# On the VC dimension of deep group convolutional neural networks

Anna Sepliarskaia[1,2], Sophie Langer[1], and Johannes Schmidt-Hieber[1]

[1]Department of Applied Mathematics, University of Twente

[2]Booking.com

October 22, 2024

### Abstract

We study the generalization capabilities of Group Convolutional Neural Networks (GCNNs) with ReLU activation function by deriving upper and lower bounds for their Vapnik-Chervonenkis (VC) dimension. Specifically, we analyze how factors such as the number of layers, weights, and input dimension affect the VC dimension. We further compare the derived bounds to those known for other types of neural networks. Our findings extend previous results on the VC dimension of continuous GCNNs with two layers, thereby providing new insights into the generalization properties of GCNNs, particularly regarding the dependence on the input resolution of the data.

## 1 Introduction

Convolutional Neural Networks (CNNs) have revolutionized the field of computer vision, achieving remarkable success in tasks such as image classification (Krizhevsky et al., 2012), object detection (Ren et al., 2016), and segmentation (Long et al., 2015). Their effectiveness can be partly attributed to their translation invariant architecture, enabling CNNs to recognize objects regardless of their position in an image. However, while CNNs are effective at capturing translation symmetries, there has been a growing interest in incorporating additional structure into neural networks to handle a wider range of transformations. These architectures aim to combine the flexibility of learning with the robustness of structure-preserving features (see, e.g., (Hinton and Wang, 2011; Lee et al., 2015)).

GCNNs were first introduced by Cohen and Welling (2016a) to improve statistical efficiency and enhance geometric reasoning. Since then equivariant network structures have evolved to support equivariance on Euclidean groups (Bekkers et al., 2018; Bekkers, 2019; Weiler et al., 2018), compact groups (Kondor and Trivedi, 2018) and Riemannian manifolds (Weiler et al., 2021). More recent architectures have even been generalized beyond other types of symmetry groups (Zhdanov et al., 2024; Dehmamy et al., 2021; Smets et al., 2023)

These advancements have significantly broadened the applicability of CNNs to more complex tasks including fluid dynamics (Wang et al., 2020; Vinuesa and Brunton, 2022), electrodynamics (Zhdanov et al., 2024), medical image segmentation (Bekkers et al., 2018), partial differential equation (PDE) solvers (Brandstetter et al., 2022), and video tracking (Sosnovik et al., 2021). A notable example of the success of equivariant neural networks is the AlphaFold algorithm, which achieves high accuracy in protein structure prediction through a novel roto-translation equivariant attention architecture of the neural network (Jumper et al., 2021; Fuchs et al., 2020).

Equivariant neural networks are expected to have smaller sample complexity if compared to neural networks without built-in symmetries (Bietti et al., 2021; Sannai et al., 2021; Elesedy, 2022; Shao et al., 2022). Given the symmetries are present in the data, equivariant networks should therefore achieve comparable performance with fewer examples.

To investigate this hypothesis, we study the generalization capabilities of a GCNN class that is invariant to group actions. We compute upper and lower bounds for the Vapnik-Chervonenkis (VC) dimension of this class and compare them with the VC dimension of deep feedforward neural networks (DNNs).

Previously, the VC dimension has been used to analyze the generalization properties of various network architectures, including vanilla DNNs (Bartlett et al., 2019; Anthony and Bartlett, 1999) and CNNs (Kohler and Walter, 2023). However, to the best of our knowledge, the only result on the VC dimension of GCNNs pertains to two-layer GCNNs with continuous groups, which were shown to have an infinite VC dimension (Petersen and Sepliarskaia, 2024). Our work differs from earlier research by analyzing the VC dimension of GCNNs in relation to the number of layers, weights, and input dimensions of the neural network, without restricting the architecture to just two layers (as in (Petersen and Sepliarskaia, 2024)) or DNNs (as in (Bartlett et al., 2019)).

## 2 GCNNs

### 2.1 Group theoretic concepts

GCNNs capture symmetries in data via group theory. A *group* $G$ is a set equipped with an operation $\circ$ such that $h, g \in G$ implies $h \circ g \in G$; there exists an identity element $e \in G$ such that $e \circ g = g \circ e = g$ for all $g \in G$; there exists an inverse element $g^{-1} \in G$ such that $g^{-1} \circ g = g \circ g^{-1} = e$; and for any $g, h, i \in G$, we have $(g \circ h) \circ i = g \circ (h \circ i)$.

A **group action** describes how a group interacts with another set. More specifically, an *action* of $G$ on a set $\mathcal{X}$ is a map $\circ : G \times \mathcal{X} \to \mathcal{X}$ such that for all $g_1, g_2 \in G$ and for all $x \in \mathcal{X}$

$$(g_1 g_2) \circ x = g_1 \circ (g_2 \circ x).$$

On the functions $\{f : G \to \mathbb{R}\}$, the *left regular representation* of a group $G$ is the action

$$(g \circ f)(g') = f(g^{-1}g'), \quad \text{for any } g, g' \in G. \tag{1}$$

For a so-called kernel function $\mathcal{K} : G \to \mathbb{R}$ and a function $f : G \to \mathbb{R}$, the group convolution of $\mathcal{K}$ with $f$ at an element $g \in G$ is defined as

$$(\mathcal{K} \star f)(g) := \int_G \mathcal{K}(g^{-1}g') \cdot f(g') \, d\mu(g'), \tag{2}$$

with $\mu$ being the Haar measure (Procesi, 2007, Theorem 8.1.2). Note that, the Haar measure is a unique left-invariant measure that exists for all locally compact groups (see (Stroppel, 2006, Definition 1.18)). For the integral in (2) to be well-defined, we assume that both, the signal $f$ and the kernel $\mathcal{K}$, are measurable and bounded.

The group convolution computes a weighted average over the group elements. It is equivariant with respect to the left regular representation of $G$ defined in (1), that is, $g \circ (\mathcal{K} \star f) = (\mathcal{K} \star (g \circ f))$ for any group element $g$.

In the case where $G$ is a finite group, the group convolution becomes the sum

$$(\mathcal{K} \star f)(g) = \frac{1}{|G|} \sum_{g' \in G} \mathcal{K}(g^{-1}g') \cdot f(g'). \tag{3}$$

In practice, continuous groups such as rotations or translations are often discretized, and computations are performed on finite grids. This discretization process involves approximating the group $G$ by a finite subset

$$G^r = \{g_1, g_2, \ldots, g_r\}.$$

We refer to the cardinality $r$ as the resolution of the discretization.

For example, for the group formed by all continuous rotations, the discretization selects a finite number of rotation angles. For the translation group, discretization consists of a finite number of shifts.

After discretization, the group convolution operates on this finite set. Ignoring the reweighting it is then given by

$$(\mathcal{K} * f)(g) := \sum_{j=1}^{r} \mathcal{K}(g^{-1}g_j) \cdot f(g_j). \tag{4}$$

This operation is called *G-correlation* (Cohen and Welling, 2016b). Note that $(\mathcal{K} * g) = r \cdot (\mathcal{K} \star g)$. If $\mathcal{K} = \mathbf{1}(\cdot = e)$ with $e$ the identity element of the group, then, on $G^r$, $K * f = f$. The G-correlation heavily depends on the discretization of the group and can differ significantly from the integral version (2). However, both definitions become approximately the same (up to rescaling) if the elements $g_j$ are drawn from the Haar measure. In the case of GCNNs, the discretization is determined by the structure of the data.

## 2.2 GCNN architecture

To parametrize the kernel function, let

$$\mathrm{K}_s : G \to \mathbb{R}, \quad s = 1, \ldots, k \tag{5}$$

be a set of basis functions. The kernel functions $\mathcal{K}_{\mathbf{w}}$ are then expressed as linear combinations of these basis functions, that is,

$$\mathcal{K}_{\mathbf{w}} = \sum_{s=1}^{k} w_s \mathrm{K}_s, \tag{6}$$

where $\mathbf{w} = (w_1, \ldots, w_k)$ is the vector of trainable parameters.

Given an activation function $\sigma : \mathbb{R} \to \mathbb{R}$, a group convolutional unit (GCNN unit) or G-convolution takes an input function on the discretized group $f = (f_1, \ldots, f_{m_0})^T : G^r \to \mathbb{R}^{m_0}$ as input and outputs another function on the discretized group $h : G^r \to \mathbb{R}$. The output $h$ is

$$h = \sigma\left( \sum_{i=1}^{m_0} \mathcal{K}_{\mathbf{w}_i} * f_i - b \right), \tag{7}$$

with the group convolution operation $*$ defined in (4). The weight vectors $\mathbf{w}_1, \ldots, \mathbf{w}_m$ and the bias $b$ are parameters that are learned from the data. In line with the common terminology, we refer to the output as feature map.

A GCNN layer (also called $G$-convolutional layer or $G$-conv layer) computes several GCNN units in parallel, using the same input functions but applying different kernel functions (also known as filters), each with potentially different parameters. Specifically, a GCNN layer with $M$ units and input function $f = (f_1, \ldots, f_{m_0})^T : G^r \to \mathbb{R}^{m_0}$, computes the following $M$ functions

$$h_j = \sigma\left( \sum_{i=1}^{m_0} \mathcal{K}_{\mathbf{w}_{ij}} * f_i - b_j \right), \quad j = 1, \ldots, M, \tag{8}$$

where $\mathcal{K}_{\mathbf{w}_{ij}}$ connects the $i$-th input with the $j$-th output. The trainable parameters in the GCNN layer are the weight vectors $\mathbf{w}_{ij}$ and the biases $b_j$. As other network architectures, GCNNs are typically structured hierarchically, with several GCNN layers followed by fully connected layers. We assume that the input of the first GCNN layer are already functions on the discretized group.

Let $L$ be the number of GCNN layers and assume that the respective numbers of GCNN units in the layers $1, \ldots, L$ is denoted by $m_1, \ldots, m_L$. To derive a recursive formula for the GCNN, we denote the inputs of the GCNN by $h_{0,1}, \ldots, h_{0,m_0}$. Note that they are also assumed to be functions $G^r \to \mathbb{R}$. If for a given $\ell = 0, \ldots, L-1$, $h_{\ell,1}, \ldots, h_{\ell,m_\ell}$ are the outputs (also known as feature maps) of the $\ell$-th GCNN layer, then, the $m_{\ell+1}$ outputs of the $(\ell+1)$-st GCNN layer are given by

$$h_{\ell+1,j} = \sigma\left( \sum_{i=1}^{m_\ell} \mathcal{K}_{\mathbf{w}_{ij}^{(\ell)}} * h_{\ell,i} - b_j^{(\ell)} \right) \tag{9}$$

for $= 1, \ldots, m_{\ell+1}$. The final output of the GCNN is given by

$$\sum_{i=1}^{m_L} \sum_{g \in G^r} h_{L,i}(g), \qquad (10)$$

which equals up to reweighting a global average pooling operation. As we take the sum over all group elements $g$, this operation makes the network invariant instead of equivariant to geometric transformation (see, e.g., (Kondor and Trivedi, 2018; Keriven and Peyré, 2019; Bekkers et al., 2018)).

