

A PROOFS

We begin with a simple lemma detailing the Fourier space end-to-end predictor for a G-CNN.

Lemma A.1 (G-CNN in Fourier space). *A G-CNN given by $\langle \mathbf{x} \star \mathbf{w}_1 \star \cdots \star \mathbf{w}_{L-1}, \mathbf{w}_L \rangle$ is equivalent to $\langle \mathbf{x}, \mathbf{w}_L \star \mathbf{w}_{L-1}^- \star \cdots \star \mathbf{w}_1^- \rangle$, or in Fourier space to $\frac{1}{|G|} \langle \hat{\mathbf{x}}, \hat{\mathbf{w}}_L \cdots \hat{\mathbf{w}}_1 \rangle_M$. Moreover,*

Proof.

$$\begin{aligned}
 N(x) &= \langle \mathbf{x} \star \mathbf{w}_1 \star \cdots \star \mathbf{w}_{L-1}, \mathbf{w}_L \rangle \\
 &= \frac{1}{|G|} \langle \mathcal{F}_M(\mathbf{x} \star \mathbf{w}_1 \star \cdots \star \mathbf{w}_{L-1}), \mathcal{F}_M \mathbf{w}_L \rangle_M \\
 &= \frac{1}{|G|} \text{tr}[\mathcal{F}_M(\mathbf{x} \star \mathbf{w}_1 \star \cdots \star \mathbf{w}_{L-1}) \hat{\mathbf{w}}_L^\dagger] \\
 &= \frac{1}{|G|} \text{tr}[\mathcal{F}_M(\mathbf{x} \star \mathbf{w}_1 \star \cdots \star \mathbf{w}_{L-2}) \hat{\mathbf{w}}_{L-1}^\dagger \hat{\mathbf{w}}_L^\dagger] \\
 &= \frac{1}{|G|} \text{tr}[\hat{\mathbf{x}} \hat{\mathbf{w}}_1^\dagger \cdots \hat{\mathbf{w}}_{L-1}^\dagger \hat{\mathbf{w}}_L^\dagger] \\
 &= \frac{1}{|G|} \langle \hat{\mathbf{x}}, \hat{\mathbf{w}}_L \cdots \hat{\mathbf{w}}_1 \rangle_M
 \end{aligned} \tag{A.1}$$

To see the stated equality in real space, observe that $\widehat{\mathbf{w}_1^-} = \widehat{\mathbf{w}}_1^\dagger$ by definition. Thus,

$$\mathcal{F}_M(\mathbf{w}_L \star \mathbf{w}_{L-1}^- \star \cdots \star \mathbf{w}_1^-) = \hat{\mathbf{w}}_L \cdots \hat{\mathbf{w}}_1$$

□

A.1 ABELIAN

The proof of Proposition 5.2 is included below.

First, recall the primary theorem of Yun et al. (2020):

Theorem A.2 (paraphrased from Yun et al. (2020)). *If there exists $\lambda > 0$ such that the initial directions $\bar{\mathbf{v}}_1, \dots, \bar{\mathbf{v}}_L$ of the network parameters satisfy $|\mathcal{F} \bar{\mathbf{v}}_\ell|_j^2 - |\mathcal{F} \bar{\mathbf{v}}_L|_j^2 \geq \lambda$ for all $\ell \in [L-1]$ and $j \in [m]$, i.e. of the Fourier transform magnitudes of the initial directions look sufficiently different pointwise (which is likely for e.g. a random initialization), then $\beta(\Theta(t))$ converges in a direction that aligns with $\mathbf{S}^T \boldsymbol{\rho}_\infty$ where $\boldsymbol{\rho}_\infty \in \mathbb{C}^m$ denotes a stationary point of the following optimization program:*

$$\min_{\boldsymbol{\rho} \in \mathbb{C}^m} \|\boldsymbol{\rho}\|_{2/L} \quad \text{s.t.} \quad y_i \mathbf{x}_i^T \mathcal{F}^T \boldsymbol{\rho} \geq 1, \forall i \in [n] \tag{A.2}$$

Since $\mathbf{S} = \mathcal{F}$ is invertible, then in fact $\beta(\Theta(t))$ converges in a direction that aligns with a stationary point \mathbf{z}_∞ of the following optimization program:

$$\min_{\mathbf{z} \in \mathbb{C}^m} \|\mathcal{F} \mathbf{z}\|_{2/L} \quad \text{s.t.} \quad y_i \mathbf{x}_i^T \mathbf{z} \geq 1, \forall i \in [n] \tag{A.3}$$

We proceed by showing that abelian G-CNNs can be written as a vector of parameters contracted (according to tensor operations) with an orthogonally decomposable data tensor, which is the primary condition for Theorem A.2 to hold.

For convenience, we restate Proposition A.3 before detailing the proof.

Proposition A.3 (paraphrased from Yun et al. (2020)). *Let $\mathbf{M}(\mathbf{x})$ be a function that maps data $\mathbf{x} \in \mathbb{R}^d$ to a data tensor $\mathbf{M}(\mathbf{x}) \in \mathbb{R}^{k_1 \times k_2 \times \cdots \times k_L}$. The data input into an L -layer tensorized neural network can be written in the form of an **orthogonally decomposable data tensor** if there exists a full column rank matrix $\mathbf{S} \in \mathbb{C}^{m \times d}$ and semi-unitary matrices $\mathbf{U}_1, \dots, \mathbf{U}_L \in \mathbb{C}^{k_\ell \times m}$ where $d \leq m \leq \min_\ell k_\ell$ such that $\mathbf{M}(\mathbf{x})$ can be written as:*

$$\mathbf{M}(\mathbf{x}) = \sum_{j=1}^m [\mathbf{S} \mathbf{x}]_j ([\mathbf{U}_1]_{\cdot, j} \otimes [\mathbf{U}_2]_{\cdot, j} \otimes \cdots \otimes [\mathbf{U}_L]_{\cdot, j}) \tag{A.4}$$

and such that the network output is the tensor multiplication between $\mathbf{M}(x)$ and each layer's parameters:

$$\begin{aligned} \text{NN}(x; \Theta) &= \mathbf{M}(x) \cdot (\mathbf{W}_1, \dots, \mathbf{W}_L) \\ &= \sum_{i_1=1}^d \cdots \sum_{i_L=1}^d \mathbf{M}(x)_{i_1 \dots i_L} (\mathbf{W}_1)_{i_1} \cdots (\mathbf{W}_L)_{i_L} \end{aligned}$$

For Proposition A.3. By direct manipulation:

$$\begin{aligned} f(x; \Theta) &= \mathbf{M}(x) \cdot (v_1, \dots, v_L) \\ &= \sum_{i_1=1}^d \cdots \sum_{i_L=1}^d \mathbf{M}(x)_{i_1 \dots i_L} (v_1)_{i_1} \cdots (v_L)_{i_L} \\ &= \sum_{i_1=1}^d \cdots \sum_{i_L=1}^d \left(\sum_{j=1}^d [\mathbf{S}x]_j ([U_1]_{\cdot, j} \otimes [U_2]_{\cdot, j} \otimes \cdots \otimes [U_L]_{\cdot, j}) \right)_{i_1 \dots i_L} (v_1)_{i_1} \cdots (v_L)_{i_L} \\ &= \sum_{i_1=1}^d \cdots \sum_{i_L=1}^d \left(\sum_{j=1}^d [\mathbf{S}x]_j ([U_1]_{i_1, j} [U_2]_{i_2, j} \cdots [U_L]_{i_L, j}) \right) (v_1)_{i_1} \cdots (v_L)_{i_L} \\ &= \sum_{j=1}^d [\mathbf{S}x]_j ([U_1^T v_1]_{i_1} [U_2^T v_2]_{i_2} \cdots [U_L^T v_L]_{i_L}) \\ &= d^{\frac{L-1}{2}} \sum_{j=1}^d [\mathcal{F}x]_j ([\overline{\mathcal{F}v_1}]_{i_1} [\overline{\mathcal{F}v_2}]_{i_2} \cdots [\overline{\mathcal{F}v_L}]_{i_L}) \\ &= d^{\frac{L-1}{2}} \sum_{j=1}^d [\mathcal{F}x]_j ([\overline{\mathcal{F}v_1}]_{i_1} [\overline{\mathcal{F}v_2}]_{i_2} \cdots [\overline{\mathcal{F}v_L}]_{i_L}) \\ &= d^{\frac{L-1}{2}} \sum_{j=1}^d [\mathcal{F}x]_j (\overline{[\mathcal{F}v_1]_{i_1} [\mathcal{F}v_2]_{i_2} \cdots [\mathcal{F}v_L]_{i_L}}) \\ &= \langle \mathcal{F}x, \mathcal{F}v_1 \odot \cdots \odot \mathcal{F}v_L \rangle \\ &= \langle \mathcal{F}x \odot \overline{\mathcal{F}v_1}, \mathcal{F}v_2 \odot \cdots \odot \mathcal{F}v_L \rangle \\ &= \langle \mathcal{F}(x \star v_1) \odot \overline{\mathcal{F}v_2}, \mathcal{F}v_3 \odot \cdots \odot \mathcal{F}v_L \rangle \\ &= \langle \mathcal{F}(x \star v_1 \star v_2) \odot \overline{\mathcal{F}v_3}, \mathcal{F}v_4 \odot \cdots \odot \mathcal{F}v_L \rangle \\ &= \langle \mathcal{F}(x \star v_1 \star v_2 \star \cdots \star v_{L-1}), \mathcal{F}v_L \rangle \\ &= \langle x \star v_1 \star \cdots \star v_{L-1}, v_L \rangle \end{aligned}$$

Here, we have used that the filters are real-valued. \square

Note that Theorem 5.3 then merely requires that $\|\mathcal{F}^{-T}z\| = \|\overline{\mathcal{F}z}\| = \|\mathcal{F}z\|$, for real-valued z .

A.1.1 FUNDAMENTAL THEOREM OF FINITE ABELIAN GROUPS

While the proof in the previous section is complete and correct, intuition (and/or alternate analysis) for abelian groups is aided by the important fact that all finite abelian groups are a direct product of cyclic groups.

Theorem A.4 (Fundamental theorem of finite abelian groups (Dummit & Foote, 2004)). *Any finite abelian group is a direct product of a finite number of cyclic groups whose orders are prime powers uniquely determined by the group.*

Given a decomposition of an abelian group G into k cyclic groups $C_{d_1} \times \cdots \times C_{d_k}$, one can easily construct the group Fourier transform as a Kronecker product of discrete Fourier transform matrices

which are the group Fourier transforms of the respective cyclic groups.

$$G = C_{d_1} \times \cdots \times C_{d_k} \implies \mathcal{F} = \bigotimes_{i=1}^k \mathcal{F}_{d_i}, \quad (\text{A.5})$$

where \bigotimes denotes the Kronecker product over matrices and \mathcal{F}_d is the standard (unitary) discrete Fourier transform matrix of dimension d defined as

$$\mathcal{F}_d = \frac{1}{\sqrt{d}} \begin{bmatrix} \omega_d^{0 \cdot 0} & \omega_d^{0 \cdot 1} & \cdots & \omega_d^{0 \cdot (d-1)} \\ \omega_d^{1 \cdot 0} & \omega_d^{1 \cdot 1} & \cdots & \omega_d^{1 \cdot (d-1)} \\ \vdots & \vdots & \ddots & \vdots \\ \omega_d^{(d-1) \cdot 0} & \omega_d^{(d-1) \cdot 1} & \cdots & \omega_d^{(d-1) \cdot (d-1)} \end{bmatrix}, \quad \omega_d = e^{\frac{-2\pi i}{d}}. \quad (\text{A.6})$$

From this result, it is clear that the desired properties of the Fourier transform and convolution hold.

