



HAL
open science

Understanding Equivariant Self-Supervised Learning in Musical Pitch Class Space

Vincent Lostanlen, Yuexuan Kong, Gabriel Meseguer-Brocal, Mathieu Lagrange,
Romain Hennequin

► **To cite this version:**

Vincent Lostanlen, Yuexuan Kong, Gabriel Meseguer-Brocal, Mathieu Lagrange, Romain Hennequin. Understanding Equivariant Self-Supervised Learning in Musical Pitch Class Space. IEEE Signal Processing Letters, 2025. ⟨hal-05297975⟩

HAL Id: hal-05297975

<https://hal.science/hal-05297975v1>

Submitted on 5 Oct 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons CC BY 4.0 - Attribution - International License

Understanding Equivariant Self-Supervised Learning in Musical Pitch Class Space

Vincent Lostanlen, Yuexuan Kong, Gabriel Meseguer-Brocal, Mathieu Lagrange, Romain Hennequin

Abstract—STONE, which stands for self-supervised tonality estimator, has recently demonstrated the practical feasibility of recognizing key signatures in music signals given little or no human annotation. In this article, we revisit STONE from a more theoretical standpoint. We show that cross-power spectral density (CPSD) defines a differentiable measure of harmonic discrepancy between key signature profiles (KSP). Having set the CPSD frequency to seven cycles per octave, we offer a geometric interpretation of this discrepancy via the circle of fifths and conduct an algebraic study to prove that all its local minima are global. We rely on the equivariance property of deep convolutional networks to prove that the STONE loss function is invariant to circular frequency shifts of the constant-Q transform. We conclude by identifying a phenomenon of spontaneous symmetry breaking in STONE: since modes of limited transposition (e.g., augmented, diminished, whole-tone) are multistable in the circle of fifths, the associated CPSD gradient is driven towards more “tonal” (i.e., asymmetric) scales, such as diatonic or pentatonic.

Index Terms—Convolutional networks, cross-power spectral density, equivariance, music, self-supervised learning.

I. INTRODUCTION

EQUIVARIANCE is an algebraic property which is useful in the design of deep convolutional networks (convnets) [1]. In the context of music signal processing, equivariance to pitch transposition may be accomplished by applying a constant- Q transform (CQT), i.e., a time–frequency representation in which the frequency boundaries between adjacent subbands follow a geometric progression: $\xi_n = 2^{n/Q} \xi_0$ where n and Q are integers and where ξ_0 is a positive constant [2].

Certain music information retrieval (MIR) tasks combine pitch equivariance with octave equivalence, i.e., *invariance* to pitch transposition for intervals which are multiples of the octave: $k = mQ$ for any $m \in \mathbb{Z}$. These considerations have inspired many techniques in feature engineering, e.g., chromagrams [3], spiral scattering [4], folded CQT [5]; and feature learning, e.g., deep chroma [6], spiral convnets [7], and Tonnetz-space transform [8]. Another line of research has proposed to discover pitch equivariance and octave equivalence from unlabeled CQT samples of musical sounds [9], [10]. Yet, until recently, these algebraic properties were regarded as insufficient for the modeling of high-level music information such as tonality or chords. For these tasks, the CQT–convnet pipeline and its octave-equivalent variants have historically been trained via supervised learning [11].

The situation changed in 2024 with the publication of the first self-supervised tonality estimator (STONE) [12]. Building upon prior work in self-supervised fundamental frequency estimation [13], STONE is a convnet which is trained to regress artificial pitch transpositions between any two unlabeled musical excerpts from the same audio track. It does so with two essential tools: a fully pitch-equivariant and octave-invariant convnet, named ChromaNet; and a differentiable loss function based on circular cross-power spectral density (CPSD). STONE learns a 12-dimensional representation of CQT, which we call *key signature profile* (KSP) because it correlates with key signature. A modified version of STONE, known as S-KEY, has recently surpassed supervised models in the automatic classification of major and minor keys [14]. But despite its practical success in MIR, STONE lacks a rigorous foundation.

In this article, we present a mathematical theory of equivariant self-supervised learning (SSL) in musical pitch space. The key idea is to harness the algebraic properties of the discrete Fourier transform (DFT), particularly the convolution theorem and modular arithmetics between integer exponents of the complex roots of unity. Our main result is that STONE learns to solve a well-posed inverse problem up to a global pitch class shift: the associated CPSD-based objective admits Q global minima, which coincide with paired vertices of a regular polygon in the complex plane. Optimizing convnet weights with this objective encourages equivariance to key signature versus robustness to other aspects of musical information, such as melody, rhythm, lyrics, and instrumentation.

Section II presents the STONE model, particularly the KSP operator for octave equivalence in the ChromaNet and the (ω, k) -discrepancy as a differentiable and nonconstrative objective for equivariant SSL. Section III shows that the CPSD-based objective is *sparsity-inducing*: at any of its global minima, the two ChromaNet predictions have a single nonzero entry, which may be interpreted as a key signature. Section IV shows that the objective is *invariant*: pitch-shifting both KSPs leaves the objective unchanged; in musical terms, the learned assignment of CQT to key signature is relative, not absolute. Section V shows that the objective is *spontaneously symmetry-breaking*: it is globally maximal at periodic KSPs; in particular those which describe modes of limited transposition, e.g., augmented, diminished, and whole-tone.

Together, these three properties prevent STONE from collapsing to a trivial solution, e.g., a constant KSP; explain how the lack of manual annotation may be obviated by prior knowledge about signal processing and musical harmony; and enable its interpretability in the practical context of tonality estimation.