In our work, we consider the ReLU activation function $\sigma(x) = \max\{x, 0\}$. We denote the class of ReLU GCNNs with $L$ layers, $m_i$ units in layer $i = 0, \ldots, L$, $k$ dimensional weight vectors in (6), and $r$ the cardinality of the discretized group by

$$\mathcal{H}(k, m_0, \ldots, m_L, r). \qquad (11)$$

The learnable parameters are the vectors $\mathbf{w}_{ij}^{(\ell)}$ and the biases $b_j^{(\ell)}$ for $j = 1, \ldots, m_\ell$, $\ell = 0, \ldots, L - 1$. During the training phase, they are updated through gradient-based optimization techniques, such as stochastic gradient descent (SGD). The updates aim to minimize an objective function, typically a loss function measuring the difference between the network predictions and the true labels. While in practice, GCNN architectures also include feedforward layers, we only focus in this work on the GCNN layers.

## 2.3 CNNs as a specific case of GCNNs

The convolutional layer in CNNs can be obtained as a specific case of the GCNN layer for the translation group $T$. An element $\mathbf{t}$ of the translation group corresponds to a vector $(t_1, t_2) \in \mathbb{R}^2$, and the group operation is defined as $\mathbf{t} \circ \mathbf{t}' = (t_1 + t_1', t_2 + t_2')$, meaning that one vector is shifted by the components of another. The inverse of a translation $\mathbf{t}$ is $\mathbf{t}^{-1} = (-t_1, -t_2)$, reversing the direction of the shift.

A square image can be interpreted as a function on the unit square $[0, 1] \times [0, 1]$. Setting the function values to zero outside the unit square, it can then be extended to a function on $\mathbb{R}^2$. Since $\mathbb{R}^2$ is isomorphic to the translation group, a square image can thus be viewed as a function on the translation group $T$.

To illustrate the discretization step, consider the MNIST dataset (LeCun et al., 1998), which consists of grayscale images of handwritten digits ranging from 0 to 9. Each MNIST image is represented by grayscale values. This means that $m_0$ is equal to 1, as the image is characterized solely by pixel brightness. In contrast, for an RGB image, $m_0 = 3$, corresponding to the red, green, and blue components of each pixel.

Furthermore, since the values are only on a $28 \times 28$ pixel grid, the translation group $T$ is discretized by

$$T^{784} = \left\{ \left( \frac{i}{28}, \frac{j}{28} \right) \mid i, j \in [28] \right\}.$$

In turn, we can view an MNIST image as a function $f$ on $T^{784}$. Note that the function value $f(\frac{i}{28}, \frac{j}{28})$ is the $(i, j)$-th pixel value.

In CNNs, the coefficients of the convolutional filters are typically represented by $s \times s$ weight matrices, with $s$ a prechosen integer. For simplicity, we consider $s = 3$ in the following. The convolutional filter

computes

$$\sum_{\ell,k=-1}^{1} w_{\ell+1,k+1} f\left(\frac{i+\ell}{28}, \frac{j+k}{28}\right), \quad i,j = 1, \ldots, 28.$$

For the kernel

$$\mathcal{K}_{\mathbf{w}}(u_1, u_2) = \sum_{\ell,k=-1}^{1} w_{\ell+1,k+1} \mathbf{1}\left((u_1, u_2) = \left(\frac{\ell}{28}, \frac{k}{28}\right)\right),$$

and $\mathbf{s}_{i,j} := (i/28, j/28)$, this can also be rewritten in the form (4) via

$$
\begin{aligned}
(\mathcal{K}_{\mathbf{w}} * f)(\mathbf{s}_{i,j}) &= \sum_{\mathbf{t} \in T^{784}} \mathcal{K}_{\mathbf{w}}(\mathbf{s}_{i,j}^{-1} \circ \mathbf{t}) \cdot f(\mathbf{t}) \\
&= \sum_{\mathbf{t} \in T^{784}} \sum_{\ell,k=-1}^{1} w_{\ell+1,k+1} \cdot \mathbf{1}\left(\left(t_1 - \frac{i}{28}, t_2 - \frac{j}{28}\right)\right. \\
&= \left.\left(\frac{\ell}{28}, \frac{k}{28}\right)\right) \cdot f(\mathbf{t}) \\
&= \sum_{\ell,k=-1}^{1} w_{\ell+1,k+1} f\left(\frac{i+\ell}{28}, \frac{j+k}{28}\right).
\end{aligned}
$$

Thus, in this equivalence the kernel is a linear combination of indicator functions.

## 2.4 Comparison of GCNNs and deep feedforward neural networks

GCNNs and deep feedforward neural networks (DNNs) differ in how their computational units are defined. In a DNN, each unit computes an affine transformations applied to the output of the previous layer, followed by an activation function $\sigma$. If $\mathbf{z}$ is the output of the previous layer, $\mathbf{w}$ is a weight vector, and $b$ a bias, then a unit in a DNN computes

$$\sigma\left(\mathbf{w}^{\top}\mathbf{z} - b\right). \tag{12}$$

The class of ReLU DNNs with $L$ layers (that is, $L - 1$ hidden layers and one output layer), $m_i$ units (or neurons) in layer $i = 0, \ldots, L$, and a single unit in the output layer (i.e., $m_L = 1$) is denoted by

$$\mathcal{F}(m_0, \ldots, m_L). \tag{13}$$

While both the DNN class $\mathcal{F}(m_0, \ldots, m_L)$ and the GCNN class $\mathcal{H}(k, m_0, \ldots, m_L, r)$ share the same architectural structure and apply the ReLU activation function, they differ in their unit definition, meaning that $\mathcal{F}(m_0, \ldots, m_L)$ uses DNN units (12) and $\mathcal{H}(k, m_0, \ldots, m_L, r)$, employs GCNN units (7).

## 3 VC dimension of GCNNs

We now derive upper bounds for the VC dimension of the GCNN class $\mathcal{H}(k, m_0, \ldots, m_L, r)$. We begin by formally introducing the VC dimension.

**Definition 3.1** (Growth function, VC dimension, shattering). *Let $\mathcal{H}$ denote a class of functions from $\mathcal{F}$ to $\{-1, 1\}$ (often referred to as the hypotheses class). For any non-negative integer $m$, we define the growth function of $\mathcal{H}$ as the maximum number of distinct classifications of $m$ samples that can be achieved by classifiers from $\mathcal{H}$. Specifically, it is defined as:*

$$\Pi_{\mathcal{H}}(m) := \max_{f_1, \ldots, f_m \in \mathcal{F}} |\{(h(f_1), \ldots, h(f_m)) : h \in \mathcal{H}\}|.$$

*If $\mathcal{H}$ can generate all possible $2^m$ classifications for a set of $m$ inputs, we say that $\mathcal{H}$ shatters that set. Formally, if*

$$|\{(h(f_1), \ldots, h(f_m)) : h \in \mathcal{H}\}| = 2^m,$$

*we say $\mathcal{H}$ shatters the set $\{f_1, \ldots, f_m\}$.*

*The Vapnik-Chervonenkis dimension of $\mathcal{H}$, denoted as $\mathrm{VC}(\mathcal{H})$, is the size of the largest shattered set, specifically the largest $m$ such that $\Pi_{\mathcal{H}}(m) = 2^m$. If no such largest $m$ exists, we define $\mathrm{VC}(\mathcal{H}) = \infty$.*

The VC dimension cannot be directly applied to a class of real-valued functions, such as neural networks. To address this, we follow the approach in (Bartlett et al., 2019) and use the pseudodimension as a measure of complexity. The pseudodimension is a natural extension of the VC dimension and retains similar uniform convergence properties (see (Pollard, 1990) and Theorem 19.2 in (Anthony and Bartlett, 1999)).

**Definition 3.2** (pseudodimension). *For a class $\mathcal{H}$ of real-valued functions, we define its pseudodimension as $\mathrm{VC}(\mathcal{H}) := \mathrm{VC}(\mathrm{sign}(\mathcal{H}))$, where*

$$\mathrm{sign}(\mathcal{H}) := \{\mathrm{sign}(H - b) \mid H \in \mathcal{H}, \; b \in \mathbb{R}\},$$

*where $\mathrm{sign}(x)$ is 1 for $x > 0$ and $-1$ otherwise. We write $\Pi_{\mathcal{H}}$ to denote a growth function of $\mathrm{sign}(\mathcal{H})$.*

For common parametrized function classes, the VC dimension relates to the number of parameters. For DNNs, the VC dimension further depends on the network depth, as discussed in (Bartlett et al., 2019).

The following result provides an upper bound on the VC dimension of the GCNN class $\mathcal{H}(k, m_0, \ldots, m_L, r)$. Recall that $L$ is the number of layers, $m_i$ is the number of units in layer $i = 0, \ldots, L$, weight vectors are $k$ dimensional, and $r$ is the cardinality of the discretized group.

**Theorem 3.3** (Upper Bound). *The VC dimension of the GCNN class $\mathcal{H} = \mathcal{H}(k, m_0, \ldots, m_L, r)$ is upper bounded by*

$$UB(\mathcal{H}) := L + 1 + 4 \left( \sum_{\ell=1}^{L} W_\ell \right) \log_2 \left( 8er \sum_{\ell=1}^{L} m_\ell \right), \tag{14}$$

*with $W_\ell$ the number of parameters up to the $\ell$-th layer, that is,*

$$W_\ell := \sum_{j=1}^{\ell} m_j (k m_{j-1} + 1). \tag{15}$$

To verify (15), observe that each unit in layer $i$, has $km_{i-1}$ weight parameters and one bias. A proof sketch of this theorem is provided in Section 4, with a detailed proof in the supplement.

