

Optimal Transport-driven Implicit Neural Network Design

Hongteng Xu

GSAI, RUC

Oct. 31, 2023



中國人民大學
RENMIN UNIVERSITY OF CHINA

高瓴人工智能學院
Gaoling School of Artificial Intelligence

Outline

1 Implicit Neural Network Design

- ▶ Preliminaries and representative work
- ▶ Our motivation and principle

2 Optimal Transport-driven GNN Design

- ▶ Implicit Global Pooling: Regularized OT layers for generalized global pooling
- ▶ Implicit Message-Passing: Quasi-Wasserstein loss and transductive message-passing

3 Optimal Transport-driven Transformer Design

- ▶ Implicit Attention Layer: Sliceformer, Make multi-head attention as simple as sorting

Neural Network Design: Engineering or Art?

The progress of AI is mainly attributed to the development of model architectures.

- ▶ Vision: AlexNet, VGG, ResNet, ViT, ...
- ▶ NLP: RNN, LSTM, BERT, GPT, ...
- ▶ Graph: Spatial and Spectral GNNs, ...

Neural Network Design: Engineering or Art?

The progress of AI is mainly attributed to the development of model architectures.

- ▶ Vision: AlexNet, VGG, ResNet, ViT, ...
- ▶ NLP: RNN, LSTM, BERT, GPT, ...
- ▶ Graph: Spatial and Spectral GNNs, ...

Ironically, till now, we only summarize very coarse and empirical design principle for neural networks.

- ▶ The deeper, the larger, the better.
- ▶ Tricks: dropout, batchnorm, non-smooth activations, ...

Neural Network Design: Engineering or Art?

The progress of AI is mainly attributed to the development of model architectures.

- ▶ Vision: AlexNet, VGG, ResNet, ViT, ...
- ▶ NLP: RNN, LSTM, BERT, GPT, ...
- ▶ Graph: Spatial and Spectral GNNs, ...

Ironically, till now, we only summarize very coarse and empirical design principle for neural networks.

- ▶ The deeper, the larger, the better.
- ▶ Tricks: dropout, batchnorm, non-smooth activations, ...

Lead to a chaotic research field, with few artistic masterpieces (e.g., Transformer) + many imitations without originality.

Neural Network Design: Explicit v.s. Implicit

- ▶ **Explicit NN layer:**

- ▶ Define and learn $f : \mathcal{X} \mapsto \mathcal{Z}$.
- ▶ Deploy: $z = f(x)$

Neural Network Design: Explicit v.s. Implicit

- ▶ **Explicit NN layer:**

- ▶ Define and learn $f : \mathcal{X} \mapsto \mathcal{Z}$.
- ▶ Deploy: $z = f(x)$

- ▶ **Implicit NN layer:**

- ▶ Define and learn $g : \mathcal{X} \times \mathcal{Z} \mapsto \mathbb{R}^n$
- ▶ Deploy: Find z s.t. $g(x, z) = 0$

Neural Network Design: Explicit v.s. Implicit

▶ **Explicit NN layer:**

- ▶ Define and learn $f : \mathcal{X} \mapsto \mathcal{Z}$.
- ▶ Deploy: $z = f(x)$

▶ **Implicit NN layer:**

- ▶ Define and learn $g : \mathcal{X} \times \mathcal{Z} \mapsto \mathbb{R}^n$
- ▶ Deploy: Find z s.t. $g(x, z) = 0$

Implicit neural network design revisits neural networks from a more mathematical and interpretable perspective, e.g., **optimization**, differential equations, and dynamic systems.

- ▶ In general, g is more interpretable than f because its design (NOT learning!) is often determined by an optimization problem.

Optimization-based Implicit Neural Network Design

Layer	Explicit Design	Implicit Design
ReLU	$y = \max\{0, x\}$	$y = \arg \min_{z \geq 0} z - x ^2$

Optimization-based Implicit Neural Network Design

Layer	Explicit Design	Implicit Design
ReLU	$y = \max\{0, x\}$	$y = \arg \min_{z \geq 0} z - x ^2$
Linear+ReLU	$\mathbf{y} = \max\{\mathbf{0}, \mathbf{W}\mathbf{x} + \mathbf{b}\}$	$\mathbf{y} = \arg \min_{\mathbf{z} \geq \mathbf{0}} \ \mathbf{z} - \mathbf{W}\mathbf{x} - \mathbf{b}\ ^2$

Optimization-based Implicit Neural Network Design

Layer	Explicit Design	Implicit Design
ReLU	$y = \max\{0, x\}$	$y = \arg \min_{z \geq 0} z - x ^2$
Linear+ReLU	$\mathbf{y} = \max\{\mathbf{0}, \mathbf{W}\mathbf{x} + \mathbf{b}\}$	$\mathbf{y} = \arg \min_{\mathbf{z} \geq \mathbf{0}} \ \mathbf{z} - \mathbf{W}\mathbf{x} - \mathbf{b}\ ^2$
Sigmoid	$y = \frac{1}{1 + \exp(x)}$	$y = \arg \max_{z \in [0, 1]} zx + z \log z + (1 - z) \log(1 - z)$

Optimization-based Implicit Neural Network Design

Layer	Explicit Design	Implicit Design
ReLU	$y = \max\{0, x\}$	$y = \arg \min_{z \geq 0} z - x ^2$
Linear+ReLU	$\mathbf{y} = \max\{\mathbf{0}, \mathbf{W}\mathbf{x} + \mathbf{b}\}$	$\mathbf{y} = \arg \min_{\mathbf{z} \geq \mathbf{0}} \ \mathbf{z} - \mathbf{W}\mathbf{x} - \mathbf{b}\ ^2$
Sigmoid	$y = \frac{1}{1 + \exp(x)}$	$y = \arg \max_{z \in [0, 1]} zx + z \log z + (1 - z) \log(1 - z)$
Softmax	$y_i = \frac{\exp(x_i)}{\sum_{j=1}^N \exp(x_j)}$	$\mathbf{y} = \arg \max_{\mathbf{z} \in \Delta^{N-1}} \langle \mathbf{z}, \mathbf{x} \rangle + \langle \mathbf{z}, \log \mathbf{z} \rangle$

Optimization-based Implicit Neural Network Design

Layer	Explicit Design	Implicit Design
ReLU	$y = \max\{0, x\}$	$y = \arg \min_{z \geq 0} z - x ^2$
Linear+ReLU	$\mathbf{y} = \max\{\mathbf{0}, \mathbf{W}\mathbf{x} + \mathbf{b}\}$	$\mathbf{y} = \arg \min_{\mathbf{z} \geq \mathbf{0}} \ \mathbf{z} - \mathbf{W}\mathbf{x} - \mathbf{b}\ ^2$
Sigmoid	$y = \frac{1}{1 + \exp(x)}$	$y = \arg \max_{z \in [0,1]} zx + z \log z + (1 - z) \log(1 - z)$
Softmax	$y_i = \frac{\exp(x_i)}{\sum_{j=1}^N \exp(x_j)}$	$\mathbf{y} = \arg \max_{\mathbf{z} \in \Delta^{N-1}} \langle \mathbf{z}, \mathbf{x} \rangle + \langle \mathbf{z}, \log \mathbf{z} \rangle$

- ▶ Design a neural network layer \Leftrightarrow Define an optimization problem + Select a solver
- ▶ Feed-forward computation \Leftrightarrow Solve an optimization problem
- ▶ Backward propagation \Leftrightarrow Adjust hyper-parameters in an optimization problem

Representative Implicit Neural Networks (Layers)

- ▶ Optimization-driven
 - ▶ Input Convex Neural Networks [1]
 - ▶ Deep Declarative Networks [2]
- ▶ Differential equation-driven
 - ▶ Neural ODEs [3]
 - ▶ Deep Equilibrium Models [4]

1 Amos, B., Xu, L., & Kolter, J. Z. Input convex neural networks. ICML, 2017.

2 Gould, S., Hartley, R., & Campbell, D. Deep declarative networks. IEEE TPAMI, 2021.

3 Chen, R. T., Rubanova, Y., Bettencourt, J., & Duvenaud, D. K. Neural ordinary differential equations. NeurIPS, 2018.

4 Bai, S., Kolter, J. Z., & Koltun, V. Deep equilibrium models. NeurIPS, 2019.

5 <https://github.com/bamos/thesis>

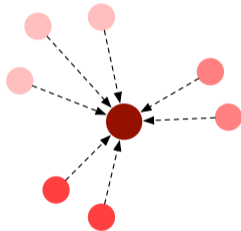
6 <http://implicit-layers-tutorial.org>

Motivation and Principle of Proposed Work

Essentially, many existing NN layers work for information fusion

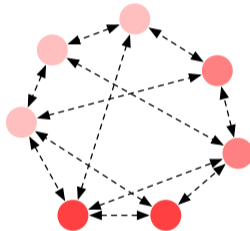
Global Pooling

$$y = f_{\theta}(X)$$



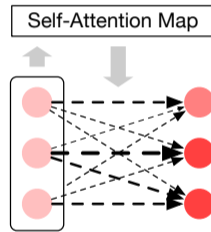
Message-Passing

$$Y = A_{\theta}X$$



Self-Attention

$$Y = A(x, \theta)X$$

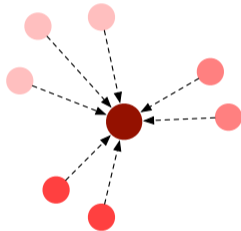


Motivation and Principle of Proposed Work

Essentially, many existing NN layers work for information fusion

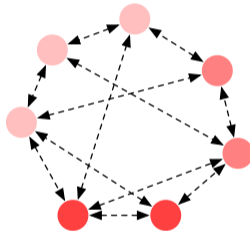
Global Pooling

$$y = f_{\theta}(X)$$



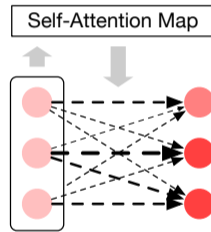
Message-Passing

$$Y = A_{\theta}X$$



Self-Attention

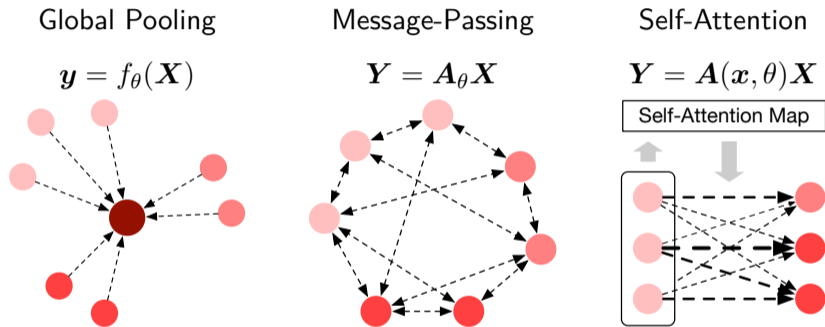
$$Y = A(x, \theta)X$$



What is the design principle of the layers?

Motivation and Principle of Proposed Work

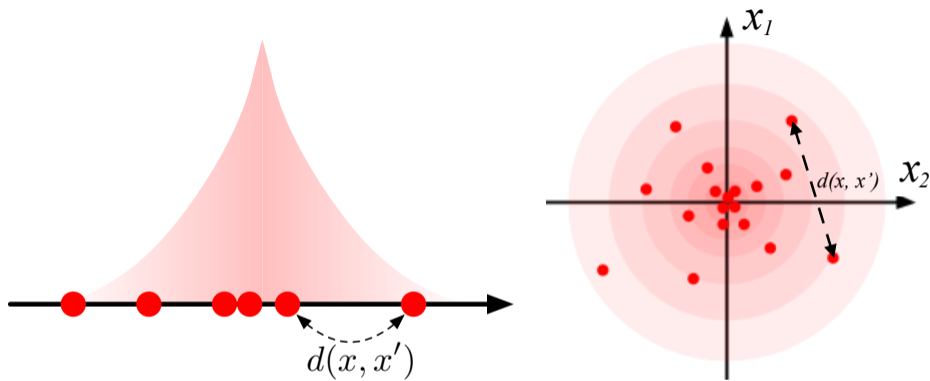
Essentially, many existing NN layers work for information fusion



What is the design principle of the layers?

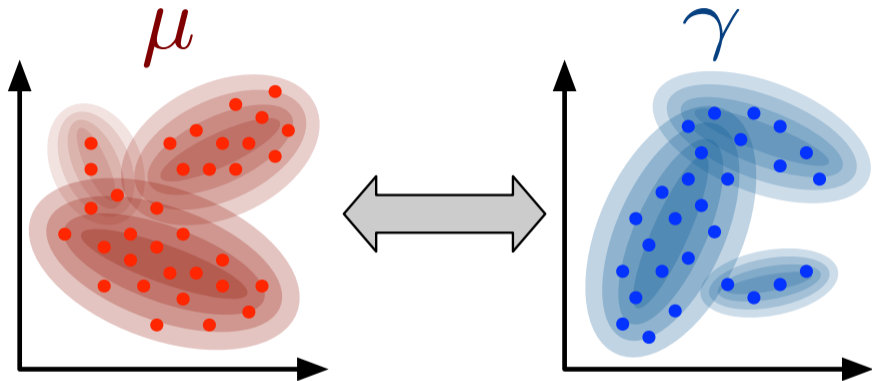
- ▶ Explore the **alignment** principle of information fusion through the lens of OT
- ▶ Develop OT-based implicit surrogates for above layers, improving interpretability and boosting performance

Preliminaries of OT: Metric-Measure Space



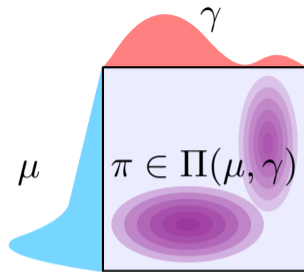
- ▶ $\mathcal{X}_{d,\mu} := (\mathcal{X}, d, \mu)$: A metric-measure space, where $x \in \mathcal{X}$ is a sample in the space.
- ▶ d : A distance metric of samples (e.g., Euclidean distance).
- ▶ \mathbb{P} : A space of (probability) measures defined on \mathcal{X} .
- ▶ $\mu \in \mathbb{P}$: a probability measure on \mathcal{X} .

