



1

<sup>2</sup> **Supporting Information for**

<sup>3</sup> **Cooperative Learning for Multiview Analysis**

<sup>4</sup> **Daisy Yi Ding, Shuangning Li, Balasubramanian Narasimhan and Robert Tibshirani**

<sup>5</sup> **Robert Tibshirani.**

<sup>6</sup> **E-mail:** tibs@stanford.edu

<sup>7</sup> **This PDF file includes:**

<sup>8</sup> Supporting text

<sup>9</sup> Figs. S1 to S9

10 **Supporting Information Text**

11 **1. Adaptive cooperative learning**

12 In this section, we outline an adaptive strategy for optimizing over  $\lambda_x$  and  $\lambda_z$  for different data views. We call this adaptive  
 13 cooperative learning. The method incorporates the values of  $\lambda_x$  and  $\lambda_z$  that have been adaptively optimized by the one-at-a-time  
 14 algorithm (Algorithm S1) as a penalty factor in the direct algorithm (Algorithm S2). In the two-dimensional grid of  $\lambda_x$  and  $\lambda_z$ ,  
 15 our proposed strategy works by iteratively searching along one axis of  $\lambda$  while fixing the other constant.

**Algorithm S1** *One-at-a-time algorithm for adaptive cooperative learning (regression).*

**Input:**  $X \in \mathcal{R}^{n \times p_x}$  and  $Z \in \mathcal{R}^{n \times p_z}$ , the response  $y \in \mathcal{R}^n$ , and a fixed hyperparameter  $\rho \in \mathcal{R}$ .

**Output:**  $\hat{\theta}_x$  and  $\hat{\theta}_z$  from the last iteration, along with the hyperparameters  $\lambda_x^*$  and  $\lambda_z^*$  and the corresponding CV errors.

1. Initialize  $\theta_x^{(0)} \in \mathcal{R}^{p_x}$  and  $\theta_z^{(0)} \in \mathcal{R}^{p_z}$ .
2. For  $k = 0, 1, 2, \dots$  until convergence:
  - (a) Set  $y_x^* = \frac{y}{1+\rho} - \frac{(1-\rho)Z\theta_z}{(1+\rho)}$ . Solve Lasso( $X, y_x^*, \lambda$ ) over a decreasing grid of  $\lambda$  values. Update  $\theta_x^{(k+1)}$  to be the solution and record the hyperparameter  $\lambda_x^*$  that minimizes the CV error.
  - (b) Set  $y_z^* = \frac{y}{1+\rho} - \frac{(1-\rho)X\theta_x}{(1+\rho)}$ . Solve Lasso( $Z, y_z^*, \lambda$ ) over a decreasing grid of  $\lambda$  values. Update  $\theta_z^{(k+1)}$  to be the solution and record the hyperparameter  $\lambda_z^*$  that minimizes the CV error.

**Algorithm S2** *Direct algorithm for adaptive cooperative learning (regression).*

**Input:**  $X \in \mathcal{R}^{n \times p_x}$  and  $Z \in \mathcal{R}^{n \times p_z}$ , the response  $y \in \mathcal{R}^n$ , and a grid of hyperparameter values  $(\rho_{\min}, \dots, \rho_{\max})$ .

**for**  $\rho \leftarrow \rho_{\min}, \dots, \rho_{\max}$  **do**

Run Algorithm S1 with both (X,Z) and (Z,X) with the same folds for CV. Select the one with the lower sum of the two CV errors. Get the corresponding  $\lambda_x^*$  and  $\lambda_z^*$ .

Set

$$\tilde{X} = \begin{pmatrix} X & Z \\ -\sqrt{\rho}X & \sqrt{\rho}Z \end{pmatrix}, \tilde{y} = \begin{pmatrix} y \\ \mathbf{0} \end{pmatrix}.$$

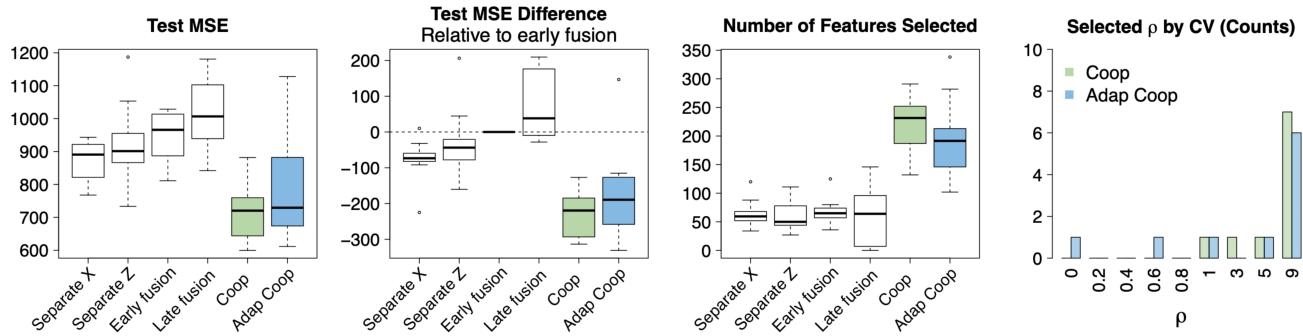
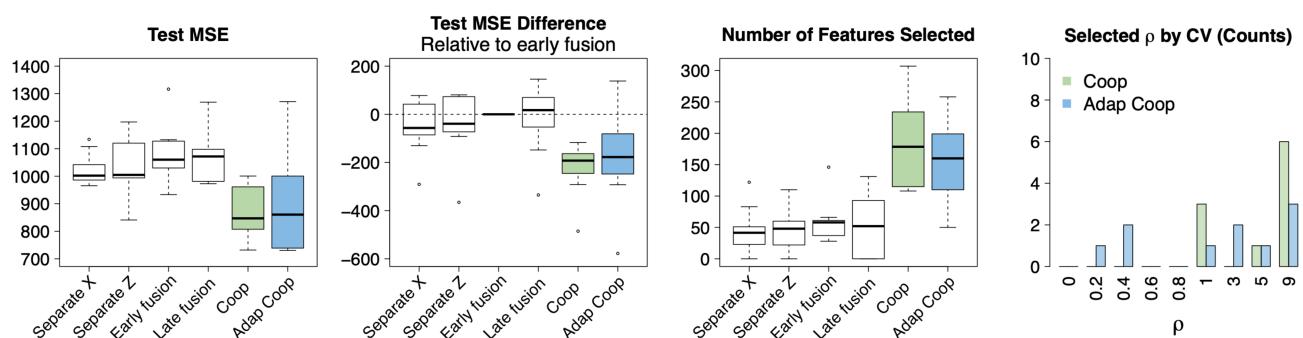
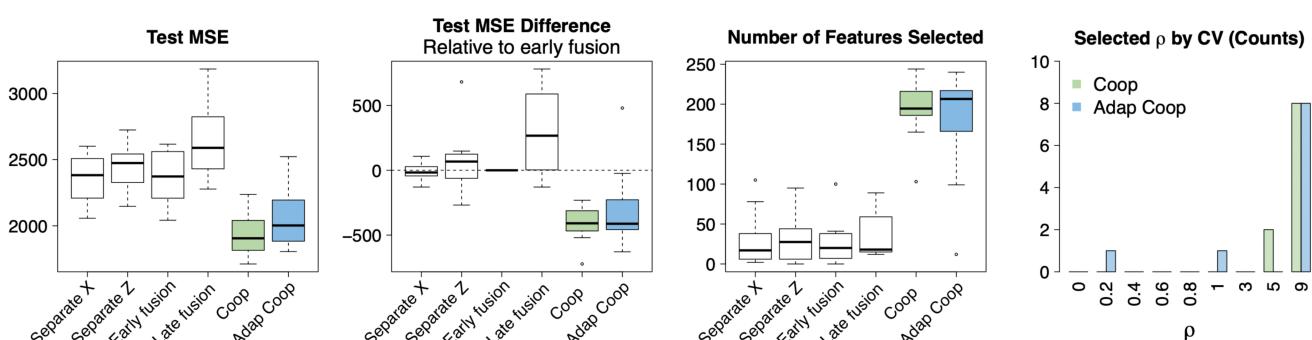
Solve Lasso( $\tilde{X}, \tilde{y}, \lambda$ ) over a decreasing grid of  $\lambda$  values, with a penalty factor of  $(1, \dots, 1, \frac{\lambda_z^*}{\lambda_x^*}, \dots, \frac{\lambda_z^*}{\lambda_x^*})$ . Note that we form folds from the rows of X and Z and then construct the corresponding  $\tilde{X}$ .

**end**

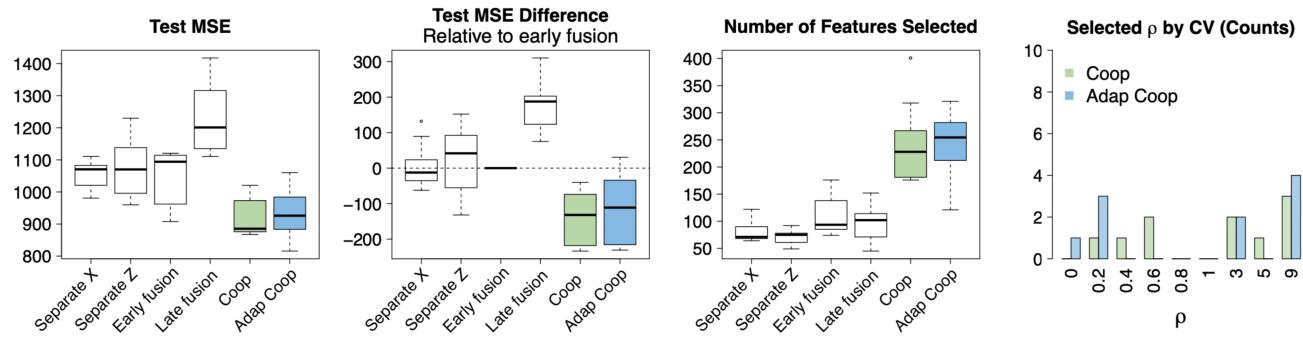
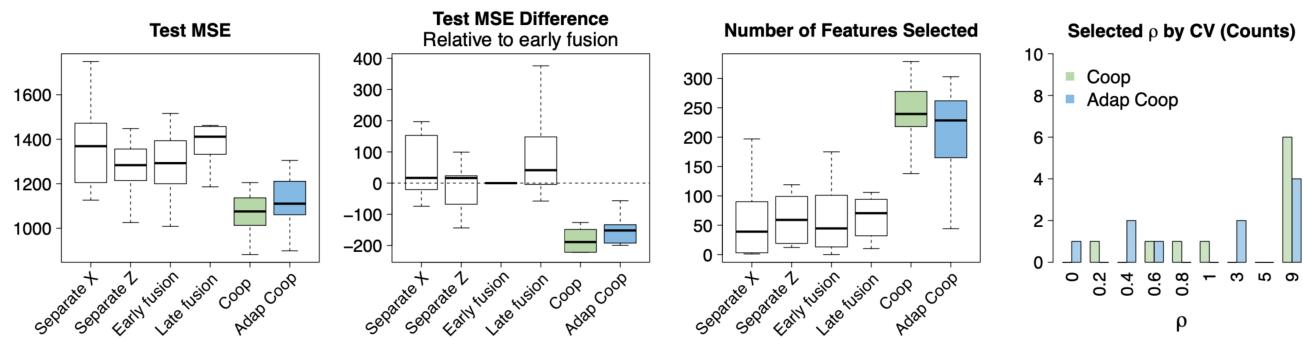
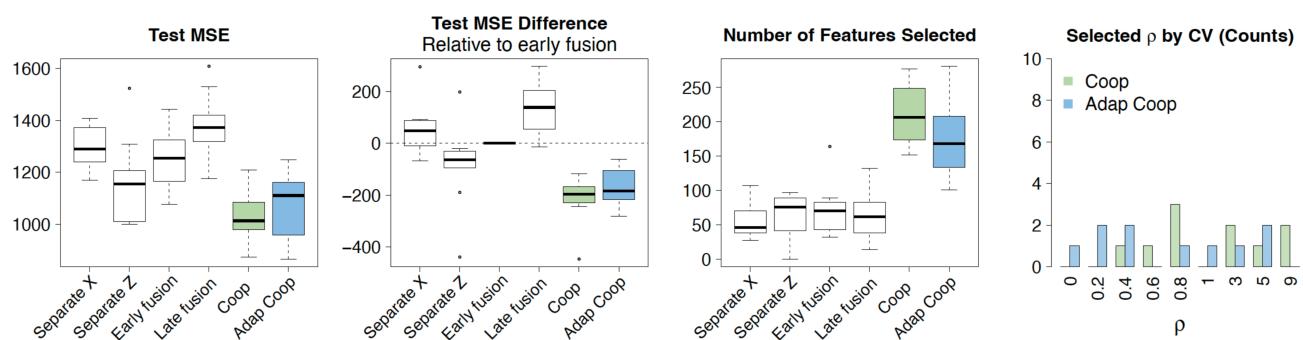
Select the optimal value of  $\rho$  based on the CV error and get the final fit.

