

3.2. Energy-based Contrastive Learning

Initially, we define p_d as the distribution of the graph data and \mathcal{T} as a set of predetermined data augmentation operators. Given a dual-attribute subgraph g^\dagger and two augmentation views $t, t' \sim \mathcal{T}$ selected uniformly at random, we propose an ECL approach to build a joint distribution $p_\theta(\nu, \nu')$ to approximate $p_d(\nu, \nu')$ over two views $\nu, \nu' = t(g^\dagger), t'(g^\dagger)$.

Definition 1. $p_\theta(\nu, \nu')$ can be defined in a contrastive paradigm $f_\theta(\nu, \nu')$ as:

$$p_\theta(\nu, \nu') = \frac{\exp(-f_\theta(\nu, \nu'))}{Z(\theta)}, \quad (5)$$

where $Z(\theta) = \int \int \exp(-f_\theta(\nu, \nu')) d\nu d\nu'$.

Building upon the assumption that semantically similar pairs (ν, ν') have nearby projections with high p_d , while dissimilar ones would correspond to distant projections with low p_d , we solve for the distance between ν and ν' through $f_\theta(\cdot)$. $z = \phi_\theta(\varphi_\theta(\nu))$ is the corresponding representation of ν . $\varphi_\theta(\cdot)$ is a GNN encoder, and $\phi_\theta(\cdot)$ is a linear projection. The term $\|z - z'\|$ is used to indicate the inverse of semantic similarity of ν and ν' . To approximate $p_\theta(\nu, \nu')$ to $p_d(\nu, \nu')$, Eq. 1 can be rephrased as:

$$\min_{\theta} \mathbb{E}_{p_d}[-\log p_\theta(\nu, \nu')]. \quad (6)$$

Proposition 1. The joint distribution $p_\theta(\nu, \nu')$ can be formulated as an EBM:

$$p_\theta(\nu, \nu') = \frac{\exp(-E_\theta(\nu, \nu'))}{Z(\theta)}, \quad (7)$$

where $E_\theta(\nu, \nu') = \|z - z'\|^2/\tau$, and τ is a temperature parameter. The gradient of the objective of Eq. 6 is expressed as:

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{p_d}[-\log p_\theta(\nu, \nu')] = \\ \mathbb{E}_{p_d}[\nabla_{\theta} E_\theta(\nu, \nu')] - \mathbb{E}_{p_\theta}[\nabla_{\theta} E_\theta(\nu, \nu')]. \end{aligned} \quad (8)$$

To avoid directly calculating $Z(\theta)$, we employ Bayes' rule (Bayes, 1763) to reformulate $\mathbb{E}_{p_d}[-\log p_\theta(\nu, \nu')]$ as:

$$\begin{aligned} \mathbb{E}_{p_d}[-\log p_\theta(\nu, \nu')] = \\ \mathbb{E}_{p_d}[-\log p_\theta(\nu'|\nu)] + \mathbb{E}_{p_d}[-\log p_\theta(\nu)], \end{aligned} \quad (9)$$

where $p_\theta(\nu)$ is the marginal distribution of $p_\theta(\nu, \nu')$ over ν' .

Theorem 1. The marginal distribution $p_\theta(\nu)$ is an EBM:

$$p_\theta(\nu) = \frac{\exp(-E_\theta(\nu))}{Z(\theta)}, \quad (10)$$

where $E_\theta(\nu) = -\log \int e^{-\|z - z'\|^2/\tau} d\nu'$, detailed in the appendix. The gradient of the objective of Eq. 10 is defined as:

$$\nabla_{\theta} \mathbb{E}_{p_d}[-\log p_\theta(\nu)] = \mathbb{E}_{p_d}[\nabla_{\theta} E_\theta(\nu)] - \mathbb{E}_{p_\theta}[\nabla_{\theta} E_\theta(\nu)]. \quad (11)$$