The Key ML Task: Distribution/Sample Matching



- ▶ Data Clustering, Domain Adaptation, Generative Modeling, Evaluation of Generative Model, ...

Distribution Matching Tool: The Kantorovich-form of OT

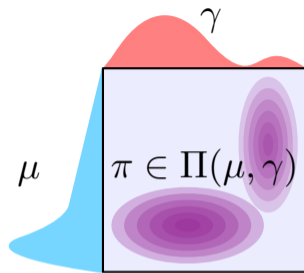


Leonid Kantorovich (1912-1986)

The Kantorovich-form of OT proposed in 1939

- ▶ $\Pi(\mu, \gamma) = \{\pi > 0 \mid \int_x \pi(x, y)dx = \gamma(y), \int_y \pi(x, y)dy = \mu(x)\}$ include all joint distributions taking μ and γ as marginals.
- ▶ $\pi \in \Pi(\mu, \gamma)$ is called **transport plan or coupling**.

Distribution Matching Tool: The Kantorovich-form of OT

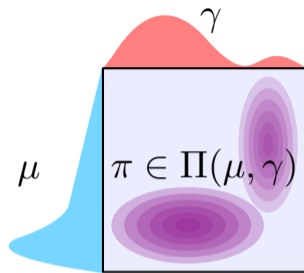


Leonid Kantorovich (1912-1986) The Kantorovich-form of OT proposed in 1939

- ▶ $\Pi(\mu, \gamma) = \{\pi > 0 \mid \int_x \pi(x, y)dx = \gamma(y), \int_y \pi(x, y)dy = \mu(x)\}$ include all joint distributions taking μ and γ as marginals.
- ▶ $\pi \in \Pi(\mu, \gamma)$ is called **transport plan or coupling**.
- ▶ Find an **optimal** transport plan to minimize the expected cost.

$$W_p(\mu, \gamma) := \left(\inf_{\pi \in \Pi(\mu, \gamma)} \int_{(x, y) \in \mathcal{X}^2} d^p(x, y) d\pi(x, y) \right)^{1/p}$$

Distribution Matching Tool: The Kantorovich-form of OT

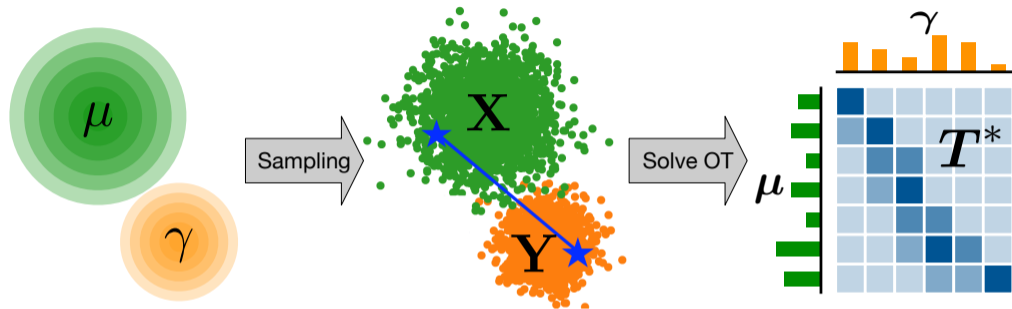


Leonid Kantorovich (1912-1986) The Kantorovich-form of OT proposed in 1939

- ▶ $\Pi(\mu, \gamma) = \{\pi > 0 \mid \int_x \pi(x, y)dx = \gamma(y), \int_y \pi(x, y)dy = \mu(x)\}$ include all joint distributions taking μ and γ as marginals.
- ▶ $\pi \in \Pi(\mu, \gamma)$ is called **transport plan or coupling**.
- ▶ Find an **optimal** transport plan to minimize the expected cost.

$$W_p(\mu, \gamma) := \left(\inf_{\pi \in \Pi(\mu, \gamma)} \int_{(x, y) \in \mathcal{X}^2} d^p(x, y) d\pi(x, y) \right)^{1/p} = \inf_{\pi \in \Pi(\mu, \gamma)} \mathbb{E}_{x, y \sim \pi}^{1/p} [d^p(x, y)]. \quad (1)$$

Empirical OT Problem Defined on Samples



Given $\mathbf{X} = \{x_m\}_{m=1}^M \sim \mu$, $\mathbf{Y} = \{y_n\}_{n=1}^N \sim \gamma$, $\boldsymbol{\mu} \in \Delta^{N-1}$ and $\boldsymbol{\gamma} \in \Delta^{M-1}$,

$$\widehat{W}_p(\mathbf{X}, \mathbf{Y}) := \left(\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\gamma})} \sum_{m=1}^M \sum_{n=1}^N d^p(x_m, y_n) t_{mn} \right)^{1/p} = \left(\min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\gamma})} \langle \mathbf{D}, \mathbf{T} \rangle \right)^{1/p} \quad (2)$$

where $\mathbf{D} = [d^p(x_m, y_n)]$, $\mathbf{T} = [t_{mn}]$, $\Pi(\boldsymbol{\mu}, \boldsymbol{\gamma}) = \{\mathbf{T} > \mathbf{0} | \mathbf{T}\mathbf{1}_M = \boldsymbol{\mu}, \mathbf{T}^\top \mathbf{1}_N = \boldsymbol{\gamma}\}$.

Wasserstein Distance between 1D Samples

- ▶ When $\dim(\mathcal{X}) = 1$, W_p has a closed form, related to **1D histogram transform and equalization**.

$$W_p(\mu, \gamma) = \left(\int_0^1 |F^{-1}(z) - G^{-1}(z)|^p dz \right)^{1/p}, \quad (3)$$

where $F, G : \mathcal{X} \mapsto [0, 1]$ are CDF's of μ and ν .

Wasserstein Distance between 1D Samples

- ▶ When $\dim(\mathcal{X}) = 1$, W_p has a closed form, related to **1D histogram transform and equalization**.

$$W_p(\mu, \gamma) = \left(\int_0^1 |F^{-1}(z) - G^{-1}(z)|^p dz \right)^{1/p}, \quad (3)$$

where $F, G : \mathcal{X} \mapsto [0, 1]$ are CDF's of μ and ν .

- ▶ Given $\mathbf{x} = \{x_n\}_{n=1}^N \sim \mu$ and $\mathbf{y} = \{y_n\}_{n=1}^N \sim \nu$:

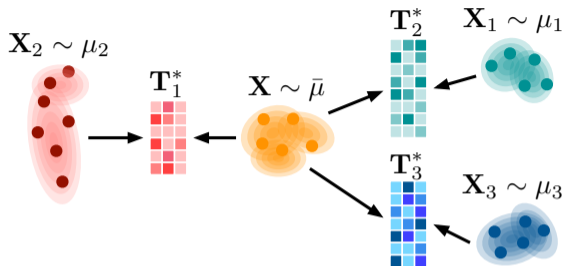
$$\widehat{W}_p(\mathbf{x}, \mathbf{y}) = \left(\sum_{n=1}^N |x_n - y_{\sigma(n)}|^p dz \right)^{1/p}, \quad (4)$$

where σ denotes the **sorting** operation.

- ▶ OT plan \mathbf{T}^* is the permutation matrix corresponding to σ .

Theorem: For one dimensional $x_1 \leq \dots \leq x_N$ and $y_1 \leq \dots \leq y_N$, identity permutation ($\sigma(n) = n$ for $n = 1, \dots, N$) leads to the optimal transport between them.

Wasserstein Barycenters



- ▶ Denote $\mathcal{P}_{\mathcal{X}_d}$ as the space of all probability measures in the metric space \mathcal{X}_d .
- ▶ $(\mathcal{P}_{\mathcal{X}_d}, W_p)$ becomes a metric space of probability measures.
- ▶ Given a set of probability measures $\{\mu_k\}_{k=1}^K \subset \mathcal{P}_{\mathcal{X}_d}$, we can define the **p -Wasserstein barycenter** [Agueh et al, 2021] as

$$\bar{\mu} := \arg \min_{\mu \in \mathcal{P}_{\mathcal{X}_d}} \sum_{k=1}^K W_p^p(\mu, \mu_k). \quad (5)$$

[Agueh et al, 2021] Agueh, M. and Carlier, G., Barycenters in the Wasserstein space. SIAM Journal on Mathematical Analysis, 2011.

Advantages of Optimal Transport

A valid metric for probability measures

- ▶ Apply to distribution comparison and fitting

Advantages of Optimal Transport

A valid metric for probability measures

- ▶ Apply to distribution comparison and fitting

The OT plan/matrix indicates the coherency of sample pairs

- ▶ Apply to point cloud matching and registration
- ▶ Achieve sample averaging and fusion

Advantages of Optimal Transport

A valid metric for probability measures

- ▶ Apply to distribution comparison and fitting

The OT plan/matrix indicates the coherency of sample pairs

- ▶ Apply to point cloud matching and registration
- ▶ Achieve sample averaging and fusion

Have potentials for designing implicit neural network and achieving interpretable information fusion.

Case 1: OT-based Implicit Global Pooling Design

Given $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$, find a global pooling $f : \mathcal{X} \mapsto \mathbb{R}^D$.

Case 1: OT-based Implicit Global Pooling Design

Given $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$, find a global pooling $f : \mathcal{X} \mapsto \mathbb{R}^D$.

- ▶ Mean-pooling

$$f(\mathbf{X}) = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \quad (6)$$

- ▶ Max-pooling

$$f(\mathbf{X}) = \left\|_{d=1}^D \max_n \{x_{dn}\}_{n=1}^N \right. \quad (7)$$

- ▶ Attention-pooling

$$f(\mathbf{X}) = \mathbf{X} \mathbf{a}_X \quad (8)$$

Case 1: OT-based Implicit Global Pooling Design

Given $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$, find a global pooling $f : \mathcal{X} \mapsto \mathbb{R}^D$.

- ▶ Mean-pooling

$$f(\mathbf{X}) = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \quad (6)$$

- ▶ Max-pooling

$$f(\mathbf{X}) = \left\|_{d=1}^D \max_n \{x_{dn}\}_{n=1}^N \right. \quad (7)$$

- ▶ Attention-pooling

$$f(\mathbf{X}) = \mathbf{X} \mathbf{a}_X \quad (8)$$

- ▶ The design and selection of f is **empirical and sub-optimal**.
- ▶ Can we build a generalized framework covering the above designs?

Case 1: OT-based Implicit Global Pooling Design

Determine the poolings by a “sample-feature dimension” distribution $\mathbf{P} = [p_{dn}]$.

Case 1: OT-based Implicit Global Pooling Design

Determine the poolings by a “sample-feature dimension” distribution $\mathbf{P} = [p_{dn}]$.

- ▶ A global pooling outputs the expectation of samples conditioned on different feature dimensions

$$f(\mathbf{X}) = \prod_{d=1}^D \mathbb{E}_{n \sim p_{n|d}} [x_{dn}]$$

Case 1: OT-based Implicit Global Pooling Design

Determine the poolings by a “sample-feature dimension” distribution $\mathbf{P} = [p_{dn}]$.

- ▶ A global pooling outputs the expectation of samples conditioned on different feature dimensions

$$f(\mathbf{X}) = \prod_{d=1}^D \mathbb{E}_{n \sim p_{n|d}} [x_{dn}] = (\mathbf{X} \odot \underbrace{\text{diag}^{-1}(\overbrace{\mathbf{P}\mathbf{1}_N}^{p=[p_d]})\mathbf{P}}_{\tilde{\mathbf{P}}=[p_{n|d}]})\mathbf{1}_N. \quad (9)$$

Case 1: OT-based Implicit Global Pooling Design

Determine the poolings by a “sample-feature dimension” distribution $\mathbf{P} = [p_{dn}]$.

- ▶ A global pooling outputs **the expectation of samples conditioned on different feature dimensions**

$$f(\mathbf{X}) = \prod_{d=1}^D \mathbb{E}_{n \sim p_{n|d}} [x_{dn}] = (\mathbf{X} \odot \underbrace{\text{diag}^{-1}(\overbrace{\mathbf{P}\mathbf{1}_N}^{p=[p_d]})\mathbf{P}}_{\tilde{\mathbf{P}}=[p_{n|d}]})\mathbf{1}_N. \quad (9)$$

Mean-pooling	$\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$	\Leftrightarrow	$\mathbf{P} = [\frac{1}{DN}]$
Max-pooling	$\prod_{d=1}^D \max_n \{x_{dn}\}_{n=1}^N$	\Leftrightarrow	$\mathbf{P} \in \{0, \frac{1}{D}\}^{D \times N}$ and $\mathbf{P}\mathbf{1}_N = \frac{1}{D}\mathbf{1}_D$
Attention-pooling	$\mathbf{X}\mathbf{a}_X$	\Leftrightarrow	$\mathbf{P} = \frac{1}{D}\mathbf{1}_D\mathbf{a}_X^\top$

Case 1: OT-based Implicit Global Pooling Design

Determine the poolings by a “sample-feature dimension” distribution $\mathbf{P} = [p_{dn}]$.