16 **2. More comprehensive simulation studies on cooperative regularized regression**

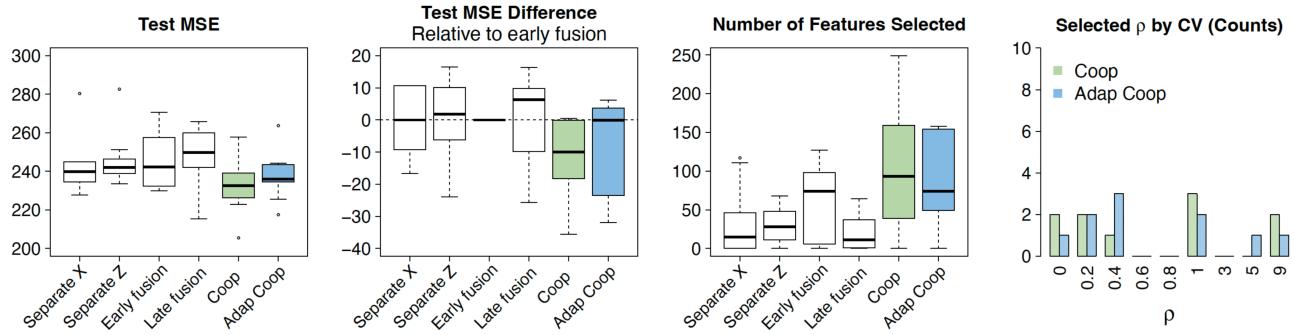
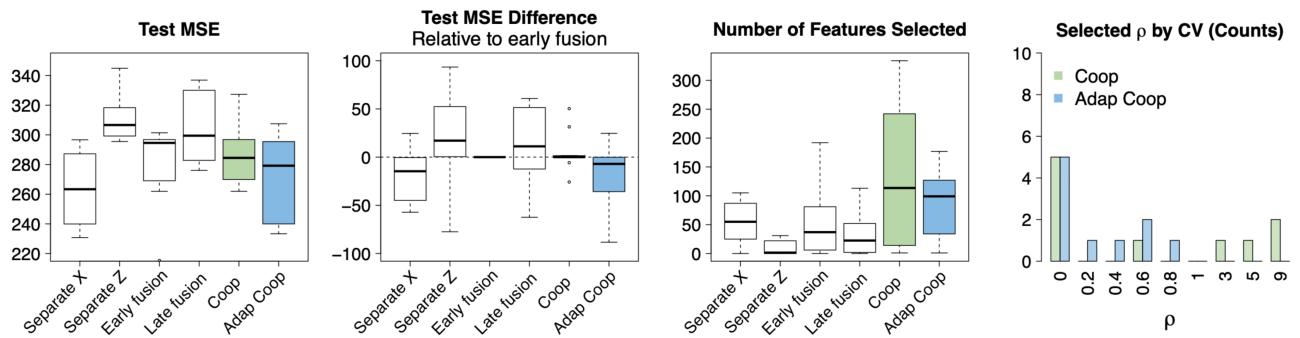
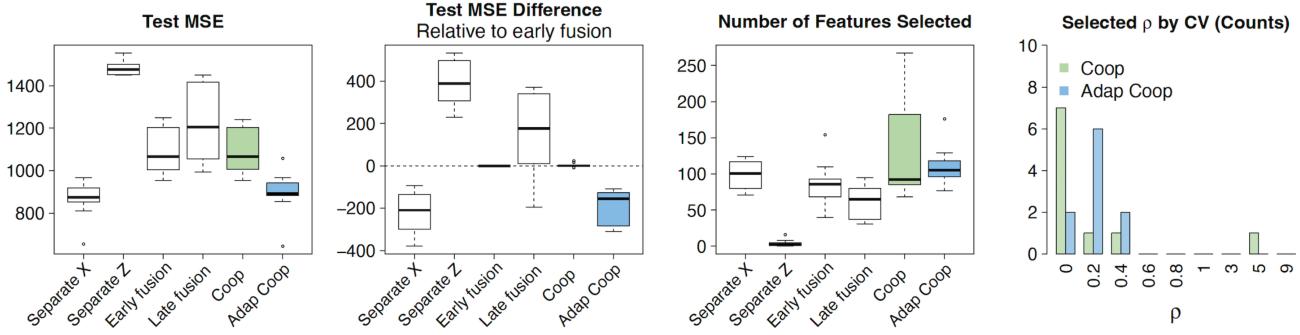
17 Fig. S1-S3 show more simulation results of the high-dimensional settings ( $p = 1000, n = 200$ ), and Fig. S4-S6 show the  
 18 simulation results of the lower-dimensional settings ( $p = 200, n = 500$ ).

**A****B****C**

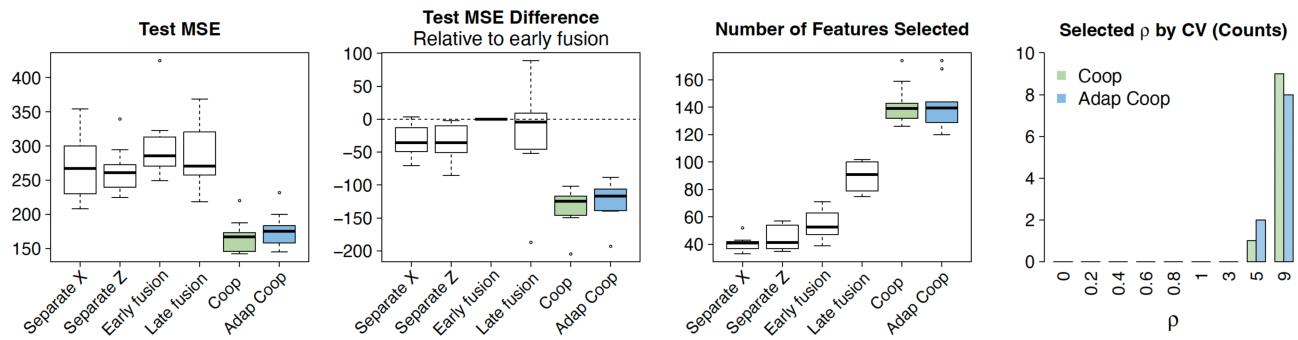
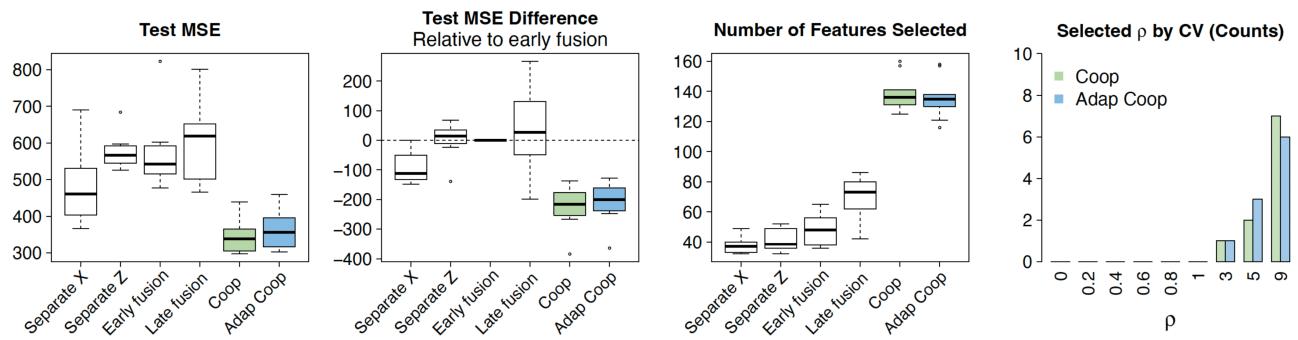
**Fig. S1.** Simulation studies on cooperative regularized linear regression when  $X$  and  $Z$  are high-dimensional and have a high level of correlation with each other. (A) Simulation results when  $X$  and  $Z$  have a high level of correlation and both contain signal ( $t_x = t_z = 6$ ),  $n = 200$ ,  $p = 1000$ , SNR = 1.0. The first panel shows MSE on a test test; the second panel shows the MSE difference on the test set relative to early fusion; the third panel shows the number of features selected; the fourth panel shows the  $\rho$  values selected by CV in cooperative learning. Here “Coop” refers to cooperative learning outlined in Algorithm 1 and “Adap Coop” refers to adaptive cooperative learning outlined in Algorithm S2. (B) Simulation results when  $X$  and  $Z$  have a high level of correlation and both contain signal ( $t_x = t_z = 6$ ),  $n = 200$ ,  $p = 1000$ , SNR = 0.6. (C) Simulation results when  $X$  and  $Z$  have a high level of correlation,  $X$  contains more signal than  $Z$  ( $t_x = 4$ ,  $t_z = 2$ ),  $n = 200$ ,  $p = 1000$ , SNR = 0.6.

**A****B****C**

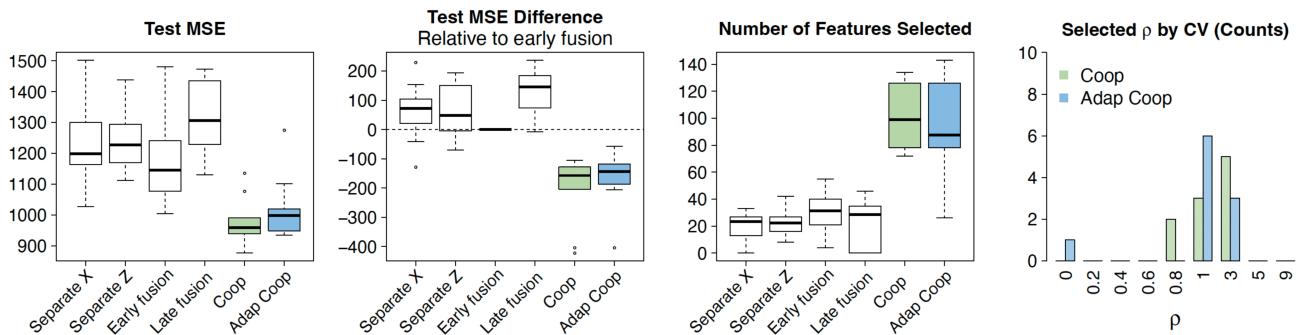
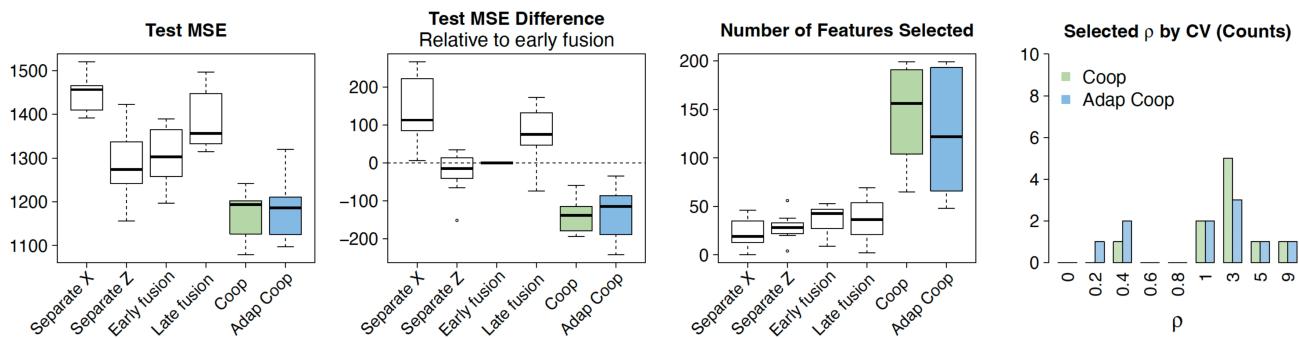
**Fig. S2.** Simulation studies on cooperative regularized linear regression when  $X$  and  $Z$  are high-dimensional and have a medium level of correlation with each other. (A) Simulation results when  $X$  and  $Z$  have a medium level of correlation and both contain signal ( $t_x = t_z = 2$ ),  $n = 200$ ,  $p = 1000$ , SNR = 3.5. The setup is the same as in Figure S1. (B) Simulation results when  $X$  and  $Z$  have a medium level of correlation and both contain signal ( $t_x = t_z = 2$ ),  $n = 200$ ,  $p = 1000$ , SNR = 1.6. (C) Simulation results when  $X$  and  $Z$  have a medium level of correlation, and  $Z$  contains more signal than  $X$  ( $t_x = 2$ ,  $t_z = 3$ ),  $n = 200$ ,  $p = 1000$ , SNR = 1.5.

