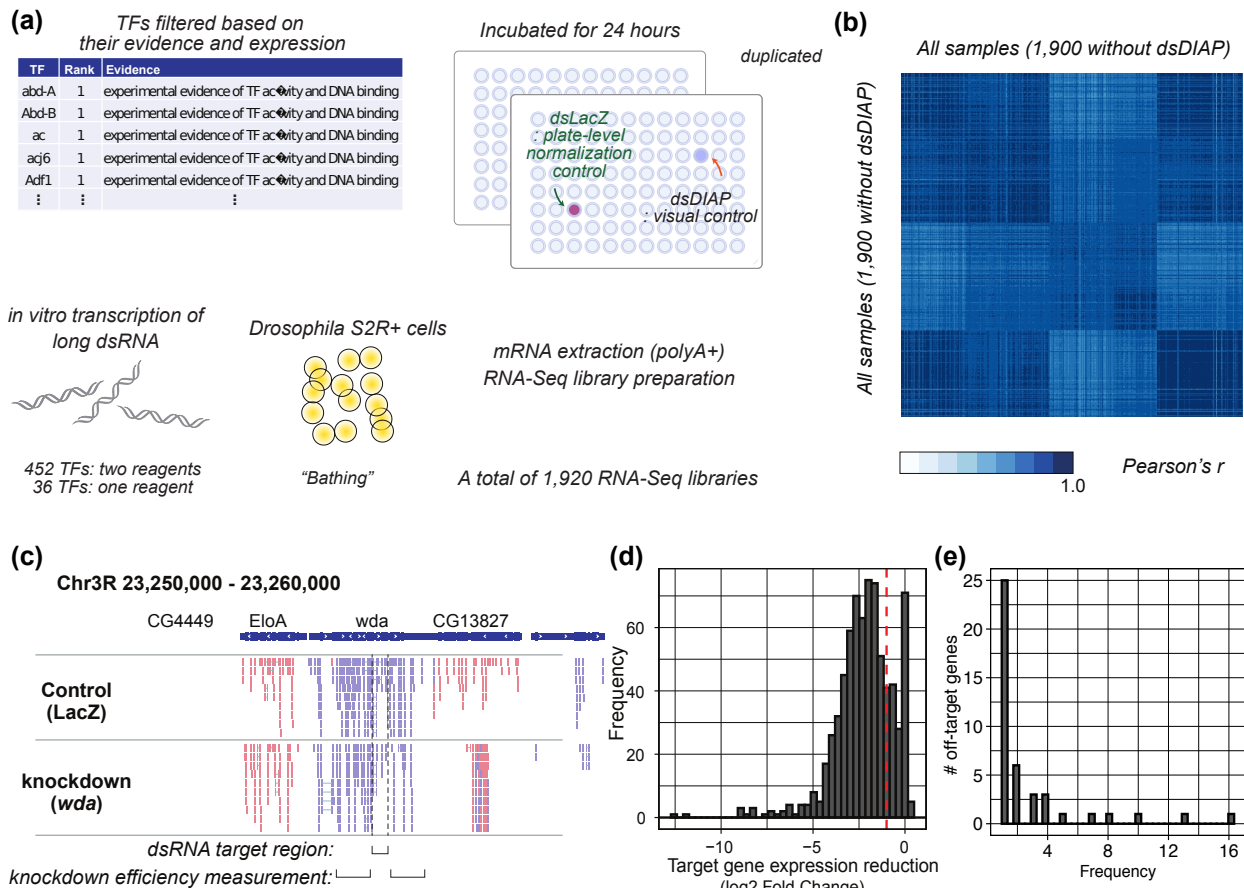
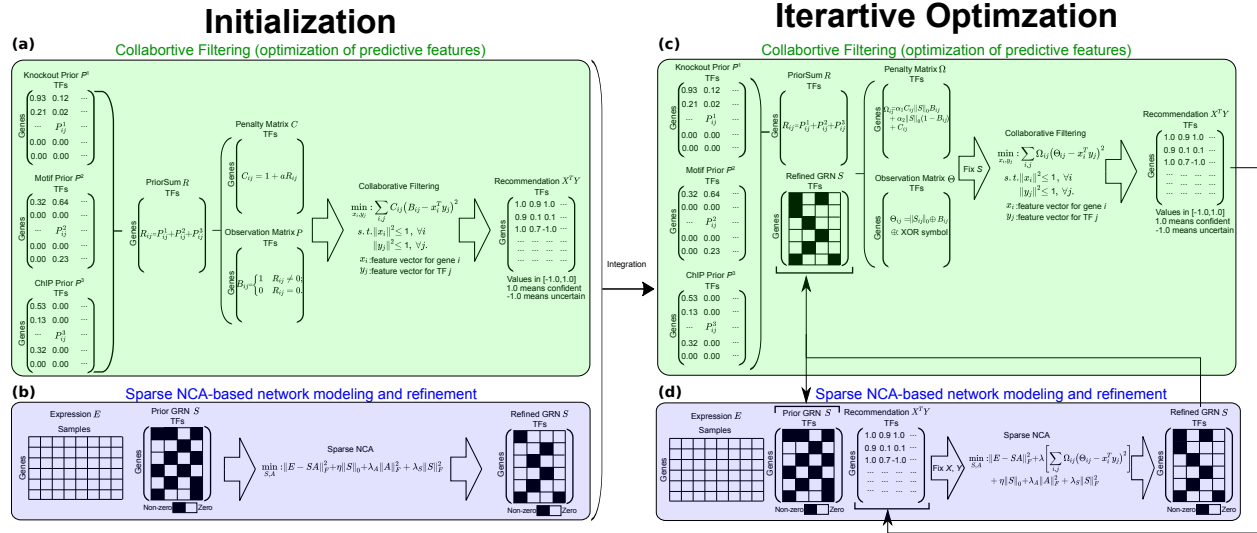