- ▶ A global pooling outputs **the expectation of samples conditioned on different feature dimensions**

$$f(\mathbf{X}) = \prod_{d=1}^D \mathbb{E}_{n \sim p_{n|d}} [x_{dn}] = (\mathbf{X} \odot \underbrace{\text{diag}^{-1}(\overbrace{\mathbf{P}\mathbf{1}_N}^{p=[p_d]})}_{\tilde{\mathbf{P}}=[p_{n|d}]}) \mathbf{1}_N. \quad (9)$$

Mean-pooling	$\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$	\Leftrightarrow	$\mathbf{P} = [\frac{1}{DN}]$
Max-pooling	$\prod_{d=1}^D \max_n \{x_{dn}\}_{n=1}^N$	\Leftrightarrow	$\mathbf{P} \in \{0, \frac{1}{D}\}^{D \times N}$ and $\mathbf{P}\mathbf{1}_N = \frac{1}{D}\mathbf{1}_D$
Attention-pooling	$\mathbf{X}\mathbf{a}_X$	\Leftrightarrow	$\mathbf{P} = \frac{1}{D}\mathbf{1}_D\mathbf{a}_X^\top$

Design global pooling is equivalent to design \mathbf{P} , leading to an OT-based implicit layer.

Regularized Optimal Transport Pooling (ROTP) Layer

$$f_{\text{rot}}(\mathbf{X}; \boldsymbol{\theta}) = (\mathbf{X} \odot \text{diag}^{-1}(\mathbf{P}_{\text{rot}}^* \mathbf{1}_N) \mathbf{P}_{\text{rot}}^*) \mathbf{1}_N, \quad (10)$$

where

$$\mathbf{P}_{\text{rot}}^* = \arg \max_{\mathbf{P} \in \Omega} \underbrace{\langle \mathbf{X}, \mathbf{P} \rangle}_{\text{OT term}} + \underbrace{\alpha_0 \langle \mathbf{C}_D \mathbf{P} \mathbf{C}_N, \mathbf{P} \rangle}_{\text{Structural Reg.}} - \underbrace{\alpha_1 \mathbf{R}(\mathbf{P})}_{\text{Smoothness}} - \underbrace{\alpha_2 \text{KL}(\mathbf{P} \mathbf{1} | \mathbf{p}_0) - \alpha_3 \text{KL}(\mathbf{P}^\top \mathbf{1} | \mathbf{q}_0)}_{\text{Marginal Prior}}$$

Regularized Optimal Transport Pooling (ROTP) Layer

$$f_{\text{rot}}(\mathbf{X}; \boldsymbol{\theta}) = (\mathbf{X} \odot \text{diag}^{-1}(\mathbf{P}_{\text{rot}}^* \mathbf{1}_N) \mathbf{P}_{\text{rot}}^*) \mathbf{1}_N, \quad (10)$$

where

$$\mathbf{P}_{\text{rot}}^* = \arg \max_{\mathbf{P} \in \Omega} \underbrace{\langle \mathbf{X}, \mathbf{P} \rangle}_{\text{OT term}} + \underbrace{\alpha_0 \langle \mathbf{C}_D \mathbf{P} \mathbf{C}_N, \mathbf{P} \rangle}_{\text{Structural Reg.}} - \underbrace{\alpha_1 \mathbf{R}(\mathbf{P})}_{\text{Smoothness}} - \underbrace{\alpha_2 \text{KL}(\mathbf{P} \mathbf{1} | \mathbf{p}_0) - \alpha_3 \text{KL}(\mathbf{P}^\top \mathbf{1} | \mathbf{q}_0)}_{\text{Marginal Prior}}$$

Fused Gromov-Wasserstein

Implementation	Principle	Connections to OT
$\max \langle \mathbf{X}, \mathbf{P} \rangle$	EM Principle (Max Output Energy)	OT loss

Regularized Optimal Transport Pooling (ROTP) Layer

$$f_{\text{rot}}(\mathbf{X}; \boldsymbol{\theta}) = (\mathbf{X} \odot \text{diag}^{-1}(\mathbf{P}_{\text{rot}}^* \mathbf{1}_N) \mathbf{P}_{\text{rot}}^*) \mathbf{1}_N, \quad (10)$$

where

$$\mathbf{P}_{\text{rot}}^* = \arg \max_{\mathbf{P} \in \Omega} \underbrace{\langle \mathbf{X}, \mathbf{P} \rangle}_{\text{OT term}} + \underbrace{\alpha_0 \langle \mathbf{C}_D \mathbf{P} \mathbf{C}_N, \mathbf{P} \rangle}_{\text{Structural Reg.}} - \underbrace{\alpha_1 \mathbf{R}(\mathbf{P})}_{\text{Smoothness}} - \underbrace{\alpha_2 \text{KL}(\mathbf{P} \mathbf{1} | \mathbf{p}_0) - \alpha_3 \text{KL}(\mathbf{P}^\top \mathbf{1} | \mathbf{q}_0)}_{\text{Marginal Prior}}$$

Fused Gromov-Wasserstein

Implementation	Principle	Connections to OT
$\max \langle \mathbf{X}, \mathbf{P} \rangle$	EM Principle (Max Output Energy)	OT loss
$+ \max \langle \mathbf{C}_D \mathbf{P} \mathbf{C}_N, \mathbf{P} \rangle$	Max Sample-Feature Correlation	FGW loss

Regularized Optimal Transport Pooling (ROTP) Layer

$$f_{\text{rot}}(\mathbf{X}; \boldsymbol{\theta}) = (\mathbf{X} \odot \text{diag}^{-1}(\mathbf{P}_{\text{rot}}^* \mathbf{1}_N) \mathbf{P}_{\text{rot}}^*) \mathbf{1}_N, \quad (10)$$

where

$$\mathbf{P}_{\text{rot}}^* = \arg \max_{\mathbf{P} \in \Omega} \underbrace{\langle \mathbf{X}, \mathbf{P} \rangle}_{\text{OT term}} + \underbrace{\alpha_0 \langle \mathbf{C}_D \mathbf{P} \mathbf{C}_N, \mathbf{P} \rangle}_{\text{Structural Reg.}} - \underbrace{\alpha_1 \text{R}(\mathbf{P})}_{\text{Smoothness}} - \underbrace{\alpha_2 \text{KL}(\mathbf{P} \mathbf{1} | \mathbf{p}_0) - \alpha_3 \text{KL}(\mathbf{P}^\top \mathbf{1} | \mathbf{q}_0)}_{\text{Marginal Prior}}$$

Fused Gromov-Wasserstein

Implementation	Principle	Connections to OT
$\max \langle \mathbf{X}, \mathbf{P} \rangle$	EM Principle (Max Output Energy)	OT loss
$+ \max \langle \mathbf{C}_D \mathbf{P} \mathbf{C}_N, \mathbf{P} \rangle$	Max Sample-Feature Correlation	FGW loss
$+ \min \text{R}(\mathbf{P})$	Avoid Over-sparse Distributions	Entropic/Quadratic FGW

Regularized Optimal Transport Pooling (ROTP) Layer

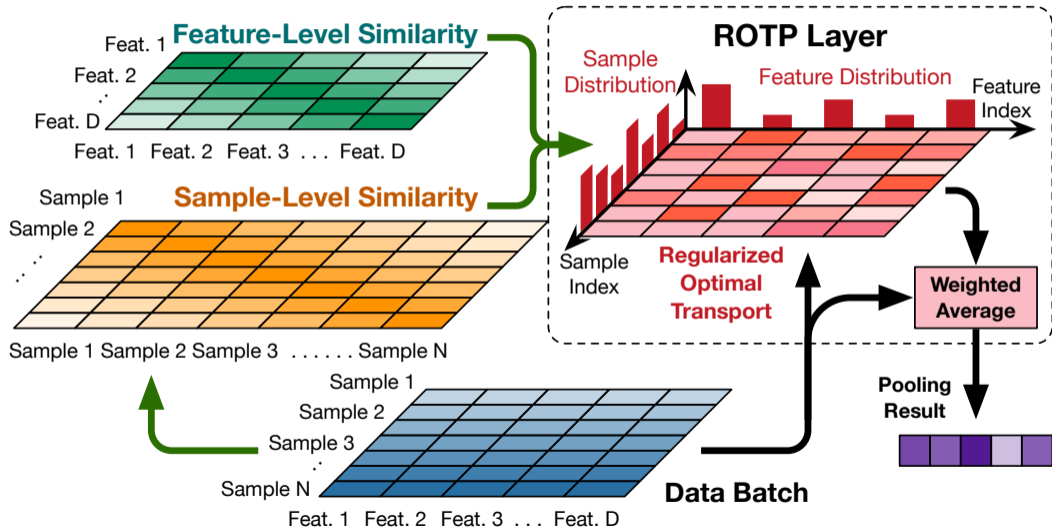
$$f_{\text{rot}}(\mathbf{X}; \boldsymbol{\theta}) = (\mathbf{X} \odot \text{diag}^{-1}(\mathbf{P}_{\text{rot}}^* \mathbf{1}_N) \mathbf{P}_{\text{rot}}^*) \mathbf{1}_N, \quad (10)$$

where

$$\mathbf{P}_{\text{rot}}^* = \arg \max_{\mathbf{P} \in \Omega} \underbrace{\langle \mathbf{X}, \mathbf{P} \rangle}_{\text{OT term}} + \underbrace{\alpha_0 \langle \mathbf{C}_D \mathbf{P} \mathbf{C}_N, \mathbf{P} \rangle}_{\text{Structural Reg.}} - \underbrace{\alpha_1 R(\mathbf{P})}_{\text{Smoothness}} - \underbrace{\alpha_2 \text{KL}(\mathbf{P} \mathbf{1} | \mathbf{p}_0) - \alpha_3 \text{KL}(\mathbf{P}^\top \mathbf{1} | \mathbf{q}_0)}_{\text{Marginal Prior}}$$

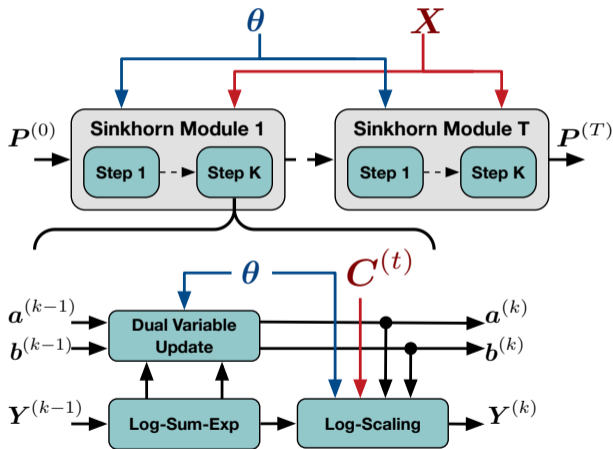
Implementation	Principle	Connections to OT
$\max \langle \mathbf{X}, \mathbf{P} \rangle$	EM Principle (Max Output Energy)	OT loss
$+ \max \langle \mathbf{C}_D \mathbf{P} \mathbf{C}_N, \mathbf{P} \rangle$	Max Sample-Feature Correlation	FGW loss
$+ \min R(\mathbf{P})$	Avoid Over-sparse Distributions	Entropic/Quadratic FGW
$+ \min \text{KL}(\mathbf{P} \mathbf{1} \mathbf{p}_0)$	Leverage Prior Knowledge	Unbalanced E/Q FGW
$+ \min \text{KL}(\mathbf{P}^\top \mathbf{1} \mathbf{q}_0)$		

Regularized Optimal Transport Pooling (ROTP) Layer



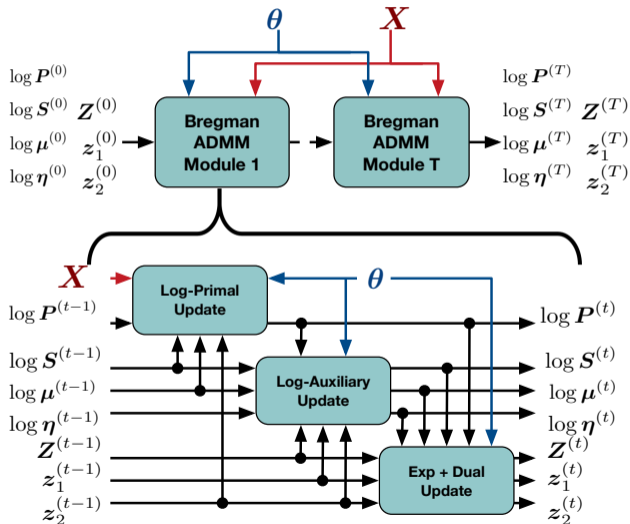
Regularized Optimal Transport Pooling (ROTP) Layer

Implement ROTP by **Unrolling Iterative Algorithms** (e.g., Sinkhorn-scaling)



Regularized Optimal Transport Pooling (ROTP) Layer

Implement ROTP by **Unrolling Iterative Algorithms** (e.g., Bregman ADMM)

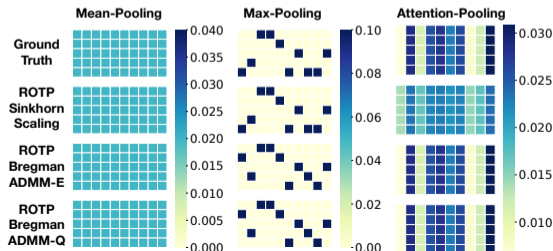


ROTP: Generalizability

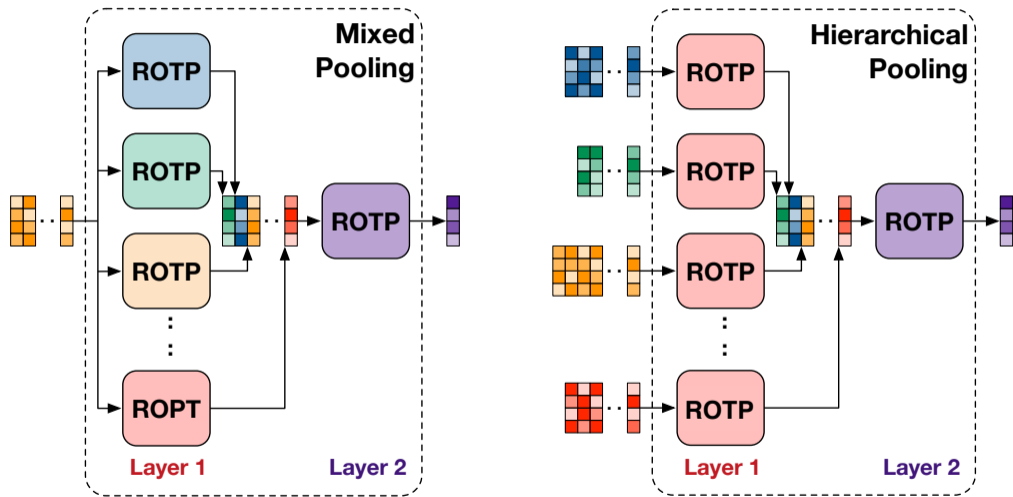
$$\mathbf{P}_{\text{rot}}^* = \arg \max_{\mathbf{P} \in \Omega} \underbrace{\langle \mathbf{X}, \mathbf{P} \rangle}_{\text{OT term}} + \underbrace{\alpha_0 \langle \mathbf{C}_D \mathbf{P} \mathbf{C}_N, \mathbf{P} \rangle}_{\text{Structural Reg.}} - \underbrace{\alpha_1 \mathbf{R}(\mathbf{P})}_{\text{Smoothness}} - \underbrace{\alpha_2 \text{KL}(\mathbf{P} \mathbf{1} | \mathbf{p}_0) - \alpha_3 \text{KL}(\mathbf{P}^\top \mathbf{1} | \mathbf{q}_0)}_{\text{Marginal Prior}}$$

