

# Automatic and Structure-Aware Sparsification of Hybrid Neural ODEs with Application to Glucose Prediction

ICLR 2026 Poster Presentation

Bob Junyi Zou<sup>1</sup> Lu Tian<sup>2</sup>

<sup>1</sup>Institute for Computational and Mathematical Engineering, Stanford University

<sup>2</sup>Department of Biomedical Data Science, Stanford University

May 2026

# The Challenge: Hybrid Models are Powerful but “Bloated”

## The Promise of Hybrid Modeling

- ▶ Combines mechanistic priors (ODEs) with neural flexibility.
- ▶ Essential for data-scarce healthcare settings (e.g., T1D glucose forecasting).

## The Reality Gap

- ▶ **Over-parameterization:** Mechanistic models are designed for biological completeness, not predictive efficiency.
- ▶ **Redundancy:** A glucose model might have 20+ latent states but only 4 observables.
- ▶ **Instability:** Cycles and feedback loops lead to stiffness and training failures.

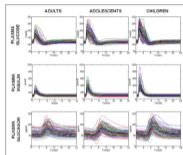


Figure 5. Simulated plasma glucose (upper), insulin (middle), and glucagon (lower panels) in the 100 in silico adults (left), adolescents (middle), and children (right panels).

new criteria for virtual subject generation to automatically discard subjects with a nonplausible physiological behavior.

Simulation results have been presented for single meal optimal open-loop therapy in adult, adolescent, and children populations. S2013 provides a more reliable framework for in silico trials for regulatory purposes, for testing glucose sensors and insulin augmented pump prediction methods, and for closed-loop single/dual hormone controller design, testing, and validation. However, it is worth noting that both S2008 and S2013 simulators have been validated and accepted by FDA for a single meal scenario only. Multiple meal scenarios can obviously be simulated, but since the simulator does not include time-varying parameters, the results would not be realistic. Inclusion of meal-by-meal and day-by-day parameter variations in the simulator is under investigation. We anticipate that the first step will be the incorporation of insulin sensitivity meal-by-meal variation as reported in Hinzaw et al.<sup>23</sup>

### Appendix

#### Model Equations

Glucose subsystem:

$$\begin{cases} \dot{G}_p(t) = EGP(t) + Rat(t) - U_p(t) - E(t) - k_p \cdot G_p(t) + k_i \cdot G(t) \\ G_p(0) = G_b \\ \dot{G}_i(t) = -U_i(t) + k_i \cdot G_p(t) - k_i \cdot G(t) \\ G_i(0) = G_b \\ G(t) = \frac{G_p}{f_c} \\ G(0) = G_i \end{cases} \quad (A1)$$

Insulin subsystem:

$$\begin{cases} \dot{I}_p(t) = -(m_1 + m_2) \cdot I_p(t) + m_1 \cdot I(t) + R_{in}(t) \\ I_p(0) = I_{p0} \\ \dot{I}_i(t) = -(m_1 + m_2) \cdot I_i(t) + m_2 \cdot I_p(t) - I_i(0) = I_{i0} \\ I_i(0) = \frac{I_p(0)}{f_c} \end{cases} \quad (A2)$$

Glucose rate of appearance:

$$\begin{cases} \dot{Q}_{in}(t) = Q_{in0}(t) + Q_{in}(t) \\ Q_{in}(0) = 0 \\ \dot{Q}_{out}(t) = -k_{out} \cdot Q_{out}(t) + D \cdot S(t) \\ Q_{out}(0) = 0 \\ \dot{Q}_{in}(t) = -k_{in}(Q_{in}(t) - Q_{in0}(t)) + k_{in} \cdot Q_{in}(t) \\ Q_{in}(0) = 0 \\ \dot{Q}_{out}(t) = -k_{out} \cdot Q_{out}(t) + k_{out} \cdot Q_{out}(t) \\ Q_{out}(0) = 0 \\ Rat(t) = \frac{f \cdot k_{in} \cdot Q_{in}(t)}{BW} \\ Rat(0) = 0 \end{cases} \quad (A3)$$

with

$$\begin{cases} k_{in}(Q_{in}) = k_{in} + \frac{k_{in} - k_{in0}}{2} \\ [\tanh[a(Q_{in} - b \cdot D)] - \tanh[b(Q_{in} - c \cdot D)]] + 2 \end{cases} \quad (A4)$$