Supplementary Materials



Supplementary Figure 1: A large-scale RNA-Seq analysis of *Drosophila* S2 cells with TF RNAi (a) A schematic illustration of the experimental design. dsLacZ. Cells treated with dsRNA against *E.coli* LacZ gene. dsDIAP. Cells treated with dsRNA against *Drosophila* Diap-1 gene (Death-associated inhibitor of apoptosis 1) (b) All sample-to-all sample pairwise correlation of gene expression is demonstrated as a heatmap. Color intensity corresponds to Pearson's correlation coefficient r . (c) An example of *wda* gene illustrates how short reads are mapped to the reference genome. Blue tiles. Reads that are aligned to 5' to 3' direction in the strand-specific RNA-Seq outcome. Red tiles. Reads for 3' to 5' direction. Top. The result from the control cells (*dsLacZ*). Bottom. Results from the cells treated with dsRNA for *wda*. Dot lines indicate the region for dsRNA targeting. Knockdown efficiencies were measured without reads mapped onto this region from both control and RNAi samples. (d) A histogram that shows the distribution of knockdown efficiency values, which are represented as reduction of target gene RNA (log2 Fold Changes between RNAi vs. Control). Dot red line. 50% reduction of mRNA (log2 Fold Change = - 1). (e) A histogram that shows an incident of significant off-target effects (adjusted p value ≤ 0.1) from all of the reagents ($n = 941$). We used BLAT [1] for the identification of any sequences that have more than 7-bp matching any of the reagents.



Supplementary Figure 2: The overview of the information flow in the NetREX-CF optimization.

Supplementary Note 1: GPALM Problem Introduction

We extend the original PALM algorithm [2] and propose the GPALM algorithm that can solve more general problems. The format of the problem that GPALM can solve is explained in this section. The actual algorithm and its convergence proof are provided in Supplementary Note 2.

The problem and basic assumptions

We are interested in solving the non-convex and non-smooth minimization problem with the following structure

$$(M) \quad \min : \Psi (X, Y, Z) := H (X, Y) + F (Y, Z), \quad (1)$$

where we have the following assumption:

Assumption 1. *The assumptions for problem (M) is as follow:*

1. $H : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ is a C^1 function.
2. $F : \mathbb{R}^m \times \mathbb{R}^l \rightarrow (-\infty + \infty]$ is a proper and lower semicontinuous (PLS) function. And $F(Y, Z)$ has the following structure $F(Y, Z) := \sum_{i=1}^m p_i(Z)g_i(Y_i) + Q(Z)$, where $p_i : \mathbb{R}^l \rightarrow \mathbb{R}$ is Lipschitz continuous with moduli $L_i(Z)$ and $p_i(Z) > 0, \forall i$ and $g_i : \mathbb{R} \rightarrow \mathbb{R}$ is lower semicontinuous and $\sup g_i(Y_i) < \lambda_i, \forall i$ and $Q : \mathbb{R}^l \rightarrow \mathbb{R}$ is Lipschitz continuous with moduli $L_Q(Z)$. $Y = [Y_1, \dots, Y_i, \dots, Y_m]$.
3. $\inf_{\mathbb{R}^n \times \mathbb{R}^m} H > -\infty$ and $\inf_{\mathbb{R}^m \times \mathbb{R}^l} F > -\infty$.
4. For any Y the function $X \rightarrow H(X, Y)$ is $C_{L_X(Y)}^{1,1}$, namely the partial gradient $\nabla_X H(X, Y)$ is globally Lipschitz with moduli $L_1(Y)$, that is

$$\|\nabla_X H(X_1, Y) - \nabla_X H(X_2, Y)\| \leq L_1(Y) \|X_1 - X_2\|. \quad (2)$$

