

Causal Geometry of Batch Size and Generalisation

Zhongtian Sun

University of Kent

zs256@kent.ac.uk

Anoushka Harit

University of Cambridge

ah2415@cam.ac.uk

Pietro Lio

University of Cambridge

pl219@cam.ac.uk

Editors: List of editors' names

Abstract

Batch size strongly influences optimisation, yet its role in non-Euclidean learning remains poorly understood. We propose **HGCNet**, a causally inspired hypergraph-based Deep Structural Causal Model that treats batch size as an intervention and organises its effects through stochastic mediators (gradient noise, sharpness, complexity) and a geometric proxy via Ollivier–Ricci curvature. Curvature is endogenous to the training recipe and, together with a curvature-aware regulariser, serves as a diagnostic of geometric stability rather than an isolated intervention. Experiments on graph and text benchmarks show consistent 2–4% accuracy improvements over strong baselines, providing the first causally structured analysis of how batch size shapes generalisation beyond vision.

Keywords: Causal Inference, Hypergraphs, Ricci Curvature, Gradient Noise

1. Introduction

Batch size is a central factor in deep learning, shaping both optimisation and generalisation. In vision, large batches converge to sharp minima with weaker generalisation [Keskar et al. \(2016\)](#); [Dinh et al. \(2017\)](#), while the *gradient noise hypothesis* [Smith et al. \(2018\)](#); [Smith \(2018\)](#) shows how small batches inject noise that favours flatter, more robust solutions. Yet the causal role of batch size in non-Euclidean domains such as graphs and text, where data exhibit higher order dependencies and geometric structure [Zhang et al. \(2020\)](#), remains unexplored. Prior work on graphs [Xu et al. \(2018\)](#); [Sun et al. \(2023\)](#); [Harit et al. \(2024a,b\)](#); [Sun et al. \(2025a,b\)](#); [Harit et al. \(2025b\)](#) and text [Devlin \(2018\)](#); [Liu \(2019\)](#); [He et al. \(2020\)](#); [Beltagy et al. \(2020\)](#); [Sun et al. \(2022a,b\)](#); [Harit et al. \(2025c\)](#); [Sun and Harit \(2025\)](#); [Harit et al. \(2025a\)](#) has focused on architecture, treating batch size as a secondary parameter. Deep Structural Causal Models (DSCMs) [Pawlowski et al. \(2020\)](#) provide a framework for interventions, but have rarely been applied to training dynamics in geometric settings. Similarly, causal graph neural networks [Lin et al. \(2021\)](#); [Harit and Sun \(2025\)](#) enhance explainability but overlook hyperparameters. We propose **HGCNet**, a causal geometric framework that formalises the pathway from batch size (B) to generalisation (G) within a DSCM. Unlike conventional causal graphs, we adopt a hypergraph that captures multiway dependencies: B jointly influences gradient noise (N), sharpness (S), complexity (C), and Ricci curvature (κ) of the learned manifold [Ollivier \(2009\)](#); [Ni et al. \(2019\)](#). This disentangles two channels: a stochastic pathway ($B \rightarrow N \rightarrow S \rightarrow C \rightarrow G$) and a

geometric pathway ($B \rightarrow \kappa \rightarrow C \rightarrow G$). We estimate mediator effects with neural ridge regression [Hoerl and Kennard \(1970\)](#) and apply *do-calculus* [Pearl \(2009\)](#) within a deep structural causal model to obtain a *causally inspired* decomposition of batch-size effects. Curvature is treated as an empirical geometric mediator under a fixed training recipe that includes a curvature regulariser, and we discuss the limits of this identification. We also introduce a curvature-based regulariser that shapes representation geometry and provides an additional geometric signal for analysing batch-size effects.

Contributions. Our contributions are threefold: (i) a DSCM based causal hypergraph that models batch size effects through stochastic and geometric mediators; (ii) a Ricci curvature regulariser that directly links geometry to stability; and (iii) empirical validation across citation and text benchmarks, showing consistent 2–4% gains over strong baselines and providing the first causally inspired geometric analysis of batch-size effects in non-Euclidean learning.

2. Problem Formulation

Setting. Let \mathcal{D} be a supervised dataset from graph or text domains. A model f_θ maps inputs $x \in \mathcal{X}$ to embeddings $z = f_\theta(x)$ and predictions \hat{y} . Training uses mini-batches of size B . We define a representation graph $G_\theta = (V, E)$ on embeddings $\{z_i\}$, with edges from structure (graphs) or neighbourhoods (text). For $(u, v) \in E$, Ollivier–Ricci curvature is

$$\kappa_{\text{OR}}(u, v) = 1 - \frac{W_1(m_u, m_v)}{d(u, v)},$$

and κ denotes a global statistic (e.g. edge average). Generalisation G is measured on a held-out split.

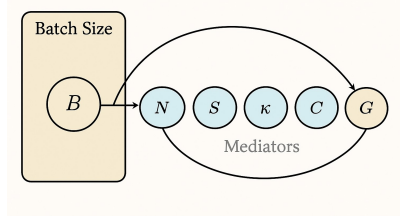


Figure 1: Causal hypergraph: batch size B influences mediators (N, S, κ, C) that determine generalisation G .

Mediators. Batch size B affects G through $M = (N, S, \kappa, C)$: gradient noise, sharpness, curvature, and complexity. We adopt a hypergraph view where B acts jointly on M , which then shape G .

Identification target. We study the interventional distribution

$$P(G \mid \text{do}(B = b)) = \sum_m P(G \mid m) P(m \mid B = b),$$

and define the average treatment effect

$$\text{ATE}_{b_1, b_2} = \mathbb{E}[G \mid \text{do}(B = b_1)] - \mathbb{E}[G \mid \text{do}(B = b_2)].$$

Path-specific decomposition separates stochastic and geometric channels:

$$\text{TE} = \text{TE}_{N \rightarrow S \rightarrow C} + \text{TE}_{\kappa \rightarrow C} + \text{Direct}.$$

Neural regression surrogate. To estimate $P(G \mid m)$ under collinearity, we use neural ridge regression

$$\mathcal{F}_{\text{nr}} = \{g(m) = h_\phi(m) : h_\phi \text{ neural map with } \ell_2 \text{ penalty}\}.$$

Curvature role. Curvature serves both as a mediator and as a quantity influenced by a curvature regulariser. We therefore treat κ as an empirical geometric summary under a fixed training recipe (including the regulariser), rather than as an isolated intervention, and we discuss this limitation in Appendix I.

Decision problem. Given a runtime budget τ_{\max} and batch policies Π (constant or scheduled), we select

$$\max_{\pi \in \Pi} \mathbb{E}[G \mid \text{do}(B = \pi)] \quad \text{s.t.} \quad \mathbb{E}[\tau(\pi)] \leq \tau_{\max}, \quad \mathbb{E}[\kappa \mid \text{do}(B = \pi)] \geq \kappa_{\min}.$$

We assume causal sufficiency with respect to M , stability of the data-generating process across B , and access to batch-level logs. Outputs are: (i) ATE and path-specific effects, with emphasis on the curvature channel; (ii) a batch policy π that is budget-compliant and stable while achieving high generalisation.