The VC bound generalizes those known for vanilla CNNs, as shown in Lemma 12 in the supplement of (Kohler and Walter, 2023), where similar dependencies between VC dimension and number of parameters and hidden layers were derived. We further compare the sample complexity of GCNNs with DNNs. For that we rely on the VC dimension bound for DNNs with piecewise-polynomial activation functions that has been derived in Theorem 7 of (Bartlett et al., 2019). Specifically, we focus on the class of DNNs with $L$ layers and $m_i$ units in layer $i$, as defined in (13). This network class corresponds to the case where $d = 1$ and $p = 1$ in Theorem 6 of (Bartlett et al., 2019). By applying the inequality $\log_2(\log_2(x)) \leq \log_2(x)$ that holds for any $x \geq 2$, the VC dimension bound for this class becomes

$$UB(\mathcal{F}) := L + 2\left(\sum_{\ell=1}^{L} W_\ell(\mathcal{F})\right)\log_2\left(4e\sum_{\ell=1}^{L}\ell m_\ell\right),$$

where $\mathcal{F}$ is used as shorthand notation for $\mathcal{F}(m_0, \ldots, m_L)$. Here $W_i(\mathcal{F})$ represents the number of parameters of the class $\mathcal{F}$ up to the $i$-th layer, that is,

$$W_\ell(\mathcal{F}) = \sum_{j=1}^{\ell} m_j(m_{j-1} + 1). \tag{16}$$

By comparing (15) and (16), and assuming an equal number of computational units per layer for both, GCNNs and DNNs, we observe that the number of parameters in GCNNs satisfies the inequality $W_i \leq kW_i(\mathcal{F})$, with $k$ the dimension of the weight vector $\mathbf{w}$ in (6).

This together with the bound in Theorem 3.3 then yields

$$UB(\mathcal{H}) \leq 2kUB(\mathcal{F}) + 4\left(\sum_{\ell=1}^{L} W_\ell(\mathcal{H})\right)\log_2(2r). \tag{17}$$

An alternative to bounding the VC dimension based on the number of layers and neurons per layer is to express the bound in terms of the total number of trainable parameters. The main advantage of this approach is that it applies to a wider range of architectures including sparsely connected GCNNs. In this context, Bartlett et al. (2019) establishes a bound on the VC dimension for the class

$$\mathcal{F}_{W,L} := \{\mathcal{F} = \mathcal{F}(m_0, \ldots, m_\ell) \mid \ell \leq L, W_L(\mathcal{F}) \leq W\}, \tag{18}$$

consisting of DNNs with at most $L$ hidden layers and an overall number $W$ of weights. In particular, they show that there exist universal constants $c$ and $C$ such that

$$c \cdot WL \log\left(\frac{W}{L}\right) \leq \text{VC}(\mathcal{F}_{W,L})$$

$$\leq \max_{\mathcal{F} \in \mathcal{F}_{W,L}} UB(\mathcal{F}) \tag{19}$$

$$\leq C \cdot WL \log W.$$

Similarly, we consider the GCNN class

$$\mathcal{H}_{W,L,r} \coloneqq \{\mathcal{H}(k, m_0, \ldots, m_\ell, r) \mid \ell \leq L, W_L \leq W\}, \tag{20}$$

consisting of all GCNN architectures with a total number of parameters bounded by $W$, depth at most $L$, and $r$ the cardinality of the discretized group.

Inequality (17), along with the lower and upper bounds on the VC dimension of DNNs in (19) leads to the following result.

**Corollary 3.4.** *Let $W$ be the total number of parameters and $L$ the depth of the network, with $L < W^{0.99}$. There exists a universal constant $C$ such that*

$$\mathrm{VC}(\mathcal{H}_{W,L,r}) \leq C\left(\mathrm{VC}(\mathcal{F}_{W,L}) + LW \log_2(r)\right).$$

This bound obtains the nearly optimal rate as shown in the next theorem.

**Theorem 3.5** (Lower bound). *If $W, L > 3$, then there exists a universal constant $c$ such that*

$$\mathrm{VC}(\mathcal{H}_{W,L,r}) \geq c \max\{\mathrm{VC}(\mathcal{F}_{W,L}), W \log_2(r)\}.$$

By combining Corollary 3.4 and Theorem 3.5, we conclude that there exist universal constants $c$ and $C$ such that

$$
\begin{aligned}
c\Big( \mathrm{VC}(\mathcal{F}_{W,L}) + W \log_2(r) \Big) \\
\leq \mathrm{VC}(\mathcal{H}_{W,L,r}) \\
\leq C\Big( \mathrm{VC}(\mathcal{F}_{W,L}) + LW \log_2(r) \Big).
\end{aligned}
$$

In practice, the depth in GCNNs is rather small, e.g., $L = 7$ in the initial article of Cohen and Welling (2016b) for the rotated MNIST dataset. In this regime, the obtained rates in the upper and lower bound of the VC dimension nearly match.

# 4    Proof sketch of Theorem 3.3

To derive an upper bound for the VC dimension of the GCNN class $\mathcal{H} = \mathcal{H}(k, m_0, \ldots, m_L, r)$, recall that the parameter $r$ is the cardinality of the discretized group $G^r \coloneqq \{g_1, g_2, \ldots, g_r\}$. The parameters $k, m_0, \ldots, m_L$ define the network architecture. The class $\mathcal{H}$ consists therefore of all functions $\mathbf{w} \mapsto h_{\mathbf{w}}$ that can be represented by a GCNN with this architecture but different network parameter $\mathbf{w} \in \mathbb{R}^W$.

The proof of Theorem 3.3 closely follows the proof of Theorem 7 of (Bartlett et al., 2019), with adjustments to account for the structure of GCNN-units. To find an upper bound of a fixed number $m$ on the VC dimension, we need to show that neural networks in $\mathcal{H}$ cannot shatter any set of $m$ functions.

Specifically, consider an input set of $m$ functions from $G^r$ to $\mathbb{R}^{m_0}$, denoted as

$$F_m \coloneqq \{f_1, \ldots, f_m\}. \tag{21}$$

We aim to bound the number of distinct sign patterns that networks in $\mathcal{H}$ can generate when applied to functions in $F_m$ i.e.,

$$|\{\text{sign}(h_{\mathbf{w}}(f_1)), \ldots, \text{sign}(h_{\mathbf{w}}(f_m)) : \mathbf{w} \in \mathbb{R}^W\}|.$$

By bounding the number of distinct sign patterns, we obtain an upper bound on the growth function $\Pi_{\mathcal{H}}$ of sign($\mathcal{H}$). If $\Pi_{\mathcal{H}} < 2^m$, then $m$ is an upper bound for the VC dimension of $\mathcal{H}$.

For fixed network architecture and fixed input, the output of the GCNN-units depends only on the network parameters. Therefore, to understand the possible sign patterns that the network can generate on $F_m$, we need to analyze the dependence of the GCNN-units on the network parameters.

Since the ReLU activation is piecewise linear, and the sum of piecewise linear functions remains piecewise linear, each GCNN unit in any layer $\ell$, as defined in (9), can be viewed as a composition of $\ell$ piecewise linear functions. This results in a piecewise polynomial function of degree $\leq \ell$ with respect to the network parameters up to layer $\ell$. Consequently, for each layer $\ell$, the parameter space $\mathbb{R}^{W_\ell}$ can be partitioned into regions $\{P_1, \ldots, P_{S(\ell)}\}$, where within each region, the GCNN units in the $(\ell+1)$-st layer behave like a fixed polynomial function in $W_\ell$ variables of $\mathbf{w}$, of degree at most $\ell$.

As a first step to prove Theorem 3.3, we recursively find a bound for $S(\ell)$ and then determine how many different sign patterns the classifiers in sign($\mathcal{H}$) can generate within each of these regions. The following lemma establishes how the upper bound for $S(\ell)$ evolves from $S(\ell-1)$. The proof is provided in the supplementary material.

**Lemma 4.1.** *Let $\mathcal{H}$ be the class of GCNNs defined in (20) with $\leq W_\ell$ network parameters up to layer $\ell \in \{1, \ldots, L\}$. If $F_m$ is the class of functions defined in (21) and $S(\ell)$ is as defined above, then, for $\ell = 0, 1, \ldots, L-1$,*

$$S(\ell+1) \leq 2 \left( \frac{2em_{\ell+1}mr(\ell+1)}{W_{\ell+1}} \right)^{W_{\ell+1}} S(\ell). \tag{22}$$

*Moreover, the GCNN units $\{h_{\ell+1,j}(g) \mid j \leq m_{\ell+1}, f \in F_m, g \in G^r\}$ with $h_{\ell+1,j}$ defined for different functions $f \in F_m$ in (8) and (9) are piecewise polynomials of degree $\leq \ell + 1$ in the network parameters.*

Next, by applying Lemma 1 in (Bartlett et al., 1998), which provides an upper bound on the number of distinct sign patterns that can be generated by polynomials of finite degree, we obtain a bound on the growth function $\Pi_{\mathcal{H}}$.

The proof is provided in the supplementary material.

**Lemma 4.2.** *Let $\mathcal{H}$ be the class of GCNNs defined in (11) with at most $W_\ell$ of parameters up to layer $\ell \leq L$ and $m_\ell$ GCNN units in layer $\ell$. Let $m > 0$ be an integer, then*

$$\Pi_{\mathcal{H}}(m) \leq 2^L \prod_{\ell=1}^{L} \left( \frac{2emrm_\ell\ell}{W_\ell} \right)^{W_\ell} 2 \left( \frac{2emL}{W_L + 1} \right)^{(W_L+1)}.$$

The relationship between VC-dimension and growth functions yields Theorem 3.3. A complete proof can be found in supplementary material.

# 5 Proof Sketch of Theorem 3.5

To establish a lower bound for the VC dimension of the GCNN class, we recall that $\mathcal{H}_{W,L,r}$ as defined in (20) represents the class of GCNNs with resolution $r$, $k$ the dimension of the kernel space, at most $L$ layers, and no more than $W$ parameters. Additionally, let $G^r := \{g_1, g_2, \ldots, g_r\}$ be a discretized group containing the identity element $e$.

To derive a lower bound, we aim to prove that the VC dimension of $\mathcal{H}_{W,L,r}$ exceeds a given integer $m$. For this, it is sufficient to find networks in $\mathcal{H}_{W,L,r}$ that can shatter a set of $m$ input functions.

In particular, to prove Theorem 3.5, one needs to show the existence of a universal constant $c > 0$ such that

$$\mathrm{VC}(\mathcal{H}_{W,L,r}) \geq c \cdot \mathrm{VC}(\mathcal{F}_{W,L}), \tag{23}$$

and

$$\mathrm{VC}(\mathcal{H}_{W,L,r}) \geq c \cdot W \lfloor \log_2 r \rfloor. \tag{24}$$

Here, $\mathcal{F}_{W,L}$ is the class of DNNs with at most $L$ hidden layers and at most $W$ weights, as defined in (18).

