
Faster Hyperparameter Search on Graphs via Calibrated Dataset Condensation

Anonymous Author(s)
Affiliation
Address
email

Abstract

1 Dataset condensation aims to reduce the computational cost of training multiple
2 models on a large dataset by condensing the training set into a small synthetic
3 one. State-of-the-art approaches rely on matching the gradients between the real
4 and synthetic data and are recently applied to condense large-scale graphs for
5 node classification tasks. Although dataset condensation may be efficient when
6 we need to train multiple models for hyperparameter optimization, there is no
7 theoretical guarantee on the generalizability of the condensed data, and it can *gen-*
8 *eralize poorly* across hyperparameters/architectures in practice; while on graphs,
9 we find and prove this overfitting is much more severe. This paper considers a
10 different condensation objective specifically for hyperparameter search. We aim
11 to generate the synthetic dataset so that the validation-performance ranking of
12 different models under different hyperparameters on the condensed and original
13 datasets are comparable. We propose a novel *hyperparameter-calibrated dataset*
14 *condensation* (HCDC) algorithm, which learns the synthetic validation data by
15 matching the *hyperparameter gradients* computed by implicit differentiation and
16 efficient inverse Hessian approximation. HCDC employs a supernet with dif-
17 ferentiable hyperparameters, making it suitable for modeling GNNs with widely
18 different convolution filters. Experiments demonstrate that the proposed framework
19 effectively maintains the validation-performance rankings of GNNs and speeds up
20 hyperparameter/architecture search on graphs.

21 1 Introduction

22 *Graph neural networks* (GNNs) have found remarkable success in tackling a variety of graph-related
23 tasks [Hamilton, 2020]. However, the prevalence of large-scale graphs in real-world contexts, such as
24 social, information, and biological networks [Hu et al., 2020], which frequently scale up to millions
25 of nodes and edges, poses significant computational issues for training GNNs. While training a single
26 model can be expensive, designing deep learning models for new tasks require substantially more
27 computations, as they involve training multiple models on the same dataset many times to verify the
28 design choices, such as architectures and hyperparameters [Elsken et al., 2019]. We ask: *how can we*
29 *reduce the computational cost for training multiple models on the same dataset, for hyperparameter*
30 *search/optimization?*

31 A natural approach is to reduce the training set size through approaches such as graph coreset
32 selection [Baker et al., 2020], graph sparsification [Batson et al., 2013], graph coarsening [Loukas,
33 2019] and graph sampling [Zeng et al., 2019]. However, these methods are restricted to selecting
34 samples from the given ones, limiting their performance upper-bound. A more effective alternative is
35 to *synthesize* informative samples rather than select from given samples. *Dataset condensation* [Zhao
36 et al., 2020] has emerged as a competent data synthesizing mechanism with promising results. It aims

37 to learn a small synthetic training set such that a model trained on the synthetic set obtains testing
38 accuracy comparable to that trained on the original training set.

39 Although *dataset condensation* achieves the state-of-the-art in terms of the performance for neural
40 networks trained on the condensed samples, it is a *unreliable and sub-optimal* solution to our
41 question, the goal of speeding up training for hyperparameter search/optimization: (1) theoretically,
42 dataset condensation learns synthetic samples that minimize the performance drop of a specific
43 model, and there is no performance guarantee when we train other models; and (2) in practice, we
44 rarely compare condensation methods with strong baselines such as various coreset methods, in
45 terms of their ability to preserve the outcome of architecture/hyperparameter optimization. In this
46 paper, we identify the poor generalizability of the condensed data on graphs [Jin et al., 2021] across
47 architectures/hyperparameters, which has been overlooked when applied for image condensation.
48 Not only we observe that graph condensation fails to preserve validation-performance ranking of
49 GNN architectures, but also we identify and prove two dominant effects causing this failure: (1) most
50 GNNs differ from each other by their design of convolution filters, thus the convolution filter used
51 during condensation is a single biased point in “the space of convolution filters”; and (2) the learned
52 adjacency matrix of the synthetic graph easily overfits the condensation objective, thus fails to
53 maintain the characteristics of the original structure and distinguish different architectures.

54 We aim to *develop a new dataset condensation framework that preserves the outcome of hyperpa-*
55 *rameter search/optimization on the condensed data.* In addition, to condense the training data, we
56 propose to learn the *validation split* of synthetic data such that the validation-performance ranking of
57 architectures on the condensed and original datasets are comparable. Similar to the standard dataset
58 condensation, this new objective can be written as a bi-level optimization problem. Inspired by the
59 gradient-matching algorithm in [Zhao et al., 2020], if assuming a continuous hyperparameter space or
60 a generic supernet which interpolates all architectures, we find and prove the validation-performance-
61 ranking-preserving goal can be realized by matching the *hyperparameter-gradients* on the synthetic
62 and real validation data. The hyperparameter-gradients (or *hypergradients* for short) can be efficiently
63 computed with constant memory overhead by the implicit function theorem (IFT) and the Neumann
64 series approximation of inverse Hessian [Lorraine et al., 2020].

65 The proposed *hyperparameter-calibrated dataset condensation* (HCDC) framework assumes *con-*
66 *tinuous* hyperparameters, which is suitable to model GNNs with different convolution matrices and
67 save the problematic generalizability of graph condensation across GNNs. Although beyond the
68 scope of this paper, HCDC also has the potential to be combined with the supernet in differentiable
69 NAS methods [Liu et al., 2018] to tackle the discrete neural architecture space, which is the primary
70 concern of NAS on image and text data. Experiments demonstrate the effectiveness of the proposed
71 framework in preserving the performance rankings of GNNs. Our distilled graph can be used as proxy
72 data for off-the-shelf graph neural architecture search algorithms to accelerate the search process.

73 Our contributions can be summarized as follows: **(1)** We formulate a new dataset condensation
74 objective specifically for hyperparameter optimization and propose the *hyperparameter-calibrated*
75 *dataset condensation* (HCDC) framework, which learns the synthetic validation data by matching
76 the hyperparameter gradients. **(2)** We prove the hardness of generalizing the condensed graph across
77 GNN architectures. **(3)** Experiments demonstrate the effectiveness of HCDC to speed up architecture
78 search on graphs when combined with off-the-shelf graph neural architecture search algorithms.

79 2 Preliminaries

80 This paper adopts graph learning notations, but HCDC is generally applicable to other data and tasks.
81 The typical downstream task on graphs is node classification. **Node classification on graph** considers
82 that we are given a graph $\mathcal{T} = (A, X, \mathbf{y})$ with adjacency matrix $A \in \{0, 1\}^{n \times n}$, node features
83 $X \in \mathbb{R}^{n \times d}$, node class labels $\mathbf{y} \in [K]^n$, and mutually disjoint node-splits $V_{train} \cup V_{val} \cup V_{test} =$
84 $[n]$. Using a GNN model $f_{\theta, \psi} : \mathbb{R}_{\geq 0}^{n \times n} \times \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times K}$, where θ is parameter and ψ is hyper-
85 parameter, we aim to find $\theta^{\mathcal{T}} = \arg \min_{\theta} \mathcal{L}_{\mathcal{T}}^{train}(\theta, \psi)$, with cross-entropy loss $\mathcal{L}_{\mathcal{T}}^{train}(\theta, \psi) =$
86 $\sum_{i \in V_{train}} \ell([f_{\theta, \psi}(A, X)]_i, y_i)$. The transductive setting can be easily generalized to the inductive
87 setting by assuming only $\{A_{ij} \mid i, j \in V_{train}\}$ and $\{X_i \mid i \in V_{train}\}$ are used during training.

88 Now, we start from reviewing the *standard dataset condensation* (SDC) and its natural *bilevel*
89 *optimization* (BL) formulation [Wang et al., 2018].

90 **Standard Dataset Condensation (SDC)** aims to find a synthetic graph $\mathcal{S} = (A', X', y')$ of size
 91 $c \ll n$, with (weighted) adjacency matrix $A' \in \mathbb{R}_{\geq 0}^{c \times c}$, node features $X' \in \mathbb{R}^{c \times d}$, node labels $y' \in$
 92 $[K]^c$, and (possibly) train/validation¹ splits $V'_{train} \cup V'_{val} = [c]$. The goal of dataset condensation is
 93 to obtain *comparable generalization performance* on real graph by training on the condensed graph,
 94 i.e., $\mathcal{L}_{\mathcal{T}}^{test}(\theta^{\mathcal{T}}, \psi) \approx \mathcal{L}_{\mathcal{T}}^{test}(\theta^{\mathcal{S}}, \psi)$ where $\theta^{\mathcal{S}} = \arg \min_{\theta} \mathcal{L}_{\mathcal{S}}^{train}(\theta, \psi)$ is optimized on the synthetic
 95 graph. By posing $\theta^{\mathcal{S}}$ as a function of the condensed graph \mathcal{S} , dataset condensation can be formulated
 96 as a *bilevel optimization* problem,

$$S^* = \arg \min_{\mathcal{S}} \mathcal{L}_{\mathcal{T}}^{train}(\theta^{\mathcal{S}}(\mathcal{S}), \psi) \quad \text{s.t.} \quad \theta^{\mathcal{S}}(\mathcal{S}) = \arg \min_{\theta} \mathcal{L}_{\mathcal{S}}^{train}(\theta, \psi) \quad (\text{SDC-BL})$$

97 However, the above problem involves a nested-loop optimization and solving the inner loop for
 98 $\theta^{\mathcal{S}}(\mathcal{S})$ at each iteration to recover the gradients for \mathcal{S} requires a computationally expensive procedure:
 99 unrolling the recursive computation graph for \mathcal{S} over multiple optimization steps of θ .

100 [Zhao et al. \[2020\]](#) alleviate the computational issue by its *gradient-matching* (GM) formulation.
 101 To start with, assuming neural network $f_{\theta, \psi}$ is a *locally smooth function*, and thus similar weights
 102 $\theta^{\mathcal{S}} \approx \theta^{\mathcal{T}}$ imply similar mappings in a local neighborhood and thus generalization performance. Then
 103 one can formulate the condensation objective as matching the optimized parameters (which depends
 104 on initialization θ_0), i.e., finding $S^* = \arg \min_{\mathcal{S}} \mathbb{E}_{\theta_0 \sim P_{\theta_0}} [D(\theta^{\mathcal{S}}(\mathcal{S}, \theta_0), \theta^{\mathcal{T}}(\theta_0))]$ s.t. $\theta^{\mathcal{S}}(\mathcal{S}, \theta_0) =$
 105 $\arg \min_{\theta} \mathcal{L}_{\mathcal{S}}^{train}(\theta(\theta_0), \psi)$ where $\theta^{\mathcal{T}}(\theta_0) = \arg \min_{\theta} \mathcal{L}_{\mathcal{T}}^{train}(\theta(\theta_0), \psi)$ and $D(\cdot, \cdot)$ is a distance
 106 function.