V. Lostanlen and M. Lagrange are with Nantes Université, École Centrale Nantes, CNRS, LS2N, UMR 6004, F-44000 Nantes, France. Y. Kong, G. Meseguer-Brocal, and R. Hennequin are with Deezer Research, Paris, 75009, France. This research was funded, in whole or in part, by the French National Research Agency (ANR) under the project MuReNN (ANR-23-CE23-0007-01).

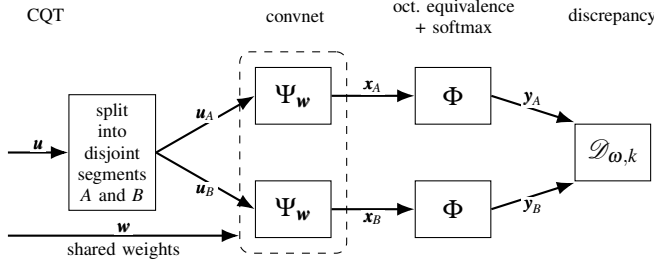


Fig. 1. Flowchart of a deep learning model for self-supervised tonality estimation (STONE). We split constant- Q transform (CQT) \mathbf{u} of an unlabeled music signal into disjoint segments A and B , yielding a pair of CQTs $(\mathbf{u}_A, \mathbf{u}_B)$ which we assume to be in the same key. We feed \mathbf{u}_A and \mathbf{u}_B to a deep convolutional network (convnet) $\Psi_{\mathbf{w}}$, followed by a non-learnable operator Φ which performs octave equivalence and softmax. This results in two key signature profiles (KSP) \mathbf{y}_A and \mathbf{y}_B ; i.e., nonnegative vectors summing to one in dimension Q . These KSPs serve as input to $\mathcal{D}_{\omega,k}$, a measure of harmonic discrepancy which is based on the cross-power spectral density between \mathbf{y}_A and \mathbf{y}_B and which is differentiable with respect to neural network weights \mathbf{w} .

Notation. Integers J and Q are the number of octaves and number of bins per octave in the CQT, respectively. The integer k is the pitch interval between the key signatures of A and B , measured in semitones if $\omega = 1$ or in perfect fifths if $\omega = 7$.

II. THE STONE MODEL FOR TONALITY ESTIMATION

Definition II.1. Given $Q \in \mathbb{N}^*$, the *standard Q -simplex* is the set of vectors in \mathbb{R}^Q with nonnegative entries summing to one:

$$S_Q = \{\mathbf{y} \in \mathbb{R}^Q \mid \|\mathbf{y}\|_1 = 1; \forall q \in \mathbb{Z}/Q\mathbb{Z}, \mathbf{y}[q] \geq 0\}. \quad (1)$$

Definition II.2. Given $J, Q \in \mathbb{N}^*$, the *key signature profile (KSP) operator* is the function $\Phi: \mathbb{R}^{JQ} \rightarrow S_Q$ such that:

$$\forall \mathbf{x} \in \mathbb{R}^{JQ}, \forall q \in \mathbb{Z}/Q\mathbb{Z}, \Phi(\mathbf{x})[q] = \frac{\exp\left(\sum_{j=0}^{J-1} \mathbf{x}[jQ+q]\right)}{\sum_{q'=0}^{Q-1} \exp\left(\sum_{j=0}^{J-1} \mathbf{x}[jQ+q']\right)}. \quad (2)$$

Definition II.3. Given $Q \in \mathbb{N}^*$, $\omega \in \mathbb{Z}/Q\mathbb{Z}$ and $k \in \mathbb{Z}$, the (ω, k) -discrepancy between two KSP vectors \mathbf{y}_A and \mathbf{y}_B is:

$$\mathcal{D}_{\omega,k}(\mathbf{y}_A, \mathbf{y}_B) = \frac{1}{2} \left| e^{2\pi i \omega k / Q} - \hat{\mathbf{y}}_A[\omega] \hat{\mathbf{y}}_B[\omega]^* \right|^2 \quad (3)$$

where $\hat{\mathbf{y}}_A[\omega] = \sum_{q=0}^{Q-1} \mathbf{y}_A[q] \exp(-2\pi i \omega q / Q)$ is the discrete Fourier transform (DFT) of \mathbf{y}_A at ω and idem for \mathbf{y}_B , and where the asterisk symbol denotes complex conjugation.

Note. The (ω, k) -discrepancy is a function of the circular cross-correlation \mathbf{R}_{AB} between \mathbf{y}_A and \mathbf{y}_B , defined by: $\mathbf{R}_{AB}[\tau] = \sum_{q=0}^{Q-1} \mathbf{y}_A[q] \mathbf{y}_B[(q+\tau) \bmod Q]$. Indeed, applying a circular change of variable $\tau = (q' - q) \bmod Q$ to $\hat{\mathbf{y}}_A[\omega] \hat{\mathbf{y}}_B[\omega]^*$ leads to a convolution theorem of the form:

$$\begin{aligned} \hat{\mathbf{y}}_A[\omega] \hat{\mathbf{y}}_B[\omega]^* &= \sum_{q=0}^{Q-1} \sum_{q'=0}^{Q-1} \mathbf{y}_A[q] \mathbf{y}_B[q'] \exp\left(2\pi i \frac{\omega(q'-q)}{Q}\right) \\ &= \sum_{q=0}^{Q-1} \sum_{\tau=0}^{Q-1} \mathbf{y}_A[q] \mathbf{y}_B[(q+\tau) \bmod Q] \exp\left(2\pi i \frac{\omega\tau}{Q}\right) \\ &= \widehat{\mathbf{R}}_{AB}[\omega]^* = \widehat{\mathbf{R}}_{BA}[\omega]. \end{aligned} \quad (4)$$

$\widehat{\mathbf{R}}_{AB}$ is sometimes called *cross-power spectral density (CPSD)*.