To prove (23), the first step is to establish a connection between DNNs and GCNNs. The next lemma demonstrates how a DNN can be associated with a GCNN with the same number of parameters and layers. Specifically, it states that when a DNN is evaluated on inputs $f(g)$, where $f \in F_m$ and $g \in G^r$ (with $F_m$ defined as in (21)), the sum of its outputs, taken over different elements of $G^r$, matches the output of the corresponding GCNN for the same function $f$.

Before stating the lemma, recall that $\mathcal{F}(m_0, \ldots, m_L)$, as defined in (13), denotes the class of fully connected feedforward ReLU networks with $L$ layers, where $m_i$ denotes the number of units in the $i$-th layer for $i = 1, \ldots, L$. The outputs of the last hidden layer of any neural network $\tilde{h}_{\mathbf{w}} \in \mathcal{F}(m_0, \ldots, m_L)$ with parameters $\mathbf{w}$ can be represented as a vector of size $m_L$, this is, $(\tilde{h}_{\mathbf{w}}^{(1)}, \ldots, \tilde{h}_{\mathbf{w}}^{(m_L)})$.

**Lemma 5.1.** *Consider GCNNs where the G-correlation in (4) uses kernels from a one-dimensional vector space with a fixed basis given by the indicator function of the identity element $e$. For every $\tilde{h}_{\mathbf{w}} \in \mathcal{F}(m_0, \ldots, m_L)$ there exists a GCNN $h_{\mathbf{w}}$, with the same number of channels in each layer, i.e., $h_{\mathbf{w}} \in \mathcal{H}(1, m_0, \ldots, m_L, r)$ and parameters $\mathbf{w}$, such that for any input function $f : G^r \to \mathbb{R}^{m_0}$*

$$\sum_{i=1}^{m_L} \sum_{j=1}^{r} \tilde{h}_{\mathbf{w}}^{(i)}(f(g_j)) = h_{\mathbf{w}}(f).$$

The proof of this lemma is provided in the supplementary material. The lemma implies that instead of finding GCNNs that shatter $F_m$, we can also find DNNs in the class $\mathcal{F}(m_0, \ldots, m_L)$ shattering $F_m$. For the construction of this DNN architecture we define sums of so-called 'indicator' neural networks, i.e., DNNs with ReLU activation functions that approximate indicator functions over a specified interval. The endpoints of this interval act as parameters of the neural networks, allowing to adjust the interval by modifying these parameters.

For $a < b$ and $\epsilon > 0$, the shallow ReLU network

$$\mathbf{1}_{(a,b,\epsilon)}(x) = \frac{1}{\epsilon}\Big(\big(x - (a - \epsilon)\big)_+ - \big(x - a\big)_+$$
$$+ \big(x - b\big)_+ - \big(x - (b + \epsilon)\big)_+\Big), \tag{25}$$

with four neurons in the hidden layer and $(x)_+ := \max(x, 0)$ the ReLU activation function, approximates the indicator function on $[a, b]$ in the sense that $\mathbf{1}_{(a,b,\epsilon)}(x)$ is 1 if $x \in [a, b]$ and 0 if $x < a - \epsilon$ or $x > b + \epsilon$.

To show (23), we construct a DNN architecture from $m_0 + 1$ smaller DNN classes. One of these DNN classes is capable of shattering a set of $\mathrm{VC}(\mathcal{F}_{W,L})$ vectors in $\mathbb{R}^{m_0}$, while the remaining $m_0$ DNN classes produce 'indicator' networks, ensuring that the combined DNNs vanish outside a certain $m_0$-dimensional hypercube. We show, that this class of networks shatters a set of $\mathrm{VC}(\mathcal{F}_{W,L})$ input functions from the set $\{f : G^r \to \mathbb{R}^{m_0}\}$. By applying Lemma 5.1, we further show that a class of GCNNs with at most $5W$ weights and $L$ layers can shatter the same set of $\mathrm{VC}(\mathcal{F}_{W,L})$ input functions. Finally, we use that $\mathrm{VC}(\mathcal{F}_{tW,L}) \geq c_t \cdot \mathrm{VC}(\mathcal{F}_{W,L})$ for a constant $t$ with $c_t < 1$.

**Lemma 5.2.** *For $L > 3$ let $\mathcal{H}_{5W,L+1,r}$ be the class of GCNNs defined as in (20) and $\mathcal{F}_{W,L}$ be the class of DNNs defined as in (18). Then*

$$\mathrm{VC}(\mathcal{H}_{5W,L+1,r}) \geq \mathrm{VC}(\mathcal{F}_{W,L}).$$

**Corollary 5.3.** *In the setting of Lemma 5.2, if the number of layers $L > 3$ and $L \leq W^{0.99}$, then there exists a constant $c$ such that*

$$\mathrm{VC}(\mathcal{H}_{W,L,r}) \geq c \cdot \mathrm{VC}(\mathcal{F}_{W,L}).$$

To prove the lower bound in (24), we first show that an 'indicator' neural network class can shatter input functions $F_m \subset \{f : G^r \to \mathbb{R}\}$ with $m := \lfloor \log_2 r \rfloor$, by adjusting its parameters, i.e., the interval endpoints to the values of the input functions from $F_m$. By using Lemma 5.1 again, we then find a corresponding GCNN, which also shatters $F_m$. This shows, that for any two numbers $A < B$, the GCNN can serve as an indicator function, outputting zero for any input function that maps outside the interval $[A, B]$. This ensures that the GCNN only shatters functions within the specified range. The next lemma provides a formal statement. The proof is provided in the supplementary material.

**Lemma 5.4.** *Let $\mathcal{H}_{4,L,r}$ be the class of GCNNs defined as in (20). Then*

$$VC(\mathcal{H}_{4,L,r}) \geq \lfloor \log_2 r \rfloor.$$

*Moreover, for any two numbers $A < B$, there exists a finite subclass of GCNNs $\mathcal{H} \subset \mathcal{H}_{4,L,r}$ that shatters a set of $\log_2 r$ input functions*

$$\{f_i : G^r \to [A, B] \mid i = 1, \ldots, \log_2 r\},$$

*and outputs zero for any input function $f : G^r \to \mathbb{R} \setminus [A, B]$.*

By using that each class $\mathcal{H}_{4,L,r}$ can shatter a set of functions

$$\{f_i : G^r \to [A_j, B_j] | i = 1, \ldots, \lfloor \log_2 r \rfloor\}$$

with disjoint intervals $[A_j, B_j]$, $j = 1, \ldots, \tilde{W}$, we can find a class of GCNNs that shatters a set of functions with $\tilde{W}\lfloor \log_2 r \rfloor$ elements. By choosing $W = 4\tilde{W}$, this shows that

$$\mathrm{VC}(\mathcal{H}_{W,L,r}) \geq \frac{1}{4}W\lfloor \log_2 r \rfloor.$$

The following corollary is a formal statement of the previous inequality. It is proved in the supplementary material.

**Corollary 5.5.** *The VC dimension of the class $\mathcal{H}_{4W,L,r}$, consisting of GCNNs with $4W$ weights, $L$ layers, and resolution $r$ satisfies the inequality*

$$\mathrm{VC}\left(\mathcal{H}_{4W,L,r}\right) \geq W\lfloor \log_2 r \rfloor.$$

By combining Corollaries 5.3 and 5.5, we obtain the lower bound

$$\mathrm{VC}(\mathcal{H}_{W,L,r}) \geq \max\left\{c\,\mathrm{VC}(\mathcal{F}_{W,L}), \frac{1}{4}W\lfloor \log_2 r \rfloor\right\},$$

thereby proving Theorem 3.5.

# 6    Conclusion

In this work, we established nearly-tight VC dimension bounds for a class of GCNNs. The bounds reveal that the complexity of GCNNs depends on the number of layers, the number of weights, and the resolution of the group acting on the input data.

While the VC dimensions of GCNNs and deep neural networks (DNNs) are similar, for GCNNs an extra term $W\lfloor \log_2 r \rfloor$ occurs that increases logarithmically in the resolution $r$ of the input data. This finding aligns with previous work (Petersen and Sepliarskaia, 2024) showing that, when the resolution of the group approaches infinity, the VC dimension of the GCNN becomes infinite as well. The logarithmic scaling with respect to the resolution highlights the sensitivity of GCNNs to the discretization of the group, providing deeper insight into how data impacts the model complexity and generalization capabilities.

# A    Supplementary material

## A.1    Proof of Theorem 3.3

Recall that

$$\mathcal{H} = \mathcal{H}(k, m_0, \ldots, m_L, r), \tag{26}$$

where $r$ is the cardinality of the discretized group $G^r := \{g_1, g_2, \ldots, g_r\}$. The parameter $k$ determines the number of basis functions

$$\mathrm{K}_s : G^r \to \mathbb{R}, \quad s = 1, \ldots, k, \tag{27}$$

in the parametrization of the kernel function

$$\mathcal{K}_{\mathbf{w}} = \sum_{s=1}^{k} w_s \mathrm{K}_s.$$

The other parameters $m_0, \ldots, m_L$ define the network architecture, and $W_\ell$ represents the number of parameters in the GCNN up to layer $\ell$. The class $\mathcal{H}$ consists of all functions that can be represented by a neural network with this architecture.

We restate Theorem 3.3 for convenience:

**Theorem A.1** (Theorem 3.3). *The VC dimension of the GCNN class $\mathcal{H} = \mathcal{H}(k, m_0, \ldots, m_L, r)$ with $r > 1$, is bounded from above by*

$$UB(\mathcal{H}) := L + 1 + 4 \left( \sum_{\ell=1}^{L} W_\ell \right) \log_2 \left( 8er \sum_{\ell=1}^{L} m_\ell \right). \tag{28}$$

For the proof, we consider an input consisting of $m$ functions from $G^r$ to $\mathbb{R}^{m_0}$, denoted by

$$F_m := \{f_1, \ldots, f_m\}. \tag{29}$$

To prove Lemma 4.1, we use the following known result:

**Lemma A.2.** *[Lemma 1, Bartlett et al. (1998)] Let $p_1, \ldots, p_{\tilde{m}}$ be polynomials of degree at most $t$ depending on $n \leq \tilde{m}$ variables. Then*

$$\Pi := |\{(\mathrm{sign}(p_1(x)), \ldots, \mathrm{sign}(p_m(x))) : x \in \mathbb{R}^n\}| \leq 2 \left( \frac{2e\tilde{m}t}{n} \right)^n.$$

Recall that $S(\ell)$ is the number of regions in the parameter space $\mathbb{R}^{W_\ell}$, such that in each region, the GCNN units in the $\ell$-th layer (denoted by $\{h_{\ell,j}(g) \mid j \leq m_\ell, f \in F_m, g \in G^r\}$) behave like a fixed polynomial of degree at most $\ell$ in the $W_\ell$ network parameters that occur up to layer $\ell$.

