

1 **Large-scale proteome inference from unpaired single-cell transcriptomic and**
2 **proteomic data by msInfer**

3 Tianyi Zhao^{1, 2, #}, Yuzhi Sun^{3, #, *}, Renjie Liu^{2, 3, #}, Liyuan Zhang³, Chengcheng Zhang³, Yuran Jia³, Liang Cheng^{4, *},
4 Guohua Wang^{3, *}, Yadong Wang^{1, *}

5

These authors contributed equally: Tianyi Zhao, Yuzhi Sun, Renjie Liu.

* Corresponding authors: Yuzhi Sun (yuzhi@stu.hit.edu.cn), Liang Cheng (liangcheng@hrbmu.edu.cn), Guohua Wang (ghwang@hit.edu.cn), Yadong Wang (ydwang@hit.edu.cn).

1 School of Medicine and Health, Harbin Institute of Technology, Harbin, China.

2 Harbin Institute of Technology Zhengzhou Research Institute, Zhengzhou, China.

3 Faculty of Computing, Harbin Institute of Technology, Harbin, China.

4 Harbin Medical University, Harbin, China.

6 Supplementary Notes

7 **Note1: msInfer reveals a more comprehensive change across different** 8 **tumour states**

9 msInfer can be used to study the microenvironments associated with different stages of breast cancer. We used
10 the single-cell RNA expression atlas measured by Bhupinder *et al.*¹, which includes samples from human breast
11 tissue in the normal, preneoplastic, and tumour states (a total of more than 300,000 cells from 8 normal samples, 4
12 BRCA1+ carrier tissue samples and 8 triple-negative breast cancer samples). We utilized msInfer, referencing Leduc
13 2022, Specht, and Khan which are from breast tissue to infer the corresponding proteomic data for the Bhupinder
14 dataset. Using transcriptomic data combined with inferred proteomics data can better perform t-SNE dimensionality
15 reduction and clustering of cells (Supplementary Fig. S5). Subsequently, we conducted a joint analysis of
16 transcriptomic and proteomic data across the different microenvironments associated with the normal, preneoplastic,
17 and cancerous states. Here, we focused on mesenchymal and epithelial cells closely associated with the breast cancer
18 microenvironment.