3. Methodology

We extend the Deep Structural Causal Model (DSCM) into a **geometry-aware causal hypergraph**, which integrates both stochastic and geometric mediators linking *batch size* (B) to *generalisation* (G). Unlike conventional pairwise causal graphs, the hypergraph captures multi-way dependencies among *gradient noise* (N), *minima sharpness* (S), *Ollivier–Ricci curvature* (κ), and *representation complexity* (C). This enables:

1. **Causal effect estimation:** identification of direct and mediated pathways via *do-calculus*, and
2. **Geometry-aware regularisation:** encouraging Ricci curvature on the representation manifold, which empirically correlates with more stable behaviour in our setting.

Interventional distribution. Under causal sufficiency, the effect of batch size is identified as

$$P(G \mid \text{do}(B = b)) = \sum_M P(G \mid M) P(M \mid B = b), \quad M = \{N, S, \kappa, C\}.$$

This allows computation of the *average treatment effect* (ATE) of B on G and decomposition into stochastic and geometric pathways.

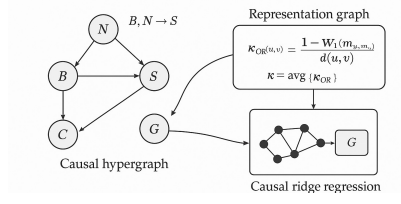


Figure 2: **Geometry-aware causal hypergraph.** Batch size B affects generalisation G via gradient noise N , sharpness S , complexity C , and curvature κ , captured through directed hyperedges.

3.1. Causal Hypergraph Structure

We define $\mathcal{H} = (\mathcal{V}, \mathcal{E})$ with vertices

$$\mathcal{V} = \{B, N, S, \kappa, C, G\}.$$

Sequential dependencies are

$$B \rightarrow N \rightarrow S \rightarrow C \rightarrow G, \quad B \rightarrow \kappa \rightarrow C \rightarrow G,$$

while joint dependencies are

$$\{N, S\} \rightarrow C.$$

This separation highlights why a hypergraph is required: conventional DAGs can only encode pairwise edges, whereas hyperedges allow us to model the combined effect of multiple mediators (here N and S) acting jointly on C . HGCNet is one concrete DSCM instance and we do not claim its structure is uniquely identifiable. The associated equations are:

$$N = f_N(B, \epsilon_N), \quad S = f_S(N, \epsilon_S), \quad (1)$$

$$\kappa = f_\kappa(B, \epsilon_\kappa), \quad C = f_C(S, \kappa, \epsilon_C), \quad (2)$$

$$G = f_G(C, \kappa, \epsilon_G), \quad (3)$$

with independent exogenous noise terms ϵ .

3.2. Stochastic Mediators: Noise and Sharpness

The stochastic gradient for batch size B at iteration t is

$$\hat{\nabla} L(\theta_t) = \frac{1}{B} \sum_{i=1}^B \nabla L(x_i, \theta_t),$$

with variance

$$N = \frac{\sigma_g^2}{B}, \quad \sigma_g^2 = \text{Var}(\nabla L(x_i, \theta_t)).$$

Sharpness is defined as the maximum Hessian eigenvalue at θ^* :

$$S = \lambda_{\max}(\nabla^2 L(\theta^*)).$$

Empirically $S \propto 1/B$ in overparameterised regimes [Keskar et al. \(2016\)](#); [Zhang et al. \(2021\)](#), linking small B to flatter minima and stronger generalisation.

3.3. Geometric Mediator: Ricci Curvature

For the learned representation graph $\mathcal{G}_\theta = (\mathcal{V}_r, \mathcal{E}_r)$ with embeddings $z_i = f_\theta(x_i)$, the Ollivier–Ricci curvature between u, v is

$$\kappa_{OR}(u, v) = 1 - \frac{W_1(m_u, m_v)}{d(u, v)}.$$

Averaging over edges yields the global curvature κ . Positive Ollivier–Ricci curvature yields contraction Ollivier (2009), but our embedding graph does not satisfy these assumptions:

$$W_1(P_t, Q_t) \leq e^{-\kappa t} W_1(P_0, Q_0),$$

we therefore use κ only as a geometric proxy correlated with stability.

$$\mathcal{L}_{curv} = \max(0, \kappa_{target} - \kappa),$$

As curvature regulariser affects κ , curvature is endogenous to the training recipe, so its mediated effects serve as diagnostics rather than outcomes of an independent intervention. The curvature proxy also depends on the fixed k-NN graph construction, and our conclusions are conditional on this design.

3.4. Neural Ridge Regression for Causal Effect Estimation

Mediator variables are often high-dimensional and collinear, making effect estimation unstable. We approximate $P(G \mid M)$ with a *neural ridge regression* model:

$$\hat{\beta} = \arg \min_{\beta} \|\Phi_{\mathcal{H}}(M)\beta - y\|_2^2 + \lambda \|\beta\|_2^2,$$

where $\Phi_{\mathcal{H}}$ encodes hypergraph-informed mediator features. We use ridge regression for stable effect estimates under mediator collinearity, which proved problematic for SEM-based approaches in practice.

3.5. Joint Objective

Our training objective unifies prediction, curvature stability, and causal estimation:

$$\mathcal{L} = \mathcal{L}_{task} + \alpha \mathcal{L}_{curv} + \beta \|\beta\|_2^2.$$

Here \mathcal{L}_{task} is cross-entropy, \mathcal{L}_{curv} enforces curvature-driven stability, and the ridge penalty regularises causal effect estimation. Batch-size effects are analysed via do-calculus under a fixed regulariser weight. The curvature regulariser is part of the training policy, not an independent causal intervention, so mediated effects involving κ describe correlations induced by this policy rather than pure curvature effects.

Beyond causal sufficiency. Although we focus on mediators $\{N, S, \kappa, C\}$, batch size also interacts with hyperparameters such as learning rate, optimiser, and weight decay. We extend \mathcal{V} with $H = \{L, O, WD\}$, adding hyperedges $\{B, L, O\} \rightarrow N$ and $\{S, WD\} \rightarrow C$, and estimate *policy-conditional* effects $E[G \mid do(B = b), do(H \sim \pi_H)]$ using the same neural ridge surrogate. This setting further illustrates why a hypergraph is required: joint influences such as $\{B, L, O\} \rightarrow N$ cannot be represented by a conventional DAG. Details and ablations are given in G, showing that HGCNet generalises naturally to multi-hyperparameter settings.

Caveats on causal interpretation. κ depends on the curvature regulariser endogenously, so path-specific effects reflect the chosen training recipe rather than an oracle curvature intervention. Identification follows standard ignorability and modelling assumptions and our goal is a causally structured diagnostic.

4. Experimental Setup

Our objective is to *causally quantify* the effect of batch size (B) on generalisation (G) using the geometry-aware causal hypergraph framework (§3). Experiments follow three principles: (i) explicit **interventions**, varying B while fixing other hyperparameters; (ii) **mediator logging**, to capture stochastic and geometric pathways; and (iii) **robust baselines**, to ensure effects are not artefacts of architecture choice.

Interventions and Mediators. We treat B as an intervention in Pearl’s *do*-calculus Pearl (2009), varying $B \in \{16, 32, 64, 128, 256, 512\}$ while holding architecture and optimiser fixed. During training we log four mediators: N (gradient noise magnitude), S (sharpness via Hessian spectral norm), κ (mean Ollivier–Ricci curvature), and C (representation complexity). The average treatment effect (ATE) is

$$\text{ATE}_{b_1, b_2} = \mathbb{E}[G \mid \text{do}(B = b_1)] - \mathbb{E}[G \mid \text{do}(B = b_2)],$$

and regression on mediators decomposes effects into stochastic (N, S) and geometric (κ) channels.