107 The *parameter-matching* problem is still a bilevel optimization but can be simplified with several
 108 approximations. Firstly, $\theta^{\mathcal{S}}(\mathcal{S}, \theta_0)$ is approximated by the output of an incomplete gradient-descent
 109 optimization, $\theta^{\mathcal{S}}(\mathcal{S}, \theta_0) \approx \theta_{t+1}^{\mathcal{S}} \leftarrow \theta_t^{\mathcal{S}} - \eta \nabla_{\theta} \mathcal{L}_{\mathcal{S}}^{train}(\theta_t^{\mathcal{S}}, \psi)$. However, the target parameter $\theta^{\mathcal{T}}(\theta_0) =$
 110 $\arg \min_{\theta} \mathcal{L}_{\mathcal{T}}^{train}(\theta(\theta_0), \psi)$ may be far away from $\theta_{t+1}^{\mathcal{S}}$. [Zhao et al. \[2020\]](#) propose to match $\theta_{t+1}^{\mathcal{S}}$ with
 111 incompletely optimized $\theta_{t+1}^{\mathcal{T}} \leftarrow \theta_t^{\mathcal{T}} - \eta \nabla_{\theta} \mathcal{L}_{\mathcal{T}}^{train}(\theta_t^{\mathcal{T}}, \psi)$ at each iteration t , and the condensation
 112 objective is now $S^* = \arg \min_{\mathcal{S}} \mathbb{E}_{\theta_0 \sim P_{\theta_0}} \left[\sum_{t=0}^{T-1} D(\theta_t^{\mathcal{S}}, \theta_t^{\mathcal{T}}) \right]$.

113 Starting from the common initialization θ_0 and up to iteration t , if $\theta_t^{\mathcal{S}}$ can always track $\theta_t^{\mathcal{T}}$
 114 by optimizing \mathcal{S} , i.e., $\theta_t^{\mathcal{S}} \approx \theta_t^{\mathcal{T}}$. For the one step update, we can replace $D(\theta_{t+1}^{\mathcal{S}}, \theta_{t+1}^{\mathcal{T}})$ by
 115 $D(\nabla_{\theta} \mathcal{L}_{\mathcal{S}}^{train}(\theta_t^{\mathcal{S}}, \psi), \nabla_{\theta} \mathcal{L}_{\mathcal{T}}^{train}(\theta_t^{\mathcal{T}}, \psi)) \approx D(\nabla_{\theta} \mathcal{L}_{\mathcal{S}}^{train}(\theta_t^{\mathcal{S}}, \psi), \nabla_{\theta} \mathcal{L}_{\mathcal{T}}^{train}(\theta_t^{\mathcal{S}}, \psi))$. Repeating this
 116 inductive argument, the condensation objective is approximated by matching the gradients at each
 117 iteration t ,

$$S^* = \arg \min_{\mathcal{S}} \mathbb{E}_{\theta_0 \sim P_{\theta_0}} \left[\sum_{t=0}^{T-1} D(\nabla_{\theta} \mathcal{L}_{\mathcal{S}}^{train}(\theta_t^{\mathcal{S}}, \psi), \nabla_{\theta} \mathcal{L}_{\mathcal{T}}^{train}(\theta_t^{\mathcal{S}}, \psi)) \right] \quad (\text{SDC-GM})$$

118 We now have a single deep network with parameters θ trained on the condensed graph \mathcal{S} . While \mathcal{S}
 119 is optimized such that the distance between the gradient vectors of $\mathcal{L}_{\mathcal{T}}^{train}$ and of $\mathcal{L}_{\mathcal{S}}^{train}$ w.r.t. the
 120 parameters θ is minimized. Cosine distance $D(\cdot, \cdot) = \cos(\cdot, \cdot)$ works well in practice.

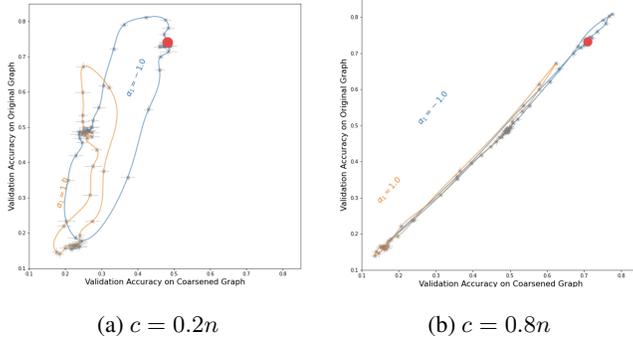
121 3 Standard Dataset Condensation Is Problematic Across GNNs

122 Despite its success in preserving the model performance when trained on the condensed dataset,
 123 the *gradient-matching* algorithm *naturally overfits* the model $f_{\theta, \psi}$ used during condensation and
 124 generalizes poorly to others. There is no guarantee that the condensed synthetic data \mathcal{S}^* which
 125 minimizes the objective (Eq. (SDC-GM)) for a specific model $f_{\theta, \psi}$ (marked by its hyperparameter ψ)
 126 can generalize well to other models $f_{\theta, \psi'}$ where $\psi' \neq \psi$.

127 We find this overfitting issue could be much more *severe on graphs*. For the ease of theoretical
 128 analysis², we consider the simple linear regression problem with *linear convolution* models in this
 129 section, $f_{\theta=[\theta_C, \theta_W]}(A, X) = C(A, \theta_C)XW(\theta_W)$, where $C(A; \theta_C)$ is the *convolution matrix* which
 130 has the same sizes as the adjacency matrix A and possibly also depends on the parameters $\theta_C \in \mathbb{R}^p$,
 131 and $W(\theta_W)$ is the learnable linear weight matrix which is a reshape of the parameters $\theta_W \in \mathbb{R}^d$. The
 132 loss is now sum of squares $\mathcal{L}(\theta, \psi) = \|\mathbf{y} - CXW\|_2^2$ (train-split subscript omitted) where the labels
 133 \mathbf{y} are continuous.

¹The validation split of synthetic data is only required by our HCDC; see Eq. (SDC-BL) vs. Eq. (DCHPO).

²Assuming convex loss and linear models still reflects the general generalization issue; see Appendix A.2.



Ratio (c/n)	A' learned	$A' = I_c$
0.05%	61.3 ± 0.5	59.2 ± 1.1
0.25%	64.2 ± 0.4	63.2 ± 0.3

(a) Graph condensation with identity adjacency $A' = I_c$.

	CVT	GCN	SGC ($K = 2$)	GIN
GCN	—	60.3 ± 0.3	59.2 ± 0.7	42.2 ± 4.3
SGC	—	59.2 ± 1.1	60.5 ± 0.6	39.0 ± 7.1
GIN	—	—	—	59.1 ± 1.1

(b) Generalization across GNNs.

Figure 1: The manifold of GNNs with convolution filters $C = I + \alpha_1(\hat{L}) + \alpha_2(2\hat{L} - I)$ (\hat{L} defined in ChebNet; see Appendix A.3) projected to the plane of validation accuracy on condensed (x-axis) and original (y-axis) graphs under two ratios c/n on Cora. $C = I + \hat{L}$ (red dot) is a biased point in this model space.

Table 1: Empirical results on Ogbn-arxiv verifying the two effects (Propositions 2 and 3) that hinders the generalization of condensed graph across GNNs. (a) Adjacency overfitting, (b) Convolution mismatch.

134 This *linear convolution* model generalizes a wide variety of GNNs [Balcilar et al., 2021, Ding
 135 et al., 2021]; see Appendix A.3. For example, the convolution matrix of graph convolution network
 136 (GCN) [Kipf and Welling, 2016] is $C(A) = \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}$ where \tilde{A} and \tilde{D} are the self-loop-added
 137 adjacency and degree matrix. It also generalizes the one-dimensional *convolution neural network*
 138 (1D-CNN) (with one channel), where the convolution matrix of kernel size $p = (2K + 1)$ is
 139 $C(\theta_C) = \sum_{k=-K}^K [\theta_C]_k P^k$ and P is the cyclic permutation matrix correspond to a unit shift.

140 We say the *gradient-matching objective is satisfied on a non-degenerate trajectory*, if there exists a
 141 fixed learning trajectory $(\theta_t^S)_{t=0}^{T-1}$ which span the whole parameter space, i.e., $\text{span}(\theta_0^S, \dots, \theta_{T-1}^S) =$
 142 \mathbb{R}^{p+d} , such that the *gradient-matching* loss on this trajectory (the objective of Eq. (SDC-GM) without
 143 expectation) is 0. The *validity of standard dataset condensation* (SDC) can be readily verified;
 144 see Proposition 4 in Appendix A.1, where we show if the gradient-matching objective is satisfied on
 145 a non-degenerate trajectory, the optimizer on the condensed dataset \mathcal{S} is also optimal on the original
 146 dataset \mathcal{T} .

147 However, as on the generalizability of condensed dataset across models, we obtain **contrary results**
 148 **for 1D-CNNs and GNNs**.

149 **Proposition 1** (Successful Generalization of SDC across 1D-CNNs). *Assuming least-square re-*
 150 *gression with one-dimensional linear convolution $f_{\theta=[\theta_C, \theta_W]}^{2K+1}(X) = (\sum_{k=-K}^K [\theta_C]_k P^k) XW(\theta_W)$,*
 151 *where kernel size is $(2K + 1), K \geq 0$, if the gradient-matching objective is satisfied on a*
 152 *non-degenerate trajectory for f^{2K+1} , then the condensed dataset \mathcal{S}^* still satisfies the gradient-*
 153 *matching objective on any trajectories $(\theta_t^S)_{t=0}^{T-1}$ for any linear convolution $f^{2K'+1}$ with kernel size*
 154 *$(2K' + 1), K \geq K' \geq 0$.*

155 The intuition behind Proposition 1 is that the 1D-CNN of kernel size $(2K + 1)$ is a “supernet” of
 156 the 1D-CNN of kernel size $(2K' + 1)$ if $K' \leq K$, and the condensed dataset via a bigger model
 157 can generalize well to smaller ones. This result suggests us to *use a sufficiently large model during*
 158 *condensation*, to enable the generalization of the condensed dataset to a wider range of models.