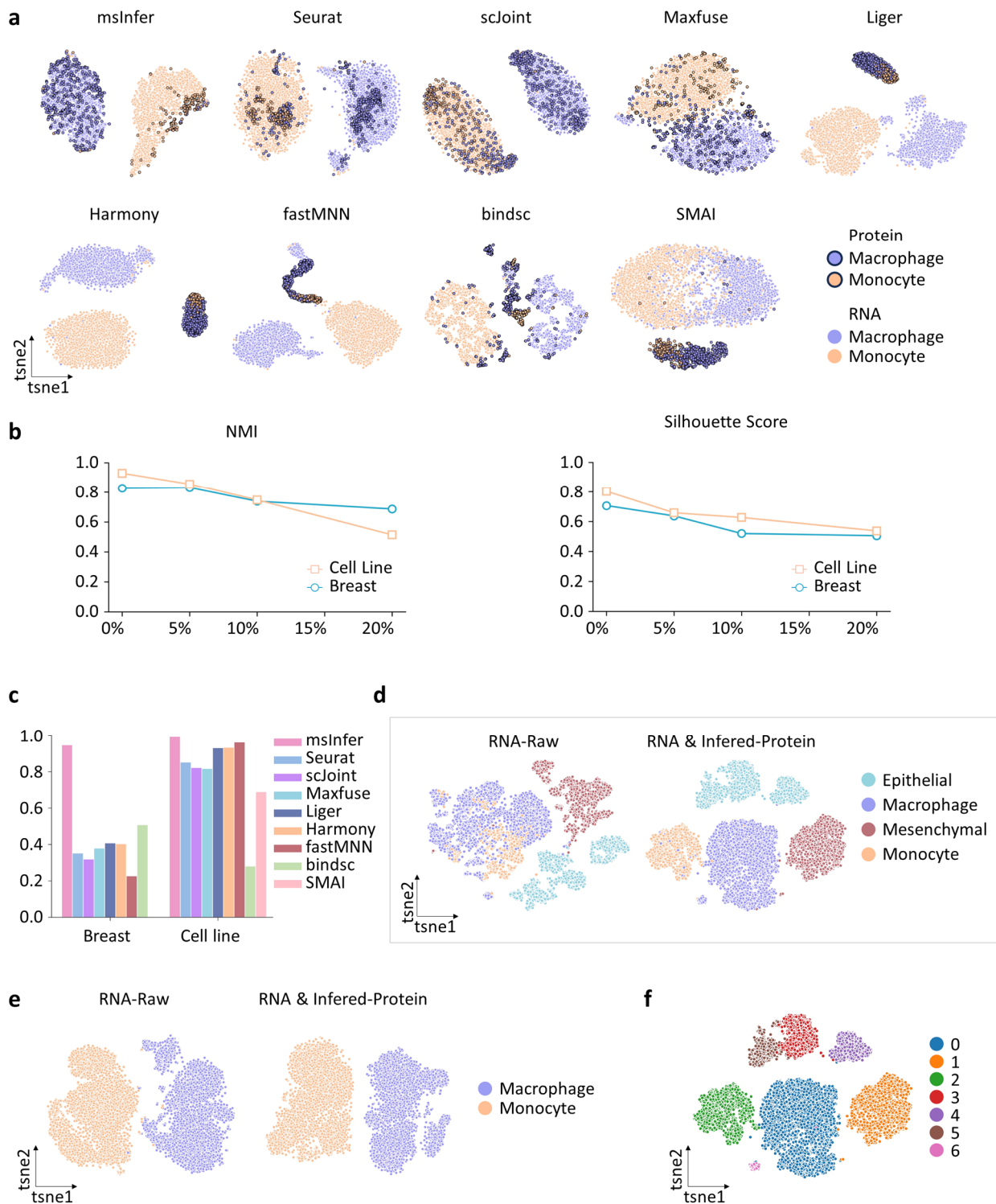
19 We conducted differential expression analysis (N-B1 and B1-TN) on the transcriptomic and proteomic data,
20 respectively. We then summarized the gene sets on the basis of upregulation and downregulation as well as cell type,
21 and created an UpSet plot. As shown in Supplementary Fig. S6a, in epithelial cells, the ATOX1, CSTB, and NUCKS1
22 genes were significantly upregulated at both the RNA and protein levels. Previous studies have shown that the ATOX1
23 gene can alter cell migration capabilities, which is closely related to the progression of breast cancer²⁻⁴. The CSTB
24 gene has been identified as crucial for the proteolytic cascade associated with tumour progression, including tumour
25 growth, invasion, and metastasis⁵. Additionally, the NUCKS1 gene has been confirmed to participate in tumour
26 suppression^{6,7}. In epithelial cells, 27 genes were significantly downregulated at both the RNA and protein levels,
27 with functions related primarily to metabolic regulation, survival, and anti-apoptosis⁸⁻¹². Similarly, as shown in
28 Supplementary Fig. S6b, 25 genes in mesenchymal cells were downregulated at both stages, with functions also
29 related to metabolic regulation and anti-apoptosis¹³⁻¹⁷. However, no genes exhibiting consistent upregulation were
30 identified in mesenchymal cells. Thus, msInfer efficiently analyses breast cancer progression at the cellular level
31 from a proteomic perspective.

32 We then analysed the RNAs and proteins that exhibited significant expression changes at both the N-B1 and B1-
33 TN stages within the same cell type. Specifically, we first selected the top 200 genes (100 upregulated and 100
34 downregulated, sorted by logarithmic fold change) that presented significant changes (p-adjusted less than 0.05) at
35 each stage. We then took the overlap of the two gene sets and created a Sankey plot to illustrate the continuous
36 changes in the expression of the RNAs and proteins at each stage. A substantial portion of RNAs in epithelial and
37 mesenchymal cells exhibited a trend of downregulation followed by upregulation, while some showed upregulation
38 followed by no change (Supplementary Fig. S6c). However, proteins predominantly displayed a continuous
39 upregulation trend in epithelial cells, whereas a continuous downregulation trend was more evident in mesenchymal
40 cells. As executors of cellular functions, proteins exhibit a more stable trend of change than do RNAs. Notably, during
41 the two stages (N-B1 and B1-TN), we identified a gene, HMGA1, in epithelial cells that exhibited a persistent decline
42 in RNA expression while protein abundance continuously increased. HMGA1 is closely related to the migration,
43 invasion, and metastasis of TNBC cells¹⁸, and its high protein expression can serve as a biomarker to predict
44 metastasis incidence¹⁹, histological grade, clinical stage²⁰, and survival time²¹. The protein abundance inferred by
45 msInfer aligns with that reported in previous studies, indicating that high expression of the HMGA1 protein is highly

46 correlated with breast cancer progression, while scRNA-seq showed contrasting results. This highlights the
47 importance of msInfer in complementing multiomics information to comprehensively investigate the role of the
48 tumour microenvironment.