Here, according to Eq. 9, the ECL objective is decomposed into the generative and discriminative terms, given by:

$$\mathcal{L}_b(\theta) = \mathbb{E}_{p_d}[-\log p_\theta(\nu'|\nu)] + \alpha \mathbb{E}_{p_d}[-\log p_\theta(\nu)], \quad (12)$$

where α is a hyperparameter to trade off the strength of two terms. According to Eq. 11, the gradient of the Eq. 12 is:

$$\begin{aligned} \nabla_{\theta} \mathcal{L}_b(\theta) = \mathbb{E}_{p_d}[-\nabla_{\theta} \log p_\theta(\nu'|\nu)] + \\ \alpha \mathbb{E}_{p_d}[\nabla_{\theta} E_\theta(\nu)] - \alpha \mathbb{E}_{p_\theta}[\nabla_{\theta} E_\theta(\nu)]. \end{aligned} \quad (13)$$

By this way, we solve $Z(\theta)$ in the discriminative term without additional calculations. For the generative term, we merely need to sample ν^* from $p_d(\nu)$ with adding noise $\mathcal{N}(0, \lambda)$ and iteratively optimize ν^* through SGLD, as indicated in Eq. 3.

Implementation 1. To implement the training of ECL, we approximate the generative and discriminative terms of Eq. 12, respectively, using the empirical mean of $p_\theta(\nu)$.

Given a mini-batch of samples $\{(\nu_n, \nu'_n)\}_{n=1}^N$, along with its representations $\{(z_n, z'_n)\}_{n=1}^N$, we have N positive and $2(N-1)$ negative samples. Therefore, the empirical mean $\hat{p}_\theta(\nu_n)$ (Kim & Ye, 2022) is defined as:

$$\hat{p}_\theta(\nu_n) = \frac{1}{2N} \sum_{\nu'_m: \nu'_m \neq \nu_n}^{2(N-1)} p_\theta(\nu_n, \nu'_m). \quad (14)$$

For the discriminative term, we utilize $\frac{p_\theta(\nu_n, \nu'_n)}{\hat{p}_\theta(\nu_n)}$ to approximate the conditional probability density $p_\theta(\nu'|\nu)$. According the SimCLR framework (Chen et al., 2020a), $\min_{\theta} \mathbb{E}_{p_d}[-\log \hat{p}_\theta(\nu'_n|\nu_n)]$ can be represented as:

$$\min_{z \in f_\theta(\nu)} -\log \left(\frac{\exp(-\|z_n - z'_n\|^2/\tau)}{\frac{1}{2N} \sum_{\nu'_m: \nu'_m \neq \nu_n}^{2(N-1)} \exp(-\|z_n - z'_m\|^2/\tau)} \right). \quad (15)$$

Considering only N positive samples in the generative term, we simplify Eq. 14 to $\hat{p}_\theta(\nu_n) = \frac{1}{N} \sum_{m=1}^N p_\theta(\nu_n, \nu'_m)$. The approximation of $\min_{\theta} \mathbb{E}_{p_d}[-\log \hat{p}_\theta(\nu, \nu_n)]$ is denoted as:

$$\min_{z \in f_\theta(\nu)} -\log \left(\sum_{n=1}^N \exp(-\|z - z_n\|^2/\tau) \right). \quad (16)$$

In summation, the final objective of ECL is:

$$\mathcal{L}_E(\theta) = \mathcal{L}_b(\theta) + \beta \mathcal{L}_r(\theta), \quad (17)$$

where $\mathcal{L}_r(\theta) = \frac{1}{2N} \sum_{m \neq n} E_\theta(\nu_n, \nu'_m)^2$ is the L_2 regularization loss that serves to prevent gradient overflow due to the excessive energy values. β is a trade-off hyperparameter.

Table 1. Node classification accuracy (mean(%) \pm std) with the standard splits on various benchmark datasets. The top three results are highlighted in **first**, **second**, and **third** best, respectively. "OOM" indicates out of memory.