Definition II.4. Given $Q \in \mathbb{N}^*$ and $\omega, k \in \mathbb{Z}/Q\mathbb{Z}$; two CQT spectrograms $\mathbf{u}_A, \mathbf{u}_B$; and a model $\Psi_{\mathbf{w}}$ with trainable weights \mathbf{w} , the *STONE loss* is the function $\mathcal{L}_{\omega,k}: W \rightarrow \mathbb{R}_{\geq 0}$ such that:

$$\mathcal{L}_{\omega,k}(\mathbf{w} \mid \mathbf{u}_A, \mathbf{u}_B) = \mathcal{D}_{\omega,k}((\Phi \circ \Psi_{\mathbf{w}})(\mathbf{u}_A), (\Phi \circ \Psi_{\mathbf{w}})(\mathbf{u}_B)), \quad (5)$$

where W is the domain of \mathbf{w} , Φ is the key signature profile operator (Definition II.2) and $\mathcal{D}_{\omega,k}$ is (ω, k) -discrepancy (Definition II.3). The vertical bar in $\mathcal{L}_{\omega,k}(\mathbf{w} \mid \mathbf{u}_A, \mathbf{u}_B)$ denotes a separation between trainable parameters and input data.

Note. In practice, we construct pairs $(\mathbf{u}_A, \mathbf{u}_B)$ by extracting two disjoint 15-second segments within the same unlabeled music signal \mathbf{u} . By setting $k = 0$ in $\mathcal{D}_{\omega,k}$, we assume that \mathbf{u}_A and \mathbf{u}_B are in the same (unknown) key. We concede that this assumption is occasionally disproven by the presence of key changes between segments A and B . However, key changes have become increasingly rare in mainstream pop music, at least at the 15-second time scale¹. Thus, the founding hypothesis of STONE is that the rate of label noise (setting $k = 0$ despite occasional key changes) is outweighed by the sheer scale of available data for SSL, of the order of one million songs [14].

Proposition II.5. The STONE loss is differentiable.

Proof. The model $\Psi_{\mathbf{w}}$ is differentiable over W , by construction. The KSP operator Φ is differentiable over \mathbb{R}^{JQ} as a composition of a linear transformation and a pointwise softmax (Equation 2). Lastly, $\mathcal{D}_{\omega,k}$ is differentiable with respect to the real and imaginary parts of both its input arguments as a polynomial of degree two (Equation 4). We conclude using Equation 5. ■

Note. We typically set $Q = 12$ to match twelve-tone equal temperament in music. For the sake of conciseness, we do not recall the dependency of Φ , $\mathcal{D}_{\omega,k}$, and $\mathcal{L}_{\omega,k}$ upon Q .

III. SPARSITY PROPERTIES

Definition III.1. Given $Q \in \mathbb{N}^*$ and $\omega \in \mathbb{Z}/Q\mathbb{Z}$, the *Kronecker delta symbol* at index $\alpha \in \mathbb{Z}/Q\mathbb{Z}$ is the vector such that, for every $q \in \mathbb{Z}/Q\mathbb{Z}$, $\delta_{\alpha}[q] = 1$ if $q = \alpha \bmod Q$ and zero otherwise.

Theorem III.2. Let $Q \in \mathbb{N}^*$ and $\omega, k \in \mathbb{Z}/Q\mathbb{Z}$. If $\gcd(\omega, Q) = 1$, the global minima of $\mathcal{D}_{\omega,k}$ over S_Q are the pairs $(\mathbf{y}_A, \mathbf{y}_B) \in S_Q$ of the form $(\delta_{\alpha}, \delta_{\alpha+k})$ for every $\alpha \in \mathbb{Z}/Q\mathbb{Z}$.

We will apply the lemma below, proven in the Appendix.

Lemma III.3. Let P and Q be two coprime integers where $Q > 0$. For all $q, q' \in \mathbb{Z}/Q\mathbb{Z}$ where $q \neq q'$, $qP \bmod Q \neq q'P \bmod Q$.

Proof of Theorem III.2. Since $\mathcal{D}_{\omega,k}$ is clearly nonnegative, it suffices to identify all antecedents of zero by $\mathcal{D}_{\omega,k}$ in S_Q . The identity $\mathcal{D}_{\omega,k}(\mathbf{y}_A, \mathbf{y}_B) = 0$ implies $\hat{\mathbf{y}}_A[\omega] \hat{\mathbf{y}}_B[\omega]^* = \exp(2\pi i \omega k)$, which in turn implies $|\hat{\mathbf{y}}_A[\omega]| |\hat{\mathbf{y}}_B[\omega]| = 1$ after applying the complex modulus. Yet, by the triangular inequality and by definition of the DFT, $\mathbf{y}_A \in S_Q$ implies $|\hat{\mathbf{y}}_A[\omega]| \leq \|\mathbf{y}_A\|_1 = 1$. Therefore, $\mathcal{D}_{\omega,k}(\mathbf{y}_A, \mathbf{y}_B) = 0$ corresponds to an equality case of the triangular inequality: all nonzero terms in the DFT of \mathbf{y}_A at ω have the same phase. By contradiction, let us

