

SUPPLEMENTARY INFORMATION FOR:

Universal Forward Training and Structure-free Learning

This PDF file includes:

1. Supplementary Texts

Hypothesis 1-3

System Matrixization

Geometric Understanding

Convergence analysis

System analysis

2. Supplementary Methods

Method S1-S4

Algorithm S1-S7

3. Supplementary Experiments S1-S6

4. Supplementary Tables S1-S6

5. Supplementary Figures S1-S10

Supplementary Texts

This paper carries out a theoretical analysis and builds a theoretical system based on Hypothesis 1-3. The Hypotheses used as the basis only need to be reasonable and applicable and do not necessarily need to be proved. However, the whole mathematical calculus process can deepen the understanding of the relationship and meaning of each hypothesis.

Hypothesis 1

Hypothesis 1: For a multi-input-single-output (MISO) function with continuous derivatives (smooth)

$$f: x \mapsto y = f(x), x = [x^1, x^2, \dots, x^P],$$

(a) (Existentiality): For any certain data point (x_0, y_0) , the function has the local linearization:

$$y_0 = f(x_0) = C_0 + Cx_0^T, \text{ where } C = [C_1, \dots, C_P].$$

(b) (Uniqueness): If C_p is the corresponding partial derivative of the output y_0 to the input x_0^p , then C is unique.

(c) (Approximation): In the sufficiently small neighbor of any certain data point (x_0, y_0) , any point (x, y) can be approximated as:

$$y = f(x) \approx C_0 + Cx^T,$$

where C and C_0 is determined by (x_0, y_0) (that is $y_0 = f(x_0) = C_0 + Cx_0^T$ in Hypothesis 1(a)(b)).

(d) (Continuity): C and C_0 are also the continuous function of input x .

Proof:

(a)-(b) (Existentiality and Uniqueness):

For a multi-input-single-output (MISO) function

$$f: x \mapsto y = f(x), x = [x^1, x^2, \dots, x^P],$$

with continuous derivatives is a continuous and derivable function.

If the mapping function is continuous and derivable, then $\frac{\partial y}{\partial x^p}$ is existent and unique.

For any certain data point (x_0, y_0) , define

$$C_p \triangleq \frac{\partial y}{\partial x^p} \big|_{x=x_0}$$

and

$$C_0 = y_0 - \sum_{p=1}^P C_p x_0^p$$

So

$$y_0 = f(x_0) = C_0 + \sum_{p=1}^P C_p x_0^p = C_0 + C x_0^T$$

where, $C = [C_1, \dots, C_P]$.

According to the existence and uniqueness of C_p , **Hypothesis 1**(a) and (b) can be proved.

(c) (Approximation):

The proof of **Hypothesis 1**(c) is obviously easy to be accepted because of continuity and “substituting curve with straight”.

According to the “substituting curve with straight”, for a smooth “curve”, we can approximate it with “straight”. For any data point (x, y) in the sufficiently small neighbor of a certain data point (x_0, y_0) , that is

$$y = f(x) \approx C'_0 + C' x^T.$$

In the neighbor area of (x_0, y_0) , all the data satisfy it, including (x_0, y_0) . And

$$y_0 = f(x_0) = C_0 + C x_0^T,$$

by Hypothesis 1(a).

So $C'_0 = C_0$, $C' = C$ is a reasonable solution (The derivation cannot guarantee the uniqueness of the solution).

According to the continuity of $f(x)$ and $g(x) = C_0 + C x^T$,

$$\lim_{x \rightarrow x_0} f(x) - g(x) = f(x_0) - g(x_0) = 0.$$

Then, for a sufficiently small $\varepsilon > 0$ and $\varepsilon \rightarrow 0$,

$\exists \delta > 0$, when $|x - x_0| < \delta$ then $|f(x) - g(x)| < \varepsilon$,

So

$$y = f(x) \approx g(x) = C_0 + Cx^T.$$

(d) (Continuity):

(1) C and C_0 , as observable properties variables, are the function of input x , according to Causality and Determinism of **Hypothesis 2**.

(2) According to the continuity of derivatives (The premise of **Hypothesis 1**), C (**Hypothesis 1(b)**, Uniqueness, C_p is the corresponding partial derivative) is the continuous function of input x .

Because

$$C_0 = y - \sum_{p=1}^P C_p x^p$$

and the continuity of the function and C , C_0 is also the continuous function of input x .

Proven.

Extension of Hypothesis 1(d):

First, try to explain that C and C_0 , as observable properties variables, are the function of input x . It is obvious because only the input can change other variables when the function is fixed, according to Causality and Determinism of **Hypothesis 2**. If not, in this situation, then C and C_0 are constant, and $y = C_0 + Cx^T$ just a linear function, which does not conform to the premise of any function. Thus, C and C_0 are the function of input x . The case of constant (linear function) can be regarded as a special case.

Secondly, for any data point (x', y') in the sufficiently small neighbor of a certain data point (x, y) , and

$$\begin{cases} y = f(x) = C_0 + Cx^T = C_s[1; x^T] = C_sX \\ y' = f(x') = C'_0 + C'x'^T = C'_s[1; x'^T] = C'_sX' \end{cases}$$

where $X = [1; x]$, $X' = [1; x'^T]$, $C_s = [C_0, C]$, $C'_s = [C'_0, C']$

$$\begin{cases} C_s = X^\dagger y \\ C'_s = X'^\dagger y' \end{cases}$$

then according to the continuity of function f , it can be seen that

$$\lim_{x' \rightarrow x} y' = y, \quad \lim_{x' \rightarrow x} X' = X$$

According to the continuity of the generalized inverse matrix (Ref. 29, 30), it can be obtained that:

$$\text{rank}(X') = \text{rank}(X) = 1,$$

then

$$\lim_{x' \rightarrow x} X'^\dagger = X^\dagger.$$

Thus

$$\lim_{x' \rightarrow x} C'_s = C_s$$

It indicates continuity.

At the same time, for sufficiently small $\Delta x = x' - x$, $C'_s \approx C_s$, same as **Hypothesis 3(b)**.

Hypothesis 2

Hypothesis 2 as follows is based on “causality and determinism” in the classical scientific view. Stable means that the output is fixed in the same case. Independence is to simplify the follow-up analysis.

Hypothesis 2: For a stable parametric system,

(a) (Causality and Determinism): Any output and observable properties (or characters) of any system are determined stably by the (i) input, (ii) system parameters (weights), and (iii) system structure (the sum of covert parameters, hidden variables, topology, and other influencing factors).

(b) (Independence and Causal Decoupling): The effect of (i) input, (ii) parameters, and (iii) structure could be considered independent of each other.

“Causality and Determinism” in Hypothesis 2 is a classical scientific view, not suitable for all systems like chaotic systems in an engineering sense. Similarly, independence is a simplified and idealized treatment.

By Hypothesis 2, the matrix $\mathbb{C}_X = \aleph(X, W, SS) = \aleph_X(SS, W) = \aleph_{X,W}(SS) = \aleph_{X,SS}(W)$. Thus \mathbb{C}_X could be an observable window to analyze the system structure and parameters, and a representative of the system. This is the ideological basis on \mathbb{C}_X of system design in the next section.

Hypothesis 3

Hypothesis 3: For two stable parametric systems processed by LL

$$\begin{cases} Y = \mathbb{H}(W) = \mathbb{C}_X W \\ Y' = \mathbb{H}'(W') = \mathbb{C}'_X W' \end{cases}$$

(a) (Invariability): The systems with the same structure and parameters ($W = W'$), for the same input vector X of the dataset, have the same output vector ($Y = Y'$) and the same constant matrix ($\mathbb{C}_X = \mathbb{C}'_X$).

(b) (Comparability): The systems with the same structure and **comparable** parameter ($W \approx W'$), for the same input vector X of the dataset, have the **comparable** constant matrix ($\mathbb{C}_X \approx \mathbb{C}'_X$), and

$$\lim_{W' \rightarrow W} \mathbb{C}'_X = \mathbb{C}_X .$$

Proof:

(a) (Invariability):

Suppose a dataset contains N_d data $\{(x_i, y_i) \mid i = 1, \dots, N_d\}$, where output $y = [y^1, y^2, \dots, y^Q]$, input $x = [x^1, x^2, \dots, x^P]$, and let $\mathbf{Y} = [y_1, y_2, \dots, y_{N_d}]^T_{(N_d \times Q) \times 1}$, $\mathbf{X} = [x_1, x_2, \dots, x_{N_d}]^T_{N_d \times 1}$,

$$\mathbb{C}_X = [C_{x_1}^{(1)}; C_{x_1}^{(2)}; \dots; C_{x_1}^{(Q)}; C_{x_2}^{(1)}; \dots; C_{x_{N_d}}^{(1)}; C_{x_{N_d}}^{(2)}; \dots; C_{x_{N_d}}^{(Q)}]_{(N_d \times Q) \times (S+1)}, \quad \mathbf{W} = [1, W]^T_{(S+1) \times 1}.$$

Because the systems with the same structure and parameters (weights) \mathbf{W} , for the same input x of the dataset, have the same mapping function and output (by **Hypothesis 2**), that is,

$$y = F(x, W) = H(W)$$

So, the systems have the same output vector $\mathbf{Y} = [y_1, y_2, \dots, y_{N_d}]^T_{(N_d \times Q) \times 1}$.

For any data (x, y) in $\{(x_i, y_i)\}$, according to the existence and uniqueness of \mathbb{C} (by **Hypothesis 1**), the existence and uniqueness of

$$[C_x^{(1)}; C_x^{(2)}; \dots; C_x^{(Q)}]_{Q \times (S+1)}$$

is correct, and

$$y = H(W) = \mathbb{C}_x \mathbf{W}$$

Make further efforts, for all the data in the dataset $\{(x_i, y_i)\}$, the existence and uniqueness of

$$\mathbb{C}_X = [C_{x_1}^{(1)}; C_{x_1}^{(2)}; \dots; C_{x_1}^{(Q)}; C_{x_2}^{(1)}; \dots; C_{x_{N_d}}^{(1)}; C_{x_{N_d}}^{(2)}; \dots; C_{x_{N_d}}^{(Q)}]_{(N_d \times Q) \times (S+1)}$$

is correct. Thus the systems with the same structure and parameter (weights) \mathbf{W} , for the same input vector \mathbf{X} of the dataset, have the same constant matrix \mathbb{C}_X .

(b) (Comparability): The systems with the same structure and **comparable** parameter ($\mathbf{W} \approx \mathbf{W}'$), for the same input vector \mathbf{X} of the dataset, have the **comparable** constant matrix ($\mathbb{C}_X \approx \mathbb{C}'_X$), and

$$\lim_{\mathbf{W}' \rightarrow \mathbf{W}} \mathbb{C}'_X = \mathbb{C}_X.$$

Define: $|\mathbf{v}| \triangleq |v_1| + |v_2| + \dots + |v_N|$, $\mathbf{v} = [v_1, v_2, \dots, v_N]$.

The weights can be considered to be the input of the current system, for the same \mathbf{X} of the dataset. Continuity (“continuous derivable” of **Hypothesis 1**) is the important prerequisite for all discussions in this paper. According to the **Hypothesis 3(a)** and the **Hypothesis 1(d)**,

$$\lim_{\mathbf{W}' \rightarrow \mathbf{W}} \mathbf{Y}' = \mathbf{Y}, \quad \lim_{\mathbf{W}' \rightarrow \mathbf{W}} \mathbb{C}'_X = \mathbb{C}_X.$$

Thus, for sufficiently small $|\Delta \mathbf{W}| = |\mathbf{W} - \mathbf{W}'|$, $\mathbb{C}'_X \approx \mathbb{C}_X$. **Proven.**

Extension of Hypothesis 3(b)

Hypothesis 3(b) is obviously easy to accept because of continuity.

Define: $|\mathbf{v}| \triangleq |v_1| + |v_2| + \dots + |v_N|$, $\mathbf{v} = [v_1, v_2, \dots, v_N]$.

If $|\mathbf{W} - \mathbf{W}'| = |\Delta\mathbf{W}| = \varepsilon > 0$, ε is sufficiently small.

then $|\mathbf{Y} - \mathbf{Y}'| = |\Delta\mathbf{Y}| = \delta$, $|\mathbb{C}_X - \mathbb{C}'_X| = \theta$ and $\lim_{\varepsilon \rightarrow 0} \delta = 0$, $\lim_{\varepsilon \rightarrow 0} \theta = 0$.

So,

$$\begin{aligned} |\Delta\mathbf{Y}| = |\mathbf{Y} - \mathbf{Y}'| &= |\mathbb{C}_X \mathbf{W} - \mathbb{C}'_X \mathbf{W}'| = |\mathbb{C}_X \mathbf{W} - \mathbb{C}'_X \mathbf{W}| = |\mathbb{C}_X \mathbf{W} - \mathbb{C}_X \mathbf{W}' + \mathbb{C}_X \mathbf{W}' - \mathbb{C}'_X \mathbf{W}'| \\ &\leq |\mathbb{C}_X| |\Delta\mathbf{W}| + |\mathbb{C}_X - \mathbb{C}'_X| |\mathbf{W}'| \end{aligned}$$

$$\Rightarrow |\mathbb{C}_X - \mathbb{C}'_X| \geq \frac{|\Delta\mathbf{Y}| - |\mathbb{C}_X| |\Delta\mathbf{W}|}{|\mathbf{W}'|}$$

The above inequality gives the lower limit of error $|\mathbb{C}_X - \mathbb{C}'_X|$. Two phenomena were observed at the same time.