49 We summarized the genes whose expression was continuously upregulated or downregulated and conducted
50 functional enrichment analysis. As shown in Supplementary Fig. S6d, we identified many pathways with similar
51 functions in both cell types, such as extracellular exosomes, which play important roles in cancer development and
52 metastasis. We also discovered functions unique to specific cell types; for example, apoptosis-related pathways were
53 found in epithelial cells, whereas mitochondrial function-related pathways, such as ATP:ADP antiporter activity, were
54 identified in mesenchymal cells. These pathways are closely related to the development and environment of cancer
55 cells.

Supplementary Figures



57

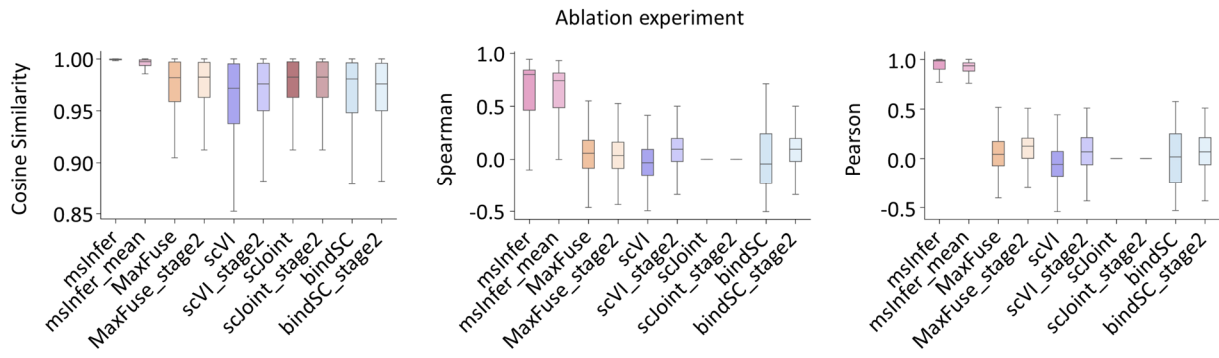
58

59

60

Supplementary Fig. S1 Validation on scRNA-seq and Mass Spectrometry-based Single-cell Proteomics Dataset. a. t-SNE visualization of transcriptomic and proteomic cells distributions with integrated embedding features on Cell Line task. **b.** Clustering performance (NMI and Silhouette Score) under disturbance. **c.** Accuracy of cell type matching. **d.** t-SNE visualization of cells with

61 RNA data and RNA & inferred-protein data. **e.** t-SNE visualization of cells with RNA data and RNA & inferred-protein data on Cell
62 Line task. **f.** Cell distribution map after Leiden clustering of transcriptome and proteome data from the breast dataset. derived from Wu
63 et al.⁵⁰ .
64