**Lemma A.3.** *[Lemma 4.1] Let $\mathcal{H}$ be the class of GCNNs defined in (26), with at most $W_\ell$ parameters up to layer $\ell \in \{1, \ldots, L\}$. If $F_m$ is the class of functions defined in (29), and $S(\ell)$ is as defined above, then for $\ell = 0, 1, \ldots, L - 1$,*

$$S(\ell+1) \leq 2 \left( \frac{2em_{\ell+1}mr(\ell+1)}{W_{\ell+1}} \right)^{W_{\ell+1}} S(\ell). \tag{30}$$

14

*Moreover, the GCNN units $\{h_{\ell+1,j}(g) \mid j \leq m_{\ell+1}, f \in F_m, g \in G^r\}$ with $h_{\ell+1,j}$ defined for different functions $f \in F_m$, are piecewise polynomials of degree $\leq \ell + 1$ in the network parameters.*

*Proof.* As a first step of the proof, we show that any GCNN unit $h_{\ell,j}$ of any layer $\ell \in \{0, \ldots, L\}$ and $j \in \{1, \ldots, m_\ell\}$ forms a piecewise polynomial of degree at most $\ell$. We proceed by induction on the layers $\ell$.

For the base case $\ell = 0$, the GCNN units $h_{0,j}$ for $j \leq m_0$ correspond to the input of the network, which is independent of the network parameters. Therefore, $h_{0,j}$ are polynomials of degree 0.

Assume the statement holds for all layers up to $\ell$. We now prove it for layer $\ell + 1$. The GCNN unit in layer $\ell + 1$ is defined by a convolution with the feature maps from the previous layer, that is,

$$h_{\ell+1,j} = \sigma \left( \sum_{i=1}^{m_\ell} \mathcal{K}_{\mathbf{w}_{ij}^{(\ell)}} * h_{\ell,i} - b_j^{(\ell)} \right),$$

where the convolutional filter is expanded in terms of the fixed basis functions $\mathsf{K}_s$ via $\mathcal{K}_{\mathbf{w}} = \sum_{s=1}^{k} w_s \mathsf{K}_s$. For fixed network parameters, $g \mapsto h_{\ell,j}(g)$ is a function of the group, with $h_{\ell,j}(g) \in \mathbb{R}$ for any $g \in G^r$. By the induction hypothesis, $h_{\ell,j}(g)$ are piecewise polynomials of degree at most $\ell$ with respect to the network parameters, with the polynomial pieces depending on the network input and the group element $g$.

Next, observe that for any input and any group element $g$, the term $(\mathsf{K}_s * h_{\ell,j})(g)$ can be written as

$$(\mathsf{K}_s * h_{\ell,j})(g) = \sum_{g' \in G^r} \mathsf{K}_s(g^{-1} \circ g') \cdot h_{\ell,j}(g').$$

Since $h_{\ell,j}(g')$ is a piecewise polynomial of degree at most $\ell$, it follows that $(\mathsf{K}s * h_{\ell,j})(g)$ is also a piecewise polynomial of degree at most $\ell$. Thus, the convolution

$$(\mathcal{K}_{\mathbf{w}} * h_{\ell,j})(g) = \sum_{s=1}^{k} w_s (\mathsf{K}_s * h_{\ell,j})(g)$$

is a weighted sum of piecewise polynomials, which remains a piecewise polynomial. However, multiplying by the weights $w_s$ increases the degree of the polynomial, making it at most $\ell + 1$. Subtracting the bias term $b_j^{(\ell)}$ and applying the ReLU activation function may increase the number of pieces, but does not increase the degree of the polynomials. Therefore, for any input and any group element $g$, the GCNN unit $h_{\ell+1,j}(g)$ remains a piecewise polynomial with degree $\leq \ell + 1$. This completes the proof by induction.

Next, we show (30). Each GCNN unit in layer $\ell + 1$ is computed by

$$\sigma \left( \sum_{i=1}^{m_\ell} \mathcal{K}_{\mathbf{w}_{ij}^{(\ell)}} * h_{\ell,i} \right),$$

with $\sigma$ being the ReLU activation function $\sigma(x) = \max\{x, 0\}$. As mentioned above, applying the ReLU function can increases the number of regions in the parameter space where the GCNN units behave as polynomials. This occurs, as the ReLU function either outputs the input itself (for positive values) or zero (otherwise). As a result its application decomposes each of the $S(\ell)$ regions of the parameter space in layer $\ell$ in multiple subregions. To bound this number of subregions, we need to count the number of possible sign pattern that can arise after applying the ReLU activation.

15

Fixing one of the $S(\ell)$ regions of layer $\ell$, by definition, all functions $h_{\ell,j}(g)$ are polynomials in the parameters of degree at most $\ell$. Each of the $m_{\ell+1}$ GCNN unit in layer $\ell+1$, are then also a polynomial of degree at most $\ell+1$, leading to an overall amount of polynomials of $m_{\ell+1}mr$, where $m$ is the number of input function defined in the (29) and $r$ is the resolution. Applying Lemma A.2 to the $\tilde{m} = m_{\ell+1}mr$ polynomials of degree $t = \ell+1$ depending on $n = W_{\ell+1}$ parameters leads to an overall number of different sign patterns of each region of $S(\ell)$ of

$$2\left(\frac{2em_{\ell+1}mr(\ell+1)}{W_{\ell+1}}\right)^{W_{\ell+1}}.$$

This shows (30) and concludes the proof.

$\square$

**Lemma A.4.** *[Lemma 4.2] Let $\mathcal{H}$ be the class of GCNNs defined in (26), with at most $W_\ell$ parameters up to layer $\ell \leq L$, and $m_\ell$ GCNN units in layer $\ell$. For any integer $m > 0$, the growth function of this class can be bounded by*

$$\Pi_{\mathcal{H}}(m) \leq 2^L \prod_{\ell=1}^{L}\left(\frac{2emrm_\ell\ell}{W_\ell}\right)^{W_\ell} 2\left(\frac{2emL}{W_L+1}\right)^{W_L+1}.$$

*Proof.* Lemma A.3 shows that after $L$ layers, there are at most

$$2^L \prod_{\ell=1}^{L}\left(\frac{2emrm_\ell\ell}{W_\ell}\right)^{W_\ell}$$

regions in the parameter space $\mathbb{R}^W$, on which the GCNN units in the last layer $\{h_{L,i}(g) \mid i \leq m_L, f \in F_m, g \in G^r\}$ behave like a fixed polynomial function of degree $\leq L$ in $W$ variables.

Recall that the final output of the neural network, is obtained by applying average pooling to the outputs of the GCNN units in the last layer. This implies that, for a fixed network architecture and input, the output of the neural network is a piecewise polynomial of degree at most $L$, depending on all $W_L$ network parameters. Since there are $m$ possible inputs $f_1, \ldots, f_m$, we get $m$ piecewise polynomials, each corresponding to one of these inputs. Bounding the growth function $\Pi_{\mathcal{H}}(m)$ now means we need to count the number of different sign patterns that arises for classifiers in $\text{sign}(\mathcal{H})$. For that, we recall that by Definition 3.2 in the main article,

$$\text{sign}(\mathcal{H}) := \{\text{sign}(h_{\mathbf{w}} - b) \mid h_{\mathbf{w}} \in \mathcal{H}, \mathbf{w} \in R^{W_L}, b \in \mathbb{R}\}.$$

Applying Lemma A.2 to $m$ polynomials of the form $h_{\mathbf{w}} - b$ of degree at most $L$ and $W_L + 1$ variables leads to no more than

$$2\left(\frac{2emL}{W_L+1}\right)^{W_L+1} \tag{31}$$

distinct sign patterns that the classifiers in $\text{sign}(\mathcal{H})$ can produce.

Thus, the growth function within each region, where the GCNN units in the last layer $\{h_{L,i}(g) \mid i \leq m_L, f \in F_m, g \in G^r\}$ behave like a fixed polynomial function in $W$ variables, is bounded by (31). As a result, we conclude that the overall growth function $\Pi_{\mathcal{H}}(m)$ is bounded by

$$S(L) \cdot 2\left(\frac{2emL}{W_L+1}\right)^{W_L+1} = 2^L \prod_{\ell=1}^{L}\left(\frac{2emrm_\ell\ell}{W_\ell}\right)^{W_\ell} \cdot 2\left(\frac{2emL}{W_L+1}\right)^{W_L+1}.$$

This completes the proof. $\qquad \square$

For the proof of the Theorem A.1 we also use the following technical lemma

**Lemma A.5.** *[Lemma 16, Bartlett et al. (1998)] Suppose $2^{\tilde{m}} \leq 2^{\kappa} \left(\frac{\tilde{m}\cdot\tilde{r}}{\tilde{w}}\right)^{\tilde{w}}$ for some $\tilde{r} \geq 16$ and $\tilde{m} \geq \tilde{w} \geq \kappa \geq 0$. Then $\tilde{m} \leq \kappa + \tilde{w}\log_2(2\tilde{r}\log_2 \tilde{r})$.*

*Proof of Theorem A.1.* Let $m := \mathrm{VC}(\mathcal{H})$. For convenience, define the sum

$$\tilde{W} := \sum_{i=1}^{L} W_i, \tag{32}$$

where $W_i$ denotes the number of parameters of a GCNN in $\mathcal{H}$ up to layer $i$.

We consider two complementary cases and prove the theorem for each of them separately.

**Case 1:** $m < \tilde{W} + W_L + 1$. In this case, we have $\tilde{W} + W_L + 1 < 3\tilde{W} < UB(\mathcal{H})$, where $UB(\mathcal{H})$ is defined in (28). For the latter inequality we use that $\log_2\left(8er\sum_{\ell=1}^{L} m_\ell\right) > 1$. Therefore, Theorem A.1 holds.