Method	Cora	Citeseer	Cornell	Texas	Wisconsin	Actor	Pubmed	Arxiv	Products	Proteins
GCN	81.46 \pm 0.58	71.36 \pm 0.31	47.84 \pm 5.55	57.83 \pm 2.76	57.45 \pm 4.30	30.01 \pm 0.77	79.18 \pm 0.29	70.77 \pm 0.19	75.64 \pm 0.21	72.51 \pm 0.35
GAT	81.41 \pm 0.77	70.69 \pm 0.58	46.22 \pm 6.33	54.05 \pm 7.35	57.65 \pm 7.75	28.91 \pm 0.83	77.85 \pm 0.42	69.90 \pm 0.25	79.45 \pm 0.59	72.02 \pm 0.44
LDS	83.01 \pm 0.41	73.55 \pm 0.54	47.87 \pm 7.14	58.92 \pm 4.32	61.70 \pm 3.58	31.05 \pm 1.31	OOM	OOM	OOM	OOM
GEN	80.21 \pm 1.72	71.15 \pm 1.81	57.02 \pm 7.19	65.94 \pm 4.13	66.07 \pm 3.72	27.21 \pm 2.05	78.91 \pm 0.69	OOM	OOM	OOM
SGSR	83.48 \pm 0.43	72.96 \pm 0.25	44.32 \pm 2.16	60.81 \pm 4.87	56.86 \pm 1.24	30.23 \pm 0.38	78.09 \pm 0.53	OOM	OOM	OOM
GRCN	83.87 \pm 0.49	72.43 \pm 0.61	54.32 \pm 8.24	62.16 \pm 7.05	56.08 \pm 7.19	29.97 \pm 0.71	78.92 \pm 0.39	OOM	OOM	OOM
IDGL	83.88 \pm 0.42	72.20 \pm 1.18	50.00 \pm 8.98	62.43 \pm 6.09	59.41 \pm 4.11	28.16 \pm 1.41	OOM	OOM	OOM	OOM
GAuG-O	82.20 \pm 0.80	71.60 \pm 1.10	57.60 \pm 3.80	56.90 \pm 3.60	54.80 \pm 5.70	25.80 \pm 1.00	79.30 \pm 0.40	OOM	OOM	OOM
SUBLIME	83.40 \pm 0.42	72.30 \pm 1.09	70.54 \pm 5.98	70.03 \pm 4.23	66.81 \pm 6.55	30.79 \pm 0.68	73.80 \pm 0.60	55.50 \pm 0.10	—	—
ProGNN	80.30 \pm 0.57	68.51 \pm 0.52	54.05 \pm 6.16	48.37 \pm 12.17	62.54 \pm 7.56	22.35 \pm 0.88	71.60 \pm 0.46	OOM	OOM	OOM
CoGSL	81.76 \pm 0.24	73.09 \pm 0.42	52.16 \pm 3.21	59.46 \pm 4.36	58.82 \pm 1.52	32.95 \pm 1.20	OOM	OOM	OOM	OOM
STABLE	80.20 \pm 0.68	68.91 \pm 1.01	44.03 \pm 4.05	55.24 \pm 6.04	53.00 \pm 5.27	30.18 \pm 1.00	OOM	OOM	OOM	OOM
NodeFormer	80.28 \pm 0.82	71.31 \pm 0.98	42.70 \pm 5.51	58.92 \pm 4.32	48.43 \pm 7.02	25.51 \pm 1.17	78.21 \pm 1.43	55.40 \pm 0.23	—	—
ECL-GSR	84.06 \pm 0.84	73.70 \pm 0.75	71.27 \pm 2.06	72.97 \pm 3.39	67.79 \pm 1.03	33.71 \pm 0.96	80.91 \pm 1.12	71.09 \pm 0.31	80.47 \pm 0.22	74.64 \pm 0.31

Table 2. Pairwise independent sample t -tests comparison of ECL-GSR with other methods on Cora, Citeseer, Actor, and Pubmed datasets.

