

# Exploiting Edge Features in Graph-based Learning with Fused Network Gromov-Wasserstein Distance

Junjie Yang<sup>1</sup> Matthieu Labeau<sup>1</sup> Florence d'Alché-Buc<sup>1</sup>

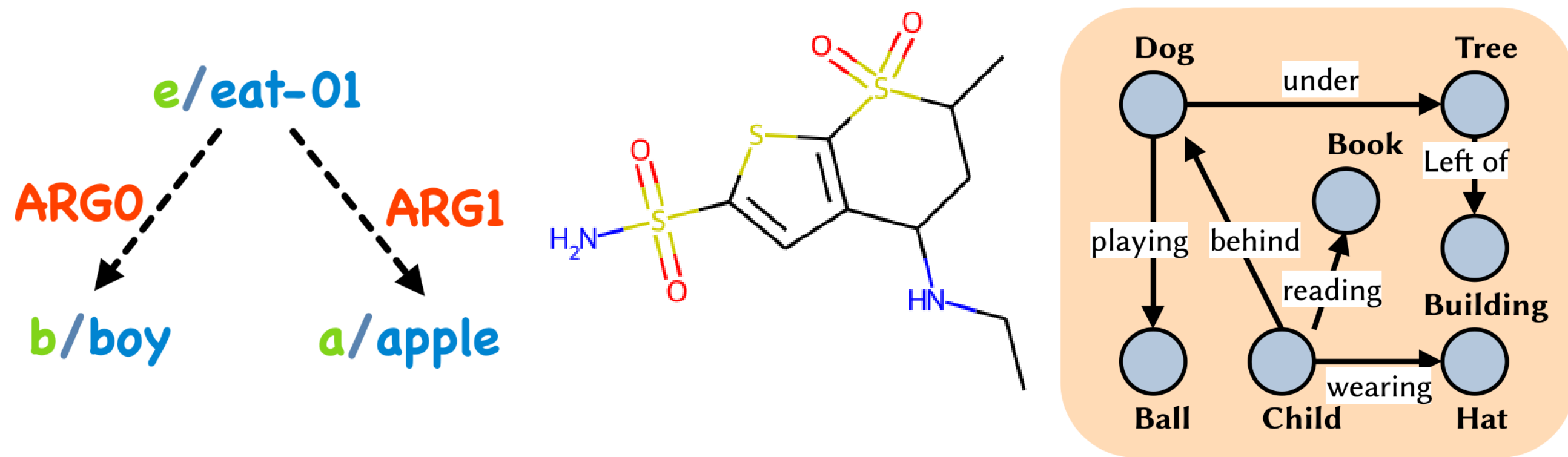
<sup>1</sup>LTCI, Télécom Paris, Institut Polytechnique de Paris, France

## Motivation

**Problem:** Graph-based Learning with Optimal Transport

**Existing Works:** Graph Classification with FGW [5], Graph Clustering with GW [4], FGW [5] or NGW [2], Graph Dictionary Learning with (F)GW [6], Supervised Graph Prediction with FGW [1]

**Goal:** We want to **unlock OT-based learning for edge featured graphs**. We target especially **Supervised Graph Prediction** task.