A.2 NON-ABELIAN

Theorem A.5. *Consider a classification task with ground-truth linear predictor β , trained via a linear G-CNN architecture $\text{NN}(\mathbf{x}) = \langle \mathbf{x} \star \mathbf{w}_1 \star \cdots \star \mathbf{w}_{L-1}, \mathbf{w}_L \rangle$ (see Section 5 for architecture details) with $L \geq 2$ layers under the exponential loss. Then for almost any datasets $\{\mathbf{x}_i, y_i\}$ separable by β , any bounded sequence of step sizes η_t , and almost all initializations, suppose that:*

- The loss $\mathcal{L}(\mathbf{W})$ converges to 0
- The gradients with respect to the end-to-end linear predictor, $\nabla_{\beta} \mathcal{L}(\mathcal{P}(\mathbf{W}^t))$, converge in direction as $t \rightarrow \infty$
- The iterates \mathbf{W}^t themselves converge in direction as $t \rightarrow \infty$ to a separator $\mathcal{P}(\mathbf{W}^t)$ with positive margin

When $L = 2$, we need an additional technical assumption, Assumption A.7. Then, the resultant linear predictor $\hat{\beta}^\infty$ is a positive scaling of a first order stationary point of the optimization problem:

$$\min_{\beta} \|\hat{\beta}\|_{2/L}^{(S)} \quad \text{s.t.} \quad \forall n, y_n \langle \beta_n, \mathbf{x}_n \rangle_M \geq 1 \quad (\text{A.7})$$

In this section, we prove the non-abelian case, Theorem A.5. The proof of our result proceeds according to the following outline:

1. By applying a general result of Gunasekar et al. (2018b), Theorem A.6, we characterize the implicit regularization in the full space of parameters, \mathbf{W} (in contrast to the end-to-end linear predictor β), as the stationary point of an optimization problem Equation A.9 in \mathbf{W} .
2. Separately, we define a *distinct* optimization problem, Equation A.7 in β . The goal is to demonstrate that stationary points of Equation A.9 are a subset of the stationary points of Equation A.7.
3. The *necessary* KKT conditions for Equation A.9 characterize its stationary points. Using this characterization, we show that the *sufficient* KKT optimality conditions for Equation A.7 are in fact also satisfied for the corresponding end-to-end predictor. Thus, we show that for any stationary point \mathbf{W}^\dagger of Equation A.9, the linear predictor $\mathcal{P}(\mathbf{W}^\dagger)$ is a stationary point of Equation A.7.

First, recall that Gunasekar et al. (2018b) prove the following general result about the implicit regularization of any homogeneous polynomial parametrization:

Theorem A.6 (Homogeneous polynomial parametrization, Theorem 4 of Gunasekar et al. (2018b)). *Let \mathbf{W} be the concatenation of all (real-valued) parameters \mathbf{W}_i . For any homogeneous polynomial map $\mathcal{P} : \mathbb{R}^P \rightarrow \mathbb{R}^{|G|}$ from parameters $\mathbf{W}_i \in \mathbb{R}^P$ to linear predictors, almost all datasets*

$\{\mathbf{x}_n, y_n\}_{n=1}^N$ separable by the ground truth predictor $\beta := \{\mathcal{P}(\mathbf{W}) : \mathbf{W} \in \mathbb{R}^P\}$, almost all initializations \mathbf{W}^0 , and any bounded sequence of step sizes $\{\eta_t\}_t$, consider the gradient descent updates:

$$\mathbf{W}^{t+1} = \mathbf{W}^t - \eta_t \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}^t) = \mathbf{W}^t - \eta_t \nabla_{\mathbf{W}} \mathcal{P}(\mathbf{W}^t) \nabla_{\beta} \mathcal{L}(\mathcal{P}(\mathbf{W}^t)) \quad (\text{A.8})$$

Suppose furthermore that the exponential loss converges to zero, that the gradients $\nabla_{\beta} \mathcal{L}(\beta^t)$ converge in direction, and that the iterates \mathbf{W}^t themselves converge in direction to yield a separator with positive margin. Then, the limit direction of the parameters $\bar{\mathbf{W}}^{\infty} = \lim_{t \rightarrow \infty} \frac{\mathbf{W}^{(t)}}{\|\mathbf{W}^{(t)}\|_2}$ is a positive scaling of a first order stationary point of the following optimization problem:

$$\min_{\mathbf{W} \in \mathbb{R}^P} \|\mathbf{W}\|_2^2 \quad \text{s.t.} \quad \forall n, y_n \langle \mathbf{x}_n, \mathcal{P}(\mathbf{W}) \rangle \geq 1. \quad (\text{A.9})$$

To keep track of constant factors, let $\mathbf{W}^{\infty} = \tau \bar{\mathbf{W}}^{\infty}$ denote the first order stationary point itself. Furthermore, let \mathbf{W}_i^{∞} denote the individual layers (or parameter blocks) comprising \mathbf{W}^{∞} , and similarly let $\bar{\mathbf{W}}_i^{\infty}$ denote the individual layers comprising $\bar{\mathbf{W}}^{\infty}$. We then have via the KKT conditions that:

$$\begin{aligned} & \exists \{\alpha_n : \alpha_n \geq 0\}_{n=1}^N \text{ s.t. } \alpha_n = 0 \text{ if } y_n \langle \mathbf{x}_n, \mathcal{P}(\mathbf{W}^{\infty}) \rangle > 1 \\ & \mathbf{W}_i^{\infty} = \nabla_{\mathbf{W}_i} \mathcal{P}(\mathbf{W}^{\infty}) \left[\sum_n \alpha_n y_n \mathbf{x}_n \right] = \nabla_{\mathbf{W}_i} \left\langle \mathcal{P}(\mathbf{W}^{\infty}), \sum_n \alpha_n y_n \mathbf{x}_n \right\rangle \end{aligned} \quad (\text{A.10})$$

While this is an interesting result alone, the goal of implicit regularization is to characterize the final linear predictor (which is some function \mathcal{P} of the complete parametrization \mathbf{W}). To that end, consider the following optimization problem in β :

$$\min_{\beta} \|\hat{\beta}\|_{2/L}^{(S)} \quad \text{s.t.} \quad \forall n, y_n \langle \beta_n, \mathbf{x}_n \rangle \geq 1 \quad (\text{A.11})$$

We will leverage the *necessary* KKT conditions from Equation A.10 to show that first-order stationary points of Equation A.9 are (up to a scaling) also first-order stationary points of Equation A.7, using the *sufficient* KKT conditions for Equation A.7.

Using standard KKT sufficiency conditions, the first-order stationary points of Equation A.7 are those vectors β such that there exist $\tilde{\alpha}_1, \dots, \tilde{\alpha}_n$ satisfying:

1. **Feasibility:** $\forall n, y_n \langle \beta, \mathbf{x}_n \rangle \geq 1$ and $\tilde{\alpha}_i \geq 0 \quad \forall i$
2. **Complementary slackness:** $\forall i, \tilde{\alpha}_i = 0$ if $y_n \langle \beta, \mathbf{x}_n \rangle > 1$
3. **Membership in subdifferential:** $\sum_n \tilde{\alpha}_n y_n \hat{\mathbf{x}}_n \in \partial^o \|\hat{\beta}\|_{2/L}^{(S)}$

In the third condition above, ∂^o is the local sub-differential of Clarke (1975): $\partial^o f(\beta) = \text{conv}\{\lim_{i \rightarrow \infty} \nabla f(\beta + \mathbf{h}_i) : \mathbf{h}_i \rightarrow 0\}^8$.

We will need the following assumption in the special case $L = 2$:

Assumption A.7 ($L = 2$ bounded subgradient). Let $\hat{\mathbf{z}} = \sum_n \tilde{\alpha}_n y_n \hat{\mathbf{x}}_n$ result from the KKT conditions of the optimization problem in \mathbf{W} , Equation A.10, as described previously. Then, we assume that $\|\hat{\mathbf{z}}\|_{\infty}^{(S)} \leq 1$.

Let $\beta^{\infty} = \mathcal{P}(\mathbf{W}^{\infty})$ and let $\tilde{\alpha}_i = \frac{1}{\gamma} \alpha_i$ for all i , where γ is equal to $\left(\|\hat{\beta}^{\infty}\|_{\frac{2}{L}}\right)$ for $L > 2$ and to 1 otherwise. (Note that by homogeneity of \mathcal{P} , $\mathcal{P}(\mathbf{W}^{\infty}) = \mathcal{P}(\tau \bar{\mathbf{W}}^{\infty}) = \tau^L \mathcal{P}(\bar{\mathbf{W}}^{\infty})$.) We will check these conditions one by one for β^{∞} and $\tilde{\alpha}_i$, with the first two following immediately from Theorem A.6 and the last one requiring the most manipulation.

Feasibility Trivially, $\tilde{\alpha}_i \geq 0 \quad \forall i$ by definition of α in Equation A.10. Similarly, $y_n \langle \mathbf{x}_n, \beta^{\infty} \rangle = y_n \langle \mathbf{x}_n, \mathcal{P}(\mathbf{W}^{\infty}) \rangle \geq 1$.

⁸ \mathbf{h}_i is a sequence of vectors in some linear space, and we take $\mathbf{h}_i \rightarrow 0$ as an entry-wise statement. This is because the vectors are finite-dimensional, so all norms are equivalent.

Complementary slackness If $y_n \langle \beta^\infty, \mathbf{x}_n \rangle > 1 = y_n \langle \mathcal{P}(\mathbf{W}^\infty), \mathbf{x}_n \rangle > 1$, then $\tilde{\alpha}_n \propto \alpha_n = 0$.