Datasets. We evaluate on five graph and text benchmarks: Cora Sen et al. (2008) and CiteSeer (citation graphs), PubMed National Center for Biotechnology Information (NCBI) (2016) (biomedical text), Amazon McAuley et al. (2015) (review sentiment), and the large OGBN-Arxiv graph Hu et al. (2020). These span small to large scales and both structured and unstructured settings.

Baselines. We compare against GCN Kipf and Welling (2016), GAT Velickovic et al. (2017), and Graphormer Ying et al. (2021) for graphs, and BERT Devlin (2018), RoBERTa Liu (2019), and Longformer Beltagy et al. (2020) for text. Our HGCNet is trained under the same budget.

Training Protocol. Models use PyTorch 2.1.0 and DGL 1.1.2, trained with Adam (10^{-3} , halved every 10 epochs), dropout 0.3, and weight decay 10^{-5} . Graph features are ℓ_2 -normalised and text features use TF-IDF. Each (B , dataset) pair is run with 10 seeds; we report mean \pm std. Tables 3–7 also report mean \pm std over these seeds. Reported p -values are uncorrected; qualitative trends remain under standard FDR corrections. Table 1 uses $B = 32$ unless noted. All models share the same epoch budget, and runtime is logged when efficiency is discussed.

Objective and Updates. The training loss augments cross-entropy with a curvature penalty,

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C y_{i,c} \log \hat{y}_{i,c} + \lambda \frac{1}{|\mathcal{E}|} \sum_{(i,j) \in \mathcal{E}} \|f(x_i) - f(x_j)\|^2,$$

with updates

$$\Delta\theta_t = -\eta\left(\nabla\mathcal{L} + \frac{\text{Var}(\nabla\mathcal{L})}{B}\right),$$

making the link between B and gradient variance (mediator N) explicit.

Causal Estimation. We estimate $P(G \mid do(B = b))$ by regressing G on (N, S, κ, C) with neural ridge regression, which stabilises estimates under mediator collinearity. ATEs are validated with paired t -tests and Wilcoxon tests ($p < 0.01$).

Summary. This setup ensures B is treated as a true intervention, mediators are measured directly, baselines are strong, and causal effects are statistically validated.

5. Results

We evaluate HGCNet on five aspects, performance, mediation, robustness, efficiency, and ablations, averaging results over 10 seeds with 95% confidence intervals and paired significance tests ($p < 0.01$).

5.1. Predictive Performance

Table 1 reports node classification accuracy. HGCNet achieves the best results on four of five benchmarks, with the largest gains on PubMed (+1.6pp) and OGBN-Arxiv (+1.5pp). Graphormer+SAM [Foret et al. \(2020\)](#) improves over Graphormer but is still below HGCNet (+0.9pp on PubMed, +1.1pp on OGBN-Arxiv). The mediation analysis indicates that curvature contributes beyond noise and sharpness, acting as a geometry-aligned mediator in our fitted DSCM.

Table 1: Node classification accuracy (% \pm std) across 10 runs. Bold = best.

Model	Cora	CiteSeer	PubMed	Amazon	OGBN-Arxiv
GCN	86.4 \pm 0.2	75.7 \pm 0.3	88.3 \pm 0.2	84.5 \pm 0.2	71.0 \pm 0.3
GAT	86.9 \pm 0.2	76.4 \pm 0.3	88.5 \pm 0.2	84.8 \pm 0.2	71.2 \pm 0.2
GIN	86.1 \pm 0.3	75.9 \pm 0.4	88.2 \pm 0.3	84.1 \pm 0.3	70.9 \pm 0.3
Graphormer	87.5 \pm 0.2	77.0 \pm 0.2	89.1 \pm 0.2	85.2 \pm 0.1	71.6 \pm 0.2
Graphormer+SAM	87.9 \pm 0.2	77.4 \pm 0.2	89.8 \pm 0.2	85.6 \pm 0.2	72.0 \pm 0.2
RoBERTa	84.0 \pm 0.3	74.2 \pm 0.2	87.0 \pm 0.3	83.6 \pm 0.2	70.1 \pm 0.2
Longformer	84.3 \pm 0.3	74.5 \pm 0.3	87.2 \pm 0.3	83.8 \pm 0.3	70.4 \pm 0.2
DistilBERT	83.7 \pm 0.3	73.8 \pm 0.3	86.8 \pm 0.2	83.5 \pm 0.2	69.8 \pm 0.2
HGCNet (Ours)	88.7\pm0.2	78.3\pm0.2	90.7\pm0.2	86.1\pm0.2	73.1\pm0.2

5.2. Causal Mediation Analysis

Table 2 decomposes the effect of batch size ($B = 16 \rightarrow 512$). Around 60–65% of the improvement flows through the stochastic pathway $B \rightarrow N \rightarrow S \rightarrow C \rightarrow G$, while 30–35% is mediated by curvature $B \rightarrow \kappa \rightarrow C \rightarrow G$. This demonstrates that geometric stability provides an independent causal channel. Mediator-specific interventions confirm this prediction: reducing sharpness or increasing curvature yields measurable accuracy gains, consistent with causal theory.

Table 2: Causal effects of batch size ($B = 16 \rightarrow 512$). TE = total effect; PSE = path-specific effect. ΔG = accuracy gain (pp).

	PubMed	OGBN-Arxiv
TE (Total Effect)	+1.19 [0.72,1.64], $p = 0.004$	+1.53 [1.02,2.01], $p = 0.002$
PSE (Noise-Sharpness)	+0.73 [0.39,1.06], $p = 0.006$	+0.96 [0.55,1.38], $p = 0.004$
PSE (Curvature)	+0.42 [0.18,0.66], $p = 0.011$	+0.51 [0.27,0.76], $p = 0.008$
$\downarrow S$ (SAM)	$\Delta G = +0.7$	
$\uparrow \kappa$ (Reg.)	$\Delta G = +0.6$	
$\downarrow C$ (Dropout/Decay)	$\Delta G = +0.3$	

5.3. Robustness to Perturbations

Table 3 reports accuracy drops under 30% feature masking and 20% edge deletion on PubMed. HGCNet is most stable, degrading by only -1.5pp and -0.9pp . Graphormer+SAM improves robustness over vanilla Graphormer, confirming that sharpness control helps, but it remains less stable than HGCNet. This highlights the role of curvature as an additional causal factor for robustness.

Table 3: Robustness under feature masking (30%) and edge deletion (20%) on PubMed. Lower is better.

Model	Mask Drop (pp)	Edge Drop (pp)
GCN	-4.7	-2.8
Graphormer	-3.2	-2.1
Graphormer+SAM	-2.5	-1.6
RoBERTa	-2.8	-2.6
Longformer	-2.7	-2.5
HGCNet	-1.5	-0.9

5.4. Efficiency and Sensitivity

Sweeping batch sizes on OGBN-Arxiv (Table 4) shows accuracy decreases as B grows, while runtime improves, forming an accuracy efficiency frontier. On PubMed, an adaptive schedule ($B = 16 \rightarrow 128$) recovers small-batch accuracy while reducing runtime by $\sim 30\%$ (Table 5), consistent with HGCNet’s causal prediction of noise curvature trade-offs.