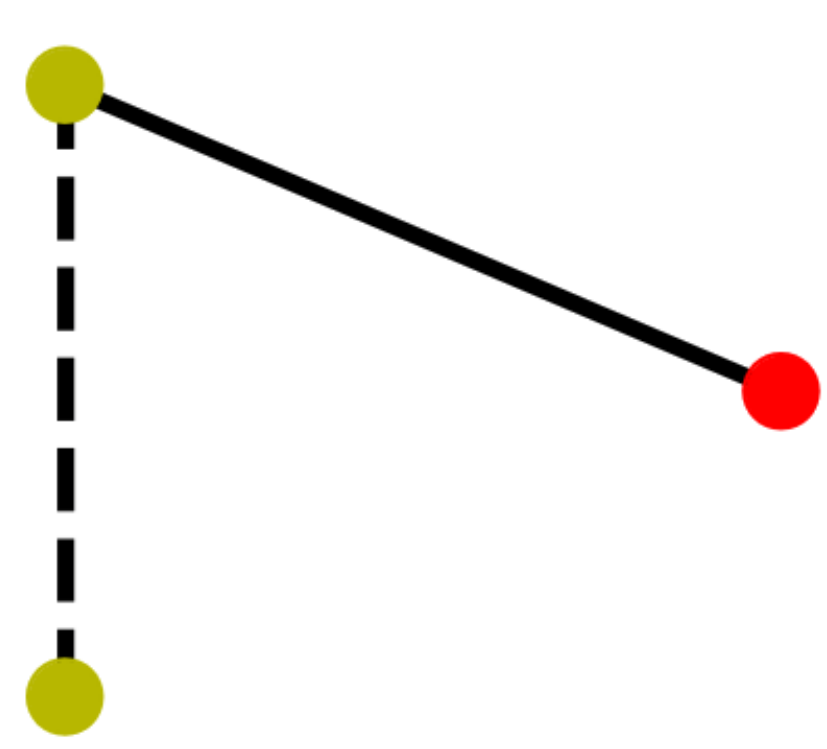


## Node and Edge Featured Graph

A node and edge featured graph of size  $m$  is a quadruple of the form  $(F, A, E, \mathbf{p})$  where

- $F \in \Psi^m$  is a tuple of points valued in a metric space  $(\Psi, d_\Psi)$
- $A \in \mathbb{R}^{m \times m}$  is a real-valued matrix
- $E \in \Omega^{m \times m}$  is a tuple of points valued in a metric space  $(\Omega, d_\Omega)$
- $\mathbf{p} \in \Sigma_m$  is a simplex histogram

We denote  $\mathcal{G}$  as a set of such quadruples.



- $\Psi = \{\text{red, yellow}\}$ : the node-color space
- $\Omega = \{\text{solid, dashed, non-edge}\}$ : the edge-type space

$$F = \begin{bmatrix} [1, 0] \\ [1, 0] \\ [0, 1] \end{bmatrix}, A = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

$$E = \begin{bmatrix} [0, 0, 1] & [0, 1, 0] & [1, 0, 0] \\ [0, 1, 0] & [0, 0, 1] & [0, 0, 1] \\ [1, 0, 0] & [0, 0, 1] & [0, 0, 1] \end{bmatrix}, \mathbf{p} = \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix}$$

## Fused Network Gromov-Wasserstein Distance

Given  $g = (F, A, E, \mathbf{p})$  of size  $m$ ,  $\tilde{g} = (\tilde{F}, \tilde{A}, \tilde{E}, \tilde{\mathbf{p}})$  of size  $\tilde{m}$  corresponding to two tuples of  $\mathcal{G}$ , and trade-off parameters  $(\alpha, \beta) \in [0, 1]^2$ , the Fused Network Gromov-Wasserstein distance between them for  $(p, q) \in [1, \infty]$  is written as:

$$\text{FNGW}_{\alpha, \beta, q, p}(g, \tilde{g}) = \min_{\pi \in \Pi(\mathbf{p}, \tilde{\mathbf{p}})} \mathcal{E}_{\alpha, \beta, q, p}(F, A, E, (\tilde{F}, \tilde{A}, \tilde{E}), \pi)$$

with

$$\mathcal{E}_{\alpha, \beta, q, p}(F, A, E, (\tilde{F}, \tilde{A}, \tilde{E}), \pi) = \left( \sum_{i, j, k, l} \left[ \alpha d_\Omega(E(i, k), \tilde{E}(j, l))^q + \beta |A(i, k) - \tilde{A}(j, l)|^q + (1 - \alpha - \beta) d_\Psi(F(i), \tilde{F}(j))^q \right]^p \pi_{k, l} \pi_{i, j} \right)^{1/p}$$

The FNGW distance satisfies the following **metric** properties: **positivity**, **symmetry**, **equality** with a corresponding notion of weak isomorphism, **relaxed triangle inequality** with a factor of  $2^{q-1}$ .