Membership in subdifferential We first characterize the set $\partial^o \|\hat{\beta}\|_{2/L}^{(S)}$ for a generic matrix β , and then show that $\hat{z} \triangleq \sum_n \tilde{\alpha}_n y_n \hat{\mathbf{x}}_n \in \partial^o \|\hat{\beta}^\infty\|_{2/L}^{(S)}$. When $L = 2$ and thus $p = 1$, the Schatten norm $\|\hat{\beta}^\infty\|_1$ is indeed a norm and its subgradient is known; see e.g. Watson (1992). We restate this result below:

Lemma A.8 (Subdifferential of p -Schatten norm, $p = 1$). *Suppose $L = 2$, such that $\frac{2}{L} = p = 1$. Let \mathbf{A} be an $n \times n$ complex-valued matrix. Then we have*

$$\partial^o \|\mathbf{A}\|_1^{(S)} = \left\{ \mathbf{G} : \|\mathbf{A}\|_1^{(S)} = \text{tr}[\mathbf{G}^\dagger \mathbf{A}], \|\mathbf{G}\|_\infty^{(S)} \leq 1 \right\} \quad (\text{A.12})$$

Lemma A.9 (Subdifferential of p -Schatten norm, $p < 1$). *Suppose $L > 2$ and let $\frac{2}{L} = p$, such that $0 < p < 1$. Let \mathbf{A} be an $n \times n$ complex-valued matrix with singular value decomposition $\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^\dagger$. Let Π_U project onto the column space of \mathbf{A} , i.e. $\Pi_U(\mathbf{M}) = \mathbf{U} \mathbf{U}^\dagger \mathbf{M}$. Let Π_V project onto the row space of \mathbf{A} , i.e. $\Pi_V(\mathbf{M}) = \mathbf{M} \mathbf{V} \mathbf{V}^\dagger$. Then we have*

$$\partial^o \|\mathbf{A}\|_p^{(S)} = \left\{ \mathbf{G} : \mathbf{A}^\dagger \mathbf{G} = \frac{1}{\|\mathbf{A}\|_p^{(S)}} \sqrt{\mathbf{A}^\dagger \mathbf{A}}^p \text{ and } \mathbf{G} \mathbf{A}^\dagger = \frac{1}{\|\mathbf{A}\|_p^{(S)}} \sqrt{\mathbf{A} \mathbf{A}^\dagger}^p \right\} \quad (\text{A.13})$$

Proof. Suppose \mathbf{A} is rank r and has singular value decomposition $\mathbf{A} = \sum_{i=1}^r d_i \mathbf{u}_i \mathbf{v}_i^\dagger$, where $d_i > 0$ for all i . Consider unit vectors $\mathbf{W}_1, \dots, \mathbf{W}_{n-r}$ which are a basis for the orthogonal subspace to $\text{Span}(\mathbf{u}_1, \dots, \mathbf{u}_r)$, and unit vectors $\mathbf{s}_1, \dots, \mathbf{s}_{n-r}$ which are a basis for the orthogonal subspace to $\text{Span}(\mathbf{v}_1, \dots, \mathbf{v}_r)$. Treating the space of $n \times n$ complex matrices as a n^2 -dimensional linear space, we see that $\{\mathbf{u}_i \mathbf{v}_i^\dagger\}_{i=1}^r$ form a basis for an r -dimensional subspace. Let Π_A denote the projector onto this space, $\Pi_A = \sum_{i=1}^r \mathbf{u}_i \mathbf{v}_i^\dagger$. Note that $\Pi_A \mathbf{A} = \mathbf{A}$. For any set of $n-r$ sequences $\{\epsilon_{m,i}\}_{m=1}^\infty$ such that $\lim_{m \rightarrow \infty} \epsilon_{m,i} = 0$ for all $i = 1, \dots, n-r$, consider the particular sequence of matrices $\{\mathbf{H}_m\}_{m=1}^\infty$ defined by

$$\mathbf{H}_m = \sum_{i=1}^{n-r} \epsilon_{m,i} \mathbf{W}_i \mathbf{s}_i^\dagger$$

By definition, $\lim_{m \rightarrow \infty} \|\mathbf{H}_m\|_{\text{Fro}} = 0$, where $\|\mathbf{H}_m\|_{\text{Fro}} \triangleq \|\mathbf{H}_m\|_2^{(S)}$. Also, if \mathbf{M} is a full rank matrix with singular value decomposition $\mathbf{A} \mathbf{D} \mathbf{B}^\dagger$, $\|\mathbf{M}\|_p^{(S)}$ is differentiable at \mathbf{M} and $\nabla \|\mathbf{M}\|_p^{(S)} = \frac{1}{\|\mathbf{M}\|_p^{(S)}} \mathbf{A} \mathbf{D}^{p-1} \mathbf{B}^\dagger$. For convenience of notation, let \mathbf{U} be the matrix with i^{th} column \mathbf{u}_i , \mathbf{W} the matrix with i^{th} column \mathbf{W}_i , and similarly for \mathbf{V} and \mathbf{S} with respect to vectors \mathbf{v}_i and \mathbf{s}_i respectively. Also, let \mathbf{D} be the diagonal matrix with d_i on the i^{th} diagonal.

Combining this fact with the construction of \mathbf{H}_m , we have that

$$\nabla \|\mathbf{A} + \mathbf{H}_m\|_p^{(S)} = \nabla \left\| \sum_{i=1}^r d_i \mathbf{u}_i \mathbf{v}_i^\dagger + \sum_{i=1}^{n-r} \epsilon_{m,i} \mathbf{W}_i \mathbf{s}_i^\dagger \right\|_p^{(S)} \quad (\text{A.14})$$

$$= \nabla \left\| \begin{bmatrix} \mathbf{D} & 0 & \dots & 0 \\ 0 & \epsilon_{m,1} & \dots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \dots & \epsilon_{m,n-r} \end{bmatrix} \begin{bmatrix} \mathbf{V}^\dagger \\ \mathbf{S}^\dagger \end{bmatrix} \right\|_p^{(S)} \quad (\text{A.15})$$

$$= \frac{1}{(\sum_{i=1}^r d_i^p + \sum_{i=1}^{n-r} \epsilon_{m,i}^p)^{\frac{1}{p}}} \left(\mathbf{U} \mathbf{D}^{p-1} \mathbf{V}^\dagger + \mathbf{W} \begin{bmatrix} \epsilon_{m,1}^{p-1} & \dots & 0 \\ \vdots & & \ddots \\ 0 & \dots & \epsilon_{m,n-r}^{p-1} \end{bmatrix} \mathbf{S}^\dagger \right) \quad (\text{A.16})$$

In the limit as m goes to infinity, $\sum_{i=1}^{n-r} \epsilon_{m,i}^p$ approaches 0. However, $p - 1 < 0$ implies that $\lim_{m \rightarrow \infty} \epsilon_{m,i}^{p-1} = \pm\infty$. By taking convex combinations, one can create any matrix with left and right singular vectors \mathbf{W} and \mathbf{S}^\dagger . Formally, we have:

$$\partial^o \|A\|_p^{(S)} = \text{conv}\left\{\lim_{m \rightarrow \infty} \nabla \|A + H_m\|_p : H_m \rightarrow 0\right\} \quad (\text{A.17})$$

$$= \frac{1}{\|A\|_p^{(S)}} U D^{p-1} V^\dagger + \frac{1}{\|A\|_p^{(S)}} \text{conv}\left\{\lim_{m \rightarrow \infty} W \begin{bmatrix} \epsilon_{m,1}^{p-1} & \dots & 0 \\ \vdots & & \ddots \\ 0 & \dots & \epsilon_{m,n-r}^{p-1} \end{bmatrix} S^\dagger : \epsilon_{m,i} \rightarrow_m 0\right\} \quad (\text{A.18})$$

$$= \frac{1}{\|A\|_p^{(S)}} U D^{p-1} V^\dagger + \{W \Sigma S^\dagger : \Sigma \text{ is real and diagonal}\} \quad (\text{A.19})$$

Note that for any rank-one matrix $M = ab^\dagger$, if a is not orthogonal to each column of U , then $UU^\dagger M \neq 0$. Similarly, if b is not orthogonal to each column of V , then $MVV^\dagger \neq 0$. Thus for an arbitrary matrix M , by decomposing it into a sum of rank-one matrices via its SVD, we see that $\Pi_U M = \Pi_V M = 0$ implies that both the row and column spaces of M are orthogonal to those of A , respectively. Thus, we can project via Π_U and Π_V to disregard the second term of Equation A.19, and obtain the following expression for the subgradient:

$$\partial^o \|A\|_p^{(S)} = \left\{ G : \Pi_U G = \Pi_V G = \left(\|A\|_p^{(S)}\right)^{-p} U D^{p-1} V^\dagger \right\} \quad (\text{A.20})$$

Consider first the equality

$$\Pi_U G = \left(\|A\|_p^{(S)}\right)^{-p} U D^{p-1} V^\dagger \quad (\text{A.21})$$

We can left-multiply by A^\dagger without changing the set of matrices G satisfying this relation. To see this, one can check that if $A^\dagger \Pi_U G = \left(\|A\|_p^{(S)}\right)^{-p} A^\dagger U D^{p-1} V^\dagger$, left-multiplying on both sides by $U D^{-1} V^\dagger$ recovers Equation A.21.

Similarly, consider the second equality:

$$\Pi_V G = \left(\|A\|_p^{(S)}\right)^{-p} U D^{p-1} V^\dagger \quad (\text{A.22})$$

We can right-multiply by A^\dagger without changing the set of matrices G satisfying this relation. To see this, one can check that if

$$(\Pi_V G) A^\dagger = \left(\|A\|_p^{(S)}\right)^{-p} U D^{p-1} V^\dagger A^\dagger$$

Then right-multiplying on both sides by $U D^{-1} V^\dagger$ recovers Equation A.22. Finally, observe that $A^\dagger \Pi_U G = V D U^\dagger U U^\dagger G = A^\dagger G$ and $(\Pi_V G) A^\dagger = G V V^\dagger V D U^\dagger = G V D U^\dagger = G A^\dagger$. Furthermore, $A^\dagger \frac{1}{\|A\|_p^{(S)}} U D^{p-1} V^\dagger = V D^p V^\dagger = \frac{1}{\|A\|_p^{(S)}} \sqrt{A^\dagger A}^p$ and $\frac{1}{\|A\|_p^{(S)}} U D^{p-1} V^\dagger A^\dagger = \frac{1}{\|A\|_p^{(S)}} \sqrt{A A^\dagger}^p$. This completes the proof of the lemma. \square

Lemmas A.8 and A.9 characterize the subdifferential of the Schatten norm. Now, we show that $\sum_n \tilde{\alpha}_n y_n \hat{x}_n$ satisfies Equation A.13.

Lemma A.10. Recall that our G-CNN is given by $\text{NN}(x) = \langle x \star w_1 \star \dots \star w_{L-1}, w_L \rangle$, with the vector of parameters $W = [w_1 \dots w_L]$ and end-to-end linear predictor given by $\mathcal{P}(W) = w_1 \star \dots \star w_L$. Consider an arbitrary such vector of real-valued parameters. Also, we have assumed that the filters in real-space are real-valued, i.e. $w_i \in \mathbb{R}^{|G|}$. Then the following relation holds:

$$\mathcal{F}_M \nabla_{w_\ell} \langle \mathcal{P}(W), e_i \rangle = \hat{w}_{\ell+1}^\dagger \dots \hat{w}_L^\dagger \hat{e}_i \hat{w}_1^\dagger \dots \hat{w}_{\ell-1}^\dagger \quad (\text{A.23})$$

where e_i is the i^{th} standard basis vector.