Method	Cora		Citeseer		Actor		Pubmed	
	T -statistic	P -value	T -statistic	P -value	T -statistic	P -value	T -statistic	P -value
GCN	12.33	4.24×10^{-9}	7.85	4.18×10^{-6}	11.45	1.40×10^{-8}	5.46	3.66×10^{-3}
GAT	11.15	2.12×10^{-8}	11.00	2.62×10^{-8}	14.54	2.83×10^{-10}	6.89	2.51×10^{-5}
LDS	6.41	6.37×10^{-5}	11.62	8.97×10^{-7}	10.14	3.96×10^{-3}	—	—
GEN	5.66	2.97×10^{-4}	8.22	4.65×10^{-4}	14.72	2.31×10^{-10}	8.51	3.22×10^{-4}
SGSR	10.21	6.84×10^{-4}	5.24	1.55×10^{-2}	13.35	1.15×10^{-9}	4.71	2.29×10^{-3}
GRCN	6.95	1.12×10^{-3}	13.06	2.04×10^{-5}	11.71	9.69×10^{-9}	4.21	6.86×10^{-3}
IDGL	7.58	2.47×10^{-3}	3.99	8.05×10^{-2}	13.40	1.08×10^{-9}	—	—
GAuG-O	5.92	1.71×10^{-4}	4.22	6.67×10^{-3}	16.83	2.42×10^{-11}	9.10	8.02×10^{-3}
SUBLIME	8.09	8.85×10^{-4}	7.49	4.86×10^{-4}	8.33	1.77×10^{-6}	16.04	5.45×10^{-11}
ProGNN	16.33	3.99×10^{-11}	18.10	6.97×10^{-12}	25.95	1.34×10^{-14}	21.63	3.22×10^{-13}
CoGSL	13.06	1.66×10^{-9}	4.67	3.67×10^{-2}	12.37	3.81×10^{-5}	—	—
STABLE	15.52	9.45×10^{-11}	9.43	2.86×10^{-7}	7.58	6.86×10^{-6}	—	—
NodeFormer	12.26	4.62×10^{-9}	6.09	1.22×10^{-4}	14.36	3.46×10^{-10}	4.90	1.49×10^{-3}

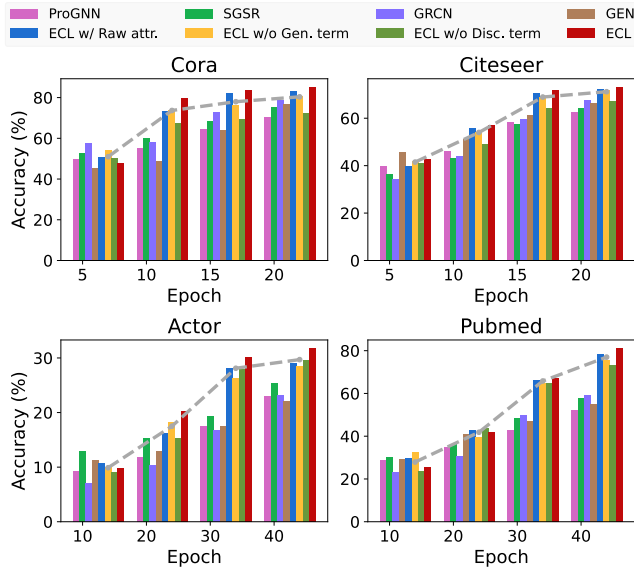


Figure 1. Performance analysis of the ECL-GSR variants and other benchmarks over a range of training epochs on four datasets.

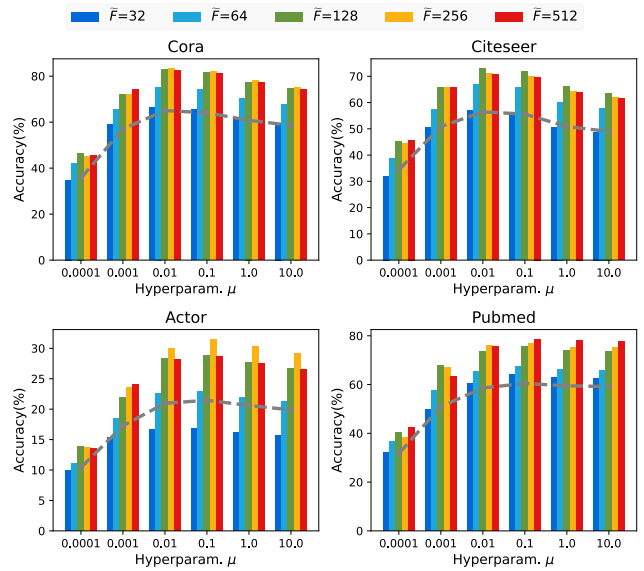


Figure 2. Hyperparameter μ and dimensionality \tilde{F} analysis of ECL-GSR on the Cora, Citeseer, Actor, and Pubmed datasets.