**A****B****C**

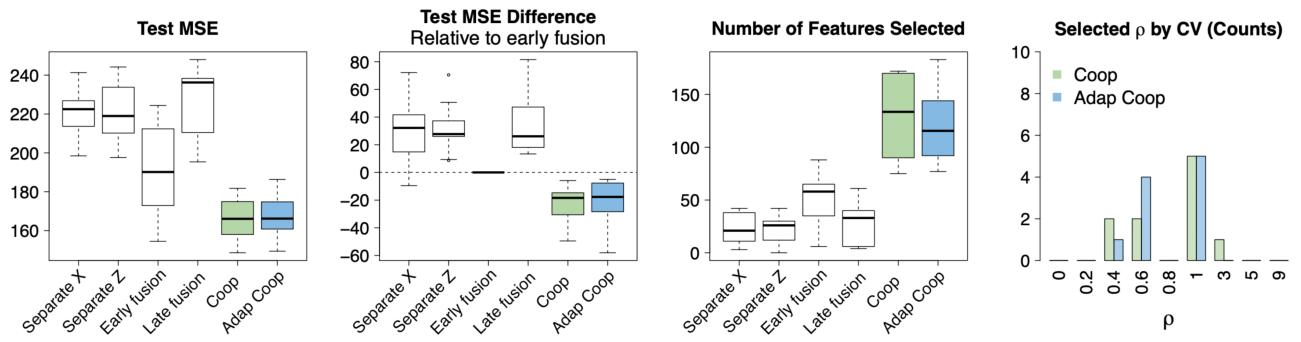
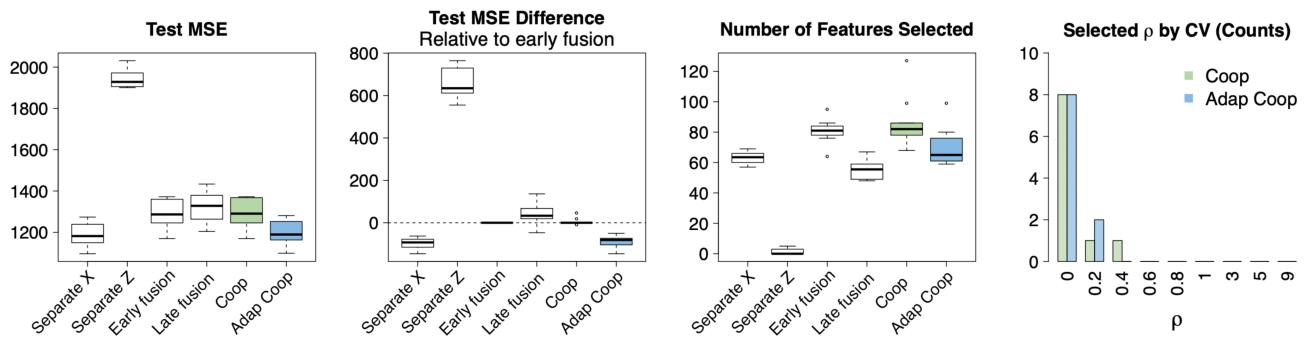
**Fig. S3.** Simulation studies on cooperative regularized linear regression when  $X$  and  $Z$  are high-dimensional and have no correlation. (A) Simulation results when  $X$  and  $Z$  have no correlation, and both  $X$  and  $Z$  contain signal (here we generated  $y$  as a linear combination of  $X$  and  $Z$  instead of the latent factors),  $n = 200$ ,  $p = 1000$ , SNR = 1.0. The setup is the same as in Figure S1. (B) Simulation results when  $X$  and  $Z$  have no correlation;  $X$  contains more signal than  $Z$  (here we generated  $y$  as a linear combination of  $X$  and  $Z$  instead of the latent factors),  $n = 200$ ,  $p = 1000$ , SNR = 1.1. (C) Simulation results when  $X$  and  $Z$  have no correlation; only  $X$  contains signal ( $t_x = 2, t_z = 0$ ),  $n = 200$ ,  $p = 1000$ , SNR = 3.5.

**A****B**

**Fig. S4.** Simulation studies on cooperative regularized linear regression when  $X$  and  $Z$  are of a lower dimension and have a high level of correlation with each other. (A) Simulation results when  $X$  and  $Z$  have a high level of correlation and both contain signal ( $t_x = t_z = 6$ ),  $n = 500$ ,  $p = 200$ , SNR = 1.2. The first panel shows MSE on a test set; the second panel shows the MSE difference on the test set relative to early fusion; the third panel shows the number of features selected; the fourth panel shows the  $\rho$  values selected by CV in cooperative learning. Here “Coop” refers to cooperative learning outlined in Algorithm 1 and “Adap Coop” refers to adaptive cooperative learning outlined in Algorithm S2. (B) Simulation results when  $X$  and  $Z$  have a high level of correlation and  $X$  contains more signal than  $Z$  ( $t_x = 5$ ,  $t_z = 3$ ),  $n = 500$ ,  $p = 200$ , SNR = 0.7.

**A****B**

**Fig. S5.** Simulation studies on cooperative regularized linear regression when  $X$  and  $Z$  are of a lower dimension and have a medium level of correlation with each other. (A) Simulation results when  $X$  and  $Z$  have a medium level of correlation and both contain signal ( $t_x = t_z = 1$ ),  $n = 500$ ,  $p = 200$ , SNR = 0.8. The setup is the same as in Figure S4. (B) Simulation results when  $X$  and  $Z$  have a medium level of correlation, and  $Z$  contains more signal than  $X$  ( $t_x = 0.6$ ,  $t_z = 0.9$ ),  $n = 500$ ,  $p = 200$ , SNR = 0.5.

**A****B**

**Fig. S6.** Simulation studies on cooperative regularized linear regression when  $X$  and  $Z$  are of a lower dimension and have no correlation with each other. (A) Simulation results when  $X$  and  $Z$  have no correlation ( $t_x = t_z = 0$ ), and both  $X$  and  $Z$  contain signal (here we generated  $y$  as a linear combination of  $X$  and  $Z$  instead of the latent factors),  $n = 500$ ,  $p = 200$ , SNR = 0.3. The setup is the same as in Figure S4. (B) Simulation results when  $X$  and  $Z$  have no correlation and only  $X$  contains signal ( $t_x = 2$ ,  $t_z = 0$ ),  $n = 500$ ,  $p = 200$ , SNR = 3.0.

19 **3. Simulation studies on cooperative learning for more than two data views**

20 Here we conduct simulation studies on cooperative learning for more than two data views. Specifically, we consider the setting  
 21 of three data views, and this generalizes easily to more data views. We generated Gaussian data with  $n = 200$  and  $p = 300$  in  
 22 each of the views  $X_1$ ,  $X_2$  and  $X_3$ , and created correlation between them using latent factors. The response  $\mathbf{y}$  was generated as  
 23 a linear combination of the latent factors, corrupted by Gaussian noise.

24 **A. Simulation procedure.** The simulation for 3 data views is set up as follows. Given values for parameters  $n, p_{x_1}, p_{x_2},$   
 25  $p_{x_3}, p_u, s_u, t_{x_1}, t_{x_2}, t_{x_3}, \beta_u, \sigma$ , we generate data according to the following procedure:

26 1.  $x_{1j} \in \mathcal{R}^n$  distributed i.i.d.  $\text{MVN}(0, I_n)$  for  $j = 1, 2, \dots, p_{x_1}$ .

27 2.  $x_{2j} \in \mathcal{R}^n$  distributed i.i.d.  $\text{MVN}(0, I_n)$  for  $j = 1, 2, \dots, p_{x_2}$ .

28 3.  $x_{3j} \in \mathcal{R}^n$  distributed i.i.d.  $\text{MVN}(0, I_n)$  for  $j = 1, 2, \dots, p_{x_3}$ .

29 4. For  $i = 1, 2, \dots, p_u$  ( $p_u$  corresponds to the number of latent factors):

30 (a)  $u_i \in \mathcal{R}^n$  distributed i.i.d.  $\text{MVN}(0, s_u^2 I_n)$ ;

31 (b)  $x_{1i} = x_{1i} + t_{x_1} * u_i$ ;

32 (c)  $x_{2i} = x_{2i} + t_{x_2} * u_i$ ;

33 (d)  $x_{3i} = x_{3i} + t_{x_3} * u_i$ .

34 5.  $X_1 = [x_{11}, x_{12}, \dots, x_{1p_{x_1}}]$ ,  $X_2 = [x_{21}, x_{22}, \dots, x_{2p_{x_2}}]$ ,  $X_3 = [x_{31}, x_{32}, \dots, x_{3p_{x_3}}]$ .

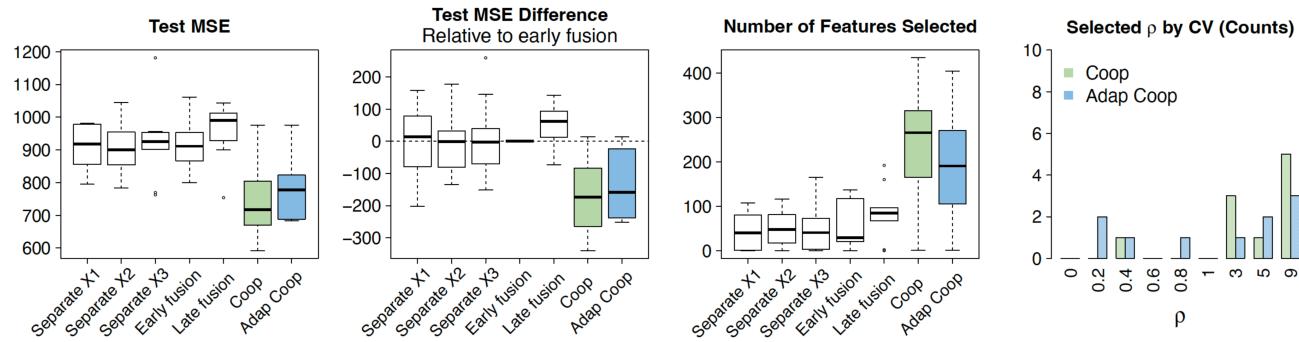
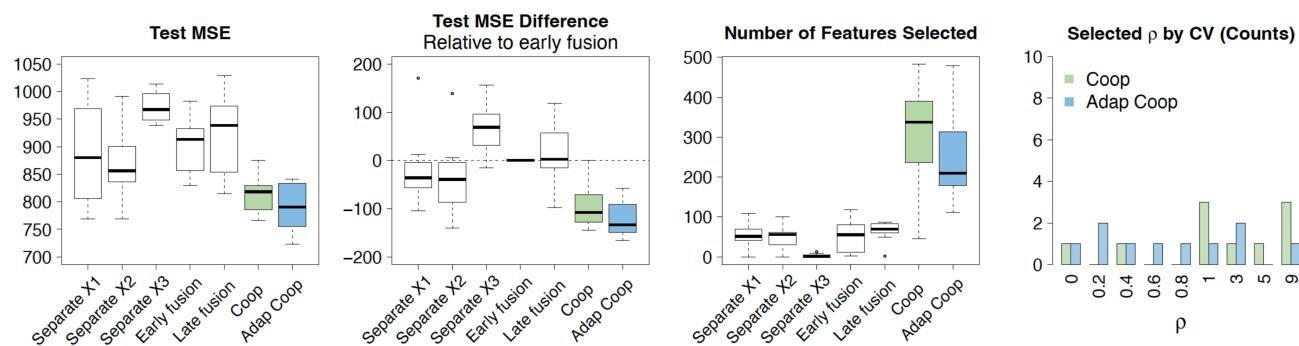
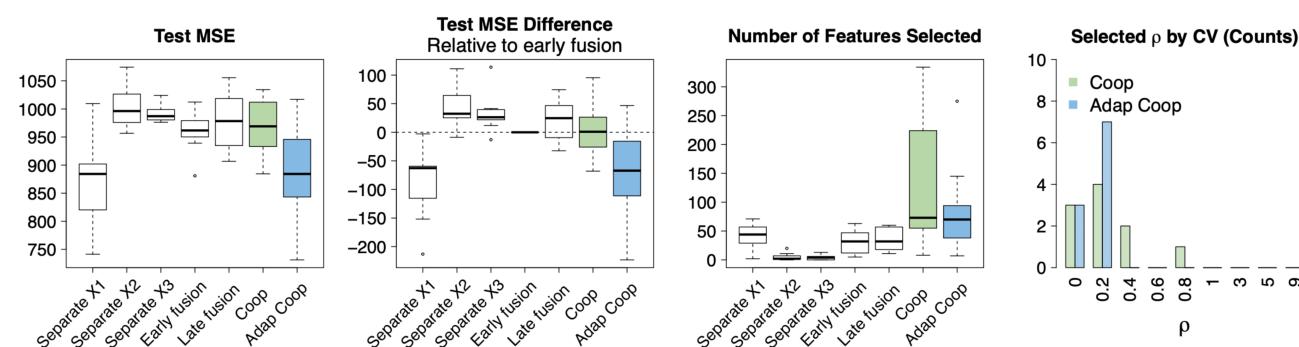
35 6.  $U = [u_1, u_2, \dots, u_{p_u}]$ ,  $\mathbf{y} = U\beta_u + \epsilon$  where  $\epsilon \in \mathcal{R}^n$  distributed i.i.d.  $\text{MVN}(0, \sigma^2 I_n)$ .