1) because $\lim_{\mathbf{W}' \rightarrow \mathbf{W}} \mathbf{Y}' = \mathbf{Y}$, and $\lim_{|\Delta\mathbf{W}| \rightarrow 0} |\Delta\mathbf{Y}| = 0$ (**Hypothesis 3(a)**)

$$\lim_{\substack{\mathbf{W}' \rightarrow \mathbf{W} \\ |\Delta\mathbf{W}| \rightarrow 0}} \frac{|\Delta\mathbf{Y}| - |\mathbb{C}_X| |\Delta\mathbf{W}|}{|\mathbf{W}'|} = 0$$

Therefore, $\Delta\mathbf{W}$ needs to be small enough in engineering implementation.

2) The smaller $|\mathbf{W}'|$ is, the larger the lower limit of error $|\mathbb{C}_X - \mathbb{C}'_X|$ is.

So the system weight should not be too small.

End

System Matrixization

For a high-dimensional parameterized system with multi-input-multi-output (MIMO) mapping function, any output y or observable properties (P or characters C) is the function of the input x , the system parameter (weights) W and system structure (SS) in **Hypothesis 2**, labeled as $y = \aleph(x, W, SS)$ or $P = \aleph(x, W, SS)$. \aleph means a high-dimensional complex mapping space.

For a high-dimensional parameterized system with multi-input-multi-output (MIMO) mapping function, any output y is the function of the input x and the system parameter (weights) W , in **Hypothesis 2**, that is

$$y = F(x, W) = [f^1(x, W), f^2(x, W), \dots, f^Q(x, W)]$$

where

$$\text{output } y = [y^1, y^2, \dots, y^Q],$$

$$\text{input } x = [x^1, x^2, \dots, x^P],$$

$$\text{weights } W = [w_1, w_2, \dots, w_S],$$

P, Q, and S are the dimensions of input, output, and weights, respectively,

and $f^\alpha = \pi^\alpha \circ F$ is the α th component function of F .

According to the function linearization LL discussed in **Hypothesis 1**,

$$f^\alpha(x, W) = C_0^\alpha + C^\alpha[x, W]^T = C^{(\alpha)}[1, x, W]^T,$$

then

$$y^T = C_S[1, x, W]^T,$$

where $C_S = [C^{(1)}; C^{(2)}; \dots; C^{(Q)}]$.

Obviously, in the case of x and W fixed, C_S is only related to the structure of the system. This processing also means a new approximate description method of complex mapping space's local characteristics by a hyperplane (determined by C_S). On the other hand, C_S varies with x and W . Therefore, C_S can only describe the local system characteristics of current x and W . In order to describe the global characteristics of the system more comprehensively, more different x and W need to be added. In engineering applications, x and W usually do not seek to change at the same time, so the changes of x and W are discussed respectively.

If the input of the system remains fixed, the output can be changed into a function of the system parameters:

$$y = H(W) = [h^1(W), h^2(W), \dots, h^Q(W)] = C_x[1, W]^T$$

where $C_x = [C_x^{(1)}, C_x^{(2)}, \dots, C_x^{(Q)}]$. The linearization of the function h^α is realized in the neighborhood of W for the data point (x, y^α) , $y^\alpha = C_x^{(\alpha)}[1, W]^T$. If C_x can be obtained, the local spatial distribution characteristics for the current W in (x, y) can be obtained, as shown in **Hypothesis 1(c)**. Of course, only one data point is considered, which obviously cannot describe the whole plant space globally, so more data points need to be added.

Suppose a dataset $(X, Y) = \{(x_i, y_i) | i = 1, \dots, N_d\}$ contains N_d data, and let

$$\mathbf{Y} = [y_1, y_2, \dots, y_{N_d}]^T_{(N_d \times Q) \times 1},$$

$$\mathbf{X} = [x_1, x_2, \dots, x_{N_d}]^T_{N_d \times 1},$$

$$\mathbb{C}_X = [C_{x_1}^{(1)}; C_{x_1}^{(2)}; \dots; C_{x_1}^{(Q)}; C_{x_2}^{(1)}; \dots; C_{x_{N_d}}^{(1)}; C_{x_{N_d}}^{(2)}; \dots; C_{x_{N_d}}^{(Q)}]_{(N_d \times Q) \times (S+1)},$$

$$\mathbf{W} = [1, W]^T,$$

then

$$\mathbf{Y} = \mathbb{H}(\mathbf{W}) = \mathbb{C}_X \mathbf{W}.$$

\mathbb{C}_X (Constant matrix or Character matrix) is a character variables matrix of function, remains constant for the fixed \mathbf{X} . On the contrary, it will change with the change of \mathbf{X} . \mathbb{C} has a clear mathematical meaning, and each element is a corresponding partial derivative of the output to the input, as shown in **Hypothesis 1(b)**. The global characteristics of the function are obtained by all local characteristics of the fitted hyperplanes for all the data of the dataset, when the dataset is large enough, due to **Hypothesis 1(c)**. Thus, the effect, accuracy, or granularity described by the constant matrix is related to the size and distribution of the dataset. Constant matrix, as a novel type of function description method, is different from the traditional mathematical analytic formula, drawing, modeling, or black box method.

Similarly obtains

$$y = \mathfrak{N}(x, W, SS) = \mathfrak{N}_{SS}(x, W) = \mathfrak{N}_{W, SS}(x),$$

$$P = \mathfrak{N}(x, W, SS) = \mathfrak{N}_x(SS, W) = \mathfrak{N}_{x, W}(SS),$$

or other functions and their corresponding matrix form, respectively, like

$$\mathbf{Y} = \mathbb{G}(\mathbf{X}) = \mathbb{C}_W \mathbf{X},$$

$$\text{or } \mathbf{Y} = \mathbb{F}(\mathbf{W}, \mathbf{X}) \approx \mathbb{C}_S[\mathbf{W}; \mathbf{X}],$$

which have obviously different physical meanings and application prospects. In spite of these, they have the same mathematical form. Hence \mathbb{C}_X , \mathbb{C}_S or \mathbb{C}_W are collectively labeled by \mathbb{C} (constant matrix), and the corresponding functional relationship are collectively labeled by

$$\mathbf{Y} = \mathbb{F}(\mathbf{X}) = \mathbb{C}\mathbf{X}.$$

In the paper, $\mathbf{Y} = \mathbb{H}(\mathbf{W}) = \mathbb{C}_X\mathbf{W}$ is taken into consideration to illustrate the solution and application method of the constant matrix.

As a result of local linearization processing, for a function $y = f(x)$, $\mathbf{Y} = \mathbb{F}(\mathbf{X}) = \mathbb{C}\mathbf{X}$ in the case of a given sampling dataset (\mathbf{X}, \mathbf{Y}) . When input matrix \mathbf{X} is a serial fixed sampling data, the constant matrix \mathbb{C} only related to the function f , realizes (i) numerization of the function and (ii) local decoupling of multivariate, by linearization processing, as shown in Hypothesis 1(a). Therefore, the processing method has four characteristics: (i) linearization, (ii) matrixization, (iii) numerization, and (iv) decoupling.

Geometric Understanding

Using the classical concepts of linearization and isomorphism, and three proposed hypotheses, LL-IC system theory is constructed, which is further explained as follows.

System matrixization and system feature matrix \mathbb{C}_X is obtained by the LL process. Local linearity is used to describe the global nonlinearity, which reflects the dialectical unity relations between the whole and the part, linearity and nonlinearity. Visualized understanding of LL is that at each sampling data point, the local linear hyperplane is used to expand its surrounding areas, and the approximate description of the global space is completed like a tent or membrane structure building (e.g. Olympiastadion München and The Shed in NewYork), as shown in Fig. S10b-d. System matrixization has characteristics of linearization, matrixization, quantification, and decoupling.

IC stemming from causal determinism solves the solution of \mathbb{C}_X and W . The core of IC is to ensure that all influences are kept fixed except for the changes (like weights) we care about. The isomorphic system spaces constructed by different parameters W can be transformed by the continuous change of W , which could be imagined as the continuous stretching and torsion of the system space, see Fig. S10e,f. In the process of stretching and torsion, the change of gradient is often limited. This is also the reason of reuse and why \mathbb{C}_X is effective.—

S1-S6 are algorithms constructed directly from the concepts of LL and IC, and the geometric understanding can be found in Fig. S10. In contrast, S7 no longer requires complex matrices, is simpler, less computation and resource requirement. Convergence is achieved by repeated iterations according to $\Delta W = kW_r$.

The geometric understanding of LIFT S7 is presented in Fig.1 e-f. W is the current parameters; W^* is the expected parameters; the black dashed lines WW^* represent the ideal tuning direction; the red solid lines represent the weight-adjusted trajectory ΔW . $\Delta W = kW_r$, the direction of adjustment is determined by randomly generated W_r , which has an obvious error with the ideal gradient direction WW^* . While the systematic error will decrease, and the active ingredient is extracted by scalar product calculation.

Convergence analysis

Obviously, in Fig.1 e-f, the LIFT S7 convergence condition can be obtained as:

$$l_r k \|W_r\| < \|WW'\| \text{ and } \langle E', E \rangle \neq 0.$$

where

$$k = (Y'' - Y')^\dagger (Y - Y') = \frac{\langle E', E \rangle}{\langle E, E \rangle}$$

l_r is the learning rate, WW^* is the ideal gradient direction, WW' has the same direction with ΔW , and $\Delta W^* WW'$ is a isosceles triangle (Idealized abstraction, due to the complex spatial distribution of the system, this abstract result does not necessarily hold.).

According to the geometry understanding, the weight adjustment $l_r \Delta W$ should be in the range of WW' , that is $l_r k \|W_r\| < \|WW'\|$.

Additional Notes:

(1) If $\langle E', E \rangle = 0$, $\Delta W = 0$, no adjustment for current W and training ineffective. The probability of orthogonality of two vectors in a multidimensional space is very low, thus ensuring the success of the algorithm, which converges in the overall process.

(2) Due to $l_r k \|W_r\| < \|WW'\|$, Although $\|WW'\|$ is difficult to obtain in practical engineering, it does give us some hints, such as: $\|W_r\|$ and the learning rate l_r should be appropriately small.

System analysis

System matrixization is a geometric algebraization method, has characteristics of linearization, matrixization, quantification, and decoupling, and enables two-dimensionalized representation of high-dimensional systems. To some extent, \mathbb{C}_X can be used as a tool for systems analysis, like system function.

Here's a list of some pre-thinking.

(1) Similarity

According to the Hypothesis 3: For two stable parametric systems processed by LL. (a) (Invariability): The systems with the same structure and parameters ($W = W'$), for the same input vector X of the dataset, have the same output vector ($Y = Y'$) and the same constant matrix ($\mathbb{C}_X = \mathbb{C}'_X$). (b) (Comparability): The systems with the same structure and **comparable** parameter ($W \approx W'$), for the same input vector X of the dataset, have the **comparable** constant matrix ($\mathbb{C}_X \approx \mathbb{C}'_X$), and $\lim_{W' \rightarrow W} \mathbb{C}'_X = \mathbb{C}_X$.

$$S_i = \|\mathbb{C}_X - \mathbb{C}'_X\|$$

S_i measures the direct similarity between the two systems. The smaller the S_i , the higher the similarity.

(2) Sensitivity

With the change of parameters ($\mathbf{W} \rightarrow \mathbf{W}'$), the system will change ($\mathbb{C}_X \rightarrow \mathbb{C}'_X$). Se is calculated by the following formula to show the sensitivity of the system to parameter changes.

$$Se = \|\mathbb{C}_X - \mathbb{C}'_X\|$$

$\mathbf{W} = [w_1, w_2, \dots, w_s]$, for different parameter components w_s , the sensitivity of the system to different parameters can be measured by Se_s .

$$Se_s = \|\mathbb{C}_X - \mathbb{C}'_X\|_{\Delta w_s}$$

Solving $\max \{Se_s\}$, the maximum corresponds to the most sensitive parameter w_s of the current system.

(3) smoothability

\mathbb{C}_X has a clear mathematical meaning, and each element is a corresponding partial derivative of the output to the input, as shown in **Hypothesis 1(b)**.

$$S_m = \|\mathbb{C}_X\|$$

The smaller the S_m , the smoother the system space.

Solving $\max \{\mathbb{C}_X\}$, the maximum maybe correspond to the most steepest place of the current system.

Notes:

(1) The distribution of (\mathbf{X}, \mathbf{Y}) is known.

(2) Only considering the local space, we can get local similarity, sensitivity and smoothability.

Supplementary Methods

Method S1

Direct Method:

Two completely isomorphic systems with different parameters, which are stable and smooth parametric systems, satisfy

$$\begin{cases} Y = \mathbb{H}(W) = \mathbb{C}_X W \\ Y' = \mathbb{H}'(W') = \mathbb{C}'_X W' \end{cases}$$

respectively according to **Hypothesis 1**, and

$$\mathbb{C}'_X \approx \mathbb{C}_X$$

by **Hypothesis 3**. Then,

$$\mathbb{C}'_X \approx \mathbb{C}_X = YW^\dagger,$$

where W^\dagger is the Moore-Penrose generalized inverse of W .

