

# G-SpNN: GPU-Accelerated Passivity Enforcement for S-Parameter Modeling with Neural Networks

Lijie Zeng<sup>1</sup>, Jiatai Sun<sup>1</sup>, Xiao Wu<sup>2</sup>, Dan Niu<sup>3</sup>, Tianshi Wang<sup>4</sup>, Yibo Lin<sup>5</sup>, Zuochang Ye<sup>6</sup>, Zhou Jin<sup>7</sup>

<sup>1</sup>SSSLab, China University of Petroleum-Beijing, China, <sup>2</sup>Huada Emphyrean Software Co. Ltd, China,

<sup>3</sup>School of Automation, Southeast University, China, <sup>4</sup>HiSilicon Technologies Co. Ltd,

<sup>5</sup>School of Integrated Circuits, Peking University, Beijing, <sup>6</sup>School of Integrated Circuits, Tsinghua University,

<sup>7</sup>College of Integrated Circuits, Zhejiang University, China

Email: [z.jin@zju.edu.cn](mailto:z.jin@zju.edu.cn)



SPONSORED BY



# OUTLINE

- 1 Background
- 2 G-SpNN
- 3 Experiment
- 4 Conclusion

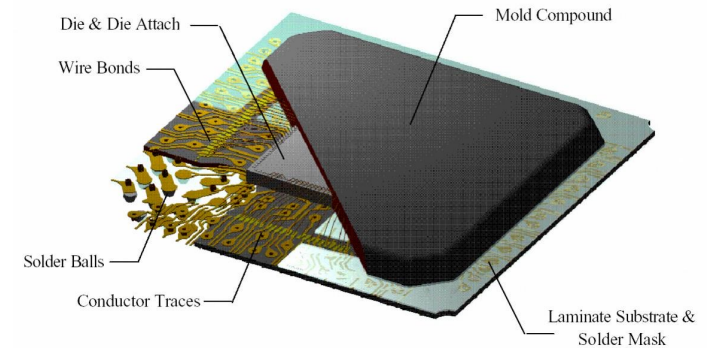
# OUTLINE

- 1 **Background**
- 2 **G-SpNN**
- 3 **Experiment**
- 4 **Conclusion**

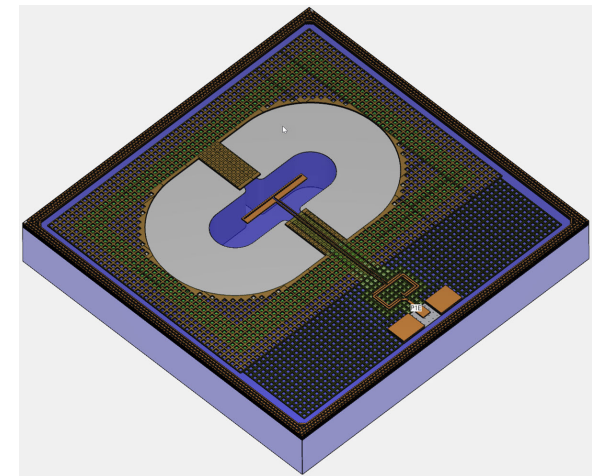


# Background

- At high frequencies, **S-parameters** (scattering parameters) are commonly used to describe the performance of microwave and RF devices.
- Running time-domain analyses with them is computationally intensive and often leads to convergence issues.
- **Macromodeling** techniques are applied to simplify these behaviors, enhancing simulation efficiency and stability.



Chip Packaging Simulation



RFIC Simulation

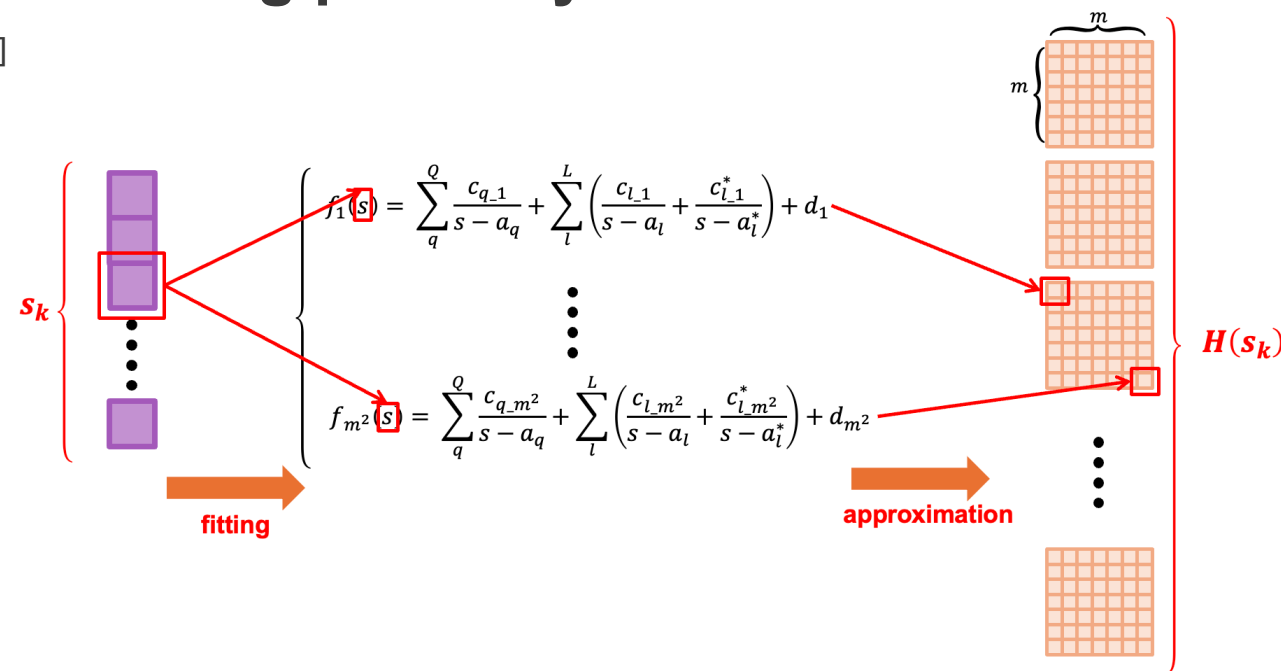
# Background

Mainstream methods typically adopt a **two-step approach**:

① **Generate a macromodel without considering passivity constraints**

Mainstream method such as: Vector Fitting (VF)<sup>[1]</sup>

Using the **VF method**, a set of functions with a given form is fitted to the known frequency data  $s_k$  and system response  $H(s_k) \in \mathbb{C}^{m \times m}$ .



# Background

## ② Applying specialized algorithms to restore passivity

Such as: Eigenvalue Perturbation (EPM)<sup>[2]</sup>, Residue Perturbation (RPM)<sup>[3]</sup>, Local Compensation (LC)<sup>[4]</sup>

Transform the rational function form into a state-space form and apply the **EPM/RPM/LC** to restore the passivity of model  $G(s)$ .

$$\underbrace{f(s_k) = \sum_{n=1}^{N_q} \frac{c_n}{s_k - a_n} + d + s_k h}_{\text{rational function form}} \quad \longrightarrow \quad \underbrace{G(s_k) = C(s_k \cdot I - A)^{-1}B + D}_{\text{state-space form}} \quad \longrightarrow \quad \underbrace{G(s_k) + G^H(s_k) \geq 0}_{\text{passivity condition}}$$

[2] S. Grivet-Talocia, "An adaptive sampling technique for passivity characterization and enforcement of large interconnect macromodels," IEEE Transactions on Advanced Packaging, vol. 30, no. 2, pp. 226–237, 2007.