**Case 2:** $m \geq \tilde{W} + W_L + 1$. Since $m$ represents the VC dimension of $\mathcal{H}$, it follows from the definition of the VC dimension (see Definition 3.1 in the main article) that $\Pi_{\mathcal{H}}(m) = 2^m$. Applying Lemma A.4 gives us

$$\Pi(m) = 2^m \leq 2^{L+1} \prod_{\ell=1}^{L} \left(\frac{2emrm_\ell\ell}{W_\ell}\right)^{W_\ell} \left(\frac{2emL}{W_L+1}\right)^{W_L+1}. \tag{33}$$

Next, we apply the weighted arithmetic-geometric mean (AM-GM) inequality to the right side of (33), using weights $W_\ell/(\tilde{W}+W_L+1)$ for $\ell = 1, 2, \ldots, L$, and $W_L/(\tilde{W}+W_L+1)$, where $\tilde{W}$ is defined in (32). This yields

$$2^m \leq 2^{L+1} \left(\frac{2em(r\sum_{\ell=1}^{L}\ell m_\ell + L)}{\tilde{W}+W_L+1}\right)^{\tilde{W}+W_L+1}.$$

The last step of the proof involves applying Lemma A.5 in Bartlett et al. (1998) to this inequality, which provides an upper bound for $m$. Before doing so, we must verify that all conditions of the lemma are satisfied. In our case, $\tilde{m}$ corresponds to $m$, $\kappa$ to $L+1$, $\tilde{w}$ to $\tilde{W}+W_L+1$, and $\tilde{r}$ to $2e(r\sum_{\ell=1}^{L}\ell m_\ell + L)$. Since $r\sum_{\ell=1}^{L}\ell m_\ell + L > 2$, we have $\tilde{r} > 16$. Moreover, we are considering the case where $m \geq \tilde{W}+W_L+1$, and it is straightforward to verify that $\tilde{W}+W_L+1 \geq L+1 > 0$. Therefore all conditions of Lemma A.5 in Bartlett et al. (1998) are indeed satisfied and we obtain

$$m \leq (L+1) + 2\tilde{W}\log_2\left(4e\left(r\sum_{\ell=1}^{L}\ell m_\ell + L\right) \cdot \log_2\left(2e\left(r\sum_{\ell=1}^{L}\ell m_\ell + L\right)\right)\right).$$

To simplify this inequality, we use that for all $a \geq 1$, $\log_2(2a\log_2 a) = \log_2(2a) + \log_2(\log_2 a) \leq 2\log_2(2a)$. Substituting $a = 2e\left(r\sum_{\ell=1}^{L}\ell m_\ell + L\right)$, we note that $a \leq 4er\sum_{\ell=1}^{L}\ell m_\ell$ and obtain

$$m \leq (L+1) + 4\tilde{W}\log_2\left(8er\sum_{\ell=1}^{L}\ell m_\ell\right),$$

completing the proof of the theorem. $\qquad \square$

## A.2 Proof of Theorem 3.5

In this section, we provide the detailed proof of Theorem 3.5, along with the proofs for Lemmas 5.1, 5.2, 5.4, and their corresponding Corollaries 5.3 and 5.5.

The class

$$\mathcal{H}_{W,L,r} := \{\mathcal{H}(k, m_0, \ldots, m_\ell, r) \mid \ell \leq L,\, W_L \leq W\}, \tag{34}$$

includes all GCNN architectures with a total number of parameters bounded by $W$, a maximum depth of $L$, and $r$ representing the cardinality of the discretized group $G^r := \{g_1, g_2, \ldots, g_r\}$ containing the identity element $e$.

Next, we recall that $\mathcal{F}(m_0, \ldots, m_L)$ represents the class of fully connected feedforward ReLU networks with $L$ layers, where $m_i$ denotes the number of units in the $i$-th layer for $i = 1, \ldots, L$. The output of the last hidden layer of any neural network $\tilde{h}_{\mathbf{w}} \in \mathcal{F}(m_0, \ldots, m_L)$, with parameters $\mathbf{w}$, can be written as a vector of size $m_L$, that is, $(\tilde{h}_{\mathbf{w}}^{(1)}, \ldots, \tilde{h}_{\mathbf{w}}^{(m_L)})$.

Finally, we define the class

$$\mathcal{F}_{W,L} := \{\mathcal{F} = \mathcal{F}(m_0, \ldots, m_\ell) \mid \ell \leq L,\, W_L(\mathcal{F}) \leq W\}, \tag{35}$$

consisting of DNNs with at most $L$ hidden layers and a total number of weights not exceeding $W$.

**Lemma A.6.** *(Lemma 5.1) Consider GCNNs where the G-correlation uses kernels from a one-dimensional vector space with a fixed basis given by the indicator function of the identity element $e$. For every $\tilde{h}_{\mathbf{w}} \in \mathcal{F}(m_0, \ldots, m_L)$, there exists a GCNN $h_{\mathbf{w}}$ with the same number of channels in each layer, i.e., $h_{\mathbf{w}} \in \mathcal{H}(1, m_0, \ldots, m_L, r)$, and parameters $\mathbf{w}$, such that for any input function $f : G^r \to \mathbb{R}^{m_0}$,*

$$\sum_{i=1}^{m_L} \sum_{j=1}^{r} \tilde{h}_{\mathbf{w}}^{(i)}(f(g_j)) = h_{\mathbf{w}}(f).$$

*Proof.* Write $\mathcal{H} := \mathcal{H}(1, m_0, \ldots, m_L, r)$. Consider a fixed input function $f : G^r \to \mathbb{R}^{m_0}$ and a weight vector $\mathbf{w} \in \mathbb{R}^W$. Recall that the number of parameters in a GCNN is given by

$$W_L := \sum_{j=1}^{L} m_j(km_{j-1} + 1), \tag{36}$$

where $k$ is the dimension of the kernel space. In our case $k = 1$ and the number of parameters for a GCNN with architecture $\mathcal{H}$ is

$$W_L = \sum_{j=1}^{L} m_j(m_{j-1} + 1). \tag{37}$$

This coincides with the number of parameters in a DNN with architecture $\mathcal{F}(m_0, \ldots, m_L)$. Consequently, the same weight vector $\mathbf{w} \in \mathbb{R}^W$ defines both, a DNN function $\tilde{h}_{\mathbf{w}}$ and a GCNN function $h_{\mathbf{w}} \in \mathcal{H}$ when the input is fixed to $f$.

We now show that the outputs of the computational units in $\tilde{h}_{\mathbf{w}}$ and $h_{\mathbf{w}}$ are equal when applied to $f(g)$ and $g$, respectively. Specifically, we aim to prove that

$$\tilde{h}_{\ell,i}\big(f(g_j)\big) = h_{\ell,i}(g_j),$$

where $\tilde{h}_{\ell,i}$ denotes a DNN computational unit in layer $\ell$ of $\tilde{h}_{\mathbf{w}}$, with parameters fixed to $\mathbf{w}$, and $h_{\ell,i}$ represents a GCNN computational unit in layer $\ell$, with parameters fixed to $\mathbf{w}$ and input set to $f$. We prove this by induction on the layer $\ell$.

The statement holds trivially for the input layer, as $\tilde{h}_{0,i}(f(g_j)) = h_{0,i}(g_j)$ for any $g_j \in G^r$. Assuming it holds for all layers up to $\ell - 1$, we now prove it for layer $\ell$.

Let K denote the indicator of the identity element $e \in G^r$. By calculating the G-correlation between K and $f$, we obtain $K * f = f$. Combining this with the definition of the GCNN unit (see (9) in the main article) and the induction hypothesis, we have

$$
\begin{aligned}
\tilde{h}_{\ell,i}(g_j) &:= \sigma\left(\sum_{t=1}^{m_{\ell-1}} \mathbf{w}_{t,i}^{(\ell-1)} \tilde{h}_{\ell-1,t}(g_j) - b_i^{(\ell)}\right) \\
&= \sigma\left(\sum_{t=1}^{m_{\ell-1}} \mathbf{w}_{t,i}^{(\ell-1)} h_{\ell-1,t}(g_j) - b_i^{(\ell)}\right) && \text{(induction assumption)} \\
&= \sigma\left(\sum_{t=1}^{m_{\ell-1}} \left(\mathbf{w}_{t,i}^{(\ell-1)} K * h_{\ell-1,t}\right)(g_j) - b_i^{(\ell)}\right) && \text{(property of K)} \\
&= \sigma\left(\sum_{t=1}^{m_{\ell-1}} \left(\mathcal{K}_{\mathbf{w}_{t,i}^{(\ell-1)}} * h_{\ell-1,t}\right)(g_j) - b_i^{(\ell)}\right) && \text{(definition of learned kernel)} \\
&= h_{\ell,i}(g_j) && \text{(definition of GCNN unit).}
\end{aligned}
$$

This shows that the outputs of the computational units in $\tilde{h}_{\mathbf{w}}$ and $h_{\mathbf{w}}$ are equal when applied to $f(g)$ and $g$, respectively.

Finally, the outputs of $\tilde{h}_{\mathbf{w}} := (\tilde{h}_{\mathbf{w}}^{(1)}, \ldots, \tilde{h}_{\mathbf{w}}^{(m_L)})$ can be rewritten into the form

$$\sum_{i=1}^{m_L} \sum_{j=1}^{r} \tilde{h}_{\mathbf{w}}^{(i)}(f(g_j)) = \sum_{i=1}^{m_L} \sum_{j=1}^{r} \tilde{h}_{L,i}(g_j) = \sum_{i=1}^{m_L} \sum_{j=1}^{r} h_{L,i}(g_j) = h_{\mathbf{w}}(f),$$

concluding the proof of the lemma. $\qquad\square$

Next, we prove Lemma 5.2. The key ideas and steps of the proof have already been outlined in the main article, so here we will focus on the formal statements that still needs to be established.

Recall that the indicator neural network

$$\mathbf{1}_{(a,b,\epsilon)} \tag{38}$$

is a shallow ReLU network with four neurons in the hidden layer (see (25) in the main article). It approximates the indicator function on the interval $[a, b]$ in the sense that $\mathbf{1}_{(a,b,\epsilon)}(x) = 1$ if $x \in [a, b]$, and $\mathbf{1}_{(a,b,\epsilon)}(x) = 0$ if $x < a - \epsilon$ or $x > b + \epsilon$.

**Lemma A.7.** *[Lemma 5.2] For $L > 3$ let $\mathcal{H}_{6W,L+1,r}$ be the class of GCNNs defined as in (34) and $\mathcal{F}_{W,L}$ be the class of DNNs defined as in (35). Then*

$$\mathrm{VC}(\mathcal{H}_{6W,L+1,r}) \geq \mathrm{VC}(\mathcal{F}_{W,L}).$$

*Proof.* Let $m$ be the VC dimension of the class of DNNs $\mathcal{F}_{W,L}$. There are $2^m$ possible binary classifications for a set of $m$ elements, subsequently denoted by $d = 2^m$.