The performance of the known system is used to directly solve the feature matrix of the unknown system.

Method S2

Difference Method:

Two completely isomorphic systems with different parameters, which are stable and smooth parametric systems, satisfy

$$\begin{cases} Y = \mathbb{H}(W) = \mathbb{C}_X W \\ Y' = \mathbb{H}'(W') = \mathbb{C}'_X W' \end{cases}$$

respectively according to **Hypothesis 1**, and

$$\mathbb{C}'_X \approx \mathbb{C}_X$$

by **Hypothesis 3**. Then,

$$\mathbb{C}'_X \approx \mathbb{C}_X \approx (Y - Y')(W - W')^\dagger = \Delta Y \Delta W^\dagger,$$

where

$$\begin{cases} \Delta Y = Y - Y' \\ \Delta W = W - W' \end{cases},$$

and $\Delta \mathbf{W}^\dagger$ is the Moore-Penrose generalized inverse of $\Delta \mathbf{W}$.

Solve for the feature matrix of the unknown system using the difference in performance between the known and unknown systems.

Method S3

Direct Method:

For a stable and smooth parametric system with training data (\mathbf{X}, \mathbf{Y}) and current weights \mathbf{W}' , its current output vector

$$\mathbf{Y}' = \mathbb{H}(\mathbf{W}') = \mathbb{C}_X \mathbf{W}'.$$

So the system feature matrix

$$\mathbb{C}_X = \mathbf{Y}' \mathbf{W}'^\dagger$$

And the object output and weights are

$$\mathbf{Y} = \mathbb{H}(\mathbf{W}) = \mathbb{C}_X \mathbf{W}.$$

Then update the weights by

$$\mathbf{W} = \mathbb{C}_X^\dagger \mathbf{Y} = (\mathbf{Y}' \mathbf{W}'^\dagger)^\dagger \mathbf{Y} = \mathbf{W}' \mathbf{Y}'^\dagger \mathbf{Y}$$

Let

$$\mathbf{k} = \mathbf{Y}'^\dagger \mathbf{Y}$$

Then

$$\mathbf{W} = \mathbf{W}' \mathbf{k}$$

Thus it can be seen that \mathbf{k} is the scale coefficient of \mathbf{W}' .

Method S4

For a stable and smooth parametric system with training data (\mathbf{X}, \mathbf{Y}) and current weights \mathbf{W}' , its current output vector

$$\mathbf{Y}' = \mathbb{H}(\mathbf{W}') = \mathbb{C}_X \mathbf{W}'.$$

And the object output and weights are

$$\mathbf{Y} = \mathbb{H}(\mathbf{W}) = \mathbb{C}_X \mathbf{W}.$$

Then

$$\Delta \mathbf{W} = \mathbf{W} - \mathbf{W}' = \mathbb{C}_X^\dagger (\mathbf{Y} - \mathbf{Y}') = \mathbb{C}_X^\dagger \mathbf{E},$$

where $\mathbf{E} = \mathbf{Y} - \mathbf{Y}'$.

Because the gap between \mathbf{W} and \mathbf{W}' is always too large to obtain a sufficiently accurate \mathbb{C}_X^\dagger . So we need to reduce this gap to a sufficiently small one. Given a little disturbance parameter \mathbf{W}_r , and the output changes to

$$\mathbf{Y}'' = \mathbb{H}(\mathbf{W}' + \mathbf{W}_r) = \mathbb{C}_X(\mathbf{W}' + \mathbf{W}_r).$$

So the system feature matrix

$$\mathbb{C}_X = (\mathbf{Y}'' - \mathbf{Y}')(\mathbf{W}' + \mathbf{W}_r - \mathbf{W}')^\dagger = (\mathbf{Y}'' - \mathbf{Y}')\mathbf{W}_r^\dagger$$

Then update the weights by

$$\begin{cases} \mathbf{W} = \mathbf{W}' + l_r \Delta \mathbf{W} \\ \Delta \mathbf{W} = \mathbb{C}_X^\dagger \mathbf{E} = ((\mathbf{Y}'' - \mathbf{Y}')\mathbf{W}_r^\dagger)^\dagger (\mathbf{Y} - \mathbf{Y}') = \mathbf{W}_r (\mathbf{Y}'' - \mathbf{Y}')^\dagger (\mathbf{Y} - \mathbf{Y}') \end{cases}.$$

where, l_r is the learning rate, $\mathbf{E}' = \mathbf{Y}'' - \mathbf{Y}'$.

Let

$$k = (\mathbf{Y}'' - \mathbf{Y}')^\dagger (\mathbf{Y} - \mathbf{Y}') = \frac{\langle \mathbf{E}', \mathbf{E} \rangle}{\langle \mathbf{E}, \mathbf{E} \rangle}$$

Then

$$\begin{cases} \mathbf{W} = \mathbf{W}' + l_r \Delta \mathbf{W} \\ \Delta \mathbf{W} = k \mathbf{W}_r \end{cases}$$

It is an iterative training method and trains all training data synchronously in an iteration.

Algorithm S1

Algorithm S1 Direct method

Input: Training data $\{(x_i, y_i)\}$, $i = 1, \dots, N_d$;

System structure and forward mapping $y = F(x, W)$.

1 W' with random initialization;

2 $y'_i = F(x_i, W')$;

3 $Y' = [y'_1, y'_2, \dots, y'_{N_d}]^T$;

4 $C_X \approx C'_X = Y'W'^\dagger$;

5 $Y = [y_1, y_2, \dots, y_{N_d}]^T$;

6 $W = C_X^\dagger Y$.

Output: W .

Algorithm S2

Algorithm S2 Difference method

Input: Training data $\{(x_i, y_i)\}$, $i = 1, \dots, N_d$; r and l_r ;
System structure and forward mapping $y = F(x, W)$.

```
1   $W$  with random initialization in  $[-1, 1]$ ;  
2   $Y = [y_1, y_2, \dots, y_{N_d}]^T$ ;  
3  for  $i = 1$  to  $N_d$  do  
4       $y_i^{(1)} = F(x_i, W)$ ;  
5  end for  
6   $Y^{(1)} = [y_1^{(1)}, y_2^{(1)}, \dots, y_{N_d}^{(1)}]^T$ ;  
7   $E = Y - Y^{(1)}$ ;  
8  while ( $epoch < Max$  and  $MSE(E) > \varepsilon$ )  
9       $W_r$  with random initialization in  $[-r, r]$ ;  
10     for  $i = 1$  to  $N_d$  do  
11          $y_i^{(2)} = F(x_i, W + W_r)$ ;  
12     end for  
13      $Y^{(2)} = [y_1^{(2)}, y_2^{(2)}, \dots, y_{N_d}^{(2)}]^T$ ;  
14      $C_X \approx C'_X = (Y^{(2)} - Y^{(1)})W_r^\dagger$ ;  
15      $\Delta W = C_X^\dagger E$ ;  
16      $W = W + l_r \Delta W$ ;  
17     for  $i = 1$  to  $N_d$  do  
18          $y_i^{(1)} = F(x_i, W)$ ;  
19     end for  
20      $Y^{(1)} = [y_1^{(1)}, y_2^{(1)}, \dots, y_{N_d}^{(1)}]^T$ ;  
21      $E = Y - Y^{(1)}$ ;  
22      $epoch++$ ;  
23 end while
```

Output: W .

Algorithm S3

Algorithm S2 Difference method with reuse of C_x^\dagger .

Input: Training data $\{(x_i, y_i)\}$, $i = 1, \dots, N_d$; r and l_r ;
System structure and forward mapping $y = F(x, W)$.

```

1   $W$  with random initialization in  $[-1, 1]$ ;
2   $Y = [y_1, y_2, \dots, y_{N_d}]^T$ ;
3  for  $i = 1$  to  $N_d$  do
4       $y_i^{(1)} = F(x_i, W)$ ;
5  end for
6   $Y^{(1)} = [y_1^{(1)}, y_2^{(1)}, \dots, y_{N_d}^{(1)}]^T$ ;
7   $E = Y - Y^{(1)}$ ;
8  while ( $epoch < Max$  and  $MSE(E) > \varepsilon$ )
9       $W_r$  with random initialization in  $[-r, r]$ ;
10     for  $i = 1$  to  $N_d$  do
11          $y_i^{(2)} = F(x_i, W + W_r)$ ;
12     end for
13      $Y^{(2)} = [y_1^{(2)}, y_2^{(2)}, \dots, y_{N_d}^{(2)}]^T$ ;
14      $C_x \approx C'_x = (Y^{(2)} - Y^{(1)})W_r^\dagger$ ;
15     while  $MSE(E)$  decreasing
16          $\Delta W = C_x^\dagger E$ ;
17          $W = W + l_r \Delta W$ ;
18     for  $i = 1$  to  $N_d$  do
19          $y_i^{(1)} = F(x_i, W)$ ;
20     end for
21      $Y^{(1)} = [y_1^{(1)}, y_2^{(1)}, \dots, y_{N_d}^{(1)}]^T$ ;
22      $E = Y - Y^{(1)}$ ;
23 end while
24  $epoch++$ ;
25 end while

```

Output: W .

Algorithm S4

Algorithm S4 Difference method with batch process of training data

Input: Training data $\{(x_i, y_i)\}$, $i = 1, \dots, N_d$; r , l_r , and M_b ;
System structure and forward mapping $y = F(x, W)$.

```
1   $W$  with random initialization in  $[-1, 1]$ ;  
2   $Y = [y_1, y_2, \dots, y_{N_d}]^T$ ;  
3  while (iteration  $< Max$  and  $MSE(E) > \varepsilon$ )  
4    for  $k = 0$  to  $(N_d/M_b - 1)$  do  
5       $Y = [y_{k+1}, y_{k+2}, \dots, y_{k+M_b}]^T$ ;  
6      for  $i = (k + 1)$  to  $(k + M_b)$  do  
7         $y_i^{(1)} = F(x_i, W)$ ;  
8      end for  
9       $Y^{(1)} = [y_{k+1}^{(1)}, y_{k+2}^{(1)}, \dots, y_{k+M_b}^{(1)}]^T$ ;  
10      $E = Y - Y^{(1)}$ ;  
11      $W_r$  with random initialization in  $[-r, r]$ ;  
12     for  $i = (k + 1)$  to  $(k + M_b)$  do  
13        $y_i^{(2)} = F(x_i, W + W_r)$ ;  
14     end for  
15      $Y^{(2)} = [y_{k+1}^{(2)}, y_{k+2}^{(2)}, \dots, y_{k+M_b}^{(2)}]^T$ ;  
16      $C_X \approx C'_X = (Y^{(2)} - Y^{(1)})W_r^\dagger$ ;  
17      $\Delta W = C_X^\dagger E$ ;  
18      $W = W + l_r \Delta W$ ;  
19     for  $i = (k + 1)$  to  $(k + M_b)$  do  
20        $y_i^{(1)} = F(x_i, W)$ ;  
21     end for  
22      $Y^{(1)} = [y_{k+1}^{(1)}, y_{k+2}^{(1)}, \dots, y_{k+M_b}^{(1)}]^T$ ;  
23      $E = Y - Y^{(1)}$ ;  
24     end for  
25     iteration  $++$ ;  
26 end while
```

Output: W .

Algorithm S5

Algorithm S5 Difference method with batch process of training data and weights.

Input: Training data $\{(x_i, y_i)\}$, $i = 1, \dots, N_d$; r , l_r , M_b and N_b ;
System structure and forward mapping $y = F(x, W)$.

```
1   $W$  with random initialization in  $[-1, 1]$ ;  
2   $Y = [y_1, y_2, \dots, y_{N_d}]^T$ ;  
3  while (iteration  $< Max$  and  $MSE(E) > \varepsilon$ )  
4    for  $k = 0$  to  $(N_d/M_b - 1)$  do  
5       $Y = [y_{k+1}, y_{k+2}, \dots, y_{k+M_b}]^T$ ;  
6      for  $i = (k + 1)$  to  $(k + M_b)$  do  
7         $y_i^{(1)} = F(x_i, W)$ ;  
8      end for  
9       $Y^{(1)} = [y_{k+1}^{(1)}, y_{k+2}^{(1)}, \dots, y_{k+M_b}^{(1)}]^T$ ;  
10      $E = Y - Y^{(1)}$ ;  
11     for  $p = 0$  to  $(S/N_b - 1)$  do  
12        $W_b = [W_{p+1}, W_{p+2}, \dots, W_{p+N_b}]$ ;  
13        $W_r$  with random initialization in  $[-r, r]$ ;  
14       for  $i = (k + 1)$  to  $(k + M_b)$  do  
15          $y_i^{(2)} = F(x_i, W_b + W_r)$ ;  
16       end for  
17        $Y^{(2)} = [y_{k+1}^{(2)}, y_{k+2}^{(2)}, \dots, y_{k+M_b}^{(2)}]^T$ ;  
18        $C_p \approx C'_{X, W_t} = (Y^{(2)} - Y^{(1)})W_r^\dagger$ ;  
19     end for  
20      $C_x^\dagger = [C_1^\dagger; C_2^\dagger; \dots; C_{S/N_b-1}^\dagger]$ ;  
21      $\Delta W = C_x^\dagger E$ ;  
22      $W = W + l_r \Delta W$ ;  
23     for  $i = (k + 1)$  to  $(k + M_b)$  do  
24        $y_i^{(1)} = F(x_i, W)$ ;  
25     end for
```

26 $Y^{(1)} = [y_{k+1}^{(1)}, y_{k+2}^{(1)}, \dots, y_{k+M_b}^{(1)}]^T;$

27 $E = Y - Y^{(1)};$

28 **end for**

29 iteration ++;

30 **end while**

Output: W .