159 However, the story for GNNs is vastly different. We find there are *two dominant effects* causing the
 160 condensed graph to fail to generalize across GNNs. Firstly, the learned adjacency A' of the synthetic
 161 graph \mathcal{S} can easily *overfit* the condensation objective, thus failing to maintain the characteristics of
 162 the original structure and distinguish different architectures; see Table 1 for relevant experiments.

163 **Proposition 2** (Condensed Adjacency Overfits SDC Objective). *Assuming least-square regression*
 164 *with a linear GNN, $f_{\theta}(A, X) = C(A)XW(\theta)$. For any synthetic node features $X^t \in \mathbb{R}^{c \times d}$, there*
 165 *exists a synthetic adjacency matrix $A' \in \mathbb{R}_{\geq 0}^{c \times c}$ such that the gradient-matching objective is satisfied*
 166 *on any trajectories.*

167 Secondly, GNNs differ from each other mostly on the design of convolution $C(A)$, i.e., how the con-
 168 volution weights C depend on the adjacency information A . The convolution filter $C(A)$ used during

condensation is a single biased point in “the space of convolutions”; see Fig. 1 for a visualization, thus there is a *mismatch* of inductive bias when transferring to a different GNN. These two effects lead to the following hardness results when transferring the condensed graph across GNNs.

Proposition 3 (Failed Generalization of SDC across GNNs). *Assuming least square regression and linear GNN, $f_{\theta}^C(A, X) = C(A)XW(\theta)$, there always exists a condensed graph \mathcal{S}^* , such that the gradient-matching objective is satisfied on any trajectories for f^C . However, if we train a new linear GNN $f_{\theta}^{\mathcal{C}}(A, X)$ with convolution matrix $\mathcal{C}(A')$ on \mathcal{S}^* , the relative error between the optimized model parameters of $f^{\mathcal{C}}$ on the real and condensed graphs is $\|\theta_{\mathcal{C}}^{\mathcal{S}} - \theta_{\mathcal{C}}^T\|/\|\theta_{\mathcal{C}}^T\| \geq \max\{\sigma_{\max}(Q) - 1, 1 - \sigma_{\min}(Q)\}$, where $\theta_{\mathcal{C}}^T = \arg \min_{\theta} \|\mathbf{y} - f_{\theta}^{\mathcal{C}}(A, X)\|_2^2$, $\theta_{\mathcal{C}}^{\mathcal{S}} = \arg \min_{\theta} \|\mathbf{y}' - f_{\theta}^{\mathcal{C}}(A', X')\|_2^2$, and $Q = (X^{\top}[C(A)]^{\top}[C(A)]X)(X^{\top}[\mathcal{C}(A)]^{\top}[\mathcal{C}(A)]X)^{-1}$.*

Proposition 3 provides a effective lower-bound on the relative estimation error of optimal model parameter, when a different convolution filter $\mathcal{C}(\cdot) \neq C(\cdot)$ is used³. According to the spectral characterization of convolution filters of GNNs (Table 1 of [Balcilar et al., 2021]), we can approximately compute the maximum eigenvalue of Q for some GNNs. For example, if we condense with f^C graph isomorphism network (GIN-0) [Xu et al., 2018] but train $f^{\mathcal{C}}$ GCN on the condensed graph, we have $\|\theta_{\mathcal{C}}^{\mathcal{S}} - \theta_{\mathcal{C}}^T\|/\|\theta_{\mathcal{C}}^T\| \gtrsim \text{deg} + 1$ where deg is the average node degree of the original graph. This large lower bound hints the catastrophic failure when transfer across GIN and GCN; see Table 1.

Although the results above are obtained for least squares loss and linear convolution model, it *still reflects the nature of general non-convex losses and non-linear models*. Since dataset condensation is effectively matching the local minima $\{\theta^T\}$ of the original loss $\mathcal{L}_{\mathcal{T}}^{\text{train}}(\theta, \psi)$ with the local minima $\{\theta^{\mathcal{S}}\}$ of the condensed loss $\mathcal{L}_{\mathcal{S}}^{\text{train}}(\theta, \psi)$, within the small neighborhoods surrounding the pair of local minima $(\theta^T, \theta^{\mathcal{S}})$, we can approximate the non-convex loss and non-linear model with a convex/linear one respectively. Hence the generalizability issues with convex loss and liner model may hold.

4 Hyperparameter-calibrated Dataset Condensation Objective

Our goal is to develop an optimal and reliable condensation method for architecture/hyperparameter search. Standard Dataset Condensation objective (Eq. (SDC-BL)/Eq. (SDC-GM)) does not accomplish this goal since it does not generalize across GNNs, as proven in Section 3. In this section, we propose a new condensation objective specifically for *preserving the outcome of hyperparameter optimization* (HPO) on the condensed dataset.

HPO Objective HPO finds the optimal hyperparameter ψ^T such that the corresponding model f_{θ, ψ^T} minimizes the validation loss after training, i.e.,

$$\psi^T = \arg \min_{\psi \in \Psi} \mathcal{L}_{\mathcal{T}}^{\text{val}}(\theta^T(\psi), \psi) \quad \text{s.t.} \quad \theta^T(\psi) = \arg \min_{\theta} \mathcal{L}_{\mathcal{T}}^{\text{train}}(\theta, \psi) \quad (\text{HPO})$$

We see HPO itself is a bilevel optimization, where the optimal parameter $\theta^T(\psi)$ is posed as a function of the hyperparameter ψ .

Dataset Condensation for HPO Objective If both the train and validation splits are defined on the condensed dataset \mathcal{S} , the optimal hyperparameter $\psi^{\mathcal{S}}$ is well-defined. Our goal is to find the synthetic dataset \mathcal{S} such that we can obtain *comparable validation performance* if the hyperparameters are optimized on the condensed dataset, i.e., $\mathcal{L}_{\mathcal{T}}^{\text{val}}(\theta^T(\psi^T), \psi^T) \approx \mathcal{L}_{\mathcal{T}}^{\text{val}}(\theta^T(\psi^{\mathcal{S}}), \psi^{\mathcal{S}})$. Clearly, this goal looks very similar to the goal of standard dataset condensation, preserving generalization performance $\mathcal{L}_{\mathcal{T}}^{\text{test}}(\theta^T, \psi) \approx \mathcal{L}_{\mathcal{T}}^{\text{test}}(\theta^{\mathcal{S}}, \psi)$, which hints us to formulate the new objective as a bilevel optimization problem too,

$$\mathcal{S}^* = \arg \min_{\mathcal{S}} \mathcal{L}_{\mathcal{T}}^{\text{val}}(\theta^T(\psi^{\mathcal{S}}(\mathcal{S})), \psi^{\mathcal{S}}(\mathcal{S})) \quad \text{s.t.} \quad \psi^{\mathcal{S}}(\mathcal{S}) = \arg \min_{\psi \in \Psi} \mathcal{L}_{\mathcal{S}}^{\text{val}}(\theta^{\mathcal{S}}(\psi), \psi) \quad (\text{DCHPO})$$

where $\theta^T(\psi)$ and $\theta^{\mathcal{S}}(\psi)$ are defined following Eq. (HPO).

However, this formulation (Eq. (DCHPO)) is a nested optimization (for dataset condensation) over another nested optimization (for HPO) which necessitates very high order gradients and is challenging to solve. Moreover, we have largely overlooked another important factor of hyperparameter optimization, the search space/feasible set of hyperparameters Ψ .

³If $\mathcal{C}(\cdot) = C(\cdot)$ Proposition 4 guarantees $\theta_{\mathcal{C}}^{\mathcal{S}} = \theta_{\mathcal{C}}^T$ and the lower bound in Proposition 3 is 0.

214 In contrast to parameter optimization, where the search space is usually assumed to be the continuous
 215 and unbounded Euclidean space, the search space of hyperparameters Ψ can be either a discrete
 216 set $\Psi = \{\psi_1, \dots\} = \Psi_1 \times \dots \times \Psi_p$ (where each hyperparameter vector $\psi \in \Psi$ consists of p
 217 discrete hyperparameters of various types, for example, neural network type, width, depth, batch
 218 size, etc) or a small bounded set of continuous hyperparameters around its optimum (for example,
 219 learning rate, dropout rate, sample weights, etc). Often we face compositions of these discrete- and
 220 continuous-natured hyperparameters, and we can either model them all as discrete ones and search
 221 by grid search, Bayesian optimization, and reinforcement learning; or relax the discrete search space
 222 to a continuous one.

223 **Hyperparameter-Calibration: A Sufficient Alternative to Dataset Condensation for HPO** The
 224 finiteness/boundedness nature of the search space Ψ cast another challenge to the dataset condensation
 225 for HPO. To avoid the complex combinatorial/constrained optimization in Eq. (HPO), we ask: *Is it*
 226 *possible to preserve the outcome of HPO without solving HPO (Eq. (HPO)) directly?* In this spirit,
 227 we consider a sufficient alternative condition to preserve the outcome of HPO on Ψ .

228 **Definition 1** (Hyperparameter-Calibration). *Given original dataset \mathcal{T} , generic model $f_{\theta, \psi}$, and*
 229 *hyperparameter search space Ψ , we say a condensed dataset \mathcal{S} is hyperparameter-calibrated, if for*
 230 *any $\psi_1 \neq \psi_2 \in \Psi$, it holds that,*

$$(\mathcal{L}_{\mathcal{T}}^{val}(\theta^{\mathcal{T}}(\psi_1), \psi_1) - \mathcal{L}_{\mathcal{T}}^{val}(\theta^{\mathcal{T}}(\psi_2), \psi_2))(\mathcal{L}_{\mathcal{S}}^{val}(\theta^{\mathcal{S}}(\psi_1), \psi_1) - \mathcal{L}_{\mathcal{S}}^{val}(\theta^{\mathcal{S}}(\psi_2), \psi_2)) > 0 \quad (\text{HC})$$

231 *that is, changes of validation loss on \mathcal{T} and \mathcal{S} always have the same sign, where $\theta^{\mathcal{T}}(\psi) =$
 232 $\arg \min_{\theta} \mathcal{L}_{\mathcal{T}}^{train}(\theta, \psi)$ denotes the parameters optimized on the training split of \mathcal{T} with hyperparam-*
 233 *eter ψ , similar for $\theta^{\mathcal{S}}(\psi)$.*

234 It is clear that if hyperparameter calibration (HC) is satisfied, HPO on the original and condensed
 235 datasets yields the same result. Therefore, our mission changes to *how to ensure hyperparameter-*
 236 *calibration for a single pair of hyperparameters (ψ_1, ψ_2) ?*

237 **HCDC: Hypergradient-alignment Objective** To proceed, we make an important extra assumption
 238 that the (possibly discrete) search space Ψ can be extended to a compact and connected set $\Psi' \supset \Psi$,
 239 where we can define continuation of the generic model $f_{\theta, \psi}$ on Ψ' so that $f_{\theta, \psi}$ is differentiable
 240 anywhere in Ψ' . Such a continual extension naturally exists on graphs (see Section 5) or can be
 241 provided by differentiable NAS approaches; see Section 6.