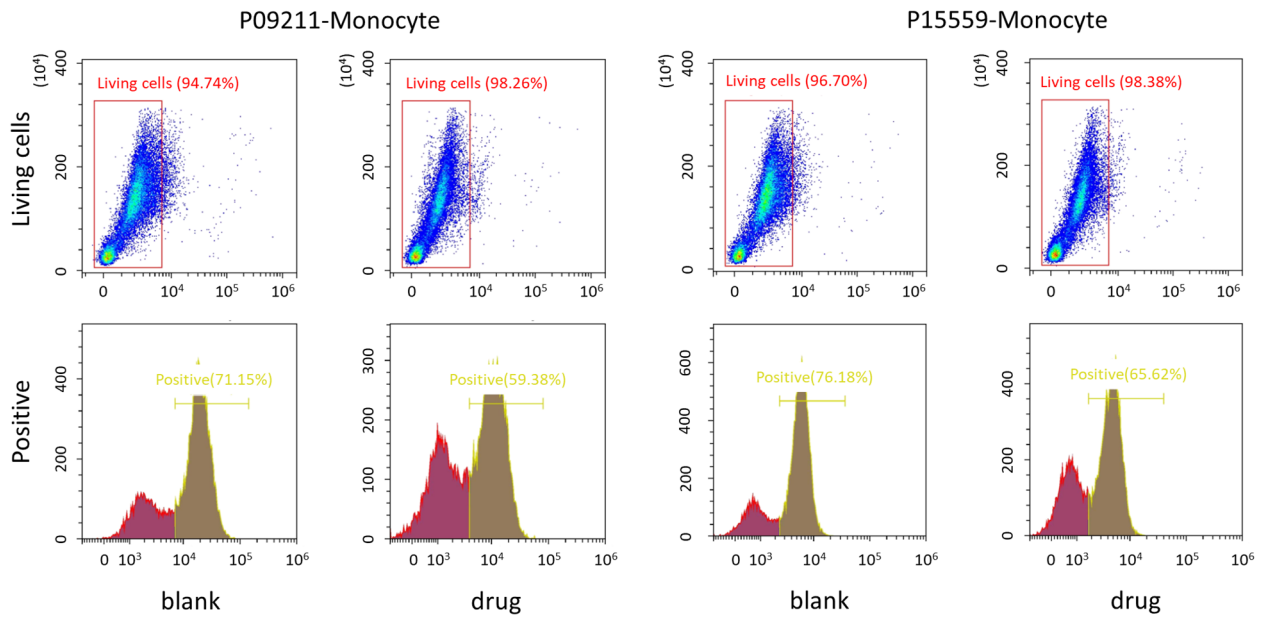
65

66

Supplementary Fig. S2 Ablation experiment on benchmark. Use mean calculation to replace stage 2 of msInfer for prediction, and use integrated methods combined with stage 2 for prediction.

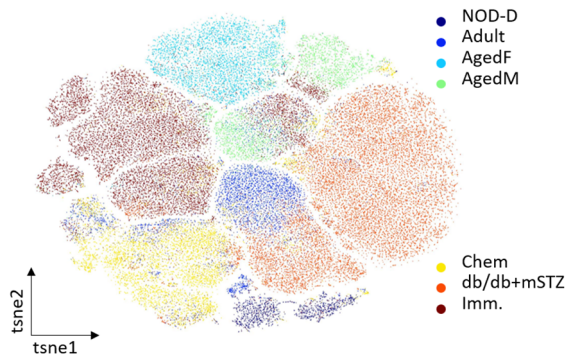
67

68



69
70
71

Supplementary Fig. S3 Changes in the levels of proteins P09211 and P15559 in monocytes before and after treatment with cisplatin measured by flow cytometry.