Proof. Since $\langle \mathcal{P}(\mathbf{W}), \mathbf{e}_i \rangle$ is real, we have $\langle \mathcal{P}(\mathbf{W}), \mathbf{e}_i \rangle = \langle \mathbf{e}_i, \mathcal{P}(\mathbf{W}) \rangle$. Plugging this in,

$$\begin{aligned}
\mathcal{F}_M \nabla_{\mathbf{w}_\ell} \mathcal{P}(\mathbf{W})[\mathbf{e}_i] &= \mathcal{F}_M \nabla_{\mathbf{w}_\ell} \langle \mathcal{P}(\mathbf{W}), \mathbf{e}_i \rangle \\
&= \mathcal{F}_M \nabla_{\mathbf{w}_\ell} \langle \mathbf{e}_i, \mathcal{P}(\mathbf{W}) \rangle \\
&= \mathcal{F}_M \nabla_{\mathbf{w}_\ell} \frac{1}{|G|} \langle \mathcal{F}_M \mathbf{e}_i, \mathcal{F}_M \mathcal{P}(\mathbf{W}) \rangle_M \\
&= \mathcal{F}_M \nabla_{\mathbf{w}_\ell} \frac{1}{|G|} \langle \hat{\mathbf{e}}_i, \hat{\mathbf{w}}_L \cdots \hat{\mathbf{w}}_1 \rangle_M \\
&= \mathcal{F}_M \nabla_{\mathbf{w}_\ell} \frac{1}{|G|} \text{tr}[\hat{\mathbf{e}}_i (\hat{\mathbf{w}}_L \cdots \hat{\mathbf{w}}_1)^\dagger] \\
&= \mathcal{F}_M \nabla_{\mathbf{w}_\ell} \frac{1}{|G|} \text{tr}[\hat{\mathbf{w}}_{\ell+1}^\dagger \cdots \hat{\mathbf{w}}_L^\dagger \hat{\mathbf{e}}_i \hat{\mathbf{w}}_1^\dagger \cdots \hat{\mathbf{w}}_\ell^\dagger] \\
&= \mathcal{F}_M \nabla_{\mathbf{w}_\ell} \frac{1}{|G|} \langle \hat{\mathbf{w}}_{\ell+1}^\dagger \cdots \hat{\mathbf{w}}_L^\dagger \hat{\mathbf{e}}_i \hat{\mathbf{w}}_1^\dagger \cdots \hat{\mathbf{w}}_{\ell-1}^\dagger, \hat{\mathbf{w}}_\ell \rangle_M \\
&= \mathcal{F}_M \nabla_{\mathbf{w}_\ell} \left\langle \mathcal{F}_M^{-1}(\hat{\mathbf{w}}_{\ell+1}^\dagger \cdots \hat{\mathbf{w}}_L^\dagger \hat{\mathbf{e}}_i \hat{\mathbf{w}}_1^\dagger \cdots \hat{\mathbf{w}}_{\ell-1}^\dagger), \mathbf{w}_\ell \right\rangle \\
&= \hat{\mathbf{w}}_{\ell+1}^\dagger \cdots \hat{\mathbf{w}}_L^\dagger \hat{\mathbf{e}}_i \hat{\mathbf{w}}_1^\dagger \cdots \hat{\mathbf{w}}_{\ell-1}^\dagger
\end{aligned} \tag{A.24}$$

□

Letting $\mathbf{z} = \sum_n \tilde{\alpha}_n y_n \mathbf{x}_n$ and $\mathbf{W} = \mathbf{W}^\infty$, Lemma A.10 implies that

$$\mathcal{F}_M \nabla_{\mathbf{w}_\ell^\infty} \langle \mathcal{P}(\mathbf{W}^\infty), \mathbf{z} \rangle = \hat{\mathbf{w}}_{\ell+1}^\infty{}^\dagger \cdots \hat{\mathbf{w}}_L^\infty{}^\dagger \hat{\mathbf{z}} \hat{\mathbf{w}}_1^\infty{}^\dagger \cdots \hat{\mathbf{w}}_{\ell-1}^\infty{}^\dagger \tag{A.25}$$

By combining Equation A.10 with Lemma A.10, we have

$$\frac{1}{\gamma} \hat{\mathbf{w}}_\ell^\infty = \frac{1}{\gamma} \nabla_{\mathbf{w}_\ell^\infty} \left\langle \mathcal{P}(\mathbf{W}^\infty), \sum_n \alpha_n y_n \mathbf{x}_n \right\rangle \tag{A.26}$$

$$= \nabla_{\mathbf{w}_\ell^\infty} \left\langle \mathcal{P}(\mathbf{W}^\infty), \sum_n \tilde{\alpha}_n y_n \mathbf{x}_n \right\rangle \tag{A.27}$$

$$= \hat{\mathbf{w}}_{\ell+1}^\infty{}^\dagger \cdots \hat{\mathbf{w}}_L^\infty{}^\dagger \hat{\mathbf{z}} \hat{\mathbf{w}}_1^\infty{}^\dagger \cdots \hat{\mathbf{w}}_{\ell-1}^\infty{}^\dagger \tag{A.28}$$

As a result:

$$\begin{aligned}
\hat{\mathbf{w}}_\ell^\infty &= \gamma \hat{\mathbf{w}}_{\ell+1}^\infty{}^\dagger \cdots \hat{\mathbf{w}}_L^\infty{}^\dagger \hat{\mathbf{z}} \hat{\mathbf{w}}_1^\infty{}^\dagger \cdots \hat{\mathbf{w}}_{\ell-1}^\infty{}^\dagger \\
\hat{\mathbf{w}}_\ell^\infty \hat{\mathbf{w}}_\ell^\infty{}^\dagger &= \gamma \hat{\mathbf{w}}_{\ell+1}^\infty{}^\dagger \cdots \hat{\mathbf{w}}_L^\infty{}^\dagger \hat{\mathbf{z}} \hat{\mathbf{w}}_1^\infty{}^\dagger \cdots \hat{\mathbf{w}}_\ell^\infty{}^\dagger
\end{aligned} \tag{A.29}$$

Applying this relation with $\ell = L$, we have that

$$\hat{\mathbf{w}}_L^\infty \hat{\mathbf{w}}_L^\infty{}^\dagger = \gamma \hat{\mathbf{z}} \hat{\mathbf{w}}_1^\infty{}^\dagger \cdots \hat{\mathbf{w}}_L^\infty{}^\dagger \tag{A.30}$$

$$= \gamma \hat{\mathbf{z}} \hat{\boldsymbol{\beta}}^\infty{}^\dagger \tag{A.31}$$

Taking adjoints of both sides implies that $\hat{\mathbf{z}} \hat{\boldsymbol{\beta}}^\infty{}^\dagger$ is Hermitian, which will be useful later.

Let $\hat{\boldsymbol{\beta}}^\infty = \mathcal{F}_M \mathcal{P}(\mathbf{W}^\infty)$, from which we can derive the following recursion:

$$\begin{aligned}
\hat{\boldsymbol{\beta}}^\infty \hat{\boldsymbol{\beta}}^\infty{}^\dagger &= \hat{\mathbf{w}}_L^\infty \cdots \hat{\mathbf{w}}_1^\infty \hat{\mathbf{w}}_1^\infty{}^\dagger \cdots \hat{\mathbf{w}}_L^\infty{}^\dagger \\
&= \gamma^1 \hat{\mathbf{w}}_L^\infty \cdots \hat{\mathbf{w}}_2^\infty \hat{\mathbf{w}}_2^\infty{}^\dagger \cdots \hat{\mathbf{w}}_L^\infty{}^\dagger \hat{\mathbf{z}} \hat{\mathbf{w}}_1^\infty{}^\dagger \cdots \hat{\mathbf{w}}_L^\infty{}^\dagger \text{ by Equation A.29} \\
&= \gamma^2 \hat{\mathbf{w}}_L^\infty \cdots \hat{\mathbf{w}}_3^\infty \hat{\mathbf{w}}_3^\infty{}^\dagger \cdots \hat{\mathbf{w}}_L^\infty{}^\dagger \hat{\mathbf{z}} \hat{\mathbf{w}}_1^\infty{}^\dagger \hat{\mathbf{w}}_2^\infty{}^\dagger \cdots \hat{\mathbf{w}}_L^\infty{}^\dagger \hat{\mathbf{z}} \hat{\mathbf{w}}_1^\infty{}^\dagger \cdots \hat{\mathbf{w}}_L^\infty{}^\dagger \text{ again, by Equation A.29} \\
&= \gamma^2 \hat{\mathbf{w}}_L^\infty \cdots \hat{\mathbf{w}}_3^\infty \hat{\mathbf{w}}_3^\infty{}^\dagger \cdots \hat{\mathbf{w}}_L^\infty{}^\dagger \hat{\mathbf{z}} \hat{\boldsymbol{\beta}}^\infty{}^\dagger \hat{\mathbf{z}} \hat{\boldsymbol{\beta}}^\infty \text{ by definition of } \boldsymbol{\beta}^\infty \\
&= \gamma^2 \hat{\mathbf{w}}_L^\infty \cdots \hat{\mathbf{w}}_3^\infty \hat{\mathbf{w}}_3^\infty{}^\dagger \cdots \hat{\mathbf{w}}_L^\infty{}^\dagger (\hat{\mathbf{z}} \hat{\boldsymbol{\beta}}^\infty{}^\dagger)^2 \text{ by repeated application of Equation A.29} \\
&= \gamma^L (\hat{\mathbf{z}} \hat{\boldsymbol{\beta}}^\infty{}^\dagger)^L
\end{aligned} \tag{A.32}$$

Similarly to Equation A.29, we have:

$$\hat{\mathbf{w}}_\ell^{\infty \dagger} \hat{\mathbf{w}}_\ell^\infty = \gamma \hat{\mathbf{w}}_\ell^{\infty \dagger} \hat{\mathbf{w}}_{\ell+1}^{\infty \dagger} \dots \hat{\mathbf{w}}_L^{\infty \dagger} \hat{\mathbf{z}} \hat{\mathbf{w}}_1^{\infty \dagger} \dots \hat{\mathbf{w}}_{\ell-1}^{\infty \dagger} \quad (\text{A.33})$$

By considering $\ell = 1$, we have $\hat{\mathbf{w}}_1^{\infty \dagger} \hat{\mathbf{w}}_1^\infty = \hat{\beta}^{\infty \dagger} \hat{\mathbf{z}}$, which shows that $\hat{\beta}^{\infty \dagger} \hat{\mathbf{z}}$ is Hermitian as well.

