1    **STACCato: Supervised Tensor Analysis tool for studying Cell-cell Communication using scRNA-seq**

2    **data across multiple samples and conditions**

3    Qile Dai[1,2], Michael P. Epstein[2*], Jingjing Yang[2*]

4        1.  Department of Biostatistics and Bioinformatics, Emory University School of Public Health, Atlanta,
5            Georgia 30322, United States of America.
6        2.  Center for Computational and Quantitative Genetics, Department of Human Genetics, Emory
7            University School of Medicine, Atlanta, Georgia 30322, United States of America.
8    *Correspondence Authors: M.P.E. (mpepste@emory.edu) and J.Y. (jingjing.yang@emory.edu)

9

10   **Abstract**

11       Research on cell-cell communication (CCC) is crucial for understanding biology and diseases. Many

12   existing CCC inference tools neglect potential confounders, such as batch and demographic variables, when

13   analyzing multi-sample, multi-condition scRNA-seq datasets. To address this significant gap, we introduce

14   STACCato, a **S**upervised **T**ensor **A**nalysis tool for studying **C**ell-cell **C**ommunication, that identifies CCC

15   events and estimates the effects of biological conditions (e.g., disease status, tissue types) on such events,

16   while adjusting for potential confounders. Application of STACCato to both simulated data and real scRNA-

17   seq data of lupus and autism studies demonstrate that incorporating sample-level variables into CCC inference

18   consistently provides more accurate estimations of disease effects and cell type activity patterns than existing

19   methods that ignore sample-level variables. A computational tool implementing the STACCato framework is

20   available on GitHub.

21   **Introduction**

22       Cell-cell communication (CCC) involves cells exchanging signals to coordinate physiological and

23   developmental functions in multicellular organisms. The study of CCC events, which involves interactions

24   between one ligand-receptor pair from one sender cell type to one receiver cell type, is important for

25   elucidating biological processes, exploring disease mechanisms, and inspiring advancements in drug

26   discovery. Using gene expression data produced by single-cell RNA sequencing (scRNA-seq) technology,

27   multiple computational tools are now available to infer CCC events[1–9].

28    Recently, high-throughput sequencing technology advancements have significantly reduced the cost of

29    scRNA-seq, allowing researchers to gather scRNA-seq data from multiple biological samples under multiple

30    biological conditions[10–13], such as disease versus healthy control samples or samples from multiple tissue

31    types. Most existing computational tools developed for CCC inference were originally designed for analyzing

32    single-sample scRNA-seq data[1–7]. When attempting to apply these tools to multi-sample multi-condition

33    scRNA-seq datasets, a three-step procedure is typically necessary. First, data from all samples within the same

34    condition are combined to create an aggregated "sample" per condition. Second, communication scores are

35    calculated for CCC events using the aggregated "samples", one per condition. Last, CCC events with

36    significantly different communication scores across conditions are identified as condition-related CCC events.

37    Another proposed strategy to handle such multi-sample multi-condition single-cell data is to use the tensor

38    decomposition technique, which has been used to extract underlying lower-dimensional patterns from high-

39    dimensional genomic data[8,9,14,15]. For example, the recently developed tool Tensor-cell2cell [11] constructs a 4-

40    dimensional communication score tensor, with 4 dimensions corresponding to samples, ligand-receptor pairs

41    sender cell types, and receiver cell types. Tensor-cell2cell applies unsupervised tensor decomposition to

42    identify underlying communication patterns, and then tests if the communication patterns are significantly

43    different across conditions.

44    An important drawback of both the three-step procedure and the Tensor-cell2cell tool for analyzing

45    multi-sample and multi-condition scRNAseq data is that they ignore important sample-level variables (such as

46    processing batch, age, gender, and ancestry) that are typically collected in such studies. These variables can

47    have substantial impacts on both biological conditions and CCC, likely confounding the identification of

48    condition-related CCC events. Neglecting these confounding variables may mask true biological associations

49    between CCC events and conditions, or, even more concerning, lead to false positive associations that could

50    result in misguided interpretations of CCC events. Therefore, the development of a CCC inference tool to

51    effectively incorporate sample-level variables and adjust for potential confounding variables in multi-sample

52    multi-condition scRNA-seq data becomes increasingly important.

2

53        To bridge this gap, we introduce the **S**upervised **T**ensor **A**nalysis tool for studying **C**ell-cell

54        **C**ommunication (STACCato), that uses multi-sample multi-condition scRNA-seq dataset to identify CCC events

55        significantly associated with conditions while adjusting for potential sample-level confounders. STACCato

56        considers the same 4-dimentional communication score tensor as the Tensor-cell2cell tool, with 4 dimensions

57        corresponding to samples, ligand-receptor pairs, sender cell types, and receiver cell types. Different from the

58        Tensor-cell2cell tool, STACCato employs supervised tensor decomposition[16] to fit a regression model that

59        considers the 4-dimensional communication score tensor as the outcome variable while treating the biological

60        conditions (e.g., disease status, time points, tissue types) and other sample-level covariates (e.g., batch and

61        demographic variables) as independent variables. Through this supervised tensor-based regression model,

62        STACCato can identify CCC events and estimate the impact of conditions on CCC events, while effectively

63        controlling for potential confounding variables.

64        In subsequent sections, we first introduce the analytical framework of STACCato. We then apply

65        STACCato to two real datasets: the Systemic Lupus Erythematosus (SLE) dataset[10,11] consisting of scRNA-seq

66        data of peripheral blood mononuclear cells (PBMC) samples from 154 SLE patients and 97 healthy controls,

67        and the Autism Spectrum Disorder (ASD) dataset[12] consisting of snRNA-seq data of prefrontal cortex (PFC)

68        samples from 13 ASD patients and 10 controls. Notably, the SLE dataset exhibits an unbalanced study design,

69        resulting in batch effects being highly confounded with the disease effect. We observed dramatic changes in

70        estimated disease effects for CCC events before and after adjusting for batch effects, leading to contrasting

71        conclusions regarding the associations between these CCC events and SLE. These findings underscore the

72        substantial impact of confounding variables on CCC inference, emphasizing the necessity of accounting for

73        confounding variables in CCC studies. We further validate these observations through a simulation study

74        considering various study designs. Finally, we conclude with a discussion.

75        **Results**

76        *STACCato framework*

77        We propose STACCato, a powerful tool that utilizes multi-sample multi-condition scRNA-seq data to

78        identify condition-related CCC events while accounting for potential confounding variables. Briefly, STACCato

79   first generates a 4D communication score tensor with four dimensions representing samples, ligand-receptor

80   pairs, sender cell types, and receiver cell types (Figure 1A-1C). Next, STACCato employs a supervised tensor

81   decomposition method that incorporates sample-level information (such as biological conditions or batches) to

82   estimate a coefficient tensor, representing the effects of sample-level variables on CCC events (Figure 1C).

83   Finally, we conduct parametric bootstrapping to assess the significance of the estimated coefficients. We

84   describe the general supervised tensor decomposition framework below and relegate the technical details to the

85   Methods section.

86   *Supervised tensor decomposition of communication score tensor*

87   With respect to an CCC event involving the interaction of ligand-receptor pair $j$ from sender cell type

88   $k$ to receiver cell type $l$, we consider the following regression model to assess the association between the CCC

89   event and the condition adjusting for other covariates,

$$y_{ijkl} = \beta_{1jkl}x_{i1} + \cdots + \beta_{qjkl}x_{iq} + \epsilon_{ijkl};$$

91   $$i = 1, \cdots I; \ j = 1, \cdots J; \ k = 1, \cdots K; \ l = 1, \cdots L; \ q = 1, \cdots Q. \quad \text{(Equation 1)}$$

92   Here, $I, J, K, L$, and $Q$ are the total number of samples, ligand-receptor pairs, sender cell types, receiver cell

93   types, and sample-level variables, respectively. In Equation 1, $y_{ijkl}$ denotes the communication score

94   representing the communication level of the CCC event involving the interaction of ligand-receptor pair $j$ from

95   sender cell type $k$ to receiver cell type $l$ in sample $i$ (see Methods for details about communication score

96   calculation); $x_{iq}$ denotes the sample-level variable $q$, such as biological condition or batch, for sample $i$; $\beta_{qjkl}$

97   denotes the effect of variable $q$ on the communication score of the CCC event involving the interaction of ligand-

98   receptor pair $j$ from sender cell type $k$ to receiver cell type $l$; and $\epsilon_{ijkl} \sim N(0, \sigma^2)$ denotes the random error that

99   follows a Gaussian distribution with mean 0 and standard deviation $\sigma$.

100   A straightforward way to estimate $\beta_{qjkl}$ is to fit a regression model with $\boldsymbol{y}_{jkl} = \left[y_{1jkl}, \cdots, y_{Ijkl}\right]^T$ as the

101   values of the dependent variable and sample-level information matrix $\boldsymbol{X} \in \mathbb{R}^{I \times Q}$ as the design matrix for

102   independent variables. The major limitation of this strategy is that it estimates $\boldsymbol{\beta}_{jkl} = \left[\beta_{1jkl}, \cdots, \beta_{Qjkl}\right]^T, j =$

103    $1, \cdots J$, $k = 1, \cdots K$, $l = 1, \cdots L$ separately for each CCC event and ignores the correlations among CCC events.