Algorithm S6

Algorithm S6 Difference method with batch process of training data and weights and reuse of C_x^\dagger .

Input: Training data $\{(x_i, y_i)\}$, $i = 1, \dots, N_d$; r , l_r , M_b and N_b ;
System structure and forward mapping $y = F(x, W)$.

```

1    $W$  with random initialization in  $[-1, 1]$ ;
2    $Y = [y_1, y_2, \dots, y_{N_d}]^T$ ;
3   for  $k = 0$  to  $(N_d/M_b - 1)$  do
4        $Y = [y_{k+1}, y_{k+2}, \dots, y_{k+M_b}]^T$ ;
5       for  $i = (k + 1)$  to  $(k + M_b)$  do
6            $y_i^{(1)} = F(x_i, W)$ ;
7       end for
8        $Y^{(1)} = [y_{k+1}^{(1)}, y_{k+2}^{(1)}, \dots, y_{k+M_b}^{(1)}]^T$ ;
9        $E = Y - Y^{(1)}$ ;
10      while (iteration  $< Max$  and  $MSE(E) > \varepsilon$ )
11           $W$  random reordering
12          for  $p = 0$  to  $(S/N_b - 1)$  do
13               $W_b = [W_{p+1}, W_{p+2}, \dots, W_{p+N_b}]$ ;
14               $W_r$  with random initialization in  $[-r, r]$ ;
15              for  $i = (k + 1)$  to  $(k + M_b)$  do
16                   $y_i^{(2)} = F(x_i, W_b + W_r)$ ;
17              end for
18               $Y^{(2)} = [y_{k+1}^{(2)}, y_{k+2}^{(2)}, \dots, y_{k+M_b}^{(2)}]^T$ ;
19               $C_p \approx C'_{x, W_t} = (Y^{(2)} - Y^{(1)})W_r^\dagger$ ;
20          end for
21           $C_x^\dagger = [C_1^\dagger; C_2^\dagger; \dots; C_{S/N_b-1}^\dagger]$ ;
22          while  $MSE(E)$  decreasing
23               $\Delta W = C_x^\dagger E$ ;
24               $W = W + l_r \Delta W$ ;
25          for  $i = (k + 1)$  to  $(k + M_b)$  do

```

```

26          $y_i^{(1)} = F(x_i, W);$ 
27     end for
28      $Y^{(1)} = [y_{k+1}^{(1)}, y_{k+2}^{(1)}, \dots, y_{k+M_b}^{(1)}]^T;$ 
29      $E = Y - Y^{(1)};$ 
30 end while
31 end for
32     iteration  ++;
33 end while

```

Output: W .

Algorithm S7

Algorithm S7 Improved LIFT without matrix.

Input: Training data $\{(x_i, y_i)\}$, $i = 1, \dots, N_d$; r and l_r ;
System structure and forward mapping $y = F(x, W)$.

```
1   $W$  with random initialization in  $[-1, 1]$ ;  
2   $Y = [y_1, y_2, \dots, y_{N_d}]^T$ ;  
3  for  $i = 1$  to  $N_d$  do  
4       $y_i^{(1)} = F(x_i, W)$ ;  
5  end for  
6   $Y^{(1)} = [y_1^{(1)}, y_2^{(1)}, \dots, y_{N_d}^{(1)}]^T$ ;  
7   $E = Y - Y^{(1)}$ ;  
8  while ( $epoch < Max$  and  $MSE(E) > \varepsilon$ )  
9       $W_r$  with random initialization in  $[-r, r]$ ;  
10     for  $i = 1$  to  $N_d$  do  
11          $y_i^{(2)} = F(x_i, W + W_r)$ ;  
12     end for  
13      $Y^{(2)} = [y_1^{(2)}, y_2^{(2)}, \dots, y_{N_d}^{(2)}]^T$ ;  
14      $E' = Y^{(2)} - Y^{(1)}$ ;  
15      $k = \frac{\langle E', E \rangle}{\langle E, E \rangle}$   
16      $\Delta W = kW_r$ ;  
17      $W = W + l_r \Delta W$ ;  
18 end while
```

Output: W .

Supplementary Experiments

Experiment S1: IIR

The N th-order IIR(Infinite Impulse Response) Filter has a system function:

$$H(z) = \frac{Y(z)}{X(z)} = \frac{\sum_{k=0}^N b_k z^{-k}}{1 - \sum_{k=1}^N a_k z^{-k}},$$

and

$$H(e^{j\omega}) = H(z) \big|_{z=e^{j\omega}} = |H(e^{j\omega})| e^{j\beta(j\omega)},$$

where $|H(e^{j\omega})|$ is amplitude-frequency characteristic, and $\beta(j\omega)$ is phase-frequency characteristic.

For this IIR filter system, input is angular frequency $\omega = k2\pi/1000$, $k = 0, 1, \dots, 500$, corresponding output is the amplitude-frequency characteristic $|H(e^{j\omega})|$, that is

$$|H(e^{j\omega})| = f(\omega) = \left| \sum_{k=0}^N b_k z^{-k} / (1 - \sum_{k=1}^N a_k z^{-k}) \right| \bigg|_{z=e^{j\omega}} = \left| \sum_{k=0}^N b_k e^{-jk\omega} / (1 - \sum_{k=1}^N a_k e^{-jk\omega}) \right|$$

And the system parameter $W = [a_1, a_2, \dots, a_N, b_0, b_1, \dots, b_N]$. It is a single-input-single-output (SISO) system. The IIR filter system is designed by the solution of $\{a_k\}$ and $\{b_k\}$.

The difference equation corresponding to the IIR system is

$$y(n) = \sum_{k=0}^N b_k x(n-k) + \sum_{k=1}^N a_k y(n-k)$$

It means that there is feedback between the input and output of the IIR system.

Experiment 1: Design a 100-order IIR filter so that its amplitude-frequency characteristics meet specific requirements (WL-type filter).

Specific model parameters are given in Table S1.

Experiment S2: DNN

Neural networks can be used for classification and regression prediction and are a fundamental element of modern AI solutions. Classical Multi-Layer Feedforward Networks (MLFNs) consist of an input layer, multiple hidden layers, and an output layer. Suppose the classical structure is P-

L - Q , and the specific parameters are shown in Table S2. The neural network adopts a fully connected structure. The feedforward neural networks perform forward computation by iterating layer by layer:

$$\begin{cases} \mathbf{z}^{(l)} = \mathbf{W}^{(l)} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)} \\ \mathbf{a}^{(l)} = \phi(\mathbf{z}^{(l)}) \end{cases}$$

Experiment 2: Design a fully-connected DNN with **40 hidden layers**, to fit a two-input-two-output nonlinear system (two mixed benchmark functions: Ackley and Adjiman).

Specific model parameters are given in Table S2.

Experiment S3: SaNN

To verify the adaptability of nonclassical nonlinear components in ANN, this section selects the **MIMO Sampling Neural Network (SaNN)**. MIMO SaNN uses single-input-single-output (SISO) Sampling neural network (SaNN) as the neuron activation function and accomplishes neural network learning by neuron activation functions training. The feedforward neural networks perform forward computation by iterating layer by layer:

$$\begin{cases} \mathbf{z}^{(l)} = \mathbf{W}^{(l)} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)} \\ \mathbf{a}^{(l)} = \phi(\mathbf{z}^{(l)}) \end{cases}$$

where the neurons' input and output satisfy the following formula:

$$\begin{cases} t = \frac{N_s}{2} (\text{Softsign}(\mathbf{z}^{(l)}) + 1) = \frac{N_s}{2} \left(\frac{\mathbf{z}^{(l)}}{1 + |\mathbf{z}^{(l)}|} + 1 \right) \\ \mathbf{a}^{(l)} = \sum_{n=1}^{N_s} W_n Sa \left[\frac{\pi}{T} (t - nT) \right] \end{cases}$$

The main differences between SaNN and traditional multi-layer neural networks are i) The connection weights between hidden layers do not need to be trained. The inter-layer link weights were randomly initialized with 0 or ± 1 , and then remained fixed, in Experiment 3. ii) The number of training weights for each layer in SaNN is $M_l \times N_s$, rather than $M_l \times M_l$ in MLFNs. Since usually $N_s < M_l$, it is beneficial to reduce the network size. iii) The activation function of each neuron is completely independent and trainable, while the traditional network uses a pre-set and fixed activation function, like sigmoid, tanh, ReLU, ELU, and so on.

Experiment 3: Design a SaNN with **10 hidden layers**, to fit a two-input-two-output nonlinear system (two mixed benchmark functions: Ackley and Adjiman).

Specific model parameters are given in Table S3.

Experiment S4: SaKAN

KAN(**K**olmogorov-**A**rnold **N**etwork) is a novel and hot network architecture based on Kolmogorov-Arnold Representation Theorem, that has emerged in nearly a month. It improves the network performance and interpretability by replacing the weight parameters with learnable univariate functions, and has the potential to become an important direction for driving the development of deep learning models. Its main features are that the learnable activation and B-spline function replacing the traditional linear weights.

As a new type of network architecture, there are still many areas to be improved, such as the difficulty of network training. To further verify the universality of the LIFT algorithm, we constructed a SaKAN (KAN with SaNN) network that replaced the B-spline function with the SaNN. The basic structure of SaKAN is as follows:

$$\left\{ \begin{array}{l} W_{i,j}^{(l)} = \varphi_{SaNN}(x_i^{(l)}) \\ y_j^{(l)} = \sum_i W_{i,j}^{(l)} \\ x_j^{(l)} = Softsign(y_j^{(l+1)}) = \frac{y_j^{(l+1)}}{|y_j^{(l+1)}| + 1} \end{array} \right.$$

where the function φ_{SaNN} is the learnable input-output relationship of SaNN, which weights are the system parameters to be sought.

Experiment 4: Design a SaKAN with **1** hidden layers, to fit a two-input-two-output nonlinear system (two mixed benchmark functions: Ackley and Adjiman).

Specific model parameters are given in Table S4.

Experiment S5: RAN

RAN: Neural Network with Randomly selected Activation functions.

Here employed is still the MLP(Multilayer Perceptron) or MLFNs (Multi-Layer Feedforward Networks), classical neural network structure. The neural networks perform forward computation by iterating layer by layer:

$$\begin{cases} \mathbf{z}^{(l)} = \mathbf{W}^{(l)} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)} \\ \mathbf{a}^{(l)} = \phi(\mathbf{z}^{(l)}) \end{cases}$$

Where the activation functions $\phi(\cdot)$ are no longer a fixed function, but the specific functions randomly selected at initialization from 10 functions. Details in Table S6. These functions include traditional activation functions, non-continuous segmented functions, and non-derivative functions.

Experiment 5: Design a RAN with **3** hidden layers, to fit a two-input-two-output nonlinear system (two mixed benchmark functions: Ackley and Adjiman).

Specific model parameters are given in Table S5.

Experiment S6: MLP

MLP(Multilayer Perceptron) or MLFNs (Multi-Layer Feedforward Networks) is a classical neural network structure. DNN in Experiment S2 is also a MLP with more hidden layers.

In order to compare LIFT with BP, a three-layer network structure is selected. Because basic BP is not suitable for deep networks.

The neural networks perform forward computation by iterating layer by layer:

$$\begin{cases} \mathbf{z}^{(l)} = \mathbf{W}^{(l)} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)} \\ \mathbf{a}^{(l)} = \phi(\mathbf{z}^{(l)}) \end{cases}$$

Experiment 6: Design a MLP with **3** hidden layers, to fit a two-input-two-output nonlinear system (two mixed benchmark functions: Ackley and Adjiman).

Specific model parameters are given in Table S6.

Supplementary Tables

Table S1.

Some parameters of **Experiment 1** for IIR design.

Items	Descriptions or Parameters
Input P	1
Output Q	1
Weights S	201 where 100 for $\{a_n\}$, $n=1, \dots, 100$, and 101 for $\{b_n\}$, $n=0, 1, \dots, 100$.
Object	WL-type filter, the object as shown in Fig. 2, S1 and Data S1.
Input	$\omega \in (0, 2\pi)$
Output	$ H(e^{i\omega}) $
Nth-order	100
N_d	500
W initialization	random in $[-1, 1]$
Pre-training	Yes, Algorithm S1
Iterative training algorithm	Algorithm S3
Learn rate l	$\begin{cases} 0.1, \text{epoch} < 500000 \\ 0.01, \text{others} \end{cases}$
W_r	Random in $[-0.01, 0.01]$
Epoch Max	100 Million
Training upper limit per Epoch	100

Table S2.