By definition, there exists a natural number $m_0$ and a set of $m$ vectors

$$\mathcal{Y} := \{\mathbf{y}_1, \ldots, \mathbf{y}_m\} \subset \mathbb{R}^{m_0}, \tag{39}$$

that can be shattered by a subset of networks $\tilde{\mathcal{H}} \subset \mathcal{F}_{W,L}$. Since there are no more than $d$ distinct classifiers for $\mathcal{Y}$, the class $\tilde{\mathcal{H}}$ consists of at most $d$ DNN functions.

Next, we construct a DNN architecture using $m_0 + 1$ smaller DNN classes. One of these classes is $\tilde{\mathcal{H}}$, while the remaining $m_0$ classes consist of "indicator" networks, as described in (38). These indicator networks ensure that the combined DNN vanishes outside a certain $m_0$-dimensional hypercube. To define this hypercube, we use the set $\mathcal{Y}$ from above.

Specifically, we choose numbers $A > \max_{\mathbf{y} \in \mathcal{Y}} \|\mathbf{y}\|_\infty + 1$ and $B > A$, and define the $m_0$-dimensional hypercube

$$\Pi := \{\mathbf{y} = (y_1, \ldots, y_{m_0}) \mid A \leq y_i \leq B\}.$$

To construct a DNN that vanishes on $\Pi$, we define an approximate indicator function for $\Pi$, using a DNN with $m_0$-dimensional input $\mathbf{y} = (y_1, \ldots, y_{m_0})$:

$$I : \mathbb{R}^{m_0} \to \mathbb{R}, \quad I(\mathbf{y}) := \frac{1}{m_0} \sum_{i=1}^{m_0} \mathbf{1}_{(A,B,0.5)}(y_i),$$

where $\mathbf{1}_{(A,B,0.5)}(y_i)$ is an indicator network that approximates the indicator function for values within $(A, B)$.

The final DNN is formed by combining functions from $\tilde{\mathcal{H}}$ with the indicator function $I$. Since DNNs can be summed if they have the same depth, we adjust the depth of $I$ to match the depth of the functions from $\tilde{\mathcal{H}}$ while ensuring that $I$ remains constant on $\mathcal{Y}$ and $\Pi$. Specifically, we use the fact that for $I(\mathbf{y}) > 0$, $\sigma(I(\mathbf{y})) = I(\mathbf{y})$ for the ReLU activation function $\sigma(x) = \max\{x, 0\}$ (for any $\mathbf{y}$ from $\Pi$ or $\mathcal{Y}$). This means that by composing $I$ with the required number of ReLU functions, we can construct a DNN that satisfies the desired properties. This construction requires at most $L < W$ additional weights.

To complete the proof, we need to show that there are $m$ input functions $F_m := \{f_1, \ldots, f_m\} \subset \{f : G^r \to \mathbb{R}^{m_0}\}$ that can be shattered by GCNNs from $\mathcal{H}_{6W,L+1,r}$. As the set $F_m$, we choose functions defined by $f_i(e) = \mathbf{y}_i$ and $f_i(g) \in \Pi$ for $g \in G^r$ and $g \neq e$.

By the definition of shattering (see Definition 3.1 in the main article), to prove that $F_m$ is shattered, it is sufficient to show that for any binary classifier $\mathcal{C} : F_m \to \{0, 1\}$, there exists a corresponding function in $\mathrm{sign}(\mathcal{H}_{6W,L+1,r})$ whose values coincide with those of $\mathcal{C}$ on $F_m$.

Choose $\tilde{h} \in \tilde{\mathcal{H}}$ such that for some $b \in \mathbb{R}$, $\mathrm{sign}(\tilde{h}(\mathbf{y}_i) - b) = \mathcal{C}(f_i)$ for $i = 1, \ldots, m$.

Since $\Pi$ is compact, we define

$$T := \max_{\mathbf{y} \in \Pi} |\tilde{h}(\mathbf{y})|.$$

The final DNN $\tilde{h}_{\mathcal{C}}$ adjusts $\tilde{h}$ such that it vanishes on $\Pi$ but coincides with $\text{sign}(\tilde{h} - b)$ on $\mathcal{Y}$,

$$\tilde{h}_{\mathcal{C}} := \sigma\big(\tilde{h} - (T - b)I - b\big),$$

with $\sigma(x) = \max\{x, 0\}$

Thus, for any $f_i \in F_m$,

$$\text{sign}\left(\sum_{j=1}^{r} \tilde{h}_{\mathcal{C}}(f_i(g_j))\right) = \text{sign}(\tilde{h}(\mathbf{y}_i) - b) = \mathcal{C}(f_i).$$

By Lemma A.6, we can define a GCNN $h_{\mathcal{C}}$ such that $h_{\mathcal{C}}(f) = \sum_{s=1}^{r} \tilde{h}_{\mathcal{C}}\big(f(g_s)\big)$ for any $f \in F_m$. This implies that $\text{sign}(h_{\mathcal{C}}(f_i - b)) = \mathcal{C}(f_i)$ for any $f_i \in F_m$.

As the number of weights in $\tilde{h}_{\mathcal{C}}$ is $W + L + 4m_0 < 6W$, this shows that our GCNN is in the class $\mathcal{H}_{6W, L+1, r}$, completing the proof of the lemma. $\qquad\square$

**Corollary A.8.** *[Corollary 5.3] In the setting of Lemma 5.2, if the number of layers $L > 3$ and $L \leq W^{0.99}$, then there exists a constant $c$ such that*

$$\text{VC}(\mathcal{H}_{W,L,r}) \geq c \cdot \text{VC}(\mathcal{F}_{W,L}).$$

*Proof.* From Equation (2) in Bartlett et al. (2019) , we know that for the class of fully connected neural networks $\mathcal{F}_{W,L}$ with $L$ layers and at most $W$ overall parameters, there exist constants $c_0$ and $C_0$ such that

$$c_0 \cdot WL \log\left(\frac{W}{L}\right) \leq \text{VC}(\mathcal{F}_{W,L}) \leq C_0 \cdot WL \log W. \tag{40}$$

Moreover, by Lemma A.7, we have

$$\text{VC}(\mathcal{H}_{6W', L'+1, r}) \geq \text{VC}(\mathcal{F}_{W', L'}).$$

By choosing $W = 6W'$ and $L = L' + 1$, this shows that

$$\text{VC}(\mathcal{H}_{W,L,r}) \geq \text{VC}(\mathcal{F}_{\lfloor \frac{1}{6} W \rfloor, L-1}).$$

To obtain the statement in the lemma, we combine this bound with the left inequality in (40), leading to

$$\text{VC}(\mathcal{H}_{W,L,r}) \geq \text{VC}(\mathcal{F}_{\lfloor \frac{1}{6} W \rfloor, L-1}) \geq c_0 \cdot \left(\frac{1}{6} W - 1\right)(L-1) \log\left(\frac{\frac{1}{6} W - 1}{L - 1}\right).$$

For some constant $c_1 > 0$, the right-hand side of this inequality is bounded from below by

$$c_1 \cdot WL \log W.$$

By using the right inequality in (40), this can be further bounded,

$$c_1 \cdot WL \log W \geq c \cdot \text{VC}(\mathcal{F}_{W,L}),$$

showing the assertion. $\qquad\square$

Next, we provide the proof for the second part of Theorem 3.5, which states that for some universal constant $c > 0$, the VC dimension $VC(\mathcal{H}_{W,L,r})$ is bounded by $c \cdot W \log_2(r)$. As mentioned in the main article, the first step of the proof is Lemma 5.4.

**Lemma A.9.** *[Lemma 5.4] Let $\mathcal{H}_{4,L,r}$ be the class of GCNNs defined in (34). Then*

$$VC(\mathcal{H}_{4,L,r}) \geq \lfloor \log_2 r \rfloor.$$

*Moreover, for any two numbers $A < B$, there exists a finite subclass of GCNNs $\mathcal{H} \subset \mathcal{H}_{4,L,r}$ that shatters a set of $\lfloor log_2 r \rfloor$ input functions*

$$F_m := \{f_i : G^r \to [A, B] \mid i = 1, \ldots, \lfloor \log_2 r \rfloor\},$$

*and outputs zero for any input function $f : G^r \to \mathbb{R} \setminus [A, B]$.*

*Proof.* To simplify the notation, let $m := \lfloor \log_2 r \rfloor$. It will be enough to show that a subclass of GCNNs $\mathcal{H} \subset \mathcal{H}_{4,L,r}$ shatters the set of $m$ input functions as this immediately implies that

$$VC(\mathcal{H}_{4,L,r}) \geq \lfloor \log_2 r \rfloor.$$

The proof involves selecting $d := 2^m$ distinct points in the interval $[A, B]$ and defining "indicator" neural networks of the form (38) that output 1 at exactly one of these points. By adjusting the parameters of these networks, we can control the intervals of our indicator networks and ensure that each network outputs 1 at the desired point.

Specifically, define $\delta := \frac{B-A}{2(d+2)}$ and select the $d$ points

$$\mathcal{Y} := \{y_i := A + i\delta \mid i = 1, \ldots, d\}.$$

The input functions $F_m$ are now chosen from $\{f : G^r \to \mathcal{Y} \cup \{B - \delta\}\}$.

There are $d = 2^m$ different binary classifiers for the set of $m$ elements. Each binary classifier is defined by the elements for which it outputs 1, and we can index these classifiers by the subsets of $\{1, 2, \ldots, m\}$, denoted by $S_1, \ldots, S_d$. In our construction, each $y_i \in \mathcal{Y}$ corresponds to the binary classifier determined by $S_i$. More formally, the set of $m$ input functions $F_m := \{f_1, \ldots, f_m\}$ is defined by

$$f_j(g_i) := \begin{cases} y_i, & \text{if } j \in S_i, \\ B - \delta, & \text{otherwise.} \end{cases}$$

Next, we define the finite subclass in $\mathcal{H}_{4,L,r}$ that shatters $F_m$ and outputs zero for any function $f : G^r \to \mathbb{R} \setminus [A, B]$.

By the definition of shattering (Definition 3.1 in the main article), for any binary classifier $\mathcal{C} : F_m \to \{-1, 1\}$, we need to find a function in $\text{sign}(\mathcal{H}_{4,L,r})$ matching $\mathcal{C}$ on $F_m$.

