

# Glassoformer: a query-sparse transformer for post-fault power grid voltage prediction

Yunling Zheng<sup>1</sup>   Carson Hu<sup>1</sup>   Guang Lin<sup>2</sup>   Meng Yue<sup>3</sup>  
Bao Wang<sup>4</sup>   Jack Xin<sup>1</sup>

<sup>1</sup>Department of Mathematics, University of California, Irvine

<sup>2</sup>Department of Mathematics and School of Mechanical Engineering, Purdue University

<sup>3</sup>Interdisciplinary Science Department, Brookhaven National Laboratory

<sup>4</sup>Department of Mathematics and Scientific Computing and Imaging Institute,  
The University of Utah

April 17, 2022



# Table of Contents

## 1 Introduction

## 2 Our approach

- Transformer
- Self-attention
- Query sparsification

- RGSM

## 3 Experiment and result

- Experiment
- Result

## 4 Acknowledgments



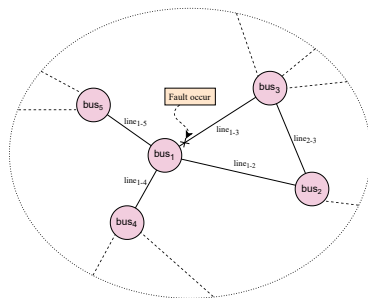
# Section 1

## Introduction



# Background

- Power grid system is a lifeline networks but fragile when interruptions occur.
- Have an online prediction for power grid stability after fault happens.
- Traditional physical simulation approach:
  - time-consuming
  - details of fault information
  - not suitable for online assessment



Power grid system with a fault on line<sub>1-3</sub>.

# Related work

- Data driven solutions
  - predict stability by binary indicator
  - estimate stability margin
  - not enough
- Classical signal processing: Prony's Method
  - linear time-invariant filters
  - limited to response with a rational function
  - require manual observation and parameter choice
  - not robust
- Improved with spatial coupling: 1D-CNN
  - deep feature extraction
- Our approach: Glassoformer
  - equip encoder-decoder with attentions
  - propose efficient sparse self-attention to reduce the cost
  - outperforms existing methods



# Table of Contents

- 1 Introduction
- 2 Our approach
  - Transformer
  - Self-attention
  - Query sparsification
- RGSM
- 3 Experiment and result
  - Experiment
  - Result
- 4 Acknowledgments



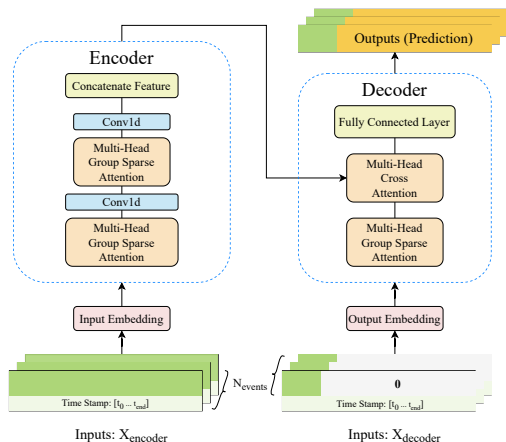
## Section 2

### Our approach



# Overview of the model

- Based on encoder-decoder structure.
- Adjusted embedding layers for post-fault prediction problems.
- Different inputs for encoder and decoder.



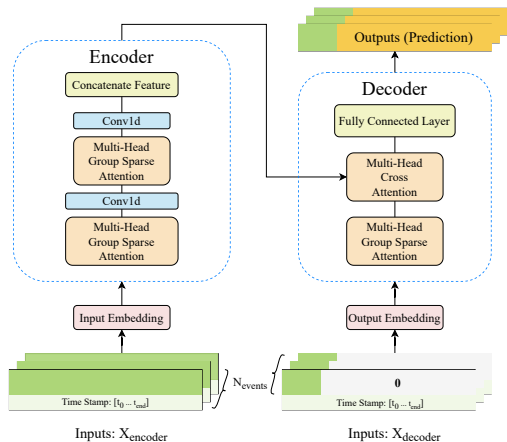
Auto-encoder architecture





# Embedding

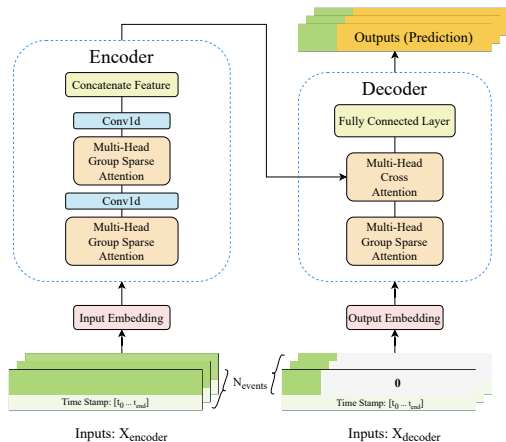
- In transient prediction, full period time stamps ( $\mathbf{x}^s$ ) and temporal time-series features  $\mathbf{x}$  are required.
- Input embedding:
  - positional embedding
  - 1D convolution
- Output embedding:
  - zero padding
  - 1D convolution