Some parameters of **Experiment 2** of DNN:

Items	Descriptions or Parameters
Input	$\mathbf{X} = [x_1, x_2, \dots, x_P]_{1 \times P}$
Output	$\mathbf{Y} = [y_1, y_2, \dots, y_Q]_{1 \times Q}$
The weights matrix between layer $l-1$ and layer l	$\mathbf{W}^{(l)}$
The bias of neurons in layer l	$\mathbf{b}^{(l)}$
Input dimension, P	2
Output dimension, Q	2
The number of hidden layers, L	40
Weights dimension, S	15680
The number of neurons in layer l , MI	$MI=20$ for all layers
Activation function	‘tanh’ for all neurons: $a_i^{(l)} = \phi_i^{(l)}(z_i^{(l)}) = \frac{e^{z_i^{(l)}} - e^{-z_i^{(l)}}}{e^{z_i^{(l)}} + e^{-z_i^{(l)}}}$
Object function	Benchmark functions: Normalized Ackley and Adjiman $\begin{cases} y_1 = f(x_1, x_2) = e^{-0.2 \sqrt{\frac{x_1^2 + x_2^2}{2}}} + 3 \cos 2x_1 + \sin 2x_2 \\ y_2 = f(x_1, x_2) = \cos(x_1) \sin(x_2) - \frac{x_1}{x_2^2 + 1} \end{cases}$ $x_1 \in [0, 10], x_2 \in [0, 10]$ See Fig. S3 .
Training data set N_d	10000
Test data set N_t	10000
Batch process of weights M_b	98
Batch process of training data N_b	98
W initialization	random in $[-1, 1]$
Pre-training	No
Iterative training algorithm	Algorithm S6
Learning rate lr	0.005
W_r	Random in $[-0.01, 0.01]$
Update batch of training data	Randomly selected (shuffling) every 25 iterations (Batch)
Iteration maximum	10000
Training upper limit per iteration (reuse)	10000

Table S3.

Some parameters of **Experiment 3** of SaNN.

Items	Parameters
Input dimension, P	2
Output dimension, Q	2
The number of hidden layers, L	10
Weights dimension, S	2000
The number of neurons in layer l , M_l	$M_l=20$ for all layers
Activation function	SISO SaNN neurons (details in Experiment S3).
The number of weights of a single SISO SaNN neuron, N_s	10
Object function	<p>Benchmark functions: Normalized Ackley and Adjiman</p> $\begin{cases} y_1 = f(x_1, x_2) = e^{-0.2\sqrt{x_1^2 + x_2^2}} + 3\cos 2x_1 + \sin 2x_2 \\ y_2 = f(x_1, x_2) = \cos(x_1)\sin(x_2) - \frac{x_1}{x_2^2 + 1} \end{cases}$ <p>$x_1 \in [0, 10], x_2 \in [0, 10]$</p> <p>See Fig. S5.</p>
Training data set N_d	10000
Test data set N_t	<p>A total of 10000.</p> <p>In Fig. S4, test 100 data after each iteration.</p> <p>In Fig. S5, test the whole 10000 data.</p>
Batch process of weights M_b	100
Batch process of training data N_b	100
W initialization of SISO SaNN for all neurons	Random in $[-1, 1]$
Inter-layer link weights W initialization	<p>Randomly initialized with 0 or ± 1.</p> <p>As shown in Fig. S4h, a total of 3680, initialized as “0” has 1802(48.97%), “1” has 905(24.59%), and “-1” has 973 (26.44%).</p>
Pre-training	Yes, Algorithm S1 by randomly selected 100 train data.
Iterated operation	Algorithm S6
Learn rate lr	0.01
W_r	Random in $[-0.01, 0.01]$
Update batch of training data	Randomly selected (shuffling) every 25 iterations (Batch)
Iteration maximum	50000
Training upper limit per iteration (reuse)	10000

Table S4.

Some parameters of **Experiment 4** of SaKAN.

Items	Parameters
Input dimension, P	2
Output dimension, Q	2
The number of hidden layers, L	1
Weights dimension, S	400
The number of neurons in layer l , M_l	$M_l=10$ for all layers
Activation function	SISO SaNN neurons (details in Experiment S4).
The number of weights of a single SISO SaNN neuron, N_s	10
Object function	Benchmark functions: Normalized Ackley and Adjiman $\begin{cases} y_1 = f(x_1, x_2) = e^{-0.2\sqrt{x_1^2 + x_2^2}} + 3\cos 2x_1 + \sin 2x_2 \\ y_2 = f(x_1, x_2) = \cos(x_1)\sin(x_2) - \frac{x_1}{x_2^2 + 1} \end{cases}$ $x_1 \in [0, 10], x_2 \in [0, 10]$
Training data set N_d	10000
Test data set N_t	A total of 10000.
Batch process of weights M_b	400 (No batch process)
Batch process of training data N_b	1000
W initialization of SISO SaNN for all neurons	Zeros
Pre-training	None
Iterated operation	Algorithm S7
Learn rate lr	1
W_r	Random in $[-0.01, 0.01]$
Update batch of training data	Randomly selected (shuffling) every 25 iterations (Batch)
Iteration maximum	50000
Training upper limit per iteration (reuse)	1

Table S5.

Some parameters of **Experiment 5** of RAN.

Items	Parameters		
Input dimension, P	2		
Output dimension, Q	2		
The number of hidden layers, L	3		
Weights dimension, S	5200		
The number of neurons in layer l , M_l	$M_l=50$ for all layers		
Activation functions	Randomly selected in 10 functions.		
	No.	Functions	Description
	1	$\phi(x) = \frac{1}{1 + e^{-x}}$	sigmoid
	2	$\phi(x) = \tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$	tanh
	3	$\phi(x) = \tanh(x/2) = \frac{e^x - 1}{e^x + 1}$	tanh (x/2)
	4	$\phi(x) = \begin{cases} x^2, & -1 < x < 1 \\ x, & \text{others} \end{cases}$	x^2 (Discontinuity)
	5	$\phi(x) = \ln(1 + x)$	ln (Non derivable)
	6	$\phi(x) = \max(0, x) = \begin{cases} 0, & x < 0 \\ x, & x \geq 0 \end{cases}$	ReLU (Non derivable)
	7	$\phi(x) = \sin(0.4x)$	sin
	8	$\phi(x) = \text{atan}(x)$	arctan
	9	$\phi(x) = \frac{x}{1 + x }$	softsign
	10	$\phi(x) = \begin{cases} \frac{x}{1+x}, & x > 4 \\ 0.2x, & -4 \leq x \leq 4 \\ \frac{x}{1-x}, & x < -4 \end{cases}$	Piecewise function (Non derivable)
Object function	Benchmark functions: Normalized Ackley and Adjiman $\begin{cases} y_1 = f(x_1, x_2) = e^{-0.2\sqrt{\frac{x_1^2 + x_2^2}{2}}} + 3\cos 2x_1 + \sin 2x_2 \\ y_2 = f(x_1, x_2) = \cos(x_1)\sin(x_2) - \frac{x_1}{x_2^2 + 1} \end{cases}$ $x_1 \in [0, 10], x_2 \in [0, 10]$		
Training data set N_d	10000		
Test data set N_t	A total of 10000.		

Batch process of weights M_b	5200 (No batch process)
Batch process of training data N_b	100
W initialization	Random
Pre-training	None
Iterated operation	Algorithm S7
Learn rate lr	1
W_r	Random in [-0.001, 0.001]
Update batch of training data	Randomly selected (shuffling) every 25 iterations (Batch)
Iteration maximum	2,000,000
Training upper limit per iteration (reuse)	1

Table S6.

Some parameters of **Experiment 6** of MLP.

Items	Parameters
Input dimension, P	2
Output dimension, Q	2
The number of hidden layers, L	3
Weights dimension, S	5200
The number of neurons in layer l , M_l	$M_l=50$ for all layers
Activation functions	tanh functions: $\phi(x) = \tanh(x) = \frac{e^{2x}-1}{e^{2x}+1}$
Object function	Benchmark functions: Normalized Ackley and Adjiman $\begin{cases} y_1 = f(x_1, x_2) = e^{-0.2} \sqrt{x_1^2 + x_2^2} + 3 \cos 2x_1 + \sin 2x_2 \\ y_2 = f(x_1, x_2) = \cos(x_1) \sin(x_2) - \frac{x_1}{x_2^2 + 1} \end{cases}$ $x_1 \in [0, 10], x_2 \in [0, 10]$
Training data set N_d	10000
Test data set N_t	A total of 10000.
Batch process of weights M_b	5200 (No batch process)
Batch process of training data N_b	100
W initialization	Random
Pre-training	None
Iterated operation	BP Algorithm S7
Learn rate lr	1 for S7 0.001 for BP(The learning rate is too large, the network is not easy to converge.)
W_r	Random in [-0.001, 0.001]
Update batch of training data	Randomly selected (shuffling) every 25 iterations (Batch)
Iteration maximum	2,000,000
Training upper limit per iteration (reuse)	1

Supplementary Figures

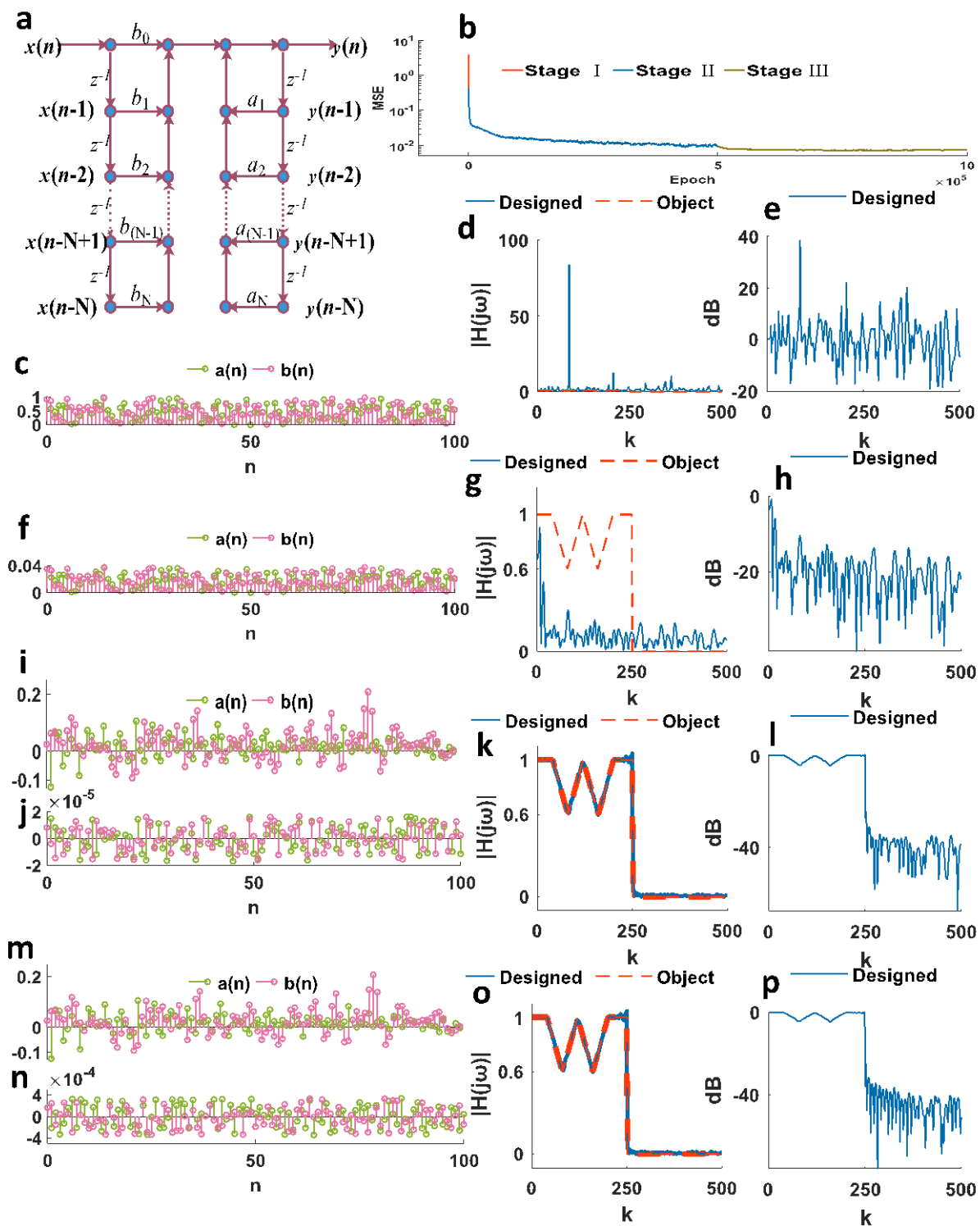


Fig. S1. The results of Experiment 1 on IIR design. Stage I is pre-training in Algorithm S1; Stage II is the iterative training in epoch 500,000, learn rate 0.1, in Algorithm S3; Stage III is the iterative training in epoch 1000,000, learn rate 0.01, in Algorithm S3. **a**, The system structure of classical direct-type N th-order IIR Filter, input $x(n)$ and output $y(n)$. **b**, The convergence results of Experiment 1 on IIR design. **c**, Randomly initialized system parameters $\{a_k\}$ and $\{b_k\}$. **d**, Amplitude-frequency characteristic curve $|H(e^{j\omega})|$ of randomly initialized system. **e**, Bode plot for the magnitude of the frequency response ($20\lg|H(e^{j\omega})|$) of randomly initialized system. **f**, System parameters after pre-training (Stage I). **g**, $|H(e^{j\omega})|$ after Stage I. **h**, $20\lg|H(e^{j\omega})|$ after Stage I. **i**, System parameters after Stage II. **j**, Adjusted value of system parameters in epoch 500,000. **k**, $|H(e^{j\omega})|$ after Stage II. **l**, $20\lg|H(e^{j\omega})|$ after Stage II. **m**, System parameters after Stage III. **n**, the Adjusted value of system parameters in epoch 1000,000. **o**, $|H(e^{j\omega})|$ after Stage III. **p**, $20\lg|H(e^{j\omega})|$ after Stage III.