Using Equation A.33, we can similarly reason about $\hat{\beta}^{\infty \dagger} \hat{\beta}^\infty$:

$$\begin{aligned} \hat{\beta}^{\infty \dagger} \hat{\beta}^\infty &= \hat{\mathbf{w}}_1^{\infty \dagger} \dots \hat{\mathbf{w}}_L^{\infty \dagger} \hat{\mathbf{w}}_L^\infty \dots \hat{\mathbf{w}}_1^\infty \\ &= \gamma \hat{\mathbf{w}}_1^{\infty \dagger} \dots \hat{\mathbf{w}}_{L-1}^{\infty \dagger} (\hat{\mathbf{w}}_L^{\infty \dagger} \hat{\mathbf{z}} \hat{\mathbf{w}}_1^{\infty \dagger} \dots \hat{\mathbf{w}}_{L-1}^{\infty \dagger}) \hat{\mathbf{w}}_{L-1}^\infty \dots \hat{\mathbf{w}}_1^\infty \\ &= \gamma (\hat{\beta}^{\infty \dagger} \hat{\mathbf{z}}) \hat{\mathbf{w}}_1^{\infty \dagger} \dots \hat{\mathbf{w}}_{L-1}^{\infty \dagger} \hat{\mathbf{w}}_{L-1}^\infty \dots \hat{\mathbf{w}}_1^\infty \\ &= \gamma (\hat{\beta}^{\infty \dagger} \hat{\mathbf{z}}) \hat{\mathbf{w}}_1^{\infty \dagger} \dots \hat{\mathbf{w}}_{L-1}^{\infty \dagger} \hat{\mathbf{w}}_L^{\infty \dagger} \hat{\mathbf{z}} \hat{\mathbf{w}}_1^{\infty \dagger} \dots \hat{\mathbf{w}}_{L-2}^{\infty \dagger} \hat{\mathbf{w}}_{L-2}^\infty \dots \hat{\mathbf{w}}_1^\infty \\ &= \gamma^2 (\hat{\beta}^{\infty \dagger} \hat{\mathbf{z}})^2 \hat{\mathbf{w}}_1^{\infty \dagger} \dots \hat{\mathbf{w}}_{L-2}^{\infty \dagger} \hat{\mathbf{w}}_{L-2}^\infty \dots \hat{\mathbf{w}}_1^\infty \\ &= \gamma^L (\hat{\beta}^{\infty \dagger} \hat{\mathbf{z}})^L \end{aligned} \quad (\text{A.34})$$

For $L > 2$, we have shown in Lemma A.9 that

$$\partial^o \|\hat{\beta}^\infty\|_p^{(S)} = \left\{ \mathbf{G} : \hat{\beta}^{\infty \dagger} \mathbf{G} = \frac{1}{\|\hat{\beta}^\infty\|_p^{(S)}} \sqrt{\hat{\beta}^{\infty \dagger} \hat{\beta}^\infty}^p \text{ and } \mathbf{G} \hat{\beta}^{\infty \dagger} = \frac{1}{\|\hat{\beta}^\infty\|_p^{(S)}} \sqrt{\hat{\beta}^\infty \hat{\beta}^{\infty \dagger}}^p \right\} \quad (\text{A.35})$$

Let's check that setting $\mathbf{G} = \hat{\mathbf{z}}$ satisfies this relation. Since $p = \frac{2}{L}$, $\frac{p}{2} = \frac{1}{L}$ and, by Equation A.32 and that $\hat{\mathbf{z}} \hat{\beta}^{\infty \dagger}$ is Hermitian:

$$(\hat{\beta}^\infty \hat{\beta}^{\infty \dagger})^{p/2} = \gamma \hat{\mathbf{z}} \hat{\beta}^{\infty \dagger} \quad (\text{A.36})$$

Similarly, by Equation A.34 and that $\hat{\beta}^{\infty \dagger} \hat{\mathbf{z}}$ is Hermitian:

$$(\hat{\beta}^{\infty \dagger} \hat{\beta}^\infty)^{\frac{p}{2}} = \gamma \hat{\beta}^{\infty \dagger} \hat{\mathbf{z}} \quad (\text{A.37})$$

By choice of γ , $\gamma = \|\hat{\beta}^\infty\|_{\frac{2}{L}}^{(S)}$. Thus, $\hat{\mathbf{z}} \in \partial^o \|\hat{\beta}^\infty\|_p^{(S)}$ as desired for $L > 2$.

For $L = 2$, $p = 1$ and by Lemma A.8 we had the following expression for the subgradient:

$$\partial^o \|\mathbf{A}\|_1^{(S)} = \left\{ \mathbf{G} : \|\mathbf{A}\|_1^{(S)} = \text{tr}[\mathbf{G}^\dagger \mathbf{A}], \|\mathbf{G}\|_\infty^{(S)} \leq 1 \right\} \quad (\text{A.38})$$

For this special case of $L = 2$, we have the required assumption:

As a technical condition, we had to assume in Assumption A.7 that $\|\hat{\mathbf{z}}\|_\infty^{(S)} \leq 1$. (We believe that with a refined analysis in future work, this assumption can be shown to be true given only our existing assumptions.) By the previous reasoning for $L = 2$, we had that $\hat{\beta}^{\infty \dagger} \hat{\beta}^\infty = \gamma^2 (\hat{\beta}^{\infty \dagger} \hat{\mathbf{z}})^2$. Since $\hat{\beta}^{\infty \dagger} \hat{\beta}^\infty = \mathbf{U} \mathbf{D}^2 \mathbf{U}^\dagger$ is positive semi-definite and symmetric, and since $\hat{\beta}^{\infty \dagger} \hat{\mathbf{z}}$ is Hermitian, we can take the square root of both sides and obtain that $\mathbf{U} \mathbf{D} \mathbf{U}^\dagger = \gamma \hat{\beta}^{\infty \dagger} \hat{\mathbf{z}} = \gamma \hat{\mathbf{z}}^\dagger \hat{\beta}^\infty$. Thus, $\text{tr}[\hat{\mathbf{z}}^\dagger \hat{\beta}^\infty] = \text{tr}[\mathbf{D} \mathbf{U}^\dagger \mathbf{U}] = \text{tr}[\mathbf{D}] = \|\hat{\beta}^\infty\|_1^{(S)}$, which is what was needed (since $\gamma = 1$ for the $L = 2$ case).

A.2.1 INFINITE GROUPS WITH BAND-LIMITED INPUTS

Here we consider the case of more general, infinite-dimensional compact Lie groups. Such groups admit a Fourier transform which is an operator between infinite-dimensional spaces, rather than a finite matrix as before, but which has the same key properties: a convolution theorem, and preservation of inner products. To be concrete, the Fourier transform is now defined as $\hat{f}(\rho) = \int_{g \in G} \rho(g) f(g) \mu(g)$, where $\mu(g)$ denotes the Haar measure for the group.

Since it is impossible to store a general function $x : G \rightarrow \mathbb{R}$, one must make simplifying assumptions on the input x . A common assumption is that x is band-limited in Fourier space, i.e. \hat{x} is supported only on a small (and finite) subset of Fourier coefficients contained within the irreps $\rho \in S$. For many such groups, there is a natural hierarchy of irreps (analogous to low frequencies and high frequencies in classical Fourier analysis), and so practical architectures typically assume only those corresponding to low frequencies are non-zero.

For ease of analysis, we will assume that the input functions and all convolutional filters are real-valued in Fourier space. The architecture of our G-CNN is the same as in the finite-dimensional group setting except that we will apply functions entirely in the finite-dimensional Fourier space over the irreps in S . Given a function \hat{x} supported only on S in Fourier space, we run gradient descent over only the Fourier coefficients on S of filters \hat{w}_i , and assume they are zero elsewhere. Before we proceed, we define \mathcal{F}_M in the natural way after restricting the irreps to the band-limiting space S :

$$\underline{\hat{f}} = \mathcal{F}_M \mathbf{f} = \bigoplus_{\rho \in S} \hat{f}(\rho)^{\oplus d_\rho} \in \text{GL} \left(\sum_{\rho \in S} d_\rho^2, \mathbb{C} \right). \quad (\text{A.39})$$

In this band-limited space, for each filter, there are $\sum_{\rho \in S} d_\rho^2$ trainable parameters corresponding to the entries of the irreps in S . Note that these entries are also orthogonal to each other with respect to the inner product $\langle \cdot, \cdot \rangle_M$ and thus form a vector space of dimension $\sum_{\rho \in S} d_\rho^2$.

Recall that we previously applied Theorem A.6 to the homogeneous polynomial parametrization from $(w_1, \dots, w_L) \rightarrow \langle x * w_1 * \dots * w_{L-1}, w_L \rangle$. Since we will operate in Fourier space, instead consider the homogeneous polynomial parametrization $\widehat{\mathbf{W}} \in \mathbb{R}^p$ containing the parameters stored in the matrices $\{\widehat{w}_1, \dots, \widehat{w}_L\}$. In other words, there are $p = L(\sum_{\rho \in S} d_\rho^2)$ parameters stored in the matrices contained in the set $\widehat{\mathbf{W}}$.

$$\widehat{\mathbf{W}}^{t+1} = \widehat{\mathbf{W}}^t - \eta_t \nabla_{\widehat{\mathbf{W}}} \mathcal{L}(\widehat{\mathbf{W}}^t) = \widehat{\mathbf{W}}^t - \eta_t \nabla_{\widehat{\mathbf{W}}} \mathcal{P}(\widehat{\mathbf{W}}^t) \nabla_{\widehat{\beta}} \mathcal{L}(\mathcal{P}(\widehat{\mathbf{W}}^t)) \quad (\text{A.40})$$

Note that here, iterates are only allowed to vary over the finite subset S of Fourier coefficients and are assumed to be equal to 0 elsewhere. If we assume further that the exponential loss converges to zero, that the gradients $\nabla_{\widehat{\beta}} \mathcal{L}(\widehat{\beta}^t)$ converge in direction, and that the iterates $\widehat{\mathbf{W}}^t$ themselves converge in direction to yield a separator with positive margin, then the limit direction of the parameters $\widehat{\mathbf{W}}^\infty = \lim_{t \rightarrow \infty} \frac{\widehat{\mathbf{W}}^t}{\|\widehat{\mathbf{W}}^t\|_2}$ is a positive scaling of a first order stationary point of the following optimization problem:

$$\min_{\widehat{\mathbf{W}} \in \mathbb{R}^p} \|\widehat{\mathbf{W}}\|_2^2 \quad \text{s.t.} \quad \forall n, y_n \langle \widehat{\mathbf{x}}_n, \mathcal{P}(\widehat{\mathbf{W}}) \rangle_M \geq 1. \quad (\text{A.41})$$

Again letting $\widehat{\mathbf{W}}^\infty = \tau \widehat{\mathbf{W}}^\infty$ denote the first order stationary point itself, $\widehat{\mathbf{W}}^\infty_i$ the individual layers (or parameter blocks) comprising $\widehat{\mathbf{W}}^\infty$, and $\widehat{\mathbf{W}}^\infty_i$ the individual layers comprising $\widehat{\mathbf{W}}^\infty$, we again have via the KKT conditions that:

$$\begin{aligned} & \exists \{\alpha_n : \alpha_n \geq 0\}_{n=1}^N \text{ s.t. } \alpha_n = 0 \text{ if } y_n \langle \widehat{\mathbf{x}}_n, \mathcal{P}(\widehat{\mathbf{W}}^\infty) \rangle_M > 1 \\ & \widehat{\mathbf{W}}_i^\infty = \nabla_{\widehat{\mathbf{W}}_i} \mathcal{P}(\widehat{\mathbf{W}}^\infty) \left[\sum_n \alpha_n y_n \widehat{\mathbf{x}}_n \right] = \nabla_{\widehat{\mathbf{W}}_i} \left\langle \mathcal{P}(\widehat{\mathbf{W}}^\infty), \sum_n \alpha_n y_n \widehat{\mathbf{x}}_n \right\rangle_M \end{aligned} \quad (\text{A.42})$$