[3] B. Gustavsen, "Fast passivity enforcement for pole-residue models by perturbation of residue matrix eigenvalues," IEEE Transactions on Power Delivery, vol. 23, no. 4, pp. 2278–2285, 2008.

[4] T. Wang and Z. Ye, "Robust passive macro-model generation with local compensation," IEEE Transactions on Microwave Theory and Techniques, vol. 60, no. 8, pp. 2313–2328, 2012.

# Background

**Two-Step Approach**

**VF + EPM/RPM/LC**

existing issues

**compromise model accuracy**

the most critical issue

address this issue

**Domain-Alternated Optimization (DAO) [5]**

the proposal of a new method



[5] Ye Z, Wang T, Li Y. Domain-alternated optimization for passive macromodeling[J]. IEEE Transactions on Very Large Scale Integration (VLSI) Systems, 2014

# Background

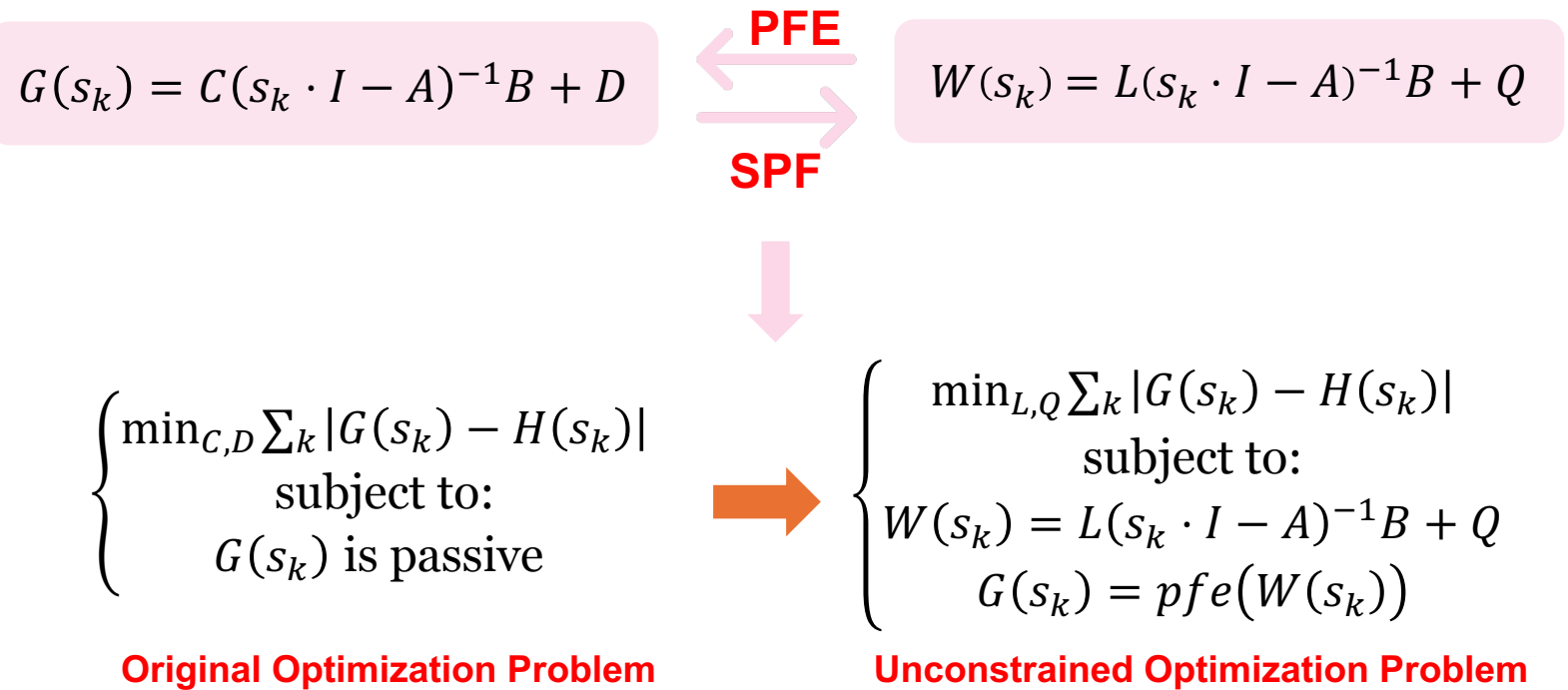
DAO introduces **two key transformation operators**. Based on the transformation using **SPF** and **PFE**, the system  $G(s_k)$  derived from  $W(s_k)$  is guaranteed to **preserve passivity**, and the original optimization problem can be further converted into an **unconstrained optimization problem**.

## Spectral Factorization (SPF)

$$\begin{aligned} R &= Q^T Q = D + D^T \\ \tilde{A} &= BR^{-1}C - A, \tilde{B} = BQ^{-1}, \tilde{C} = C^T R^{-1}C \\ K\tilde{A} + \tilde{A}^T K - K\tilde{B}\tilde{B}^T K - \tilde{C} &= 0 \\ L &= Q^{-T}C - Q^{-T}B^T K \\ W(s_k) &= L(s_k \cdot I - A)^{-1}B + Q \end{aligned}$$

## Partial Fractional Expansion (PFE)

$$\begin{aligned} M &= \text{kron}(A^T, I_n) + \text{kron}(I_n, A^T) \\ \text{vec}(K) &= -M^{-1}\text{vec}(L^T L) \\ C &= B^T K + Q^T L, D = \frac{1}{2} Q^T Q \end{aligned}$$





# Background

Three-Step Approach

VF + EPM/RPM/LC



Domain-Alternated Optimization

The essence of DAO is to transform the passivity-constrained problem into an **unconstrained optimization problem**, thereby improving macromodel accuracy while maintaining passivity throughout the process.

Stage	Method	Passivity	Time (s)	Total Error
1	VF	non-passive	544	1.69%
2	LC	passive	253	8.38%
3	DAO	passive	12238	2.63%

case with 21-port system



# Background

Three-Step Approach

VF + EPM/RPM/LC



Domain-Alternated Optimization

However, as the number of ports in complex integrated circuits continues to grow, the three-step approach increasingly reveals additional issues.

For the previously mentioned 21-port system, the iteration time of DAO is **15 times longer** than the total time of the first two steps. For a 64-port system, each iteration consumes an average of **22GB of memory**, while for a 138-port system, memory usage **exceeds 31.8GB**.

Stage	Method	Passivity	Time (s)	Total Error
1	VF	non-passive	544	1.69%
2	LC	passive	253	8.38%
3	DAO	passive	12238	2.63%

case with 21-port system



# Background

Three-Step Approach

VF + EPM/RPM/LC



Domain-Alternated Optimization

However, as the number of ports in complex integrated circuits continues to grow, the three-step approach faces several challenges.

## the Two Primary Challenges

High Memory Consumption

Slow Convergence

For the previously mentioned 21-port system, the iteration time of DAO is 12238s. For a 138-port system, memory usage exceeds 31.8GB.