Fused Gromov-Wasserstein

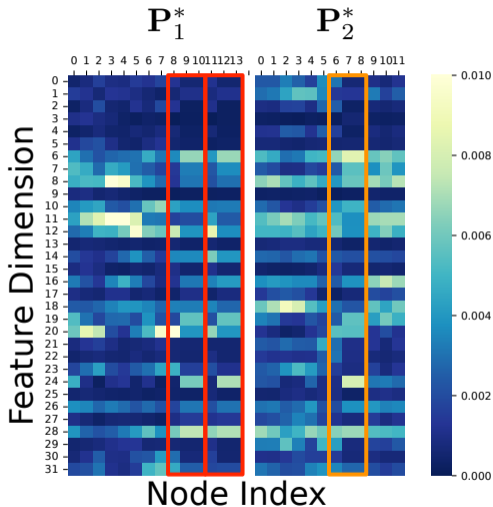
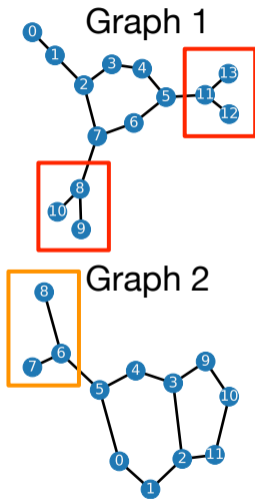
$f_{\text{rot}}(\mathbf{X}; \boldsymbol{\theta})$	α_0	α_1	α_2	α_3	\mathbf{p}_0	\mathbf{q}_0
Mean-pooling	0	$\rightarrow \infty$	$\rightarrow \infty$	$\rightarrow \infty$	$\frac{1}{D} \mathbf{1}_D$	$\frac{1}{N} \mathbf{1}_N$
Max-pooling	0	0	$\rightarrow \infty$	0	$\frac{1}{D} \mathbf{1}_D$	—
Attention-pooling	0	$\rightarrow \infty$	$\rightarrow \infty$	$\rightarrow \infty$	$\frac{1}{D} \mathbf{1}_D$	\mathbf{a}_X



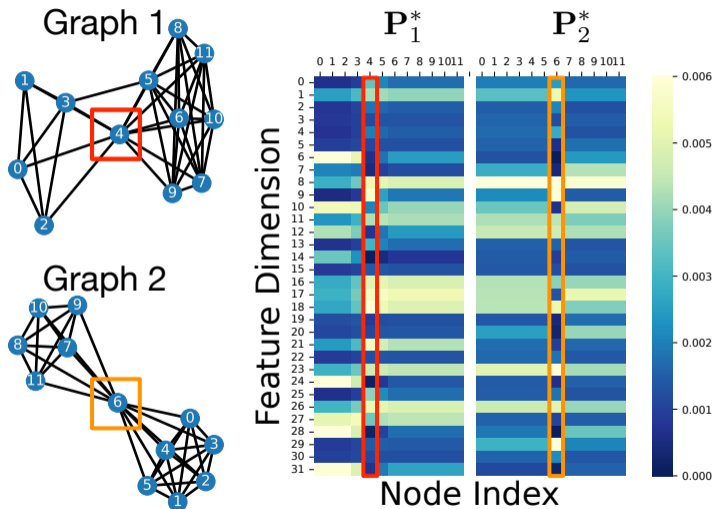
ROTP: Flexibility



ROTP: Interpretability



ROTP: Interpretability



ROTP: Superiority on Data Fitting

Dataset	NCII	PROTEINS	MUTAG	COLLAB	RDT-B	RDT-M5K	IMDB-B	IMDB-M
#Graphs	4,110	1,113	188	5,000	2,000	4,999	1,000	1,500
Average #Nodes	29.87	39.06	17.93	74.49	429.63	508.52	19.77	13.00
Average #Edges	32.30	72.82	19.79	2,457.78	497.75	594.87	96.53	65.94
#Classes	2	2	2	3	2	5	2	3
Add	67.96 \pm 0.43	72.97 \pm 0.54	89.05\pm0.86	71.06 \pm 0.43	80.00 \pm 1.49	50.16 \pm 0.97	70.18 \pm 0.87	47.56 \pm 0.56
Mean	64.82 \pm 0.52	66.09 \pm 0.64	86.53 \pm 1.62	72.35 \pm 0.44	83.62 \pm 1.18	52.44 \pm 1.24	70.34 \pm 0.38	48.65 \pm 0.91
Max	65.95 \pm 0.76	72.27 \pm 0.33	85.90 \pm 1.68	73.07 \pm 0.57	82.62 \pm 1.25	44.34 \pm 1.93	70.24 \pm 0.54	47.80 \pm 0.54
DeepSet	66.28 \pm 0.72	73.76\pm0.47	87.84 \pm 0.71	69.74 \pm 0.66	82.91 \pm 1.37	47.45 \pm 0.54	70.84 \pm 0.71	48.05 \pm 0.71
Mixed	66.46 \pm 0.74	72.25 \pm 0.45	87.30 \pm 0.87	73.22 \pm 0.35	84.36 \pm 2.62	46.67 \pm 1.63	71.28 \pm 0.26	48.07 \pm 0.25
GatedMixed	63.86 \pm 0.76	69.40 \pm 1.93	87.94 \pm 1.28	71.94 \pm 0.40	80.60 \pm 3.89	44.78 \pm 4.53	70.96 \pm 0.60	48.09 \pm 0.44
Set2Set	65.10 \pm 1.12	68.61 \pm 1.44	87.77 \pm 0.86	72.31 \pm 0.73	80.08 \pm 5.72	49.85 \pm 2.77	70.36 \pm 0.85	48.30 \pm 0.54
Attention	64.35 \pm 0.61	67.70 \pm 0.95	88.08 \pm 1.22	72.57 \pm 0.41	81.55 \pm 4.39	51.85 \pm 0.66	70.60 \pm 0.38	47.83 \pm 0.78
GatedAtt	64.66 \pm 0.52	68.16 \pm 0.90	86.91 \pm 1.79	72.31 \pm 0.37	82.55 \pm 1.96	51.47 \pm 0.82	70.52 \pm 0.31	48.67 \pm 0.35
DynamicP	62.11 \pm 0.27	65.86 \pm 0.85	85.40 \pm 2.81	70.78 \pm 0.88	67.51 \pm 1.82	32.11 \pm 3.85	69.84 \pm 0.73	47.59 \pm 0.48
GNP	68.20\pm0.48	73.44\pm0.61	88.37 \pm 1.25	72.80 \pm 0.58	81.93 \pm 2.23	51.80 \pm 0.61	70.34 \pm 0.83	48.85 \pm 0.81
ASAP	68.09 \pm 0.42	70.42 \pm 1.45	87.68 \pm 1.42	68.20 \pm 2.37	73.91 \pm 1.50	44.58 \pm 0.44	68.33 \pm 2.50	43.92 \pm 1.13
SAGP	67.48 \pm 0.65	72.63 \pm 0.44	87.88 \pm 2.22	70.19 \pm 0.55	74.12 \pm 2.86	46.00 \pm 1.74	70.34 \pm 0.74	47.04 \pm 1.22
OTK	67.96 \pm 0.55	69.52 \pm 0.76	86.90 \pm 1.83	71.35 \pm 0.91	74.28 \pm 1.39	50.57 \pm 1.20	70.94 \pm 0.79	48.41 \pm 0.89
SWE	68.06 \pm 0.98	70.09 \pm 1.22	85.68 \pm 2.07	72.17 \pm 1.29	79.30 \pm 3.94	51.11 \pm 1.55	70.34 \pm 1.05	48.93 \pm 1.34
WEGL	68.16\pm0.62	71.58 \pm 0.94	88.68\pm1.66	72.55 \pm 0.69	82.80 \pm 1.73	52.03 \pm 0.60	71.94\pm0.75	49.20 \pm 0.87
ROTP _S	68.27\pm1.06	73.10\pm0.22	88.84\pm1.21	71.20 \pm 0.55	81.54 \pm 1.38	51.00 \pm 0.61	70.74 \pm 0.80	47.87 \pm 0.43
ROTP _{B-E} ($\alpha_0 = 0$)	66.23 \pm 0.50	67.71 \pm 1.70	86.82 \pm 2.02	73.86 \pm 0.44	86.80\pm1.19	52.25 \pm 0.75	71.72 \pm 0.88	50.48\pm0.14
ROTP _{B-Q} ($\alpha_0 = 0$)	66.18 \pm 0.76	69.88 \pm 0.87	85.42 \pm 1.10	74.14\pm0.24	87.72\pm1.03	52.79\pm0.60	72.34\pm0.50	49.36 \pm 0.52
ROTP _{B-E} (learn α_0)	65.90 \pm 0.94	70.19 \pm 0.66	88.01 \pm 1.51	74.05\pm0.34	86.78 \pm 1.14	52.77\pm0.69	71.76 \pm 0.62	50.28\pm0.86
ROTP _{B-Q} (learn α_0)	65.96 \pm 0.32	70.12 \pm 1.17	86.79 \pm 1.81	74.27\pm0.47	88.67\pm0.99	52.84\pm0.60	71.78\pm1.00	49.44\pm0.46

ROTP: Superiority on Data Fitting

Dataset	Messidor	Component	Function
D	687	200	200
#Positive bags	654	423	443
#Negative bags	546	2,707	4,799
#Instances	12,352	36,894	55,536
Min. bag size	8	1	1
Max. bag size	12	53	51
Add	74.33 \pm 2.56	93.35 \pm 0.98	96.26 \pm 0.48
Mean	74.42 \pm 2.47	93.32 \pm 0.99	96.28 \pm 0.66
Max	73.92 \pm 3.00	93.23 \pm 0.76	95.94 \pm 0.48
DeepSet	74.42 \pm 2.87	93.29 \pm 0.95	96.45 \pm 0.51
Mixed	73.42 \pm 2.29	93.45\pm0.61	96.41 \pm 0.53
GatedMixed	73.25 \pm 2.38	93.03 \pm 1.02	96.22 \pm 0.65
Set2Set	73.58 \pm 3.74	93.19 \pm 0.95	96.43 \pm 0.56
Attention	74.25 \pm 3.67	93.22 \pm 1.02	96.31 \pm 0.66
GatedAtt	73.67 \pm 2.23	93.42\pm0.91	96.51\pm0.77
DynamicP	73.16 \pm 2.12	93.26 \pm 1.30	96.47\pm0.58
GNP	73.54 \pm 3.68	92.86 \pm 1.96	96.10 \pm 1.03
OTK	74.78 \pm 2.89	93.19 \pm 0.93	96.31 \pm 1.02
SWE	74.46 \pm 3.72	93.32 \pm 1.26	96.42 \pm 0.88
ROTP _S	75.42\pm2.96	93.29 \pm 0.83	96.62\pm0.48
ROTP _{B-E} ($\alpha_0 = 0$)	74.83 \pm 2.07	93.16 \pm 1.02	96.17 \pm 0.43
ROTP _{B-Q} ($\alpha_0 = 0$)	75.08 \pm 2.06	93.13 \pm 0.94	96.09 \pm 0.46
ROTP _{B-E} (learn α_0)	75.33\pm1.96	93.16 \pm 1.08	96.22 \pm 0.44
ROTP _{B-Q} (learn α_0)	75.17\pm2.45	93.45\pm0.96	96.22 \pm 0.48

Dataset	DECAGON DiBr-APND	DECAGON Anae-Fati	DECAGON PleuP-Diar	FEARS
#Graph sets	6,309	2,922	2,842	6,338
#Positive sets	3,189	1,526	1,422	3,169
Positive label	Difficulty breathing	Anaemia	Pleural pain	Non- myopathy
#Negative sets	3,120	1,396	1,420	3,169
Negative label	Pressure decreased	Fatigue	Diarrhea	Myopathy
Set size	2	2	2	2~52
Add	50.86 \pm 0.97	63.15 \pm 1.79	62.32\pm1.08	75.89 \pm 1.33
Mean	51.10 \pm 1.09	61.95 \pm 2.60	61.30 \pm 2.68	72.42 \pm 1.51
Max	50.59 \pm 0.77	61.88 \pm 2.03	60.11 \pm 2.03	82.02\pm0.72
DeepSet	49.83 \pm 1.07	56.24 \pm 5.20	51.78 \pm 3.10	82.40\pm1.56
Mixed	51.13 \pm 0.99	63.83\pm1.19	60.91 \pm 2.12	81.54 \pm 1.13
GatedMixed	51.39 \pm 0.63	61.50 \pm 1.61	59.12 \pm 2.12	81.88 \pm 1.14
Set2Set	50.72 \pm 1.71	59.35 \pm 2.04	55.01 \pm 3.59	79.29 \pm 0.84
Attention	50.52 \pm 1.10	61.40 \pm 2.03	61.33 \pm 2.40	75.98 \pm 0.74
GatedAtt	50.74 \pm 0.61	62.15 \pm 0.77	58.80 \pm 1.18	75.84 \pm 1.29
DynamicP	51.01 \pm 1.88	55.93 \pm 1.56	52.58 \pm 2.91	74.00 \pm 1.61
GNP	50.00 \pm 1.88	53.98 \pm 6.34	52.58 \pm 4.68	62.71 \pm 15.55
ASAP	50.89 \pm 0.82	63.66 \pm 1.81	60.67 \pm 2.69	77.15 \pm 1.13
SAGP	49.87 \pm 0.77	63.62 \pm 1.28	59.86 \pm 2.43	77.29 \pm 1.04
OTK	50.96 \pm 1.11	63.68 \pm 1.59	61.66\pm2.39	79.40 \pm 1.08
SWE	51.05 \pm 2.15	63.21 \pm 2.02	61.37 \pm 3.13	80.64 \pm 1.86
WEGL	51.67\pm0.85	63.79\pm2.54	61.36 \pm 2.30	81.98 \pm 0.77
ROTP _S	51.96\pm0.71	62.91 \pm 1.13	59.40 \pm 0.90	79.75 \pm 0.71
ROTP _{B-E}	51.26 \pm 0.84	63.86\pm2.41	62.57\pm1.34	82.55\pm0.42
ROTP _{B-Q}	52.72\pm0.66	63.15 \pm 1.27	60.88 \pm 1.65	81.43 \pm 1.12

Case 2: OT-based Implicit Message-Passing Design

Given a graph $G(\mathcal{V}, \mathcal{E}, \mathbf{A}, \mathbf{X})$, where $\mathbf{Y}_{\mathcal{V}_L}$ are labels of partial nodes.