Likewise, for any fixed X the function $Y_i \rightarrow H(X, Y_i)$ is assumed to be $C_{L_{Y_i}(X)}^{1,1}$.

5. For any fixed Y the function $Z \rightarrow F(Y, Z)$ is assumed to be $C_{L_Z(Y)}^{1,1}$.
6. ∇H is Lipschitz continuous on bounded subsets of $\mathbb{R}^n \times \mathbb{R}^m$. In other words, for each bounded subsets $T_1 \times T_2$ of $\mathbb{R}^n \times \mathbb{R}^m$ there exist $M > 0$ such that any (X_1, Y_1) and (X_2, Y_2) :

$$\|(\nabla_X H(X_1, Y_1) - \nabla_X H(X_2, Y_2), \nabla_Y H(X_1, Y_1) - \nabla_Y H(X_2, Y_2))\| \leq M \|(X_1 - X_2, Y_1 - Y_2)\|. \quad (3)$$

Subdifferentials of nonconvex and nonsmooth functions

Definition 1. Let $\sigma : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ be a PLS function. For a given $x \in \text{dom } \sigma$, the Frechet subdifferential of σ at x , written $\hat{\partial}\sigma(x)$, is the set of all vectors $u \in \mathbb{R}^d$ which satisfy

$$\liminf_{\substack{y \neq x \\ y \rightarrow x}} \frac{\sigma(y) - \sigma(x) - \langle u, y - x \rangle}{\|y - x\|} \geq 0. \quad (4)$$

When $x \in \text{dom } \sigma$, we set $\hat{\partial}\sigma(x) = \emptyset$.

Proposition 1. $\partial(\lambda f(x)) = \lambda \partial f(x)$ for any $\lambda > 0$.

The proposition can be proved based on Definition 1.

Proximal map

Let $\sigma : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ be a PLS function. Given $x \in \mathbb{R}^d$ and $t > 0$, the proximal map associate to σ is defined by:

$$\text{prox}_\lambda^\sigma(x) := \arg \min \left\{ \sigma(u) + \frac{\lambda}{2} \|u - x\|^2, u \in \mathbb{R}^d \right\} \quad (5)$$

The proximal map has the following important property (Lemma 3.2 in [2]).

Lemma 1. Let $h : \mathbb{R}^d \rightarrow \mathbb{R}$ be a continuously differentiable function with gradient ∇h assumed L_h Lipschitz continuous and let $\sigma : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ be a proper and lower semi-continuous function with $\inf_{\mathbb{R}^d} \sigma > -\infty$. Fix any $t > L_h$, then for any $u \in \text{dom } \sigma$ and any $u^+ \in \mathbb{R}^d$ defined by

$$u^+ \in \text{prox}_t^\sigma \left(u - \frac{1}{t} \nabla h(u) \right), \quad (6)$$

we have

$$h(u^+) + \sigma(u^+) \leq h(u) + \sigma(u) - \frac{1}{2}(t - L_h) \|u^+ - u\|^2. \quad (7)$$

Supplementary Note 2: GPALM Algorithm and its Convergence Analysis

The Algorithm

Here we first write out the algorithm that is able to solve the problem (M) with a convergence guarantee.

Convergence analysis

The proof procedure has followed the proofs introduced in the original PALM algorithm [2].

Supplementary Table 1: The GPALM algorithm.

Algorithm 1: The algorithm for the problem (M).

Initialization: X^0, Y^0 , and Z^0 .