Stage	Method	Passivity	Time (s)	Total Error
1	VF	non-passive	544	1.69%
2	LC	passive	253	8.38%
3	DAO	passive	12238	2.63%

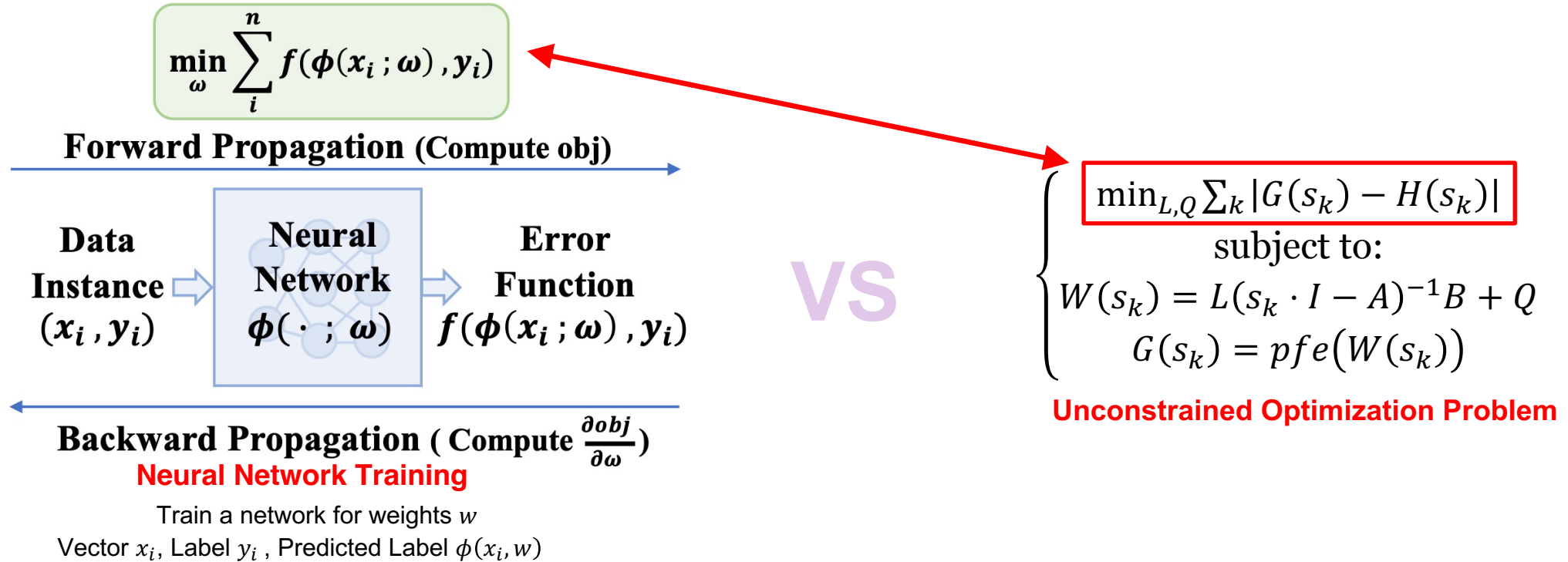
case with 21-port system



# OUTLINE

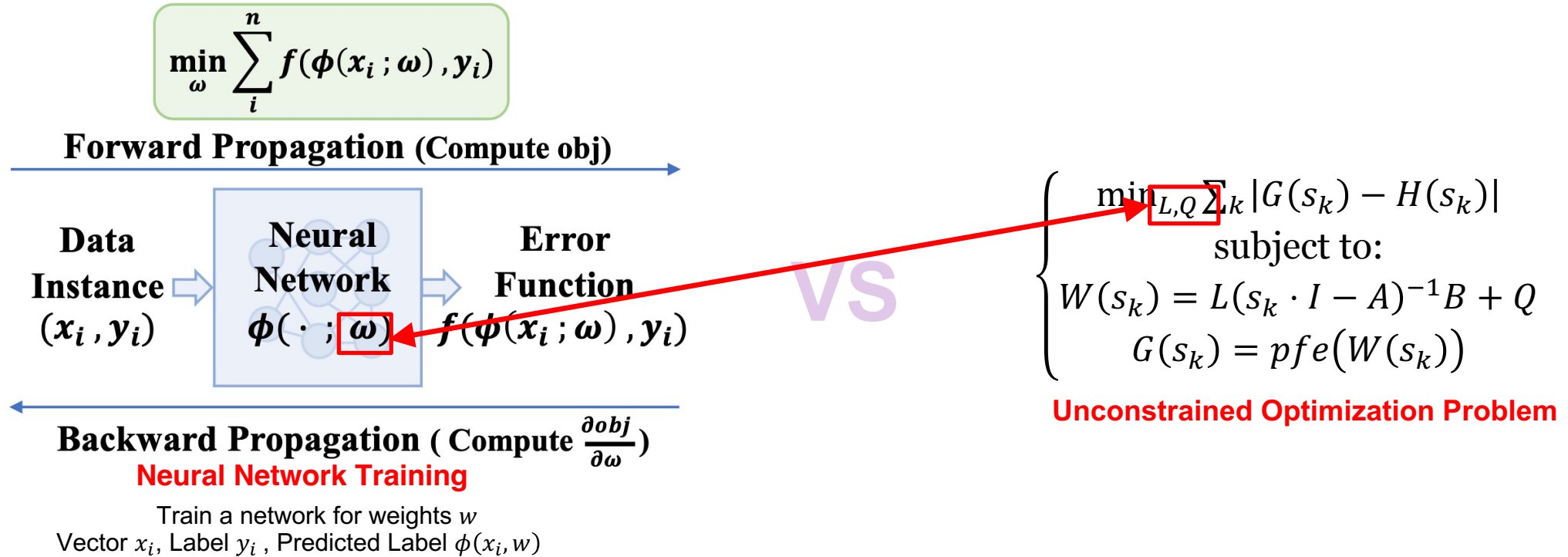
- 1 Background
- 2 **G-SpNN**
- 3 Experiment
- 4 Conclusion



