
E²Gen: Evidential Energy-Based Generation for Fair Federated Graph Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Federated Graph Learning is rapidly evolving as a privacy-preserving collaborative
2 approach for decentralized graph data. However, severe fairness challenges are
3 increasingly undermining federated systems by systematically degrading predic-
4 tions for structurally disadvantaged minority nodes. The inherent vulnerabilities
5 and missing topological contexts in Federated Graph Learning are deeply entan-
6 gled, making traditional federated fairness methods and simple oversampling less
7 effective. In our work, we propose an effective Evidential Energy-based Gen-
8 eration framework for Fair Federated Graph Learning (*E²Gen*). At the client
9 level, it explicitly identifies structurally deficient nodes using a multi-axis metric,
10 synthesizing targeted representations via conditional energy-based models, and
11 selects reliable samples through an evidential quality gate. At the server level, the
12 local performance disparities uploaded by each client are evaluated to construct
13 a fairness gap assessment, making the global model absorb equitable improve-
14 ments by further adjusting the aggregation weights. Our method can handle high
15 topological heterogeneity, does not require strict generative normalization, and is
16 effective under both homophilic and heterophilic graph structures. Extensive results
17 on various settings of federated graph scenarios under severe fairness challenges
18 validate the effectiveness of this approach. The code is anonymously available at
19 <https://anonymous.4open.science/r/E-Gen-A2C6>.

20 1 Introduction

21 Federated Learning (FL) (McMahan et al., 2017) has fundamentally transformed decentralized
22 machine learning. It allows multiple clients to collaboratively train a shared model while keeping
23 sensitive data strictly local, making it a highly effective solution for privacy-preserving applications
24 (Kairouz and McMahan, 2021). This paradigm naturally extends to graph-structured data, forming the
25 basis of Federated Graph Learning (FGL) (Zhang et al., 2021a). While FGL successfully leverages
26 distributed graphs in domains like healthcare, finance, and social networks, its decentralized nature
27 introduces severe fairness challenges. In the federated context, fairness requires the global model to
28 deliver equitable predictive performance across diverse demographic groups. When this principle is
29 violated, the model systematically degrades predictions for minority entities, potentially leading to
30 discriminatory outcomes in high-stakes risk assessments.

31 With the objective of mitigating these disparities, fairness-aware mechanisms in traditional FL
32 have been widely studied (Li et al., 2019; Mohri et al., 2019). Certain methods prioritize minority
33 groups by directly reweighting local objective functions. Some approaches dynamically adjust server
34 aggregation weights or utilize personalized learning to balance performance across heterogeneous
35 clients (Fallah et al., 2020; Li et al., 2021b). However, these methods often struggle to perform
36 effectively in FGL environments due to the unique properties of graph data. Unlike traditional
37 machine learning that assumes independent data instances, graph learning fundamentally relies on