¹Visit: <https://ethanhein.substack.com/p/identifying-modulations>

assume that there are two such nonzero terms, indexed by q and q' respectively. Since ω is coprime with Q , we may apply Lemma III.3 and deduce that $\omega q \bmod Q \neq \omega q' \bmod Q$: in other words, $\exp(2\pi i \omega q / Q)$ and $\exp(2\pi i \omega q' / Q)$ have different phases modulo 2π , which leads to a contradiction. Furthermore, let us assume that there are no nonzero terms: this implies $\mathbf{y}_A \notin S_Q$ and leads to another contradiction. Having excluded both of these contradictions, we deduce that there exists a unique nonzero term in the DFT of \mathbf{y}_A at ω : let us denote its index by $\alpha \in \mathbb{Z}/Q\mathbb{Z}$. Since $\mathbf{y}_A \in S_Q$, it follows that $\mathbf{y}_A[\alpha] = 1$ and $\mathbf{y}_A[q] = 0$ for $q \neq \alpha \bmod Q$. The same reasoning for \mathbf{y}_B leads us to claim that there exists $\beta \in \mathbb{Z}/Q\mathbb{Z}$ such that $\mathbf{y}_B[\beta] = 1$ and $\mathbf{y}_B[q] = 0$ for $q \neq \beta \bmod Q$. Thus $\mathbf{y}_A = \delta_\alpha$ and $\mathbf{y}_B = \delta_\beta$. We compute the following half-angle factorization:

$$\begin{aligned} \mathcal{D}_{\omega,k}(\delta_\alpha, \delta_\beta) &= \left| e^{2\pi i \omega k / Q} - \hat{\delta}_\alpha[\omega] \hat{\delta}_\beta[\omega]^* \right|^2 \\ &= \left| e^{2\pi i \omega k / Q} - e^{2\pi i \omega (\beta - \alpha) / Q} \right|^2 \\ &= 4 \left| \sin \left(2\pi \frac{\omega(k + \alpha - \beta)}{Q} \right) \right|^2. \end{aligned} \quad (6)$$

Solving $\mathcal{D}_{\omega,k}(\delta_\alpha, \delta_\beta) = 0$ leads to $\omega(k + \alpha - \beta) = 0 \bmod Q$, which we rewrite as $\omega(\beta - \alpha) = \omega k \bmod Q$. Since ω is coprime with Q , we may apply Lemma III.3 again and deduce that $\beta = (\alpha + k) \bmod Q$, which completes the proof. ■

Note. The condition $\gcd(\omega, Q) = 1$ is necessary to guarantee that $\mathcal{D}_{\omega,k}$ has a countable (*a fortiori*, finite) number of global minima. Assuming $\gcd(\omega, Q) = m$ with $m > 1$, and letting a parameter ρ take any real value between zero and one, we may observe that $\mathcal{D}_{\omega,k}(\rho \delta_\alpha + (1 - \rho) \delta_{\alpha+(Q/m)}, \delta_{\alpha+k}) = 0$.

Theorem III.4. *Let $Q \in \mathbb{N}^*$ and $\omega, k \in \mathbb{Z}/Q\mathbb{Z}$. If $\gcd(\omega, Q) = 1$, all local minima of the (ω, k) -discrepancy operator are global.*

Proof. Let the bilinear map \mathcal{B}_ω denote the CPSD operator at frequency ω : for any $\mathbf{y}_A, \mathbf{y}_B$, $\mathcal{B}_\omega(\mathbf{y}_A, \mathbf{y}_B) = \hat{\mathbf{y}}_A[\omega] \hat{\mathbf{y}}_B[\omega]^*$ (see Equation 4). We define $\mathcal{E}_{\omega,k} : z \in \mathbb{C} \mapsto \frac{1}{2} |e^{2\pi i \omega k / Q} - z|^2$ and decompose $\mathcal{D}_{\omega,k}$ as $(\mathcal{E}_{\omega,k} \circ \mathcal{B}_\omega)$. Up to a factor of one half, $\mathcal{E}_{\omega,k}$ is proportional to the squared Euclidean distance to the point $z_{\omega,k} = e^{2\pi i \omega k / Q}$ in the complex plane. Thus, $\mathcal{E}_{\omega,k}$ is strongly convex and has a single local minimum over \mathbb{C} : $\mathcal{E}_{\omega,k}(z_{\omega,k}) = 0$. Analogously to the proof of Theorem III.2, we apply Lemma III.3 and an equality case of the triangular inequality to show that the antecedents of $z_{\omega,k}$ by the bilinear map \mathcal{B}_ω are the pairs $(\mathbf{y}_A, \mathbf{y}_B) \in S_Q$ of the form $(\delta_\alpha, \delta_{\alpha+k})$ for every $\alpha \in \mathbb{Z}/Q\mathbb{Z}$. By Theorem III.2, all of these pairs are global minima for $\mathcal{D}_{\omega,k}$, which completes the proof. ■

Note. The result above may or may not translate to the loss function $\mathcal{L}_{\omega,k}$, depending on the choice of neural network architecture $\Psi_{\mathbf{w}}$ and initialization. If the Jacobian of $\mathbf{w} \mapsto (\Psi_{\mathbf{w}}(\mathbf{u}_A), \Psi_{\mathbf{w}}(\mathbf{u}_B))$ is singular at some $\mathbf{w} \in W$, the STONE gradient $\nabla \mathcal{L}_{\omega,k}$ may vanish even though the partial derivatives of $\mathcal{D}_{\omega,k}$ with respect to \mathbf{y}_A and \mathbf{y}_B are nonzero. To reduce the risk of a singular Jacobian, [12] have designed $\Psi_{\mathbf{w}}$ as a residual network (ResNet); more specifically, a ConvNeXT [15].