► **Traditional Node-Level Learning Paradigm:**

$$\max_{\theta} \prod_{v \in \mathcal{V}_L} p(\mathbf{y}_v | \mathbf{X}, \mathbf{A}; \theta) \Leftrightarrow \min_{\theta} \sum_{v \in \mathcal{V}_L} \psi(g_v(\mathbf{X}, \mathbf{A}; \theta), \mathbf{y}_v). \quad (11)$$

- p is Gaussian $\Leftrightarrow \psi$ is MSE.
- p is Sigmoid/Softmax $\Leftrightarrow \psi$ is cross entropy loss.

Case 2: OT-based Implicit Message-Passing Design

Given a graph $G(\mathcal{V}, \mathcal{E}, \mathbf{A}, \mathbf{X})$, where $\mathbf{Y}_{\mathcal{V}_L}$ are labels of partial nodes.

▶ **Traditional Node-Level Learning Paradigm:**

$$\max_{\theta} \prod_{v \in \mathcal{V}_L} p(\mathbf{y}_v | \mathbf{X}, \mathbf{A}; \theta) \Leftrightarrow \min_{\theta} \sum_{v \in \mathcal{V}_L} \psi(g_v(\mathbf{X}, \mathbf{A}; \theta), \mathbf{y}_v). \quad (11)$$

▶ p is Gaussian $\Leftrightarrow \psi$ is MSE.

▶ p is Sigmoid/Softmax $\Leftrightarrow \psi$ is cross entropy loss.

▶ **What we really want to do is maximizing the probability of $\mathbf{Y}_{\mathcal{V}_L}$, i.e.,**

$$\max_{\theta} p(\mathbf{Y}_{\mathcal{V}_L} | \mathbf{X}, \mathbf{A}; \theta) \quad (12)$$

Case 2: OT-based Implicit Message-Passing Design

Given a graph $G(\mathcal{V}, \mathcal{E}, \mathbf{A}, \mathbf{X})$, where $\mathbf{Y}_{\mathcal{V}_L}$ are labels of partial nodes.

- ▶ **Traditional Node-Level Learning Paradigm:**

$$\max_{\theta} \prod_{v \in \mathcal{V}_L} p(\mathbf{y}_v | \mathbf{X}, \mathbf{A}; \theta) \Leftrightarrow \min_{\theta} \sum_{v \in \mathcal{V}_L} \psi(g_v(\mathbf{X}, \mathbf{A}; \theta), \mathbf{y}_v). \quad (11)$$

- ▶ p is Gaussian $\Leftrightarrow \psi$ is MSE.
- ▶ p is Sigmoid/Softmax $\Leftrightarrow \psi$ is cross entropy loss.
- ▶ **What we really want to do is maximizing the probability of $\mathbf{Y}_{\mathcal{V}_L}$, i.e.,**

$$\max_{\theta} p(\mathbf{Y}_{\mathcal{V}_L} | \mathbf{X}, \mathbf{A}; \theta) \quad (12)$$

- ▶ **The above two equations are equivalent iff the labels are conditional independent, which is questionable in practice.**

$$p(\mathbf{y}_v | \mathbf{x}_v) \neq p(\mathbf{y}_v | \mathbf{x}_v, \mathbf{y}_{v'}),$$

$$p(\mathbf{y}_v, \mathbf{y}_{v'} | \mathbf{x}_v, \mathbf{x}_{v'}) = p(\mathbf{y}_v | \mathbf{x}_v, \mathbf{x}_{v'}, \mathbf{y}_{v'}) p(\mathbf{y}_{v'} | \mathbf{x}_v, \mathbf{x}_{v'}) \neq p(\mathbf{y}_v | \mathbf{x}_v, \mathbf{x}_{v'}) p(\mathbf{y}_{v'} | \mathbf{x}_v, \mathbf{x}_{v'}),$$

$$p(\mathbf{Y}_{\mathcal{V}} | \mathbf{X}, \mathbf{A}) \neq \prod_{v \in \mathcal{V}} p(\mathbf{y}_v | \mathbf{X}, \mathbf{A}).$$

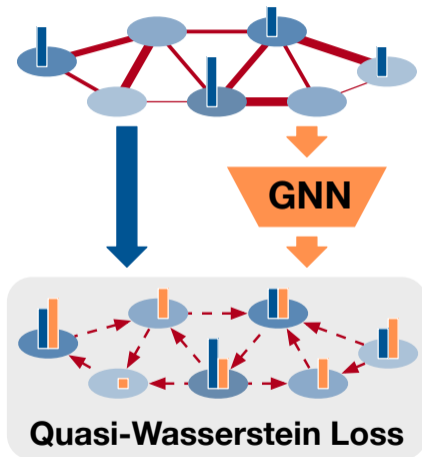
Proposed Learning Strategy

Ideally, we would like to define a loss for label sets, rather than individual labels, e.g.,

$$\min_{\theta} \text{Loss}(\underbrace{g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta)}_{\hat{\mathcal{Y}}_{\mathcal{V}_L}}, \mathbf{Y}_{\mathcal{V}_L}), \quad (13)$$

which can be implemented as an **OT problem of labels defined on a graph**.

- ▶ **Moreover, lead to an implicit message-passing layer (with a QW loss) encoding label transport.**



Proposed Quasi-Wasserstein (QW) Loss of GNNs

- ▶ 1-Wasserstein distance between signals on a graph

$$W_1(\boldsymbol{\mu}, \boldsymbol{\gamma}) := \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\gamma})} \langle \mathbf{D}, \mathbf{T} \rangle = \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\gamma})} \sum_{v, v' \in \mathcal{V} \times \mathcal{V}} t_{vv'} d_{vv'}, \quad (14)$$

where \mathbf{D} is the shortest-path distance matrix.

Proposed Quasi-Wasserstein (QW) Loss of GNNs

- ▶ 1-Wasserstein distance between signals on a graph

$$W_1(\boldsymbol{\mu}, \boldsymbol{\gamma}) := \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\gamma})} \langle \mathbf{D}, \mathbf{T} \rangle = \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\gamma})} \sum_{v, v' \in \mathcal{V} \times \mathcal{V}} t_{vv'} d_{vv'}, \quad (14)$$

where \mathbf{D} is the shortest-path distance matrix.

- ▶ W_1 is equivalent to a minimum-cost flow problem:

$$W_1(\boldsymbol{\mu}, \boldsymbol{\gamma}) = \min_{\mathbf{f} \in \Omega(\mathbf{S}_{\mathcal{V}}, \boldsymbol{\mu}, \boldsymbol{\gamma})} \|\text{diag}(\mathbf{w}) \mathbf{f}\|_1, \quad (15)$$

where $\mathbf{f} = [f_e] \in \mathbb{R}^{|\mathcal{E}|}$ is the flow on each edge.

$$\mathbf{S}_{\mathcal{V}} = [s_{ve}] \in \{0, \pm 1\}^{|\mathcal{V}| \times |\mathcal{E}|}, \quad s_{ve} = \begin{cases} 1 & \text{if } v \text{ is "head" of edge } e \\ -1 & \text{if } v \text{ is "tail" of edge } e \\ 0 & \text{otherwise.} \end{cases} \quad (16)$$

Proposed Quasi-Wasserstein (QW) Loss of GNNs

- ▶ 1-Wasserstein distance between signals on a graph

$$W_1(\boldsymbol{\mu}, \boldsymbol{\gamma}) := \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\gamma})} \langle \mathbf{D}, \mathbf{T} \rangle = \min_{\mathbf{T} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\gamma})} \sum_{v, v' \in \mathcal{V} \times \mathcal{V}} t_{vv'} d_{vv'}, \quad (14)$$

where \mathbf{D} is the shortest-path distance matrix.

- ▶ W_1 is equivalent to a minimum-cost flow problem:

$$W_1(\boldsymbol{\mu}, \boldsymbol{\gamma}) = \min_{\mathbf{f} \in \Omega(\mathbf{S}_{\mathcal{V}}, \boldsymbol{\mu}, \boldsymbol{\gamma})} \|\text{diag}(\mathbf{w}) \mathbf{f}\|_1, \quad (15)$$

where $\mathbf{f} = [f_e] \in \mathbb{R}^{|\mathcal{E}|}$ is the flow on each edge.

$$\mathbf{S}_{\mathcal{V}} = [s_{ve}] \in \{0, \pm 1\}^{|\mathcal{V}| \times |\mathcal{E}|}, \quad s_{ve} = \begin{cases} 1 & \text{if } v \text{ is "head" of edge } e \\ -1 & \text{if } v \text{ is "tail" of edge } e \\ 0 & \text{otherwise.} \end{cases} \quad (16)$$

$$\Omega(\mathbf{S}_{\mathcal{V}}, \boldsymbol{\mu}, \boldsymbol{\gamma}) = \mathcal{U}^{|\mathcal{E}|} \cap \{\mathbf{f} \mid \mathbf{S}_{\mathcal{V}} \mathbf{f} = \boldsymbol{\gamma} - \boldsymbol{\mu}\}, \quad \text{where } \mathcal{U} = \begin{cases} [0, \infty), & \text{Directed } G, \\ \mathbb{R}, & \text{Undirected } G. \end{cases}$$

Proposed Quasi-Wasserstein (QW) Loss of GNNs

- ▶ When only partial signals on \mathcal{V}_L are given, we have **partial Wasserstein distance**: for $\mu, \gamma \in \mathbb{R}^{|\mathcal{V}_L|}$,

$$W_1^{(P)}(\mu, \gamma) = \min_{\mathbf{f} \in \Omega(\mathcal{S}_{\mathcal{V}_L}, \mu, \gamma)} \|\text{diag}(\mathbf{w}) \mathbf{f}\|_1, \quad (17)$$

- ▶ For partially-observed multi-dimensional labels, i.e., $\mathbf{Y}_{\mathcal{V}_L} = [\mathbf{y}_{\mathcal{V}_L}^{(c)}] \in \mathbb{R}^{|\mathcal{V}_L| \times C}$, **aggregating each dimension's $W_1^{(P)}$ leads to the QW loss**:

$$\begin{aligned} QW(\hat{\mathbf{Y}}_{\mathcal{V}_L}, \mathbf{Y}_{\mathcal{V}_L}) &= \sum_{c=1}^C W_1^{(P)}(\hat{\mathbf{y}}_{\mathcal{V}_L}^{(c)}, \mathbf{y}_{\mathcal{V}_L}^{(c)}) \\ &= \sum_{c=1}^C \min_{\mathbf{f}^{(c)} \in \Omega(\mathcal{S}_{\mathcal{V}_L}, \hat{\mathbf{y}}_{\mathcal{V}_L}^{(c)}, \mathbf{y}_{\mathcal{V}_L}^{(c)})} \|\text{diag}(\mathbf{w}) \mathbf{f}^{(c)}\|_1 \end{aligned}$$

Proposed Quasi-Wasserstein (QW) Loss of GNNs

- ▶ When only partial signals on \mathcal{V}_L are given, we have **partial Wasserstein distance**: for $\boldsymbol{\mu}, \boldsymbol{\gamma} \in \mathbb{R}^{|\mathcal{V}_L|}$,

$$W_1^{(P)}(\boldsymbol{\mu}, \boldsymbol{\gamma}) = \min_{\mathbf{f} \in \Omega(\mathcal{S}_{\mathcal{V}_L}, \boldsymbol{\mu}, \boldsymbol{\gamma})} \|\text{diag}(\mathbf{w}) \mathbf{f}\|_1, \quad (17)$$

- ▶ For partially-observed multi-dimensional labels, i.e., $\mathbf{Y}_{\mathcal{V}_L} = [\mathbf{y}_{\mathcal{V}_L}^{(c)}] \in \mathbb{R}^{|\mathcal{V}_L| \times C}$, **aggregating each dimension's $W_1^{(P)}$ leads to the QW loss**:

$$\begin{aligned} QW(\hat{\mathbf{Y}}_{\mathcal{V}_L}, \mathbf{Y}_{\mathcal{V}_L}) &= \sum_{c=1}^C W_1^{(P)}(\hat{\mathbf{y}}_{\mathcal{V}_L}^{(c)}, \mathbf{y}_{\mathcal{V}_L}^{(c)}) \\ &= \sum_{c=1}^C \min_{\mathbf{f}^{(c)} \in \Omega(\mathcal{S}_{\mathcal{V}_L}, \hat{\mathbf{y}}_{\mathcal{V}_L}^{(c)}, \mathbf{y}_{\mathcal{V}_L}^{(c)})} \|\text{diag}(\mathbf{w}) \mathbf{f}^{(c)}\|_1 \\ &= \min_{\mathbf{F} \in \Omega_C(\mathcal{S}_{\mathcal{V}_L}, g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta), \mathbf{Y}_{\mathcal{V}_L})} \|\text{diag}(\mathbf{w}) \mathbf{F}\|_1, \end{aligned} \quad (18)$$

- ▶ $\mathbf{F} = [\mathbf{f}^{(c)}]$ is the flow matrix with size $|\mathcal{E}| \times C$
- ▶ Feasible domain $\Omega_C = \mathcal{U}^{|\mathcal{E}| \times C} \cap \{\mathbf{F} \mid \mathcal{S}_{\mathcal{V}_L} \mathbf{F} = \mathbf{Y}_{\mathcal{V}_L} - g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta)\}$.