36 We compare the following methods: (1) separate  $X_1$ , separate  $X_2$ , and separate  $X_3$ : the standard lasso is applied on the  
 37 separate data views of  $X_1$ ,  $X_2$  and  $X_3$  with 10-fold CV; (2) early fusion: the standard lasso is applied on the concatenated  
 38 data views of  $X_1$ ,  $X_2$  and  $X_3$  with 10-fold CV (note that this is equivalent to cooperative learning with  $\rho = 0$ ); (3) late fusion:  
 39 separate lasso models are first fitted on  $X_1$ ,  $X_2$  and  $X_3$  independently with 10-fold CV, and the three resulting predictors  
 40 are then combined through linear least squares for the final prediction; (4) cooperative learning (regression) and adaptive  
 41 cooperative learning.

42 We evaluated the performance based on the mean-squared error (MSE) on a test set and conducted each simulation  
 43 experiment 10 times.

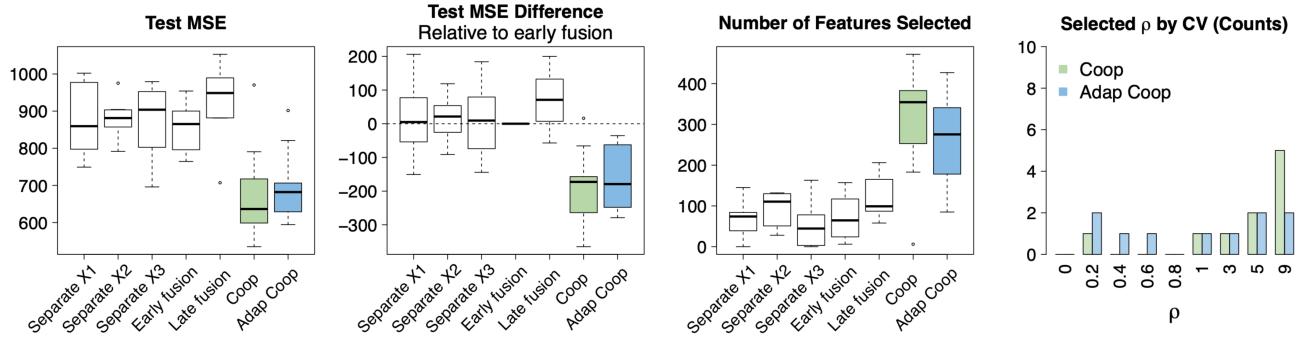
44 **B. Simulation results.** Fig. S7-S8 show the simulation results for 3 data views. Overall, the simulation results can be summarized  
 45 as follows:

- 46 • Cooperative learning performs the best in terms of test MSE across the range of SNR and correlation settings. It is most  
 helpful when the data views are correlated and contain signal. When the correlation between data views is higher, higher  
 values of  $\rho$  are more likely to be selected.
- 49 • When only two data views are correlated and contain signal (as in Fig. S7B and Fig. S8C), cooperative learning also gives  
 50 performance gains by leveraging the correlation through the agreement penalty, while early fusion can be outperformed  
 51 by the separate models fit on the data views containing the signal.
- 52 • When only one view contains signal and the views are not correlated (as in Fig. S7C), cooperative learning is outperformed  
 53 by the separate model fit on the view containing the signal, but adaptive cooperative learning is able to perform on par  
 54 with the separate model, outperforming early and late fusion.

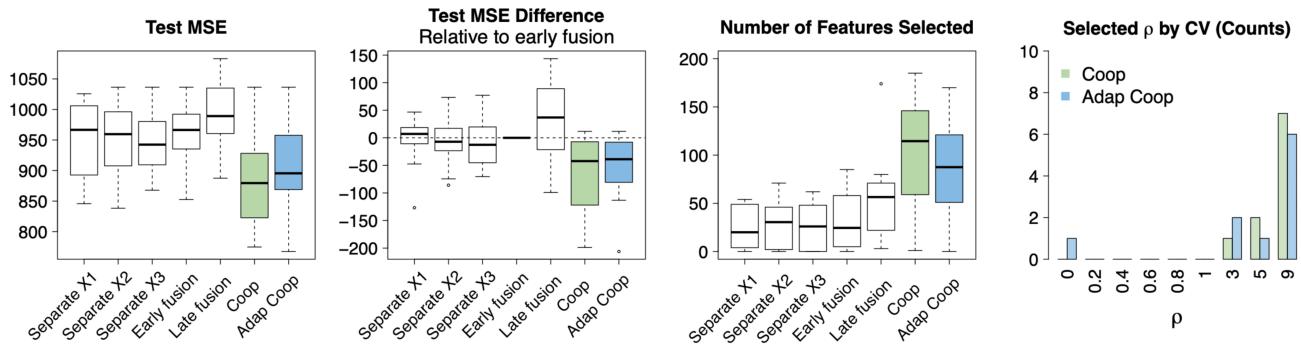
**A****B****C**

**Fig. S7.** Simulation studies on cooperative regularized linear regression for more than two data views. (A) Simulation results when  $X_1$ ,  $X_2$  and  $X_3$  are correlated and all contain signal ( $t_{x_1} = t_{x_2} = t_{x_3} = 2$ ),  $n = 200$ ,  $p = 900$ , SNR = 1.5. The first panel shows MSE on a test set; the second panel shows the MSE difference on the test set relative to early fusion; the third panel shows the number of features selected; the fourth panel shows the  $\rho$  values selected by CV in cooperative learning. Here “Coop” refers to cooperative learning outlined in Algorithm 1 and “Adap Coop” refers to adaptive cooperative learning outlined in Algorithm S2. (B) Simulation results when only  $X_1$  and  $X_2$  are correlated and contain signal ( $t_{x_1} = t_{x_2} = 2$ ,  $t_{x_3} = 0$ ),  $n = 200$ ,  $p = 900$ , SNR = 1.5. (C) Simulation results when  $X_1$ ,  $X_2$  and  $X_3$  are uncorrelated, and only  $X_1$  contains signal ( $t_{x_1} = 2$ ,  $t_{x_2} = t_{x_3} = 0$ ),  $n = 200$ ,  $p = 900$ , SNR = 1.5.

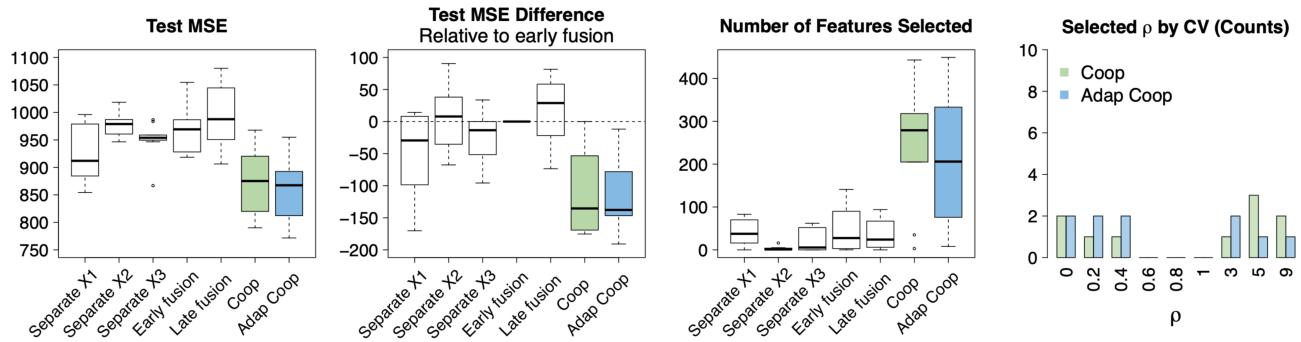
A



B



C



**Fig. S8.** Simulation studies on cooperative regularized linear regression for more than two data views. (A) Simulation results when  $X_1, X_2$  and  $X_3$  are correlated and all contain signal ( $t_{x_1} = t_{x_2} = t_{x_3} = 2$ ),  $n = 200$ ,  $p = 900$ , SNR = 2.5. The first panel shows MSE on a test test; the second panel shows the MSE difference on the test set relative to early fusion; the third panel shows the number of features selected; the fourth panel shows the  $\rho$  values selected by CV in cooperative learning. Here “Coop” refers to cooperative learning outlined in Algorithm 1 and “Adap Coop” refers to adaptive cooperative learning outlined in Algorithm S2. (B) Simulation results when  $X_1, X_2$  and  $X_3$  are correlated and all contain signal ( $t_{x_1} = t_{x_2} = t_{x_3} = 2$ ),  $n = 200$ ,  $p = 900$ , SNR = 0.6. (c) Simulation results when only  $X_1$  and  $X_3$  are correlated;  $X_1$  contains more signal than  $X_3$ ,  $X_2$  does not contain signal ( $t_{x_1} = 2, t_{x_2} = 0, t_{x_3} = 1.5$ ),  $n = 200$ ,  $p = 900$ , SNR = 1.0.

55 **4. Theoretical analysis under the factor model**

56 To understand the role of the agreement penalty from a theoretical perspective, we consider the following latent factor model.  
 57 Let  $\mathbf{u} = (U_1, U_2, \dots, U_n)$  be a vector of  $n$  i.i.d. random variables with  $U_i \sim \mathcal{N}(0, 1)$ . Let  $\mathbf{y} = (Y_1, \dots, Y_n)$ ,  $\mathbf{x} = (X_1, \dots, X_n)$ ,  
 58 and  $\mathbf{z} = (Z_1, \dots, Z_n)$ , with

59 
$$Y_i = \gamma_y U_i + \varepsilon_{yi}, \quad X_i = \gamma_x U_i + \varepsilon_{xi}, \quad \text{and} \quad Z_i = \gamma_z U_i + \varepsilon_{zi}, \quad [1]$$

60 where  $\varepsilon_{yi} \sim \mathcal{N}(0, \sigma_y^2)$ ,  $\varepsilon_{xi} \sim \mathcal{N}(0, \sigma_x^2)$ ,  $\varepsilon_{zi} \sim \mathcal{N}(0, \sigma_z^2)$  independently.