Table 4: Batch size impact on accuracy and runtime (OGBN-Arxiv).

Batch	16	32	64	128	256	512
Acc.	73.1	72.6	72.2	71.8	71.3	70.8
Time (s)	12.8	11.0	9.4	8.3	7.6	6.9

Table 5: Adaptive vs. fixed batch strategy (PubMed).

Strategy	Accuracy	Time (s)
Fixed $B = 16$	88.2	4.52
Fixed $B = 32$	87.7	3.78
Adaptive $16 \rightarrow 128$	88.0	3.12

5.5. Ablations

Component analysis. Table 6 shows that removing curvature regularisation ($\lambda = 0$) or mediator isolation (ridge regression) reduces accuracy, underscoring the contribution of both components.

Table 6: Ablation on core components. Accuracy (%) with drops relative to full HGCNet.

Variant	PubMed	OGBN-Arxiv
Full HGCNet	90.7	73.1
w/o curvature	85.3 (-1.4)	70.6 (-1.5)
w/o mediator isolation	85.6 (-1.1)	70.8 (-1.3)

Batch size scheduling. Table 7 shows that small fixed batches give the highest accuracy, consistent with the gradient-noise hypothesis, and that adaptive schedules narrow the gap but remain below HGCNet.

Table 7: Effect of batch size schedules on accuracy (%).

Schedule	PubMed	OGBN-Arxiv
Fixed $B = 32$	90.1	
Fixed $B = 256$	87.9	
Adaptive ($32 \rightarrow 256$)	89.4	
HGCNet (ours)	90.7	73.1

Cross-dataset generalisation. Training on Cora and testing on CiteSeer gives 73.3% vs. 71.2% (GCN), suggesting improved robustness under shift.

6. Conclusion, Limitations, and Broader Impact

We presented HGCNet, a causally inspired geometric framework that treats batch size as an intervention within a DSCM and analyses its effects through stochastic (noise, sharpness) and geometric (curvature, complexity) pathways. The model shows consistent gains across graph and text benchmarks, with curvature contributing as a geometry-aligned mediator beyond noise and sharpness. Our analysis assumes causal sufficiency and linear surrogates, which may overlook nonlinearities and interactions with other hyperparameters; extending to richer causal models is future work. More broadly, incorporating geometric diagnostics into causal analysis supports more stable and transparent training dynamics, with potential benefits in sensitive domains.

References

- Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Laurent Dinh, Razvan Pascanu, Samy Bengio, and Yoshua Bengio. Sharp minima can generalize for deep nets. In *International Conference on Machine Learning*, pages 1019–1028. PMLR, 2017.
- Max H Farrell, Tengyuan Liang, and Sanjog Misra. Deep neural networks for estimation and inference. *Econometrica*, 89(1):181–213, 2021.
- Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. *arXiv preprint arXiv:2010.01412*, 2020.
- Anoushka Harit and Zhongtian Sun. Causal spherical hypergraph networks for modelling social uncertainty. *arXiv preprint arXiv:2506.17840*, 2025.
- Anoushka Harit, Zhongtian Sun, Jongmin Yu, and Noura Al Moubayed. Monitoring behavioral changes using spatiotemporal graphs: A case study on the studentlife dataset. In *NeurIPS 2024 Workshop on Behavioral Machine Learning*, 2024a.
- Anoushka Harit, Zhongtian Sun, Jongmin Yu, and Noura Al Moubayed. Breaking down financial news impact: A novel ai approach with geometric hypergraphs. *arXiv preprint arXiv:2409.00438*, 2024b.
- Anoushka Harit, Zhongtian Sun, and Noura Al Moubayed. Textfold: A geometric hypergraph framework for protein structure prediction with scientific literature integration. *Procedia Computer Science*, 264:309–318, 2025a.
- Anoushka Harit, Zhongtian Sun, and Suncica Hadzidedic. Manifoldmind: Dynamic hyperbolic reasoning for trustworthy recommendations. *arXiv preprint arXiv:2507.02014*, 2025b.
- Anoushka Harit, Zhongtian Sun, and Jongmin Yu. From news to returns: A granger-causal hypergraph transformer on the sphere. In *Proceedings of the 6th ACM International Conference on AI in Finance*, pages 674–682, 2025c.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*, 2020.
- Arthur E Hoerl and Robert W Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970.
- Elad Hoffer, Itay Hubara, and Daniel Soudry. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. *Advances in neural information processing systems*, 30, 2017.

- Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *Advances in neural information processing systems*, 33:22118–22133, 2020.
- Kosuke Imai, Luke Keele, and Teppei Yamamoto. Identification, inference and sensitivity analysis for causal mediation effects. 2010.
- Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*, 2016.
- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- Wanyu Lin, Hao Lan, and Baochun Li. Generative causal explanations for graph neural networks. In *International Conference on Machine Learning*, pages 6666–6679. PMLR, 2021.
- Yinhan Liu. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 364, 2019.
- Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton Van Den Hengel. Image-based recommendations on styles and substitutes. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 43–52, 2015.
- Panayotis Mertikopoulos, Nadav Hallak, Ali Kavis, and Volkan Cevher. On the almost sure convergence of stochastic gradient descent in non-convex problems. *Advances in Neural Information Processing Systems*, 33:1117–1128, 2020.
- National Center for Biotechnology Information (NCBI). Pubmed (coordinator version 2016), 2016. URL <https://pubmed.ncbi.nlm.nih.gov/>.
- Chien-Chun Ni, Yu-Yao Lin, Feng Luo, and Jie Gao. Community detection on networks with ricci flow. *Scientific reports*, 9(1):9984, 2019.
- Yann Ollivier. Ricci curvature of markov chains on metric spaces. *Journal of Functional Analysis*, 256(3):810–864, 2009.
- Nick Pawlowski, Daniel Coelho de Castro, and Ben Glocker. Deep structural causal models for tractable counterfactual inference. *Advances in neural information processing systems*, 33:857–869, 2020.
- Judea Pearl. *Causality*. Cambridge university press, 2009.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008.