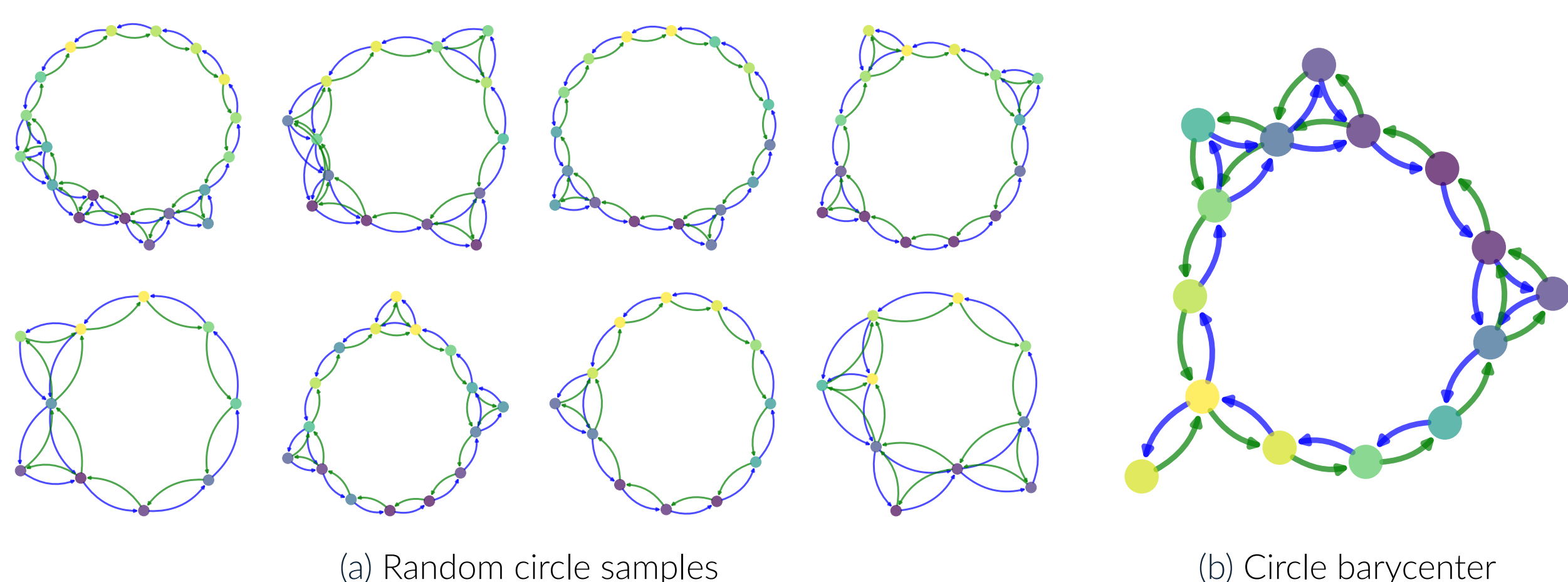


## FNGW Barycenter

Given a set  $\{g_k\}_{k=1}^K$  and a set of weights  $\{\lambda_k\}_{k=1}^K$  such that  $\sum_k \lambda_k = 1$ , the FNGW Barycenter for a pre-defined histogram  $\mathbf{p} \in \Sigma_n$  is defined as follows:

$$\text{Bary}(\{\lambda_k\}_k, \{g_k\}_k, \mathbf{p}) = \arg \min_{F \in \mathbb{R}^{n \times S}, A \in \mathbb{R}^{n \times n}, E \in \mathbb{R}^{n \times n \times T}} \sum_k \lambda_k \text{FNGW}_{\alpha, \beta}(F, A, E, \mathbf{p}, g_k)$$

We employ the Block Coordinate Descent (BCD) algorithm to obtain the FNGW barycenter where the tensor  $E$  can be updated by  $E = \frac{1}{\mathcal{I}_{n \times T \times 2} \mathbf{p} \mathbf{p}^T} \sum_k \lambda_k (E_k \times_2 \pi_k) \times_1 \pi_k$ .