61 In this section, we study the mean squared error (MSE) of the cooperative learning algorithm. More precisely, let

62 
$$\hat{\theta} = \operatorname{argmin}_{\theta} \sum_{i=1}^n \left[ \frac{1}{2} (Y_i - X_i \theta_x - Z_i \theta_z)^2 + \frac{\rho}{2} (X_i \theta_x - Z_i \theta_z)^2 \right]. \quad [2]$$

63 Let  $U_{\text{new}}$ ,  $X_{\text{new}}$ ,  $Z_{\text{new}}$  be some new random variables generated from Eq. (1) independently of the previous data, i.e.,

64 
$$Y_{\text{new}} = \gamma_y U_{\text{new}} + \varepsilon_{y \text{ new}}, \quad X_{\text{new}} = \gamma_x U_{\text{new}} + \varepsilon_{x \text{ new}}, \quad \text{and} \quad Z_{\text{new}} = \gamma_z U_{\text{new}} + \varepsilon_{z \text{ new}}, \quad [3]$$

65 where  $U_{\text{new}} \sim \mathcal{N}(0, 1)$ ,  $\varepsilon_{y \text{ new}} \sim \mathcal{N}(0, \sigma_y^2)$ ,  $\varepsilon_{x \text{ new}} \sim \mathcal{N}(0, \sigma_x^2)$ ,  $\varepsilon_{z \text{ new}} \sim \mathcal{N}(0, \sigma_z^2)$  independently. We focus on the MSE  
 66 conditioning on  $\mathbf{x}$  and  $\mathbf{z}$ :

67 
$$\operatorname{MSE}(\mathbf{x}, \mathbf{z}; \rho) = E \left[ (Y_{\text{new}} - (X_{\text{new}} \hat{\theta}_x + Z_{\text{new}} \hat{\theta}_z))^2 | \mathbf{x}, \mathbf{z} \right]. \quad [4]$$

68 The case of  $\rho = 0$  corresponds to the linear regression with no agreement penalty. We will study the behavior of  $\operatorname{MSE}(\mathbf{x}, \mathbf{z}; \rho)$   
 69 when  $\rho$  is around 0.

70 **Proposition 1.** *The derivative of  $\operatorname{MSE}(\mathbf{x}, \mathbf{z}; \rho)$  satisfies*

71 
$$\frac{d}{d\rho} [\operatorname{MSE}(\mathbf{x}, \mathbf{z}; \rho)] |_{\rho=0} = \sigma^{*2} (C_2 B_1 - 2C_1 B_2) / C_2^3, \quad [5]$$

72 where

73 
$$\begin{aligned} \sigma^{*2} &= \frac{\gamma_y^2}{1 + \gamma_x^2/\sigma_x^2 + \gamma_z^2/\sigma_z^2} + \sigma_y^2, \\ C_1 &= [(\gamma_x^2 + \sigma_x^2)(\mathbf{z}^\top \mathbf{z}) + (\gamma_z^2 + \sigma_z^2)(\mathbf{x}^\top \mathbf{x}) - 2\gamma_x \gamma_z (\mathbf{x}^\top \mathbf{z})] ((\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2), \\ B_1 &= 2 [(\gamma_x^2 + \sigma_x^2)(\mathbf{z}^\top \mathbf{z}) + (\gamma_z^2 + \sigma_z^2)(\mathbf{x}^\top \mathbf{x}) + 2\gamma_x \gamma_z (\mathbf{x}^\top \mathbf{z})] ((\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2), \\ C_2 &= (\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2, \\ B_2 &= 2 ((\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) + (\mathbf{x}^\top \mathbf{z})^2). \end{aligned} \quad [6]$$

74 **Proposition 2.** *The derivative of  $\operatorname{MSE}(\mathbf{x}, \mathbf{z}; \rho)$  at  $\rho = 0$  satisfies*

75 
$$\frac{d}{d\rho} [\operatorname{MSE}(\mathbf{x}, \mathbf{z}; \rho)] |_{\rho=0} = -\frac{4}{n} \left( 1 + \frac{2\gamma_x^2 \gamma_z^2}{\sigma_x^2 \gamma_z^2 + \sigma_z^2 \gamma_x^2 + \sigma_x^2 \sigma_z^2} \right) \left( \sigma_y^2 + \frac{\gamma_y^2 \sigma_x^2 \sigma_z^2}{\sigma_x^2 \gamma_z^2 + \sigma_z^2 \gamma_x^2 + \sigma_x^2 \sigma_z^2} \right) + \mathcal{O}_p(n^{-\frac{3}{2}}). \quad [7]$$

76 Here the notation  $\mathcal{O}_p(\cdot)$  is used with the following meaning:  $X_n = \mathcal{O}_p(a_n)$  as  $n \rightarrow \infty$  means that for any  $\varepsilon > 0$ , there exists a  
 77 finite  $M > 0$  and a finite  $N > 0$  such that  $P[|X_n/a_n| > M] < \varepsilon, \forall n > N$ .

78 The proposition establishes that the MSE is a decreasing function of  $\rho$  around 0 with high probability. In other words, the  
 79 agreement penalty is helpful in reducing the mean squared error. To further interpret the above results, we study the ratio of  
 80 the derivative to the MSE itself.

81 **Proposition 3.** *The ratio of the derivative of  $\operatorname{MSE}(\mathbf{x}, \mathbf{z}; \rho)$  at  $\rho = 0$  and  $\operatorname{MSE}(\mathbf{x}, \mathbf{z}; 0)$  satisfies*

82 
$$\frac{\frac{d}{d\rho} [\operatorname{MSE}(\mathbf{x}, \mathbf{z}; \rho)] |_{\rho=0}}{\operatorname{MSE}(\mathbf{x}, \mathbf{z}; 0)} = -\frac{4}{n} \left( 1 + \frac{2\gamma_x^2 \gamma_z^2}{\sigma_x^2 \gamma_z^2 + \sigma_z^2 \gamma_x^2 + \sigma_x^2 \sigma_z^2} \right) + \mathcal{O}_p(n^{-\frac{3}{2}}). \quad [8]$$

83 Here the notation  $\mathcal{O}_p(\cdot)$  is used with the same meaning as in Proposition 2.

84 Proposition 3 presents a simple form of the ratio of the derivative to the MSE itself. The ratio quantifies by what percentage  
 85 the “agreement” penalty decreases the MSE. It can be seen from this representation that this ratio depends on the structure  
 86 of the factor model, and that the agreement penalty is more helpful when the sample size  $n$  is smaller. In the extreme case,  
 87 when we have infinite data, i.e., when  $n = \infty$ , the derivative of the MSE is 0; in this case, we learn all the signals from the data  
 88 even without the agreement penalty.

89 **A. Proof of Proposition 1.** Here we present a lemma that is used in the proof of Proposition 1.

90 **Lemma 1.** Assume that  $c_2 > 0$ . Let  $g(\rho) = (a_1\rho^2 + b_1\rho + c_1)/(a_2\rho^2 + b_2\rho + c_2)^2$ , then

$$91 \quad g'(\rho)|_{\rho=0} = \frac{b_1c_2 - 2b_2c_1}{c_2^3}. \quad [9]$$

92 *Proof.* We compute the derivative of the function  $g$ :

$$93 \quad g'(\rho) = \frac{(2a_1\rho + b_1)(a_2\rho^2 + b_2\rho + c_2) - 2(2a_2\rho + b_2)(a_1\rho^2 + b_1\rho + c_1)}{(a_2\rho^2 + b_2\rho + c_2)^3}. \quad [10]$$

94 Evaluating at  $\rho = 0$ , we get

$$95 \quad g'(\rho)|_{\rho=0} = \frac{b_1c_2 - 2b_2c_1}{c_2^3}. \quad [11]$$

96  $\square$

97 With this lemma, we are ready to prove the proposition. We start with writing down an explicit expression for the estimator  
98  $\hat{\theta}$ . Let

$$99 \quad \tilde{X} = \begin{pmatrix} \mathbf{x} & \mathbf{z} \\ -\sqrt{\rho}\mathbf{x} & \sqrt{\rho}\mathbf{z} \end{pmatrix}, \quad \tilde{y} = \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix}. \quad [12]$$

100 Then Eq. (2) implies that  $\hat{\theta} = (\hat{\theta}_x, \hat{\theta}_z)^\top$  takes the following form:

$$101 \quad \hat{\theta} = (\tilde{X}^\top \tilde{X})^{-1} (\tilde{X}^\top \tilde{y}). \quad [13]$$

102 In particular,

$$103 \quad \tilde{X}^\top \tilde{X} = \begin{pmatrix} (1+\rho)\mathbf{x}^\top \mathbf{x} & (1-\rho)\mathbf{x}^\top \mathbf{z} \\ (1-\rho)\mathbf{x}^\top \mathbf{z} & (1+\rho)\mathbf{z}^\top \mathbf{z} \end{pmatrix}, \quad [14]$$

104 and

$$105 \quad \tilde{X}^\top \tilde{y} = (\mathbf{x}^\top \mathbf{y}, \mathbf{z}^\top \mathbf{y})^\top. \quad [15]$$

106 Therefore,

$$107 \quad \begin{aligned} \hat{\theta} &= \begin{pmatrix} \hat{\theta}_x \\ \hat{\theta}_z \end{pmatrix} = \begin{pmatrix} (1+\rho)\mathbf{x}^\top \mathbf{x} & (1-\rho)\mathbf{x}^\top \mathbf{z} \\ (1-\rho)\mathbf{x}^\top \mathbf{z} & (1+\rho)\mathbf{z}^\top \mathbf{z} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{x}^\top \mathbf{y} \\ \mathbf{z}^\top \mathbf{y} \end{pmatrix} = \frac{1}{\det} \begin{pmatrix} (1+\rho)\mathbf{z}^\top \mathbf{z} & -(1-\rho)\mathbf{x}^\top \mathbf{z} \\ -(1-\rho)\mathbf{x}^\top \mathbf{z} & (1+\rho)\mathbf{x}^\top \mathbf{x} \end{pmatrix} \begin{pmatrix} \mathbf{x}^\top \mathbf{y} \\ \mathbf{z}^\top \mathbf{y} \end{pmatrix} \\ &= \frac{1}{\det} \begin{pmatrix} (1+\rho)(\mathbf{z}^\top \mathbf{z})(\mathbf{x}^\top \mathbf{y}) - (1-\rho)(\mathbf{x}^\top \mathbf{z})(\mathbf{z}^\top \mathbf{y}) \\ (1+\rho)(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{y}) - (1-\rho)(\mathbf{x}^\top \mathbf{z})(\mathbf{x}^\top \mathbf{y}) \end{pmatrix}, \end{aligned} \quad [16]$$

108 where  $\det = (1+\rho)^2(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (1-\rho)^2(\mathbf{x}^\top \mathbf{z})^2$ .

109 We then move on to analyze the conditional distribution of  $\mathbf{u}$  and  $\mathbf{y}$  on  $\mathbf{x}$  and  $\mathbf{z}$ . By Eq. (1), we can write down a joint  
110 distribution of  $(U_i, X_i, Z_i)$ :

$$111 \quad \begin{pmatrix} U_i \\ X_i \\ Z_i \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \gamma_x & \gamma_z \\ \gamma_x & \gamma_x^2 + \sigma_x^2 & \gamma_x \gamma_z \\ \gamma_z & \gamma_x \gamma_z & \gamma_z^2 + \sigma_z^2 \end{pmatrix} \right). \quad [17]$$

112 Using formulas from conditional distribution of multivariate gaussian, we get that

$$113 \quad U_i | X_i, Z_i \sim \mathcal{N}(\text{E}[U_i | X_i, Z_i], \text{Var}[U_i | X_i, Z_i]), \quad [18]$$

114 where

$$115 \quad \begin{aligned} \text{E}[U_i | X_i, Z_i] &= (\gamma_x \quad \gamma_z) \begin{pmatrix} \gamma_x^2 + \sigma_x^2 & \gamma_x \gamma_z \\ \gamma_x \gamma_z & \gamma_z^2 + \sigma_z^2 \end{pmatrix}^{-1} \begin{pmatrix} X_i \\ Z_i \end{pmatrix} \\ &= \frac{1}{\sigma_x^2 \sigma_z^2 + \gamma_x^2 \sigma_z^2 + \gamma_z^2 \sigma_x^2} (\gamma_x \quad \gamma_z) \begin{pmatrix} \gamma_z^2 + \sigma_z^2 & -\gamma_x \gamma_z \\ -\gamma_x \gamma_z & \gamma_x^2 + \sigma_x^2 \end{pmatrix} \begin{pmatrix} X_i \\ Z_i \end{pmatrix} \\ &= \frac{\gamma_x X_i}{\sigma_x^2 + \gamma_x^2 + \gamma_x^2 \sigma_x^2 / \sigma_z^2} + \frac{\gamma_z Z_i}{\sigma_z^2 + \gamma_z^2 + \gamma_z^2 \sigma_z^2 / \sigma_x^2}, \end{aligned} \quad [19]$$