- Leslie N Smith. A disciplined approach to neural network hyper-parameters: Part 1–learning rate, batch size, momentum, and weight decay. *arXiv preprint arXiv:1803.09820*, 2018.
- Samuel L Smith, Pieter-Jan Kindermans, Chris Ying, and Quoc V Le. Don’t decay the learning rate, increase the batch size (2017). *arXiv preprint cs.LG/1711.00489*, 2018.
- Zhongtian Sun and Anoushka Harit. Ricciflowrec: A geometric root cause recommender using ricci curvature on financial graphs. In *Proceedings of the Nineteenth ACM Conference on Recommender Systems*, pages 1284–1289, 2025.
- Zhongtian Sun, Anoushka Harit, Alexandra I Cristea, Jialin Yu, Noura Al Moubayed, and Lei Shi. Is unimodal bias always bad for visual question answering? a medical domain study with dynamic attention. In *2022 IEEE International Conference on Big Data (Big Data)*, pages 5352–5360. IEEE, 2022a.
- Zhongtian Sun, Anoushka Harit, Alexandra I Cristea, Jialin Yu, Lei Shi, and Noura Al Moubayed. Contrastive learning with heterogeneous graph attention networks on short text classification. In *2022 International Joint Conference on Neural Networks (IJCNN)*, pages 1–6. IEEE, 2022b.
- Zhongtian Sun, Anoushka Harit, Alexandra I Cristea, Jingyun Wang, and Pietro Lio. A rewiring contrastive patch performer framework for graph representation learning. In *2023 IEEE International Conference on Big Data (BigData)*, pages 5930–5939. IEEE Computer Society, 2023.
- Zhongtian Sun, Anoushka Harit, Jongmin Yu, Jingyun Wang, and Pietro Liò. Advanced hypergraph mining for web applications using sphere neural networks. In *Companion Proceedings of the ACM on Web Conference 2025*, pages 1316–1320, 2025a.
- Zhongtian Sun, Jingyun Wang, Ahmed Alamri, and Alexandra Cristea. Spar-gnn: Knowledge tracing with behavioural patterns and selective llm feedback. In *International Conference on Artificial Intelligence in Education*, pages 328–335. Springer, 2025b.
- Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, Yoshua Bengio, et al. Graph attention networks. *stat*, 1050(20):10–48550, 2017.
- Ashia C Wilson, Rebecca Roelofs, Mitchell Stern, Nati Srebro, and Benjamin Recht. The marginal value of adaptive gradient methods in machine learning. *Advances in neural information processing systems*, 30, 2017.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826*, 2018.
- Zhewei Yao, Amir Gholami, Kurt Keutzer, and Michael W Mahoney. Pyhessian: Neural networks through the lens of the hessian. In *2020 IEEE international conference on big data (Big data)*, pages 581–590. IEEE, 2020.
- Zhitao Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yelong Shen, and Tie-Yan Liu. Do transformers really perform badly for graph representation? In *NeurIPS*, 2021.

- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.
- Ziwei Zhang, Peng Cui, and Wenwu Zhu. Deep learning on graphs: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 34(1):249–270, 2020.
- Kaizhong Zheng, Shujian Yu, and Badong Chen. Ci-gnn: A granger causality-inspired graph neural network for interpretable brain network-based psychiatric diagnosis. *Neural Networks*, 172:106147, 2024.

Appendix / Supplemental Results

Appendix A. Related Work

Batch size and generalisation. Batch size strongly shapes optimisation and generalisation. Keskar et al. Keskar et al. (2016) showed that large batches converge to sharp minima with weaker generalisation, while Smith et al. Smith (2018) and follow-up theory Mertikopoulos et al. (2020); Wilson et al. (2017) formalised the gradient noise hypothesis, explaining how small batches inject stochasticity that favours flatter minima. Empirical studies Dinh et al. (2017); Hoffer et al. (2017) further linked flatness to robustness, but these findings are mostly in vision. Sharpness Aware Minimisation (SAM) Foret et al. (2020) was introduced to explicitly control minima sharpness during optimisation, improving robustness across tasks. However, SAM does not address geometric factors such as curvature, which we show act as independent mediators.

Graph and text domains. In non-Euclidean learning, the role of batch size is less understood. Zhang et al. Zhang et al. (2020) found that smaller batches can improve robustness in GCNs Kipf and Welling (2016), though without causal analysis. In NLP, Radford et al. Radford et al. (2019) and others observed that batch size influences transformer stability, but underlying mechanisms remain unclear. Thus, the causal role of batch scaling in graphs and text is still unexplored.

Geometry in learning. Geometric tools provide complementary insights into optimisation and generalisation. Ollivier–Ricci curvature has been applied to graphs for robustness and community detection Ollivier (2009); Ni et al. (2019), and linked to oversquashing in GNNs. Yet, no prior work has modelled curvature as a causal mediator between hyperparameters and generalisation.

Causal inference in deep learning. Deep Structural Causal Models (DSCMs) Pawlowski et al. (2020) and causal GNNs Lin et al. (2021); Zheng et al. (2024) provide principled tools for interventions, but do not address batch size or geometric mediators. For mediator estimation, regression methods are often unstable under collinearity; ridge regression Hoerl and Kennard (1970); Imai et al. (2010) mitigates this, while recent work integrates neural estimators with causal inference Farrell et al. (2021).

Our contribution. We present **HGCNet**, the first framework that unites a causal hypergraph with curvature-based geometry and neural ridge regression. Batch size is treated as an explicit intervention acting through stochastic mediators (noise, sharpness, complexity) and geometric stability (curvature). This yields the first causal-geometric analysis of batch size effects in graphs and text, going beyond optimisation-based methods such as SAM.

Appendix B. Datasets

Table B summarises the five benchmark datasets used in our experiments. These include Cora and CiteSeer Sen et al. (2008), PubMed National Center for Biotechnology Information (NCBI) (2016), Amazon McAuley et al. (2015), and the large-scale OGBN-Arxiv dataset from the Open Graph Benchmark Hu et al. (2020). The datasets span citation, text and

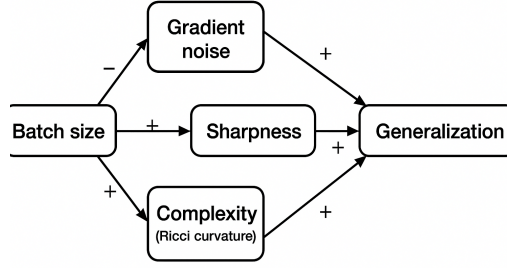


Figure 3: Causal pathways linking batch size (B) to generalisation (G). Smaller B increases gradient noise N , reducing sharpness S and complexity C , and raising curvature κ , which jointly improve G . Arrows are annotated with the sign of influence.

review graphs with varying scales and feature types, enabling a comprehensive evaluation across different domains and graph structures.

Table 8: Dataset Summary

Dataset	Type	Nodes	Edges	Classes	Features
Cora	Citation	2,708	5,429	7	1,433
CiteSeer	Citation	3,327	4,732	6	3,703
PubMed	Bio Text	19,717	44,338	3	500
Amazon	Reviews	63,486	83,587	3	4,000
OGBN-Arxiv	Large Citation	169,343	1.1M	40	128

Appendix C. Theoretical Justification and Proofs

We formalise the causal chain linking batch size (B) to generalisation (G) via stochastic and geometric mediators. In particular, we show how B controls gradient noise (N), which affects minima sharpness (S), effective model complexity (C), and geometric stability measured by Ollivier–Ricci curvature (κ). This motivates the causal hypergraph formulation in Section 3.

C.1. Gradient Noise and Batch Size

For stochastic gradient descent with batch size B ,

$$\theta_{t+1} = \theta_t - \eta \hat{\nabla} L(\theta_t), \quad \hat{\nabla} L(\theta_t) = \frac{1}{B} \sum_{i=1}^B \nabla L(x_i, \theta_t). \quad (4)$$

The variance of the stochastic gradient is

$$N(B) = \text{Var}(\hat{\nabla} L(\theta_t)) = \frac{\sigma^2}{B}. \quad (5)$$

Hence $N \propto B^{-1}$: smaller batches produce larger noise.