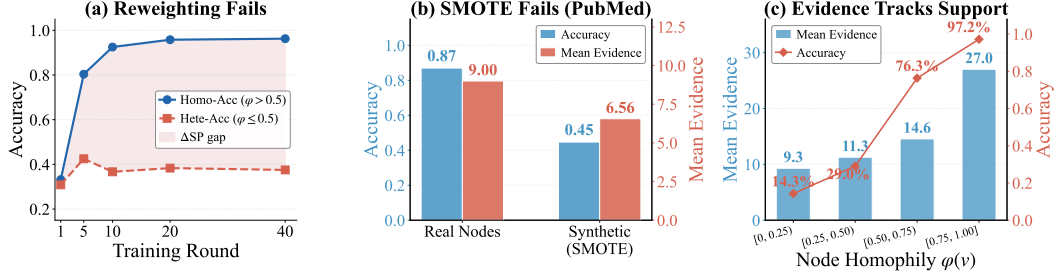


Figure 2: Empirical observations on the inherent vulnerabilities of disadvantaged nodes in FGL. **(a) Reweighting Fails:** Using the class-level reweighting method BoostFGL, the accuracy gap between advantaged (homophilous, $\varphi > 0.5$) and disadvantaged ($\varphi \leq 0.5$) nodes still persists and widens to 0.59 at round 40. **(b) SMOTE Fails:** Traditional latent interpolation (e.g., FedSig) generates unreliable representations, causing accuracy to plummet from 0.87 to 0.45 and mean evidence to drop from 9.00 to 6.56. **(c) Evidence Tracks Reliability:** The quantified mean evidence monotonically increases alongside node homophily (φ) and accuracy, confirming it as a valid metric for support deficiency. For detailed experimental setups and parameter configurations, please refer to Appendix B.

38 message passing mechanisms (Gilmer et al., 2017), where a node’s embedding is recursively updated
 39 by aggregating information from its neighbors. This interconnected nature dictates that a node’s
 40 learning quality is deeply entangled with its local context. Consequently, within this message-passing
 41 paradigm, disadvantaged nodes face critical inherent vulnerabilities. They lack the high-quality local
 42 patterns and informative connections necessary to aggregate meaningful messages. Since standard
 43 FL fairness techniques operate solely by adjusting local loss functions or server aggregation weights
 44 and remain inherently agnostic to these underlying node-level deficits, they cannot mitigate these
 45 vulnerabilities, thus proving ineffective for graph-based fairness issues.

46 Based on the aforementioned discussion, we review the challenges existing in FGL
 47 fairness. First of all, to address performance disparities, many existing fairness-aware
 48 FGL methods heavily modify the local objective function by reweighting
 49 losses or assigning higher gradient penalties to underperforming nodes (Wu et al.,
 50 2025; Chen et al., 2026). However, amplifying a poor learning signal cannot yield
 51 meaningful features if the underlying informative context is fundamentally missing.
 52 As empirically demonstrated in Figure 2(a), relying solely on reweighting mechanisms
 53 (e.g., BoostFGL) fails to bridge the performance gap between advantaged and
 54 disadvantaged nodes. This inability to create new meaningful representations inherently
 55 limits the effectiveness of reweighting strategies. This leads to our first question: 1) *How can we explicitly identify disadvantaged nodes and
 56 synthetically generate targeted representations to compensate for their inherent vulnerabilities?*

67 To address this first question, we introduce the **Deficit Score**, a comprehensive diagnostic metric that
 68 decomposes a node’s inherent vulnerabilities along three complementary axes: **Local Homophily**,
 69 **Node Degree**, and **Class Frequency**. Instead of relying on a single coarse indicator like local training
 70 loss (Tang et al., 2020; Zhu et al., 2020), this metric evaluates the multi-dimensional context along
 71 these axes to produce a calibrated aggregate score. By elevating fairness diagnosis from macro performance
 72 tracking to explicit, fine-grained evaluation, the Deficit Score provides a stable, reproducible
 73 signal to pinpoint precisely which disadvantaged nodes require targeted augmentation. We empirically
 74 validate the necessity and effectiveness of this fused multi-axis formulation through comprehensive
 75 ablation studies, as detailed in Section 5.3. Subsequently, to generate higher-quality synthetic rep-
 76 resentations for these identified nodes, we turn to **energy-based models (EBMs)** and explore their

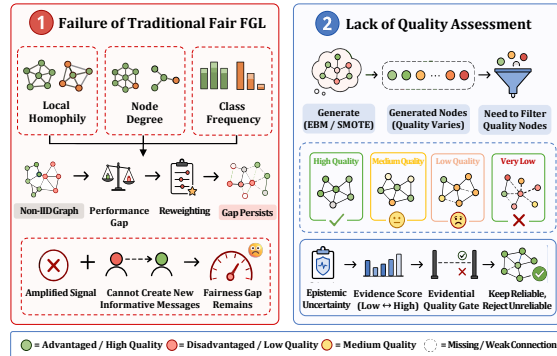


Figure 1: **Problem Illustration.** We describe the challenges in fair FGL: **I)** Traditional reweighting methods merely amplify existing signals without creating new informative messages, failing to close the fairness gap. **II)** Generation-based methods lack quality assessment, necessitating an evidential quality gate to filter out unreliable generated nodes.