Equivalently, writing the above in terms of the matrices \widehat{w}_ℓ^∞ , and defining $\widehat{\mathbf{z}} = \sum_n \alpha_n y_n \widehat{\mathbf{x}}_n$, we have an equivalence to Lemma A.10.

$$\begin{aligned}
\nabla_{\mathbf{w}_\ell} \mathcal{P}(\widehat{\mathbf{W}}^\infty)[\widehat{\mathbf{z}}] &= \nabla_{\widehat{\mathbf{w}}_\ell} \langle \mathcal{P}(\widehat{\mathbf{W}}^\infty), \widehat{\mathbf{z}} \rangle_M \\
&= \nabla_{\widehat{\mathbf{w}}_\ell} \langle \widehat{\mathbf{z}}, \mathcal{P}(\widehat{\mathbf{W}}^\infty) \rangle_M \\
&= \nabla_{\widehat{\mathbf{w}}_\ell} \langle \widehat{\mathbf{z}}, \widehat{\mathbf{w}}_L^\infty \cdots \widehat{\mathbf{w}}_1^\infty \rangle_M \\
&= \nabla_{\widehat{\mathbf{w}}_\ell} \text{tr}[\widehat{\mathbf{z}}(\widehat{\mathbf{w}}_L^\infty \cdots \widehat{\mathbf{w}}_1^{\infty\dagger})] \\
&= \nabla_{\widehat{\mathbf{w}}_\ell} \text{tr}[\widehat{\mathbf{w}}_{\ell+1}^{\infty\dagger} \cdots \widehat{\mathbf{w}}_L^{\infty\dagger} \widehat{\mathbf{z}} \widehat{\mathbf{w}}_1^{\infty\dagger} \cdots \widehat{\mathbf{w}}_\ell^{\infty\dagger}] \\
&= \nabla_{\widehat{\mathbf{w}}_\ell} \langle \widehat{\mathbf{w}}_{\ell+1}^{\infty\dagger} \cdots \widehat{\mathbf{w}}_L^{\infty\dagger} \widehat{\mathbf{z}} \widehat{\mathbf{w}}_1^{\infty\dagger} \cdots \widehat{\mathbf{w}}_{\ell-1}^{\infty\dagger}, \widehat{\mathbf{w}}_\ell^\infty \rangle_M \\
&= \widehat{\mathbf{w}}_{\ell+1}^{\infty\dagger} \cdots \widehat{\mathbf{w}}_L^{\infty\dagger} \widehat{\mathbf{z}} \widehat{\mathbf{w}}_1^{\infty\dagger} \cdots \widehat{\mathbf{w}}_{\ell-1}^{\infty\dagger}
\end{aligned} \tag{A.43}$$

Using the KKT conditions above and this fact, all of the manipulations demonstrating that a positive scaling of $\sum_n \alpha_n y_n \widehat{\mathbf{x}}_n$ is a first-order stationary point of the optimization problem below carry over exactly as they do in the previous part. This yields the following formal result:

Theorem A.11. *Consider a classification task with ground-truth linear predictor $\widehat{\beta}$, trained via a real-valued, Fourier-space, band-limited linear G-CNN architecture $\text{NN}(\widehat{\mathbf{x}}) = \langle \widehat{\mathbf{x}}, \widehat{\mathbf{w}}_L \cdots \widehat{\mathbf{w}}_1 \rangle$ with $L \geq 2$ layers under the exponential loss. Then for almost any datasets $\{\mathbf{x}_i, y_i\}$ separable by β , any bounded sequence of step sizes η_t , and almost all initializations, suppose that:*

- The loss converges to 0
- The gradients with respect to the end-to-end linear predictor $\widehat{\beta}$ converge in direction
- The iterates $\widehat{\mathbf{w}}_i$ themselves converge in direction to a separator with positive margin

When $L = 2$, we need an additional technical assumption, Assumption A.7. Then, the resultant linear predictor $\widehat{\beta}^\infty$ is a positive scaling of a first order stationary point of the optimization problem:

$$\min_{\widehat{\beta}} \left\| \widehat{\beta} \right\|_{2/L}^{(S)} \quad \text{s.t.} \quad \forall n, y_n \langle \widehat{\mathbf{x}}_n, \widehat{\beta} \rangle_M \geq 1. \tag{A.44}$$

B GROUP FOURIER TRANSFORMS

To aid the reader in understanding the notation and structure behind the group Fourier transform, the following exposition is given for reference and convenience. Here, we introduce important concepts from representation theory and from there, provide explicit constructions the group Fourier transform.

A representation of a group G is a vector space V together with a G -linear map $\rho : G \rightarrow \text{GL}(V)$. Of particular interest is the *regular representation* which we construct as follows. Let G be a group of order n and choose $V = \mathbb{C}^n$. Consider an element $u \in \mathbb{C}[G]$, so $u = a_1 g_1 + \cdots + a_n g_n$ where $a_i \in \mathbb{C}$ and $g_n \in G$, and the associated vector $\mathbf{u} \in \mathbb{C}^n$ such that $\mathbf{u} = [a_1 \cdots a_n]$.

The action of left-multiplication of u for any $h \in G$ yields $hu = a_1(hg_1) + \cdots + a_n(hg_n)$, which is equivalent to a permutation of the coefficients, so there is a unique matrix $H \in \text{GL}(\mathbb{C}^n)$ such that $H\mathbf{u}$ is equivalent to the associated vector for hu . The G -linear map $L_h : h \mapsto H$ is the (left) *regular representation*.

The direct sum of two G -representations (ρ_1, V_1) and (ρ_2, V_2) can be constructed by

$$(\rho_1 \oplus \rho_2)(g) = \begin{bmatrix} \rho_1(g) & \\ & \rho_2(g) \end{bmatrix} \tag{B.1}$$

The dimension d_ρ of a representation (ρ, V) is defined to be $\dim(V)$. A finite-dimensional representation is *irreducible* if it cannot be written as the direct sum of two nontrivial representations.

Denote \widehat{G} to be the set of isomorphism classes of irreducible subrepresentations of L_u , and let d_{ρ_i} denote the dimension of $\rho_i \in \widehat{G}$. As it turns out, there is an isomorphism

$$L_u \cong \bigoplus_{\rho \in \widehat{G}} \rho^{\oplus d_\rho}$$

Where ρ ranges over one representative from each of the isomorphism classes of \widehat{G} , repeated according to its multiplicity. In general, this decomposition is not uniquely determined as it depends on the choice of representatives. Throughout this paper, we choose representatives such that each ρ is *unitary*, meaning that every $\rho(g)$ is a unitary matrix.

Every function $f : G \rightarrow \mathbb{C}$ can be considered as equivalent to an element $u_f \in \mathbb{C}[G]$ by setting $u_f = f(g_1)g_1 + \dots + f(g_n)g_n$. And as we have already seen, every $u \in \mathbb{C}[G]$ can naturally be considered a subrepresentation of L_u . Then the *Fourier transform* of f at a representation ρ , denoted $\widehat{f}(\rho) \in \text{GL}(d_\rho, \mathbb{C})$ where $\widehat{f}(\rho) = \sum_{u \in G} f(u)\rho(u)$, can be considered as a projection of $L_u(u_f)$ onto the orthogonal subspace described by ρ . Throughout the paper we use slightly different notations and characterizations of the Fourier transform depending on the context, but all share this projection as the fundamental operation.

Recalling Equation B.1, there is a representation isomorphic to L_u , which we will suggestively call \mathcal{F}_M , that block-diagonalizes the image of L_u into orthogonal subspaces of unitary irreducible representations (each $\rho \in \widehat{G}$ repeated with multiplicity d_ρ):

$$\mathcal{F}_M(u) = \begin{bmatrix} \rho_1(u) & & \\ & \ddots & \\ & & \rho_j(g) \end{bmatrix} \in \text{GL}(\mathbb{C}^n) \quad (\text{B.2})$$

For the last piece of the puzzle, extend the domain of \mathcal{F}_M to all functions $f : G \rightarrow \mathbb{C}$ (by considering f as an element of $\mathbb{C}[G]$). Then we obtain

$$\mathcal{F}_M f = \begin{bmatrix} \widehat{f}(\rho_1) & & \\ & \ddots & \\ & & \widehat{f}(\rho_j) \end{bmatrix} \in \text{GL}(\mathbb{C}^n) \quad (\text{B.3})$$

Which we call the *matrix Fourier transform* of f .

We also make use of the *Fourier basis matrix* \mathcal{F} , which depends only on the group G and not the function f . To construct it, we first need the operation *Flatten* which vertically stacks the columns of a matrix. Then define the transform

$$\phi(f) = \begin{bmatrix} \text{Flatten}(\widehat{f}(\rho_1)) \\ \vdots \\ \text{Flatten}(\widehat{f}(\rho_j)) \end{bmatrix} \quad (\text{B.4})$$

Letting $e_i : G \rightarrow \mathbb{C}$ be the indicator function $e_i(g_j) = \mathbb{1}_{i=j}$ we can describe the unitary Fourier basis matrix \mathcal{F} for a group G as a row-scaling of

$$\mathcal{F} \propto [\phi(e_1) \quad \phi(e_2) \quad \dots \quad \phi(e_n)] \quad (\text{B.5})$$

In other words, given a column vectorization \mathbf{f} of a function f such that $\mathbf{f}_i = f(g_i)$, then \mathcal{F} is the matrix such that the ‘unflattening’ of $\mathcal{F}\mathbf{f}$ is equal to $\mathcal{F}_M f$ up to the row-scaling constants. Thus we can treat the group Fourier transform either as an abstract isomorphism or as a concrete matrix-vector multiplication, depending on the application.