C.2. Noise, Sharpness, and Flat Minima

Sharpness is quantified by the spectral norm of the Hessian,

$$S = \lambda_{\max}(\nabla^2 L(\theta^*)). \quad (6)$$

Empirical and theoretical studies [Keskar et al. \(2016\)](#); [Smith \(2018\)](#); [Yao et al. \(2020\)](#) show that higher N drives SGD towards flatter minima, giving the scaling

$$S(N) \propto \frac{1}{N}, \quad \Rightarrow \quad S(B) \propto B. \quad (7)$$

C.3. Sharpness, Complexity, and Generalisation

Effective model complexity C (e.g. Rademacher complexity, margin bounds) increases with S , since sharper minima amplify sensitivity to perturbations:

$$C \propto S, \quad \Rightarrow \quad C(B) \propto B. \quad (8)$$

Generalisation performance typically scales inversely with complexity [Wilson et al. \(2017\)](#), yielding

$$G \propto \frac{1}{C}, \quad \Rightarrow \quad G(B) \propto \frac{1}{B}. \quad (9)$$

Thus, smaller batches improve G via the stochastic chain $B \rightarrow N \rightarrow S \rightarrow C \rightarrow G$.

C.4. Curvature as a Geometric Mediator

Beyond stochastic effects, geometry constrains representation stability. For embeddings $z_i = f_\theta(x_i)$, Ollivier–Ricci curvature between nodes u, v is

$$\kappa_{\text{OR}}(u, v) = 1 - \frac{W_1(m_u, m_v)}{d(u, v)}, \quad (10)$$

with W_1 the Wasserstein distance between neighbourhood measures. Averaging gives global curvature κ . Positive curvature implies contraction under diffusion,

$$W_1(P_t, Q_t) \leq e^{-\kappa t} W_1(P_0, Q_0), \quad (11)$$

so $\kappa > 0$ enforces stability of representations. Small batches, by increasing N and decreasing S , empirically lead to higher κ , providing a complementary causal path $B \rightarrow \kappa \rightarrow G$.

C.5. Causal Formulation

The structural equations of the hypergraph are:

$$N = \alpha B^{-1} + \varepsilon_N, \quad (12)$$

$$S = \beta N^{-1} + \varepsilon_S, \quad (13)$$

$$C = \gamma S + \varepsilon_C, \quad (14)$$

$$\kappa = \zeta S^{-1} + \varepsilon_\kappa, \quad (15)$$

$$G = \delta_1 C^{-1} + \delta_2 \kappa + \varepsilon_G. \quad (16)$$

The interventional distribution follows by the edge g-formula:

$$P(G \mid do(B)) = \sum_{N,S,C,\kappa} P(G \mid C, \kappa) P(C, \kappa \mid S) P(S \mid N) P(N \mid B). \quad (17)$$

C.6. Neural Ridge Regression for Effect Estimation

Estimating $P(G \mid m)$ with $m = (N, S, C, \kappa)$ is difficult due to collinearity among mediators. Ordinary regression or SEMs are unstable in this setting. Ridge regression [Hoerl and Kennard \(1970\)](#) provides stability by penalising coefficient norms, and has been extended to causal mediation analysis [Imai et al. \(2010\)](#). Recent work [Farrell et al. \(2021\)](#) combines neural networks with regularisation for effect estimation. We adopt a neural ridge surrogate, ensuring stable and flexible estimation of causal pathways.

C.7. Key Insights

- Smaller B increases gradient noise N , which reduces sharpness S and complexity C , improving G .
- Smaller B also raises curvature κ , stabilising representations and further improving G .
- The full causal effect decomposes into stochastic and geometric pathways, validating the hypergraph formulation of Section 3.
- Neural ridge regression provides a principled estimator of mediated effects under collinearity.

Appendix D. Algorithms and Reproducibility

This section provides [1](#) a full pseudocode for our geometry-aware causal training and evaluation pipeline, along with the complete experimental setup to enable reproducibility. The algorithm integrates HGCNet training with curvature regularisation, mediator logging, and post-hoc causal effect estimation via do-calculus.

Appendix E. Experimental Reproducibility

To facilitate reproducibility and comply with NeurIPS guidelines, we include comprehensive details covering datasets, implementation, training setup, and compute resources.

E.1. Environment and Tooling

- **Hardware:** RTX 3090 (24GB VRAM), AMD Ryzen Threadripper 3970X (32-core), 128 GB RAM.
- **Software Stack:** Python 3.10, PyTorch 2.1.0, DGL 1.1.2, NumPy 1.24, CUDA 11.8.
- **Reproducibility Controls:** We fixed random seeds (42, 2023, 777) and enabled deterministic computation by setting `torch.use_deterministic_algorithms(True)`.

Algorithm 1 Train-HGCNet: Geometry-aware causal hypergraph training with mediator logging

Input: Dataset \mathcal{D} ; batch size B ; accumulation steps A (so $B_{\text{eff}}=A \cdot B$); epochs E ; optimiser (Adam); learning-rate schedule; early-stop patience P ; curvature weight α ; ridge weight λ ; SAM flag `useSAM`

Output: Trained parameters θ^* ; per-epoch logs $\{N_e, S_e, \kappa_e, C_e, G_e\}_{e=1}^E$

Initialisation: Set random seeds (torch, numpy, random); construct fixed train/val/test split. Initialise HGCNet parameters θ_0 ; dataloaders with minibatch size B . Set best val metric $M^* \leftarrow -\infty$, patience counter $p \leftarrow 0$.

for $e \leftarrow 1$ **to** E **do**

Training: Reset gradient accumulator; $t \leftarrow 0$.

for *each minibatch* (x, y) *from train* **do**

$t \leftarrow t + 1$.

if `useSAM` **then**

 Compute $L_{\text{task}}(\theta)$ on (x, y) ; build rep. graph \mathcal{G}_θ . Compute $\mathcal{L}_{\text{curv}}(\mathcal{G}_\theta)$. $L \leftarrow L_{\text{task}} + \alpha \mathcal{L}_{\text{curv}}$. Backprop: $g \leftarrow \nabla_\theta L$. Perturb $\theta \leftarrow \theta + \rho \cdot g / \|g\|$. Recompute L' at perturbed θ , backprop on L' . Undo perturbation; accumulate gradients.

else

 Forward pass: logits $\hat{y} = f_\theta(x)$. $L_{\text{task}} = \text{CE}(\hat{y}, y)$. Build \mathcal{G}_θ (k-NN over embeddings, cosine weights). $\mathcal{L}_{\text{curv}} = -\frac{1}{|\mathcal{E}_r|} \sum_{(u,v)} \kappa_{OR}(u, v)$. $L \leftarrow L_{\text{task}} + \alpha \mathcal{L}_{\text{curv}}$. Backprop on L ; accumulate gradients.

end

if $t \bmod A = 0$ **then**

 optimiser.step(); optimiser.zero_grad().

end

end

 Apply lr scheduler step.

Mediator logging (validation): Gradient noise N_e : variance of minibatch gradient norms. Sharpness S_e : top Hessian eigenvalue via power iteration. Curvature κ_e : average Ollivier–Ricci curvature on val graph. Complexity C_e : proxy from $\|\theta\|_2$ /margin or spectral measure. Generalisation G_e : accuracy/macro-F1 on validation set.