242 Now, if we consider the special case where ψ_1 is within the neighborhood of ψ_2 , i.e., $\psi_1 \in B_r(\psi_2)$
 243 for some $r > 0$, and reparameterize $\psi_1 = \psi + \Delta\psi$, $\psi_2 = \psi$ with $r \geq \|\Delta\psi\|_2 \rightarrow 0^+$. The
 244 change in validation loss is approximated *up to first-order* by the hyperparameter-gradients (hy-
 245 pergradients for short) $(\mathcal{L}_{\mathcal{T}}^{val}(\theta^{\mathcal{T}}(\psi_1), \psi_1) - \mathcal{L}_{\mathcal{T}}^{val}(\theta^{\mathcal{T}}(\psi_2), \psi_2)) \approx \nabla_{\psi} \mathcal{L}_{\mathcal{T}}^{val}(\theta^{\mathcal{T}}(\psi), \psi) \cdot \Delta\psi$. The
 246 hyperparameter-calibration condition within this tiny neighborhood $B_r(\psi)$ is then simplified to
 247 $\nabla_{\psi} \mathcal{L}_{\mathcal{T}}^{val}(\theta^{\mathcal{T}}(\psi), \psi) / \nabla_{\psi} \mathcal{L}_{\mathcal{S}}^{val}(\theta^{\mathcal{S}}(\psi), \psi)$, i.e., the two hypergradient vectors are aligned (i.e., point-
 248 ing same direction).

249 Considering the extended search space Ψ' can be covered by the union of many small neighborhoods,
 250 we derive the hypergradient-alignment condition: $\nabla_{\psi} \mathcal{L}_{\mathcal{T}}^{val}(\theta^{\mathcal{T}}(\psi), \psi) / \nabla_{\psi} \mathcal{L}_{\mathcal{S}}^{val}(\theta^{\mathcal{S}}(\psi), \psi)$ for any
 251 $\psi \in \Psi'$. It is not hard to show that the condition above is equivalent to hyperparameter-calibration
 252 (Definition 1) on a connected and compact set Ψ' . (1) *Necessity* proved by contradiction. If
 253 there exists $\psi_0 \in \Psi'$ such that the two gradient vectors are not aligned at ψ_0 , then there exists
 254 small perturbation $\Delta\psi_0$ such that $(\mathcal{L}_{\mathcal{T}}^{val}(\theta^{\mathcal{T}}(\psi_0 + \Delta\psi_0), \psi_0 + \Delta\psi_0) - \mathcal{L}_{\mathcal{T}}^{val}(\theta^{\mathcal{T}}(\psi_0), \psi_0))$ and
 255 $(\mathcal{L}_{\mathcal{S}}^{val}(\theta^{\mathcal{S}}(\psi_0 + \Delta\psi_0), \psi_0 + \Delta\psi_0) - \mathcal{L}_{\mathcal{S}}^{val}(\theta^{\mathcal{S}}(\psi_0), \psi_0))$ have different signs. (2) *Sufficiency* proved
 256 by integration. For any pair $\psi_1 \neq \psi_2 \in \Psi'$, if we have a continuous path connecting ψ_1 and ψ_2 , then
 257 integrating hypergradients $\nabla_{\psi} \mathcal{L}_{\mathcal{T}}^{val}(\theta^{\mathcal{T}}(\psi), \psi)$ and $\nabla_{\psi} \mathcal{L}_{\mathcal{S}}^{val}(\theta^{\mathcal{S}}(\psi), \psi)$ through the path recovers the
 258 hyperparameter-calibration condition.

259 In this regard, enforcing hypergradient-alignment on Ψ' is sufficient to hyperparameter calibration on
 260 Ψ , thus ensuring the outcome of HPO over Ψ is preserved. The hypergradient-alignment objective
 261 below realizes hyperparameter-calibrated dataset condensation (HCDC).

$$\mathcal{S}^* = \arg \min_{\mathcal{S}} \sum_{\psi \in \Psi'} D(\nabla_{\psi} \mathcal{L}_{\mathcal{T}}^{val}(\theta^{\mathcal{T}}(\psi), \psi), \nabla_{\psi} \mathcal{L}_{\mathcal{S}}^{val}(\theta^{\mathcal{S}}(\psi), \psi)) \quad (\text{HCDC})$$

262 where cosine distance $D(\cdot, \cdot) = \cos(\cdot, \cdot)$ is used.

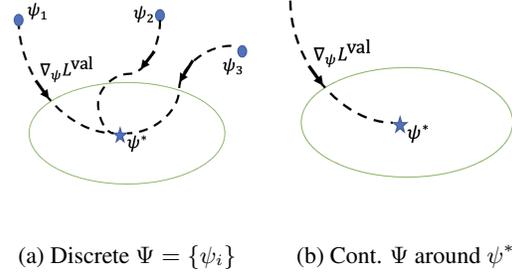
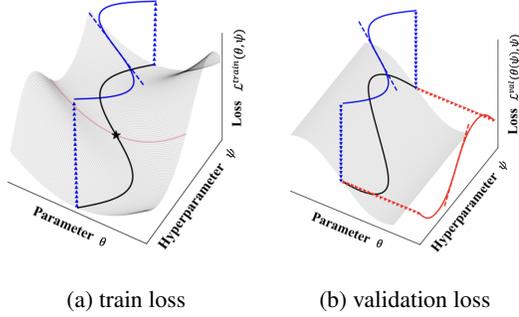


Figure 2: The parameter-hyperparameter manifolds and IFT. The blue solid line is the best response $\theta^*(\psi)$. The red dashed line is hypergradients $\nabla_{\psi} \mathcal{L}^{val}(\theta(\psi), \psi)$.

Figure 3: Where to align the hypergradients in HCDC (Eq. (HCDC)); see Section 5 for explanations.

263 5 Implementation of HCDC

264 Finally, we work on implementing and simplifying the hyperparameter-calibrated dataset condensa-
 265 tion (HCDC) objective and apply it to the graph architecture/hyperparameter search problem.

266 **How is HCDC connected to standard dataset condensation (SDC) (Eq. (SDC-GM))?** Theoretically
 267 speaking, the objective of HCDC, preserving the outcome of hyperparameter optimization (HPO), is
 268 orthogonal to the objective of SDC, preserving generalization performance. Therefore, we can limit
 269 the part of the synthetic dataset they optimize to make the two algorithms completely independent.
 270 While SDC only learns the training split of \mathcal{S} , we restrict HCDC to only optimize the validation split
 271 of \mathcal{S} in Eq. (HCDC) and keep the training split fixed⁴. Nevertheless, we need to find the condensed
 272 training data \mathcal{S}^{train} before HCDC, and this can be done by all kinds of approaches, from uniform
 273 sampling to SDC.

274 **How to compute hypergradients and optimize the hypergradient-alignment loss in Eq. (HCDC)?**
 275 The efficient computation of hypergradients $\nabla_{\psi} \mathcal{L}_{\mathcal{T}}^{val}(\theta^{\mathcal{T}}(\psi), \psi)$ and $\nabla_{\psi} \mathcal{L}_{\mathcal{S}}^{val}(\theta^{\mathcal{S}}(\psi), \psi)$ uses
 276 the implicit function theorem (IFT) (see Section 4 for visualization), $\nabla_{\psi} \mathcal{L}_{\mathcal{T}}^{val}(\theta^{\mathcal{T}}(\psi), \psi) =$
 277 $-\left[\frac{\partial^2 \mathcal{L}_{\mathcal{T}}^{train}(\theta, \psi)}{\partial \psi \partial \theta^T}\right] \left[\frac{\partial^2 \mathcal{L}_{\mathcal{T}}^{train}(\theta, \psi)}{\partial \theta \partial \theta^T}\right]^{-1} \nabla_{\theta} \mathcal{L}_{\mathcal{T}}^{val}(\theta, \psi) + \nabla_{\psi} \mathcal{L}_{\mathcal{T}}^{val}(\theta, \psi)$, where $\nabla_{\psi} \mathcal{L}_{\mathcal{T}}^{val}(\theta, \psi)$ is the di-
 278 rect gradient and often identically 0. The first term is the product of inverse training Hesse-
 279 sian $\left[\frac{\partial^2 \mathcal{L}_{\mathcal{T}}^{train}(\theta, \psi)}{\partial \theta \partial \theta^T}\right]^{-1}$, the training mixed partials $\left[\frac{\partial^2 \mathcal{L}_{\mathcal{T}}^{train}(\theta, \psi)}{\partial \psi \partial \theta^T}\right]$ and the validation gradients
 280 $\nabla_{\theta} \mathcal{L}_{\mathcal{T}}^{val}(\theta, \psi)$. While the other parts can be computed by back-propagation, the inverse Hessian
 281 needs to be approximated. Instead of using conjugate gradient method, Lorraine et al. [2020]
 282 propose a stable, tractable and efficient Neumann series approximation, $\left[\frac{\partial^2 \mathcal{L}_{\mathcal{T}}^{train}(\theta, \psi)}{\partial \theta \partial \theta^T}\right]^{-1} =$
 283 $\lim_{i \rightarrow \infty} \sum_{j=0}^i \left[I - \frac{\partial^2 \mathcal{L}_{\mathcal{T}}^{train}(\theta, \psi)}{\partial \theta \partial \theta^T} \right]^j$ with constant memory constraint. To optimize the validation
 284 part of \mathcal{S} w.r.t. the cosine hypergradient-matching loss in Eq. (HCDC), note that we only need to
 285 take gradients of $\nabla_{\theta} \mathcal{L}_{\mathcal{S}}^{val}(\theta, \psi)$ and $\nabla_{\psi} \mathcal{L}_{\mathcal{S}}^{val}(\theta, \psi)$ w.r.t. \mathcal{S}^{val} . This can be handled by the same
 286 back-propagation technique in SDC, where we take gradients of $\nabla_{\theta} \mathcal{L}_{\mathcal{S}}^{train}(\theta, \psi)$ w.r.t \mathcal{S}^{train} .