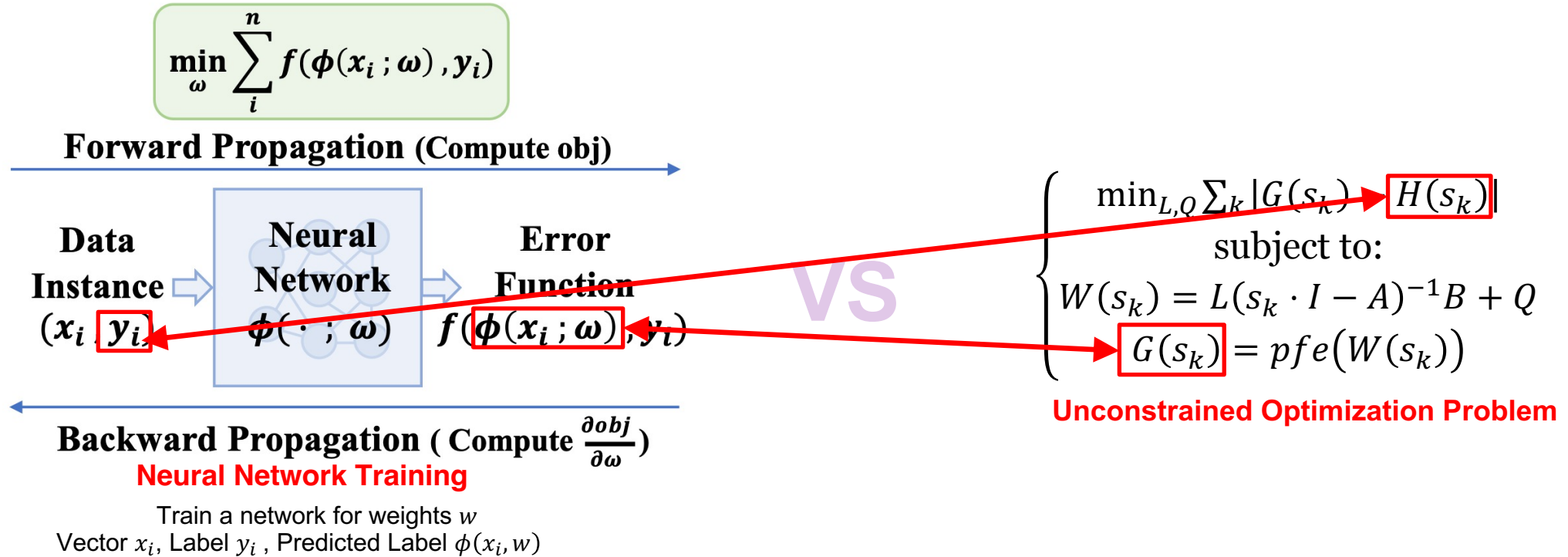


The **error function in unconstrained optimization** can be viewed as analogous to **the prediction error** encountered in neural network training.





The **optimization variables  $L, Q$**  in unconstrained optimization can be viewed as analogous to **the trainable weight** encountered in neural network training.



$H(s_k)$ ,  $G(s_k)$  in unconstrained optimization can be viewed as analogous to the label and predicted label encountered in neural network training.

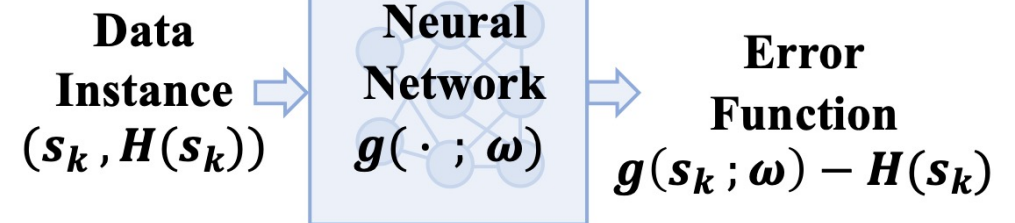
$$\begin{cases} \min_{L,Q} \sum_k |G(s_k) - H(s_k)| \\ \text{subject to:} \\ W(s_k) = L(s_k \cdot I - A)^{-1}B + Q \\ G(s_k) = pfe(W(s_k)) \end{cases}$$

Unconstrained Optimization Problem



$$\min_{\omega} \sum_k^N |g(s_k; \omega) - H(s_k)|$$

Forward Propagation (Compute obj)



Backward Propagation (Compute  $\frac{\partial obj}{\partial \omega}$ )  
Neural Network Training

Solve a set of parameters for macromodel  
with  $w = (L, Q)$

The two optimization frameworks share a **fundamental similarity**, making it possible to **leverage neural network techniques** to accelerate macromodeling optimization.

# G-SpNN | Loss Function

To make the neural network training process more efficient, the **objective function can be reformulated** to reduce computational complexity.

The summation of errors over  $k$  terms increases computational complexity.

$$\text{Error} = \min_{C,D} \sum_k |G(s_k) - H(s_k)|$$

Original Objective Function

$$\begin{aligned} \text{Error\_vec} &= |Fy - h| \\ y &= \begin{bmatrix} \text{vec}(C) \\ \text{vec}(D) \end{bmatrix} \quad F_i = [\text{kron}([s_i I - A]^{-1}, I_n) \quad I_m^2] \\ h &= \begin{bmatrix} \text{Vec}(\text{Re}(H(s_1))) \\ \vdots \\ \text{Vec}(\text{Re}(H(s_N))) \\ \text{Vec}(\text{Im}(H(s_1))) \\ \vdots \\ \text{Vec}(\text{Im}(H(s_N))) \end{bmatrix} \quad F = \begin{bmatrix} \text{Re}(F_1) \\ \vdots \\ \text{Re}(F_N) \\ \text{Im}(F_1) \\ \vdots \\ \text{Im}(F_N) \end{bmatrix} \\ F &= Q_F R_F \\ b &= Q_F^T h, \quad \delta^2 = h^T h - b^T b \end{aligned}$$

Intermediate Computation Steps

$$\text{loss} = \text{Error} = (R_F y - b)^T (R_F y - b) + \delta^2$$

Reformulated Objective Function



# G-SpNN | Loss Function

Based on the new definition of the loss function, optimization problem can be further reformulated.

$$\begin{cases} \min_{C,D} \sum_k |G(s_k) - H(s_k)| \\ \text{subject to:} \\ G(s_k) \text{ is passive} \end{cases}$$

Original Optimization Problem



$$\begin{cases} \min_{L,Q} \sum_k |G(s_k) - H(s_k)| \\ \text{subject to:} \\ W(s_k) = L(s_k \cdot I - A)^{-1}B + Q \\ G(s_k) = pfe(W(s_k)) \end{cases}$$

Unconstrained Optimization Problem



$$\begin{cases} \text{variable: } L, Q \\ \min: f(y) = (R_F y - b)^T (R_F y - b) + \delta^2 \\ \text{subject to:} \\ x = \begin{bmatrix} \text{vec}(L) \\ \text{vec}(Q) \end{bmatrix} \\ y = pfe(x) \end{cases}$$

Reformulated Optimization Problem

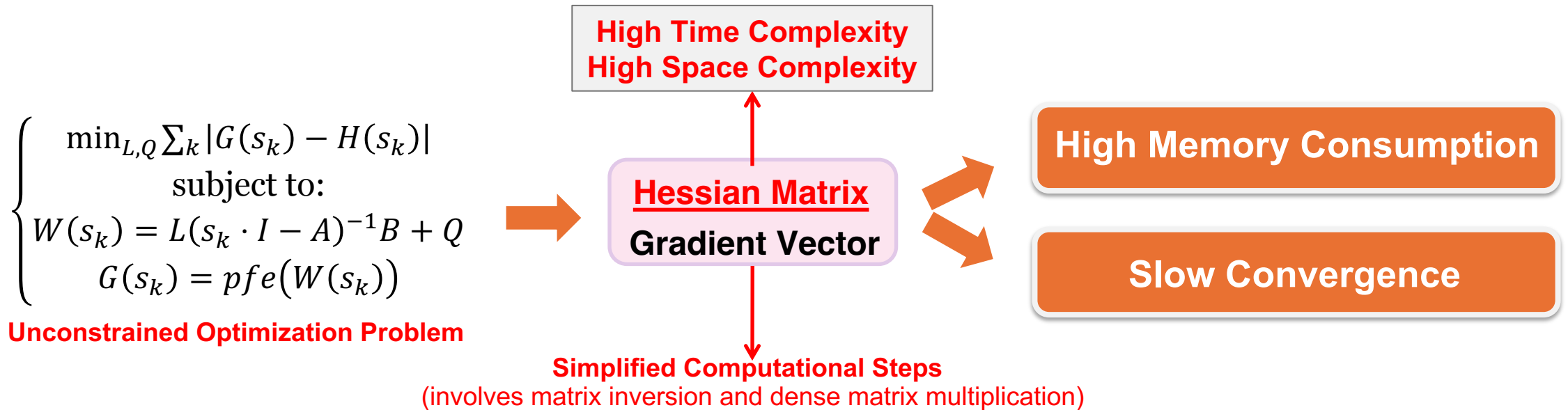
The reformulated optimization problem can also be analogized to neural network training for solution.