77 potential. Energy is fundamentally an unnormalized probability likelihood (LeCun et al., 2006; Song
 78 and Kingma, 2021), offering a flexible modeling approach that is not constrained by strict normaliza-
 79 tion. The core strength of EBMs lies in their architectural versatility, allowing them to be integrated
 80 with virtually any model architecture. In our work, we combine energy functions with Graph Convo-
 81 lutional Networks (GCN) (Kipf and Welling, 2016) to construct an Energy-based GCN. This design
 82 preserves the GCN’s capability to capture complex relational information while benefiting from the
 83 flexibility of energy-based modeling to stably generate plausible synthetic latent representations. By
 84 unifying the Deficit Score and this generative capability, we establish a robust, targeted augmentation
 85 foundation to effectively compensate for the inherent vulnerabilities of disadvantaged nodes.

86 Secondly, while generating synthetic data presents a potential solution to mitigate these inherent
 87 vulnerabilities, directly injecting these representations into local training introduces severe risks of
 88 negative transfer. For instance, some approaches adapt oversampling techniques like GraphSMOTE
 89 (Zhao et al., 2021) to the federated latent space (Bi et al., 2024). However, applying simple linear
 90 interpolation within a highly non-IID topological space often generates out-of-distribution noise
 91 or unreliable representations. As empirically shown in Figure 2(b), samples generated by simple
 92 interpolation suffer from severe semantic degradation and high epistemic uncertainty, which we
 93 quantify as a sharp decrease in predictive “evidence”. Compounding this issue, standard softmax
 94 classifiers are notoriously overconfident on borderline synthetic samples (Guo et al., 2017; Hendrycks
 95 and Gimpel, 2016), rendering them incapable of reliably assessing generation quality. This brings
 96 up our second question: 2) *How can we accurately quantify the epistemic uncertainty of generated*
 97 *graph representations to ensure only reliable samples are incorporated into the training process?*

98 To solve this second question regarding quality control, we introduce **Evidential Quality Gating** (Sen-
 99 soy et al., 2018). Motivated by the theoretical framework of subjective logic (Jsang, 2018), we adopt
 100 an evidence vector $e \in \mathbb{R}_+^K$ for a K -class classification task. Mathematically, this evidence vector is
 101 used to parameterize a Dirichlet distribution $\text{Dir}(p|\alpha)$ over the categorical probabilities, where the
 102 concentration parameters are defined as $\alpha = e + 1$. Unlike standard softmax outputs that only provide
 103 overconfident point estimates, this evidence-parameterized distribution can be utilized to explicitly
 104 quantify the model’s epistemic uncertainty—essentially capturing the model’s “lack of knowledge”
 105 regarding an unfamiliar representation. As validated in Figure 2(c), this theoretical evidence metric
 106 exhibits a strong positive correlation with a node’s actual reliability, making it an ideal diagnostic
 107 signal. Therefore, we evaluate every EBM-generated sample through this gate, dynamically down-
 108 weighting or rejecting synthetic representations that exhibit low total evidence. Finally, to align local
 109 improvements with global fairness, we introduce a fairness gap-driven server aggregation mechanism
 110 that proportionally rewards clients who successfully close their local performance disparities. By
 111 seamlessly integrating the Deficit Score, the local EBM generative foundation, and the Evidential
 112 Quality Gate, we propose **E²Gen**, an Evidential Energy-based **Generation** framework for Fair FGL.
 113 Our principal contributions are summarized as follows.