1 **for** $k = 0, 1, \dots, K$ **do**

2

$$X^{k+1} \in \text{prox}_{\mu_X^k}^{H(\cdot, Y^k)}(X^k) \quad (8)$$

3 **for** $i = 0, 1, \dots, m$ **do**

4 $Y_i^{k+1} \in \text{prox}_{\mu_{Y_i}^k}^{F(\cdot, Z^k)}\left(Y_i^k - \frac{1}{\mu_{Y_i}^k} \nabla_{Y_i} H(X^{k+1}, Y_i^k)\right)$ (9)

5 **end**

$$Z^{k+1} \in \text{prox}_{\mu_Z^k}^{F(Y^{k+1}, \cdot)}(Z^k) \quad (10)$$

6 **end**

Theorem 1. Assume $\Psi(B)$ is a PLS function with $\inf \Psi > -\infty$, the sequence $\{B^k\}_{k \in \mathbb{N}}$ is a Cauchy sequence and converges to a critical point of $\Psi(B)$, if the following four conditions hold [2]:

(i) *Sufficiently decreasing:* there exist some positive constant $\rho_1 > 0$, such that

$$\Psi(B^k) - \Psi(B^{k+1}) \geq \rho_1 \|B^{k+1} - B^k\|^2, \forall k. \quad (11)$$

(ii) *Relative error:* there exist some positive constant $\rho_2 > 0$, such that for any $w^k \in \partial\Psi(B^k)$,

$$\|w^k\| \leq \rho_2 \|B^{k+1} - B^k\|, \forall k. \quad (12)$$

(iii) *Continuity:* there exist a subsequence $\{B^{k_j}\}_{j \in \mathbb{N}}$ and B^* , such that

$$B^{k_j} \rightarrow B^*, \Psi(B^{k_j}) \rightarrow \Psi(B^*), \text{ as } j \rightarrow +\infty. \quad (13)$$

(iv) *KL property:* Ψ satisfies KL property in its effective domain.

By the theorem above, we only need to check that the sequence generated by Algorithm 6 satisfies the conditions (i) - (iv).

Proposition 2. Algorithm 6 is a global convergence algorithm.

Proof. Follow Theorem 1, we prove Algorithm 6 satisfies conditions (i)- (iv).

Condition (i). Based on (8), we know

$$X^{k+1} \in \text{prox}_{\mu_X^k}^{H(\cdot, Y^k)}(X^k) = \arg \min \left\{ H(X, Y^k) + \frac{1}{\mu_X^k} \|X - X^k\|, X \in \mathbb{R}^n \right\}, \quad (14)$$

which implies

$$H(X^{k+1}, Y^k) + F(Y^k, Z^k) \leq H(X^k, Y^k) + F(Y^k, Z^k) - \frac{\mu_X^k}{2} \|X^{k+1} - X^k\| \quad (15)$$

We then apply Lemma 1 to (9),

$$H(X^{k+1}, Y_i^{k+1}) + F(Y_i^{k+1}, Z^k) \leq H(X^{k+1}, Y_i^k) + F(Y_i^k, Z^k) - \frac{1}{2} (\mu_{Y_i}^k - L_Y(X^{k+1})) \|Y_i^{k+1} - Y_i^k\| \quad (16)$$

Similar to the derivation related to X , for Z we get

$$H(X^{k+1}, Y^{k+1}) + F(Y^{k+1}, Z^{k+1}) \leq H(X^{k+1}, Y^{k+1}) + F(Y^{k+1}, Z^k) - \frac{\mu_Z^k}{2} \|Z^{k+1} - Z^k\| \quad (17)$$

Let $B^k = (X^k, Y^k, Z^k)$ and sum over equations from (15) to (17). We have

$$\Psi(B^{k+1}) \leq \Psi(B^k) - \frac{\mu_X^k}{2} \|X^{k+1} - X^k\| - \sum_i \frac{1}{2} (\mu_{Y_i}^k - L_{Y_i}(X^{k+1})) \|Y_i^{k+1} - Y_i^k\| - \frac{\mu_Z^k}{2} \|Z^{k+1} - Z^k\|. \quad (18)$$

We know that μ_X^k , μ_Y^k , and μ_Z^k have their lower bound and $\mu_{Y_i}^k > L_Y(X^{k+1})$. Therefore, we can get $\rho_1 = \frac{\mu_Y^k}{2} + \sum_i \frac{1}{2} (\mu_{Y_i}^k - L_{Y_i}(X^{k+1})) + \frac{\mu_Z^k}{2}$. Then for $B^k = (X^k, Y^k, Z^k)$ we have

$$\Psi(B^k) - \Psi(B^{k+1}) \geq \rho_1 \|B^{k+1} - B^k\|^2, \forall k \quad (19)$$

tha proves condition (i).