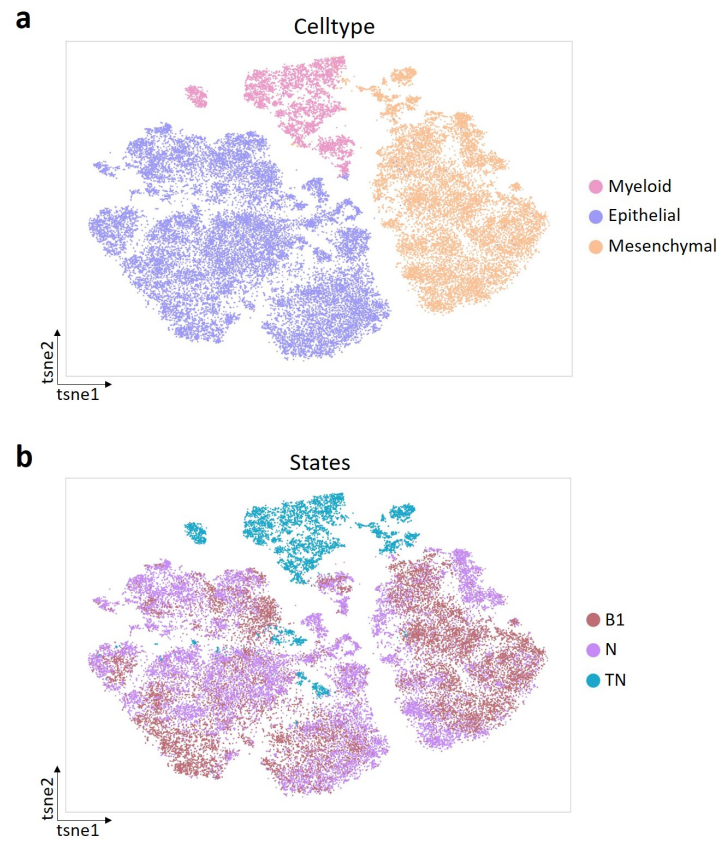
72

73

74

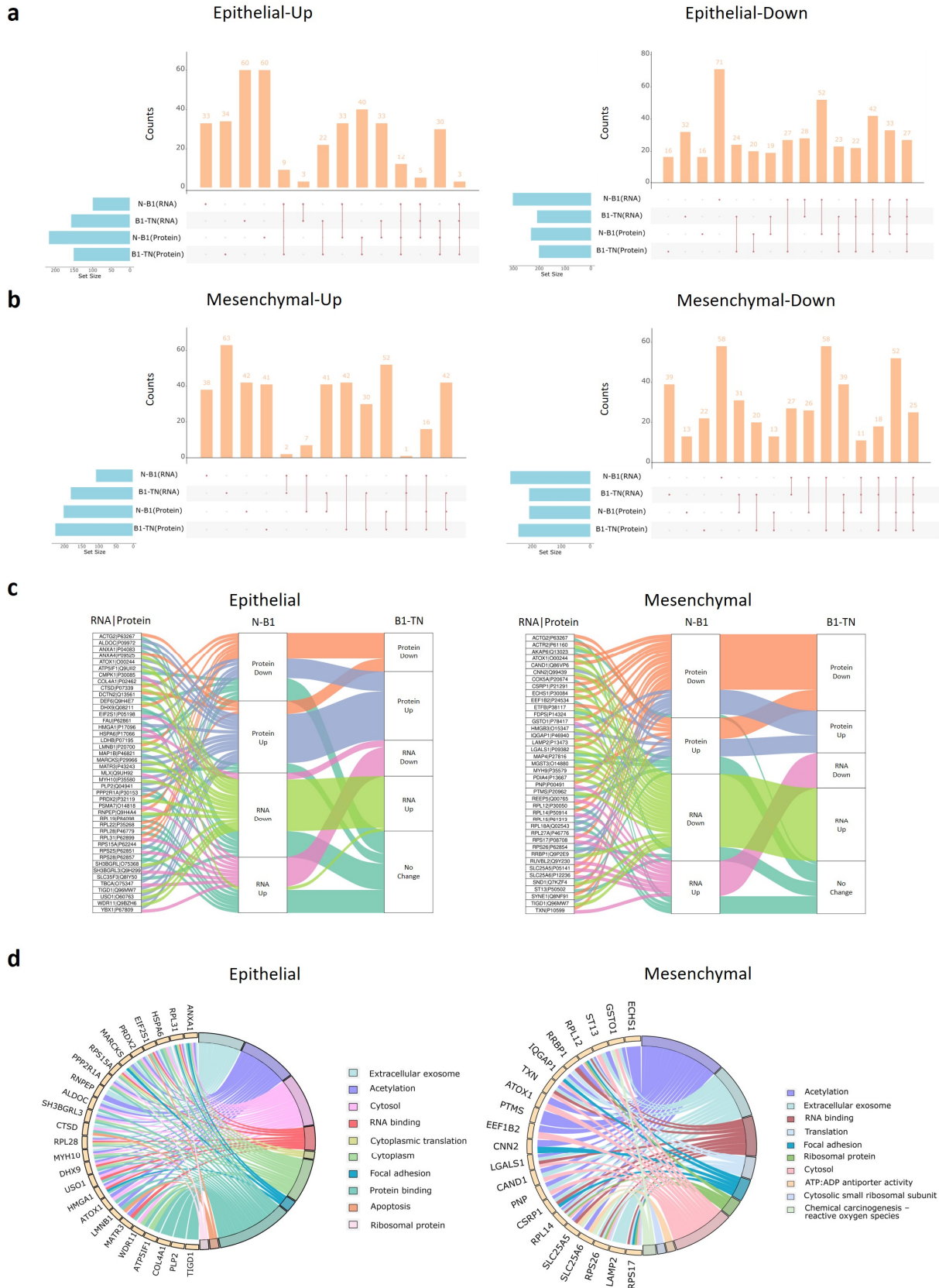
75

Supplementary Fig. S4 Scatter plot of beta atlas cells distribution only using transcriptome data



76
77
78
79

Supplementary Fig. S5 T-SNE plot of using transcriptomics and combined with proteomics predicted by msInfer in N, B1, and TN group cells.



Supplementary Fig. S6 Single-cell RNA and protein changes in different environments. a. UpSet plot showing the number of

83 differentially expressed RNAs and proteins in epithelial cells at different cancer stages. N represents the normal (healthy) state, B1
84 indicates the preneoplastic state (BRCA1+/-), and TN denotes cancer patients. For example, N-B1(RNA) indicates the differentially
85 expressed RNAs between the preneoplastic state and the normal state, with the upper bar plot showing the number of genes (33) and
86 the line connecting the points indicating the overlap of different groups. **b.** UpSet plot showing the number of differentially expressed
87 RNAs and proteins in mesenchymal cells at different cancer stages. **c.** Sankey plot illustrating the continuous changes in RNA and
88 protein levels from the normal to preneoplastic to cancerous state. The top 200 differentially expressed RNAs and proteins shared
89 between the N-B1 and B1-TN states were selected for plotting, with No change indicating no significant differences. **d.** RNAs and
90 proteins whose expression continuously changed, as shown in Fig. 5c, were selected for functional enrichment analysis. For example,
91 ALDOC was selected for enrichment analysis since its protein expression is upregulated in both N-B1 and B1-TN. The same pathways
92 are represented in the same colour in the two enrichment analysis diagrams.

93

94