An Implicit Message Passing Associated with The QW Loss

- ▶ Learning GNN with the QW loss:

$$\begin{aligned} & \min_{\theta} QW(g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta), \mathbf{Y}_{\mathcal{V}_L}) \\ &= \min_{\theta} \min_{\mathbf{F} \in \Omega_C(\mathbf{S}_{\mathcal{V}_L}, g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta), \mathbf{Y}_{\mathcal{V}_L})} \|\text{diag}(\mathbf{w})\mathbf{F}\|_1. \end{aligned} \tag{19}$$

An Implicit Message Passing Associated with The QW Loss

- ▶ Learning GNN with the QW loss:

$$\begin{aligned} & \min_{\theta} QW(g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta), \mathbf{Y}_{\mathcal{V}_L}) \\ & = \min_{\theta} \min_{\mathbf{F} \in \Omega_C(\mathcal{S}_{\mathcal{V}_L}, g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta), \mathbf{Y}_{\mathcal{V}_L})} \|\text{diag}(\mathbf{w})\mathbf{F}\|_1. \end{aligned} \tag{19}$$

- ▶ Optionally, the flow can be used to parameterize the adjacency matrix

$$\min_{\theta, \xi} \min_{\mathbf{F} \in \Omega_C(\mathcal{S}_{\mathcal{V}_L}, g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}(\mathbf{F}; \xi); \theta), \mathbf{Y}_{\mathcal{V}_L})} \|\text{diag}(\mathbf{w})\mathbf{F}\|_1, \tag{20}$$

An Implicit Message Passing Associated with The QW Loss

- ▶ Learning GNN with the QW loss:

$$\begin{aligned} & \min_{\theta} QW(g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta), \mathbf{Y}_{\mathcal{V}_L}) \\ & = \min_{\theta} \min_{\mathbf{F} \in \Omega_C(\mathbf{S}_{\mathcal{V}_L}, g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta), \mathbf{Y}_{\mathcal{V}_L})} \|\text{diag}(\mathbf{w})\mathbf{F}\|_1. \end{aligned} \quad (19)$$

- ▶ Optionally, the flow can be used to parameterize the adjacency matrix

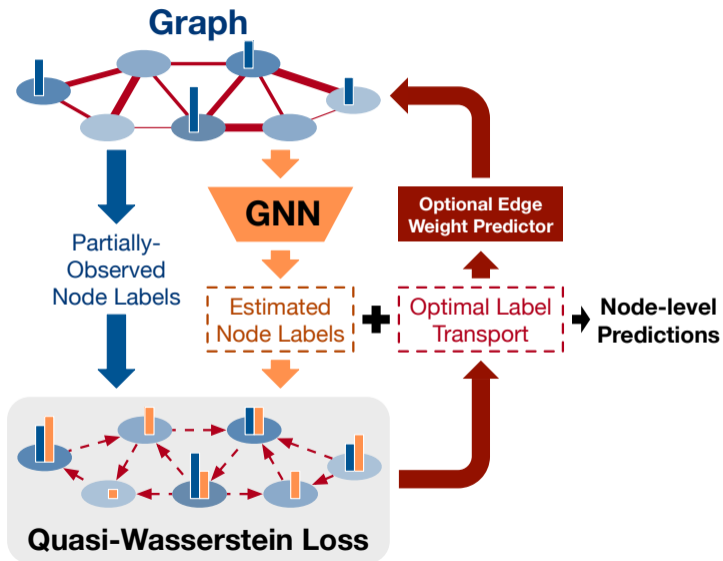
$$\min_{\theta, \xi} \min_{\mathbf{F} \in \Omega_C(\mathbf{S}_{\mathcal{V}_L}, g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}(\mathbf{F}; \xi); \theta), \mathbf{Y}_{\mathcal{V}_L})} \|\text{diag}(\mathbf{w})\mathbf{F}\|_1, \quad (20)$$

- ▶ Adding the optimal label transport in the testing phase: For $v \in \mathcal{V} \setminus \mathcal{V}_L$, we predict its label as

$$\tilde{\mathbf{y}}_v := g_v(\mathbf{X}, \mathbf{A}; \theta^*) + \mathbf{S}_v \mathbf{F}^*, \text{ s.t. } \mathbf{F}^* \leftarrow QW(g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta), \mathbf{Y}_{\mathcal{V}_L}). \quad (21)$$

An implicit message passing layer encoding training labels' optimal transport.

An Implicit Message Passing Associated with The QW Loss



Solvers of QW Loss

$$\begin{aligned} & \min_{\theta} \min_{\mathbf{F}} \|\text{diag}(\mathbf{w})\mathbf{F}\|_1. \\ & \text{s.t. } \mathbf{F} \in \mathcal{U}^{|\mathcal{E}| \times C} \cap \{\mathbf{F} \mid \mathbf{S}_{\mathcal{V}_L}\mathbf{F} = \mathbf{Y}_{\mathcal{V}_L} - g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta)\} \end{aligned} \tag{22}$$

where $B_{\phi}(x, y) = \phi(x) - \phi(y) - \langle \nabla \phi(y), x - y \rangle$.

Solvers of QW Loss

$$\begin{aligned} & \min_{\theta} \min_{\mathbf{F}} \|\text{diag}(\mathbf{w})\mathbf{F}\|_1. \\ & \text{s.t. } \mathbf{F} \in \mathcal{U}^{|\mathcal{E}| \times C} \cap \{\mathbf{F} \mid \mathbf{S}_{\mathcal{V}_L}\mathbf{F} = \mathbf{Y}_{\mathcal{V}_L} - g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta)\} \end{aligned} \quad (22)$$

where $B_{\phi}(x, y) = \phi(x) - \phi(y) - \langle \nabla \phi(y), x - y \rangle$.

- ▶ An inexact solver based on Bregman Divergence-based Relaxation:

$$\min_{\theta, \mathbf{F} \in \mathcal{U}^{|\mathcal{E}| \times C}} \|\text{diag}(\mathbf{w})\mathbf{F}\|_1 + \lambda B_{\phi}(g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta) + \mathbf{S}_{\mathcal{V}_L}\mathbf{F}, \mathbf{Y}_{\mathcal{V}_L}). \quad (23)$$

Solvers of QW Loss

$$\begin{aligned} & \min_{\theta} \min_{\mathbf{F}} \|\text{diag}(\mathbf{w})\mathbf{F}\|_1. \\ & \text{s.t. } \mathbf{F} \in \mathcal{U}^{|\mathcal{E}| \times C} \cap \{\mathbf{F} \mid \mathbf{S}_{\mathcal{V}_L}\mathbf{F} = \mathbf{Y}_{\mathcal{V}_L} - g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta)\} \end{aligned} \quad (22)$$

where $B_{\phi}(x, y) = \phi(x) - \phi(y) - \langle \nabla \phi(y), x - y \rangle$.

- ▶ An inexact solver based on Bregman Divergence-based Relaxation:

$$\min_{\theta, \mathbf{F} \in \mathcal{U}^{|\mathcal{E}| \times C}} \|\text{diag}(\mathbf{w})\mathbf{F}\|_1 + \lambda B_{\phi}(g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta) + \mathbf{S}_{\mathcal{V}_L}\mathbf{F}, \mathbf{Y}_{\mathcal{V}_L}). \quad (23)$$

- ▶ An exact solver based on Bregman ADMM:

$$\begin{aligned} & \min_{\theta, \mathbf{F} \in \mathcal{U}^{|\mathcal{E}| \times C}} \max_{\mathbf{Z} \in \mathbb{R}^{|\mathcal{V}_L| \times C}} \|\text{diag}(\mathbf{w})\mathbf{F}\|_1 \\ & \quad + \langle \mathbf{Z}, g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta) + \mathbf{S}_{\mathcal{V}_L}\mathbf{F} - \mathbf{Y}_{\mathcal{V}_L} \rangle \\ & \quad + \lambda B_{\phi}(g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta) + \mathbf{S}_{\mathcal{V}_L}\mathbf{F}, \mathbf{Y}_{\mathcal{V}_L}). \end{aligned} \quad (24)$$

Compare with Traditional Learning Paradigm

Method	Setting	Node Classification	Node Regression
Apply the Traditional loss	ψ	Cross-entropy or KL	MSE
	Predicted \mathbf{y}_v	$g_v(\mathbf{X}, \mathbf{A}; \theta), \forall v \in \mathcal{V} \setminus \mathcal{V}_L$	
Apply the QW loss	ϕ	Entropy	$\frac{1}{2} \ \cdot\ _2^2$
	$B_\phi(= \psi)$	KL	MSE
	Predicted \mathbf{y}_v	$g_v(\mathbf{X}, \mathbf{A}; \theta) + \mathbf{S}_v \mathbf{F}^*, \forall v \in \mathcal{V} \setminus \mathcal{V}_L$	

Compare with Traditional Learning Paradigm

Method	Setting	Node Classification	Node Regression
Apply the Traditional loss	ψ	Cross-entropy or KL	MSE
	Predicted \mathbf{y}_v	$g_v(\mathbf{X}, \mathbf{A}; \theta), \forall v \in \mathcal{V} \setminus \mathcal{V}_L$	
Apply the QW loss	ϕ	Entropy	$\frac{1}{2} \ \cdot\ _2^2$
	$B_\phi(= \psi)$	KL	MSE
	Predicted \mathbf{y}_v	$g_v(\mathbf{X}, \mathbf{A}; \theta) + \mathbf{S}_v \mathbf{F}^*, \forall v \in \mathcal{V} \setminus \mathcal{V}_L$	

Theorem

Let $\{\theta^*, \mathbf{F}^*, \mathbf{Z}^*\}$ be the optimal solution of (24), $\{\theta^\dagger, \mathbf{F}^\dagger\}$ be the optimal solution of (23), and θ^\ddagger be the optimal solution of $\min_{\theta} B_\phi(g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta), \mathbf{Y}_{\mathcal{V}_L})$, we have

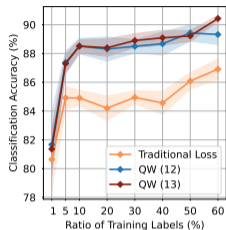
$$\begin{aligned}
 B_\phi(g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta^*) + \mathbf{S}_{\mathcal{V}_L} \mathbf{F}^*, \mathbf{Y}_{\mathcal{V}_L}) &\leq B_\phi(g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta^\dagger) + \mathbf{S}_{\mathcal{V}_L} \mathbf{F}^\dagger, \mathbf{Y}_{\mathcal{V}_L}) \\
 &\leq B_\phi(g_{\mathcal{V}_L}(\mathbf{X}, \mathbf{A}; \theta^\ddagger), \mathbf{Y}_{\mathcal{V}_L})
 \end{aligned}$$

QW Loss and Implicit MP: Superiority

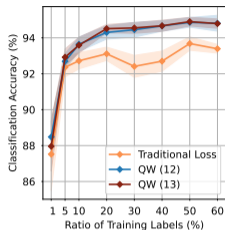
Node classification tasks:

Model	Method	Homophilic graphs					Heterophilic graphs					Overall Improve
		Cora	Citeseer	Pubmed	Computers	Photo	Squirrel	Chameleon	Actor	Texas	Cornell	
	#Nodes ($ \mathcal{V} $)	2,708	3,327	19,717	13,752	5,201	7,650	2,277	7,600	183	183	
	#Features (D)	1,433	3,703	500	767	754	2,089	2,325	932	1,703	1,703	
	#Edges ($ \mathcal{E} $)	5,278	4,552	44,324	245,861	119,081	198,358	31,371	26,659	279	277	
	Intra-edge rate	81.0%	73.6%	80.2%	77.7%	82.7%	22.2%	23.0%	21.8%	6.1%	12.3%	
	#Classes (C)	7	6	5	10	8	5	5	5	5	5	
GCN	$\textcircled{2}$ +LPA	87.44 \pm 0.96	79.98 \pm 0.84	86.93 \pm 0.29	88.42 \pm 0.45	93.24 \pm 0.43	46.55 \pm 1.15	63.57 \pm 1.16	34.00 \pm 1.28	77.21 \pm 3.28	61.91 \pm 5.11	—
	QW	86.34 \pm 1.45	78.51 \pm 1.22	84.72 \pm 0.70	82.48 \pm 0.69	88.10 \pm 1.31	44.81 \pm 1.81	60.90 \pm 1.63	32.43 \pm 1.59	78.69 \pm 6.47	68.72 \pm 5.95	-1.36
GAT	$\textcircled{2}$	87.88\pm0.79	81.36\pm0.41	87.89\pm0.40	89.20\pm0.41	93.81\pm0.36	52.62\pm0.49	68.10\pm1.01	38.09\pm0.50	84.10\pm2.95	84.26\pm2.98	+4.81
	QW	89.20\pm0.79	80.75\pm0.78	87.42 \pm 0.33	90.08 \pm 0.36	94.38 \pm 0.25	48.20 \pm 1.67	64.31 \pm 2.01	35.68\pm0.60	80.00 \pm 3.11	68.09 \pm 2.13	—
GIN	$\textcircled{2}$	89.11 \pm 0.66	80.19 \pm 0.64	88.38\pm0.23	90.41\pm0.28	94.65\pm0.24	55.03\pm1.35	67.35\pm1.42	33.86 \pm 2.13	80.33\pm1.80	70.21\pm2.13	+1.14
	QW	86.22 \pm 0.95	76.18\pm0.78	87.87\pm0.23	80.87 \pm 1.43	89.83 \pm 0.72	39.11 \pm 2.23	64.29 \pm 1.51	32.37\pm1.56	72.79 \pm 4.92	62.55 \pm 4.80	—
GraphSAGE	$\textcircled{2}$	86.24\pm0.90	76.13 \pm 1.09	87.53 \pm 0.34	89.28\pm0.45	92.60\pm0.44	65.29\pm0.68	73.26\pm1.12	32.32 \pm 1.93	77.54\pm2.60	64.04\pm3.62	+5.22
	QW	88.24\pm0.95	79.81 \pm 0.80	88.14 \pm 0.25	89.71 \pm 0.38	95.08 \pm 0.26	43.79 \pm 0.59	63.26 \pm 1.09	38.99\pm0.85	90.00 \pm 2.30	84.26 \pm 2.98	—
APNP	$\textcircled{2}$	87.59 \pm 0.77	80.52\pm0.68	88.61\pm0.32	90.17\pm0.24	95.25\pm0.25	54.37\pm0.89	68.32\pm0.68	37.82 \pm 0.45	90.33\pm1.97	86.38\pm2.13	+1.18
	QW	88.14 \pm 0.73	80.47 \pm 0.74	88.12 \pm 0.31	85.32 \pm 0.37	88.51 \pm 0.31	36.15 \pm 0.75	52.93 \pm 1.71	40.46 \pm 0.64	91.31 \pm 1.97	87.66 \pm 2.13	—
BernNet	$\textcircled{2}$	88.74\pm0.84	80.94\pm0.61	89.48\pm0.28	86.95\pm0.82	94.43\pm0.24	38.73\pm1.06	53.76\pm1.25	40.78\pm0.74	91.48\pm2.30	87.87\pm2.34	+1.41
	QW	88.28 \pm 1.00	79.81 \pm 0.79	88.87 \pm 0.38	87.61 \pm 0.46	93.68 \pm 0.28	51.15 \pm 1.09	67.96 \pm 1.05	40.72 \pm 0.80	93.28 \pm 1.48	90.21 \pm 2.35	—
ChebNetII	$\textcircled{2}$	89.03\pm0.76	81.35\pm0.71	89.03\pm0.38	89.58\pm0.47	94.55\pm0.39	55.22\pm0.64	71.66\pm1.18	40.91\pm0.71	93.44\pm1.80	90.85\pm2.34	+1.41
	QW	88.26 \pm 0.89	80.00\pm0.74	88.57 \pm 0.36	86.58 \pm 0.71	93.50 \pm 0.34	57.78 \pm 0.84	71.71 \pm 1.40	40.70 \pm 0.77	92.79 \pm 1.48	88.94\pm2.78	—
	QW	88.54\pm0.76	79.47 \pm 0.70	89.47\pm0.36	90.43\pm0.22	94.84\pm0.37	60.55\pm0.64	74.05\pm0.68	41.37\pm0.67	93.93\pm0.98	87.23 \pm 3.62	+1.11

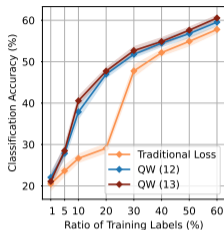
QW Loss and Implicit MP: Superiority



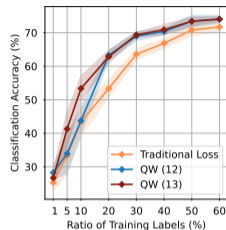
(a) Computers



(b) Photo



(c) Squirrel



(d) Chameleon

Figure: Illustrations of the methods' performance given different amounts of labeled nodes.

Node regression tasks:

Model	Method	Homophilic graphs		Heterophilic graphs	
		Computers	Photo	Actor	Cornell
GIN	(2)	0.0605 \pm 0.0018	0.0459 \pm 0.0044	0.1570 \pm 0.0014	0.1609 \pm 0.0359
	QW	0.0244\pm0.0028	0.0203\pm0.0012	0.1564\pm0.0012	0.1524\pm0.0043
BernNet	(2)	0.0871 \pm 0.0002	0.0488 \pm 0.0009	0.1661\pm0.0020	0.0989 \pm 0.0076
	QW	0.0364\pm0.0038	0.0297\pm0.0014	0.1671 \pm 0.0008	0.0753\pm0.0024

Case 3: OT-based Implicit Structured Attention Layer Design

Empirically, the power of Transformer is attributed to its Multi-head Attention (MHA) Architecture:

$$\text{MHA}_\theta(\mathbf{X}) := \text{Concat}_{\text{col}}(\{\text{Att}(\mathbf{V}_m; \mathbf{Q}_m, \mathbf{K}_m)\}_{m=1}^M), \quad (25)$$

For $m = 1, \dots, M$, we have $\mathbf{V}_m = \mathbf{X}\mathbf{W}_{V,m}$, $\mathbf{Q}_m = \mathbf{X}\mathbf{W}_{Q,m}$, and $\mathbf{K}_m = \mathbf{X}\mathbf{W}_{K,m}$.

$$P(\mathbf{Q}, \mathbf{K}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right), \quad (26)$$

Case 3: OT-based Implicit Structured Attention Layer Design

Empirically, the power of Transformer is attributed to its Multi-head Attention (MHA) Architecture:

$$\text{MHA}_\theta(\mathbf{X}) := \text{Concat}_{\text{col}}(\{\text{Att}(\mathbf{V}_m; \mathbf{Q}_m, \mathbf{K}_m)\}_{m=1}^M), \quad (25)$$

For $m = 1, \dots, M$, we have $\mathbf{V}_m = \mathbf{X}\mathbf{W}_{V,m}$, $\mathbf{Q}_m = \mathbf{X}\mathbf{W}_{Q,m}$, and $\mathbf{K}_m = \mathbf{X}\mathbf{W}_{K,m}$.

$$P(\mathbf{Q}, \mathbf{K}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right), \quad (26)$$

Challenges:

- ▶ **Over-smoothness** when dealing with long sequences.
- ▶ **Quadratic complexity** w.r.t. the sequence length.

Development Tendency of Attention Layer

Model	Attention($\mathbf{V}; \mathbf{Q}, \mathbf{K}$)	Complexity	Attention Structure	Faster?	Better?
Transformer	$\text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(DN^2)$	Dense + Row-wisely normalized	—	—

Development Tendency of Attention Layer

Model	Attention($V; Q, K$)	Complexity	Attention Structure	Faster?	Better?
Transformer	$\text{Softmax}\left(\frac{QK^\top}{\sqrt{D}}\right)V$	$\mathcal{O}(DN^2)$	Dense + Row-wisely normalized	—	—
SparseTrans	Local2D-Softmax $\left(\frac{QK^\top}{\sqrt{D}}\right)V$	$\mathcal{O}(DN^{1.5})$	Sparse + Row-wisely normalized	Yes	Comparable
Longformer	Local1D-Softmax $\left(\frac{QK^\top}{\sqrt{D}}\right)V$	$\mathcal{O}(DNL)$	Sparse + Row-wisely normalized	Yes	Comparable
Reformer	LSH-Softmax $\left(\frac{QK^\top}{\sqrt{D}}\right)V$	$\mathcal{O}(DN \log N)$	Sparse + Row-wisely normalized	Yes	Comparable
CosFormer	$(Q_{\cos}K_{\cos}^\top + Q_{\sin}K_{\sin}^\top)V$	$\mathcal{O}(DE_{QK})$	Sparse	Yes	Comparable

Development Tendency of Attention Layer

Model	Attention($\mathbf{V}; \mathbf{Q}, \mathbf{K}$)	Complexity	Attention Structure	Faster?	Better?
Transformer	$\text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(DN^2)$	Dense + Row-wisely normalized	—	—
SparseTrans	Local2D-Softmax $\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(DN^{1.5})$	Sparse + Row-wisely normalized	Yes	Comparable
Longformer	Local1D-Softmax $\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(DNL)$	Sparse + Row-wisely normalized	Yes	Comparable
Reformer	LSH-Softmax $\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(DN \log N)$	Sparse + Row-wisely normalized	Yes	Comparable
CosFormer	$(\mathbf{Q}_{\cos}\mathbf{K}_{\cos}^\top + \mathbf{Q}_{\sin}\mathbf{K}_{\sin}^\top)\mathbf{V}$	$\mathcal{O}(DE_{QK})$	Sparse	Yes	Comparable
Performer	$\phi_r(\mathbf{Q})\phi_r(\mathbf{K})^\top\mathbf{V}$	$\mathcal{O}(DNr)$	Low-rank	Yes	No
Linformer	$\text{Softmax}\left(\frac{\mathbf{Q}\psi_r(\mathbf{K})^\top}{\sqrt{D}}\right)\psi_r(\mathbf{V})$	$\mathcal{O}(DNr)$	Low-rank + Row-wisely normalized	Yes	No

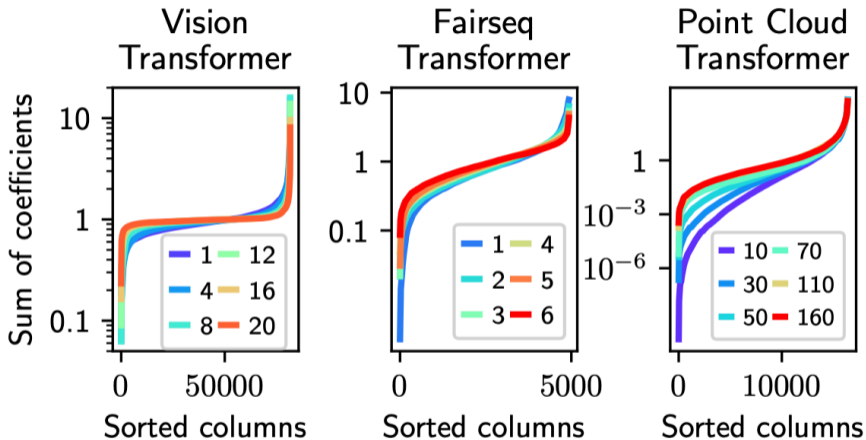
Development Tendency of Attention Layer

Model	Attention($\mathbf{V}; \mathbf{Q}, \mathbf{K}$)	Complexity	Attention Structure	Faster?	Better?
Transformer	$\text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(DN^2)$	Dense + Row-wisely normalized	—	—
SparseTrans	Local2D-Softmax $\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(DN^{1.5})$	Sparse + Row-wisely normalized	Yes	Comparable
Longformer	Local1D-Softmax $\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(DNL)$	Sparse + Row-wisely normalized	Yes	Comparable
Reformer	LSH-Softmax $\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(DN \log N)$	Sparse + Row-wisely normalized	Yes	Comparable
CosFormer	$(\mathbf{Q}_{\cos}\mathbf{K}_{\cos}^\top + \mathbf{Q}_{\sin}\mathbf{K}_{\sin}^\top)\mathbf{V}$	$\mathcal{O}(DE_{QK})$	Sparse	Yes	Comparable
Performer	$\phi_r(\mathbf{Q})\phi_r(\mathbf{K})^\top\mathbf{V}$	$\mathcal{O}(DNr)$	Low-rank	Yes	No
Linformer	$\text{Softmax}\left(\frac{\mathbf{Q}\psi_r(\mathbf{K})^\top}{\sqrt{D}}\right)\psi_r(\mathbf{V})$	$\mathcal{O}(DNr)$	Low-rank + Row-wisely normalized	Yes	No
Sinkformer	Sinkhorn $_K\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(KDN^2)$	Dense + Doubly stochastic	No	Yes

Development Tendency of Attention Layer

Model	Attention($\mathbf{V}; \mathbf{Q}, \mathbf{K}$)	Complexity	Attention Structure	Faster?	Better?
Transformer	$\text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(DN^2)$	Dense + Row-wisely normalized	—	—
SparseTrans	Local2D-Softmax $\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(DN^{1.5})$	Sparse + Row-wisely normalized	Yes	Comparable
Longformer	Local1D-Softmax $\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(DNL)$	Sparse + Row-wisely normalized	Yes	Comparable
Reformer	LSH-Softmax $\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(DN \log N)$	Sparse + Row-wisely normalized	Yes	Comparable
CosFormer	$(\mathbf{Q}_{\cos}\mathbf{K}_{\cos}^\top + \mathbf{Q}_{\sin}\mathbf{K}_{\sin}^\top)\mathbf{V}$	$\mathcal{O}(DE_{QK})$	Sparse	Yes	Comparable
Performer	$\phi_r(\mathbf{Q})\phi_r(\mathbf{K})^\top\mathbf{V}$	$\mathcal{O}(DNr)$	Low-rank	Yes	No
Linformer	$\text{Softmax}\left(\frac{\mathbf{Q}\psi_r(\mathbf{K})^\top}{\sqrt{D}}\right)\psi_r(\mathbf{V})$	$\mathcal{O}(DNr)$	Low-rank + Row-wisely normalized	Yes	No
Sinkformer	Sinkhorn $_K\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(KDN^2)$	Dense + Doubly stochastic	No	Yes
Sliceformer	Sort$_{\text{col}}(\mathbf{V})$	$\mathcal{O}(DN \log N)$	Full-rank+Sparse+Doubly stochastic	Yes	Yes

Doubly Stochastic Tendency of Attention Map



Sander, Michael E., et al. Sinkformers: Transformers with doubly stochastic attention. AISTATS, 2022.