# G-SpNN | Further Analysis of the Main Challenges

For the transformed unconstrained optimization problem, it can be solved by computing the **gradient vector** and **Hessian matrix**, but this is also the main source of the challenges.



$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y_C}{\partial x_L} & \frac{\partial y_C}{\partial x_Q} \\ \frac{\partial y_D}{\partial x_L} & \frac{\partial y_D}{\partial x_Q} \end{bmatrix} = \begin{bmatrix} -J_{CK}M^{-1}\mathcal{N}_L(L) + \mathcal{K}_Q(Q) & \mathcal{K}_L(L) \\ 0 & \mathcal{N}_Q(Q) \end{bmatrix}$$

$$\frac{\partial^2 y}{\partial x_k^2} = \frac{\partial^2 y}{\partial x^2} \mathbf{e}_k = \frac{\partial}{\partial x} \left( \frac{\partial y}{\partial x} \right) \mathbf{e}_k = \begin{bmatrix} -J_{CK}M^{-1}\mathcal{N}_L(L_k) + K_Q(Q_k) & K_L(L_k) \\ 0 & \mathcal{N}_Q(Q_k) \end{bmatrix}$$

$$\frac{\partial f}{\partial y} = 2E^T R \quad \frac{\partial^2 f}{\partial y^2} = 2R^T R \quad \frac{\partial f}{\partial x} = \frac{\partial f}{\partial y} \frac{\partial y}{\partial x}$$

$$\frac{\partial^2 f}{\partial x^2} = \frac{\partial}{\partial x} \left( \frac{\partial f}{\partial x} \right)^T = \frac{\partial}{\partial x} \left( \frac{\partial f}{\partial y} \frac{\partial y}{\partial x} \right)^T = \left( \frac{\partial y}{\partial x} \right)^T \frac{\partial^2 f}{\partial y^2} \frac{\partial y}{\partial x} + \frac{\partial^2 y}{\partial x^2} \left( \frac{\partial f}{\partial y} \right)^T$$

As previously mentioned, the full computation of **the Hessian matrix incurs significant memory and time overhead**; therefore, we adopt the **LBFGS method** to approximate the inverse of the Hessian matrix.

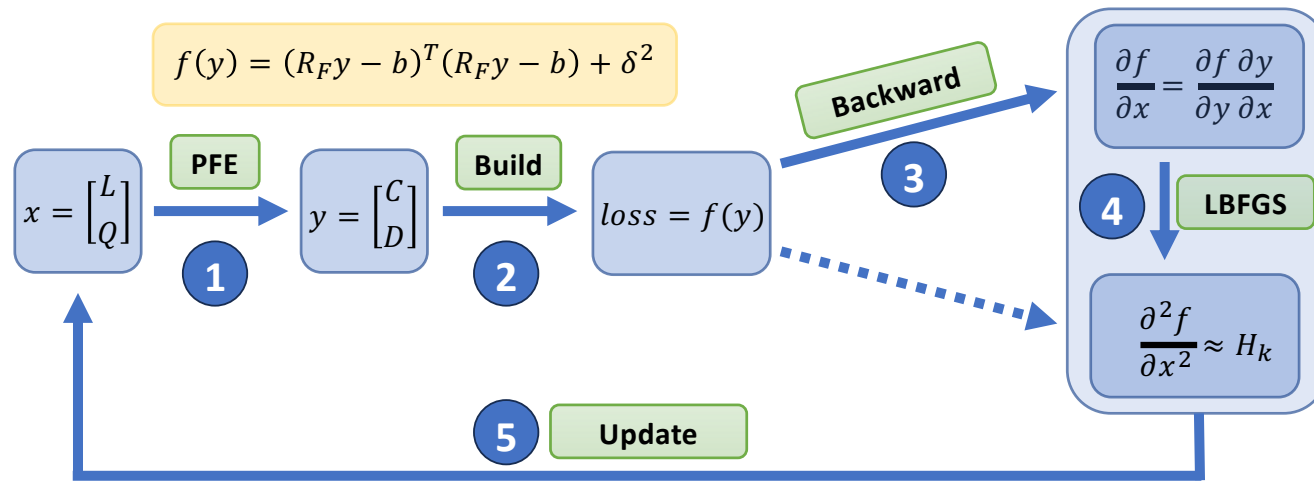
**The update formula in LBFGS method**

$$H_{k+1} = H_k - \frac{H_k \Delta g_k \Delta g_k^T H_k}{\Delta g_k^T \Delta x_k} + \frac{\Delta x_k \Delta x_k^T}{\Delta x_k^T \Delta g_k}$$

Where:

- 1)  $H_k$  is the Hessian inverse approximation at iteration  $k$
- 2)  $\Delta g_k = \left(\frac{\partial f}{\partial x}\right)_{k+1} - \left(\frac{\partial f}{\partial x}\right)_k$  is the change in the gradient
- 3)  $\Delta x_k = x_{k+1} - x_k$  is the change in the parameter

The essence of LBFGS is to efficiently approximate the inverse of the Hessian matrix by **retaining information from the most recent iterations**, thereby accelerating the convergence of large-scale optimization problems.



Computational Graph with LBFGS Method

**Step 1:**  $x$  undergoes the *PFE* operation to generate parameters  $y$  (the passive system  $G(s)$ ).

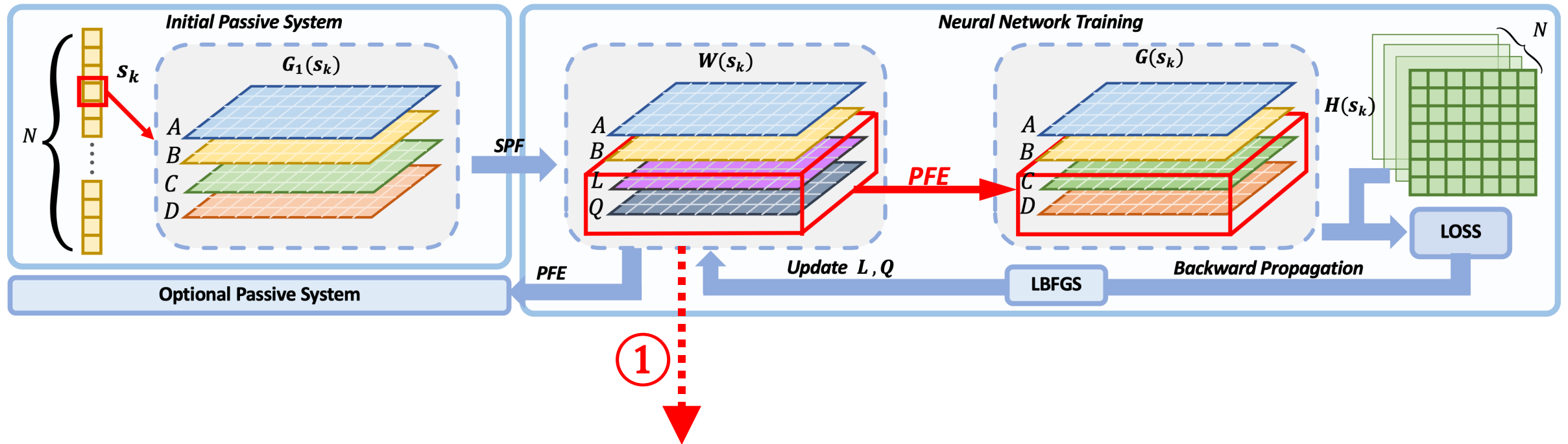
**Step 2:** Construct the loss function.