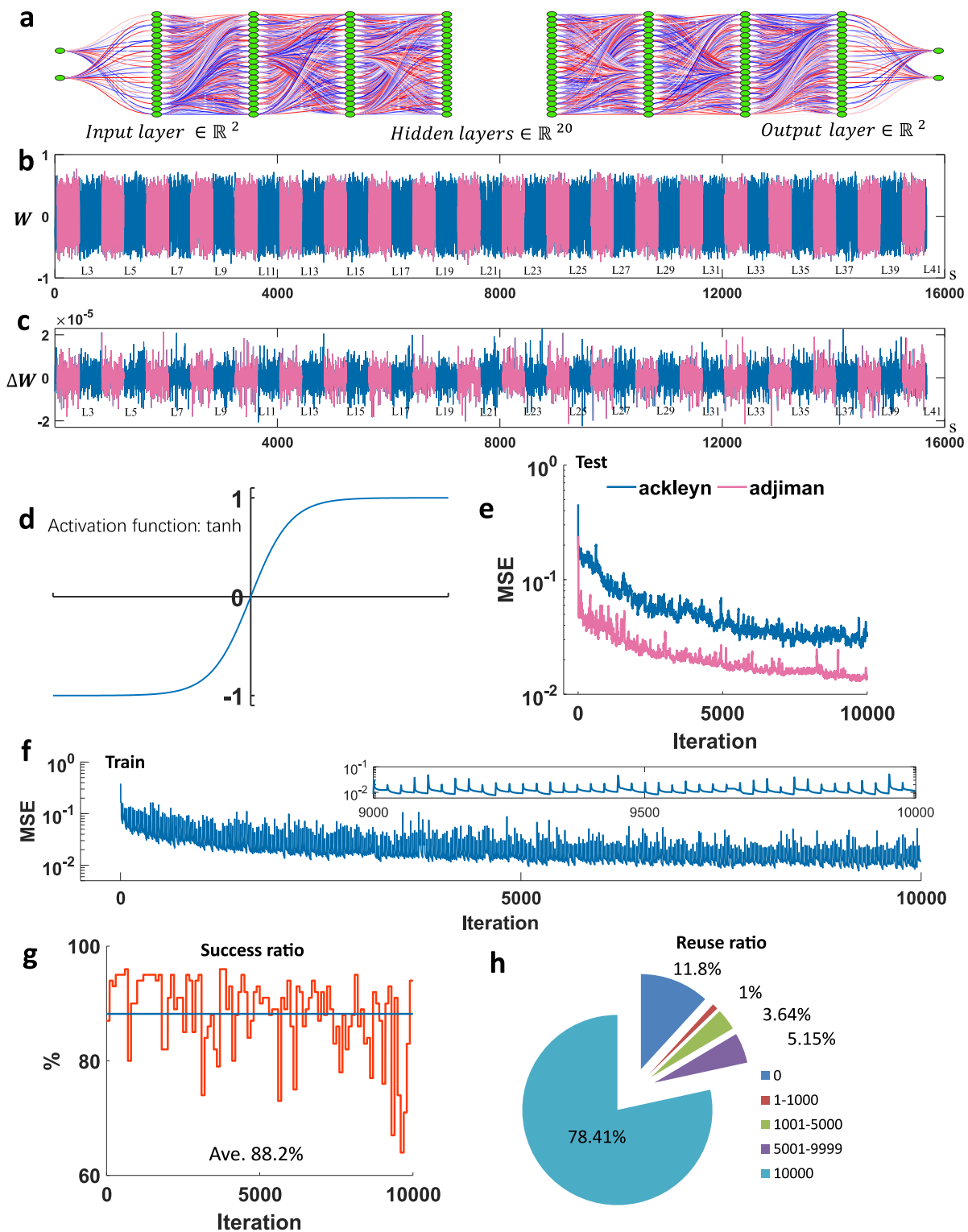


Fig. S2. The results of Experiment 2 on DNN design, in 10,000 iterations, learning rate 0.005, and Algorithm S6. a, The structure of DNN. **b,** The network weighs in 10,000 iterations. **c,** The adjusted value of network weights in 10,000 iterations. **d,** The activation function of neurons is tanh. **e,** Two benchmark function output convergence in MSE of the test. **f,** Training MSE. MSE continuously reduced for each update training data batch per 25 iterations. **g,** The success ratio with iteration. The ratio is calculated once per 100 iterations. The whole average success ratio is 88.2% as shown in the blue line. **h,** The reuse situation of \mathbb{C}_X . “0” means the current \mathbb{C}_X is fail to train. “10000 times” is the preset up-limit of reuse.

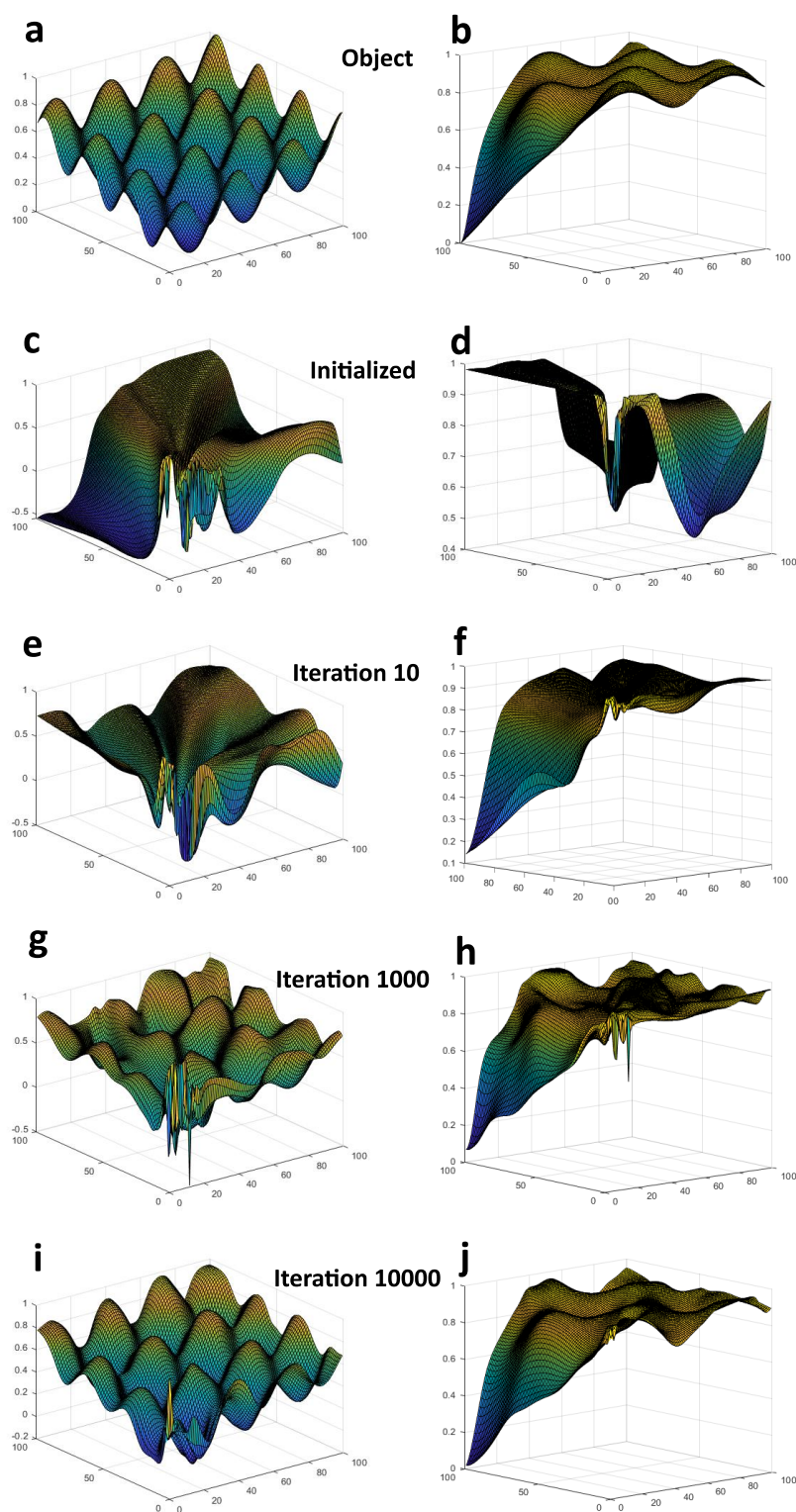


Fig. S3. The results of Experiment 2 on DNN design (The supplement of Fig.S2.). **a**, The benchmark function: Ackley. **b**, The benchmark function: Adjiman. **c**, Initialized network output

of Ackley with the MSE 0.4536. **d**, Initialized network output of Adjiman with the MSE 0.2399. **e**, Network output of Ackley with the MSE 0.1906 in iteration 10. **f**, Network output of Adjiman with the MSE 0.0503 in iteration 10. **g**, Network output of Ackley with the MSE 0.0922 in 1000 iterations. **h**, Network output of Adjiman with the MSE 0.0305 in 1000 iterations. **i**, Network output of Ackley with the MSE 0.0323 in 10000 iterations. **j**, Network output of Adjiman with the MSE 0.0130 in 10000 iterations.

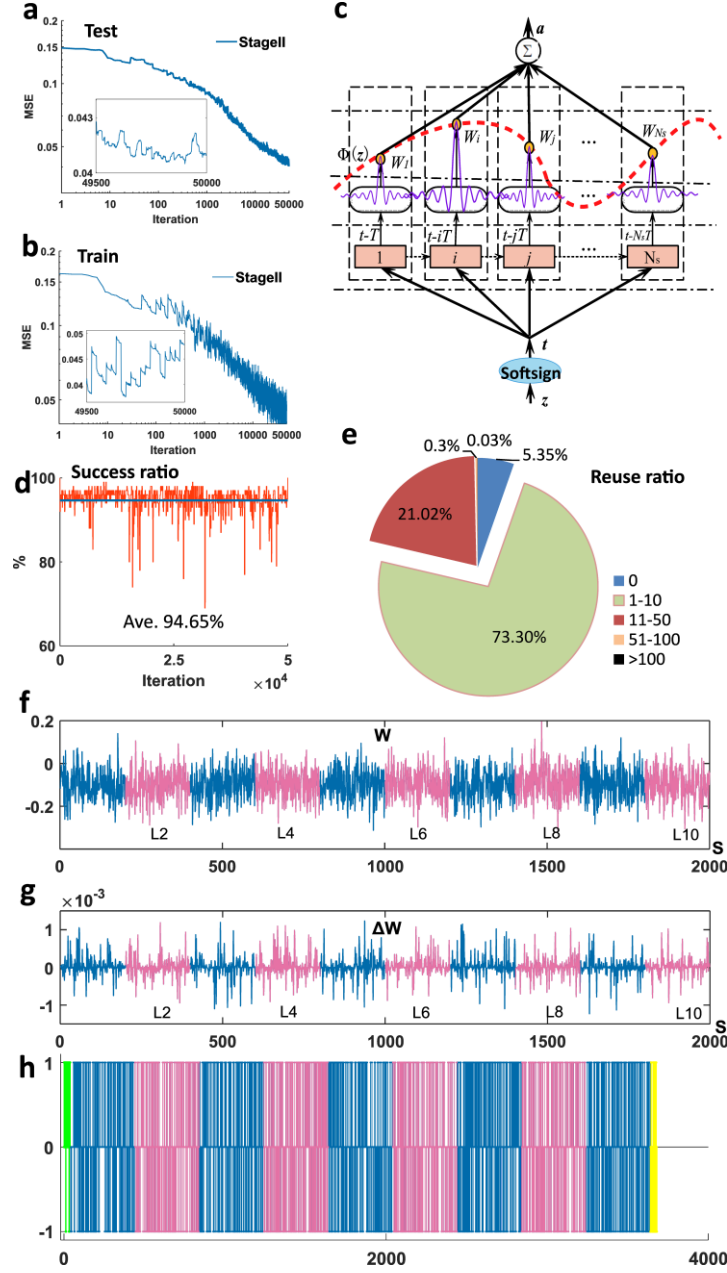


Fig. S4. The results of Experiment 3 on SaNN design, in iteration 50,000, learning rate 0.01, Algorithm S1, and Algorithm S6. a, benchmark function output convergence in MSE of the test. **b,** MSE of training. Update training data batch per 25 iterations. **c,** The structure of SISO SaNN. **d,** The success ratio with iteration. The ratio is calculated once per 100 iterations. The whole average success ratio is 88.2% as shown in the blue line. **e,** The reuse situation of \mathbb{C}_X . “0” means the failure calculation of the current \mathbb{C}_X . “10000 times” is the preset up limit of reuse. **f,** The network weights in iteration 10,000. **g,** The adjusted value of network weights in 10,000 iterations. **h,** Inter-layer link weights between different layers. A total of 3680, initialized as “0” has 1802(48.97%), “1” has 905(24.59%), and “-1” has 973 (26.44%).