(a) Random circle samples

(b) Circle barycenter

## Supervised Graph Prediction with FNGW

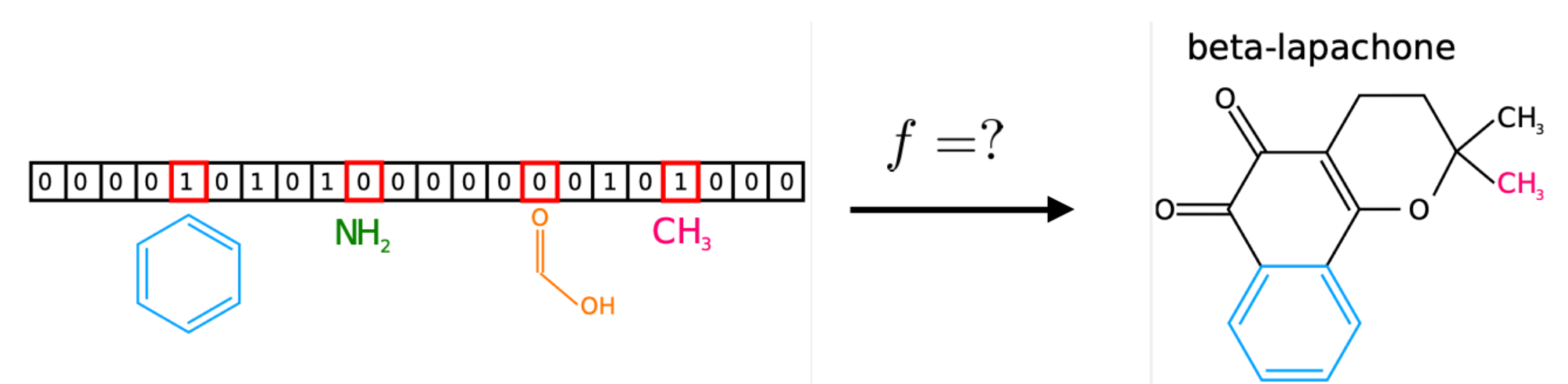
Given input space  $\mathcal{X}$ , output graph space  $\mathcal{G}$ , relaxed graph space  $\mathcal{G}_m = \{(F, A, E, \mathbf{p}) \mid C \in [0, 1]^{m \times m}, F = (F_i)_{i=1}^m \in \text{Conv}(\mathcal{F})^m, E = (E_{ij}) \in \text{Conv}(\mathcal{T})^{m \times m}, \mathbf{p} = m^{-1} \mathbf{1}_m\}$  where  $\mathcal{F} \subset \mathbb{R}^S$  and  $\mathcal{T} \subset \mathbb{R}^T$  are finite node and edge feature spaces, and training samples  $\{(x_i, g_i)\}_{i=1}^n$ , Supervised Graph Prediction requires finding an estimator  $f: \mathcal{X} \rightarrow \mathcal{G}_m$  of the minimizer  $f^*$  of the expected risk  $\mathcal{R}(f) = \mathbb{E}_\rho[\text{FNGW}_{\alpha, \beta}(f(X), G)]$ . Based on the work of [3, 1], we propose an estimator of the form

$$\hat{f}(x) = \arg \min_{g \in \mathcal{G}_m} \sum_{i=1}^n \xi(x)_i \text{FNGW}_{\alpha, \beta}(g, g_i)$$

with  $\xi(x) = \mathbf{K} S^T (\mathbf{S} \mathbf{K}^2 S^T + n \lambda \mathbf{S} \mathbf{K} S^T)^{\dagger} S \kappa_x$  where  $\mathbf{K} \in \mathbb{R}^{n \times n}$  is the input kernel Gram matrix,  $\kappa_x = (k(x, x_1), \dots, k(x, x_n))^T \in \mathbb{R}^n$ , and  $S \in \mathbb{R}^{s \times n}$  with  $s \ll n$  is a sketching matrix.