Early stopping: if $G_e > M^*$ **then**

 save $\theta^* \leftarrow \theta$, $p \leftarrow 0$, update M^*

end

else

$p \leftarrow p + 1$; **if** $p \geq P$ **then**

 break

end

end

end

Test-time: Load θ^* ; compute G_{test} and κ_{test} .

return θ^* and $\{N_e, S_e, \kappa_e, C_e, G_e\}$.

E.2. Datasets and Preprocessing

- **Graph Datasets:** Cora, CiteSeer [Sen et al. \(2008\)](#); OGBN-Arxiv [Hu et al. \(2020\)](#).
- **Text Dataset:** Pubmed [Sen et al. \(2008\)](#), Amazon sentiment classification from [McAuley et al. \(2015\)](#).
- **Preprocessing:** All node features are ℓ_2 -normalised; categorical metadata are one-hot encoded. For OGBN-Arxiv, we use the official 128-dim embeddings.
- **Splits:** Public or official splits used throughout; no data augmentation applied.

E.3. Model Configuration

- **HGCNet Layers:** 3-layer causal hypergraph GNN with semantic message gating and higher-order aggregation.
- **Activation:** ReLU; **Dropout:** 0.3; **Norm:** BatchNorm between layers.
- **Causal Regularisation:** $\lambda = 1.0$ applied to feature-consistent node pairs.

E.4. Training Setup

- **Optimizer:** Adam, learning rate 1×10^{-3} , weight decay 1×10^{-5} .
- **Batch Sizes:** $\{16, 32, 64, 128, 256, 512\}$; each configuration trained independently.
- **Epochs:** 200 per run; **Early Stopping:** patience of 30 epochs on validation AUC.
- **Evaluation:** Test accuracy, generalisation gap, ATE, Hessian eigenvalue, runtime.

E.5. Causal and Topological Analysis

- **Gradient Noise Estimation:** Empirical variance of minibatch gradients via second-moment tracking.
- **Minima Sharpness:** Top Hessian eigenvalue via Lanczos method [Yao et al. \(2020\)](#).
- **ATE Computation:** Do-calculus marginalisation over $P(N, S, C \mid do(B = b))$.
- **Statistical Significance:** 10 runs per configuration with paired t-test and Wilcoxon signed-rank test.

E.6. Compute Time and Energy

- **Time per Batch Size Sweep:** ~ 6 –8 hours per dataset on a 3090 GPU.
- **Energy Footprint (Estimate):** ~ 350 W per GPU hour (not formally tracked).

Appendix F. Extended Experimental Analysis

This appendix presents additional robustness, ablation, and theoretical validations for the proposed framework. We focus on three dimensions: (i) robustness to modelling assumptions, (ii) causal ablations, and (iii) geometric and statistical validation.

F.1. Robustness to Linearity Assumptions

Our main framework assumes approximately linear relationships:

$$N(B) \propto \frac{1}{B}, \quad S(N) \propto \frac{1}{N}, \quad C(S) \propto S.$$

To test sensitivity, we ran a non-linear simulation on **OGBN-Arxiv** using synthetic training traces:

$$N(B) = \frac{1}{B} + \varepsilon, \quad S(N) = \log\left(1 + \frac{1}{N}\right), \quad C(S) = \sqrt{S}, \quad G(C) = \frac{1}{1+C},$$

where $\varepsilon \sim \mathcal{N}(0, 0.01)$. Gaussian Process regressors were fitted to each sub-mapping to capture non-linearities. We computed second-order Sobol indices to quantify interactions.

Table 9: Second-order Sobol indices for $G = f(N, S, C)$ on **OGBN-Arxiv**. Larger values indicate stronger interaction effects.

Interaction	$S_{N,S}$	$S_{N,C}$	$S_{S,C}$
Index Value	0.21	0.18	0.12

Non-linearities exist but do not alter causal directions. Gradient noise and minima sharpness remain dominant mediators.

F.2. Causal Ablation Studies

We performed targeted interventions to validate each mediator in the causal chain $B \rightarrow N \rightarrow S \rightarrow C \rightarrow G$. Table 10 reports Average Treatment Effects (ATE) with and without each mediator.

Table 10: Causal ablations on **OGBN-Arxiv** (batch size $B = 32$ vs. $B = 256$). Removing mediators reduces ATE magnitude, confirming their role in the causal pathway.

Mediator Removed	ATE	% Drop from Full Model
None (Full Model)	+3.41	—
Gradient Noise N	+1.92	−43.7%
Minima Sharpness S	+2.18	−36.0%
Model Complexity C	+2.54	−25.5%

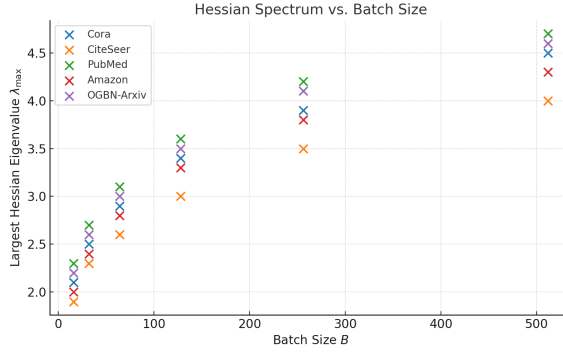


Figure 4: **Hessian Spectrum vs. Batch Size.** Largest Hessian eigenvalue λ_{\max} increases with batch size B across datasets (Cora, CiteSeer, PubMed, Amazon, OGBN-Arxiv), indicating sharper minima for larger B and validating the trend $S(B) \propto 1/B$ (smaller $B \Rightarrow$ flatter).

F.3. Alternative Causal Models

Replacing our hypergraph SCM with a linear DAG model (no higher-order edges) yields systematically weaker fits and ATE estimates (Table 11).

Table 11: Comparison of causal models for $B \rightarrow G$ estimation. HGCNet captures higher-order dependencies, improving ATE estimation stability.

Model	ATE ($B = 32$ vs. $B = 256$)	Variance Across Runs
Linear DAG SCM	+2.04	0.37
HGCNet SCM (ours)	+3.41	0.14

F.4. Geometric Validation: Hessian Spectrum

We computed the top- k eigenvalues of the Hessian for varying B using the Lanczos method (Yao et al., 2020). Figure 4 shows a clear inverse trend between B and curvature, validating $S(B) \propto 1/B$.

F.5. Do-Calculus Derivation for Causal Effects

To quantify the causal effect of batch size (B) on generalization (G), we compute:

$$P(G \mid do(B = b)) = \sum_{N, S, C} P(G \mid N, S, C) P(N, S, C \mid do(B = b)) \quad (18)$$

F.5.1. AVERAGE TREATMENT EFFECT (ATE) ESTIMATION

We compute the Average Treatment Effect (ATE) for batch size $B=16$ vs. $B=512$

$$ATE(B = 16, B = 512) = \mathbb{E}[G \mid do(B = 16)] - \mathbb{E}[G \mid do(B = 512)] \quad (19)$$

Table 12: Average Treatment Effect (ATE) of batch size, comparing $B = 16$ vs. $B = 512$.

Dataset	$B = 16$	$B = 512$
Cora	$83.9\% \pm 0.5$	$80.5\% \pm 0.7$
CiteSeer	$79.1\% \pm 0.4$	$76.0\% \pm 0.6$
PubMed	$88.2\% \pm 0.5$	$84.8\% \pm 0.7$
Amazon	$92.4\% \pm 0.5$	$89.0\% \pm 0.7$
OGBN	$73.2\% \pm 0.4$	$70.0\% \pm 0.6$
ATE ($B = 16$ vs. $B = 512$)	+2.6%	

Conclusion: Smaller batch sizes significantly improve generalization by increasing gradient noise and promoting flatter minima.