The matrix \mathcal{F} can be explicitly constructed as described in Definition 4.1. Denoting $e_{[\rho, i, j]}$ as the column-major vectorized basis for element ρ_{ij} in the group Fourier transform, then we can form the matrix

$$\mathcal{F} = \sum_{u \in G} \sum_{\rho \in \widehat{G}} \frac{\sqrt{d_\rho}}{\sqrt{|G|}} \sum_{i, j=1}^{d_\rho} \rho(u)_{ij} e_{[\rho, i, j]} e_u^T. \quad (\text{B.6})$$

For visualization, consider the dihedral group $D_6 = \{1, r, r^2, a, ar, ar^2\}$, which has three irreducible representations ρ_1, ρ_2, ρ_3 (up to isomorphism) of dimensions 1, 1, and 2 respectively, and let $f : D_6 \rightarrow \mathbb{C}$. Using colors instead of values at first to avoid numerical clutter:

$$\hat{f}(\rho_1) = \begin{bmatrix} \text{blue} \end{bmatrix} \quad \hat{f}(\rho_2) = \begin{bmatrix} \text{orange} \end{bmatrix} \quad \hat{f}(\rho_3) = \begin{bmatrix} \text{green} & \text{green} \\ \text{green} & \text{green} \end{bmatrix}$$

Since $L_{D_6} \cong \rho_1 \oplus \rho_2 \oplus \rho_3^2$, we can get something like

$$\mathcal{F}_M f = \begin{bmatrix} \text{blue} & & & & \\ & \text{orange} & & & \\ & & \text{green} & & \\ & & & \text{green} & \\ & & & & \text{green} \end{bmatrix} \in \text{GL}(\mathbb{C}^n)$$

Whereas for the unitary Fourier basis matrix we have the form

$$\mathcal{F} \propto \begin{array}{c} \begin{array}{c} \rho_1 \\ \rho_2 \\ [\rho_3]_{11} \\ [\rho_3]_{21} \\ [\rho_3]_{12} \\ [\rho_3]_{22} \end{array} \begin{array}{c} 1 \quad r \quad r^2 \quad a \quad ar \quad ar^2 \\ \begin{bmatrix} \text{blue} & \text{blue} & \text{blue} & \text{blue} & \text{blue} & \text{blue} \\ \text{orange} & \text{orange} & \text{orange} & \text{orange} & \text{orange} & \text{orange} \\ \text{green} & \text{green} & \text{green} & \text{green} & \text{green} & \text{green} \\ \text{green} & \text{green} & \text{green} & \text{green} & \text{green} & \text{green} \\ \text{green} & \text{green} & \text{green} & \text{green} & \text{green} & \text{green} \\ \text{green} & \text{green} & \text{green} & \text{green} & \text{green} & \text{green} \end{bmatrix} \end{array} \end{array}$$

Note that we do not yet include the row-scaling constants. Now explicitly, choosing ρ_1 the trivial irrep, ρ_2 the sign irrep, and ρ_3 the representation sending

$$\rho_3(a) = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad \rho_3(r) = \begin{bmatrix} \omega & 0 \\ 0 & \omega^2 \end{bmatrix} \quad (\text{B.7})$$

Where $\omega = e^{2\pi i/3}$, then the matrix \mathcal{F} is exactly

$$\mathcal{F} = \frac{1}{\sqrt{6}} \begin{array}{c} \begin{array}{c} \rho_1 \\ \rho_2 \\ [\rho_3]_{11} \\ [\rho_3]_{21} \\ [\rho_3]_{12} \\ [\rho_3]_{22} \end{array} \begin{array}{c} 1 \quad r \quad r^2 \quad a \quad ar \quad ar^2 \\ \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & -1 & -1 & -1 \\ \sqrt{2} & \sqrt{2}\omega & \sqrt{2}\omega^2 & 0 & 0 & 0 \\ 0 & 0 & 0 & \sqrt{2} & \sqrt{2}\omega & \sqrt{2}\omega^2 \\ 0 & 0 & 0 & \sqrt{2}\omega^2 & \sqrt{2}\omega & \sqrt{2} \\ \sqrt{2} & \sqrt{2}\omega^2 & \sqrt{2}\omega & 0 & 0 & 0 \end{bmatrix} \end{array} \end{array} \quad (\text{B.8})$$

Note that the above is only one possible way of writing \mathcal{F} since the 2-dimensional irreducible representation of D_6 is unique only up to conjugation by a unitary matrix.

C UNCERTAINTY PRINCIPLES FOR GROUPS

In mathematics, uncertainty principles categorize trade-offs of the “amount of information” stored in a function between canonically conjugate regimes, *e.g.*, position (real regime) and momentum (Fourier regime) of a physical particle. More generically, uncertainty principles show that a function and its Fourier transform cannot both be very localized or concentrated. In the context of group theory, uncertainty principles specifically show that sparse support in either the real or Fourier regime of a group necessarily implies dense support in the conjugate regime (Wigderson & Wigderson, 2021). Such results are directly relevant when interpreting implicit regularization of linear G-CNNs which bias gradient descent towards sparse solution in the Fourier basis of the group. In this section, we formally state and summarize these group theoretic uncertainty principles to provide intuition into the properties of functions which linear group convolutional networks are likely to learn.

Abelian groups For abelian groups, the fundamental uncertainty principle details a trade-off between the norms of a function in its real and Fourier bases.

Theorem C.1 (Generalization of Donoho-Stark Theorem (Wigderson & Wigderson, 2021)). *Given a finite abelian group G , let \hat{G} be the set of irreducible representations of G and $f : G \rightarrow \mathbb{C}$ be a function mapping group elements to complex numbers. Let \mathbf{f} be the vectorized function and \mathcal{F} be the unitary group Fourier transform (see Definition 4.1), then*

$$\frac{\|\mathbf{f}\|_1}{\|\mathbf{f}\|_\infty} \frac{\|\mathcal{F}\mathbf{f}\|_1}{\|\mathcal{F}\mathbf{f}\|_\infty} \geq |G| \quad (\text{C.1})$$

Remark. *Since the size of the support of a vector is bounded by $|\text{supp}(\mathbf{v})| \geq \frac{\|\mathbf{v}\|_1}{\|\mathbf{v}\|_\infty}$, this directly implies that $|\text{supp}(\mathbf{f})| |\text{supp}(\mathcal{F}\mathbf{f})| \geq |G|$ recovering the Donoho-Stark theorem (Donoho & Stark, 1989; Matusiak et al., 2004).*

Non-abelian groups Since non-abelian groups have matrix valued irreducible representations, uncertainty theorems must account for norms and notions of support in the context of matrices. Here, we will provide two different uncertainty theorems for the non-abelian setting – one via the rank of irreducible representations and another via the Schatten norm of irreducible representations. Uncertainty relations for non-abelian groups make use of the matrix group Fourier transform detailed in Definition 4.1.

Theorem C.2 (Meshulam uncertainty theorem (Meshulam, 1992)). *Given a finite non-abelian group G , let \hat{G} be the set of irreducible representations of G and $f : G \rightarrow \mathbb{C}$ be a function mapping group elements to complex numbers. Let \mathbf{f} be the vectorized function and \mathcal{F}_M be the matrix group Fourier transform (see Definition 4.1), then*

$$|\text{supp}(\mathbf{f})| \text{rank}(\mathcal{F}_M \mathbf{f}) = |\text{supp}(\mathbf{f})| \left[\sum_{\rho \in \hat{G}} d_\rho \text{rank}(\hat{f}(\rho)) \right] \geq |G| \quad (\text{C.2})$$

The above theorem shows that the rank of the matrix Fourier transformed function is the proper notion of support in the uncertainty theorem for a non-abelian group. A stronger uncertainty principle which is a more direct corollary to Theorem C.1 can be obtained via the Schatten norms of the irreducible representations as shown below.

Theorem C.3 (Kuperberg uncertainty theorem (Wigderson & Wigderson, 2021)). *Given a finite non-abelian group G , let \hat{G} be the set of irreducible representations of G and $f : G \rightarrow \mathbb{C}$ be a function mapping group elements to complex numbers. Let \mathbf{f} be the vectorized function and \mathcal{F}_M be the matrix group Fourier transform (see Definition 4.1), then*

$$\frac{\|\mathbf{f}\|_1}{\|\mathbf{f}\|_\infty} \frac{\|\mathcal{F}_M \mathbf{f}\|_1^{(S)}}{\|\mathcal{F}_M \mathbf{f}\|_\infty^{(S)}} = \frac{\|\mathbf{f}\|_1}{\|\mathbf{f}\|_\infty} \frac{\sum_{\rho \in \hat{G}} d_\rho \|\hat{f}(\rho)\|_1^{(S)}}{\max_{\rho \in \hat{G}} \|\hat{f}(\rho)\|_\infty^{(S)}} \geq |G|. \quad (\text{C.3})$$

D VISUALIZING THE IMPLICIT BIAS

Implicit biases induced by the G-CNN architectures studied here can be readily observed by analyzing coefficients of the linearized transformation in the Real or Fourier regimes. Here, we visualize the linearized outputs a 3-layer linear G-CNN over the Dihedral group D_{60} which has 4 scalar irreps and 14 irreps of dimension 2 (hence 2×2 matrices). Figure 6 shows these linearized coefficients of the G-CNN, CNN, and FC at initialization and training.

As evident in Figure 6, the learned coefficients of the G-CNN are sparse in the Fourier regime of the group. This sparsity pattern appears over blocks of irreps of length four, corresponding to coefficients of the 2×2 irreps of D_{60} . Furthermore, the values of the coefficients within a block are roughly constant, highlighting the bias towards low rank irreps. The relative denseness of coefficients of the trained G-CNN in the real regime is inherent due to the uncertainty principles of group functions. Unlike the G-CNN, the fully connected network (FC) appears to have no bias towards sparseness in its coefficients in either the real or Fourier regime. On a related note, the cyclic group of CNNs share some of the same irreps as those of the G-CNN studied here. This may be one explanation for the partial sparsity patterns observed in the coefficients of the CNN in the Fourier regime.

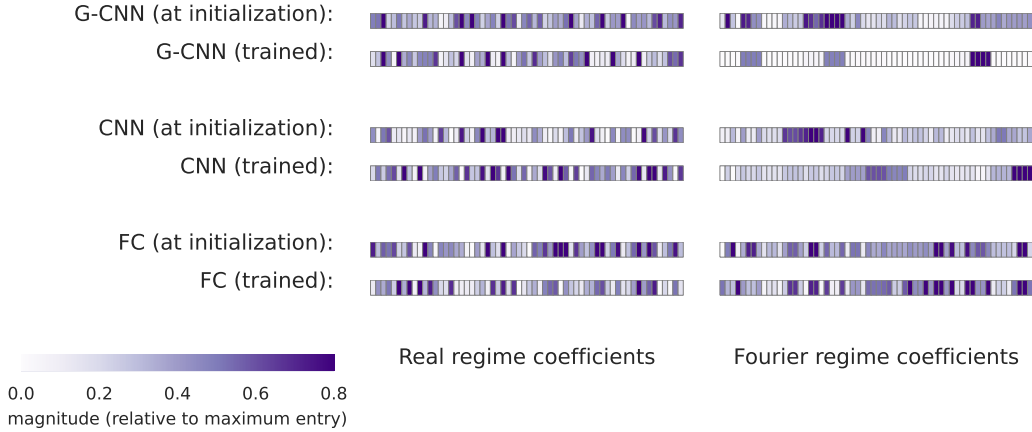


Figure 6: Linearized functions of the G-CNN (over the dihedral group D_{60}) are sparse in the Fourier regime of the group. Furthermore, the sparsity pattern shows up in blocks of 4 which are the blocks containing coefficients of individual irreps of the dihedral group D_{60} . Linearized functions in the real regime of the G-CNN, in contrast, are rather dense, highlighting the uncertainty principles inherent in the implicit bias of G-CNNs. Sparsity in the group Fourier regime is not evident in CNNs or fully connected (FC) networks.