**Step 3:** Use automatic differentiation to perform backpropagation via the chain rule and compute the first-order derivative of the loss function with respect to the network weights  $x$ .

**Step 4:** Use the LBFGS method to approximate the inverse of the Hessian matrix.

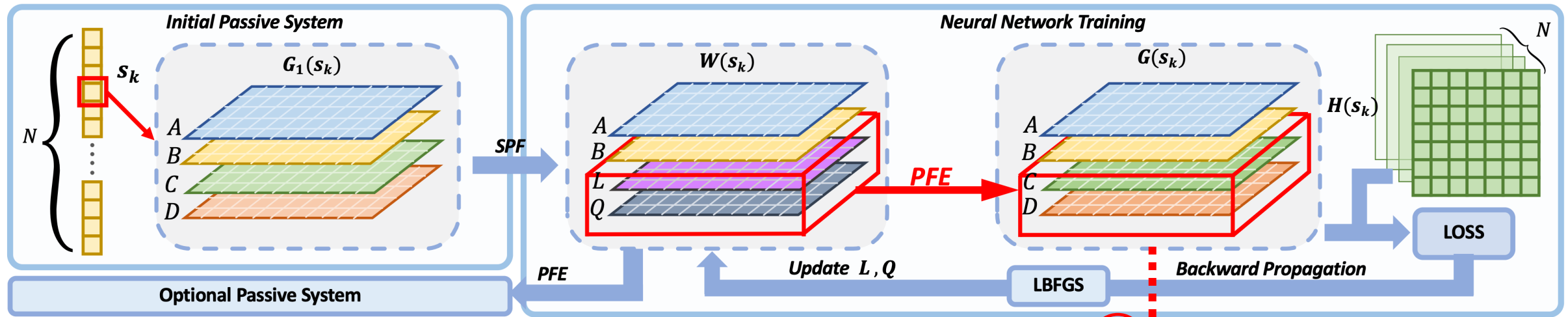
**Step 5:** Update the parameters  $L$  and  $Q$  of the initial passive system.

# G-SpNN | The Overall Framework



Starting with a given passive system  $G_1(s_k)$ , an unconstrained system  $W(s_k)$  is first derived by SPF transformation and **represented as a layer in the neural network.**

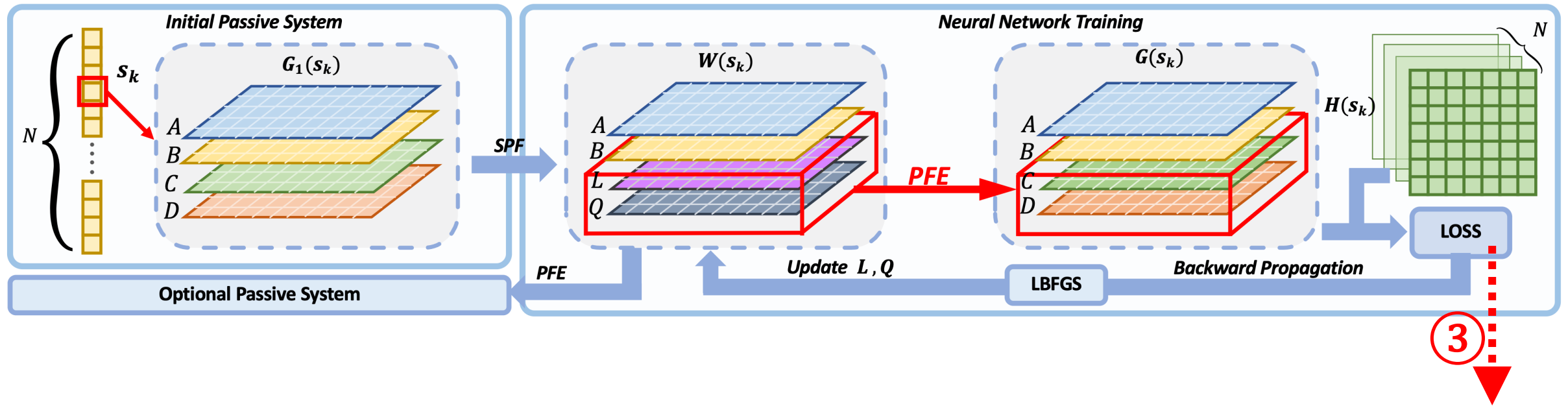
# G-SpNN | The Overall Framework



Next, the PFE transformation is applied to generate a new network layer  $G(s_k)$ , corresponding to a passive system. This reformulates the problem as a neural network training task.

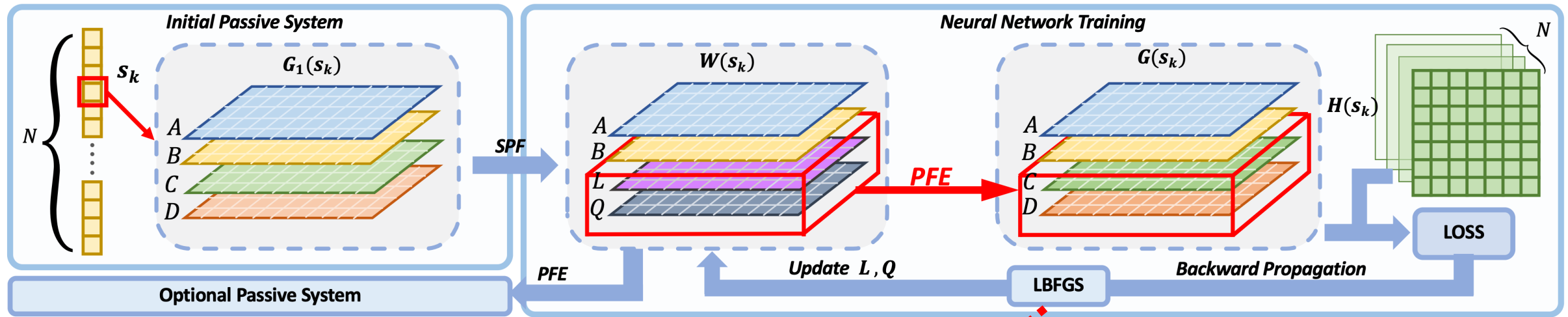


# G-SpNN | The Overall Framework



During forward propagation, the system, together with the tabulated data  $H(s_k)$ , is used to **compute the loss value**.

# G-SpNN | The Overall Framework



For efficient training, the LBFGS method is further incorporated with backpropagation to **compute gradients and update the network parameters**. Once the iteration stopping criteria are satisfied, the optimized passive system is obtained.

# OUTLINE

- 1 Background
- 2 G-SpNN
- 3 Experiment
- 4 Conclusion

# Experiment | Experimental Setup

## Experimental Environment

- We implement and test the proposed G-SpNN on an i7-14700KF @5.6GHz CPU with 32GB of memory, and a GeForce RTX 4070 SUPER GPU with 12GB of VRAM.
- G-SpNN is implemented based on PyTorch and compared against the framework DAO, which is implemented in MATLAB and is open-sourced on GitHub.