287 **Where to align the hypergradients in Eq. (HCDC)?** The hypergradient-alignment condition, as a
 288 sufficient condition for preserving the outcome of HPO, is often too strong. For a discrete search
 289 space Ψ , we can preserve the order of any $\psi_1 \neq \psi_2 \in \Psi$, as long as there exists a continuous
 290 path connecting ψ_1 and ψ_2 on which the hypergradients $\nabla_{\psi} \mathcal{L}_{\mathcal{T}}^{val}(\theta^{\mathcal{T}}(\psi), \psi)$ and $\nabla_{\psi} \mathcal{L}_{\mathcal{S}}^{val}(\theta^{\mathcal{S}}(\psi), \psi)$
 291 are aligned. To further avoid the $O(p^2)$ paths, we propose to align the hyperparameters on the p
 292 continuous-HPO trajectories. The i -th continuous-HPO trajectory starts from $\psi_{i,0}^S = \psi_i \in \Psi$ and
 293 update through $\psi_{i,t+1}^S \leftarrow \psi_{i,t}^S - \eta \nabla_{\psi} \mathcal{L}_{\mathcal{S}}^{val}(\theta^{\mathcal{S}}(\psi_{i,t}^S), \psi_{i,t}^S)$. All of the p trajectories will approach the
 294 optima ψ^S which form a “connected” path between any pair of hyperparameters $\psi_i \neq \psi_j \in \Psi$. For a
 295 continuous search space $\Psi = \Psi'$, since it is often bounded narrowly around the optima ψ^S , we again
 296 align the hypergradients along the optimization trajectories $(\psi_{i,t}^S)_{t=0}^{T-1}$ despite that the starting points
 297 $\psi_i \in \Psi$ is now randomly sampled; see Section 4.

⁴It is also possible for graph condensation when the train and validation subgraphs are not connected.

298 **What graph architecture/hyperparameter search problem can HCDC solve?** We illustrate how to
 299 tackle the two types of search spaces: (1) discrete and finite Ψ and (2) continuous and bounded Ψ with
 300 two typical examples originating from the problem of searching for the best convolution matrix $C(A)$
 301 on a large graph $\mathcal{T} = (A, X, y)$. **(1) Discrete and finite search space Ψ :** often the most important
 302 question of architecture search on large graphs is *what design of convolution filter performs best on the*
 303 *given graph?* One may simply train the set of p prior-defined GNNs $\{f_{[\theta_{C_i}, \theta_w]}^{C_i} \mid i = 1, \dots, p\}$ whose
 304 convolution matrices are $\mathcal{C} = \{C_1(A; \theta_{C_1}), \dots, C_p(A; \theta_{C_p})\}$ and compare their validation perfor-
 305 mance. We can formulate this problem as HPO, by defining an “interpolated” model $f_{[\theta_C, \theta_w], \psi}^C$ whose
 306 convolution matrix is $C(A; \theta_C, \psi) = \varphi_1 C_1(A; \theta_{C_1}) + \dots + \varphi_p C_p(A; \theta_{C_p})$, where hyperparameters
 307 $\psi = [\varphi_1, \dots, \varphi_p] \in \Psi$ and $\theta_C = [\theta_{C_1}, \dots, \theta_{C_p}]$. The feasible set $\Psi = \{\psi_1 = \mathbf{e}_1^p, \dots, \psi_p = \mathbf{e}_p^p\}$ is
 308 the set of unit vectors in \mathbb{R}^p and the extended search space can be defined as $\Psi' = [0, 1]^p$. **(2) Contin-**
 309 **uous and bounded search space Ψ :** one may also use a continuous generic formula, e.g., truncated
 310 series, to model a wide-range of convolution filters, i.e., $C(A; \psi) = \sum_{i=1}^p \varphi_i C_i(A)$, for example in
 311 ChebNet [Defferrard et al., 2016] or SGC [Wu et al., 2019] (see Appendix A.3). The only difference
 312 to the previous case is that the search space $\Psi = \Psi'$ can be larger than $[0, 1]^p$.

313 6 Related Work

314 **Graph condensation** [Jin et al., 2021] achieves the state-of-the-art performance for preserving GNNs’
 315 performance on the simplified graph. However, Jin et al. [2021] only adapt the gradient-matching
 316 algorithm of dataset distillation Zhao et al. [2020] to graph data, together with a MLP-based generative
 317 model for edges [Anand and Huang, 2018, Simonovsky and Komodakis, 2018], leaving out several
 318 major issues on efficiency, performance, and generalizability (discussed in Section 1). Subsequent
 319 work aims to apply the more efficient distribution-matching algorithm of dataset distillation [Zhao
 320 and Bilen, 2021a, Wang et al., 2022] to graph or speed up gradient-matching graph condensation by
 321 reducing the number of gradient-matching-steps [Jin et al., 2022]. While the efficiency issue of graph
 322 condensation is mitigated [Jin et al., 2022], the performance degradation on medium- and large-sized
 323 graphs still renders graph condensation practically meaningless. Our hyperparameter-calibrated graph
 324 distillation is specifically designed for repeated training in architecture search, which is, in contrast,
 325 well-motivated.

326 **Implicit differentiation** methods apply the implicit function theorem (IFT) to the nested-optimization
 327 problems [Ochs et al., 2015, Wang et al., 2019]. The IFT requires inverting the training Hessian
 328 with respect to the network weights. Lorraine et al. [2020] approximates the inverse Hessian by
 329 the Neumann series, which is a stable alternative to conjugate gradients [Shaban et al., 2019] and
 330 successfully scales gradient-based bilevel-optimization to large networks with constant memory
 331 constraint. It is shown that unrolling differentiation around locally optimal weights for i steps is
 332 equivalent to approximating the Neumann series inverse approximation up to the first i terms.

333 **Differentiable NAS** methods, e.g., DARTS [Liu et al., 2018] explore the possibility of transforming
 334 the discrete neural architecture space into a continuously differentiable form and further uses gradient
 335 optimization to search the neural architecture. DARTS follows a cell-based search space [Zoph et al.,
 336 2018] and continuously relaxes the original discrete search strategy. Differentiable NAS techniques
 337 have also been applied to graphs to automatically design data-specific GNN architectures [Wang
 338 et al., 2021, Huan et al., 2021].

339 In addition, we summarize graph reduction methods (including graph coreset selection, graph
 340 sampling, graph sparsification, and graph coarsening), as well as the more dataset condensation and
 341 coreset selection methods beyond graphs in Appendix B.

342 7 Experiments

343 In this section we validate the effectiveness of hyperparameter-calibrated dataset condensation
 344 (HCDC) when applied to speed up graph architecture/hyperparameter search. In this section, correla-
 345 tion refers to the Spearman’s rank correlation coefficient r_s between two rankings of the ordered list of
 346 hyperparameters on the original and condensed datasets. Please refer to ?? for more implementation
 347 details.

Hyperpara. Method	Cross Validation	
	Correlation	Performance
Random	-0.04	84.0
DC	0.68	82.6
DM	0.76	82.8
Early-Stopping	0.11	84.3
HCDC	0.91	84.7

Table 2: The rank correlation and validation performance on the real dataset of the M -fold cross validation ranked/selected on the condensed dataset.

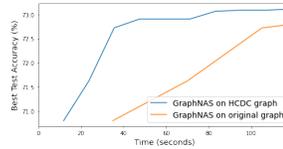


Figure 4: Speed-up to the search process of graph NAS when combined with HCDC on Ogbn-arxiv, best test performance so far vs. time spent.

348 **Synthetic experiments on CIFAR-10.** We first consider a synthetically created set of hyperpara-
349 rameters on image dataset, CIFAR-10. Consider the M -fold cross validation, where a fraction of
350 $1/M$ samples are use as the validation split each time. The M -fold cross-validation process can be
351 modeled by a set of M hyperparameters $\{\varphi_i \in \{0, 1\} \mid i = 1, \dots, M\}$, where $\varphi_i = 1$ if and only if
352 the i -th fold is used for validation. The problem of finding the best validation performance among
353 the M results can be modeled as a hyperparameter optimization problem with a discrete search
354 space $|\Psi| = M$. We compare HCDC with the gradient-matching [Zhao et al., 2020] and distribution
355 matching [Zhao and Bilen, 2021a] baselines. We also consider a uniform random sampling baseline,
356 and an early-stopping baseline where we train only $c/n * 500$ epochs but on the original dataset. The
357 results of $M = 20$ and $c/n = 1\%$ is reported in Table 2, where we see HCDC achieves the highest
358 rank correlation.

359 **Finding best convolution filter on (large) graphs.** One easy application of HCDC we analyzed
360 in Section 5 is to speed up the selection of the best suited convolution filter design on (large) graphs.
361 Following the method discussed in Section 5, we test HCDC against (1) Random: random uniform
362 sampling of nodes and find their induced subgraph, (2) GCond-X: graph condensation [Jin et al.,
363 2021] but fix the synthetic adjacency to identity, (3) GCond: the graph condensation algorithm in [Jin
364 et al., 2021], and (4) Whole Graph: when the model selection is performed on the original dataset.
365 We use random uniform sampling to find the training synthetic subgraph before we apply HCDC. For
366 the other coreset/condensation methods which does not define the validation split, we random split
367 the train and validation nodes according to the original split ratio. We not only report the Spearman’s
368 rank correlation, but also the test performance (on real dataset) of model selected by the condensed
369 dataset. The results are summarized in Table 3.

370 **Speeding up off-the-shelf graph architecture search algorithms.** Finally we test HCDC on how
371 much speed-up it can provides to the off-the-shelf graph architecture search methods. We use graph
372 NAS [Gao et al., 2019] on Ogbn-arxiv with a condensation ratio of $c/n = 0.5\%$. The search space of
373 architectures is the same as the set used in Table 3 with a focus on graphs with different convolution
374 filters. We plot the best test performance of searched architecture (so far) versus the time spent for
375 searching (in seconds) in Fig. 4. We see HCDC, as a dataset condensation approach, can further
376 speed up the search process of graph NAS and is orthogonal to the efficient search algorithms like
377 Bayesian optimization or reinforcement learning used by NAS methods.