Condition (ii). Writing down the optimality condition for (8), we have

$$\nabla_X H(X^{k-1}, Y^{k-1}) + \mu_X^{k-1} (X^k - X^{k-1}) = 0. \quad (20)$$

Let $w_X^k := -\mu_X^{k-1} (X^k - X^{k-1}) - \nabla_X H(X^{k-1}, Y^{k-1}) + \nabla_X H(X^k, Y^k)$. It is easy to prove that $w_X^k \in \partial_X \Psi(X^k, Y^k, Z^k)$. Then

$$\begin{aligned} \|w_X^k\| &\leq \mu_X^{k-1} \|X^k - X^{k-1}\| + \|\nabla_X H(X^k, Y^k) - \nabla_X H(X^{k-1}, Y^{k-1})\| \\ &\leq \mu_X^{k-1} \|X^k - X^{k-1}\| + M (\|X^k - X^{k-1}\| + \|Y^k - Y^{k-1}\|) \\ &\leq (\mu_X^{k-1} + 2M) \|B^k - B^{k-1}\|. \end{aligned} \quad (21)$$

The first inequality comes from the fact that ∇H is Lipschitz continuous on bounded subset $\mathbb{R}^n \times \mathbb{R}^m$ as assumed in Assumption 1 (6).

With the optimality condition for (9), we have

$$\nabla_{Y_i} H(X^k, Y_i^{k-1}) + \mu_{Y_i}^{k-1} (Y_i^k - Y_i^{k-1}) + \partial_{Y_i} F(Y_i^k, Z^{k-1}) = 0. \quad (22)$$

Let $w_{Y_i}^k := -\mu_{Y_i}^k (Y_i^{k+1} - Y_i^k) - \nabla_{Y_i} H(X^k, Y_i^{k-1}) + \nabla_{Y_i} H(X^k, Y_i^k) - \partial_{Y_i} F(Y_i^k, Z^{k-1}) + \partial_{Y_i} F(Y_i^k, Z^k)$. Clearly, $w_{Y_i}^k \in \partial_Y \Psi(X^k, \dots, Y_{i-1}^k, Y_i^k, Y_{i+1}^k, \dots, Z^k)$, then we have

$$\begin{aligned} \|w_{Y_i}^k\| &\leq \mu_{Y_i}^k \|Y_i^{k+1} - Y_i^k\| + \|\nabla_{Y_i} H(X^k, Y_i^k) - \nabla_{Y_i} H(X^k, Y_i^{k-1})\| + \|\partial_{Y_i} F(Y_i^k, Z^k) - \partial_{Y_i} F(Y_i^k, Z^{k-1})\| \\ &\leq \mu_{Y_i}^k \|Y_i^{k+1} - Y_i^k\| + M_{Y_i} \|Y_i^{k+1} - Y_i^k\| + \|\partial_{Y_i} (p_i(Z^k)g_i(Y_i)) - \partial_{Y_i} (p_i(Z^{k-1})g_i(Y_i))\| \\ &\leq (\mu_{Y_i}^k + M_{Y_i}) \|Y_i^{k+1} - Y_i^k\| + \|\partial g_i(Y_i) (p_i(Z^k) - p_i(Z^{k-1}))\| \\ &\leq (\mu_{Y_i}^k + M_{Y_i}) \|Y_i^{k+1} - Y_i^k\| + M_i^Z \|\partial g_i(Y_i)\| \|Z^k - Z^{k-1}\| \\ &\leq (\mu_{Y_i}^k + M_{Y_i} + M_i^Z U_{Y_i}) \|B^k - B^{k-1}\|. \end{aligned} \tag{23}$$

The second inequality utilizes the structure of $F(Y, X)$ introduced in Assumption 1 (2). The third inequality uses Proposition 1. We set $M_{Y_i} > L_{Y_i}(X)$, $M_i^Z > L_i(Z)$, and $U_{Y_i} > U_i$. Similar to things related to X , writing down the optimality condition for (10),

$$\nabla_Z F(Y^k, Z^{k-1}) + \mu_Z^{k-1} (Z^k - Z^{k-1}) = 0. \tag{24}$$

Let $w_Z^k := -\mu_Z^{k-1} (Z^k - Z^{k-1}) - \nabla_Z F(Y^k, Z^{k-1}) + \nabla_Z F(Y^k, Z^k)$. We find that $w_Z^k \in \partial_Z \Psi(X^k, Y^k, Z^k)$ and we have

$$\|w_Z^k\| \leq (\mu_Z^{k-1} + M_Z) \|B^{k+1} - B^k\|, \tag{25}$$

where $M_Z > L_Z(Y)$.