Auto-encoder architecture

# Encoder and Decoder

- Encoder consists of two identical layers:
  - multi-head group sparse attention sub-layer
  - 1D-convolution sub-layer
- Decoder:
  - multi-head group sparse self attention
  - multi-head cross attention



Auto-encoder architecture



# Self-attention

- long-range dependency and enable parallel processing
- Self-attention transforms input sequence  $\tilde{\mathbf{X}} \in \mathcal{R}^{N, D_x}$  into an output sequence  $\hat{\mathbf{V}}$  by two steps:
  - Step 1: Project input into 3 matrices:

$$\mathbf{Q} = \mathbf{W}_Q^T \tilde{\mathbf{X}}^T; \mathbf{K} = \mathbf{W}_K^T \tilde{\mathbf{X}}^T; \mathbf{V} = \mathbf{W}_V^T \tilde{\mathbf{X}}^T,$$

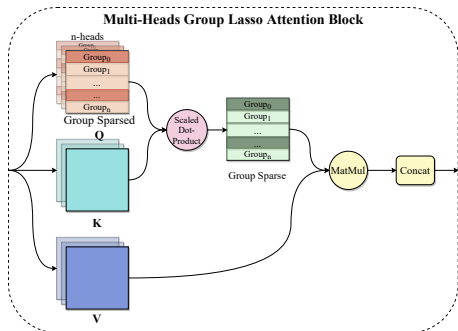
- Step 2: Compute output with softmax:

$$\hat{\mathbf{V}} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{D}}\right)\mathbf{V} := \mathbf{AV}$$



# Query sparsification

- For long sequence, the computational and memory cost are both  $\mathcal{O}(N^2)$ .
- Efficient transformers have been proposed: Informer, cluster transformer, etc.
- Our approach: Query sparsification
  - apply group Lasso (GLasso) penalty to sparsify  $\mathbf{W}_Q$
  - reduce number of nonzero rows of  $\mathbf{Q} = \mathbf{W}_Q^T \tilde{\mathbf{X}}^T$  to  $\mathcal{O}(1)$  (row-wise group sparsity)
  - complexity of  $\mathbf{AV}$  goes down to  $\mathcal{O}(N)$



Glassoformer attention block

# Algorithm

- Propose Relaxed Group-wise Splitting Method (RGSM) to achieve row-wise sparsity in query matrix  $\mathbf{Q}$ :
- Consider group lasso penalty  $\|\cdot\|_{GL}$ , the loss is:

$$L(\theta, \mathbf{W}_Q) = f(\theta, \mathbf{W}_Q) + \lambda \|\mathbf{W}_Q\|_{GL}$$

- Solve proximal problem in closed-form:

$$\mathbf{y}_g^* = \arg \min_{\mathbf{y}_g} \lambda \|\mathbf{y}_g\|_2 + \sum_{i \in I_g} \frac{1}{2} \|\mathbf{y}_{g,i} - \mathbf{w}_{g,i}\|_2^2$$

here  $\mathbf{W}_Q = \{\dots, \mathbf{w}_g, \dots\}$

- Obtain GLasso proximal operator:

$$\mathbf{y}_g^* = \mathbf{Prox}_{GL, \lambda}(\mathbf{w}_g) \triangleq \mathbf{w}_g \max(\|\mathbf{w}_g\|_2 - \lambda, \mathbf{0}) / \|\mathbf{w}_g\|_2.$$



## RGSM algorithm

Update for each iteration:

$$\mathbf{u}_g^t = \mathbf{Prox}_\lambda(\mathbf{w}_g^t), \quad \text{for } g = 1, \dots, N$$

$$(\theta, \mathbf{w})^{t+1} = (\theta, \mathbf{w})^t - \eta \nabla f(\theta^t, \mathbf{w}^t) - (0, \eta \beta(\mathbf{w}^t - \mathbf{u}^t))$$

# Convergence

## Lemma

$$L_{\beta}(\mathbf{v}^{t+1}, \mathbf{u}^t) \leq L_{\beta}(\mathbf{v}^t, \mathbf{u}^t) + \left( \frac{L_{ip}}{2} + \frac{\beta}{2} - \frac{1}{\eta} \right) \|\mathbf{v}^{t+1} - \mathbf{v}^t\|^2$$

