

### 3.2. Energy-based Contrastive Learning

Initially, we define  $p_d$  as a 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  $(\nu, \nu')$  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  $\nu$  and  $\nu'$  through  $f_\theta(\cdot)$ .  $z = \phi_\theta(\varphi_\theta(\nu))$  is the corresponding representation of  $\nu$ .  $\varphi_\theta(\cdot)$  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  $\nu$  and  $\nu'$ . 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  $\tau$  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  $\nu'$ .

**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 objective of ECL 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  $\alpha$  is a hyperparameter to trade off the strength of two terms. According to Eq. 11, the gradient of 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,  $Z(\theta)$  ingeniously cancels itself out in the discriminative term without additional calculations. For the generative term, we merely need to sample  $\nu^*$  from  $p_d(\nu)$  with adding noise  $\mathcal{N}(0, \lambda)$  and iteratively optimize  $\nu^*$  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.  $\beta$  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	<b>79.18 <math>\pm</math> 0.29</b>	<b>70.77 <math>\pm</math> 0.19</b>	<b>75.64 <math>\pm</math> 0.21</b>	<b>72.51 <math>\pm</math> 0.35</b>
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	<b>69.90 <math>\pm</math> 0.25</b>	<b>79.45 <math>\pm</math> 0.59</b>	<b>72.02 <math>\pm</math> 0.44</b>
LDS	83.01 $\pm$ 0.41	<b>73.55 <math>\pm</math> 0.54</b>	47.87 $\pm$ 7.14	58.92 $\pm$ 4.32	61.70 $\pm$ 3.58	<b>31.05 <math>\pm</math> 1.31</b>	OOM	OOM	OOM	OOM
GEN	80.21 $\pm$ 1.72	71.15 $\pm$ 1.81	<b>57.02 <math>\pm</math> 7.19</b>	<b>65.94 <math>\pm</math> 4.13</b>	<b>66.07 <math>\pm</math> 3.72</b>	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	<b>83.87 <math>\pm</math> 0.49</b>	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	<b>83.88 <math>\pm</math> 0.42</b>	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	<b>79.30 <math>\pm</math> 0.40</b>	OOM	OOM	OOM
SUBLIME	83.40 $\pm$ 0.42	72.30 $\pm$ 1.09	<b>70.54 <math>\pm</math> 5.98</b>	<b>70.03 <math>\pm</math> 4.23</b>	<b>66.81 <math>\pm</math> 6.55</b>	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	<b>73.09 <math>\pm</math> 0.42</b>	52.16 $\pm$ 3.21	59.46 $\pm$ 4.36	58.82 $\pm$ 1.52	<b>32.95 <math>\pm</math> 1.20</b>	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	<b>84.06 <math>\pm</math> 0.84</b>	<b>73.70 <math>\pm</math> 0.75</b>	<b>71.27 <math>\pm</math> 2.06</b>	<b>72.97 <math>\pm</math> 3.39</b>	<b>67.79 <math>\pm</math> 1.03</b>	<b>33.71 <math>\pm</math> 0.96</b>	<b>80.91 <math>\pm</math> 1.12</b>	<b>71.09 <math>\pm</math> 0.31</b>	<b>80.47 <math>\pm</math> 0.22</b>	<b>74.64 <math>\pm</math> 0.31</b>

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

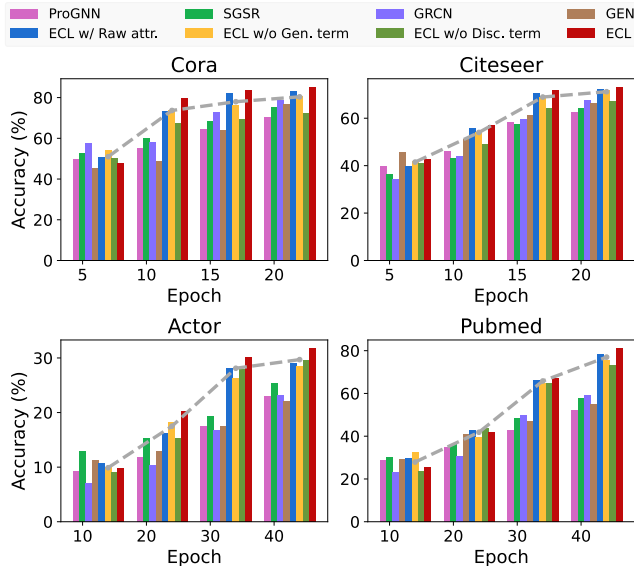


Figure 1. Performance analysis of the ECL-GSR variants and other benchmarks over a range of training epochs on four datasets.

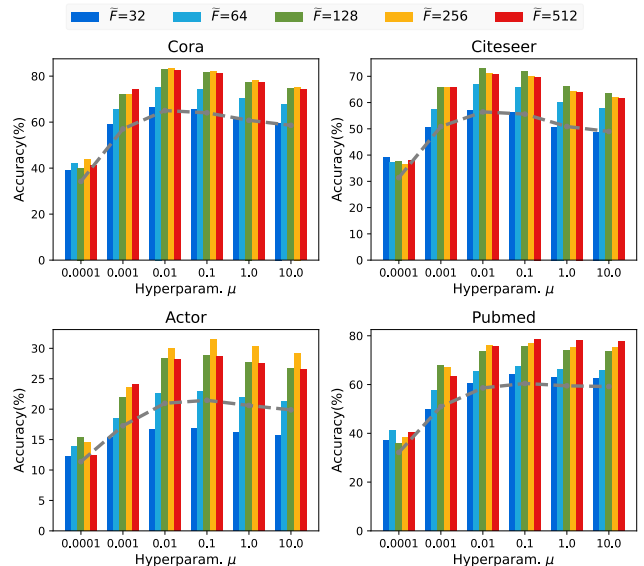


Figure 2. Hyperparameter  $\mu$  and dimensionality  $\tilde{F}$  analysis of ECL-GSR on the Cora, Citeseer, Actor, and Pubmed datasets.