Let $\rho_2 = \max\{\mu_X^{k-1} + 2M, \mu_{Y_i}^k + M_{Y_i} + M_i^Z U_{Y_i}, \mu_Z^{k-1} + M_Z\}$ and sum (21), (23), (25), we have

$$\|w^k\| \leq \rho_2 \|B^{k+1} - B^k\|, \tag{26}$$

where $w^k = (w_X^k, \dots, w_Y^k, \dots, w_Z^k) = (\partial_X^k \Psi, \dots, \partial_{Y_1^k} \Psi, \dots, \partial_{Z^k} \Psi) = \partial \Psi(X^k, Y^k, X^k) \in \partial \Psi(B^k)$.

Condition (iii). Summing (19) from $k = 0$ to $N - 1$ we have

$$\rho_1 \sum_k^{N-1} \|B^{k+1} - B^k\|^2 \leq \Psi(B^0) - \Psi(B^N) \tag{27}$$

Since $\{\Psi(B^N)\}$ is decreasing and $\inf \Psi > -\infty$, there exist some $\bar{\Psi}$ such that $\Psi(B^N) \rightarrow \bar{\Psi}$ as $N \rightarrow +\infty$. Therefore, let $N \rightarrow +\infty$ in (27), we have

$$\rho_1 \sum_k^{+\infty} \|B^{k+1} - B^k\|^2 \leq \Psi(B^0) - \bar{\Psi}, \tag{28}$$

which implies that $\lim \|B^k - B^{k-1}\| = 0$. Let $B^* = (X^*, Y^*, Z^*)$ be a limit point of $\{B^k\}_{k \in \mathbb{N}} = \{(X^k, Y^k, Z^k)\}_{k \in \mathbb{N}}$. Then (28) indicates that there is a subsequence $\{(X^{k_j}, Y^{k_j}, Z^{k_j})\}_{j \in \mathbb{N}}$ such that $(X^{k_j}, Y^{k_j}, Z^{k_j}) \rightarrow (X^*, Y^*, Z^*)$ as $j \rightarrow +\infty$.

From (9), we know

$$Y_i^{k+1} \in \arg \min \left\{ \langle Y - Y_i^k, \nabla_{Y_i} H(X^k, Y_i^k) \rangle + \frac{\mu_{Y_i}^k}{2} \|Y - Y_i^k\|^2 + F(Y, Z^k) \right\} \tag{29}$$

Let $Y = Y_i^*$ the limiting point of $\{Y_i^k\}_{k \in \mathbb{N}}$, we have

$$\begin{aligned} & \langle Y_i^{k+1} - Y_i^k, \nabla_{Y_i} H(X^k, Y_i^k) \rangle + \frac{\mu_{Y_i}^k}{2} \|Y_i^{k+1} - Y_i^k\|^2 + F(Y_i^{k+1}, Z^k) \\ & \leq \langle Y_i^* - Y_i^k, \nabla_{Y_i} H(X^k, Y_i^k) \rangle + \frac{\mu_{Y_i}^k}{2} \|Y_i^* - Y_i^k\|^2 + F(Y_i^*, Z^k) \end{aligned} \quad (30)$$

Set $k = k_j - 1$, we obtain

$$\begin{aligned} & \langle Y_i^{k_j} - Y_i^{k_j-1}, \nabla_{Y_i} H(X^{k_j-1}, Y_i^{k_j-1}) \rangle + \frac{\mu_{Y_i}^{k_j-1}}{2} \|Y_i^{k_j} - Y_i^{k_j-1}\|^2 + F(Y_i^{k_j}, Z^{k_j-1}) \\ & \leq \langle Y_i^* - Y_i^{k_j-1}, \nabla_{Y_i} H(X^{k_j-1}, Y_i^{k_j-1}) \rangle + \frac{\mu_{Y_i}^{k_j-1}}{2} \|Y_i^* - Y_i^{k_j-1}\|^2 + F(Y_i^*, Z^{k_j-1}) \end{aligned} \quad (31)$$

Let $j \rightarrow +\infty$, we get

$$\limsup_{j \rightarrow +\infty} F(Y_i^{k_j}, Z^{k_j-1}) \leq F(Y_i^*, Z^*) \quad (32)$$

From the fact that F is a PLS function, we also have

$$\limsup_{j \rightarrow +\infty} F(Y_i^{k_j}, Z^{k_j-1}) \geq F(Y_i^*, Z^*) \quad (33)$$

Based on (32) and (33), we know $\lim_{j \rightarrow +\infty} F(Y_i^{k_j}, Z^{k_j-1}) = F(Y_i^*, Z^*)$. Arguing similarly with X , we finally have

$$\lim_{j \rightarrow +\infty} \Psi(X^{k_j}, Y^{k_j}, Z^{k_j}) = \lim_{j \rightarrow +\infty} H(X^{k_j}, Y^{k_j}) + F(Y^{k_j}, Z^{k_j}) = \Psi(X^*, Y^*, Z^*). \quad (34)$$

Condition (iv). The function Ψ is a semi-algebraic function, which automatically satisfies the Kurdyka-Lojasiewicz property [2]. \square

Supplementary Note 3: Parameter Selection for Algorithms used in the study

In this section, we introduce how we select parameters for the competing algorithms.