104    For example, the interactions of the same ligand-receptor pair $j$ across different sender and receiver cell types

105    are dependent, and thus $\beta_{qjkl}$ is dependent of $\beta_{qjk'l'}$ with $k \neq k'$ and $l \neq l'$. To consider such correlations

106    among CCC events, we employ a supervised tensor technique to jointly estimate $\boldsymbol{\beta}_{jkl}$ for all $j = 1, \cdots J$, $k =$

107    $1, \cdots K$, $l = 1, \cdots L$. To do so, we note that Equation 1 is equivalent to the tensor model,

108    $$\mathcal{Y} = \mathcal{B} \times_1 \boldsymbol{X} + \mathcal{E} \qquad \text{(Equation 2)}$$

109    where $\mathcal{Y} \in \mathbb{R}^{I \times J \times K \times L}$ denotes the 4-dimensional communication score tensor with dimensions of $I$ samples, $J$

110    ligand-receptor pairs, $K$ sender cell types, and $L$ receiver cell types, with the $(i, j, k, l)$ entry corresponding to

111    $y_{ijkl}$ in Equation 1 (see Figure 1A – 1C for an example communication score tensor; see Methods for details

112    about constructing communication score tensor); $\mathcal{B} \in \mathbb{R}^{Q \times J \times K \times L}$ denotes a 4-dimensional coefficient tensor with

113    dimensions of $Q$ sample-level variables, $J$ ligand-receptor pairs, $K$ sender cell types, and $L$ receiver cell types,

114    with the $(q, j, k, l)$ entry corresponding to $\beta_{qjkl}$ in Equation 1; $\boldsymbol{X} \in \mathbb{R}^{I \times Q}$ in Equation 2 denotes sample-level

115    design matrix for $Q$ variables of $I$ samples, with the $(i, q)$ entry corresponding to $x_{iq}$ in Equation 1; $\times_1$ denotes

116    multiplying a tensor by a matrix in the tensor's first dimension; and $\mathcal{E} \in \mathbb{R}^{I \times J \times K \times L}$ denotes a 4-dimensional

117    tensor with the $(i, j, k, l)$ entry corresponding to $\epsilon_{ijkl}$ in Equation 1. The graphic representation of an example

118    tensor model as in Equation 2 is shown in Figure 1C, with disease, age, and batch as example sample-level

119    variables. The detailed illustration of how this supervised tensor technique can incorporate correlations among

120    CCC events is described in the Methods section.

121         To estimate $\mathcal{B}$ in Equation 2, we employ the supervised tensor decomposition technique[16] that considers

122    $\mathcal{B}$ in Equation 2 as a core tensor $\mathcal{G}$ multiplied by 4 factor matrices $\boldsymbol{M}_Q, \boldsymbol{M}_J, \boldsymbol{M}_K, \boldsymbol{M}_L$,

123    $$\mathcal{B} = \mathcal{G} \times_1 \boldsymbol{M}_Q \times_2 \boldsymbol{M}_J \times_3 \boldsymbol{M}_K \times_4 \boldsymbol{M}_L.$$

124    where $\times_d, d = 1,2,3,4$ denotes multiplying a tensor by a matrix in the tensor's $d$ th dimension. For the

125    convenience of presentation, we use $\mathcal{G} \times \{\boldsymbol{M}_Q, \boldsymbol{M}_J, \boldsymbol{M}_K, \boldsymbol{M}_L\}$ to denote the above tensor-by-matrix product.

126    Then the full supervised tensor decomposition model is given by:

127
$$\mathcal{Y} = \mathcal{B} \times_1 X = \mathcal{G} \times \{M_Q, M_J, M_K, M_L\} \times_1 X + \mathcal{E}, \quad \text{(Equation 3)}$$

128　where $M_Q \in \mathbb{R}^{Q \times r_Q}$, $M_J \in \mathbb{R}^{J \times r_J}$, $M_K \in \mathbb{R}^{K \times r_K}$, $M_L \in \mathbb{R}^{L \times r_L}$ are factor matrices. These factor matrices have

129　orthonormal columns (i.e., factors), which can be thought of as the principal components for each dimension.

130　Under the context of cell-cell communication, $M_Q \in \mathbb{R}^{Q \times r_Q}$ contains $r_Q$ factors, representing $r_Q$ effect patterns

131　of $Q$ covariates; $M_J \in \mathbb{R}^{J \times r_J}$ contains $r_J$ factors, representing $r_J$ activity patterns of $J$ ligand-receptor pairs;

132　$M_K \in \mathbb{R}^{K \times r_K}$ contains $r_K$ factors, representing $r_k$ activity patterns of $K$ sender cell type; $M_L$ contains $r_L$

133　factors, represents $r_L$ activity patterns of $L$ receiver cell type; $\mathcal{G} \in \mathbb{R}^{r_Q \times r_J \times r_K \times r_L}$ in Equation 3 denotes the core

134　tensor whose entries show the level of interaction among the factors from different dimensions. We define the

135　decomposition rank $r = (r_Q, r_J, r_K, r_L)$. Details regarding the determination of $r$ are described in the Methods

136　section.

137　　　　We use the QR-adjusted optimization algorithm proposed by Hu et al.[16] to estimate $\mathcal{B}, \mathcal{G}, M_Q, M_J, M_K$

138　$M_L$. The significance level of estimated coefficients in $\mathcal{B}$ are assessed using parametric bootstrap[17]. The details

139　about the optimization algorithm and bootstrap procedure are described in Methods.

140　*Applying STACCato to identify CCC events associated with SLE*

141　　　　We applied STACCato to a scRNA-seq dataset of PBMC samples from 154 SLE subjects and 97 healthy

142　controls[10,11] to identify CCC events associated with SLE while adjusting for age, gender, self-reported ancestry,

143　and processing batch (see Methods for details). The constructed 4-dimensional communication score tensor is a

144　$251 \times 55 \times 9 \times 9$ tensor containing the communication scores of CCC events for 251 samples across 55

145　ligand-receptor pairs, 9 sender cell types, and 9 receiver cell types. The 9 cell types are B cells, natural killer

146　cells (NK), proliferating T and NK cells (Prolif), CD4$^+$ T cells, CD8$^+$ T cells, CD14$^+$ classical monocytes (cM),

147　CD16$^+$ nonclassical monocytes (ncM), conventional dendritic cells (cDC), and plasmacytoid dendritic cells

148　(pDC). We used the decomposition rank $r = (r_Q = 8, r_J = 7, r_K = 4, r_L = 4)$. We used 4,999 iterations of

149　bootstrapping resampling to assess the significance levels of the estimated SLE disease effects. We identified

150　disease effects with p-value $< 0.05$ and magnitude $> 0.015$ as significant disease effects (Supplementary Figure

151　1).

152    Figure 2A displays the estimated factor matrices of the sender and receiver cell type dimension, which

153    represent the activity patterns of sender cell types and receiver cell types. The contribution of each factor to the

154    decomposition is shown in Supplementary Figure 2 (see Methods for details about the calculation of

155    contributions). In both sender and receiver cell type dimension, for factor 1 with the largest contribution, all cell

156    types display scores in the same direction, indicating a critical systematic biological process that involves all cell

157    types. Factor 2 highlights a notable contrast between the lymphocyte group (encompassing B, NK, Prolif, CD4$^+$

158    T, and CD8$^+$ T cells) and the monocyte group (comprising cM, nCM, cDC, and pDC cells), demonstrating

159    opposite activities of these two groups. Factor 3 and Factor 4 unveil distinct activity patterns specific to pDC

160    cells and B cells, respectively, shedding light on the unique roles of these two cell types.

161    Figure 2B displays significant disease effects corresponding to CCC events with B, CD8$^+$ T, cM, and

162    pDC cells as the receiver cell type. The significant effects of CCC events in other receiver cell types are shown

163    in Supplementary Figure 3. Notably, multiple ligand-receptor pairs consistently exhibit positive associations

164    with SLE across sender and receiver cell types. For instance, ligand-receptor pairs LGALS9 – PTPRC and

165    LGALS9 – CD44 consistently show positive associations with SLE across cell types (Figure 2B). This discovery

166    aligns with our earlier findings that the factors representing the systematic biological process involving all cell

167    types have the largest contributions to the decomposition.

168    STACCato also effectively identified CCC events with cell type specific disease effects. For instance,

169    ligand-receptor pair CD99 – PILRA showed negative associations with SLE only with B cells and pDC cells as

170    the receiver cell types (Figure 2B). ligand-receptor pair CD22 – PTPRC demonstrated an significant association

171    with SLE only with B cells as the sender cell type (Figure 2B), which is consistent with the knowledge that

172    CD22 is a B-cell-specific glycoprotein[18].

173    One noteworthy aspect of this SLE dataset is its highly unbalanced study design, where batch 1 included

174    only healthy controls while batch 2 included SLE patients predominantly (Supplementary Table 1).

175    Consequently, batch confounded the association of CCC events with SLE. We applied Tensor-cell2cell[8], which

176    does not consider confounding variables, to the same 4-dimensional communication score tensor of the SLE

177    dataset (Supplementary Figure 4A) and identified three factors (factor 3, 5, 7) significantly associated with SLE

178      disease (Supplementary Figure 4B). However, we found that these factors were also strongly associated with

179      batch (Supplementary Figure 5), suggesting that the disease effect was confounded by the batch effect in these

180      factors (Supplementary Figure 6). For instance, healthy controls exhibited significantly larger loadings in factor

181      3 (Supplementary Figure 4B), indicating a negative association between factor 3 and SLE. However, when

182      excluding batch 1 samples, the difference between SLE patients and healthy controls in other batches became

183      minimal in factor 3 (Supplementary Figure 6). These results demonstrated that batch 1 distorted the association

184      between factor 3 and disease in Tensor-cell2cell, leading to misleading interpretations of factor 3's role in SLE.