Touchstone files information

Num	Case	Ports	n	m	N
1	Telluride	11	56	11	258
2	test_5	30	199	30	2000
3	sp125_uniform_2	64	342	64	400
4	CKDIST_TUNEDBUF	64	338	64	2000
5	pll_testcase	138	727	138	300

$$H^{m \times m}(s_k) \quad (1 \leq k \leq N) \quad A^{n \times n}, B^{n \times m}, C^{m \times n}, \text{ and } D^{m \times m}$$

Fitting results of VF and LC

Num	VF			LC		
	SS Error	Time(s)	Passivity	SS Error	Time(s)	Passivity
1	1.52e-4	0.8906	non-passive	5.91e-1	0.1875	passive
2	1.72e-7	17.78	non-passive	6.56e-5	7.64	passive
3	6.29e-4	25.59	non-passive	6.56e-4	6.906	passive
4	4.27e-3	96.89	non-passive	5.16e-3	17.21	passive
5	5.49e-4	75.23	non-passive	5.70e-4	91.79	passive

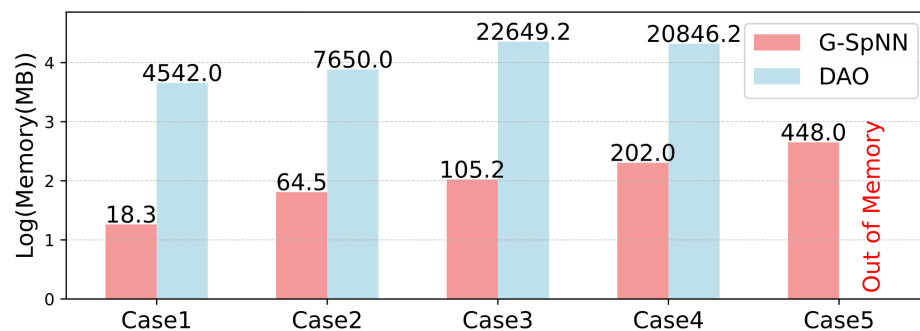


# Experiment | Convergence Speed and Memory Usage

Comparison of G-SpNN and DAO. The “–” indicates memory overrun during execution.

Num	Initial Loss	DAO					G-SpNN				
		Time (s)	#Iteration	Final Loss	SS Error	Passivity	Time (s)	#Iteration	Final Loss	SS Error	Passivity
1	6.17e12	17.47	93	4.8656	2.26e-3	passive	94.19	232	4.19e-2	2.46e-4	passive
2	6.50e8	104.65	7	3.86e-2	2.65e-5	passive	97.93	328	3.89e-2	2.65e-5	passive
3	17.36	2116.48	4	17.21	6.51e-4	passive	300.46	803	17.15	6.47e-4	passive
4	2.31e3	3923.14	9	2.24e3	4.97e-3	passive	176.56	469	2.14e3	4.69e-3	passive
5	78.32	–	–	–	–	–	403.35	456	77.62	5.58e-4	passive
Average		1540.44	28	565.53	1.98e-3		214.49	458	446.97	1.23e-3	

➤ G-SpNN achieves an average speedup of **7.63×** compared to DAO.



➤ DAO's average memory consumption is **171.3×** that of G-SpNN.

➤ Keeping the memory usage **almost constant** with an increasing number of ports.



# Experiment | Convergence Speed and Memory Usage

Comparison of G-SpNN and DAO. The “–” indicates memory overrun during execution.

Num	Initial Loss	DAO					G-SpNN				
		Time (s)	#Iteration	Final Loss	SS Error	Passivity	Time (s)	#Iteration	Final Loss	SS Error	Passivity
1	6.17e12	17.47	93	4.8656	2.26e-3	passive	94.19	232	4.19e-2	2.46e-4	passive
2	6.50e8	104.65	7	3.86e-2	2.65e-5	passive	97.93	328	3.89e-2	2.65e-5	passive
3	17.36	2116.48	4	17.21	6.51e-4	passive	300.46	803	17.15	6.47e-4	passive
4	2.31e3	3923.14	9	2.24e3	4.97e-3	passive	176.56	469	2.14e3	4.69e-3	passive
5	78.32	–	–	–	–	–	403.35	456	77.62	5.58e-4	passive
Average		1540.44	28	565.53	1.98e-3		214.49	458	446.97	1.23e-3	

## More Detailed Explanation

For Case 1, although the runtime of DAO appears shorter, the comparison of the final loss and steady-state error shows that DAO actually experiences **pseudo-convergence** and does not reach the optimal solution.

# Experiment | Convergence Speed and Memory Usage

Comparison of G-SpNN and DAO. The “–” indicates memory overrun during execution.

Num	Initial Loss	DAO					G-SpNN				
		Time (s)	#Iteration	Final Loss	SS Error	Passivity	Time (s)	#Iteration	Final Loss	SS Error	Passivity
1	6.17e12	17.47	93	4.8656	2.26e-3	passive	94.19	232	4.19e-2	2.46e-4	passive
2	6.50e8	104.65	7	3.86e-2	2.65e-5	passive	97.93	328	3.89e-2	2.65e-5	passive
3	17.36	2116.48	4	17.21	6.51e-4	passive	300.46	803	17.15	6.47e-4	passive
4	2.31e3	3923.14	9	2.24e3	4.97e-3	passive	176.56	469	2.14e3	4.69e-3	passive
5	78.32	–	–	–	–	–	403.35	456	77.62	5.58e-4	passive
Average		1540.44	28	565.53	1.98e-3		214.49	458	446.97	1.23e-3	

## More Detailed Explanation

- For Case 3 and Case 4, it should be noted that the DAO method is **forcibly terminated** during the iteration process due to **memory overflow** and does not reach the predefined convergence criterion.
- For Case 5, DAO experiences a **memory overflow** during the first iteration and could not complete the iteration.

# Experiment | Convergence Speed and Memory Usage

Comparison of G-SpNN and DAO. The “–” indicates memory overrun during execution.

Num	Initial Loss	DAO					G-SpNN				
		Time (s)	#Iteration	Final Loss	SS Error	Passivity	Time (s)	#Iteration	Final Loss	SS Error	Passivity
1	6.17e12	17.47	93	4.8656	2.26e-3	passive	94.19	232	4.19e-2	2.46e-4	passive
2	6.50e8	104.65	7	3.86e-2	2.65e-5	passive	97.93	328	3.89e-2	2.65e-5	passive
3	17.36	2116.48	4	17.21	6.51e-4	passive	300.46	803	17.15	6.47e-4	passive
4	2.31e3	3923.14	9	2.24e3	4.97e-3	passive	176.56	469	2.14e3	4.69e-3	passive
5	78.32	–	–	–	–	–	403.35	456	77.62	5.58e-4	passive
Average		1540.44	28	565.53	1.98e-3		214.49	458	446.97	1.23e-3	