Parameter Selection for PriroSum

PriorSum constructs a predicted GRN by summing overweights from all prior networks $P = \{P^1, \dots, P^d\}$. Therefore, PriorSum builds a GRN $\mathfrak{P}_{ij} = \sum_k P_{ij}^k$ and does not need to select any parameters.

Parameter Selection for LassoStARS

LassoStARS [3] is the latest version of Inferelator, it takes an unweighted prior and gene expression data as input. Because LassoStARS needs an unweighted prior network and the prior networks we have are weighted prior networks, we choose different cutoffs to construct prior networks for LassoStARS. We generate prior networks by assigning each gene the top N TFs based on the \mathfrak{P}_{ij} . For N , we set $N = \{10, 20, 30, 40\}$ and we find that $N = 10$ performs the best and report the results in the main paper. For other parameters used in LassoStARS, LassoStARS proposed a way to select the optimal parameters, therefore, we do not need to select other parameters.

Parameter Selection for MerlinP

For reconstructing the GRN for yeast, MerlinP [4] uses the same prior networks and gene expression to build a GRN and reported in the repository <https://github.com/Roy-lab/merlin-p>. We directly download the GRN they build and compared it with other methods. For reconstructing the S2 cell GRN and cell-specific GRNs, we follow the instruction provided in <https://github.com/Roy-lab/merlin-p>.

Parameter Selection for NetREX

NetREX [5] is similar to LassoStARS, taking an unweighted prior and gene expression as input. So similarly, we generate prior networks for NetREX by assigning each gene the top N TFs based on the \mathfrak{P}_{ij} . We set $N = \{10, 20, 30, 40\}$ and we find that $N = 20$ performs the best and report the results in the main paper. For the other parameters, we selected based on the suggestion provided in <https://github.com/ncbi/NetREX>.

Parameter Selection for CF

We input CF [6] with $\mathfrak{P}_{ij} = \sum_k P_{ij}^k$. The dimension of the hidden feature vector we set it to be 100, 200, and 300. The regulation term used by CF is set to be 0.1, 1, 10, 100. We try all those combinations and report the result with the best performance.

Parameter Selection for NetREX-CF

Based on the formulaiton of NetREX-CF (??), we know that we need to select h , λ_A , λ_S , η_{ij} , λ , and \bar{C}_{ij} . h is the dimension of the hidden feature vector. We find that $h = \{100, 200, 300\}$ does not change the performance much. For computational consideration, we set $h = 100$. Because λ_A and λ_S are used as standard regulation to avoid over-fitting, we set $\lambda_A = 1.0$ and $\lambda_S = 1.0$ by default. We introduce the selection of η_{ij} and \bar{C}_{ij} in the following subsection.

Selection of η_{ij}

We need to make sure $F(S, X, Y)$ is lower semi-continuous. We can first simplify the equation into

$$\begin{aligned}
F(S, X, Y) &= \lambda \left[\sum_{i,j} \Omega_{ij} (\|S_{ij}\|_0 + (1 - \|S_{ij}\|_0) B_{ij} - x_i^T y_j)^2 \right] + \sum_{i,j} \eta_{ij} \|S_{ij}\|_0 \\
&= \lambda \left[\sum_{i,j} (C_{ij} + (\bar{C}_{ij} - C_{ij}) \|S_{ij}\|_0) (\|S_{ij}\|_0 + (1 - \|S_{ij}\|_0) B_{ij} - x_i^T y_j)^2 \right] + \sum_{i,j} \eta_{ij} \|S_{ij}\|_0 \\
&= \sum_{i,j} \left\{ \lambda [\bar{C}_{ij}(1 - B_{ij})(1 + 2(B_{ij} - x_i^T y_j)) + (\bar{C}_{ij} - C_{ij})(B_{ij} - x_i^T y_j)^2] + \eta_{ij} \right\} \|S_{ij}\|_0 \\
&\quad + \sum_{ij} C_{ij} (B_{ij} - x_i^T y_j)^2
\end{aligned} \tag{35}$$