IV. INVARIANCE PROPERTIES

Definition IV.1. Let G be a group with identity element 1_G and X be a set. A (*left*) group action of G on X is a function $\mathcal{A} : G \times X \rightarrow X$ such that, for all $\mathbf{x} \in X$, $\mathcal{A}(1_G, \mathbf{x}) = \mathbf{x}$; and such that, for all $\mathbf{x} \in X$, for all $g, g' \in G$, $\mathcal{A}(g, \mathcal{A}(g', \mathbf{x})) = \mathcal{A}(gg', \mathbf{x})$.

Definition IV.2. A function $\Phi : X \rightarrow Y$ is *equivariant* to G if there exist group actions $\mathcal{A}_X : G \times X \rightarrow X$ and $\mathcal{A}_Y : G \times Y \rightarrow Y$ such that, for all $\mathbf{x} \in X$ and $g \in G$, $\Phi(\mathcal{A}_X(g, \mathbf{x})) = \mathcal{A}_Y(g, \Phi(\mathbf{x}))$.

Proposition IV.3. The key signature profile (KSP) is equivariant to circular frequency shifts in the constant- Q transform.

Proof. Let $X = \mathbb{R}^{JQ}$ and $Y = \mathbb{R}^Q$. Consider the group $G = \mathbb{Z}$ and the group action $\mathcal{A}_X : G \times X \rightarrow X$ such that:

$$\forall g \in G, \forall q \in \mathbb{Z}/JQ\mathbb{Z}, \quad \mathcal{A}_X(g, \mathbf{x})[q] = \mathbf{x}[(q - g) \bmod JQ]. \quad (7)$$

Similarly, consider the group action $\mathcal{A}_Y : G \times Y \rightarrow Y$ such that:

$$\forall g \in G, \forall q \in \mathbb{Z}/Q\mathbb{Z}, \quad \mathcal{A}_Y(g, \mathbf{y})[q] = \mathbf{y}[(q - g) \bmod Q]. \quad (8)$$

Consider the sets $C_q = \{jQ + q \bmod JQ \mid 0 \leq j < J\}$ for each $q \in \mathbb{Z}/Q\mathbb{Z}$. Definition II.2 rewrites as:

$$\forall \mathbf{x} \in \mathbb{R}^{JQ}, \forall q \in \mathbb{Z}/Q\mathbb{Z}, \quad \Phi(\mathbf{x})[q] = \frac{\prod_{n \in C_q} \exp \mathbf{x}[n]}{\prod_{q'=0}^{Q-1} \prod_{n' \in C_{q'}} \exp \mathbf{x}[n']}. \quad (9)$$

Given $g \in G$, consider $\mathcal{A}_G(g, C_q) = \{(n - g) \bmod JQ \mid n \in C_q\}$. By property of the Euclidean division: for every $n \in C_q$, there exists a unique $j \in \mathbb{Z}$ such that $n = jQ + q$. Thus:

$$\mathcal{A}_G(g, C_q) = \{(jQ + q - g) \bmod JQ \mid 0 \leq j < J\} = C_{(q-g)}. \quad (10)$$

Returning to Equation 9, a change of variables yields:

$$\begin{aligned} \mathcal{A}_Y(g, \Phi(\mathbf{x}))[q] &= \frac{\prod_{n \in C_{(q-g)}} \exp \mathbf{x}[n]}{\prod_{q'=0}^{Q-1} \prod_{n' \in C_{(q'-g)}} \exp \mathbf{x}[n']} \\ &= \frac{\prod_{n \in C_q} \exp \mathcal{A}_X(g, \mathbf{x})[n]}{\prod_{q'=0}^{Q-1} \prod_{n' \in C_{q'}} \exp \mathbf{x}[n']} \end{aligned} \quad (11)$$

Thus $\mathcal{A}_Y(g, \Phi(\mathbf{x})) = \Phi(\mathcal{A}_X(g, \mathbf{x}))$, concluding the proof. ■

Proposition IV.4. Given $Q \in \mathbb{N}^*$ and $\omega, k \in \mathbb{Z}/Q\mathbb{Z}$, consider the group $G = \mathbb{Z}$, the standard Q -simplex $Y = S_Q$ (Definition II.1), and the group action $\mathcal{A}_Y : G \times Y \rightarrow Y$ such that: $\forall g \in G, \forall q \in \mathbb{Z}/Q\mathbb{Z}, \mathcal{A}_Y(g, \mathbf{y})[q] = \mathbf{y}[(q - g) \bmod Q]$. For all $\mathbf{y}_A, \mathbf{y}_B \in S_Q$ and all $g_A, g_B \in G$, the (ω, k) -discrepancy (Definition II.3) satisfies:

$$\mathcal{D}_{\omega,k}(\mathcal{A}_Y(g_A, \mathbf{y}_A), \mathcal{A}_Y(g_B, \mathbf{y}_B)) = \mathcal{D}_{\omega, (k+g_A-g_B)}(\mathbf{y}_A, \mathbf{y}_B). \quad (12)$$