- The FNGW loss admits an **Implicit Loss Embedding (ILE)**  $\rightarrow \hat{f}$  is universally consistent and its learning rate is of order  $n^{-1/4}$  with additional assumptions.
- Sketched ILE** enables the supervised graph prediction with more than 100,000 training points.

## Experiment: Fingerprint to Molecule



	GED w/o edge feature ↓	GED w/ edge feature ↓
NNBary-FGW	5.000 ± 0.140	-
NNBary-FNGW	5.311 ± 0.090	5.756 ± 0.073
Sketched ILE-FGW	3.037 ± 0.111	-
Sketched ILE-FNGW	<b>1.449 ± 0.034</b>	<b>1.534 ± 0.029</b>

Table. Graph edit distances of different methods on the Fin2Mol test set.

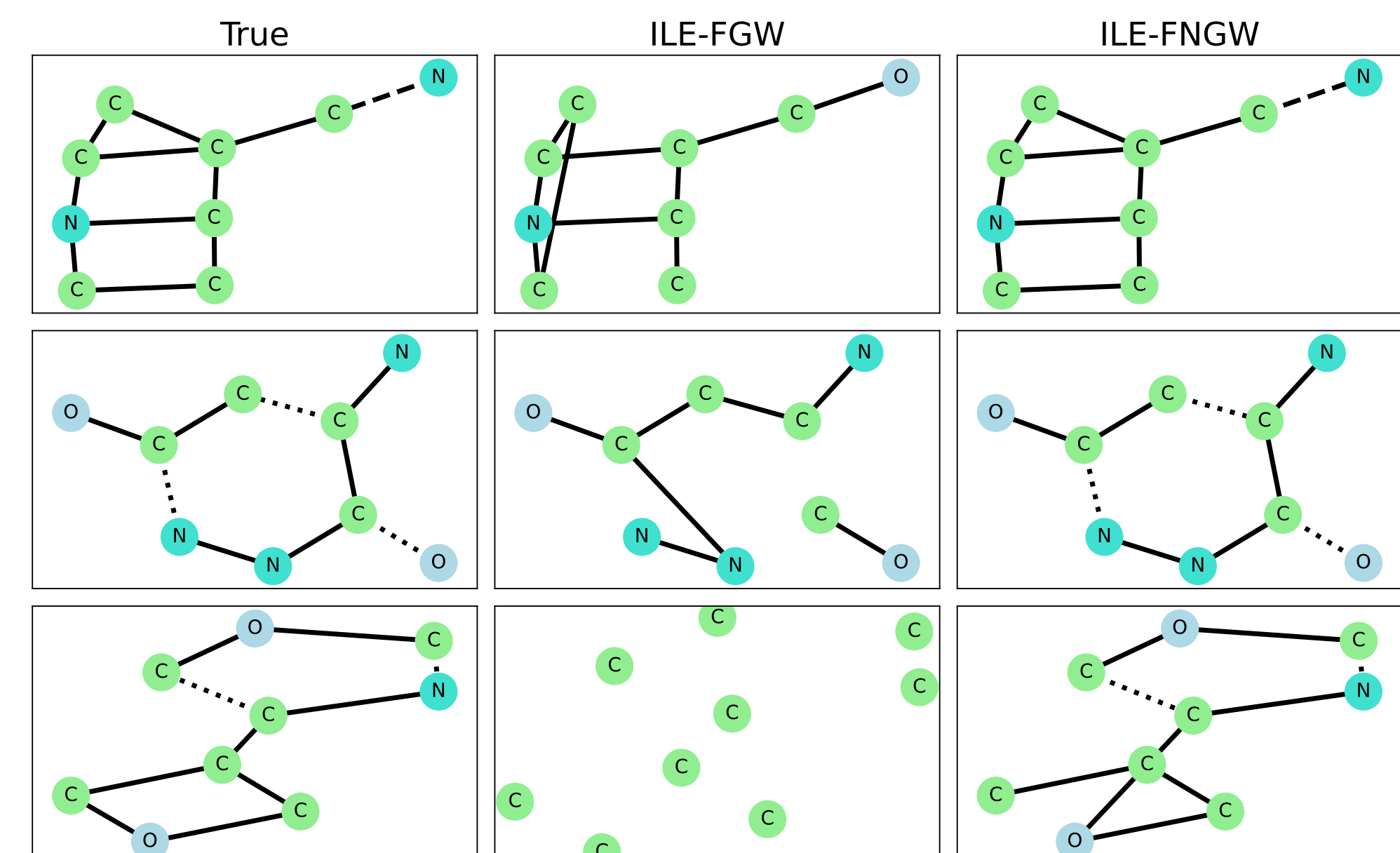


Figure. Qualitative comparison of the predicted QM9 molecules.

## Experiment: Metabolite Identification

To solve Metabolite Identification problem, the learning algorithm is expected to predict the metabolite (small molecules) given a tandem mass spectra. For each input spectra, a known set of metabolite candidates is provided.

	Top-1	Top-10	Top-20
WL kernel	9.8%	29.1%	37.4%
Fingerprint with linear kernel	28.6%	54.5%	59.9%
Fingerprint with gaussian kernel	<b>41.0%</b>	<b>62.0%</b>	<b>67.8%</b>
FGW diffuse	28.1%	53.6%	59.9%
FNGW diffuse + Bond stereo	27.7%	55.2%	60.9%
FNGW diffuse + Bond type	34.6%	55.1%	60.0%
FNGW diffuse + Mix	36.2%	58.2%	61.9%

Table. Top-k accuracies on the metabolite identification test set.