- 114 • We introduce the **Deficit Score**, a fused diagnostic metric that explicitly quantifies node-level
 115 inherent vulnerabilities across complementary multi-dimensional axes. Furthermore, we empirically
 116 demonstrate that traditional loss reweighting and SMOTE-like latent interpolation fail to reliably
 117 mitigate these fundamental vulnerabilities.
- 118 • We propose **E²Gen**, a novel generative framework for fair FGL. It leverages conditional EBMs to
 119 synthesize targeted representations for disadvantaged nodes and employs an evidential quality gate
 120 to strictly filter out unreliable generations, ensuring robust and fair local training.
- 121 • We evaluate **E²Gen** on seven benchmarks covering both homophilic and heterophilic graphs.
 122 The results demonstrate that **E²Gen** establishes new state-of-the-art performance, significantly
 123 outperforming existing standard and fairness-aware FGL baselines across various non-IID settings.

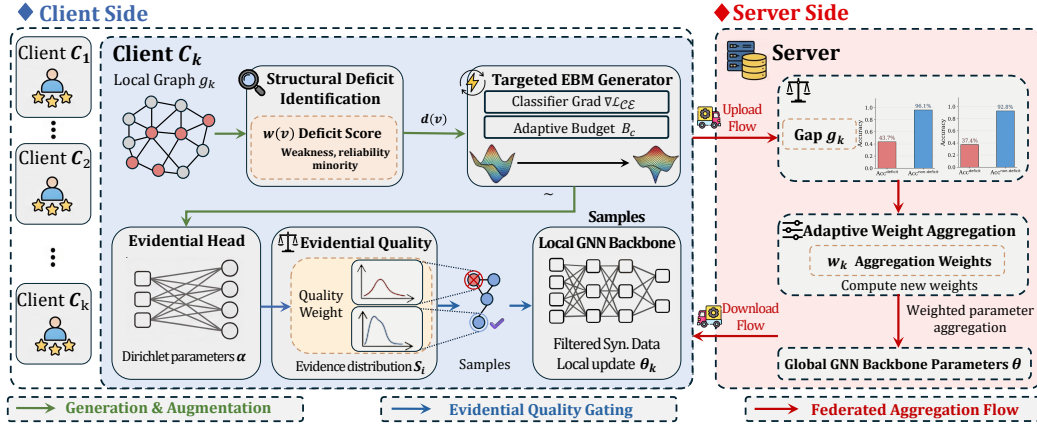


Figure 3: Architecture illustration of E²Gen. We use green, blue, and red arrows to represent the three main flows of our method: Generation & Augmentation, Evidential Quality Gating, and Federated Aggregation Flow, respectively. Best viewed in color. Zoom in for details.

Table 1: Comparison with **state-of-the-art** methods for fair federated graph learning over seven benchmark datasets under Non-IID Louvain partitioning. The upper table presents overall predictive performance, and the lower table reports fairness metrics. The best and second-best results are highlighted with bold and underline, respectively. Please see additional analysis in Section 5.2.