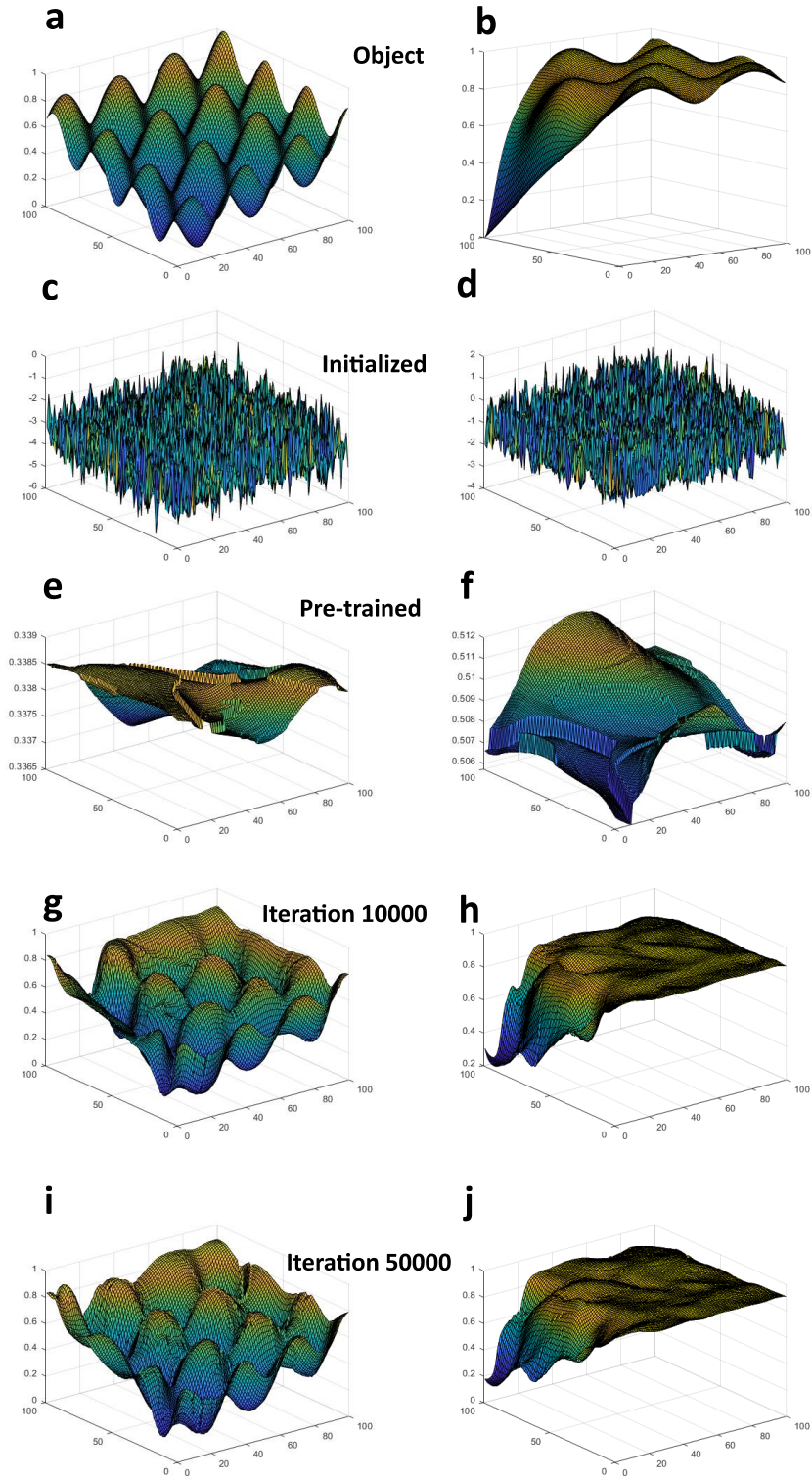


Fig. S5. The results of Experiment 3 on SaNN design (The supplement of Fig.S4.). Test the whole 10000 test dataset, different from the test in Fig.3 just 100 randomly selected test data. **a,**

The benchmark function: Ackley. **b**, The benchmark function: Adjiman. **c**, Initialized network output of Ackley with the MSE 0.7417. **d**, Initialized network output of Adjiman with the MSE 0.7154. **e**, Network output of Ackley with the MSE 0.1844 by pre-train in Algorithm S1. **f**, Network output of Adjiman with the MSE 0.1380 by pre-train in Algorithm S1. **g**, Network output of Ackley with the MSE 0.0696 in 10000 iterations in Algorithm S6. **h**, Network output of Adjiman with the MSE 0.0467 in 10000 iterations in Algorithm S6. **i**, Network output of Ackley with the MSE 0.0499 in 50000 iterations in Algorithm S6. **j**, Network output of Adjiman with the MSE 0.0410 in 50000 iterations in Algorithm S6.

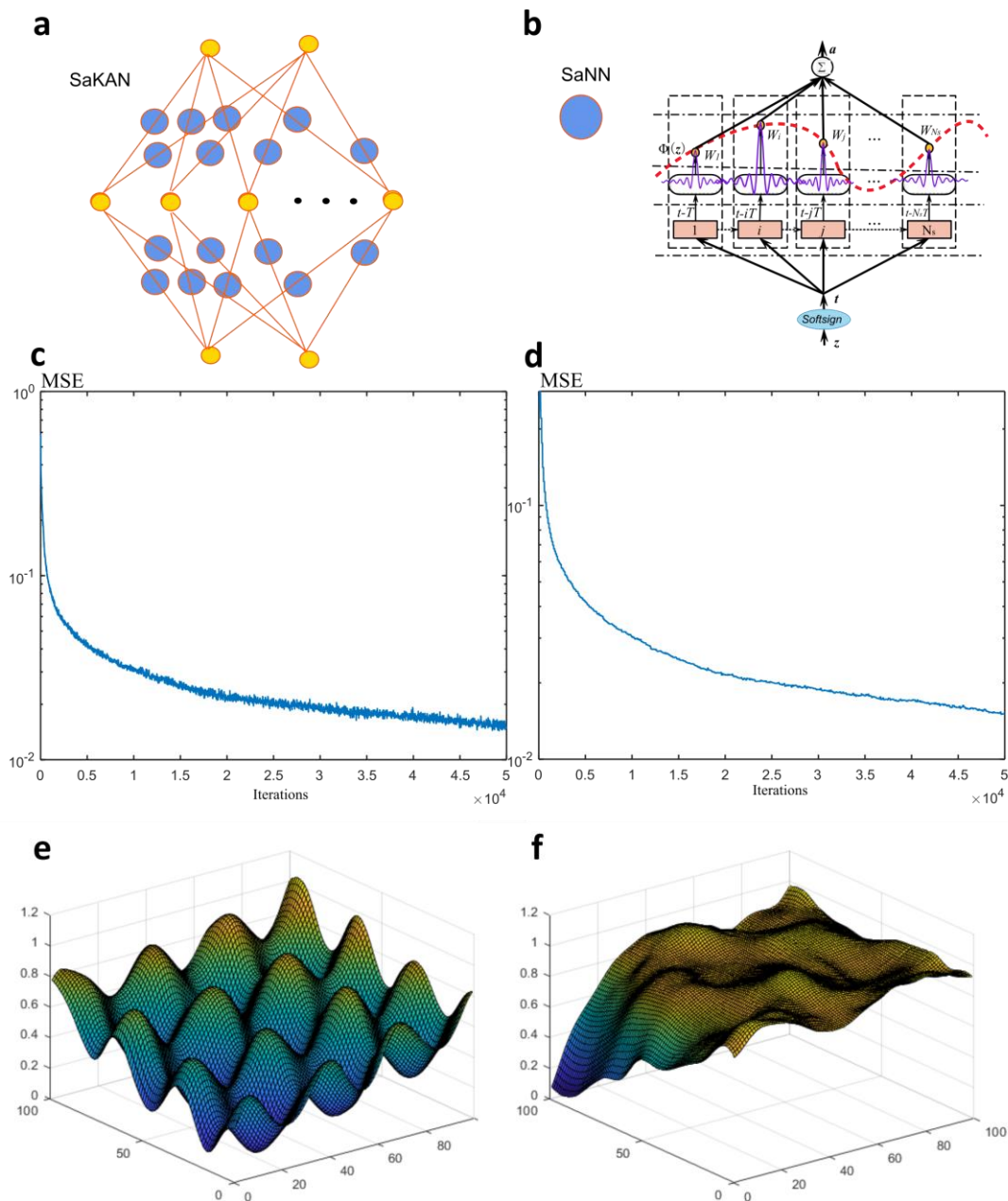


Fig. S6. The results of Experiment 4 on SaKAN design. a, The structure of SaKAN. **b,** The structure of SISO SaNN. **c,** MSE of training. **d,** MSE of testing. **e,** Network output of Ackley with the MSE 0.0155 in 50000 iterations in Algorithm S7. **f,** Network output of Adjiman with the MSE 0.0142 in 50000 iterations in Algorithm S7.

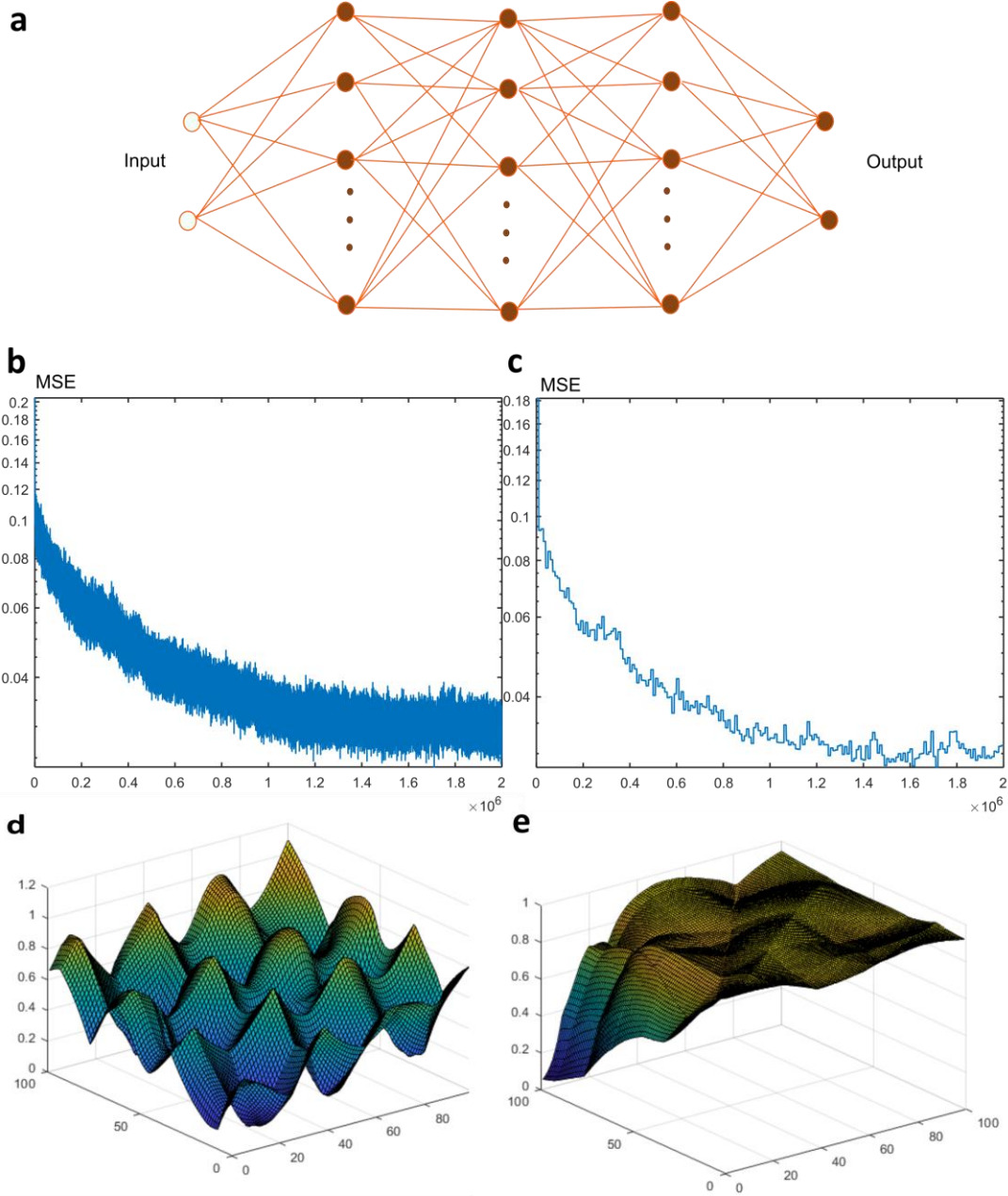


Fig. S7. The results of Experiment 5 on RAN design. a, The structure of RAN. The solid circle represents neurons with randomly selected activation functions. **b,** MSE of training. **c,** MSE of testing. **d,** Network output of Ackley with the MSE 0.0324 in 2,000,000 iterations in Algorithm S7. **e,** Network output of Adjiman with the MSE 0.0265 in 2,000,000 iterations in Algorithm S7.