Endogenous glucose production:

$$EGP(t) = k_{e1} - k_{e2} \cdot G_p(t) - k_{e3} \cdot X^1(t) + \xi \cdot X^2(t) \quad (A5)$$

$$X^1(t) = -k_1 \cdot [X^1(t) - I(t)] \quad X^1(0) = I_0 \quad (A6)$$

$$\begin{cases} \dot{I}(t) = -k_1 \cdot [I(t) - I(t)] \\ I(0) = I_0 \end{cases} \quad (A7)$$

$$\begin{cases} \dot{X}^2(t) = -k_{e2} \cdot X^2(t) + k_{e2} \cdot \max[(I(t) - H), 0] \\ X^2(0) = 0 \end{cases} \quad (A8)$$

Glucose utilization:

$$U_p(t) = F_{in} \quad (A9)$$

$$U_i(t) = \frac{[V_{in} + V_{in} \cdot X(t) \cdot (1 + \tau_{in} \cdot \text{risk})] \cdot G_i(t)}{K_{in} + G_i(t)} \quad (A10)$$

with

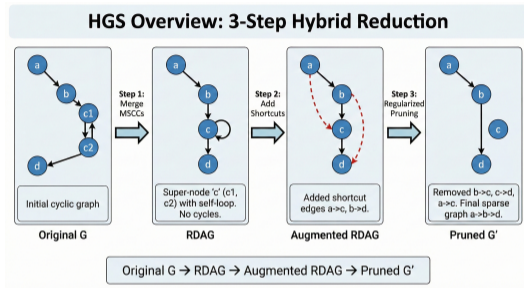
# Why not use standard pruning?

- ▶ **Stability:** Feedback loops can cause training instabilities.
- ▶ **Timescale Separation:** Standard pruning doesn't account for "fast" variables that can be effectively skipped.
- ▶ **Gradient-Free Methods:** Greedy search or genetic algorithms are too slow for training complex ODEs.
- ▶ **Pure Data-Driven GNN Pruning:** Ignores physical laws; risks breaking reachability or creating nonsensical pathways.

## Our Solution: Hybrid Graph Sparsification (HGS)

A gradient-based, structure-aware pipeline.

# Methodology: The HGS Pipeline



- 1 **Merge Cycles (Stability):** Collapse Maximal Strongly Connected Components (MSCCs) into super-nodes.
- 2 **Add Shortcuts (Expressivity):** Augment pathways via partial transitive closure to model timescale separation.
- 3 **Regularized Pruning (Sparsity):** Train with  $L_1$  penalty on edge weights to select the best subgraph.

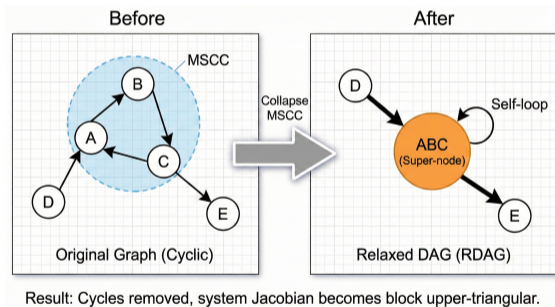
# Step 1: Merging MSCCs for Stability

## Problem:

- ▶ Physiological models are full of feedback loops.
- ▶ Cycles causes stiffness and exploding gradients.

## HGS Action:

- ▶ Collapse cycles into “super-nodes”.
- ▶ **Result:** A Relaxed DAG (RDAG) that preserves the high-level mechanistic causal structure. System Jacobian becomes block upper-triangular.



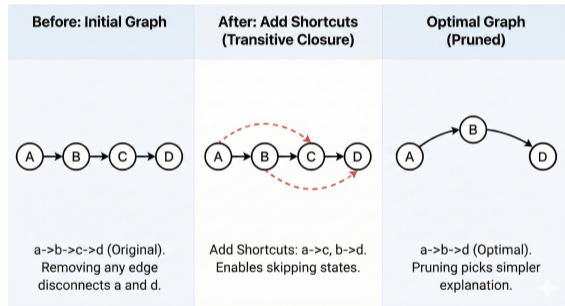
## Step 2: Mechanistic Shortcuts

### Intuition:

- ▶ In biology,  $A \rightarrow B \rightarrow C$  might happen so fast that  $A \rightarrow C$  is a better predictive model (Timescale Separation).
- ▶ Standard pruning can only *remove* edges, not *skip* steps.

### HGS Action:

- ▶ Transitive Closure: Add shortcut edges ( $u \rightarrow v$ ) if  $u$  is an ancestor of  $v$ .
- ▶ Allows the optimizer to "bypass" redundant latent states without violating mechanistic reachability.