Methods	Cora		CiteSeer		PubMed		Coauthor-CS		Coauthor-Phy		Chameleon		Squirrel	
	\mathcal{A}	\mathcal{F}	\mathcal{A}	\mathcal{F}	\mathcal{A}	\mathcal{F}	\mathcal{A}	\mathcal{F}	\mathcal{A}	\mathcal{F}	\mathcal{A}	\mathcal{F}	\mathcal{A}	\mathcal{F}
<i>Standard FL baselines</i>														
FedAvg	<u>82.64</u>	82.13	72.06	65.82	86.05	85.64	<u>91.11</u>	88.14	94.88	93.27	53.64	52.46	40.60	39.54
MOON	82.37 _{-0.27}	81.84 _{-0.29}	70.66 _{-1.40}	66.38 _{-0.56}	85.67 _{-0.38}	85.23 _{-0.41}	91.02 _{-0.09}	87.96 _{-0.18}	94.82 _{-0.06}	93.12 _{-0.15}	54.05 _{-0.41}	52.64 _{-0.18}	41.73 _{-1.13}	40.68 _{-1.14}
<i>Federated Graph Learning baselines</i>														
FedGTA	82.28 _{-0.36}	81.77 _{-0.36}	<u>73.24</u> _{-1.18}	<u>69.28</u> _{-3.46}	85.99 _{-0.06}	85.53 _{-0.11}	91.08 _{-0.03}	88.04 _{-0.10}	94.93 _{-0.05}	93.32 _{-0.05}	54.05 _{-0.41}	52.08 _{-0.38}	41.17 _{-0.57}	39.97 _{-0.43}
FedSage+	81.38 _{-1.26}	80.49 _{-1.64}	70.73 _{-1.33}	64.88 _{-0.94}	77.14 _{-8.91}	74.28 _{-11.36}	89.46 _{-1.65}	84.13 _{-4.01}	94.37 _{-0.51}	92.68 _{-0.59}	54.47 _{-0.83}	54.49 _{-2.03}	35.15 _{-5.45}	29.37 _{-10.17}
AdaFGL	81.92 _{-0.72}	81.54 _{-0.59}	71.25 _{-0.81}	65.49 _{-0.33}	86.30 _{-0.25}	85.79 _{-0.15}	90.71 _{-0.40}	87.60 _{-0.54}	94.92 _{-0.04}	93.31 _{-0.04}	55.09 _{-1.45}	52.93 _{-0.47}	41.26 _{-0.66}	40.26 _{-0.72}
FedTAD	82.46 _{-0.18}	81.65 _{-0.48}	70.29 _{-1.77}	61.20 _{-4.62}	85.87 _{-0.18}	85.36 _{-0.28}	90.73 _{-0.38}	87.68 _{-0.46}	94.54 _{-0.34}	92.89 _{-0.38}	55.30 _{-1.66}	52.27 _{-0.19}	<u>42.39</u> _{-1.79}	41.29 _{-1.75}
<i>Fairness-aware FGL baselines</i>														
FairFGL	76.80 _{-5.84}	76.84 _{-5.29}	67.11 _{-4.95}	61.89 _{-3.93}	83.15 _{-2.90}	82.85 _{-2.79}	88.96 _{-2.15}	85.64 _{-2.50}	93.47 _{-1.41}	91.36 _{-1.91}	52.60 _{-1.04}	51.16 _{-1.30}	40.79 _{-0.19}	40.64 _{-1.10}
BoostFGL	82.28 _{-0.36}	81.73 _{-0.40}	70.44 _{-1.62}	64.07 _{-1.75}	<u>86.59</u> _{-0.54}	<u>86.13</u> _{-0.49}	91.08 _{-0.03}	<u>88.20</u> _{-0.06}	<u>94.99</u> _{-0.11}	<u>93.41</u> _{-0.14}	52.18 _{-1.46}	52.02 _{-0.44}	40.98 _{-0.38}	39.06 _{-0.48}
<i>Oversampling baselines</i>														