Dataset	Ratio	Random		GCond-X		GCond		HCDC		Whole Graph Perf. (%)
		Corr.	Perf.	Corr.	Perf.	Corr.	Perf.	Corr.	Perf.	
Cora	0.9%	0.29	81.2	0.14	79.5	0.61	81.9	0.80	83.2	83.8
	1.8%	0.40	81.9	0.21	80.3	0.76	83.2	0.89	83.8	
	3.6%	0.51	82.2	0.22	80.9	0.81	83.2	0.92	83.8	
Citeseer	1.3%	0.38	71.9	0.15	70.8	0.68	71.3	0.90	73.1	73.7
	2.6%	0.56	72.2	0.29	70.8	0.79	71.5	0.93	73.7	
	5.2%	0.71	73.0	0.35	70.2	0.83	71.1	0.97	73.7	
Ogbn-arxiv	0.25%	0.59	70.1	0.39	69.8	0.59	70.3	0.77	71.9	73.4
	0.5%	0.63	70.3	0.44	70.1	0.64	70.5	0.85	72.2	
	1.0%	0.68	70.9	0.47	70.0	0.67	70.1	0.88	72.2	
Reddit	0.1%	0.42	92.1	0.39	90.9	0.53	90.9	0.88	92.1	94.3
	0.2%	0.50	93.1	0.41	90.9	0.61	91.2	0.92	92.7	
	0.4%	0.58	93.1	0.42	91.5	0.66	92.1	0.96	92.7	

Table 3: Spearman’s rank correlation and test performance of the convolution filter selected on the condensed graph.

378 **References**

- 379 William L Hamilton. *Graph Representation Learning*. Morgan & Claypool Publishers, 2020.
- 380 Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta,
381 and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *Advances in*
382 *neural information processing systems*, 33:22118–22133, 2020.
- 383 Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. Neural architecture search: A survey. *The*
384 *Journal of Machine Learning Research*, 20(1):1997–2017, 2019.
- 385 Daniel Baker, Vladimir Braverman, Lingxiao Huang, Shaofeng H-C Jiang, Robert Krauthgamer, and
386 Xuan Wu. Coresets for clustering in graphs of bounded treewidth. In *International Conference on*
387 *Machine Learning*, pages 569–579. PMLR, 2020.
- 388 Joshua Batson, Daniel A Spielman, Nikhil Srivastava, and Shang-Hua Teng. Spectral sparsification
389 of graphs: theory and algorithms. *Communications of the ACM*, 56(8):87–94, 2013.
- 390 Andreas Loukas. Graph reduction with spectral and cut guarantees. *J. Mach. Learn. Res.*, 20(116):
391 1–42, 2019.
- 392 Hanqing Zeng, Hongkuan Zhou, Ajitesh Srivastava, Rajgopal Kannan, and Viktor Prasanna. Graph-
393 saint: Graph sampling based inductive learning method. In *International Conference on Learning*
394 *Representations*, 2019.
- 395 Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. Dataset condensation with gradient matching. In
396 *International Conference on Learning Representations*, 2020.
- 397 Wei Jin, Lingxiao Zhao, Shichang Zhang, Yozen Liu, Jiliang Tang, and Neil Shah. Graph condensation
398 for graph neural networks. In *International Conference on Learning Representations*, 2021.
- 399 Jonathan Lorraine, Paul Vicol, and David Duvenaud. Optimizing millions of hyperparameters by
400 implicit differentiation. In *International Conference on Artificial Intelligence and Statistics*, pages
401 1540–1552. PMLR, 2020.
- 402 Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. In
403 *International Conference on Learning Representations*, 2018.
- 404 Tongzhou Wang, Jun-Yan Zhu, Antonio Torralba, and Alexei A Efros. Dataset distillation. *arXiv*
405 *preprint arXiv:1811.10959*, 2018.
- 406 Muhammet Balcilar, Renton Guillaume, Pierre H eroux, Benoit Ga uz ere, S ebastien Adam, and Paul
407 Honeine. Analyzing the expressive power of graph neural networks in a spectral perspective. In
408 *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- 409 Mucong Ding, Kezhi Kong, Jingling Li, Chen Zhu, John Dickerson, Furong Huang, and Tom
410 Goldstein. Vq-gnn: A universal framework to scale up graph neural networks using vector
411 quantization. *Advances in Neural Information Processing Systems*, 34:6733–6746, 2021.
- 412 Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks.
413 In *International Conference on Learning Representations*, 2016.
- 414 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural
415 networks? In *International Conference on Learning Representations*, 2018.
- 416 Micha el Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on
417 graphs with fast localized spectral filtering. In *Advances in neural information processing systems*,
418 volume 29, 2016.
- 419 Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Sim-
420 plifying graph convolutional networks. In *International conference on machine learning*, pages
421 6861–6871. PMLR, 2019.
- 422 Namrata Anand and Possu Huang. Generative modeling for protein structures. *Advances in neural*
423 *information processing systems*, 31, 2018.

- 424 Martin Simonovsky and Nikos Komodakis. Graphvae: Towards generation of small graphs using
425 variational autoencoders. In *International conference on artificial neural networks*, pages 412–422.
426 Springer, 2018.
- 427 Bo Zhao and Hakan Bilen. Dataset condensation with distribution matching. *arXiv preprint*
428 *arXiv:2110.04181*, 2021a.
- 429 Kai Wang, Bo Zhao, Xiangyu Peng, Zheng Zhu, Shuo Yang, Shuo Wang, Guan Huang, Hakan
430 Bilen, Xinchao Wang, and Yang You. Cafe: Learning to condense dataset by aligning features. In
431 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
432 12196–12205, 2022.
- 433 Wei Jin, Xianfeng Tang, Haoming Jiang, Zheng Li, Danqing Zhang, Jiliang Tang, and Bing Yin.
434 Condensing graphs via one-step gradient matching. In *Proceedings of the 28th ACM SIGKDD*
435 *Conference on Knowledge Discovery and Data Mining*, pages 720–730, 2022.
- 436 Peter Ochs, René Ranftl, Thomas Brox, and Thomas Pock. Bilevel optimization with nonsmooth
437 lower level problems. In *International Conference on Scale Space and Variational Methods in*
438 *Computer Vision*, pages 654–665. Springer, 2015.
- 439 Yuanhao Wang, Guodong Zhang, and Jimmy Ba. On solving minimax optimization locally: A
440 follow-the-ridge approach. In *International Conference on Learning Representations*, 2019.
- 441 Amirreza Shaban, Ching-An Cheng, Nathan Hatch, and Byron Boots. Truncated back-propagation
442 for bilevel optimization. In *The 22nd International Conference on Artificial Intelligence and*
443 *Statistics*, pages 1723–1732. PMLR, 2019.
- 444 Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures
445 for scalable image recognition. In *Proceedings of the IEEE conference on computer vision and*
446 *pattern recognition*, pages 8697–8710, 2018.
- 447 Zhili Wang, Shimin Di, and Lei Chen. Autogel: An automated graph neural network with explicit
448 link information. *Advances in Neural Information Processing Systems*, 34:24509–24522, 2021.
- 449 ZHAO Huan, YAO Quanming, and TU Weiwei. Search to aggregate neighborhood for graph neural
450 network. In *2021 IEEE 37th International Conference on Data Engineering (ICDE)*, pages
451 552–563. IEEE, 2021.
- 452 Yang Gao, Hong Yang, Peng Zhang, Chuan Zhou, and Yue Hu. Graphnas: Graph neural architecture
453 search with reinforcement learning. *arXiv preprint arXiv:1904.09981*, 2019.
- 454 Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs.
455 *Advances in neural information processing systems*, 30, 2017.
- 456 Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua
457 Bengio. Graph attention networks. In *International Conference on Learning Representations*,
458 2018.
- 459 Johannes Klicpera, Stefan Weissenberger, and Stephan Günnemann. Diffusion improves graph
460 learning. In *Advances in neural information processing systems*. PMLR, 2019.
- 461 Yu Rong, Yatao Bian, Tingyang Xu, Weiyang Xie, Ying Wei, Wenbing Huang, and Junzhou Huang.
462 Self-supervised graph transformer on large-scale molecular data. In *Advances in neural information*
463 *processing systems*, volume 33, 2020.
- 464 Omri Puny, Heli Ben-Hamu, and Yaron Lipman. From graph low-rank global attention to 2-fwl
465 approximation. In *International Conference on Machine Learning*. PMLR, 2020.
- 466 Jiawei Zhang, Haopeng Zhang, Congying Xia, and Li Sun. Graph-bert: Only attention is needed for
467 learning graph representations. *arXiv preprint arXiv:2001.05140*, 2020.
- 468 Ondrej Bohdal, Yongxin Yang, and Timothy Hospedales. Flexible dataset distillation: Learn labels
469 instead of images. *arXiv preprint arXiv:2006.08572*, 2020.