185      These findings highlighted the importance of adjusting for confounding effects in CCC inference.

186      *Evaluating the impact of confounding variables on CCC inference with the SLE dataset*

187      To evaluate the impact of confounding variables on CCC inference, we applied STACCato to the SLE

188      dataset with three distinct models, each incorporating different sample-level variables: Model 1, whose results

189      were shown in Figure 2 and described above, considers sample-level variables of disease status, batches, and all

190      other available covariates including age, gender, and ancestry; Model 2 considers disease status and batches

191      only; and Model 3 considers disease status only. When comparing Model 1 and Model 2 to Model 3, we observed

192      substantial changes in the estimated disease effects before and after adjusting for batch effects (Supplementary

193      Figure 7). For example, the ligand-receptor pairs macrophage migration inhibitory factor (MIF) –

194      CD74&CXCR4 and MIF – CD74&CD44 showed negative associations with SLE before batch adjustment but

195      positive associations with SLE after accounting for batch effects. Monoclonal antibodies like imalumab (anti-

196      MIF) and milatuzumab (anti-CD74) have been assessed in early phase clinical trials, demonstrating efficacy in

197      SLE treatment[19]. This suggests a positive association between MIF – CD74 and SLE, which is consistent with

198      the results adjusting for batch effects. These findings underscore how confounding variables can distort true

199      associations and emphasize the importance of considering confounding variables like batches in CCC inference.

200      We also compared the factor matrices estimated with and without adjustment of batch effects by

201      calculating the normalized chordal distance between the estimated factor matrices. Normalized chordal distance

202      is a metric ranging from 0 to 1 for measuring distances between subspaces. A larger chordal distance indicates

203      a greater difference between the subspaces of the estimated factor matrices (see Methods for details about chordal

204    distance). The normalized chordal distances between the factor matrices estimated before (Model 3) and after

205    adjusting for batches (Model 2) were 0.009 for sender cell types and 0.013 for receiver cell types, indicating

206    minor differences. These results illustrate that confounding variables can significantly influence the estimation

207    of disease effects in CCC events while having a relatively minor impact on the estimation of factor matrices.

208    *Applying STACCato to identify CCC events associated with ASD*

209    We applied STACCato on the snRNA-seq dataset of postmortem tissue samples of prefrontal cortex

210    from 13 ASD patients and 10 controls[12] to identify CCC events associated with ASD (see Methods for details).

211    We considered 16 sender/receiver cell types: fibrous astrocytes (AST-FB), protoplasmic astrocytes (AST-PP),

212    Endothelial, parvalbumin interneurons (IN-PV), somatostatin interneurons (IN-SST), SV2C interneurons (IN-

213    SV2C), VIP interneurons (IN-VIP), layer 2/3 excitatory neurons (L2/3), layer 4 excitatory neurons (L4), layer

214    5/6 corticofugal projection neurons (L5/6), layer 5/6 cortico-cortical projection neurons (L5/6-CC), maturing

215    neurons (Neu-mat), NRGN-expressing neurons (Neu-NRGN-I), NRGN-expressing neurons (Neu-NRGN-II),

216    Oligodendrocyte precursor cells (OPC), and oligodendrocytes. We applied STACCato to a $23 \times 749 \times 16 \times$

217    $16$ communication score tensor (consisting of 23 samples, 749 ligand-receptor pairs, 16 sender cell types, 16

218    receiver cell types) to examine associations between CCC events and ASD, while adjusting for age, gender, and

219    processing batch. We used the decomposition rank $\boldsymbol{r} = (r_Q = 5, r_J = 5, r_K = 5, r_L = 5)$. We used 4,999 iterations

220    of bootstrapping resampling to assess the significance levels of the estimated ASD disease effects. We identified

221    estimated disease effects with p-value $< 0.05$ and magnitude $> 0.015$ as significant disease effects

222    (Supplementary Figure 8).

223    In Figure 3A, we present the estimated factor matrices of the sender and receiver cell type dimension,

224    which depict the activity patterns of sender and receiver cell types. The contributions of all factors are shown in

225    Supplementary Figure 9A – 9B. Similar to our findings in the SLE dataset, we observed that factor 1 contributed

226    the most and reflected a systematic process involving all cell types. Factors 2 through 5 for both sender and

227    receiver cell types successfully revealed 6 cell type groups with distinct activity patterns: (1) astrocytes group

228    including AST-FB and AST-PP; (2) Endothelial; (3) inhibitory neurons group including IN-PV, IN-SST, IN-

229    SV2C, IN-VIP; (4) excitatory neurons group including L2/3, L4, L5/6, L5/6-CC; (5) expressing neurons group

230  including Neu-mat, Neu-NRGN-I, and Neu-NRGN-II; (6) neuroglia group including oligodendrocytes and OPC

231  (Figure 3A).

232     For each pair of sender cell type and receiver cell type, we ranked the ligand-receptor pairs by the

233  estimated ASD disease effects and performed preranked Gene Set Enrichment Analysis (GSEA)[20] to determine

234  if ligand-receptor pairs belonging to a particular pathway are more likely to be clustered at the top or bottom of

235  the ranked list, and thereby identifying pathways associated with ASD (see details of pathway enrichment

236  analysis in the Methods section). Figure 3B shows significantly enriched KEGG pathways[21] across AST-PP,

237  Endothelial, IN-PV, L2/3, and Neu-NRGN-I cells. A total of 10 significantly enriched pathways were identified,

238  including the axon guidance, cell adhesion molecules (CAMs), cytokine-cytokine receptor interaction,

239  extracellular matrix-receptor (ECM-receptor) interaction, ErbB signaling, focal adhesion, MAPK signaling,

240  notch signaling, regulation of actin cytoskeleton, and small cell lung cancer. Importantly, 8 out of these 10

241  pathways (axon guidance, CAMs, ECM-receptor interaction, ErbB signaling, focal adhesion, MAPK signaling,

242  regulation of actin cytoskeleton, small cell lung cancer) have been previously identified as significantly enriched

243  pathways with p-values $< 5 \times 10^{-7}$ for ASD[22]. The molecules related to the notch signaling pathway have

244  been shown to have increased expression in the PFC in an animal model of autism[23], which is consistent with

245  our observation of a positive association of the notch signaling pathway with ASD between AST-FB and L2/3

246  cells.

247  *Evaluating the impact of confounding variables on CCC inference with the ASD dataset*

248     We also examined the impact of batch information on our ASD results by fitting three distinct

249  STACCato models with Model 1 considering disease status and all available covariates including batches, age,

250  and gender (as shown in Figure 3), Model 2 considering disease status and batches only, and Model 3 considering

251  disease status only. Unlike the SLE dataset, the ASD dataset exhibits a fairly balanced design (Supplementary

252  Table 2). Consequently, batch is no longer a confounding factor. As anticipated, the estimated disease effects

253  remain consistent before and after adjusting for batch effects (Supplementary Figure 10). Interestingly, the

254  chordal distances between the factor matrices estimated before (Model 3) and after adjusting for batch (Model

255  2) were 0.384 for sender cell types and 0.438 for receiver cell types, indicating substantial discrepancies in the

256   estimated factor matrices before and after batch adjustment. We further evaluated the relative contributions of

257   all sample-level variables and found that batch contributed substantially to the communication tensor, indicating

258   a non-negligible batch effect on the communication scores (Supplementary Figure 9C). This underscores a

259   crucial point — even in datasets with balanced designs, failing to account for variables with significant impacts

260   on the CCC can significantly impact the estimation of factor matrices and, consequently, the interpretations of

261   cell type activity patterns.

262   *Simulation Study*

263   We conducted simulations to investigate how sample-level variables affect the CCC inference in

264   different study designs. We simulated the communication score tensor $\mathcal{Y} \in \mathbb{R}^{I \times J \times K \times L}$ from the supervised tensor

265   decomposition model as in Equations 2 and 3. We set $\mathcal{G}, M_Q, M_J, M_K, M_L$ in Equation 3 as the core tensor and

266   factor matrices estimated from the ASD dataset and simulated $X$ for 60 subjects with intercept, disease status,

267   and batch variables. The elements of $\mathcal{E}$ were independently simulated from a normal distribution with mean 0

268   and variance $\hat{\sigma}^2$, where $\hat{\sigma} = 0.05$ was taken as the standard error of the estimation residuals from ASD data. We

269   considered a study with 30 disease subjects and 30 healthy controls processed in two batches. We considered

270   three study designs: (1) balanced design with 15 controls and 15 disease subjects in both batches; (2) moderate

271   unbalanced design with 20 controls and 10 disease subjects in batch 1, and 10 controls and 20 disease subjects

272   in batch 2; (3) extreme unbalanced design with 30 controls and 5 disease subjects in batch 1, and batch 2 only

273   contains 25 disease subjects.