RandomOS	80.04 _{-2.60}	79.74 _{-2.39}	69.70 _{-2.36}	64.13 _{-1.69}	85.97 _{-0.08}	85.60 _{-0.04}	90.37 _{-0.74}	87.34 _{-0.80}	94.25 _{-0.63}	92.35 _{-0.92}	<u>56.13</u> _{-2.49}	<u>55.59</u> _{-3.13}	42.20 _{-1.60}	41.24 _{-1.70}
BalancedOS	81.56 _{-1.08}	81.08 _{-1.05}	72.80 _{-0.74}	66.90 _{-1.08}	85.31 _{-0.74}	84.96 _{-0.68}	90.90 _{-0.21}	87.98 _{-0.16}	94.75 _{-0.13}	93.14 _{-0.13}	54.89 _{-1.25}	53.32 _{-0.86}	42.01 _{-1.41}	<u>41.49</u> _{-1.95}
FedSig	82.10 _{-0.54}	81.60 _{-0.53}	73.10 _{-1.04}	69.23 _{-3.41}	86.45 _{-0.40}	86.00 _{-0.36}	90.88 _{-0.23}	87.70 _{-0.44}	94.74 _{-0.14}	93.09 _{-0.18}	53.43 _{-0.21}	51.08 _{-1.38}	40.70 _{-0.10}	40.51 _{-0.97}
E ² Gen (Ours)	<u>82.91</u> _{-0.27}	<u>82.25</u> _{-0.12}	<u>73.76</u> _{-1.70}	<u>70.24</u> _{-4.42}	<u>86.97</u> _{-0.92}	<u>86.57</u> _{-0.93}	<u>91.25</u> _{-0.14}	<u>88.78</u> _{-0.64}	<u>95.01</u> _{-0.13}	<u>93.58</u> _{-0.31}	<u>58.21</u> _{-4.57}	<u>57.78</u> _{-5.32}	<u>43.14</u> _{-2.54}	<u>43.04</u> _{-3.50}
Methods	Cora		CiteSeer		PubMed		Coauthor-CS		Coauthor-Phy		Chameleon		Squirrel	
	\mathcal{H}	\mathcal{E}	\mathcal{H}	\mathcal{E}	\mathcal{H}	\mathcal{E}	\mathcal{H}	\mathcal{E}	\mathcal{H}	\mathcal{E}	\mathcal{H}	\mathcal{E}	\mathcal{H}	\mathcal{E}
<i>Standard FL baselines</i>														
FedAvg	34.19	79.79	32.56	67.59	51.81	85.84	41.44	86.08	39.85	92.87	57.83	58.57	40.91	42.28
MOON	34.71 _{-0.52}	80.07 _{-0.28}	32.72 _{-0.16}	67.78 _{-0.19}	51.49 _{-0.32}	85.79 _{-0.05}	42.82 _{-1.38}	86.01 _{-0.07}	39.52 _{-0.33}	93.21 _{-0.34}	56.70 _{-1.13}	57.57 _{-1.00}	40.67 _{-0.24}	41.42 _{-0.86}
<i>Federated Graph Learning baselines</i>														
FedGTA	34.99 _{-0.80}	80.04 _{-0.25}	32.69 _{-0.13}	67.71 _{-0.12}	51.74 _{-0.07}	85.68 _{-0.16}	42.87 _{-1.43}	85.54 _{-0.54}	39.35 _{-0.50}	92.90 _{-0.03}	55.18 _{-2.65}	56.02 _{-2.55}	43.07 _{-2.16}	44.14 _{-1.86}
FedSage+	36.81 _{-2.62}	79.32 _{-0.47}	32.61 _{-0.05}	67.59 _{-0.00}	49.07 _{-2.74}	83.51 _{-2.33}	<u>43.51</u> _{-2.07}	85.97 _{-3.01}	39.42 _{-0.43}	91.84 _{-1.03}	50.75 _{-7.08}	52.29 _{-6.28}	30.46 _{-10.45}	34.36 _{-7.92}
AdaFGL	36.38 _{-2.19}	79.46 _{-0.33}	34.52 _{-1.96}	68.24 _{-0.65}	52.00 _{-0.19}	<u>85.88</u> _{-0.04}	42.13 _{-0.69}	<u>86.36</u> _{-0.28}	41.69 _{-1.84}	93.24 _{-0.37}	46.14 _{-11.69}	48.76 _{-9.81}	38.10 _{-2.81}	39.91 _{-2.37}