116 and

$$117 \quad \begin{aligned} \text{Var}[U_i | X_i, Z_i] &= 1 - (\gamma_x \quad \gamma_z) \begin{pmatrix} \gamma_x^2 + \sigma_x^2 & \gamma_x \gamma_z \\ \gamma_x \gamma_z & \gamma_z^2 + \sigma_z^2 \end{pmatrix}^{-1} \begin{pmatrix} \gamma_x \\ \gamma_z \end{pmatrix} \\ &= 1 - \frac{1}{\sigma_x^2 \sigma_z^2 + \gamma_x^2 \sigma_z^2 + \gamma_z^2 \sigma_x^2} (\gamma_x \quad \gamma_z) \begin{pmatrix} \gamma_z^2 + \sigma_z^2 & -\gamma_x \gamma_z \\ -\gamma_x \gamma_z & \gamma_x^2 + \sigma_x^2 \end{pmatrix} \begin{pmatrix} \gamma_x \\ \gamma_z \end{pmatrix} \\ &= 1 - \frac{\gamma_x^2 \sigma_z^2 + \gamma_z^2 \sigma_x^2}{\sigma_x^2 \sigma_z^2 + \gamma_x^2 \sigma_z^2 + \gamma_z^2 \sigma_x^2} = \frac{1}{1 + \gamma_x^2 / \sigma_x^2 + \gamma_z^2 / \sigma_z^2}. \end{aligned} \quad [20]$$

118 Since  $Y_i = \gamma_y U_i + \varepsilon_{yi}$ , the above analysis implies that

$$119 \quad Y_i | X_i, Z_i \sim \mathcal{N}(\text{E}[Y_i | X_i, Z_i], \text{Var}[Y_i | X_i, Z_i]), \quad [21]$$

120 where

$$121 \quad \text{E}[Y_i | X_i, Z_i] = \gamma_y \text{E}[U_i | X_i, Z_i] = \frac{\gamma_x \gamma_y X_i}{\sigma_x^2 + \gamma_x^2 + \gamma_z^2 \sigma_x^2 / \sigma_z^2} + \frac{\gamma_z \gamma_y Z_i}{\sigma_z^2 + \gamma_z^2 + \gamma_x^2 \sigma_z^2 / \sigma_x^2}, \quad [22]$$

122 and

$$123 \quad \text{Var}[Y_i | X_i, Z_i] = \gamma_y^2 \text{Var}[U_i | X_i, Z_i] + \sigma_y^2 = \frac{\gamma_y^2}{1 + \gamma_x^2 / \sigma_x^2 + \gamma_z^2 / \sigma_z^2} + \sigma_y^2. \quad [23]$$

124 Let

$$125 \quad \theta_x^* = \frac{\gamma_x \gamma_y}{\sigma_x^2 + \gamma_x^2 + \gamma_z^2 \sigma_x^2 / \sigma_z^2}, \quad \theta_z^* = \frac{\gamma_z \gamma_y}{\sigma_z^2 + \gamma_z^2 + \gamma_x^2 \sigma_z^2 / \sigma_x^2}, \quad \sigma^{*2} = \frac{\gamma_y^2}{1 + \gamma_x^2 / \sigma_x^2 + \gamma_z^2 / \sigma_z^2} + \sigma_y^2, \quad [24]$$

126 then the above shows that we can express  $Y_i$  as

$$127 \quad Y_i = \theta_x^* X_i + \theta_z^* Z_i + \varepsilon_i^*, \quad [25]$$

128 where  $\varepsilon_i^* \perp\!\!\!\perp (X_i, Z_i)$  and  $\varepsilon_i^* \sim \mathcal{N}(0, \sigma^{*2})$ . In words,  $Y_i$  can be decomposed into two independent terms: a linear combination of  $X_i$  and  $Z_i$ , and an error term independent of  $(X_i, Z_i)$ .

130 With the above tools, we are ready to study the MSE. Using Eq. (25), we can write

$$\begin{aligned} \text{MSE}(\mathbf{x}, \mathbf{z}; \rho) &= \text{E} \left[ (Y_{\text{new}} - (X_{\text{new}} \hat{\theta}_x + Z_{\text{new}} \hat{\theta}_z))^2 | \mathbf{x}, \mathbf{z} \right] \\ &= \text{E} \left[ (\theta_x^* X_{\text{new}} + \theta_z^* Z_{\text{new}} + \varepsilon_{\text{new}}^* - (X_{\text{new}} \hat{\theta}_x + Z_{\text{new}} \hat{\theta}_z))^2 | \mathbf{x}, \mathbf{z} \right] \\ &= \text{E} \left[ ((\theta_x^* - \hat{\theta}_x) X_{\text{new}} + (\theta_z^* - \hat{\theta}_z) Z_{\text{new}} + \varepsilon_{\text{new}}^*)^2 | \mathbf{x}, \mathbf{z} \right] \\ &= \text{E} \left[ ((\theta_x^* - \hat{\theta}_x) X_{\text{new}} + (\theta_z^* - \hat{\theta}_z) Z_{\text{new}})^2 | \mathbf{x}, \mathbf{z} \right] + \text{E} [\varepsilon_{\text{new}}^{*2} | \mathbf{x}, \mathbf{z}] \\ &= \text{E} \left[ ((\theta_x^* - \hat{\theta}_x) X_{\text{new}} + (\theta_z^* - \hat{\theta}_z) Z_{\text{new}})^2 | \mathbf{x}, \mathbf{z} \right] + \sigma^{*2}. \end{aligned} \quad [26]$$

132 Here the cross terms vanish because  $\varepsilon_{\text{new}}^* \perp\!\!\!\perp (X_{\text{new}}, Z_{\text{new}})$ . Since the new dataset is independent of the training dataset, we can further simply the above:

$$\begin{aligned} \text{MSE}(\mathbf{x}, \mathbf{z}; \rho) &= \text{E} \left[ ((\theta_x^* - \hat{\theta}_x) X_{\text{new}})^2 | \mathbf{x}, \mathbf{z} \right] + \text{E} \left[ ((\theta_z^* - \hat{\theta}_z) Z_{\text{new}})^2 | \mathbf{x}, \mathbf{z} \right] \\ &\quad + 2\text{E} \left[ (\theta_z^* - \hat{\theta}_z) (\hat{\theta}_x - \theta_x^*) Z_{\text{new}} X_{\text{new}} | \mathbf{x}, \mathbf{z} \right] + \sigma^{*2} \\ &= \text{E} \left[ (\hat{\theta}_x - \theta_x^*)^2 | \mathbf{x}, \mathbf{z} \right] \text{E} [X_{\text{new}}^2] + \text{E} \left[ (\hat{\theta}_z - \theta_z^*)^2 | \mathbf{x}, \mathbf{z} \right] \text{E} [Z_{\text{new}}^2] \\ &\quad + 2\text{E} \left[ (\hat{\theta}_z - \theta_z^*) (\hat{\theta}_x - \theta_x^*) | \mathbf{x}, \mathbf{z} \right] \text{E} [Z_{\text{new}} X_{\text{new}}] + \sigma^{*2} \\ &= \text{E} \left[ (\hat{\theta}_x - \theta_x^*)^2 | \mathbf{x}, \mathbf{z} \right] (\gamma_x^2 + \sigma_x^2) + \text{E} \left[ (\hat{\theta}_z - \theta_z^*)^2 | \mathbf{x}, \mathbf{z} \right] (\gamma_z^2 + \sigma_z^2) \\ &\quad + 2\text{E} \left[ (\hat{\theta}_z - \theta_z^*) (\hat{\theta}_x - \theta_x^*) | \mathbf{x}, \mathbf{z} \right] (\gamma_x \gamma_z) + \sigma^{*2}. \end{aligned} \quad [27]$$

135 We can then further decompose the terms into squared bias plus variance.

$$\begin{aligned} \text{MSE}(\mathbf{x}, \mathbf{z}; \rho) &= \text{E} [\hat{\theta}_x - \theta_x^* | \mathbf{x}, \mathbf{z}]^2 (\gamma_x^2 + \sigma_x^2) + \text{Var} [\hat{\theta}_x | \mathbf{x}, \mathbf{z}] (\gamma_x^2 + \sigma_x^2) \\ &\quad + \text{E} [\hat{\theta}_z - \theta_z^* | \mathbf{x}, \mathbf{z}]^2 (\gamma_z^2 + \sigma_z^2) + \text{Var} [\hat{\theta}_z | \mathbf{x}, \mathbf{z}] (\gamma_z^2 + \sigma_z^2) \\ &\quad + 2\text{E} [\hat{\theta}_z - \theta_z^* | \mathbf{x}, \mathbf{z}] \text{E} [\hat{\theta}_x - \theta_x^* | \mathbf{x}, \mathbf{z}] (\gamma_x \gamma_z) + 2 \text{Cov} [\hat{\theta}_z, \hat{\theta}_x | \mathbf{x}, \mathbf{z}] (\gamma_x \gamma_z) + \sigma^{*2} \\ &= B^2(\mathbf{x}, \mathbf{z}; \rho) + V(\mathbf{x}, \mathbf{z}; \rho) + \sigma^{*2}, \end{aligned} \quad [28]$$

137 where  $B^2(\mathbf{x}, \mathbf{z}; \rho) = \text{E} [\hat{\theta}_x - \theta_x^* | \mathbf{x}, \mathbf{z}]^2 (\gamma_x^2 + \sigma_x^2) + \text{E} [\hat{\theta}_z - \theta_z^* | \mathbf{x}, \mathbf{z}]^2 (\gamma_z^2 + \sigma_z^2) + 2\text{E} [\hat{\theta}_z - \theta_z^* | \mathbf{x}, \mathbf{z}] \text{E} [\hat{\theta}_x - \theta_x^* | \mathbf{x}, \mathbf{z}] (\gamma_x \gamma_z)$  is  
138 the sum of bias related terms, and  $V(\mathbf{x}, \mathbf{z}; \rho) = \text{Var} [\hat{\theta}_x | \mathbf{x}, \mathbf{z}] (\gamma_x^2 + \sigma_x^2) + \text{Var} [\hat{\theta}_z | \mathbf{x}, \mathbf{z}] (\gamma_z^2 + \sigma_z^2) + 2 \text{Cov} [\hat{\theta}_z, \hat{\theta}_x | \mathbf{x}, \mathbf{z}] (\gamma_x \gamma_z)$  is  
139 the sum of variance related terms.

140 We can then use Eq. (21) - Eq. (23) to study the bias and variance of the estimators  $\hat{\theta}_x$  and  $\hat{\theta}_z$ . We start with the bias. By  
141 Eq. (25), we have  $E[\mathbf{y} | \mathbf{x}, \mathbf{z}] = \theta_x^* \mathbf{x} + \theta_z^* \mathbf{z}$ . Therefore,

$$\begin{aligned} E[\hat{\theta}_x | \mathbf{x}, \mathbf{z}] &= E\left[\frac{1}{\det} ((1+\rho)(\mathbf{z}^\top \mathbf{z})(\mathbf{x}^\top \mathbf{y}) - (1-\rho)(\mathbf{x}^\top \mathbf{z})(\mathbf{z}^\top \mathbf{y})) | \mathbf{x}, \mathbf{z}\right] \\ &= \frac{1}{\det} ((1+\rho)(\mathbf{z}^\top \mathbf{z})(\mathbf{x}^\top E[\mathbf{y} | \mathbf{x}, \mathbf{z}]) - (1-\rho)(\mathbf{x}^\top \mathbf{z})(\mathbf{z}^\top E[\mathbf{y} | \mathbf{x}, \mathbf{z}])) \\ &= \frac{1}{\det} ((1+\rho)(\mathbf{z}^\top \mathbf{z})(\mathbf{x}^\top (\theta_x^* \mathbf{x} + \theta_z^* \mathbf{z})) - (1-\rho)(\mathbf{x}^\top \mathbf{z})(\mathbf{z}^\top (\theta_x^* \mathbf{x} + \theta_z^* \mathbf{z}))) \\ &= \frac{1}{\det} ((1+\rho)[\theta_x^*(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) + \theta_z^*(\mathbf{z}^\top \mathbf{z})(\mathbf{x}^\top \mathbf{z})] - (1-\rho)[\theta_x^*(\mathbf{x}^\top \mathbf{z})^2 + \theta_z^*(\mathbf{x}^\top \mathbf{z})(\mathbf{z}^\top \mathbf{z})]) \\ &= \frac{1}{\det} (\theta_x^* [(x^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2] + \rho(\theta_x^* [(x^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) + (\mathbf{x}^\top \mathbf{z})^2] + 2\theta_z^*(\mathbf{x}^\top \mathbf{z})(\mathbf{z}^\top \mathbf{z}))). \end{aligned} \quad [29]$$