- 470 Timothy Nguyen, Zhou rong Chen, and Jaehoon Lee. Dataset meta-learning from kernel ridge-
471 regression. In *International Conference on Learning Representations*, 2020.
- 472 Timothy Nguyen, Roman Novak, Lechao Xiao, and Jaehoon Lee. Dataset distillation with infinitely
473 wide convolutional networks. *Advances in Neural Information Processing Systems*, 34:5186–5198,
474 2021.
- 475 Felipe Petroski Such, Aditya Rawal, Joel Lehman, Kenneth Stanley, and Jeffrey Clune. Generative
476 teaching networks: Accelerating neural architecture search by learning to generate synthetic
477 training data. In *International Conference on Machine Learning*, pages 9206–9216. PMLR, 2020.
- 478 Bo Zhao and Hakan Bilen. Dataset condensation with differentiable siamese augmentation. In
479 *International Conference on Machine Learning*, pages 12674–12685. PMLR, 2021b.
- 480 Jang-Hyun Kim, Jinuk Kim, Seong Joon Oh, Sangdoon Yun, Hwanjun Song, Joonhyun Jeong, Jung-
481 Woo Ha, and Hyun Oh Song. Dataset condensation via efficient synthetic-data parameterization.
482 In *International Conference on Machine Learning*, pages 11102–11118. PMLR, 2022.
- 483 Yutian Chen, Max Welling, and Alex Smola. Super-samples from kernel herding. In *Proceedings of*
484 *the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence*, pages 109–116, 2010.
- 485 Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl:
486 Incremental classifier and representation learning. In *Proceedings of the IEEE conference on*
487 *Computer Vision and Pattern Recognition*, pages 2001–2010, 2017.
- 488 Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for
489 online continual learning. *Advances in neural information processing systems*, 32, 2019.
- 490 Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set
491 approach. In *International Conference on Learning Representations*, 2018.
- 492 Zalán Borsos, Mojmír Mutny, and Andreas Krause. Coresets via bilevel optimization for continual
493 learning and streaming. *Advances in Neural Information Processing Systems*, 33:14879–14890,
494 2020.
- 495 Mariya Toneva, Alessandro Sordani, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and
496 Geoffrey J Gordon. An empirical study of example forgetting during deep neural network learning.
497 In *International Conference on Learning Representations*, 2018.
- 498 Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet: Finding
499 important examples early in training. *Advances in Neural Information Processing Systems*, 34:
500 20596–20607, 2021.
- 501 Suraj Kothawade, Vishal Kaushal, Ganesh Ramakrishnan, Jeff Bilmes, and Rishabh Iyer. Prism: A
502 rich class of parameterized submodular information measures for guided data subset selection. In
503 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10238–10246,
504 2022.
- 505 Rishabh Iyer, Ninad Khargoankar, Jeff Bilmes, and Himanshu Asanani. Submodular combinatorial
506 information measures with applications in machine learning. In *Algorithmic Learning Theory*,
507 pages 722–754. PMLR, 2021.
- 508 Chengcheng Guo, Bo Zhao, and Yanbing Bai. Deepcore: A comprehensive library for coreset
509 selection in deep learning. *arXiv preprint arXiv:2204.08499*, 2022.
- 510 Vladimir Braverman, Shaofeng H-C Jiang, Robert Krauthgamer, and Xuan Wu. Coresets for clustering
511 in excluded-minor graphs and beyond. In *Proceedings of the 2021 ACM-SIAM Symposium on*
512 *Discrete Algorithms (SODA)*, pages 2679–2696. SIAM, 2021.
- 513 Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. Cluster-gcn: An
514 efficient algorithm for training deep and large graph convolutional networks. In *Proceedings of*
515 *the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pages
516 257–266, 2019.

- 517 Venu Satuluri, Srinivasan Parthasarathy, and Yiye Ruan. Local graph sparsification for scalable
518 clustering. In *Proceedings of the 2011 ACM SIGMOD International Conference on Management*
519 *of data*, pages 721–732, 2011.
- 520 Andreas Loukas and Pierre Vandergheynst. Spectrally approximating large graphs with smaller
521 graphs. In *International Conference on Machine Learning*, pages 3237–3246. PMLR, 2018.
- 522 Zengfeng Huang, Shengzhong Zhang, Chong Xi, Tang Liu, and Min Zhou. Scaling up graph
523 neural networks via graph coarsening. In *Proceedings of the 27th ACM SIGKDD Conference on*
524 *Knowledge Discovery & Data Mining*, pages 675–684, 2021.
- 525 Chen Cai, Ding kang Wang, and Yusu Wang. Graph coarsening with neural networks. In *International*
526 *Conference on Learning Representations*, 2020.
- 527 Yoshua Bengio. Gradient-based optimization of hyperparameters. *Neural computation*, 12(8):
528 1889–1900, 2000.
- 529 Jan Larsen, Lars Kai Hansen, Claus Svarer, and M Ohlsson. Design and regularization of neural
530 networks: the optimal use of a validation set. In *Neural Networks for Signal Processing VI.*
531 *Proceedings of the 1996 IEEE Signal Processing Society Workshop*, pages 62–71. IEEE, 1996.
- 532 Jelena Luketina, Mathias Berglund, Klaus Greff, and Tapani Raiko. Scalable gradient-based tuning
533 of continuous regularization hyperparameters. In *International conference on machine learning*,
534 pages 2952–2960. PMLR, 2016.
- 535 Fabian Pedregosa. Hyperparameter optimization with approximate gradient. In *International*
536 *conference on machine learning*, pages 737–746. PMLR, 2016.
- 537 James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate
538 curvature. In *International conference on machine learning*, pages 2408–2417. PMLR, 2015.
- 539 Liam Li and Ameet Talwalkar. Random search and reproducibility for neural architecture search. In
540 *Uncertainty in artificial intelligence*, pages 367–377. PMLR, 2020.
- 541 Arber Zela, Thomas Elsken, Tonmoy Saikia, Yassine Marrakchi, Thomas Brox, and Frank Hutter.
542 Understanding and robustifying differentiable architecture search. In *International Conference on*
543 *Learning Representations*, 2019.
- 544 Sirui Xie, Hehui Zheng, Chunxiao Liu, and Liang Lin. Snas: stochastic neural architecture search. In
545 *International Conference on Learning Representations*, 2018.
- 546 Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. In
547 *International Conference on Learning Representations*, 2017.
- 548 Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous
549 relaxation of discrete random variables. In *International Conference on Learning Representations*,
550 2017.

551 Supplementary Material

552 A Proofs and Extended Theoretical Results

553 Recall we consider the linear regression problem with linear convolution model $f_{\theta=[\theta_C, \theta_W]}(A, X) =$
554 $C(A; \theta_C)XW(\theta_W)$, where $C(A; \theta_C)$ is the *convolution matrix* which has the same sizes as the
555 adjacency matrix A and possibly also depends on the parameters $\theta_C \in \mathbb{R}^p$, and $W(\theta_W)$ is the
556 learnable linear weight matrix which is a reshape of the parameters $\theta_W \in \mathbb{R}^d$. The loss is now sum of
557 squares $\mathcal{L}(\theta, \psi) = \|\mathbf{y} - CXW\|_2^2$ (train-split subscript omitted) where the labels \mathbf{y} are continuous.

558 Recall that we say the gradient-matching objective is satisfied for a non-degenerate trajectory,
559 if there exists a fixed learning trajectory $(\theta_t^S)_{t=0}^{T-1}$ which span the whole parameter space, i.e.,
560 $\text{span}(\theta_0^S, \dots, \theta_{T-1}^S) = \mathbb{R}^{p+d}$, such that the *gradient-matching* loss on this trajectory (objective of
561 Eq. (SDC-GM) without expectation) is 0.

562 A.1 Validity of Standard Dataset Condensation

563 **Proposition 4.** (*Validity of SDC*) Assuming least square regression with linear convolution model,
564 if gradient-matching objective is satisfied for a non-degenerate trajectory, then the optimizer on
565 the condensed dataset \mathcal{S} , i.e., $\theta^S = \arg \min_{\theta} \mathcal{L}_{\mathcal{S}}(\theta, \psi)$ is also optimal for the original dataset, i.e.,
566 $\mathcal{L}_{\mathcal{T}}(\theta^S, \psi) = \min_{\theta} \mathcal{L}_{\mathcal{T}}(\theta, \psi)$.

567 *Proof.* If the gradient-matching objective is satisfied for a non-degenerate trajectory, one can derive
568 $(X^{\top} C^{\top} C X) = (X'^{\top} C'^{\top} C' X')$ and $X^{\top} C^{\top} \mathbf{y} = X'^{\top} C'^{\top} \mathbf{y}'$. These directly lead to that the
569 optimizer $\theta_W^S = \theta_W^T$ and thus $\mathcal{L}_{\mathcal{T}}(\theta^S, \psi) = \min_{\theta} \mathcal{L}_{\mathcal{T}}(\theta, \psi)$. \square

570 A.2 Generalization Issue of Standard Dataset Condensation

571 *Proof of Proposition 1:* If the gradient-matching objective is satisfied for a non-degenerate trajectory,
572 similarly, one can derive $(X^{\top} P^{k\top} P^k X) = (X'^{\top} P^{k\top} P^k X')$ and $X^{\top} P^{k\top} \mathbf{y} = X'^{\top} P^{k\top} \mathbf{y}'$ for any
573 $k = -K, \dots, 0, \dots, K$. Thus for a 1D-CNN with smaller $K' \leq K$, the equations above readily lead
574 to the gradient-matching objective for the new model. \square

575 *Proof of Proposition 2:* This proposition easily follow from the fact that the gradient-
576 matching objective is minimized for any trajectory if the optimized parameter matches $\theta^S =$
577 $(X'^{\top} C'^{\top} C' X') X'^{\top} C'^{\top} \mathbf{y}' = (X^{\top} C^{\top} C X) X^{\top} C^{\top} \mathbf{y}$. From which we see we only require the
578 product $C' X'$ to match $C X$. \square

579 *Proof of Proposition 3:* The lower bound follows from the optimized parameter θ_W^S and θ_W^T with the
580 new GNN with convolution filter \mathcal{C} and the inequality $\|AB\|_F \leq \sigma_{max}(A)\|B\|_F$. \square

581 Although the results above are obtained for least squares loss and linear convolution model, it still
582 reflects the nature of general non-convex losses and non-linear models. Since dataset condensation is
583 effectively matching the local minima $\{\theta^T\}$ of the original loss $\mathcal{L}_{\mathcal{T}}^{train}(\theta, \psi)$ with the local minima
584 $\{\theta^S\}$ of the condensed loss $\mathcal{L}_{\mathcal{S}}^{train}(\theta, \psi)$, within the small neighborhoods surrounding the pair of local
585 minima (θ^T, θ^S) , we can approximate the non-convex loss and non-linear model with a convex/linear
586 one respectively. And hence the generalizability issues with convex loss and linear model may hold.