274   We applied STACCato with two models: Model 1 considers disease status and batch variables, and

275   Model 2 considers only disease status. We calculated the mean squared errors (MSEs) of the estimated disease

276   effects across 100 simulations. Figure 4A shows that neglecting confounders in an unbalanced design can

277   generate larger estimation errors, and the MSEs of the disease effect dramatically increased as the degree of

278   imbalance became more extreme. We also assessed the proportion of estimated disease effects with opposite

279   directions to the assumed one (Supplementary Figure 11). We found that, before adjusting for batch, 14.7% of

280   the disease effects had incorrect estimated directions in the extremely unbalanced design, which was

281   significantly higher than the proportion 3.1% after adjusting for batch. Additionally, we assessed the accuracy

282  of the estimated factor matrices by calculating the chordal distance between the estimated factor matrices and

283  the assumed factor matrices. We observed that neglecting the batch variable resulted in decreased accuracy in

284  estimating the factor matrices (Figure 4B), especially in balanced and moderate unbalanced design. Failing to

285  account for the batch variable prevents the identification of factors that are solely batch-associated and not

286  disease-associated, resulting in inaccuracies in the estimated factor matrices. Conversely, in extreme unbalanced

287  designs where batch and disease are strongly correlated, batch-associated factors are also strongly linked to the

288  disease. In such scenarios, neglecting the batch variable did not significantly impact the accuracy of estimating

289  the factor matrices. These observations align with our real-data analysis findings, suggesting that regardless of

290  whether the dataset originates from a balanced or unbalanced design, incorporating information of sample-level

291  variables into CCC inference consistently leads to more accurate estimations of disease effects or activity

292  patterns of cell types.

293  We also compared STACCato to the separate regression procedure (Equation 1), where a regression

294  model was fitted with communication scores as dependent variables and sample-level variables as independent

295  variables separately for each CCC event. In contrast, STACCato employs the tensor technique to incorporate the

296  correlations among CCC events and jointly estimates the effects of considered variables for all CCC events.

297  Across all study designs, STACCato consistently achieved significantly lower MSE compared to the separate

298  regression approach (Supplementary Figure 12), justifying the advantage of using the tensor technique to account

299  for correlations among CCC events.

300  *Computational Considerations*

301  While a single STACCato decomposition only takes seconds, assessing the significance level of

302  estimated effects by bootstrapping requires performing decompositions for a substantial number of bootstrapping

303  iterations and takes hours of CPU time. We conducted the computational benchmarks using one Intel(R) Xeon(R)

304  processor (2.10 GHz). For a simulated dataset comprising 100 samples, 10 sender and receiver cell types, 600

305  ligand-receptor pairs, and 10 sample-level covariates, 99 iterations of bootstrap resampling took around 11

306  minutes and ~1.3 GB memory usage on the upper-bound.

307    Considering that the numbers of cell types and sample-level covariates generally do not vary much in

308    practice, we investigated how bootstrapping time and upper-bound memory usage vary with the number of

309    samples and the number of ligand-receptor pairs. We simulated datasets with 10 sender and receiver cell types,

310    10 sample-level covariates, and various numbers of samples (ranging from 25 to 100) and ligand-receptor pairs

311    (ranging from 150 to 600). With 99 iterations of bootstrap resampling, our simulation results revealed that

312    computational time increased linearly with the number of samples (Supplementary Figure 13A) and

313    quadratically with the number of ligand-receptor pairs (Supplementary Figure 14A). The upper bound memory

314    usage changed approximately linearly with both the number of samples and ligand-receptor pairs

315    (Supplementary Figures 13B, 14B).

316    **Discussion**

317    We present STACCato, a computational tool that utilizes multi-sample multi-condition scRNA-seq data

318    to identify CCC events associated with conditions (e.g., disease status, multiple time points, different tissue

319    types). STACCato utilizes supervised tensor decomposition to estimate the influence of the condition of interest

320    on CCC events, while adjusting for potential confounding variables. Furthermore, it facilitates the identification

321    of activity patterns among cell types involved in CCC. We applied STACCato to analyze a SLE dataset with an

322    extremely unbalanced design[10,11] and an ASD dataset with a balanced design[12]. Additionally, we conducted

323    simulation studies to mimic real data with different study designs. Our real data application and simulation

324    results demonstrated STACCato's capability to incorporate available sample-level variables, thereby enabling

325    more reliable inference regarding the associations between CCC events and conditions, as well as more robust

326    estimations of activity patterns among cell types.

327    In practice, a common approach to address batch effects in scRNA-seq data is to remove batch effects

328    before downstream analysis. This approach involves the estimation of batch effects, followed by the removal of

329    these estimated batch effects to generate "batch-effect-free" data for downstream analysis. However, as noted

330    by Nygaard et al.[24], this two-step procedure has a severe drawback: it relies on point estimates of batch effects

331    while disregarding estimation errors. In this two-step process, even when the original batch effects could be

332    eliminated, the estimation errors may introduce new batch effects. In contrast, STACCato incorporates potential

333    confounding variables, such as batch effects, into the design matrix, and jointly estimates the effects of these

334    confounders along with other variables in a single step. Moreover, although our application and simulation

335    studies focused on addressing batch effects, STACCato can adjust for all potential confounding variables in

336    biomedical research. For instance, age is often considered as a confounding factor in the identification of CCC

337    events associated with Alzheimer's disease. By incorporating all potential confounding variables into the model,

338    STACCato offers a comprehensive solution, allowing for simultaneous handling of multiple confounders and

339    facilitating more accurate CCC inference.

340         In contrast to Tensor-cell2cell, which also employs the tensor decomposition technique for CCC

341    inference, STACCato stands out in several key aspects. First, STACCato directly assesses the relationship

342    between each CCC event and the condition of interest. In contrast, Tensor-cell2cell primarily provides insights

343    into the association between the decomposed factors and conditions, without offering explicit interpretations

344    regarding individual CCC events. Second, STACCato goes a step further by not only identifying associations

345    but also estimating the condition effect for each CCC event and assessing the statistical significance of such an

346    effect. In contrast, Tensor-cell2cell focuses on determining the significance of the association between factors

347    and the condition, without providing detailed information on the magnitude of condition effects. Last, as

348    highlighted throughout our paper, STACCato has the capability to account for confounding variables, a feature

349    lacking in Tensor-cell2cell. Through our application of Tensor-cell2cell to the SLE dataset, we demonstrated its

350    inability to effectively disentangle confounding effects from disease effects in the study of CCC events.

351         It is important to note that STACCato is a highly adaptable framework that can be seamlessly

352    integrated with various existing CCC inference tools, each with its unique methods of constructing

353    communication scores. Researchers have the flexibility to select any tool of interest to calculate

354    communication scores. For example, one can use the LIANA tool[25], which incorporates a wide range of tools

355    and resources to calculate cell-cell communication scores, to calculate communication scores for all CCC

356    events and arrange the scores into a 3-dimensional communication score tensor per sample. The 3-dimensional

357    tensors of all samples can subsequently be combined into the 4-dimensional communication score tensor,

358    allowing STACCato to be applied for inferring CCC events associated with the specific condition of interest.

14

359     The STACCato framework does have its limitations. First, in scRNA-seq data, many genes may not

360     be actively expressed in single cells, resulting in a significant proportion of zero values in the cell-cell

361     communication score tensor. A future extension of STACCato involving sparse tensor decomposition, which

362     imposes sparsity constraints on the ligand-receptor pairs, may inherently address this zero-inflation problem.

363     Second, STACCato relies on a literature-curated database to perform CCC inference, limiting the identified

364     condition-related CCC events to those documented in previous literature. Extending STACCato to identify

365     novel ligand-receptor pairs is part of our ongoing research but falls outside the scope of this work.

366     To enable the use of STACCato by the public, we provide an integrated tool (see Code availability) to:

367     (1) perform supervised tensor decomposition to estimate the effects of conditions on CCC events adjusting for

368     covariates and infer activity patterns of cell types; (2) use bootstrapping resampling to assess the significance

369     level of the estimated effects; (3) conduct downstream analyses including comparing significant CCC events

370     across cell types and identifying pathways significantly associated with conditions. In conclusion, we present

371     STACCato as a valuable tool to effectively incorporate sample-level variables and adjust for possible

372     confounding variables in CCC inference using multi-sample multi-condition scRNA-seq data.

373     **Methods**

374     *Construction of a 4-dimensional communication score tensor*

375     With the matrix of gene expressions of multiple cell types from a scRNA-seq sample and the

376     literature-curated list of ligand-receptor pairs, we can calculate the communication score for the CCC event

377     involving the interaction of ligand-receptor pair $j$ from sender cell type $k$ to receiver cell type $l$ as

$$y_{jkl} = f(\text{ligand}_k, \text{receptor}_l)$$

379     where $y_{jkl}$ denotes the communication score; $\text{ligand}_k$ denotes the expression of the ligand in sender cell type

380     $k$; $\text{receptor}_l$ denotes the expression of the receptor in receiver cell type $l$; and $f$ denotes the scoring function

381     (Figure 1A). In this study, we used the scoring function $y_{jkl} = \sqrt{\text{ligand}_k \times \text{receptor}_l}$. Other available scoring

382     functions have been previously summarized by Armingol et al.[26] and Dimitrov et al[25].

383     Once we compute communication scores for a specific ligand-receptor pair $j$ across all $K$ sender cell

384     types and $L$ receiver cell types, we can create a communication score matrix (Figure 1B). In this matrix, the rows

385     represent $K$ sender cell types; the columns represent $L$ receiver cell types; and the element located in the $k^{th}$ row

386     and $l^{th}$ column corresponds to the value of $y_{jkl}$. By repeating this process for all $J$ ligand-receptor pairs, we will

387     get $J$ matrices, which can be arranged into a sample-specific 3-dimensional tensor with dimensions $J \times K \times L$

388     (Figure 1B). Then the 3-dimensional tensor of all samples can be arranged into a 4-dimensional tensor with

389     dimensions of $I$ samples, $J$ ligand-receptor pairs, $K$ sender cell types, and $L$ receiver cell types (Figure 1C). In

390     the application studies of the SLE dataset and ASD dataset, we constructed the 4-dimensional tensor using the

391     Tensor-cell2cell  package[8] (see Code availability). In the final tensor, we only included ligand-receptor pairs

392     with both ligands and receptors shared across all samples.