143 Note that

$$\begin{aligned} \det &= (1+\rho)^2(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (1-\rho)^2(\mathbf{x}^\top \mathbf{z})^2 \\ &= (\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2 + 2\rho[(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) + (\mathbf{x}^\top \mathbf{z})^2] + \rho^2[(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2] \end{aligned} \quad [30]$$

144 Therefore,

$$\begin{aligned} E[\hat{\theta}_x - \theta_x^* | \mathbf{x}, \mathbf{z}] &= \frac{\theta_x^*[(x^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2] + \rho(\theta_x^*[(x^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) + (\mathbf{x}^\top \mathbf{z})^2] + 2\theta_z^*(\mathbf{x}^\top \mathbf{z})(\mathbf{z}^\top \mathbf{z}))}{(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2 + 2\rho[(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) + (\mathbf{x}^\top \mathbf{z})^2] + \rho^2[(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2]} - \theta_x^* \\ &= \frac{\rho(\theta_x^*[-(\mathbf{z}^\top \mathbf{z})(\mathbf{x}^\top \mathbf{x}) - (\mathbf{x}^\top \mathbf{z})^2] + 2\theta_z^*(\mathbf{x}^\top \mathbf{z})(\mathbf{z}^\top \mathbf{z})) - \rho^2\theta_x^*[(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2]}{(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2 + 2\rho[(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) + (\mathbf{x}^\top \mathbf{z})^2] + \rho^2[(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2]} \\ &= \frac{b_1\rho + a_1\rho^2}{c_2 + b_2\rho + a_2\rho^2}, \end{aligned} \quad [31]$$

145 where  $b_1, a_1, c_2, b_2, a_2$  are expressions depending on  $\mathbf{x}$  and  $\mathbf{z}$  but not on  $\rho$ . We can then clearly see that when  $\rho = 0$ ,  
146  $E[\hat{\theta}_x - \theta_x^* | \mathbf{x}, \mathbf{z}] = 0$ . By symmetry, we have the same property for  $E[\hat{\theta}_z - \theta_z^* | \mathbf{x}, \mathbf{z}]$ . Therefore,

$$\begin{aligned} \frac{d}{d\rho} B^2(\mathbf{x}, \mathbf{z}; \rho) &= (\gamma_x^2 + \sigma_x^2) \frac{d}{d\rho} E[\hat{\theta}_x - \theta_x^* | \mathbf{x}, \mathbf{z}]^2 |_{\rho=0} + (\gamma_z^2 + \sigma_z^2) \frac{d}{d\rho} E[\hat{\theta}_z - \theta_z^* | \mathbf{x}, \mathbf{z}]^2 |_{\rho=0} \\ &\quad + 2(\gamma_x \gamma_z) \frac{d}{d\rho} (E[\hat{\theta}_z - \theta_z^* | \mathbf{x}, \mathbf{z}] E[\hat{\theta}_x - \theta_x^* | \mathbf{x}, \mathbf{z}]) |_{\rho=0} \\ &= 2(\gamma_x^2 + \sigma_x^2) (E[\hat{\theta}_x - \theta_x^* | \mathbf{x}, \mathbf{z}] |_{\rho=0}) \frac{d}{d\rho} E[\hat{\theta}_x - \theta_x^* | \mathbf{x}, \mathbf{z}] |_{\rho=0} \\ &\quad + 2(\gamma_z^2 + \sigma_z^2) (E[\hat{\theta}_z - \theta_z^* | \mathbf{x}, \mathbf{z}] |_{\rho=0}) \frac{d}{d\rho} E[\hat{\theta}_z - \theta_z^* | \mathbf{x}, \mathbf{z}] |_{\rho=0} \\ &\quad + 2(\gamma_x \gamma_z) (E[\hat{\theta}_z - \theta_z^* | \mathbf{x}, \mathbf{z}] |_{\rho=0}) \frac{d}{d\rho} E[\hat{\theta}_x - \theta_x^* | \mathbf{x}, \mathbf{z}] |_{\rho=0} \\ &\quad + 2(\gamma_x \gamma_z) (E[\hat{\theta}_x - \theta_x^* | \mathbf{x}, \mathbf{z}] |_{\rho=0}) \frac{d}{d\rho} E[\hat{\theta}_z - \theta_z^* | \mathbf{x}, \mathbf{z}] |_{\rho=0}. \end{aligned} \quad [32]$$

150 Since  $E[\hat{\theta}_x - \theta_x^* | \mathbf{x}, \mathbf{z}] |_{\rho=0} = E[\hat{\theta}_z - \theta_z^* | \mathbf{x}, \mathbf{z}] |_{\rho=0} = 0$ , we have

$$\frac{d}{d\rho} B^2(\mathbf{x}, \mathbf{z}; \rho) = 0. \quad [33]$$

152 It remains to study  $\frac{d}{d\rho} V(\mathbf{x}, \mathbf{z}; \rho)$ . From the form of  $\hat{\theta}_x$ , we can get that

$$\begin{aligned} \text{Var}[\hat{\theta}_x | \mathbf{x}, \mathbf{z}] &= \frac{\sigma^{*2}}{\det^2} \|((1+\rho)(\mathbf{z}^\top \mathbf{z})\mathbf{x} - (1-\rho)(\mathbf{x}^\top \mathbf{z})\mathbf{z})\|_2^2 \\ &= \frac{\sigma^{*2}}{\det^2} [(1+\rho)^2(\mathbf{z}^\top \mathbf{z})^2(\mathbf{x}^\top \mathbf{x}) + (1-\rho)^2(\mathbf{x}^\top \mathbf{z})^2(\mathbf{z}^\top \mathbf{z}) - 2(1+\rho)(1-\rho)(\mathbf{x}^\top \mathbf{z})^2(\mathbf{z}^\top \mathbf{z})] \\ &= \frac{\sigma^{*2}(\mathbf{z}^\top \mathbf{z})}{\det^2} ([(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2] + 2[(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2]\rho + a_{v1}\rho^2), \end{aligned} \quad [34]$$

154 for some  $a_{v1}$  depending on  $\mathbf{x}$  and  $\mathbf{z}$  but not on  $\rho$ . Similarly, we get that

$$\text{Var}[\hat{\theta}_z | \mathbf{x}, \mathbf{z}] = \frac{\sigma^{*2}(\mathbf{x}^\top \mathbf{x})}{\det^2} ([(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2] + 2[(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2]\rho + a_{v2}\rho^2), \quad [35]$$

156 for some  $a_{v2}$  depending on  $\mathbf{x}$  and  $\mathbf{z}$  but not on  $\rho$ . For the covariance term,

$$\begin{aligned} \text{Cov} [\hat{\theta}_x, \hat{\theta}_z | \mathbf{x}, \mathbf{z}] &= \frac{\sigma^{*2}}{\det^2} [(1+\rho)(\mathbf{z}^\top \mathbf{z})\mathbf{x} - (1-\rho)(\mathbf{x}^\top \mathbf{z})\mathbf{z}]^\top [(1+\rho)(\mathbf{x}^\top \mathbf{x})\mathbf{z} - (1-\rho)(\mathbf{z}^\top \mathbf{x})\mathbf{x}] \\ &= \frac{\sigma^{*2}}{\det^2} ((-1+3\rho)(1+\rho)(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z})(\mathbf{x}^\top \mathbf{z}) + (1-\rho)^2(\mathbf{x}^\top \mathbf{z})^3) \\ &= \frac{\sigma^{*2}(\mathbf{x}^\top \mathbf{z})}{\det^2} (-[(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2] + 2[(\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2]\rho + a_{v3}\rho^2), \end{aligned} \quad [36]$$

158 for some  $a_{v3}$  depending on  $\mathbf{x}$  and  $\mathbf{z}$  but not on  $\rho$ . Combining the three terms, we get

$$\begin{aligned} V(\mathbf{x}, \mathbf{z}; \rho) &= \text{Var} [\hat{\theta}_x | \mathbf{x}, \mathbf{z}] (\gamma_x^2 + \sigma_x^2) + \text{Var} [\hat{\theta}_z | \mathbf{x}, \mathbf{z}] (\gamma_z^2 + \sigma_z^2) + \text{Cov} [\hat{\theta}_z, \hat{\theta}_x | \mathbf{x}, \mathbf{z}] (\gamma_x \gamma_z) \\ &= \sigma^{*2} \frac{C_1 + B_1 \rho + A_1 \rho^2}{\det^2} = \sigma^{*2} \frac{C_1 + B_1 \rho + A_1 \rho^2}{(C_2 + B_2 \rho + A_2 \rho^2)^2}, \end{aligned} \quad [37]$$

160 where

$$\begin{aligned} C_1 &= [(\gamma_x^2 + \sigma_x^2)(\mathbf{z}^\top \mathbf{z}) + (\gamma_z^2 + \sigma_z^2)(\mathbf{x}^\top \mathbf{x}) - 2\gamma_x \gamma_z (\mathbf{x}^\top \mathbf{z})] ((\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2), \\ B_1 &= 2[(\gamma_x^2 + \sigma_x^2)(\mathbf{z}^\top \mathbf{z}) + (\gamma_z^2 + \sigma_z^2)(\mathbf{x}^\top \mathbf{x}) + 2\gamma_x \gamma_z (\mathbf{x}^\top \mathbf{z})] ((\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2), \\ C_2 &= (\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) - (\mathbf{x}^\top \mathbf{z})^2, \\ B_2 &= 2((\mathbf{x}^\top \mathbf{x})(\mathbf{z}^\top \mathbf{z}) + (\mathbf{x}^\top \mathbf{z})^2). \end{aligned} \quad [38]$$

162 By Lemma 1,  $\frac{d}{d\rho} V(\mathbf{x}, \mathbf{z}; \rho)|_{\rho=0} = \sigma^{*2}(C_2 B_1 - 2C_1 B_2)/C_2^3$ .

163 Finally by Eq. (28)

$$\begin{aligned} \frac{d}{d\rho} \text{MSE}(\mathbf{x}, \mathbf{z}; \rho)|_{\rho=0} &= \frac{d}{d\rho} B^2(\mathbf{x}, \mathbf{z}; \rho)|_{\rho=0} + \frac{d}{d\rho} V(\mathbf{x}, \mathbf{z}; \rho)|_{\rho=0} \\ &= \frac{d}{d\rho} V(\mathbf{x}, \mathbf{z}; \rho)|_{\rho=0} = \sigma^{*2}(C_2 B_1 - 2C_1 B_2)/C_2^3. \end{aligned} \quad [39]$$

165 **B. Proof of Proposition 2.** By the central limit theorem, we have that

$$\mathbf{x}^\top \mathbf{x} = n(\gamma_x^2 + \sigma_x^2) + \mathcal{O}_p(\sqrt{n}), \quad \mathbf{z}^\top \mathbf{z} = n(\gamma_z^2 + \sigma_z^2) + \mathcal{O}_p(\sqrt{n}), \quad \mathbf{x}^\top \mathbf{z} = n(\gamma_x \gamma_z) + \mathcal{O}_p(\sqrt{n}). \quad [40]$$

166 Plugging into Eq. (6) gives

$$\begin{aligned} C_1 &= 2n^3 [(\gamma_x^2 + \sigma_x^2)(\gamma_z^2 + \sigma_z^2) - \gamma_x^2 \gamma_z^2]^2 + \mathcal{O}_p(n^{5/2}), \\ B_1 &= 4n^3 [(\gamma_x^2 + \sigma_x^2)(\gamma_z^2 + \sigma_z^2) + \gamma_x^2 \gamma_z^2] [(\gamma_x^2 + \sigma_x^2)(\gamma_z^2 + \sigma_z^2) - \gamma_x^2 \gamma_z^2] + \mathcal{O}_p(n^{5/2}), \\ C_2 &= n^2 [(\gamma_x^2 + \sigma_x^2)(\gamma_z^2 + \sigma_z^2) - \gamma_x^2 \gamma_z^2] + \mathcal{O}_p(n^{3/2}), \\ B_2 &= 2n^2 [(\gamma_x^2 + \sigma_x^2)(\gamma_z^2 + \sigma_z^2) + \gamma_x^2 \gamma_z^2] + \mathcal{O}_p(n^{3/2}). \end{aligned} \quad [41]$$