Implement Attention Maps Implicitly via Slicing-Sorting

The attention layer in Sliceformer is implemented by slicing-sorting:

$$\text{SliceSort}(\mathbf{X}) := \text{Sort}_{\text{col}}(\underbrace{\mathbf{X}\mathbf{W}_V}_{\mathbf{V}=[\mathbf{v}_i]}) = \text{Concat}_{\text{col}}(\{\mathbf{P}_i\mathbf{v}_i\}_{i=1}^{MD}), \quad (27)$$

\mathbf{P}_i is a permutation matrix corresponding to the sorting of \mathbf{v}_i .

- ▶ Full-rank, sparse, and doubly stochastic.

Implement Attention Maps Implicitly via Slicing-Sorting

The attention layer in Sliceformer is implemented by slicing-sorting:

$$\text{SliceSort}(\mathbf{X}) := \text{Sort}_{\text{col}}(\underbrace{\mathbf{X}\mathbf{W}_V}_{\mathbf{V}=[\mathbf{v}_i]}) = \text{Concat}_{\text{col}}(\{\mathbf{P}_i\mathbf{v}_i\}_{i=1}^{MD}), \quad (27)$$

\mathbf{P}_i is a permutation matrix corresponding to the sorting of \mathbf{v}_i .

- ▶ Full-rank, sparse, and doubly stochastic.

Slicing-sorting operation = computing the Wasserstein barycenter of $\{\mathbf{v}_i\}_{i=1}^{MD}$:

$$\min_{\mathbf{v}} \sum_{i=1}^{MD} W_1(\mathbf{v}, \mathbf{v}_i)$$

Implement Attention Maps Implicitly via Slicing-Sorting

The attention layer in Sliceformer is implemented by slicing-sorting:

$$\text{SliceSort}(\mathbf{X}) := \text{Sort}_{\text{col}}(\underbrace{\mathbf{X}\mathbf{W}_V}_{\mathbf{V}=[\mathbf{v}_i]}) = \text{Concat}_{\text{col}}(\{\mathbf{P}_i\mathbf{v}_i\}_{i=1}^{MD}), \quad (27)$$

\mathbf{P}_i is a permutation matrix corresponding to the sorting of \mathbf{v}_i .

- ▶ Full-rank, sparse, and doubly stochastic.

Slicing-sorting operation = computing the Wasserstein barycenter of $\{\mathbf{v}_i\}_{i=1}^{MD}$:

$$\min_{\mathbf{v}} \sum_{i=1}^{MD} W_1(\mathbf{v}, \mathbf{v}_i) \xrightarrow{\text{1D case}} \sum_{i=1}^{MD} W_1(\mathbf{v}_1, \mathbf{v}_i)$$

Implement Attention Maps Implicitly via Slicing-Sorting

The attention layer in Sliceformer is implemented by slicing-sorting:

$$\text{SliceSort}(\mathbf{X}) := \text{Sort}_{\text{col}}(\underbrace{\mathbf{X}\mathbf{W}_V}_{\mathbf{V}=[\mathbf{v}_i]}) = \text{Concat}_{\text{col}}(\{\mathbf{P}_i\mathbf{v}_i\}_{i=1}^{MD}), \quad (27)$$

\mathbf{P}_i is a permutation matrix corresponding to the sorting of \mathbf{v}_i .

- ▶ Full-rank, sparse, and doubly stochastic.

Slicing-sorting operation = computing the Wasserstein barycenter of $\{\mathbf{v}_i\}_{i=1}^{MD}$:

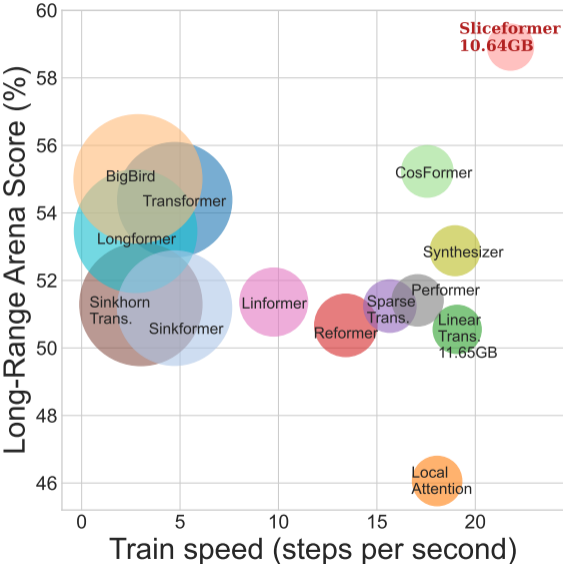
$$\min_{\mathbf{v}} \sum_{i=1}^{MD} W_1(\mathbf{v}, \mathbf{v}_i) \xrightarrow{\text{1D case}} \sum_{i=1}^{MD} W_1(\mathbf{v}_1, \mathbf{v}_i) \xrightarrow{\text{Alignment}} \sum_{i=1}^{MD} \sum_{n=1}^N |v_{1n} - v_{i\sigma_i(n)}|^2 \quad (28)$$

$$\sigma_i \Leftrightarrow \mathbf{P}_i, \quad \forall i = 1, \dots, MD.$$

Sliceformer: Superiority on LRA

Model	ListOps	Text	Retrieval	Image	Path	Path-X	Avg.
Transformer	36.37	64.27	57.46	42.44	71.40	FAIL	54.39
Local Att.	15.82	52.98	53.39	41.46	66.63	FAIL	46.06
Linear Trans.	16.13	65.90	53.09	42.34	75.30	FAIL	50.55
Reformer	37.27	56.10	53.40	38.07	68.50	FAIL	50.67
Sinkformer	30.70	64.03	55.45	41.08	64.65	FAIL	51.18
SparseTrans	17.07	63.58	59.59	44.24	71.71	FAIL	51.24
SinkhornTrans	33.67	61.20	53.83	41.23	67.45	FAIL	51.29
Linformer	35.70	53.94	52.27	38.56	76.34	FAIL	51.36
Performer	18.01	<u>65.40</u>	53.82	42.77	77.05	FAIL	51.41
Synthesizer	36.99	61.68	54.67	41.61	69.45	FAIL	52.88
Longformer	35.63	62.85	56.89	42.22	69.71	FAIL	53.46
BigBird	36.05	64.02	59.29	40.83	74.87	FAIL	55.01
Cosformer	37.90	63.41	61.36	43.17	70.33	FAIL	55.23
Sliceformer	<u>37.65</u>	64.60	62.23	48.02	82.04	FAIL	58.91

Sliceformer: Superiority on LRA



Sliceformer: Generalizability on Other Tasks

Image Classification: Model size ($\times 10^6$) and Top-1 accuracy (%)

Data Metric	Dogs vs. Cats		MNIST		CIFAR-10		CIFAR-100		Tiny-ImageNet	
	Size	Top-1	Size	Top-1	Size	Top-1	Size	Top-1	Size	Top-1
ViT	1.90	79.03	9.60	98.78	9.60	80.98	9.65	53.99	22.05	52.74
Sliceformer	1.11	79.71	6.50	98.81	6.46	82.16	6.50	54.24	18.50	51.77

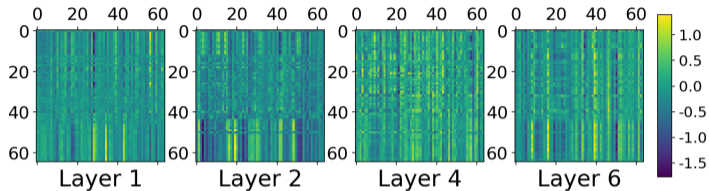
Text Classification:

Data Metric	IMDB	
	Size	Top-1
Transformer	8.84	83.05
Sliceformer	8.05	84.91

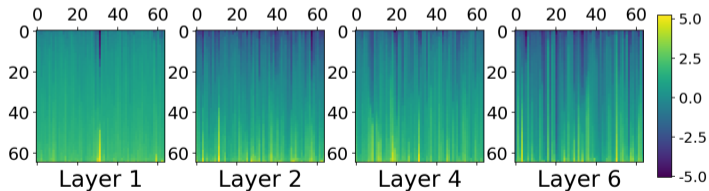
Molecular Property Prediction:

Data Metric	PCQM4M-LSC	
	Size	MAE
Graphormer	47.09	0.1287
Sliceformer	32.91	0.1308

Empirical Rationality of Sliceformer: Suppressing Mode Collapse



(a) ViT



(b) Sliceformer

Figure: $MHA_{\theta}(\mathbf{X})$ of ViT v.s. $SliceSort(\mathbf{X})$ of Sliceformer.

Empirical Rationality of Sliceformer: Suppressing Mode Collapse

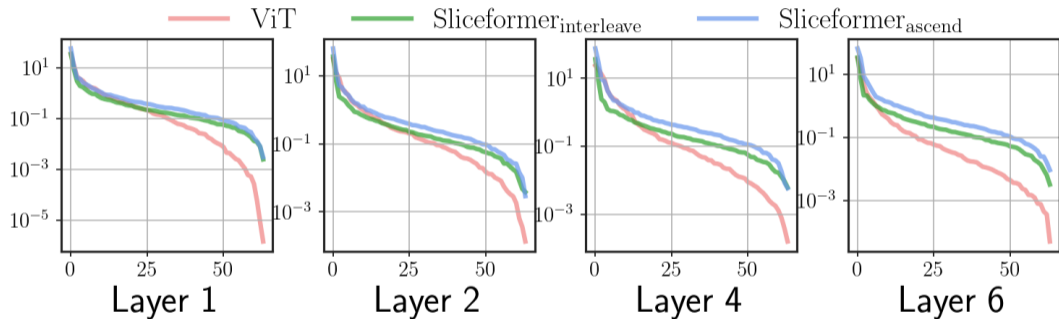


Figure: The comparisons on the singular spectrum.

Summary

Study the implicit design of information fusion layers and propose OT-driven solutions:

Layer	Design	OT Problem
Global Pooling	$(\mathbf{X} \odot \text{diag}^{-1}(\mathbf{P}_{\text{rot}}^*(\theta)\mathbf{1}_N)\mathbf{P}_{\text{rot}}^*(\theta))\mathbf{1}_N$	ROTP
Message-Passing	$g_{\mathcal{V}}(\mathbf{X}, \mathbf{A}; \theta) + \mathbf{S}_{\mathcal{V}}\mathbf{F}^*$	W_1 on graph
MHA Layer	$\text{Concat}_{\text{col}}(\{\mathbf{P}_i\mathbf{X}\mathbf{w}_i\}_{i=1}^{MD})$	1D WB

Summary

Study the implicit design of information fusion layers and propose OT-driven solutions:

Layer	Design	OT Problem
Global Pooling	$(\mathbf{X} \odot \text{diag}^{-1}(\mathbf{P}_{\text{rot}}^*(\theta)\mathbf{1}_N)\mathbf{P}_{\text{rot}}^*(\theta))\mathbf{1}_N$	ROTP
Message-Passing	$g_{\mathcal{V}}(\mathbf{X}, \mathbf{A}; \theta) + \mathbf{S}_{\mathcal{V}}\mathbf{F}^*$	W_1 on graph
MHA Layer	$\text{Concat}_{\text{col}}(\{\mathbf{P}_i\mathbf{X}\mathbf{w}_i\}_{i=1}^{MD})$	1D WB

Their common advantages:

- ▶ Simplify the design, even computation
- ▶ Feed-forward computation with better interpretability
- ▶ Advantages on data fitting

Summary

Study the implicit design of information fusion layers and propose OT-driven solutions:

Layer	Design	OT Problem
Global Pooling	$(\mathbf{X} \odot \text{diag}^{-1}(\mathbf{P}_{\text{rot}}^*(\theta)\mathbf{1}_N)\mathbf{P}_{\text{rot}}^*(\theta))\mathbf{1}_N$	ROTP
Message-Passing	$g_{\mathcal{V}}(\mathbf{X}, \mathbf{A}; \theta) + \mathbf{S}_{\mathcal{V}}\mathbf{F}^*$	W_1 on graph
MHA Layer	$\text{Concat}_{\text{col}}(\{\mathbf{P}_i\mathbf{X}\mathbf{w}_i\}_{i=1}^{MD})$	1D WB

Their common advantages:

- ▶ Simplify the design, even computation
- ▶ Feed-forward computation with better interpretability
- ▶ Advantages on data fitting

Current limitations:

- ▶ Over-fitting when distribution shifts
- ▶ The potentials in generative tasks are not investigated
- ▶ Some designs have limitations on model capacity (e.g., The Top-1 Acc. of Sliceformer on full-sized ImageNet is $< 70\%$.)

References and Resources

ROTP

- ▶ Paper: <https://ieeexplore.ieee.org/abstract/document/10247589/>
- ▶ Code: <https://github.com/SDS-Lab/ROT-Pooling>

QW Loss

- ▶ Paper: <https://arxiv.org/abs/2310.11762>

Sliceformer

- ▶ Paper: <https://arxiv.org/abs/2310.17683>
- ▶ Code: <https://github.com/SDS-Lab/sliceformer>

More OT-driven Work

- ▶ AAIL'22 Tutorial on Gromov-Wasserstein Learning
- ▶ IJCAI'23 Tutorial on Optimal Transport-based Machine Learning
- ▶ <https://hongtengxu.github.io/talks.html>

Thank you!

`https://hongtengxu.github.io`

`https://github.com/HongtengXu`

`hongtengxu@ruc.edu.cn`