Proof. Let us denote the discrete Fourier transform (DFT) of $\mathcal{A}_Y(g_A, \mathbf{y}_A)$ by $\widehat{\mathcal{A}_Y(g_A, \mathbf{y}_A)}$. The DFT is equivariant to \mathcal{A}_Y : $\widehat{\mathcal{A}_Y(g_A, \mathbf{y}_A)}[\omega] = \exp(-2\pi i \omega g_A / Q) \widehat{\mathbf{y}}_A[\omega]$ and idem for $\widehat{\mathcal{A}_Y(g_B, \mathbf{y}_B)}$. Thus, by definition of (ω, k) -discrepancy:

$$\begin{aligned} \mathcal{D}_{\omega,k}(\mathcal{A}_Y(g_A, \mathbf{y}_A), \mathcal{A}_Y(g_B, \mathbf{y}_B)) &= \\ \frac{1}{2} \left| e^{2\pi i \omega k / Q} - e^{2\pi i \omega (g_B - g_A) / Q} \widehat{\mathbf{y}}_A[\omega] \widehat{\mathbf{y}}_B[\omega]^* \right|^2. \end{aligned} \quad (13)$$

After multiplying the subtraction by $\exp(2\pi i \omega (g_A - g_B)/Q)$, we recognize $\mathcal{D}_{\omega, (k+g_A-g_B)}(\mathbf{y}_A, \mathbf{y}_B)$ and conclude the proof. ■

Theorem IV.5. *The STONE loss function is invariant to circular frequency shifts in the constant- Q transform: given $Q \in \mathbb{N}^*$ and $\omega, k \in \mathbb{Z}/Q\mathbb{Z}$, for all $\mathbf{w} \in W$, $g \in \mathbb{Z}$, and $\mathbf{u}_A, \mathbf{u}_B \in \mathbb{R}^{JQ \times T}$,*

$$\mathcal{L}_{\omega, k}(\mathbf{w} \mid \mathcal{A}_U(g, \mathbf{u}_A), \mathcal{A}_U(g, \mathbf{u}_B)) = \mathcal{L}_{\omega, k}(\mathbf{w} \mid \mathbf{u}_A, \mathbf{u}_B), \quad (14)$$

where $\mathcal{A}_U(g, \mathbf{u}_A)[q, t] = \mathbf{u}_A[(q-g) \bmod JQ, t]$ is a cyclic group action of \mathbb{Z} over $\mathbb{R}^{JQ \times T}$ and idem for $\mathcal{A}_U(g, \mathbf{u}_B)$.

Proof. By Propositions IV.3 and IV.4 with $g_A = g_B = g$. ■

Note. The equivariance of $(\Phi \circ \Psi_{\mathbf{w}})$ offers the opportunity to include artificially transposed versions of \mathbf{u}_A and \mathbf{u}_B in the STONE loss, and derive a suitable value for the parameter k . For example, in [12], \mathbf{u}_A and \mathbf{u}_B are assumed to be in the same key and the modified STONE loss combines three terms:

$$\begin{aligned} \mathcal{L}_{\omega, k}^{\text{multi}}(\mathbf{w} \mid \mathbf{u}_A, \mathbf{u}_B) &= \mathcal{D}_{\omega, 0}((\Phi \circ \Psi_{\mathbf{w}})(\mathbf{u}_A), (\Phi \circ \Psi_{\mathbf{w}})(\mathbf{u}_B)) \\ &+ \mathcal{D}_{\omega, k}((\Phi \circ \Psi_{\mathbf{w}})(\mathbf{u}_A), (\Phi \circ \Psi_{\mathbf{w}})(\mathcal{A}_U(k, \mathbf{u}_A))) \\ &+ \mathcal{D}_{\omega, k}((\Phi \circ \Psi_{\mathbf{w}})(\mathbf{u}_B), (\Phi \circ \Psi_{\mathbf{w}})(\mathcal{A}_U(k, \mathbf{u}_A))). \end{aligned} \quad (15)$$

Clearly, $\mathcal{L}_{\omega, k}^{\text{multi}}$ satisfies the same invariance property as in Theorem IV.5. At the same time, empirical evidence suggests that the pretext task $\mathcal{L}_{\omega, k}^{\text{multi}}$ is less susceptible to trivial collapse (i.e., predicting a constant KSP) than $\mathcal{D}_{\omega, 0}$, while incurring a moderate (two-fold) increase in computational cost.

V. SYMMETRY-BREAKING PROPERTIES

We have seen in Theorems III.2 and III.4 that all local minima of $\mathcal{D}_{\omega, k}$ are global minima. Moreover, we have characterized them as pairs of one-hot vectors, which are aperiodic patterns. Yet, by Proposition IV.4, $\mathcal{L}_{\omega, k}$ is invariant to circular pitch shifts. In physics, this phenomenon is known as spontaneous symmetry breaking (SSB) [16]: the stable states of a system—namely, $(\delta_{\alpha}, \delta_{\alpha+k})$ for $\alpha \in \mathbb{Z}/Q\mathbb{Z}$ —are not symmetric (i.e., invariant) under a symmetry of the underlying “Hamiltonian” $\mathcal{D}_{\omega, k}$, understood as a function of $S_Q \times S_Q$.

We propose to elaborate on this observation by studying the value of $\mathcal{D}_{\omega, k}$ over pairs of periodic signals $(\mathbf{y}_A, \mathbf{y}_B)$, which are left unchanged by the action of proper subgroups of $G_Q = \mathbb{Z}/Q\mathbb{Z}$. Since G_Q is cyclic of order Q , these subgroups are of the form G_P with $P \mid Q$; i.e., P is a divisor of Q .