169 Thus we have that

$$\begin{aligned} C_1 B_2 &= 4n^5 [(\gamma_x^2 + \sigma_x^2)(\gamma_z^2 + \sigma_z^2) + \gamma_x^2 \gamma_z^2] [(\gamma_x^2 + \sigma_x^2)(\gamma_z^2 + \sigma_z^2) - \gamma_x^2 \gamma_z^2]^2 + \mathcal{O}_p(n^{9/2}), \\ C_2 B_1 &= 4n^5 [(\gamma_x^2 + \sigma_x^2)(\gamma_z^2 + \sigma_z^2) + \gamma_x^2 \gamma_z^2] [(\gamma_x^2 + \sigma_x^2)(\gamma_z^2 + \sigma_z^2) - \gamma_x^2 \gamma_z^2]^2 + \mathcal{O}_p(n^{9/2}). \end{aligned} \quad [42]$$

171 Therefore,

$$C_2 B_1 - 2C_1 B_2 = -4n^5 [(\gamma_x^2 + \sigma_x^2)(\gamma_z^2 + \sigma_z^2) + \gamma_x^2 \gamma_z^2] [(\gamma_x^2 + \sigma_x^2)(\gamma_z^2 + \sigma_z^2) - \gamma_x^2 \gamma_z^2]^2 + \mathcal{O}_p(n^{9/2}). \quad [43]$$

173 Now we also know that

$$C_2^3 = n^6 [(\gamma_x^2 + \sigma_x^2)(\gamma_z^2 + \sigma_z^2) - \gamma_x^2 \gamma_z^2]^3 + \mathcal{O}_p(n^{11/2}). \quad [44]$$

174 Hence

$$\begin{aligned} \frac{d}{d\rho} [\text{MSE}(\mathbf{x}, \mathbf{z}; \rho)]|_{\rho=0} &= \sigma^{*2}(C_2 B_1 - 2C_1 B_2)/C_2^3 \\ &= \frac{-4n^5 [(\gamma_x^2 + \sigma_x^2)(\gamma_z^2 + \sigma_z^2) + \gamma_x^2 \gamma_z^2] [(\gamma_x^2 + \sigma_x^2)(\gamma_z^2 + \sigma_z^2) - \gamma_x^2 \gamma_z^2]^2}{n^6 [(\gamma_x^2 + \sigma_x^2)(\gamma_z^2 + \sigma_z^2) - \gamma_x^2 \gamma_z^2]^3} \sigma^{*2} + \mathcal{O}_p(n^{-3/2}) \\ &= -\frac{4}{n} \frac{(\gamma_x^2 + \sigma_x^2)(\gamma_z^2 + \sigma_z^2) + \gamma_x^2 \gamma_z^2}{(\gamma_x^2 + \sigma_x^2)(\gamma_z^2 + \sigma_z^2) - \gamma_x^2 \gamma_z^2} \sigma^{*2} + \mathcal{O}_p(n^{-3/2}) \\ &= -\frac{4}{n} \left(1 + \frac{2\gamma_x^2 \gamma_z^2}{\sigma_x^2 \gamma_z^2 + \sigma_z^2 \gamma_x^2 + \sigma_x^2 \sigma_z^2}\right) \sigma^{*2} + \mathcal{O}_p(n^{-3/2}) \\ &= -\frac{4}{n} \left(1 + \frac{2\gamma_x^2 \gamma_z^2}{\sigma_x^2 \gamma_z^2 + \sigma_z^2 \gamma_x^2 + \sigma_x^2 \sigma_z^2}\right) \left(\sigma_y^2 + \frac{\gamma_y^2 \sigma_x^2 \sigma_z^2}{\sigma_x^2 \gamma_z^2 + \sigma_z^2 \gamma_x^2 + \sigma_x^2 \sigma_z^2}\right) + \mathcal{O}_p(n^{-3/2}). \end{aligned} \quad [45]$$

177 **C. Proof of Proposition 3.** For  $\text{MSE}(\mathbf{x}, \mathbf{z}; 0)$ , we have that by Eq. (28),

$$178 \quad \text{MSE}(\mathbf{x}, \mathbf{z}; 0) = B^2(\mathbf{x}, \mathbf{z}; 0) + V(\mathbf{x}, \mathbf{z}; 0) + \sigma^{*2} = V(\mathbf{x}, \mathbf{z}; 0) + \sigma^{*2}. \quad [46]$$

179 Here we make use of the fact that when  $\rho = 0$ ,  $E[\hat{\theta}_x - \theta_x^* | \mathbf{x}, \mathbf{z}] = E[\hat{\theta}_z - \theta_z^* | \mathbf{x}, \mathbf{z}] = 0$  and that  $B^2(\mathbf{x}, \mathbf{z}; 0) = 0$ . For  
180  $V(\mathbf{x}, \mathbf{z}; 0)$ , we have that by Eq. (37) and Eq. (40),

$$181 \quad V(\mathbf{x}, \mathbf{z}; 0) = \sigma^{*2} \frac{C_1}{C_2^2} = \frac{4\sigma^{*2}}{n} + \mathcal{O}_p(n^{-3/2}) = \mathcal{O}_p\left(\frac{1}{n}\right). \quad [47]$$

182 Therefore,

$$183 \quad \text{MSE}(\mathbf{x}, \mathbf{z}; 0) = V(\mathbf{x}, \mathbf{z}; 0) + \sigma^{*2} = \sigma^{*2} + \mathcal{O}_p\left(\frac{1}{n}\right). \quad [48]$$

184 Thus, together with the result in Proposition 2, we have

$$185 \quad \frac{\frac{d}{d\rho} [\text{MSE}(\mathbf{x}, \mathbf{z}; \rho)]|_{\rho=0}}{\text{MSE}(\mathbf{x}, \mathbf{z}; 0)} = -\frac{4}{n} \left(1 + \frac{2\gamma_x^2\gamma_z^2}{\sigma_x^2\gamma_z^2 + \sigma_z^2\gamma_x^2 + \sigma_x^2\sigma_z^2}\right) + \mathcal{O}_p\left(n^{-\frac{3}{2}}\right). \quad [49]$$

## 186 5. Distribution of predicted versus true time to delivery for the labor onset prediction example

187 We show in Figure S9 the distribution of predicted and true time to delivery for each patient, which gives a better sense of the  
188 quality of the predictions for the regression task. The left plot shows the distribution of time to delivery for all patients at their  
189 first time points in the longitudinal study; the right plot shows the predicted versus true time to delivery for the training and  
190 test samples. This is based on one random split of the training and test sets of 40 and 13 patients, respectively.

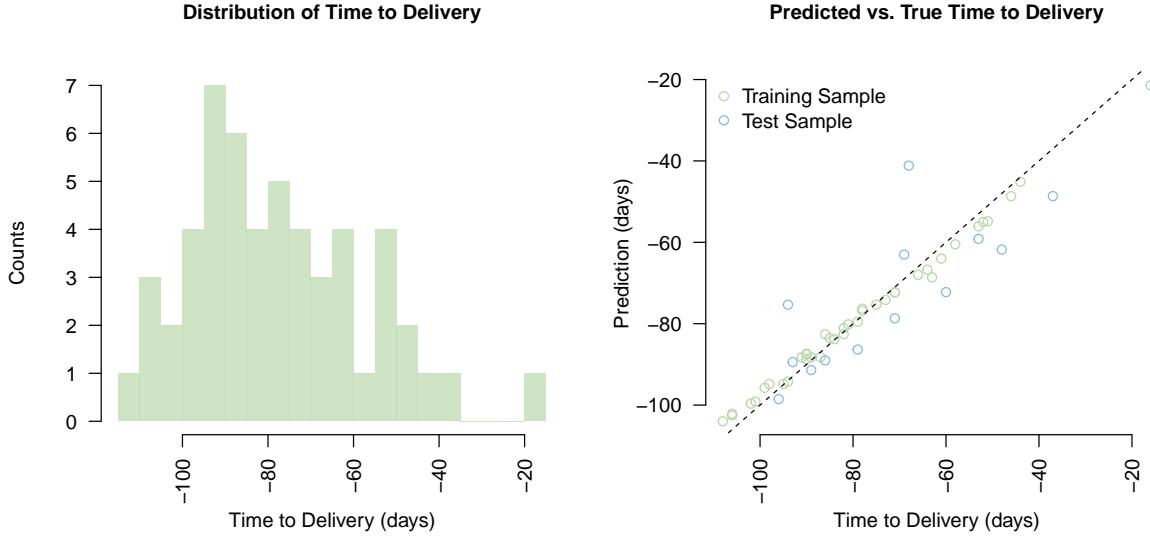


Fig. S9. Distribution of time to delivery and predicted versus true time to delivery for training and test samples. The predictions were derived from cooperative learning.

## 191 6. Procedure for generating the imaging and “omics” data

192 Here we outline the detailed procedure for data generation in the simulation study with imaging and “omics” data in Algorithm  
193 S3. The “omics” data ( $X$ ), imaging data ( $Z$ ), and the response  $y$  are generated such that there are correlations between  $X$ ,  $Z$ ,  
194 and  $y$ .

---

**Algorithm S3** *Simulation procedure for generating the imaging and “omics” data.*

---

**Input:** Parameters  $n, p_x, p_u, s_u, t, \sigma, \beta_u, I_{\max}, \text{ndim}$ , threshold.

**Output:**  $X \in \mathcal{R}^{n \times p_x}$  (omics),  $Z \in \mathcal{R}^{n \times \text{ndim} \times \text{ndim} \times 1}$  (images assuming one color channel),  $\mathbf{y} \in \mathcal{R}^n$ .

1.  $x_j \in \mathcal{R}^n$  distributed i.i.d.  $\text{MVN}(0, I_n)$  for  $j = 1, 2, \dots, p_x$
  2. For  $i = 1, 2, \dots, p_u$  ( $p_u < p_x$ , where  $p_u$  corresponds to the number of factors):
    - (a)  $u_i \in \mathcal{R}^n$  distributed i.i.d.  $\text{MVN}(0, s_u^2 I_n)$
    - (b)  $x_i = x_i + t * u_i$
  3.  $U = [u_1, u_2, \dots, u_{p_u}]$ ,  $X = [x_1, x_2, \dots, x_{p_x}]$
  4.  $\mathbf{y}_u = U\beta_u + \epsilon$  where  $\epsilon \in \mathcal{R}^n$  distributed i.i.d.  $\text{MVN}(0, \sigma^2 I_n)$
  5. For  $i = 1, 2, \dots, n$ :
    - (a)  $P_i = \frac{1}{1+\exp(\mathbf{y}_{u_i})}$ ,  $\mathbf{y}_i \sim \text{Bernoulli}(P_i)$
    - (b) Generate a 2D pixel matrix of image  $Z_i \in \mathcal{R}^{\text{dim} \times \text{dim} \times 1}$
    - (c) Generate a polygon  $\text{PG}_i$  inside  $Z_i$ , defined by 4 vertices  $[v_1, v_2, v_3, v_4]$  on the axes, i.e.  $v_1 = [0, a_1], v_2 = [0, a_2], v_3 = [a_3, 0], v_4 = [a_4, 0]$ , where  $a_1 \sim \text{Unif}(\frac{\text{ndim}}{2}, \text{ndim}), a_2 \sim \text{Unif}(-\text{ndim}, -\frac{\text{ndim}}{2}), a_3 \sim \text{Unif}(\frac{\text{ndim}}{2}, \text{ndim}), a_4 \sim \text{Unif}(-\text{ndim}, -\frac{\text{ndim}}{2})$
    - (d) Randomly sample points from  $Z_i$ : if the point  $[x', y']$  falls inside the polygon  $\text{PG}_i$ , i.e.  $[x', y'] \in \text{PG}_i$ , then  $Z_i[x', y'] \sim \text{Unif}(0, 1)$
    - (e) If  $\mathbf{y}_i = 1$ ,  $I_{\text{disease}} = I_{\max} \times \mathbf{y}_{u_i}$ , where  $I_{\max}$  is the maximum intensity of pixel values for images,
      - For  $x' = 1, 2, \dots, \text{ndim}$ :
        - For  $y' = 1, 2, \dots, \text{ndim}$ :
          - \*  $P(x', y') \sim \text{Unif}(0, 1)$
          - \* If  $[x', y'] \in \text{PG}_i$  and  $P(x', y') < \text{threshold}$ ,  $Z_i[x', y'] = I_{\text{disease}}$
-