E COMPUTATIONAL DETAILS

As described in Section 6, for all models we use three-layer networks over \mathbb{R}^G and binary classification tasks trained via standard gradient descent with exponential loss. We often train networks on isotropic Gaussian data points which are random vectors whose entries are drawn i.i.d. from a standard Normal. For the linear networks on simulated data (with a fully connected output layer) we use the groups D_8 , and $(C_{10} \times C_{10} \times C_2)$, and $(C_5 \times C_5) \rtimes Q_8$. For the linear and nonlinear networks on MNIST, we use $(C_{28} \times C_{28}) \times D_8$. For the networks with ReLU activations and a linear layer (instead of pooling) we use the groups D_8 and D_{60} . For the experiments on MNIST with non-linear networks, we have used the `e2cnn` package Weiler & Cesa (2019). The weights are initialized with the standard uniform initialization. We choose an appropriate learning rate for each task depending on the dimension and magnitude of the values - all learning rate choices are specified in the attached code. Since each problem is overparameterized the loss will almost surely converge to zero, so we choose enough training epochs to achieve satisfactory convergence—this ended up being around 500 epochs for most tasks. For each plot, we report 95% confidence intervals over 10 to 50 runs, depending on the classification task. The computational resources used are modest - commodity hardware should suffice to fully reproduce our results. For further details, please see the attached code in the supplementary materials.

E.1 ADDITIONAL EXPERIMENTS

We include here the real and Fourier space plots for an experiment on $C_{10} \times C_{10} \times C_2$ (Figure 7), an abelian group with three generators which captures more complicated group operations than cyclic shifts or standard planar convolutions. Inputs are vectors with elements drawn i.i.d. from the standard normal distribution.

E.2 OMITTED REAL-SPACE PLOTS

E.3 LOSS PLOTS

ACKNOWLEDGMENTS

Redacted until publication.

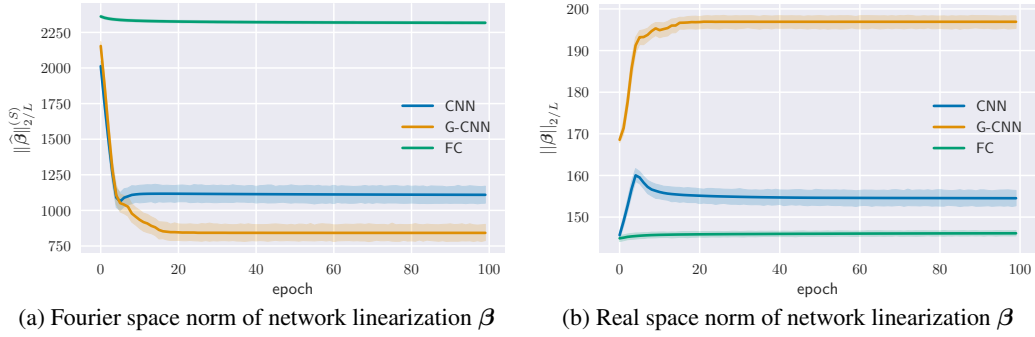


Figure 7: Norms of the linearizations of three different linear architectures for the abelian group $G = C_{10} \times C_{10} \times C_2$ trained on a binary classification task with six isotropic Gaussian data points.

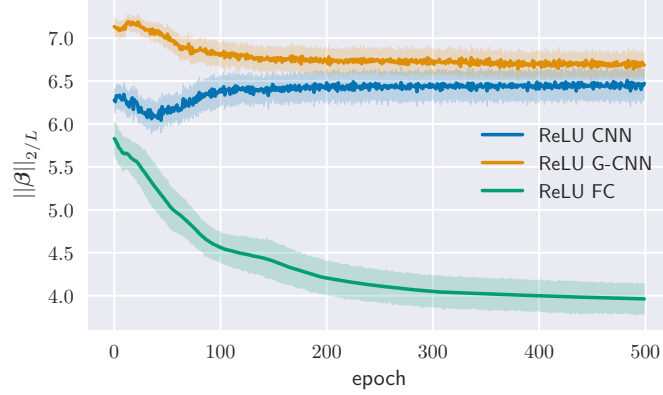


Figure 8: Real-space norm for D_8 with ReLU networks and 2 Gaussian training points. See Figure 5a for comparable plot of norms in Fourier space.

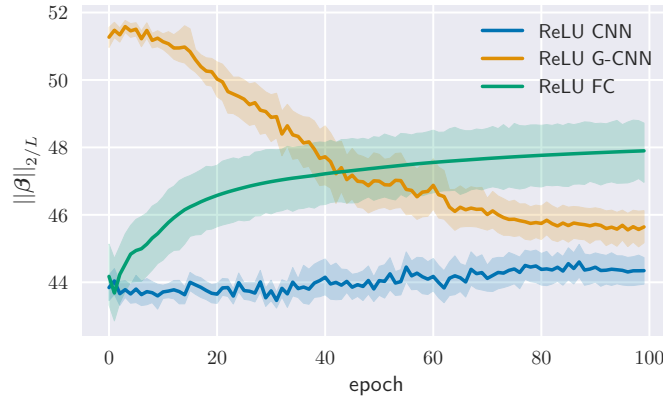


Figure 9: Real-space norm for D_{60} with ReLU networks, 10 Gaussian training points. See Figure 5b for comparable plot of norms in Fourier space.

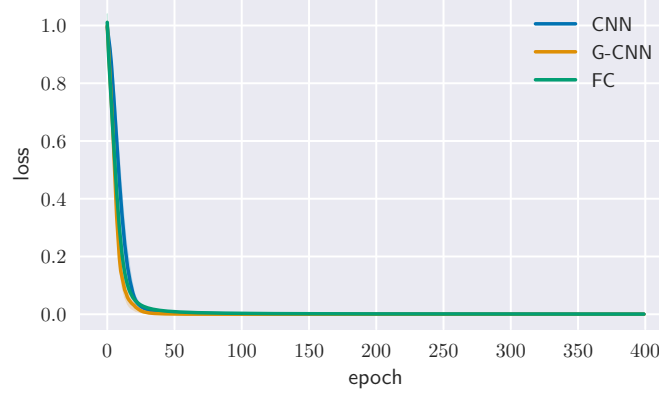


Figure 10: Loss trajectory for the setting of D_8 (see Figure 1). Networks are trained on 2 Gaussian i.i.d. data points.

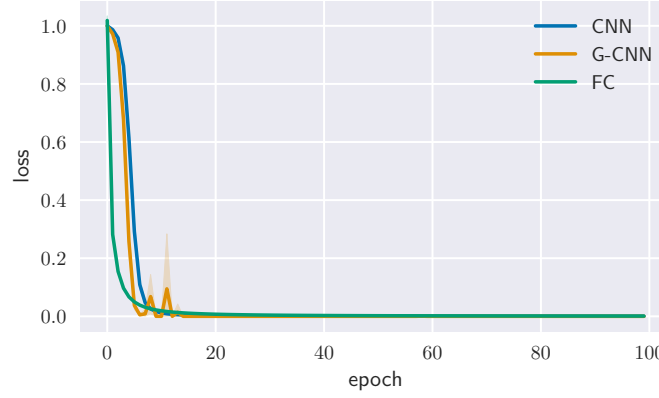


Figure 11: Loss trajectory for the setting of $C_{10} \times C_{10} \times C_2$ (see Figure 7). Networks are trained on 6 Gaussian i.i.d. data points.

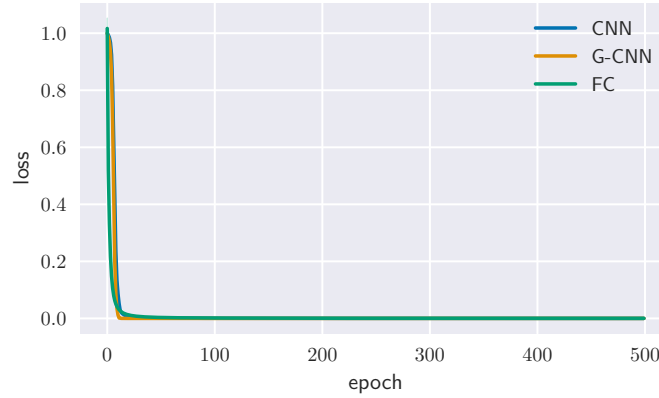


Figure 12: Loss trajectory for the setting of $(C_5 \times C_5) \times Q_8$ (see Figure 3). Networks are trained on 10 Gaussian i.i.d. data points.

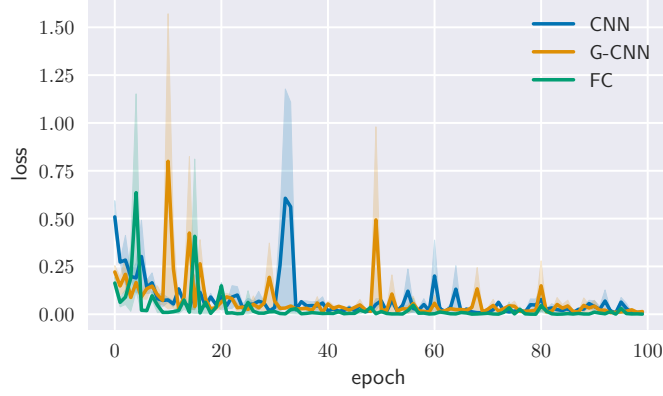


Figure 13: Loss trajectory for the setting of $(C_{28} \times C_{28}) \times D_8$ (see Figure 4). Networks are trained on the digits 1 and 5 from MNIST.

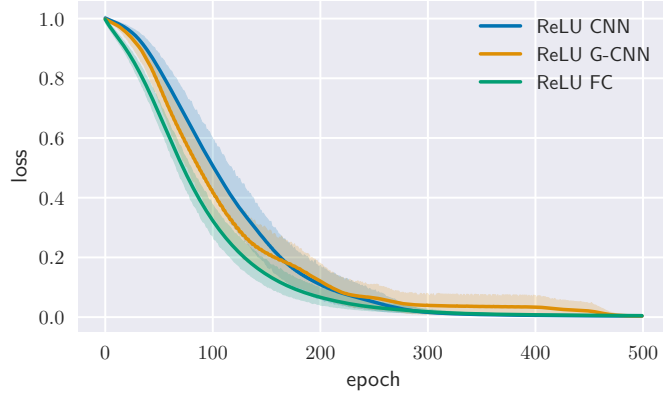


Figure 14: Loss trajectory for the setting of D_8 (see Figure 5a). Networks are nonlinear with ReLU activations and trained on 2 Gaussian i.i.d. data points.

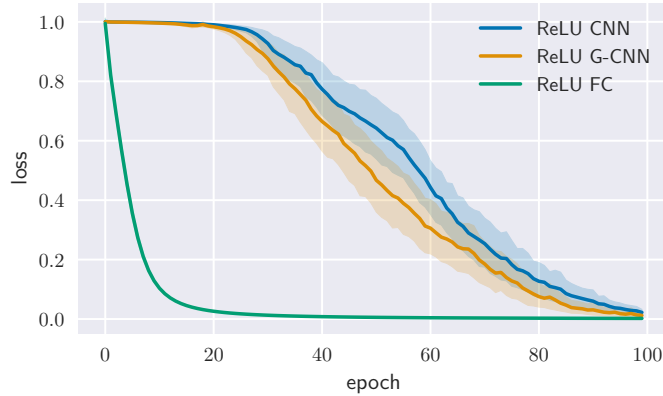


Figure 15: Loss trajectory for the setting of D_{60} (see Figure 5b). Networks are nonlinear with ReLU activations and trained on 10 distinct frequencies as data points.