Note. In music theory, the binary vectors $\mathbf{y} \in \{0, 1\}^Q$ which are fixed points of $\mathbf{y} \mapsto \mathcal{A}_Y(P, \mathbf{y})$ for some $P \mid Q$ are called “modes of limited transposition” [17]: e.g., tritone, augmented triad, diminished seventh chord, whole-tone scale, and more.

The theorem below shows that the STONE loss is at an unstable equilibrium if (but not only if) its inputs are periodic. For simplicity, we only discuss the case where both elements in the pair are equal ($\mathbf{y}_A = \mathbf{y}_B = \mathbf{y}$) and we have set $k = 0$.

Theorem V.1. *Given $Q \in \mathbb{N}^*$, $\omega \in \mathbb{Z}/Q\mathbb{Z}$ such that $\gcd(\omega, Q) = 1$, and $G = \mathbb{Z}$, consider the standard simplex $Y = S_Q$ (Definition II.1), the group action $\mathcal{A}_Y : G \times Y \rightarrow Y$ (Definition IV.1), and the $(\omega, 0)$ -discrepancy $\mathcal{D}_{\omega, 0}$ (Definition II.3). For any $\mathbf{y} \in S_Q$, if \mathbf{y} is periodic—i.e., if there exists $P \mid Q$,*

$1 \leq P < Q$ such that $\mathbf{y} = \mathcal{A}_Y(P, \mathbf{y})$ —then the gradient of “self-discrepancy” $\mathbf{y} \mapsto \mathcal{D}_{\omega, 0}(\mathbf{y}, \mathbf{y})$ vanishes at \mathbf{y} : $\nabla \mathcal{D}_{\omega, 0}(\mathbf{y}, \mathbf{y}) = 0$.

We will apply the lemma below, proven in the Appendix.

Lemma V.2. *Given $Q \in \mathbb{N}^*$ and $\omega \in \mathbb{Z}/Q\mathbb{Z}$, for all $\mathbf{y} \in S_Q$,*

$$\mathcal{D}_{\omega, 0}(\mathbf{y}, \mathbf{y}) \leq \frac{1}{2}. \quad (16)$$

Proof of Theorem V.1. Let us factorize Q as KP where P is the period of \mathbf{y} . For each $q \in \mathbb{Z}/Q\mathbb{Z}$, we perform the Euclidean division: $q = kP + r$ where $0 \leq k \leq K$ and $0 \leq r < P$. Denoting by $\mathcal{A}_Y^{(k)}$ the group action \mathcal{A}_Y applied iteratively k times, the DFT of \mathbf{y} at frequency ω rewrites as a double sum:

$$\begin{aligned} \hat{\mathbf{y}}[\omega] &= \sum_{k=0}^{K-1} \sum_{r=0}^{P-1} \mathcal{A}_Y^{(k)}(-P, \mathbf{y})[r] \exp\left(-2\pi i \frac{\omega(kP+r)}{Q}\right) \\ &= \left(\sum_{k=0}^{K-1} e^{-2\pi i \omega k P / Q} \right) \times \left(\sum_{r=0}^{P-1} \mathbf{y}[r] e^{-2\pi i \omega r / Q} \right). \end{aligned} \quad (17)$$

We recognize the first factor as the sum of K^{th} complex roots of unity. Since K is integer, this sum equals zero. Equating $\hat{\mathbf{y}}[\omega]$ to zero in Definition II.3 yields $\mathcal{D}_{\omega, k}(\mathbf{y}, \mathbf{y}) = \frac{1}{2}$. By Lemma V.2, this value is a global maximum for the application $\mathbf{y} \mapsto \mathcal{D}_{\omega, 0}(\mathbf{y}, \mathbf{y})$ over S_Q . Since this application is differentiable (see proof of Proposition II.5), its gradient vanishes at local extrema, and *a fortiori* at global maxima; which concludes the proof. ■

To recap, stochastic optimization with STONE is spontaneously symmetry-breaking because symmetric patterns have the highest self-discrepancy (Theorem V.1) while one-hot vectors have the lowest self-discrepancy (Theorem III.2).

VI. CONCLUSION

STONE, which stands for self-supervised tonality estimation, is a “hybrid deep learning” model [18]: what it lacks in human supervision, it makes up for in mathematical structure [12]. In this article, we have seen that the interpretability of learned STONE embeddings as key signature profiles (KSP) does not emerge spontaneously from large-scale training: it also depends upon a pretext task design which is informed by signal processing (CQT, CPSD), deep learning (ResNets, softmax), and basic music theory (octave equivalence and, optionally, the circle of fifths). Note that STONE is not strictly tied to the CQT: in principle, it generalizes to any filterbank with center frequencies in geometric progression: e.g., SWIPE kernels [19].

Beyond its usefulness in practice, STONE is the first proof of feasibility of a machine learning system which distinguishes key signatures from polyphonic audio without any prior concept of root, bass, triad, scale degree, or fundamental line. At the same time, its publication has led to further developments: for example, a recent paper has adapted the CPSD objective to a multi-class-token Transformer architecture for multitask SSL [20]. Lastly, from the perspective of music perception and cognition, the STONE model qualifies as “illiterate” (machine) listening [21], i.e., as remote from linguistic or symbolic determinations such as pitch class names. We leave the implications of this remark to future work.