## Theorem

*If  $f$  is coercive ( $f$  bounded implies that of its independent variables, true when standard weight decay is present in network training) and the learning rate  $\eta < 2/\beta + L_{ip}$ , then  $L_{\beta}(\mathbf{v}^t, \mathbf{u}^t)$  decreases monotonically in  $t$ , and  $(\mathbf{v}^t, \mathbf{u}^t)$  converges sub-sequentially to a limit point  $(\bar{\mathbf{v}}, \bar{\mathbf{u}})$ , from which  $\bar{\mathbf{w}}$  is extracted to speed up inference.*



# Table of Contents

- 1 Introduction
- 2 Our approach
  - Transformer
  - Self-attention
  - Query sparsification
- 3 Experiment and result
  - RGSM
  - Experiment
  - Result
- 4 Acknowledgments





## Section 3

### Experiment and result



# Experimental setting

- **Dataset:** Simulation of power system in New York/New England
  - 16 generators 68 buses
  - 2248 fault events, and each 10 seconds signal
  - voltage from bus and current from line
  - graph structure of buses and lines
- Experiment details:
  - baseline: 1D-CNN, Informer, Prony's method, Transformer with lasso
  - platform: GTX-1080Ti



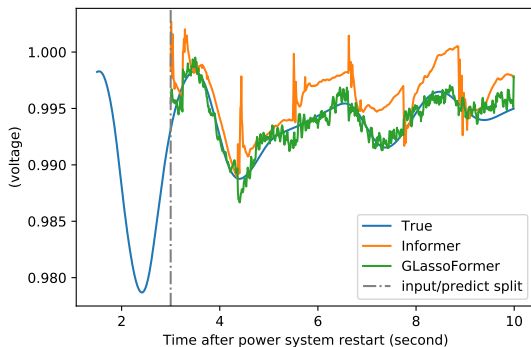
# Result

	Models	GLasso	Informer	Lasso	1DCNN	Prony
I	MSE( $\times 10^{-5}$ )	<b>3.189</b>	3.662	3.532	8.014	–
	MAE( $\times 10^{-3}$ )	<b>2.374</b>	2.684	2.543	6.087	–
II	MSE( $\times 10^{-5}$ )	<b>6.115</b>	6.599	6.501	17.21	397.6
	MAE( $\times 10^{-3}$ )	<b>3.520</b>	3.877	3.611	9.264	47.3

**Table:** Voltage prediction error comparison.

I (II): input data with (without) neighbor voltage and current features.

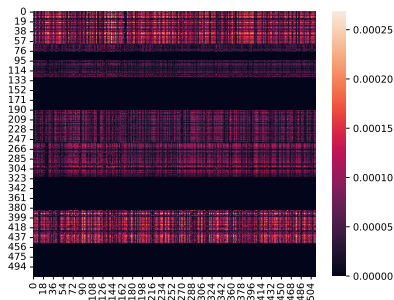
# Result



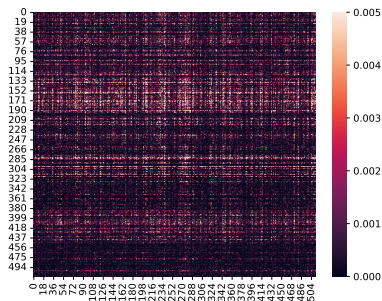
Sample post-fault voltage predictions in time (to the right of the vertical dashed line at time 3) with input from the left of the line.

# Result

Visualization of sparsity patterns (black pixel) in  $Q$ . GLassoformer obtains a better group (row-wise) sparsity.



Sample sparsity pattern of  $Q$  by GLassoformer



Sample sparsity pattern of  $Q$  by Lasso penalty

# Result

Our network: Glassoformer is faster and sparser.

Models	GLasso	Lasso	Informer	1DCNN
Num of Params (M)	5.827	5.707	7.257	0.706
Inference Time (ms)	18.98	14.92	29.76	0.5969

**Table 1:** Comparison of model parameter size (M: million), and inference time (ms: millisecond) on GTX-1080Ti.

Models	GLasso	Informer	Lasso
Pruning rate <sup>1</sup> (%)	19.09	4.220	2.674
Training time <sup>2</sup> (second)	9.215	8.975	8.987

**Table 2:** Pruning rate and training time comparison.

<sup>1</sup>: fraction of zero query vectors (threshold =  $1e-5$ ); <sup>2</sup>: time per epoch



# Table of Contents

- 1 Introduction
- 2 Our approach
  - Transformer
  - Self-attention
  - Query sparsification
- 3 Experiment and result
  - RGSM
  - Experiment
  - Result
- 4 Acknowledgments



# Acknowledgments

The work was partially supported by NSF and DOE grants DMS-1924935, DMS-1952644, DMS-1924548, DMS-1952339, DE-SC0021142.





# The End