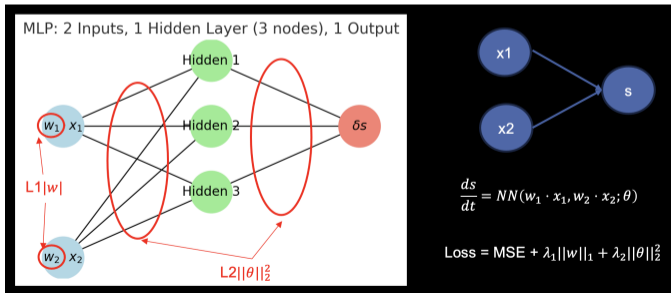


## Step 3: Regularized Pruning

We learn a sparse subgraph using a composite loss function:

$$\mathcal{L} = \text{MSE} + \lambda_1 \sum_{(u,v)} |w_{u,v}| + \lambda_2 \|\Theta\|_2^2$$

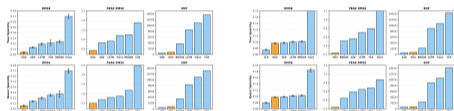
- ▶  $L_1$  on Edge Weights ( $w$ ): Encourages sparsity.
- ▶  $L_2$  on Parameters ( $\Theta$ ): Ensures stability.
- ▶ **Equivalence:** Equivalent to a first-layer Group LASSO.



# Results: Synthetic Experiments

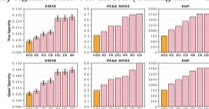
**Setup:** Ground truth sparse ODE hidden inside a redundant comprehensive graph.

- ▶ **HGS vs. Black Box:** Outperforms LSTM/TCN in low-data regimes.
- ▶ **HGS vs. Other Pruning:** Achieves lower RMSE and better sparsity (ENP) than NeuralSparse or Group LASSO.

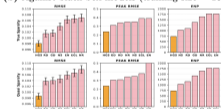


(a) Against black-box models (training size = 100)

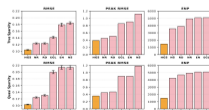
(b) Against black-box models (training size = 1000)



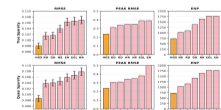
(c) Against other reduction methods  
(refined initial graph, training size = 100)



(d) Against other reduction methods  
(refined initial graph, training size = 1000)



(e) Against other reduction methods  
(comprehensive initial graph, training size = 100)



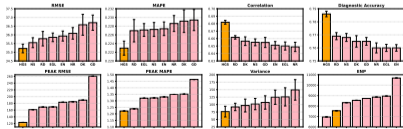
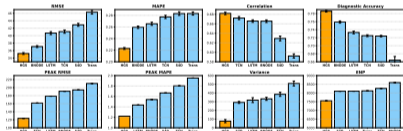
(f) Against other reduction methods  
(comprehensive initial graph, training size = 1000)

# Application: T1D Glucose Forecasting

**Data:** T1D Exercise Initiative (T1DEXI), 342 time series from 105 patients.

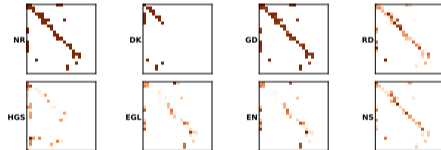
## Performance:

- ▶ Surpasses Black-box Neural ODEs and standard mechanistic models.
- ▶ **Robustness:** Significantly lower Peak RMSE (worst-case error).



## Interpretability:

- ▶ HGS pruned glucagon feedback loops.
- ▶ Suggests impaired glucagon response during exercise-induced hypoglycemia.



## Summary

We propose **HGS**, a three-step pipeline for Hybrid Neural ODEs:

- 1 **Structure-Aware:** Respects mechanistic plausibility.
- 2 **Efficient:** Gradient-based (no expensive greedy search).
- 3 **Effective:** Improves accuracy and robustness in medical time-series.

## Acknowledgment

This publication is based on research using data from Jaeb Center for Health Research Foundation that has been made available through Vivli, Inc. Vivli has not contributed to or approved, and is not in any way responsible for, the contents of this publication. We thank Alex Wang, Dessi Zeharieva, Emily Fox, Matthew Levine and Ramesh Johari for providing assistance with data acquisition and preliminary discussions on T1D and model reduction.

**Paper available at the ICLR poster session.**