FedTAD	36.43 _{-2.24}	80.02 _{-0.23}	32.49 _{-0.07}	67.46 _{-0.13}	51.41 _{-0.40}	84.71 _{-1.13}	41.69 _{-0.25}	85.15 _{-0.93}	40.39 _{-0.54}	92.89 _{-0.02}	55.68 _{-2.15}	56.64 _{-1.09}	<u>43.84</u> _{-2.93}	<u>44.52</u> _{-2.24}
<i>Fairness-aware FGL baselines</i>														
FairFGL	34.89 _{-0.70}	76.13 _{-3.66}	31.82 _{-0.74}	63.93 _{-3.66}	<u>53.01</u> _{-1.20}	82.69 _{-3.15}	41.94 _{-0.50}	83.80 _{-2.28}	37.01 _{-2.84}	89.70 _{-3.17}	54.32 _{-3.51}	53.99 _{-4.58}	42.65 _{-1.74}	43.66 _{-1.38}
BoostFGL	35.23 _{-1.04}	79.90 _{-0.11}	31.56 _{-1.00}	67.46 _{-0.13}	52.27 _{-0.46}	85.78 _{-0.06}	42.47 _{-1.03}	86.32 _{-0.24}	40.82 _{-0.97}	92.93 _{-0.06}	54.46 _{-3.37}	56.26 _{-2.31}	40.67 _{-0.24}	41.61 _{-0.67}
<i>Oversampling baselines</i>														
RandomOS	37.02 _{-2.83}	79.05 _{-0.74}	32.60 _{-0.04}	66.84 _{-0.76}	52.27 _{-0.46}	85.17 _{-0.67}	41.18 _{-0.26}	84.74 _{-1.34}	41.01 _{-1.16}	92.15 _{-0.72}	57.34 _{-3.49}	<u>58.73</u> _{-0.16}	41.30 _{-0.39}	43.66 _{-1.38}
BalancedOS	<u>37.22</u> _{-3.03}	79.71 _{-1.08}	<u>35.35</u> _{-2.79}	<u>68.57</u> _{-0.98}	50.68 _{-1.13}	84.62 _{-1.22}	42.97 _{-1.53}	86.00 _{-0.08}	40.88 _{-1.03}	93.28 _{-0.41}	56.82 _{-1.01}	58.38 _{-0.19}	40.10 _{-0.81}	42.08 _{-0.20}
FedSig	33.87 _{-0.32}	80.84 _{-1.05}	33.04 _{-0.48}	68.01 _{-0.42}	52.96 _{-1.15}	85.47 _{-0.37}	41.82 _{-0.38}	86.10 _{-0.02}	<u>42.75</u> _{-2.90}	<u>93.75</u> _{-0.88}	55.56 _{-2.27}	56.20 _{-2.37}	39.70 _{-1.21}	40.85 _{-1.43}
E ² Gen (Ours)	<u>40.69</u> _{-6.50}	<u>81.70</u> _{-1.91}	<u>36.76</u> _{-4.20}	<u>70.01</u> _{-2.42}	<u>54.40</u> _{-2.59}	<u>86.11</u> _{-0.27}	<u>45.16</u> _{-3.72}	<u>87.08</u> _{-1.00}	<u>43.28</u> _{-3.43}	<u>93.85</u> _{-0.98}	<u>57.45</u> _{-0.38}	<u>58.77</u> _{-0.20}	<u>44.78</u> _{-3.87}	<u>45.47</u> _{-3.19}

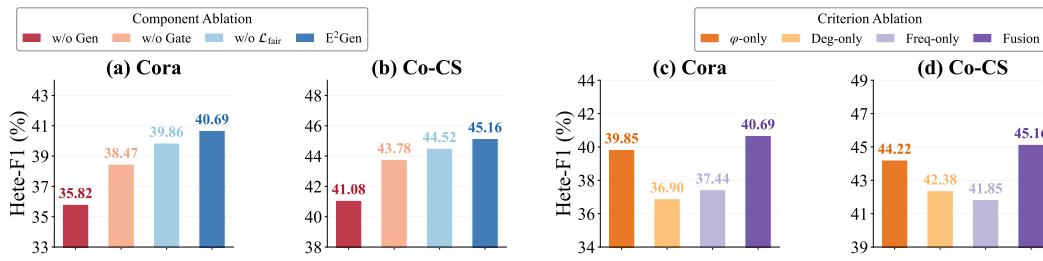


Figure 4: **Ablation Analysis** of components and criteria for E^2 Gen on the Cora and Co-CS datasets. The results illustrate the fairness performance (Hete-F1) under Component Ablation (a-b) and Criterion Ablation (c-d).