For any classifier $\mathcal{C} : F_m \to \{-1, 1\}$ we can find a subset $S \subseteq \{1, \ldots, m\}$ such that $\mathcal{C}(f_j) = 1$ if $j \in S$ and $\mathcal{C}(f_j) = -1$ if $j \in S^c$. There exists an index $i^*$ such that $S = S_{i^*}$. Using Lemma A.6, one can construct a GCNN $h_{i^*} \in \mathcal{H}_{4,L,r}$ that matches $\mathcal{C}$ on $F_m$. Indeed, for any $f_j \in F_m$,

$$h_{i^*}(f_j) := \sum_{s=1}^{r} \mathbf{1}_{(y_{i^*} - \frac{\delta}{2}, y_{i^*} + \frac{\delta}{2}, \frac{\delta}{2})}\big(f_j(g_s)\big) = \begin{cases} 1, & \text{if } j \in S_{i^*}, \\ 0, & \text{otherwise.} \end{cases} \tag{41}$$

Thus, $\text{sign}(h_{i^*}(f) - 0.5) = \mathcal{C}(f)$ for all $f \in F_m$. As an 'indicator' neural network $\mathbf{1}_{(y_{i^*} - \frac{\delta}{2}, y_{i^*} + \frac{\delta}{2}, \frac{\delta}{2})}$ has only 4 parameters and 2 layers, it is in $\mathcal{H}_{4,L,r}$.

Moreover, for any $i = 1, \ldots, d$ and any $x \in \mathbb{R} \setminus [A, B]$, $\mathbf{1}_{(y_i - \frac{\delta}{2}, y_i + \frac{\delta}{2}, \frac{\delta}{2})}(x) = 0$. Arguing as for (41), $h_{i^*}(f) = 0$ for any $f : G^r \to \mathbb{R} \setminus [A, B]$.

That means that the class $\mathcal{H} := \{h_1, \ldots, h_d\}$ shatters input functions $F_m$ and outputs 0 on the subset $\{f : G^r \to \mathbb{R} \setminus [A, B]\}$. This completes the proof. $\square$

**Corollary A.10.** *[Corollary 5.5] The VC dimension of the class $\mathcal{H}_{4W,L,r}$, consisting of GCNNs with $4W$ weights, $L$ layers, and resolution $r$ satisfies the inequality*

$$\text{VC}\big(\mathcal{H}_{4W,L,r}\big) \geq W \lfloor \log_2 r \rfloor.$$

*Proof.* To simplify notation, let $m := \lfloor \log_2 r \rfloor$.

We prove this corollary by defining $W$ disjoint intervals $[A_1, B_1], \ldots, [A_W, B_W]$, where $A_i := (m+3)i$ and $B_i := (m+2)i$. For different $i \in \{1, \ldots, W\}$ the set of input functions $\mathcal{F}_i := \{f : G^r \to [A_i, B_i]\}$ is disjoint since the values of the intervals do not overlap.

By Lemma A.9, for each $i = 1, \ldots, W$, we can find a class of GCNNs $\mathcal{H}_i \subset \mathcal{H}_{4,L,r}$ that shatters a set of $m$ input functions $F_{m,i} \subset \mathcal{F}_i$ and outputs 0 on any other set $F_{m,j}$, where $j \neq i$.

Next we show that the class of GCNNs $\mathcal{H} := \mathcal{H}_1 \oplus \mathcal{H}_2 \oplus \cdots \oplus \mathcal{H}_W \subset \mathcal{H}_{4W,L,r}$ shatters the set $F_{Wm} := \bigsqcup_{i=1}^{W} F_{m,i}$. This will prove the corollary.

By the definition of shattering, we need to find for any binary classifier $\mathcal{C} : F_{Wm} \to \{0, 1\}$, a function in $\text{sign}(\mathcal{H})$ that matches $\mathcal{C}$ on $F_{Wm}$.

For $i = 1, \ldots, W$, let $\mathcal{C}_i := \mathcal{C}\,|_{F_{m,i}}$ be the restriction of $\mathcal{C}$ to $F_{m,i}$. As $\mathcal{H}_i$ shatters $F_{m,i}$, we can choose a GCNN $h_{\mathcal{C}_i} \in \mathcal{H}_i$ such that its values match those of $\mathcal{C}_i$ on $F_{m,i}$.

Next, we show that the values of the GCNN $h_\mathcal{C} := \sum_{i=1}^{W} h_{\mathcal{C}_i}$ match $\mathcal{C}$ on $F_{Wm}$. Let $f$ be any input function from $F_{Wm}$, say $f \in F_q$. For any $i \neq q$, it holds that $h_{\mathcal{C}_i}(f) = 0$ since $h_{\mathcal{C}_i} \in \mathcal{H}_i$. Thus,

$$\sum_{i=1}^{W} h_{\mathcal{C}_i}(f) = h_{\mathcal{C}_q}(f).$$

Since $h_{\mathcal{C}_q}(f) = \mathcal{C}(f)$ by the choice of $h_{\mathcal{C}_q}$, it follows that $h_\mathcal{C}(f) = \mathcal{C}(f)$ for any $f \in F_{Wm}$.

This shows that the class $\mathcal{H}$ of GCNNs shatters $F_{Wm}$, proving the corollary. $\square$

# References

Anthony, M. and Bartlett, P. L. (1999). *Neural network learning: theoretical foundations*. Cambridge University Press, Cambridge.

Bartlett, P., Maiorov, V., and Meir, R. (1998). Almost linear VC dimension bounds for piecewise polynomial networks. *Advances in neural information processing systems*, 11.

Bartlett, P. L., Harvey, N., Liaw, C., and Mehrabian, A. (2019). Nearly-tight VC-dimension and pseudodimension bounds for piecewise linear neural networks. *The Journal of Machine Learning Research*, 20(1):2285–2301.

Bekkers, E. J. (2019). B-spline cnns on lie groups. *arXiv preprint arXiv:1909.12057*.

Bekkers, E. J., Lafarge, M. W., Veta, M., Eppenhof, K. A., Pluim, J. P., and Duits, R. (2018). Roto-translation covariant convolutional networks for medical image analysis. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2018: 21st International Conference, Granada, Spain, September 16-20, 2018, Proceedings, Part I*, pages 440–448. Springer.

Bietti, A., Venturi, L., and Bruna, J. (2021). On the sample complexity of learning under geometric stability. *Advances in Neural Information Processing Systems*, 34:18673–18684.

Brandstetter, J., Welling, M., and Worrall, D. E. (2022). Lie point symmetry data augmentation for neural pde solvers. In *International Conference on Machine Learning*, pages 2241–2256. PMLR.

Cohen, T. and Welling, M. (2016a). Group equivariant convolutional networks. In Balcan, M. F. and Weinberger, K. Q., editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 2990–2999, New York, New York, USA. PMLR.

Cohen, T. and Welling, M. (2016b). Group equivariant convolutional networks. In *International conference on machine learning*, pages 2990–2999. PMLR.

Dehmamy, N., Walters, R., Liu, Y., Wang, D., and Yu, R. (2021). Automatic symmetry discovery with lie algebra convolutional network. *Advances in Neural Information Processing Systems*, 34:2503–2515.

Elesedy, B. (2022). Group symmetry in PAC learning. In *ICLR 2022 Workshop on Geometrical and Topological Representation Learning*.

Fuchs, F., Worrall, D., Fischer, V., and Welling, M. (2020). Se (3)-transformers: 3d roto-translation equivariant attention networks. *Advances in neural information processing systems*, 33:1970–1981.

Hinton, Geoffrey, A. K. and Wang, S. (2011). Transforming auto-encoders. In *ICANN-11: International Conference on Artificial Neural Networks*.

Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool, K., Bates, R., vZídek, A., Potapenko, A., et al. (2021). Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589.

Keriven, N. and Peyré, G. (2019). Universal invariant and equivariant graph neural networks. *Advances in Neural Information Processing Systems*, 32.

Kohler, M. and Walter, B. (2023). Analysis of convolutional neural network image classifiers in a rotational symmetric model. *IEEE transactions on pattern analysis and machine intelligence*, 69(8):5203–5218.

Kondor, R. and Trivedi, S. (2018). On the generalization of equivariance and convolution in neural networks to the action of compact groups. In *International Conference on Machine Learning*, pages 2747–2755. PMLR.

Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25.

LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324.

Lee, C., Xie, S., Gallagher, P., Zhang, Z., and Tu, Z. (2015). Deeply-supervised nets. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 38, pages 562–570.

Long, J., Shelhamer, E., and Darrell, T. (2015). Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440.

Petersen, P. C. and Sepliarskaia, A. (2024). VC dimensions of group convolutional neural networks. *Neural Networks*, 169:462–474.

Pollard, D. (1990). *Empirical processes: theory and applications*, volume 2 of *NSF-CBMS Regional Conference Series in Probability and Statistics*. Institute of Mathematical Statistics, Hayward, CA; American Statistical Association, Alexandria, VA.

Procesi, C. (2007). *Lie groups: an approach through invariants and representations*, volume 115. Springer.

Ren, S., He, K., Girshick, R., and Sun, J. (2016). Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE transactions on pattern analysis and machine intelligence*, 39(6):1137–1149.

Sannai, A., Imaizumi, M., and Kawano, M. (2021). Improved generalization bounds of group invariant/equivariant deep networks via quotient feature spaces. In *Uncertainty in Artificial Intelligence*, pages 771–780. PMLR.

Shao, H., Montasser, O., and Blum, A. (2022). A theory of pac learnability under transformation invariances. *Advances in Neural Information Processing Systems*, 35:13989–14001.

Smets, B. M. N., Portegies, J., Bekkers, E. J., and Duits, R. (2023). PDE-based group equivariant convolutional neural networks. *J. Math. Imaging Vision*, 65(1):209–239.

Sosnovik, I., Moskalev, A., and Smeulders, A. W. (2021). Scale equivariance improves siamese tracking. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2765–2774.

Stroppel, M. (2006). *Locally compact groups*, volume 3. European Mathematical Society.

Vinuesa, R. and Brunton, S. L. (2022). Enhancing computational fluid dynamics with machine learning. *Nature Computational Science*, 2(6):358–366.

Wang, R., Walters, R., and Yu, R. (2020). Incorporating symmetry into deep dynamics models for improved generalization. *arXiv preprint arXiv:2002.03061*.

Weiler, M., Forré, P., Verlinde, E., and Welling, M. (2021). Coordinate independent convolutional networks– isometry and gauge equivariant convolutions on riemannian manifolds. *arXiv preprint arXiv:2106.06020*.

Weiler, M., Hamprecht, F. A., and Storath, M. (2018). Learning steerable filters for rotation equivariant cnns. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 849–858.

Zhdanov, M., Ruhe, D., Weiler, M., Lucic, A., Brandstetter, J., and Forré, P. (2024). Clifford-steerable convolutional neural networks. *arXiv preprint arXiv:2402.14730*.