587 A.3 Convolution Filters of GNNs

588 The convolution formulation of many popular GNNs [Balcilar et al., 2021] is summarized in Table 4.

589 B Extended Related Work

590 B.1 Dataset Condensation and Coreset Selection

591 Firstly, we review the two main approaches to reducing the training set size while preserving model
592 performance.

Model Name	Design Idea	Conv. Matrix Type	# of Conv.	Convolution Matrix
GCN ¹ [Kipf and Welling, 2016]	Spatial Conv.	Fixed	1	$C = \tilde{D}^{-1/2} \tilde{A} \tilde{D}^{-1/2}$
SAGE-Mean ² [Hamilton et al., 2017]	Message Passing	Fixed	2	$\begin{cases} C^{(1)} = I_n \\ C^{(2)} = D^{-1}A \end{cases}$
GAT ³ [Veličković et al., 2018]	Self-Attention	Learnable	# of heads	$\begin{cases} \mathbf{e}^{(s)} = A + I_n \text{ and} \\ h_{\mathbf{a}^{(l,s)}}^{(s)}(X_{i,:}^{(l)}, X_{j,:}^{(l)}) = \exp(\text{LeakyReLU}(\mathbf{a}^{(l,s)} \cdot (X_{i,:}^{(l)} W^{(l,s)} \parallel X_{j,:}^{(l)} W^{(l,s)}))) \end{cases}$
GIN ¹ [Xu et al., 2018]	WL-Test	Fixed + Learnable	2	$\begin{cases} C^{(1)} = A \\ C^{(2)} = I_n \text{ and } h_{\mathbf{e}^{(l)}}^{(2)} = 1 + \epsilon^{(l)} \end{cases}$
SGC! ^{1,2} [Defferrard et al., 2016]	Spectral Conv.	Learnable	order of poly.	$\begin{cases} \mathbf{e}^{(1)} = I_n, \mathbf{e}^{(2)} = 2L/\lambda_{\max} - I_n, \\ \mathbf{e}^{(s)} = 2\mathbf{e}^{(2)}\mathbf{e}^{(s-1)} - \mathbf{e}^{(s-2)} \\ \text{and } h_{\theta^{(s)}}^{(s)} = \theta^{(s)} \end{cases}$
ChebNet ² [Defferrard et al., 2016]	Spectral Conv.	Learnable	order of poly.	$\begin{cases} \mathbf{e}^{(1)} = I_n, \mathbf{e}^{(2)} = 2L/\lambda_{\max} - I_n, \\ \mathbf{e}^{(s)} = 2\mathbf{e}^{(2)}\mathbf{e}^{(s-1)} - \mathbf{e}^{(s-2)} \\ \text{and } h_{\theta^{(s)}}^{(s)} = \theta^{(s)} \end{cases}$
GDC ³ [Klicpera et al., 2019]	Diffusion	Fixed	1	$C = S$
Graph Transformers ⁴ [Rong et al., 2020]	Self-Attention	Learnable	# of heads	$\begin{cases} \mathbf{e}_{i,j}^{(s)} = 1 \text{ and } h_{(W_Q^{(l,s)}, W_K^{(l,s)})}^{(s)}(X_{i,:}^{(l)}, X_{j,:}^{(l)}) \\ = \exp(\frac{1}{\sqrt{d_{k,l}}} (X_{i,:}^{(l)} W_Q^{(l,s)})(X_{j,:}^{(l)} W_K^{(l,s)})^\top) \end{cases}$

¹ Where $\tilde{A} = A + I_n$, $\tilde{D} = D + I_n$. ² $C^{(2)}$ represents mean aggregator. Weight matrix in [Hamilton et al., 2017] is $W^{(l)} = W^{(l,1)} \parallel W^{(l,2)}$. ³ Need row-wise normalization. $C_{i,j}^{(l,s)}$ is non-zero if and only if $A_{i,j} = 1$, thus GAT follows direct-neighbor aggregation. ⁴ The weight matrices of the two convolution supports are the same, $W^{(l,1)} = W^{(l,2)}$. ⁵ Where normalized Laplacian $L = I_n - D^{-1/2}AD^{-1/2}$ and λ_{\max} is its largest eigenvalue, which can be approximated as 2 for a large graph. ⁶ Where S is the diffusion matrix $S = \sum_{k=0}^{\infty} \theta_k \mathbf{T}^k$, for example, decaying weights $\theta_k = e^{-\frac{k}{\tau}}$ and transition matrix $\mathbf{T} = \tilde{D}^{-1/2} \tilde{A} \tilde{D}^{-1/2}$. ⁷ Need row-wise normalization. Only describes the global self-attention layer, where $W_Q^{(l,s)}, W_K^{(l,s)} \in \mathbb{R}^{f_l, d_{k,l}}$ are weight matrices which compute the queries and keys vectors. In contrast to GAT, all entries of $\mathbf{e}_{i,j}^{(l,s)}$ are non-zero. Different design of Graph Transformers [Puny et al., 2020, Rong et al., 2020, Zhang et al., 2020] use graph adjacency information in different ways, and is not characterized here, see the original papers for details.

Table 4: Summary of GNNs formulated as generalized graph convolution.

593 **Dataset condensation** (or distillation) is first proposed in [Wang et al., 2018] as a learning-to-learn
594 problem by formulating the network parameters as a function of synthetic data and learning them
595 through the network parameters to minimize the training loss over the original data. However, the
596 nested-loop optimization precludes it scaling up to large-scale in-the-wild datasets. Zhao et al. [2020]
597 alleviate this issue by enforcing the gradients of the synthetic samples w.r.t. the network weights
598 to approach those of the original data, which successfully alleviates the expensive unrolling of
599 the computational graph. Based on the meta-learning formulation in [Wang et al., 2018], Bohdal
600 et al. [2020] and Nguyen et al. [2020, 2021] propose to simplify the inner-loop optimization of a
601 classification model by training with ridge regression which has a closed-form solution, while Such
602 et al. [2020] model the synthetic data using a generative network. To improve the data efficiency
603 of synthetic samples in gradient-matching algorithm, Zhao and Bilen [2021b] apply differentiable
604 Siamese augmentation, and Kim et al. [2022] introduce efficient synthetic-data parametrization.
605 Recently, a new distribution-matching framework [Zhao and Bilen, 2021a] proposes to match the
606 hidden features rather than the gradients for fast optimization, but may suffer from performance
607 degradation compared to gradient-matching [Zhao and Bilen, 2021a], where Kim et al. [2022] provide
608 some interpretation.

609 **Coreset selection** methods choose samples that are important for training based on heuristic criteria,
610 for example, minimizing the distance between coreset and whole-dataset centers [Chen et al., 2010,
611 Rebuffi et al., 2017], maximizing the diversity of selected samples in the gradient space [Aljundi
612 et al., 2019], discovering cluster centers [Sener and Savarese, 2018], and choosing samples with the
613 largest negative implicit gradient [Borsos et al., 2020]. *Forgetting* [Toneva et al., 2018] measures
614 the forgetfulness of trained samples and drops those that are not easy to forget. *GraNd* [Paul et al.,
615 2021] selects the training samples that contribute most to the training loss in the first few epochs.
616 *Prism* [Kothawade et al., 2022] select samples to maximize submodular set-functions which are
617 combinatorial generalizations of entropy measures [Iyer et al., 2021]. Recent benchmark [Guo et al.,
618 2022] of a variety of coreset selection methods for image classification indicates *Forgetting*, *GraNd*,
619 and *Prism* are among the best performing coreset methods but still evidently underperform the dataset

620 condensation baselines. Although coreset selection can be very efficient, most of the methods above
621 suffer from three major limitations: (1) their performance is upper-bounded by the information in
622 the selected samples; (2) most of them do not directly optimize the synthetic samples to preserve
623 the model performance; and (3) most of them select samples incrementally and greedily, which are
624 short-sighted.

625 B.2 Graph Reduction

626 Secondly, we summarize the traditional graph reduction method for graph neural network training.

627 **Graph coreset selection** is a non-trivial generalization of the above method coreset methods given
628 the non-*iid* nature of graph nodes and the non-linearity nature of GNNs. The very few off-the-shelf
629 graph coreset algorithms are designed for graph clustering [Baker et al., 2020, Braverman et al.,
630 2021] and are not optimal for the training of GNNs.

631 **Graph sampling** methods [Chiang et al., 2019, Zeng et al., 2019] can be as simple as uniformly
632 sampling a set of nodes and finding their induced subgraph, which is understood as a graph-counterpart
633 of uniform sampling of *iid* samples. However, most of the present graph sampling algorithms (e.g.,
634 ClusterGCN [Chiang et al., 2019] and GraphSAINT [Zeng et al., 2019]) are designed for sampling
635 multiple subgraphs (mini-batches), which forms a cover of the original graph for training GNNs
636 with memory constraint. Therefore those graph mini-batch sampling algorithms are effectively graph
637 partitioning algorithms and not optimized to find just one representative subgraph.

638 **Graph sparsification** [Batson et al., 2013, Satuluri et al., 2011] and **graph coarsening** [Loukas
639 and Vandergheynst, 2018, Loukas, 2019, Huang et al., 2021, Cai et al., 2020] algorithms are usually
640 designed to preserve specific graph properties like graph spectrum and graph clustering. Such
641 objectives are often not aligned with the optimization of downstream GNNs and are shown to be
642 sub-optimal in preserving the information to train GNNs well [Jin et al., 2021].

643 B.3 Other Related Areas

644 Lastly, we list several important relevant areas.

645 **Implicit differentiation** methods apply the implicit function theorem (IFT) to the nested-optimization
646 problems [Ochs et al., 2015, Wang et al., 2019]. The IFT requires inverting the training Hessian with
647 respect to the network weights, where early work either computes the inverse explicitly [Bengio,
648 2000, Larsen et al., 1996] or approximates it as the identity matrix [Luketina et al., 2016]. Conjugate
649 gradient (CG) is applied to invert the Hessian approximately [Pedregosa, 2016], but is difficult to
650 scale to deep networks. Several methods have been proposed to efficiently approximate Hessian
651 inverse, for example, 1-step unrolled differentiation [Luketina et al., 2016], Fisher information
652 matrix [Larsen et al., 1996], NN-structure aided Kronecker-factored inversion [Martens and Grosse,
653 2015]. Lorraine et al. [2020] use the Neumann inverse approximation, which is a stable alternative to
654 CG [Shaban et al., 2019] and successfully scale gradient-based bilevel-optimization to large networks
655 with constant memory constraint. It is shown that unrolling differentiation around locally optimal
656 weights for i steps is equivalent to approximating the Neumann series inverse approximation up to
657 the first i terms.

658 **Differentiable NAS** methods, e.g., DARTS [Liu et al., 2018] explore the possibility of transforming
659 the discrete neural architecture space into a continuously differentiable form and further uses gradient
660 optimization to search the neural architecture. DARTS follows a cell-based search space [Zoph
661 et al., 2018] and continuously relaxes the original discrete search strategy. Despite its simplicity,
662 several work cast doubt on the effectiveness of DARTS [Li and Talwalkar, 2020, Zela et al., 2019].
663 SNAS [Xie et al., 2018] points out that DARTS suffers from the unbounded bias issue towards
664 its objective, and it remodels the NAS and leverages the Gumbel-softmax trick [Jang et al., 2017,
665 Maddison et al., 2017] to learn the exact architecture parameter. Differentiable NAS techniques have
666 also been applied to graphs to automatically design data-specific GNN architectures [Wang et al.,
667 2021, Huan et al., 2021].