REFERENCES

- [1] R. Kondor and S. Trivedi, "On the generalization of equivariance and convolution in neural networks to the action of compact groups," in *Proceedings of the International Conference on Machine Learning (ICML)*. PMLR, 2018, pp. 2747–2755.
- [2] C. Schörkhuber and A. Klapuri, "Constant- q transform toolbox for music processing," in *Proceedings of the Sound and Music Computing (SMC) Conference*, 2010.
- [3] M. Müller, *Fundamentals of music processing: Audio, analysis algorithms, applications*. Springer, 2015.
- [4] V. Lostanlen and S. Mallat, "Wavelet scattering on the pitch spiral," in *Proceedings of the Digital Audio Effects (DAFX) Conference*, 2015.
- [5] J.-F. Ducher and P. Esling, "Folded CQT RCNN for real-time recognition of instrument playing techniques," in *Proceedings of the International Society for Music Information Retrieval (ISMIR) Conference*, 2019.
- [6] F. Korzeniowski and G. Widmer, "Feature learning for chord recognition: The deep chroma extractor," 2016.
- [7] V. Lostanlen and C.-E. Cella, "Deep convolutional networks on the pitch spiral for musical instrument recognition," in *Proceedings of the International Society for Music Information Retrieval (ISMIR) Conference*, 2016.
- [8] E. J. Humphrey, T. Cho, and J. P. Bello, "Learning a robust Tonnetz-space transform for automatic chord recognition," in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2012, pp. 453–456.
- [9] V. Lostanlen, S. Sridhar, B. McFee, A. Farnsworth, and J. P. Bello, "Learning the helix topology of musical pitch," in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020.
- [10] S. Sridhar and V. Lostanlen, "Helicality: An isomap-based measure of octave equivalence in audio data," in *Proceedings of the International Society for Music Information Retrieval Late-Breaking/Demo Session (ISMIR-LBD)*, 2020.
- [11] K. Choi, G. Fazekas, K. Cho, and M. Sandler, "A tutorial on deep learning for music information retrieval," *arXiv:1709.04396*, 2017.
- [12] Y. Kong, V. Lostanlen, G. Meseguer-Brocal, S. Wong, M. Lagrange, and R. Hennequin, "STONE: Self-supervised tonality estimator," in *Proceedings of the International Society for Music Information Retrieval (ISMIR) Conference*, 2024.
- [13] A. Riou, S. Lattner, G. Hadjeres, and G. Peeters, "PESTO: Pitch estimation with self-supervised transposition-equivariant objective," in *Proceedings of the International Society for Music Information Retrieval (ISMIR) Conference*, 2023.
- [14] Y. Kong, G. Meseguer-Brocal, V. Lostanlen, M. Lagrange, and R. Hennequin, "S-KEY: Self-supervised learning of major and minor keys from audio," in *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2025.
- [15] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, "A convnet for the 2020s," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 11 976–11 986.
- [16] A. Beekman, L. Rademaker, and J. Van Wezel, "An introduction to spontaneous symmetry breaking," *SciPost Physics Lecture Notes*, p. 011, 2019.
- [17] A. Barate and L. A. Ludovico, "Generalizing Messiaen's modes of limited transposition to a n -tone equal temperament," in *Proceedings of the International Conference on Sound and Music Computing (SMC)*, 2015, pp. 287–293.
- [18] G. Richard, V. Lostanlen, Y.-H. Yang, and M. Müller, "Model-based deep learning for music information research: Leveraging diverse knowledge sources to enhance explainability, controllability, and resource efficiency," *IEEE Signal Processing Magazine*, vol. 41, no. 6, pp. 51–59, 2025.
- [19] D. Marttila and J. D. Reiss, "Improving neural pitch estimation with SWIPE kernels," in *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, 2025.
- [20] Y. Kong, V. Lostanlen, R. Hennequin, M. Lagrange, and G. Meseguer-Brocal, "Multi-class-token transformer for multitask self-supervised music information retrieval," in *Proceedings of the IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, 2025.
- [21] H. Honing, *The illiterate listener: On music cognition, musicality and methodology*. Amsterdam University Press, 2011.

VII. APPENDIX

We prove Lemma III.3 in two fashions: first, with elementary modular arithmetic; second, with abstract algebra.

Proof (elementary). Given $q, q' \in \mathbb{Z}/Q\mathbb{Z}$, suppose that $qP \bmod Q = q'P \bmod Q$, i.e., that $(q - q')P = 0 \bmod Q$. Since P and Q are coprime, it follows that Q divides $(q - q')$. Yet, since $0 \leq q, q' < Q$, $-Q < (q - q') < Q$. The only multiple of Q in this range is zero; thus, $q = q'$. We conclude by contraposition. ■

Proof (advanced). Since $\gcd(P, Q) = 1$, the coset $P + Q\mathbb{Z}$ is a unit in the quotient ring $\mathbb{Z}/Q\mathbb{Z}$; i.e., $P + Q\mathbb{Z} \in (\mathbb{Z}/Q\mathbb{Z})^\times$. Therefore, $q + Q\mathbb{Z} \mapsto qP + Q\mathbb{Z}$ is a ring automorphism. ■

Proof of Lemma V.2. Definition II.3 with $\mathbf{y}_A = \mathbf{y}_B = \mathbf{y}$ yields:

$$\mathcal{D}_{\omega,0}(\mathbf{y}, \mathbf{y}) = \frac{1}{2} \left(1 - |\hat{\mathbf{y}}[\omega]|^2 \right)^2. \quad (18)$$

By the triangular inequality on the DFT: $|\hat{\mathbf{y}}[\omega]| \leq \|\mathbf{y}\|_1 \leq 1$. Thus: $0 \leq (1 - |\hat{\mathbf{y}}[\omega]|^2) \leq 1$, concluding the proof. ■