393     *STACCato incorporates correlations among CCC events*

394     Consider the full supervised tensor decomposition model in Equation 3,

395     $$\mathcal{Y} = \mathcal{B} \times_1 \boldsymbol{X} + \mathcal{E} = \mathcal{G} \times \{\boldsymbol{M}_Q, \boldsymbol{M}_J, \boldsymbol{M}_K, \boldsymbol{M}_L\} \times_1 \boldsymbol{X} + \mathcal{E}.$$

396     Elementwise, we have

397     $$\beta_{qjkl} = \sum_{r_1=1}^{r_Q} \sum_{r_2=1}^{r_J} \sum_{r_3=1}^{r_K} \sum_{r_4=1}^{r_L} g_{r_1 r_2 r_3 r_4} M_Q^{qr_1} M_J^{jr_2} M_K^{kr_3} M_L^{lr_4} \quad \text{(Equation 4)}$$

398     where $g_{r_1 r_2 r_3 r_4}$ denotes the $(r_1, r_2, r_3, r_4)$ entry of $\mathcal{G}$, $M_Q^{qr_1}$ denotes the entry in the $q^{th}$ row and $r_1^{th}$ column of

399     $M_Q$, similarly for $M_J^{jr_2}$, $M_K^{kr_3}$, and $M_L^{lr_4}$. Then for $k \neq k'$ and $l \neq l'$,

400     $$\beta_{qjk'l'} = \sum_{r_1=1}^{r_Q} \sum_{r_2=1}^{r_J} \sum_{r_3=1}^{r_K} \sum_{r_4=1}^{r_L} g_{r_1 r_2 r_3 r_4} M_Q^{qr_1} M_J^{jr_2} M_K^{k'r_3} M_L^{l'r_4} \quad \text{(Equation 5)}$$

401     Equations 4 and 5 represent the effects of covariate $q$ on two different CCC events with the same ligand-

402     receptor pair $j$ but different sender (sender cell type $k$ in Equation 4 and $k'$ in Equation 5) and receiver cell types

403     (receiver cell type $l$ in Equation 4 and $l'$ in Equation 5). These two equations share the same parameters

404　　$M_J^{jr_2}$, $r_2 = 1, \cdots r_J$. Similarly, for CCC events with the same sender cell type $k$, the effects share the same

405　　parameters $M_K^{kr_3}$, $r_3 = 1, \cdots r_K$; for CCC events with the same receiver cell type $l$, the effects share the same

406　　parameters $M_L^{lr_4}$, $r_4 = 1, \cdots r_L$. In STACCato, the effects of covariates on correlated CCC events share

407　　parameters, enabling it to effectively incorporate the complex correlation structure among these CCC events.

408　　*STACCato Optimization*

409　　　　We first determine the number of components $r_J$, $r_K$, $r_L$ for ligand-receptor pair, sender cell type, and

410　　receiver cell type dimension. For each dimension, we start by performing tensor unfolding to rearrange the

411　　elements of the communication score tensor into a matrix. For example, for the ligand-receptor pair dimension,

412　　we transform $\mathcal{Y} \in \mathbb{R}^{I \times J \times K \times L}$ into a matrix $Y_{(J)}$ with $J$ rows and $I \times K \times L$ columns. Then we set $r_J$ as the

413　　number of components that can explain more than 1% of the variation in $Y_{(J)}$. We follow the same approach to

414　　determine $r_K$ for sender cell type dimension and $r_L$ for receiver cell type dimension. We set $r_Q$ as the number of

415　　sample-level variables available in $X$.

416　　　　Denoting the supervised decomposition rank $r = (r_Q, r_J, r_K, r_L)$, we follow the optimization algorithm

417　　proposed by Hu et al. [16] to estimate $\mathcal{B}, \mathcal{G}, M_Q, M_J, M_K \, M_L$:

---

**Algorithm 1:**

Input: communication score tensor $\mathcal{Y} \in \mathbb{R}^{I \times J \times K \times L}$, sample-level design matrix $X \in \mathbb{R}^{I \times Q}$, rank $r$.

1. Normalize sample-level design matrix via QR factorization $X = QR$.
2. Project $\mathcal{Y}$ to the multilinear sample-level variable space to obtain the unconstrained coefficient tensor: $\widetilde{\mathcal{B}} = \mathcal{Y} \times_1 Q^T$.
3. Obtain rank-unconstrained coefficient tensor by performing a rank-$r$ higher-order orthogonal iteration (HOOI)[27] on $\widetilde{\mathcal{B}}$: $\widehat{\mathcal{B}}^{(0)} \leftarrow HOOI(\widetilde{\mathcal{B}}, r)$.
4. Obtain estimated coefficient tensor by re-normalizing $\widehat{\mathcal{B}}^{(0)}$ back to the original feature scales:

   $\widehat{\mathcal{B}} = \widehat{\mathcal{B}}^{(0)} \times_1 R^{-1}$.

5. Estimate $\mathcal{G}, M_Q, M_J, M_K, M_L$ by performing a rank-$r$ HOOI on $\widehat{\mathcal{B}}$: $\widehat{\mathcal{B}} \approx \widehat{\mathcal{G}} \times \{\widehat{M_Q}, \widehat{M_J}, \widehat{M_K}, \widehat{M_L}\}$.

Output: $\widehat{\mathcal{B}}, \widehat{\mathcal{G}}, \widehat{M_Q}, \widehat{M_J}, \widehat{M_K}, \widehat{M_L}$.

---

418

419     We also impose orthonormality on $M_Q, M_J, M_K\ M_L$ to ensure the uniqueness of decomposition.

420     *Parametric bootstrapping for hypothesis testing*

421     Denote the estimated communication score tensor as $\widehat{\mathcal{Y}} = \widehat{\mathcal{B}} \times_1 X$ with entry $\hat{y}_{ijkl}$ and the estimated

422     standard error of $\epsilon_{ijkl}$ as $\hat{\sigma}$, we have residual tensor $\mathcal{S} = \mathcal{Y} - \widehat{\mathcal{Y}}$ with entry $s_{ijkl} = y_{ijkl} - \hat{y}_{ijkl}$, and $\hat{\sigma}^2 =$

423     $var(\text{vec}(\mathcal{S}))$, where $\text{vec}(\mathcal{S}) = [s_{1111}, \cdots, s_{ijkl}]$ denotes the vectorized version of tensor $\mathcal{S}$.

424     For the $n^{th}$ bootstrap resampling[17], we generate a new tensor $\mathcal{S}^{(n)}$ with entries from $N(0, \hat{\sigma}^2)$ and

425     construct a new communication score tensor $\mathcal{Y}^{(n)} = \widehat{\mathcal{Y}} + \mathcal{S}^{(n)}$. We perform STACCato on $\mathcal{Y}^{(n)}$ to estimate a

426     new coefficient tensor $\widehat{\mathcal{B}}^{(n)}$. We repeat this procedure for $N$ iterations to generate $\widehat{\mathcal{B}}^{(1)}, \widehat{\mathcal{B}}^{(2)}, \cdots, \widehat{\mathcal{B}}^{(N)}$. To test

427     the null hypothesis of $H_0: b_{qjkl} = 0$, we follow the guideline suggested by Hall and Wilson[28] to define the

428     bootstrap p-value as:

429
$$p_{qjkl} = \frac{\sum_{n=1}^{N} I\left(\left|\hat{b}_{qjkl}^{(n)} - \hat{b}_{qjkl}\right| > |\hat{b}_{qjkl}|\right)}{N + 1}$$

430     where $\hat{b}_{qjkl}^{(n)}$ denotes the $(q, j, k, l)$ entry of $\widehat{\mathcal{B}}^{(n)}$; $\hat{b}_{qjkl}$ denotes the $(q, j, k, l)$ entry of $\widehat{\mathcal{B}}$, which is the estimated

431     effect of variable $q$ on the CCC events involving the ligand-receptor pair $j$ between sender cell type $k$ and

432     receiver cell type $l$; and $p_{qjkl}$ is the bootstrapping p-value for $\hat{b}_{qjkl}$.

433     *Calculation of contributions*

434     To calculate the contributions of factors of the sender and receiver cell types, we remove each factor

435     from the decomposition results and assess the changes in the estimated outcome. For example, for factor 1 in the

436     sender cell type dimension, we first remove the first column of the estimated factor matrix $\widehat{M_K}$ and construct a

437     new factor matrix $\widehat{M_K}^* \in \mathbb{R}^{K \times (r_K - 1)}$. We then eliminate the interactions between this factor and factors in other

438     dimensions from the estimated core tensor $\hat{\mathcal{G}}$, creating a new core tensor $\hat{\mathcal{G}}^* \in \mathbb{R}^{r_Q \times r_J \times (r_K - 1) \times r_L}$. With the

439     modified factor matrices and core tensor, we calculate a new predicted communication score tensor $\widehat{\mathcal{Y}}^* =$

440 $\hat{\mathcal{G}}^* \times \{\widehat{M_Q}, \widehat{M_J}, \widehat{M_K}^*, \widehat{M_L}\} \times_1 X$. The contribution of the removed factor is defined as the mean squared difference

441 between the entries of $\hat{\mathcal{Y}}^*$ and the original estimated $\hat{\mathcal{Y}} = \hat{\mathcal{G}} \times \{\widehat{M_Q}, \widehat{M_J}, \widehat{M_K}, \widehat{M_L}\} \times_1 X$.