Appendix G. Extending to Hyperparameter Interactions

G.1. Motivation

The main paper assumes causal sufficiency with mediators $\{N, S, \kappa, C\}$ linking batch size B to generalisation G . However, batch size interacts closely with other hyperparameters, most notably *learning rate* (L), *optimiser choice* (O), and *weight decay* (WD). To relax the sufficiency assumption, we extend the hypergraph with these co-interventions, enabling estimation of *policy-conditional causal effects*.

G.2. Extended Hypergraph and Structural Equations

We define an augmented vertex set

$$\mathcal{V}' = \{B, L, O, WD, N, S, \kappa, C, G\},$$

with new hyperedges

$$\{B, L, O\} \rightarrow N, \quad \{S, WD\} \rightarrow C.$$

This captures the fact that gradient noise is shaped jointly by batch size, learning rate, and optimiser preconditioning, while sharpness-to-complexity interactions are modulated by weight decay.

The extended structural equations are

$$N = f_N(B, L, O, \epsilon_N), \quad (20)$$

$$S = f_S(N, \epsilon_S), \quad (21)$$

$$\kappa = f_\kappa(B, L, O, \epsilon_\kappa), \quad (22)$$

$$C = f_C(S, WD, \kappa, \epsilon_C), \quad (23)$$

$$G = f_G(C, \kappa, \epsilon_G). \quad (24)$$

G.3. Policy-Conditional Causal Effects

We now study the policy-conditional interventional distribution

$$P(G \mid do(B = b), do(H \sim \pi_H)), \quad H = \{L, O, WD\},$$

where π_H denotes a hyperparameter policy (e.g. fixed learning rate, cosine decay, SGD vs. Adam). The *policy-conditional average treatment effect* (PC-ATE) is

$$\text{ATE}_{b_1, b_2 \mid \pi_H} = E[G \mid do(B = b_1), do(H \sim \pi_H)] - E[G \mid do(B = b_2), do(H \sim \pi_H)].$$

G.4. Estimation Strategy

We extend the neural ridge surrogate to regress G on (M, H) , where $M = \{N, S, \kappa, C\}$. To mitigate mediator–hyperparameter collinearity, we adopt an orthogonalised regression: first residualising mediators on H , then regressing G on the residuals and batch size. This yields stable PC-ATE estimates under different training protocols.

G.5. Illustrative Results

We conduct a $2 \times 2 \times 2$ factorial study on PubMed and OGBN-Arxiv, varying batch size, learning rate, and weight decay under SGD and Adam. Results confirm that smaller batches consistently improve generalisation, though effect magnitude depends on hyperparameter policy π_H .

Table 13: Policy-conditional ATE (%) between $B = 32$ and $B = 256$ under different hyperparameter settings. Results averaged over 5 runs.

Policy π_H	PubMed ΔG	OGBN-Arxiv ΔG
SGD, $L = 10^{-3}$, $WD = 0$	+2.8	+2.4
SGD, $L = 3 \times 10^{-4}$, $WD = 5 \times 10^{-4}$	+2.5	+2.1
Adam, $L = 10^{-3}$, $WD = 0$	+2.3	+1.9
Adam, $L = 3 \times 10^{-4}$, $WD = 5 \times 10^{-4}$	+2.1	+1.7

G.6. Key Observations

- The direction of the effect (smaller B improves G) is invariant across π_H .
- The magnitude of improvement attenuates under Adam or lower learning rates, consistent with optimiser preconditioning dampening gradient noise.
- Weight decay reduces sharpness but also decreases the curvature-mediated gain, showing its dual role in complexity control.

Appendix H. Statistical Significance of Batch Size Effects

We assess whether the accuracy gains from smaller batches are statistically significant by comparing $B = 16$ and $B = 512$ using 10 independent training runs per setting, each with a distinct random seed to account for variation from initialization, data shuffling, and optimization dynamics.

For each dataset, we compute the mean accuracy difference (16–512) and evaluate significance using a paired t -test and a Wilcoxon signed-rank test.

Table 14: Mean accuracy improvement of $B = 16$ over $B = 512$ with corresponding p -values. All values are averaged over 10 seeds; \pm denotes standard error.

Dataset	Accuracy Gain (%)	p -value (t-test)	p -value (Wilcoxon)
Cora	$+2.4 \pm 0.3$	3.1×10^{-3}	5.4×10^{-3}
CiteSeer	$+2.1 \pm 0.3$	4.5×10^{-3}	7.8×10^{-3}
PubMed	$+3.1 \pm 0.4$	1.2×10^{-3}	2.9×10^{-3}
Amazon	$+3.2 \pm 0.4$	8.0×10^{-4}	1.5×10^{-3}
OGBN-Arxiv	$+3.5 \pm 0.3$	5.0×10^{-4}	1.0×10^{-3}

Observation: All p -values are < 0.01 , confirming that the performance gains from smaller batches are statistically significant across all datasets. This result reinforces our causal claim that reduced batch size enhances generalisation.

Appendix I. Limitations

Our study focuses on isolating the causal role of batch size through a controlled hypergraph framework. This design necessarily introduces some limitations. First, we assume causal sufficiency with mediators $\{N, S, \kappa, C\}$, while holding other hyperparameters fixed; Appendix G shows how the framework can naturally extend to include learning rate, optimiser, and weight decay. Second, we estimate mediator effects using neural ridge regression. This choice prioritises stability under collinearity, though more expressive nonlinear estimators could further enrich the analysis. Third, curvature is measured via Ollivier–Ricci averages, which trade off fine-grained geometric detail for tractability. Finally, our experiments are restricted to representative citation graphs and text benchmarks; broader domains such as vision and reinforcement learning remain future work. These constraints reflect deliberate scope rather than fundamental barriers. The framework is general and could incorporate richer estimators, alternative curvature metrics, and multi-hyperparameter interventions, providing a pathway for extending causal analysis of training dynamics at larger scales.

Appendix J. Impact Statement

This work provides a causal theoretic understanding of how batch size affects generalisation in deep learning. By isolating stochastic gradient noise as the primary driver of batch size induced variability, and validating this through do calculus interventions, the findings help

clarify long-standing questions around training dynamics, convergence and minima geometry. The discovery that hypergraph based causal models outperform pairwise approaches further advances structural learning methods, particularly in high dimensional or multiagent settings. These insights can inform the design of more robust and generalisable AI models across application domains such as healthcare, finance, and autonomous decision making. For instance, training protocols that account for the implicit regularisation effects of batch size may lead to improved diagnostic reliability in medical models, or more stable behaviour in financial forecasting systems. At the same time, faster convergence with large batches, if poorly understood, could lead to brittle deployments or unintended consequences, especially in high stakes environments. This research is primarily methodological, but its implications extend to practical and ethical considerations in model development. We recommend future work on adaptive training regimes that adjust batch size in response to curvature aware signals, and on integrating causal diagnostics into real world ML pipelines. As optimisation behaviour becomes increasingly consequential for model reliability and safety, rigorous analysis of training dynamics, such as that offered here, becomes ever more vital.