$F(S, X, Y)$ is lower semi-continuous when the parameter before $\|S_{ij}\|_0$ in the above equation is larger than 0. After several manipulations, we find out we need to set η_{ij} as follows to make $F(S, X, Y)$ lower semi-continuous.

$$\eta_{ij} = \begin{cases} \geq 0, & B_{ij} = 1, \\ \geq \lambda \frac{C_{ij} \bar{C}_{ij}}{C_{ij} - C_{ij}}, & B_{ij} = 0. \end{cases} \tag{36}$$

Selection of \bar{C}_{ij}

C_{ij} is the penalty when we want to use $x_i^T y_j$ to learn $B_{ij} = 1$. Similarly, \bar{C}_{ij} is the penalty when we want to use $x_i^T y_j$ to learn $\|S_{ij}\|_0 = 1$. There are two situations. First, when $\|S_{ij}\|_0 = 1$ and $B_{ij} = 1$, meaning the sparse NCA-based method confirms the edge in the prior, then intuitively, we need to set $\bar{C}_{ij} = \alpha C_{ij}$, $\alpha \geq 1$. Another situation is that $\|S_{ij}\|_0 = 1$ and $B_{ij} = 0$, meaning the sparse NCA-based model confirms an edges recommended by the CF model but not appeared in the prior networks. For this case, we set $\bar{C}_{ij} \in [C_{ij}, \max(C)]$, where $\max(C)$ is the largest element in penalty matrix C . In sum, $\bar{C}_{ij} = \alpha C_{ij} \|S_{ij}\|_0 B_{ij} + \beta \|S_{ij}\|_0 (1 - B_{ij})$, where $\alpha \geq 1$ and $\beta \in [C_{ij}, \max(C)]$.

Consensus of Different Parameter Selections

As explained in the previous, for η_{ij} and \bar{C}_{ij} , we know the range of these parameters but do not know the exact optimal values. For reconstructing GRN for the yeast experiment, we set

$$\eta_{ij} = \begin{cases} \geq \theta, & B_{ij} = 1, \\ \geq \lambda \frac{C_{ij} \bar{C}_{ij}}{C_{ij} - C_{ij}} + \theta, & B_{ij} = 0, \end{cases} \tag{37}$$

where $\theta = \{0.1, 0.5, 1, 2\}$. And $\bar{C}_{ij} = \alpha C_{ij} \|S_{ij}\|_0 B_{ij} + \beta \|S_{ij}\|_0 (1 - B_{ij})$, where $\alpha = \{1, 2, 3, 10\}$ and $\beta = 10, 20, 30, 40$. For different set of parameters, we get a GRN and we get a set of GRNs $\mathfrak{G} = \{G^1, \dots\}$, where $G^i = X^T Y$ after applying all these parameters. The final prediction is the average overall predictions $G^* = \frac{\sum_i G^i}{|\mathfrak{G}|}$.

Supplementary Note 4: Data processing

ScRNA-Seq data

For all scRNA-Seq data used in the manuscript, to reduce the impact of the sparsity of the scRNA-Seq data, we eliminated cells with less than 500 genes expressed and genes that are expressed in fewer than 10% of the cells.

Supplementary References

- [1] W. J. Kent. “BLAT—the BLAST-like alignment tool”. In: *Genome Res* 12.4 (Apr. 2002), pp. 656–664.
- [2] Jérôme Bolte, Shoham Sabach, and Marc Teboulle. “Proximal alternating linearized minimization for nonconvex and nonsmooth problems”. In: *Mathematical Programming* 146.1-2 (2014), pp. 459–494.
- [3] E. R. Miraldi et al. “Leveraging chromatin accessibility for transcriptional regulatory network inference in T Helper 17 Cells”. In: *Genome Res.* 29.3 (Mar. 2019), pp. 449–463.
- [4] Alireza F Siahpirani and Sushmita Roy. “A prior-based integrative framework for functional transcriptional regulatory network inference.” In: *Nucleic acids research* 45.4 (2016), gkw963.
- [5] Y. Wang et al. “Reprogramming of regulatory network using expression uncovers sex-specific gene regulation in *Drosophila*”. In: *Nat Commun* 9.1 (Oct. 2018), p. 4061.
- [6] Y. Hu, Y. Koren, and C. Volinsky. “Collaborative filtering for implicit feedback datasets”. In: *Data Mining, 2008. ICDM’08. Eighth IEEE International Conference on*. IEEE. 2008, pp. 263–272.