442 *Chordal distance between two subspaces*

443        We use normalized chordal distance[29] to evaluate the distance between the column spaces of two factor

444 matrices. Let $A \in \mathbb{R}^{d_1 \times d_2}$, $B \in \mathbb{R}^{d_1 \times d_2}$ as two matrices whose columns are the orthonormal bases of two

445 subspaces **A** and **B**, and $A^T B = U\Sigma V^T$ as the full singular value decomposition (SVD) of $A^T B$ with $\Sigma =$

446 $diag(\sigma_1, \sigma_2, \cdots, \sigma_{d_2})$. The principal angles $\theta_1 \leq \theta_2 \leq \cdots \leq \theta_{d_2}$ between the subspaces **A** and **B** are given by:

447 $$\theta_i = \cos^{-1} \sigma_i, i = 1, \cdots, d_2$$

448 The chordal distance between the subspaces **A** and **B** is given by:

449 $$d(\mathbf{A}, \mathbf{B}) = \left( \sum_{i=1}^{d_2} \sin^2 \theta_i \right)^{\frac{1}{2}}.$$

450 Here, we use the normalized chordal distance $d^*(\mathbf{A}, \mathbf{B}) = \left( \frac{1}{d_2} \sum_{i=1}^{d_2} \sin^2 \theta_i \right)^{\frac{1}{2}}$ so that the measure is bounded

451 within [0,1]. We used the R function *chord.norm.diff* from CJIVE package[30] (see Code availability) to calculate

452 the normalized chordal distance.

453 *RNA-seq data processing*

454        For all scRNA-seq datasets used in the study, we filtered out genes expressed in fewer than 4 cells and

455 utilized the provided cell type labels from the metadata. For each sample in the dataset, we aggregate gene

456 expression from single cells/nuclei into cell types by calculating the fraction of cells with non-zero counts within

457 each cell type. Therefore, the aggregated cell-type specific gene expression is bounded within [0,1]. This

458 approach is endorsed by Tensor-cell2cell for the accurate representation of genes with low expression levels[8,31],

459 which is common among genes responsible for encoding surface proteins[32].

460 *Literature-curated lists of ligand-receptor pairs*

461  We downloaded the human list of 2,005 ligand-receptor pairs from a public available compendium of

462  lists of ligand-receptor pairs (see Data availability). This list of ligand-receptor pairs was originally curated by

463  Jin et al[1].

464  *scRNA-seq dataset of SLE patients and controls*

465  The SLE scRNA-seq dataset collects multiplexed scRNA-seq of 264 PBMC samples from 162 SLE

466  patients and 99 healthy controls[10,11]. The data in h5ad format was obtained from NCBI's Gene Expression

467  Omnibus[33] with GEO accession number 174188 (see Data availability). From the h5ad data, we extracted the

468  raw UMI counts of 32,738 genes across 1,263,676 cells from 264 samples and 99 technical replicates. We

469  reduced the dataset down to one sample per subject by selecting the sample with the largest number of cells.

470  The metadata, which was also extracted from the h5ad data, includes the information of age, processing

471  batch, ancestry, and gender of subjects. 107 (41%) subjects are Asian, 149 (57%) subjects are European, 3 (1%)

472  subjects are African American, and 2 (1%) subjects are Hispanic. We filtered out 5 samples of African American

473  or Hispanic history, and only kept samples containing 9 main cell types: B, NK, Prolif, CD4$^+$ T cells, CD8$^+$ T

474  cells, cM, ncM, cDC, and pDC cells. The remaining 251 samples include 154 SLE patients and 97 healthy

475  controls from 4 processing batches. The constructed CCC tensor for the SLE dataset resulted in a 4-dimensional

476  tensor with 251 subjects, 55 ligand-receptor pairs, 9 sender cell types, and 9 receiver cell types.

477  *scRNA-seq dataset of ASD patients and controls*

478  For the ASD dataset, we downloaded the log2-transformed UMI counts of PFC samples and the

479  corresponding metadata from the UCSC Cell Browser[34] (see Data availability). The raw dataset contains the

480  expression levels of 36,501 genes across 62,166 cells from 13 ASD patients and 10 healthy controls[12]. The

481  constructed CCC tensor for the ASD dataset resulted in a 4-dimensional tensor with 23 subjects, 749 ligand-

482  receptor pairs, 16 sender cell types, and 16 receiver cell types.

483  *Gene set enrichment analysis*

484  We follow the procedure proposed in Tensor-cell2cell to conduct the GSEA. A ligand-receptor pair is

485  considered in a pathway if all the genes participating in the ligand-receptor pair are in the pathway. We consider

486    the 22 KEGG pathways selected by Tensor-cell2cell (see Data availability). For one pair of sender cell type and

487    receiver cell type, we first rank ligand-receptor pairs by their estimated disease effects, and then use the *prerank*

488    module in the Python package GSEApy[35] (see Code availability) with 4999 permutations, gene sets with at least

489    15 elements, and a score weight of 1 to calculate the enrichment p-value and normalized enrichment score. We

490    then combined the results from all tested pairs of cell types, and performed false discovery rate (FDR) correction

491    to adjust for multiple comparisons. Pathways with FDR q-value < 0.05 were identified as pathways significantly

492    associated with disease.

**Acknowledgements**

**Data availability**

497    The human list of 2,005 ligand-receptor pairs was downloaded from

498    https://github.com/LewisLabUCSD/Ligand-Receptor-Pairs/blob/master/Human/Human-2020-Jin-LR-

499    pairs.csv. The processed data of the SLE dataset in h5ad format was downloaded from

500    https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE174188. The log2-transformed UMI counts of the

501    ASD dataset was downloaded from https://cells.ucsc.edu/autism/downloads.html. The KEGG pathways

502    selected by Tensor-cell2cell to perform GSEA was downloaded from

503    https://codeocean.com/capsule/9737314/tree/v2/data/LR-Pairs/CellChat-LR-KEGG-set.pkl.

**Code availability**

505    Source code for STACCato is available from https://github.com/daiqile96/STACCato. Source code for CJIVE

506    is available from https://cran.r-project.org/web/packages/CJIVE/index.html. Source code for Tensor-cell2cell

507    is available from https://github.com/earmingol/cell2cell. Source code for GSEApy is available from

508    https://github.com/zqfang/GSEApy.

509

510 **Reference**

511 1.  Jin, S. *et al.* Inference and analysis of cell-cell communication using CellChat. *Nat Commun* **12**, 1088
512     (2021).

513 2.  Efremova, M., Vento-Tormo, M., Teichmann, S. A. & Vento-Tormo, R. CellPhoneDB: inferring cell–cell
514     communication from combined expression of multi-subunit ligand–receptor complexes. *Nat Protoc* **15**,
515     1484–1506 (2020).

516 3.  Raredon, M. S. B. *et al.* Computation and visualization of cell–cell signaling topologies in single-cell
517     systems data using Connectome. *Sci Rep* **12**, 4187 (2022).

518 4.  Hu, Y., Peng, T., Gao, L. & Tan, K. CytoTalk: De novo construction of signal transduction networks
519     using single-cell transcriptomic data. *Science Advances* **7**, eabf1356.

520 5.  Wang, Y. *et al.* iTALK: an R Package to Characterize and Illustrate Intercellular Communication. 507871
521     Preprint at https://doi.org/10.1101/507871 (2019).

522 6.  Hou, R., Denisenko, E., Ong, H. T., Ramilowski, J. A. & Forrest, A. R. R. Predicting cell-to-cell
523     communication networks using NATMI. *Nat Commun* **11**, 5011 (2020).

524 7.  Cabello-Aguilar, S. *et al.* SingleCellSignalR: inference of intercellular networks from single-cell
525     transcriptomics. *Nucleic Acids Research* **48**, e55–e55 (2020).

526 8.  Armingol, E. *et al.* Context-aware deconvolution of cell–cell communication with Tensor-cell2cell. *Nat*
527     *Commun* **13**, 3665 (2022).

528 9.  Tsuyuzaki, K., Ishii, M. & Nikaido, I. scTensor detects many-to-many cell–cell interactions from single
529     cell RNA-sequencing data. 2022.12.07.519225 Preprint at https://doi.org/10.1101/2022.12.07.519225
530     (2022).

531 10. Thompson, M. *et al.* Multi-context genetic modeling of transcriptional regulation resolves novel disease
532     loci. *Nat Commun* **13**, 5704 (2022).

533 11. Perez, R. K. *et al.* Single-cell RNA-seq reveals cell type-specific molecular and genetic associations to
534     lupus. *Science* **376**, eabf1970 (2022).

535 12. Nassir, N. *et al.* Single-cell transcriptome identifies molecular subtype of autism spectrum disorder
536     impacted by de novo loss-of-function variants regulating glial cells. *Human Genomics* **15**, 68 (2021).

537    13. Liao, M. *et al.* Single-cell landscape of bronchoalveolar immune cells in patients with COVID-19. *Nat*

538        *Med* **26**, 842–844 (2020).

539    14. Hore, V. *et al.* Tensor decomposition for multi-tissue gene expression experiments. *Nat Genet* **48**, 1094–

540        1100 (2016).