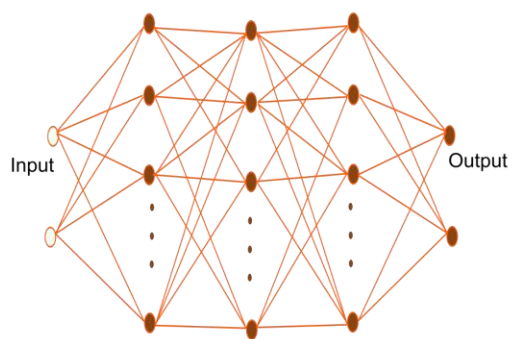
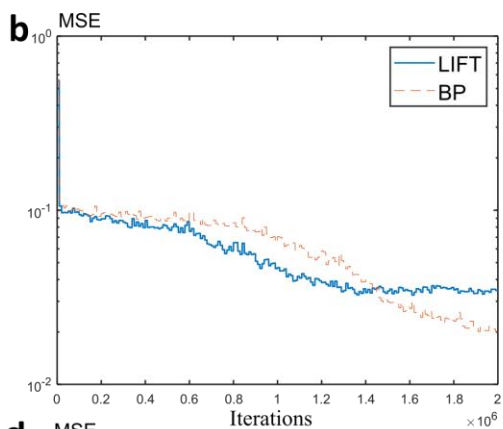
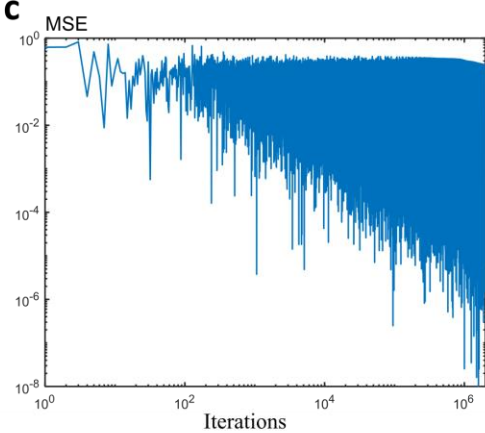
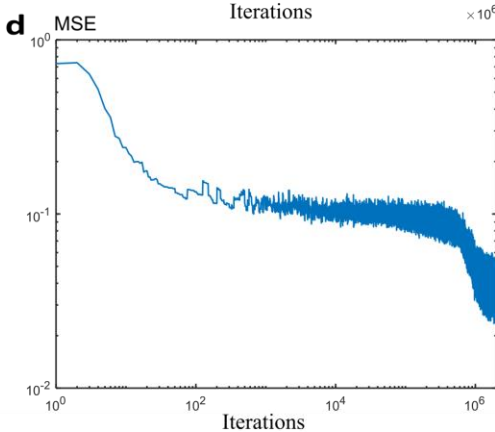
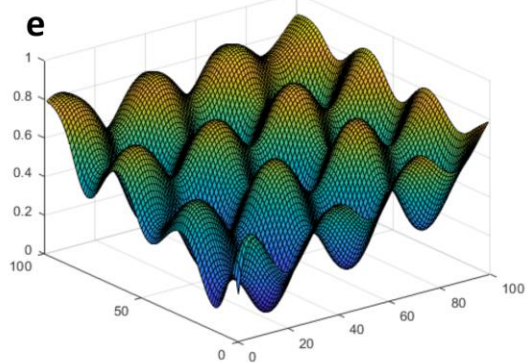
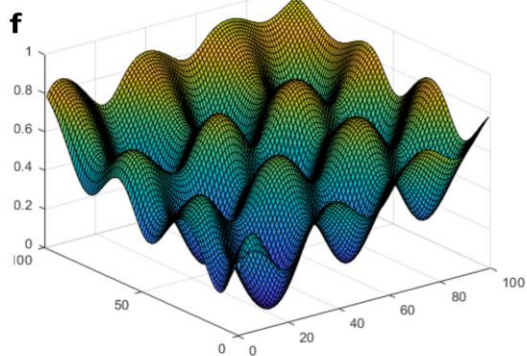
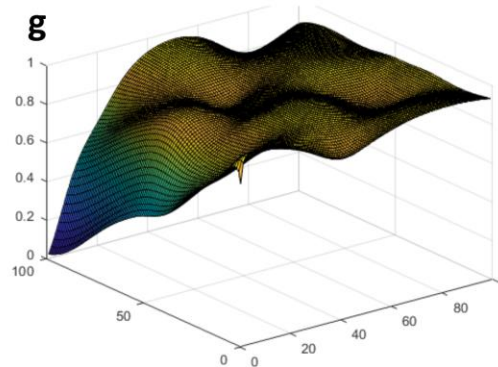
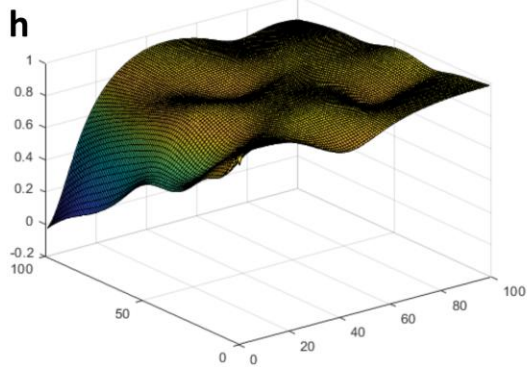
a**b****c****d****e****f****g****h**

Fig. S8. The results of Experiment 6 on MLP design. **a**, The structure of MLP with activation functions tanh . **b**, MSE of training. **c**, MSE of testing in BP. **d**, MSE of testing in LIFT. **e**, Network output of Ackley with the MSE 0.0239 in 2,000,000 iterations in BP. **f**, Network output of Ackley with the MSE 0.0391 in 2,000,000 iterations in LIFT S7. **g**, Network output of Adjiman with the MSE 0.0181 in 2,000,000 iterations in BP. **h**, Network output of Adjiman with the MSE 0.0256 in 2,000,000 iterations in LIFT S7.

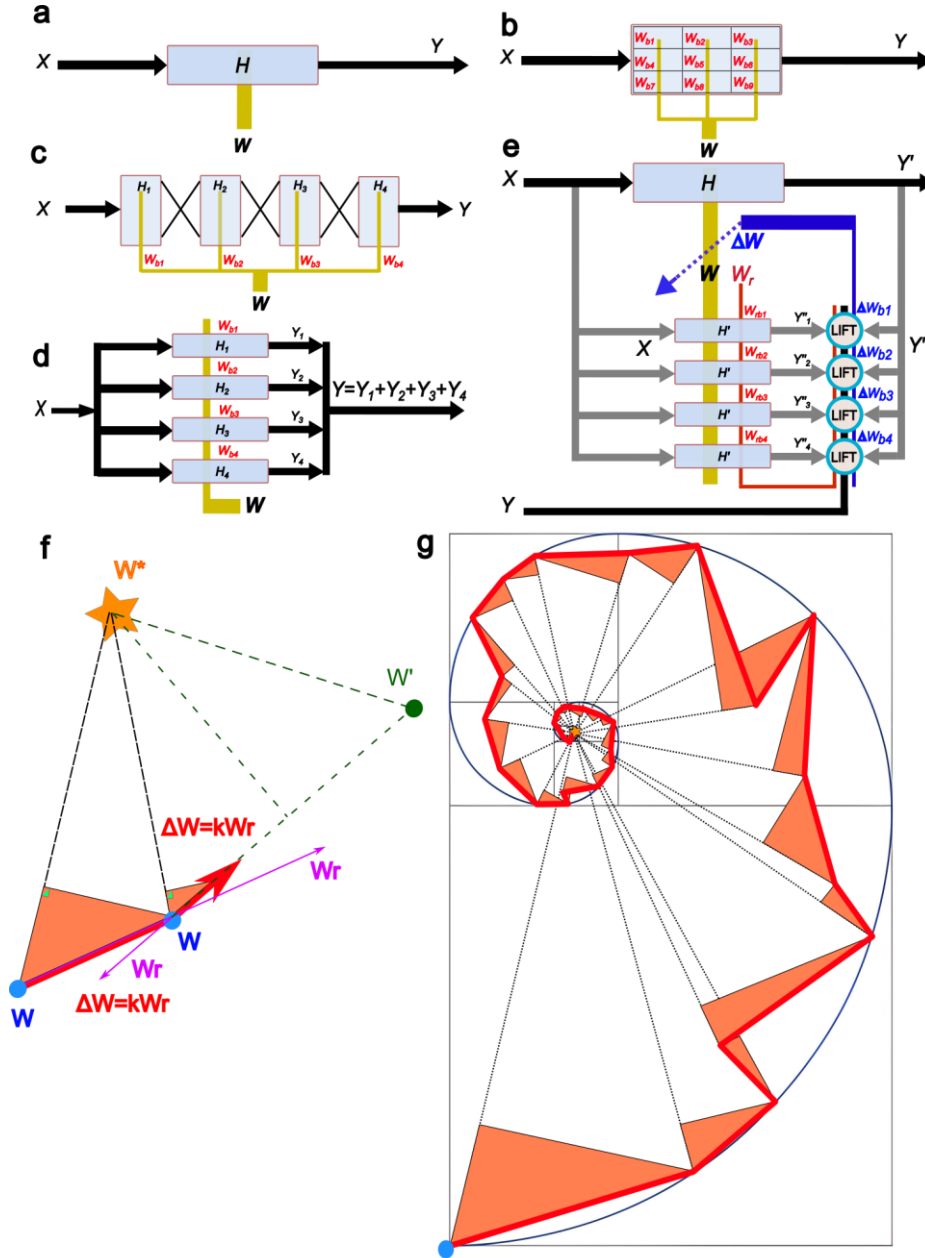


Fig. S9. Diagrammatic sketch. **a**, The structure diagram of parametric system indicated by H . **b**, The sub-region structure with batch processes(W_{b1} to W_{b9}). **c**, The hierarchical structure with batch processes(W_{b1} to W_{b4}). **d**, The parallel structure with batch processes(W_{b1} to W_{b4}). **e**, The parallel training diagram with batch processes(W_{b1} to W_{b4}). H' is the copy(reuse module) of H . See Algorithm S5, S6. **f**, Sketch of LIFT method. W is the current parameters; W^* is the expected parameters; the red solid lines represent the weight-adjusted trajectory ΔW ; the black dashed lines represent the ideal tuning direction.

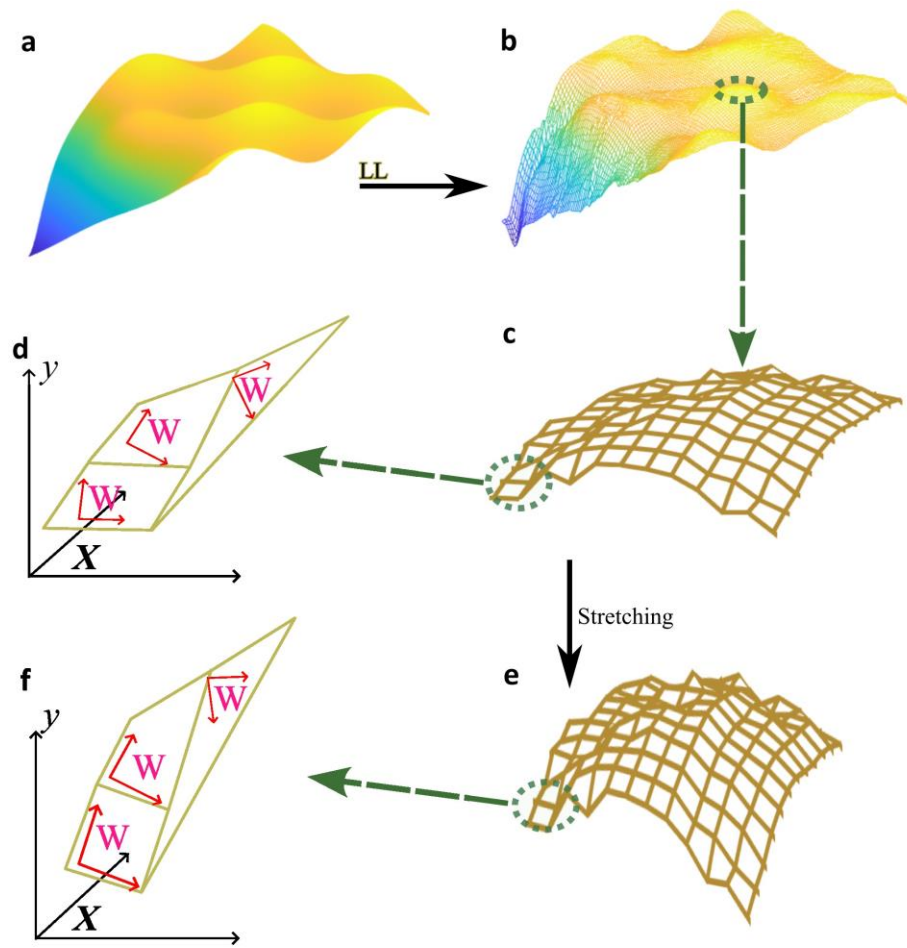


Fig. S10. The diagrammatic sketch of visualized understanding of LL-IC. **a**, System space. **b-d**, System approximation by local linearization. **e-f**, The diagrammatic sketch of isomorphism comparability. System approximation after stretching caused by weight adjustment. In the process of stretching and torsion, the change of gradient is often limited. This is also the reason why \mathbb{C}_X is effective. It's also the reason of its reuse.