## References

- L. Brogat-Motte, R. Flamary, C. Brouard, J. Rousu, and F. D'Alché-Buc. Learning to Predict Graphs with Fused Gromov-Wasserstein Barycenters. In *ICML*, volume 162 of *Proceedings of Machine Learning Research*, pages 2321–2335, July 2022.
- S. Chowdhury and F. Mémoli. The Gromov-Wasserstein Distance Between Networks and Stable Network Invariants. *Information and Inference: A Journal of the IMA*, 8(4):757–787, 2019.
- C. Ciliberto, L. Rosasco, and A. Rudi. A General Framework for Consistent Structured Prediction with Implicit Loss Embeddings. *JMLR*, 21(98):1–67, 2020.
- G. Peyré, M. Cuturi, and J. Solomon. Gromov-Wasserstein Averaging of Kernel and Distance Matrices. In *ICML*, volume 48 of *Proceedings of Machine Learning Research*, pages 2664–2672, June 2016.
- T. Vayer, L. Chapel, R. Flamary, R. Tavenard, and N. Courty. Optimal Transport for structured data with application on graphs. In *ICML*, volume 97 of *Proceedings of Machine Learning Research*, pages 6275–6284, June 2019.
- C. Vincent-Cuaz, T. Vayer, R. Flamary, M. Corneli, and N. Courty. Online Graph Dictionary Learning. In *ICML*, volume 139 of *Proceedings of Machine Learning Research*, pages 10564–10574, July 2021.

## Acknowledgements

This research work is supported by the Hi! PARIS Center and the Institut Polytechnique de Paris.