541    15. Jung, I., Kim, M., Rhee, S., Lim, S. & Kim, S. MONTI: A Multi-Omics Non-negative Tensor

542        Decomposition Framework for Gene-Level Integrative Analysis. *Frontiers in Genetics* **12**, (2021).

543    16. Hu, J., Lee, C. & Wang, M. Generalized Tensor Decomposition With Features on Multiple Modes.

544        *Journal of Computational and Graphical Statistics* **31**, 204–218 (2022).

545    17. Efron, B. & Tibshirani, R. J. *An Introduction to the Bootstrap*. (CRC Press, 1994).

546    18. Kelm, S., Gerlach, J., Brossmer, R., Danzer, C.-P. & Nitschke, L. The Ligand-binding Domain of CD22 Is

547        Needed for Inhibition of the B Cell Receptor Signal, as Demonstrated by a Novel Human CD22-specific

548        Inhibitor Compound. *J Exp Med* **195**, 1207–1213 (2002).

549    19. Bilsborrow, J. B., Doherty, E., Tilstam, P. V. & Bucala, R. Macrophage migration inhibitory factor (MIF)

550        as a therapeutic target for rheumatoid arthritis and systemic lupus erythematosus. *Expert Opin Ther*

551        *Targets* **23**, 733–744 (2019).

552    20. Subramanian, A. *et al.* Gene set enrichment analysis: A knowledge-based approach for interpreting

553        genome-wide expression profiles. *Proceedings of the National Academy of Sciences* **102**, 15545–15550

554        (2005).

555    21. Kanehisa, M. & Goto, S. KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids Res* **28**, 27–

556        30 (2000).

557    22. Wen, Y., Alshikho, M. J. & Herbert, M. R. Pathway Network Analyses for Autism Reveal Multisystem

558        Involvement, Major Overlaps with Other Diseases and Convergence upon MAPK and Calcium Signaling.

559        *PLoS ONE* **11**, (2016).

560    23. Zhang, Y. *et al.* The Notch signaling pathway inhibitor Dapt alleviates autism-like behavior, autophagy

561        and dendritic spine density abnormalities in a valproic acid-induced animal model of autism. *Prog*

562        *Neuropsychopharmacol Biol Psychiatry* **94**, 109644 (2019).

563    24. Nygaard, V., Rødland, E. A. & Hovig, E. Methods that remove batch effects while retaining group

564        differences may lead to exaggerated confidence in downstream analyses. *Biostatistics* **17**, 29–39 (2016).

565    25. Dimitrov, D. *et al.* Comparison of methods and resources for cell-cell communication inference from

566        single-cell RNA-Seq data. *Nat Commun* **13**, 3224 (2022).

567    26. Armingol, E., Officer, A., Harismendy, O. & Lewis, N. E. Deciphering cell–cell interactions and

568        communication from gene expression. *Nat Rev Genet* **22**, 71–88 (2021).

569    27. Kolda, T. G. & Bader, B. W. Tensor Decompositions and Applications. *SIAM Rev.* **51**, 455–500 (2009).

570    28. Hall, P. & Wilson, S. R. Two Guidelines for Bootstrap Hypothesis Testing. *Biometrics* **47**, 757–762

571        (1991).

572    29. Ye, K. & Lim, L.-H. Schubert Varieties and Distances between Subspaces of Different Dimensions. *SIAM*

573        *J. Matrix Anal. & Appl.* **37**, 1176–1197 (2016).

574    30. Murden, R. J., Zhang, Z., Guo, Y. & Risk, B. B. Interpretive JIVE: Connections with CCA and an

575        application to brain connectivity. *Frontiers in Neuroscience* **16**, (2022).

576    31. Booeshaghi, A. S. & Pachter, L. Normalization of single-cell RNA-seq counts by $\log(x + 1)$† or $\log(1 +$

577        $x)$†. *Bioinformatics* **37**, 2223–2224 (2021).

578    32. Baccin, C. *et al.* Combined single-cell and spatial transcriptomics reveal the molecular, cellular and spatial

579        bone marrow niche organization. *Nat Cell Biol* **22**, 38–48 (2020).

580    33. Edgar, R., Domrachev, M. & Lash, A. E. Gene Expression Omnibus: NCBI gene expression and

581        hybridization array data repository. *Nucleic Acids Res* **30**, 207–210 (2002).

582    34. Speir, M. L. *et al.* UCSC Cell Browser: visualize your single-cell data. *Bioinformatics* **37**, 4578–4580

583        (2021).

584    35. Fang, Z., Liu, X. & Peltz, G. GSEApy: a comprehensive package for performing gene set enrichment

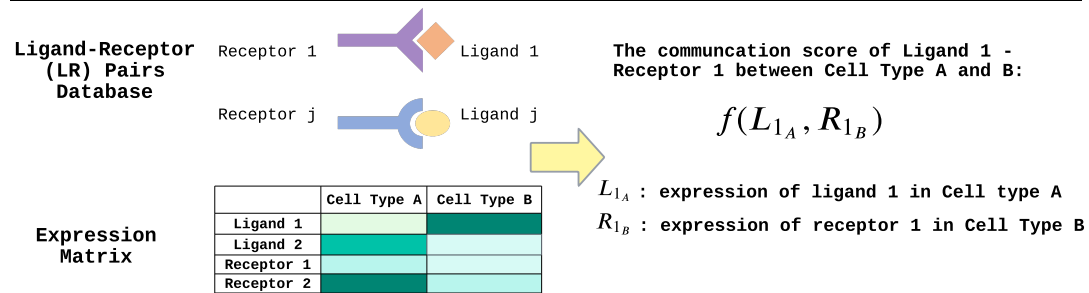585        analysis in Python. *Bioinformatics* **39**, btac757 (2023).
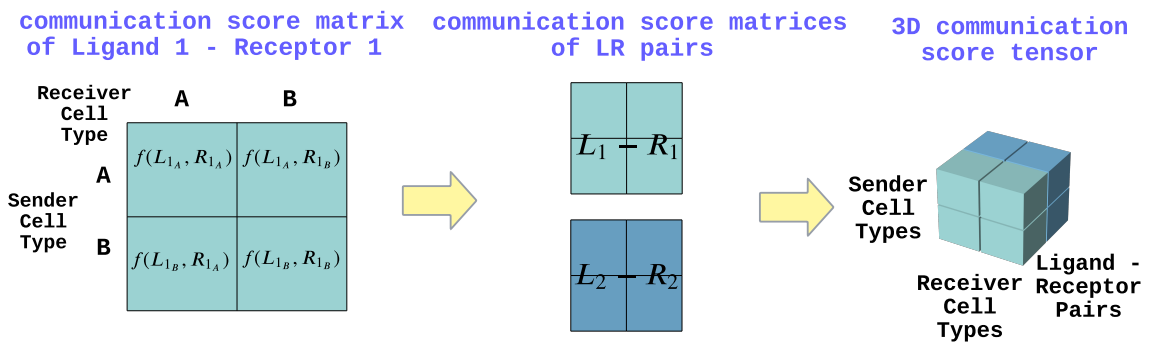
586

587

588

589 **Figures and Tables**

590 **Figure 1. STACCato analytic framework.** (A) Cell-cell communication (CCC) score is given by a function
591 of the expression levels of ligand 1 in sender cell type A ($L_{1_A}$) and receptor 1 in receiver cell type B ($R_{1_B}$). (B)
592 CCC scores are calculated for a specific ligand-receptor pair across all sender and receiver cell types. CCC
593 scores are then organized into a communication score matrix with sender cell types as rows and receiver cell
594 types as columns. Communication score matrices are repeatedly calculated for all ligand-receptor pairs and
595 organized into a 3-dimensional communication score tensor. (C) The 3-dimensional communication score
596 tensors are repeatedly constructed for all samples and then combined into a 4-dimensional communication
597 score tensor. STACCato then uses subject-level information to estimate the coefficient tensor representing the
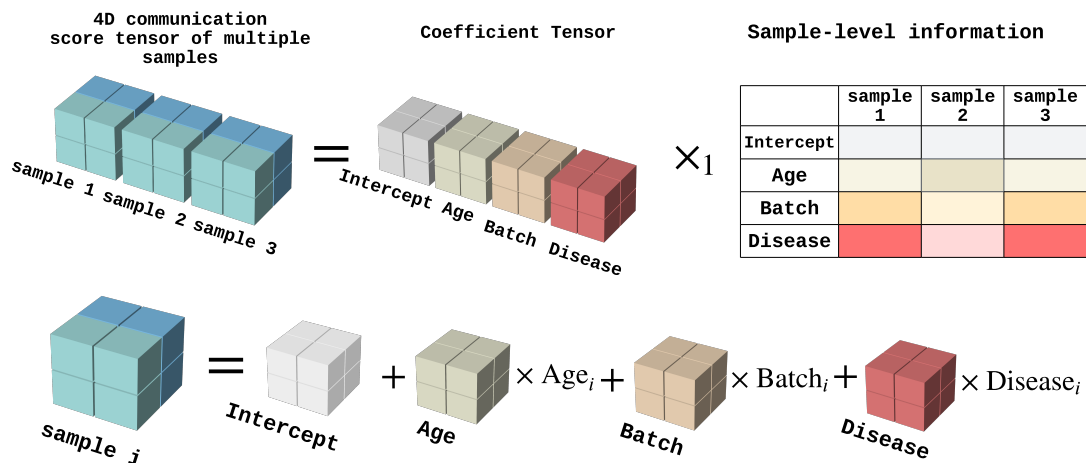