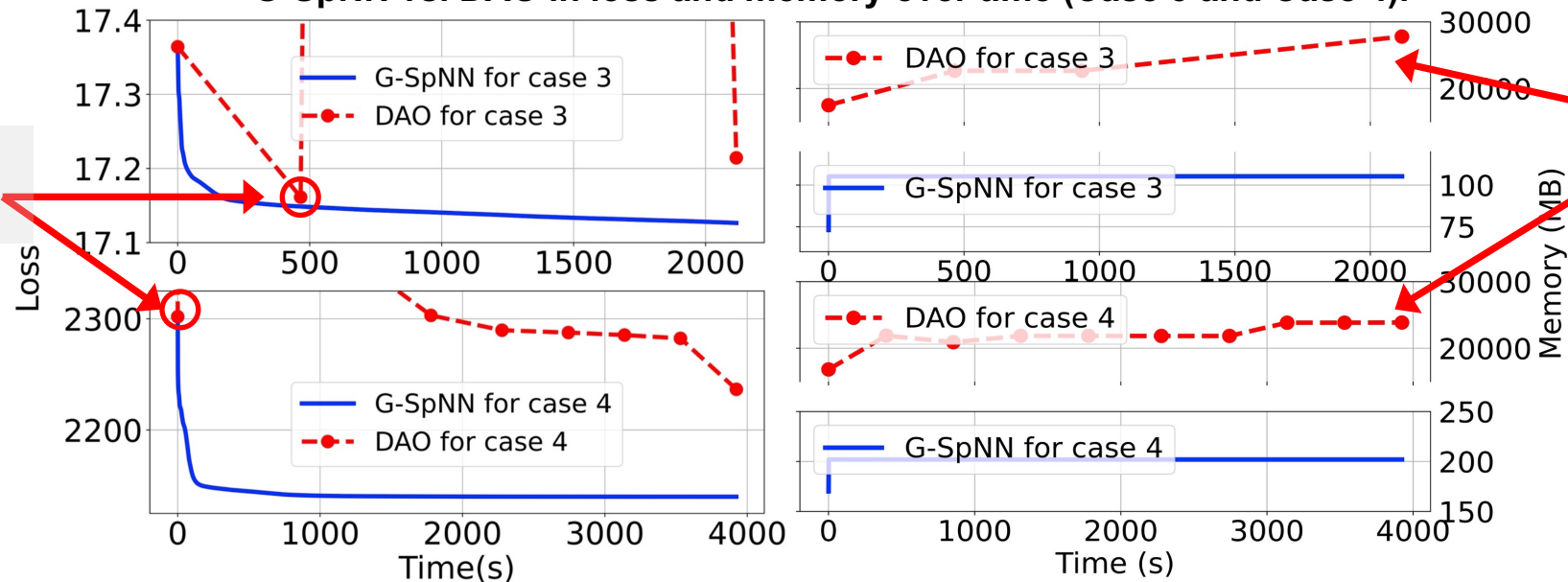
## More Detailed Explanation

Due to the high time and space complexity of the DAO method, we limit **the number of poles** in the VF method for cases 3-5 to ensure computational feasibility, which leads to higher SS Error and limits the reduction in loss.

# Experiment | Convergence Speed and Memory Usage

G-SpNN vs. DAO in loss and memory over time (Case 3 and Case 4).

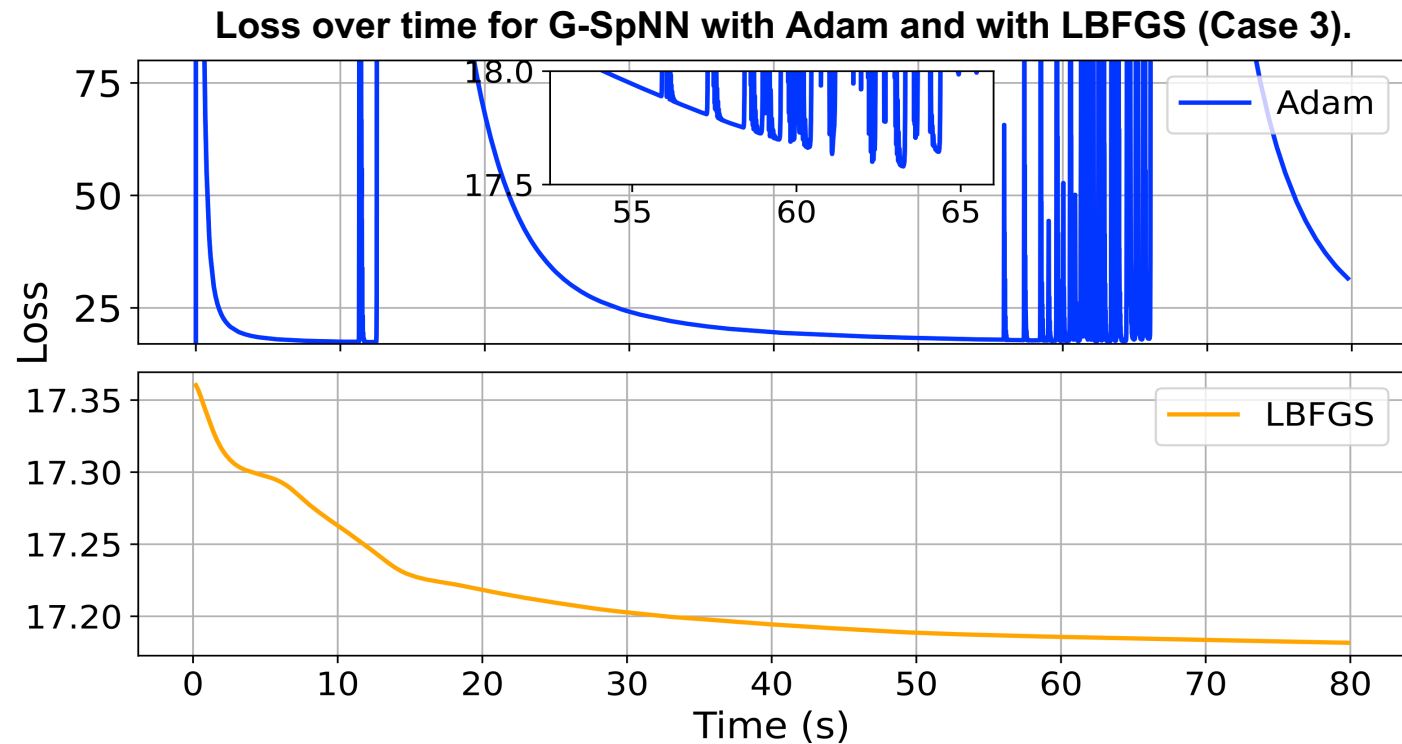
Gradient Explosion occurs in DAO



Memory Usage increases in DAO

- Using Case 3 and Case 4 as examples, figure shows the **loss and memory usage** variations during iterations.
- G-SpNN has a **smoother convergence** process with better performance, achieving **lower loss** compared to DAO.

# Experiment | Compare with Adam



The LBFGS method enables G-SpNN to **progress steadily toward convergence**, owing to second-order information guiding more effective update directions.

# OUTLINE

- 1 Background
- 2 G-SpNN
- 3 Experiment
- 4 Conclusion



# Conclusion

- Casting the passive macromodeling problem to **neural network training**, thus leveraging **GPU acceleration**.
- Using the LBFGS method to efficiently **approximate the Hessian inverse matrix**, efficiently decrease the memory cost and time overhead. Keeping the memory usage almost constant with an increasing number of ports.
- Experimental results show that G-SpNN not only converges more stably and quickly than DAO, with an average speedup of **7.63×**, its memory usage can be **reduced by two orders of magnitude** in test cases.



# Thanks for Listening Any Questions?

AI

Security

Systems

EDA

Design

Email: [z.jin@zju.edu.cn](mailto:z.jin@zju.edu.cn)