRIC samples assigned to each TNBC type. \*: Significant differences in scores between subtypes (Holm-adjusted Kruskal-Wallis test p value  $< 0.05$ ) with significantly higher scores in mesenchymal TNBC than luminal or basal TNBC (Dunn test p value  $< 0.05$ , both comparisons). Ant.: Anitmicrobials, Ant. PP: Antigen\_Processing\_and\_Presentation, BSP: BCRSignalingPathway, Ch.: Chemokines, Ch. R.: Chemokine\_Receptors, Cyt.: Cytokines, Cyt. R.: Cytokine\_Receptors, Interf.: Interferons, Interl.: Interleukins, Interl. R.: Interleukins\_Receptor, IR: Interferon\_Receptor, NKCC: NaturalKiller\_Cell\_Cytotoxicity, TCR SP: TCRsignalingPathway, TCTF: T\_cell\_TF, TF TCD: TF\_important\_for\_T\_cellll\_differentiation, TFCTCE: TF\_important\_for\_CD8\_T\_cell\_exhaustion, TGFM: TGFb\_Family\_Member, TGFMR: TGFb\_Family\_Member\_Receptor, Th HD: Th\_helper\_diff, TNFM: TNF\_Family\_Members, TNFMR: TNF\_Family\_Members\_Receptors, Treg STF: Treg\_specific\_T.

(I) TGF- $\beta$  Family Member, Chemokine Receptors, and Cytokine Receptors signature scores in TCGA and METABRIC samples assigned to each of the four refined TNBC types. Mann-Whitney *U* test p values are shown.