## A. Calculate communication score of one ligand-receptor pair

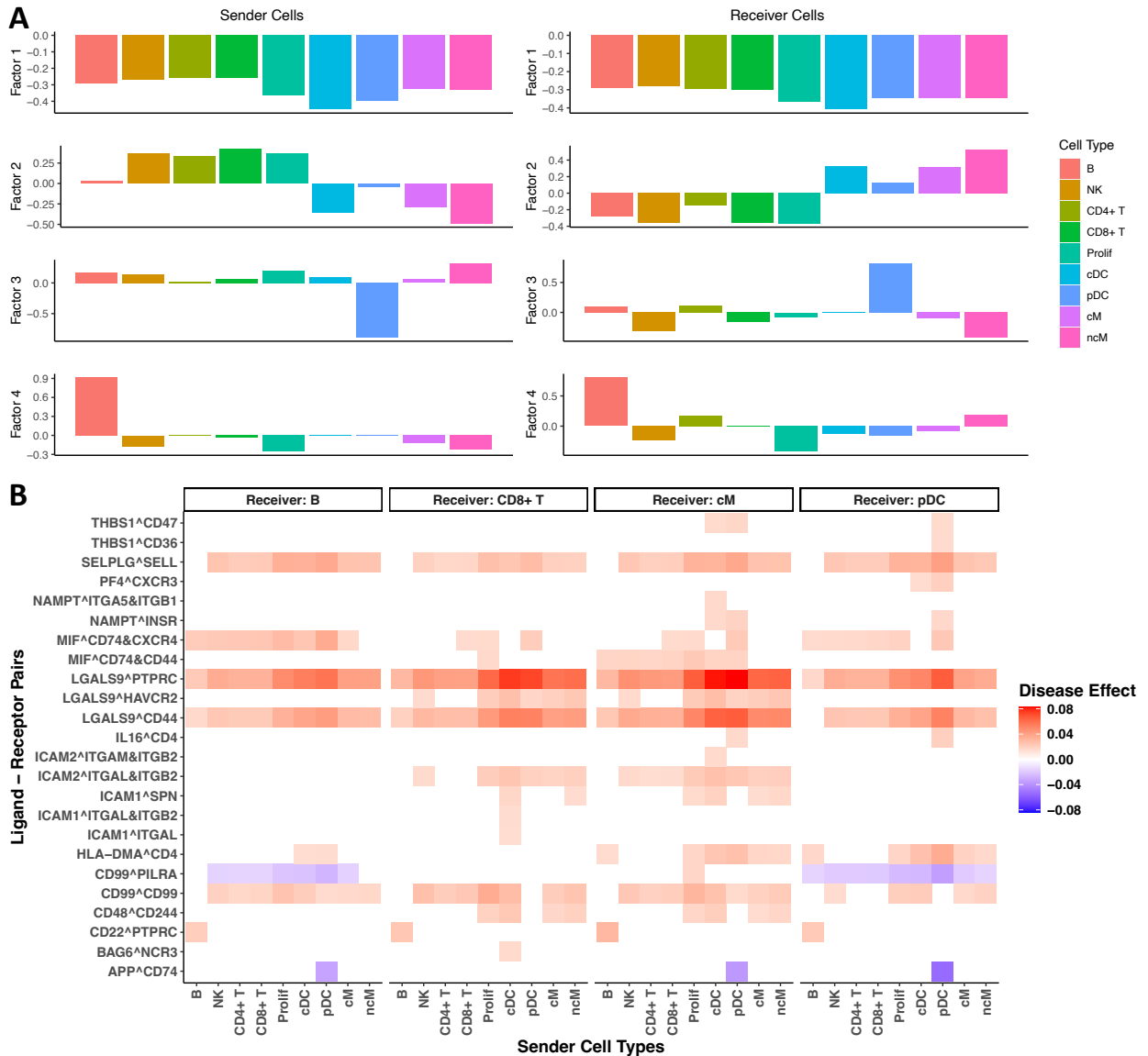## B. Construct 3D communication score tensor for one sample

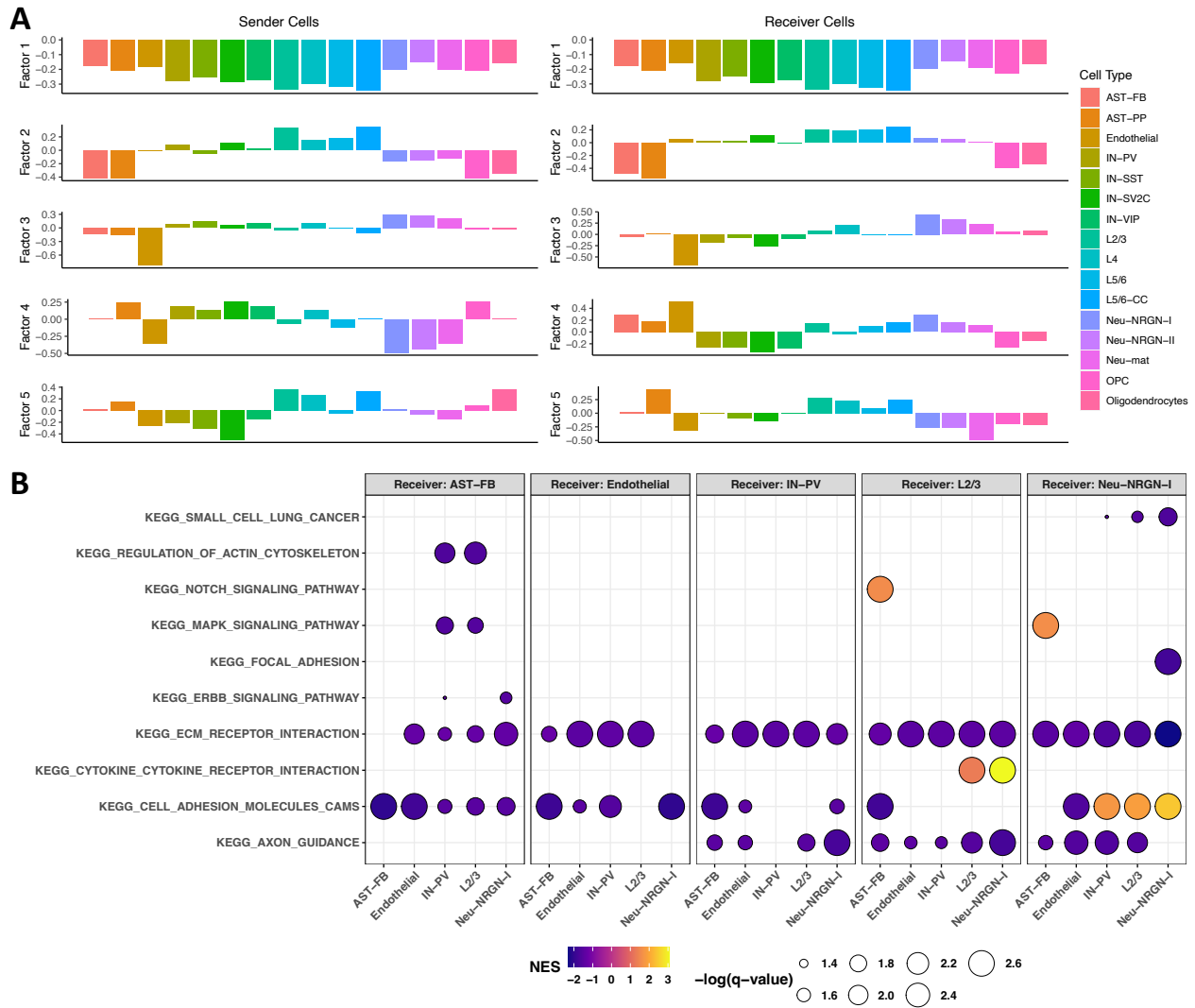## C. Supervised tensor decomposition

598   effects of subject-level variables on CCC events. While this example tensor contains only 2 cell types and 2
599   ligand-receptor pairs, the framework is generalizable to any number of cell types and ligand-receptor pairs.
600   **Figure 2. STACCato results with the SLE dataset.** (A) Bar plots of the estimated values in the factor matrices
601   of sender and receiver cell types. Each color represents one cell type. (B) Estimated significant disease effects
602   with p-values $< 0.05$ and magnitudes $> 0.015$ for communication events with B, CD8$^+$ T, cM, and pDC cells as
603   receiver cell types. Positive disease effects are colored in red while negative disease effects are colored in blue.
604   Positive disease effects indicate positive associations between CCC events and SLE, while negative disease
605   effects indicate negative associations.

606

607 **Figure 3. STACCato results with the ASD dataset.** (A) Bar plots of the estimated values in the factor matrices
608 of sender and receiver cell types. Each color represents one cell type. (B) Significantly enriched KEGG pathways
609 with false discovery rate (FDR) adjusted p-value (q-value) $< 0.05$ across AST-PP, Endothelial, IN-PV, L2/3,
610 and Neu-NRGN-I sender and receiver cell types. Colors represent the normalized enrichment scores. Positive
611 enrichment scores indicate positive associations with ASD, while negative enrichment score indicate negative
612 associations with ASD.

613

**Figure 4. STACCato simulation results:** MSE of estimated disease effects (A) and chordal distance of estimated factor matrices (B) in balanced, moderate unbalanced, and extreme unbalanced scenarios. The bar plot shows the average MSEs across 100 simulations from Model 1 considering disease status and batch (red bars) and Model 2 considering disease status only (green bars) with black error bars showing standard errors.