
Towards Isotropic Deep Learning

A New Default Inductive Bias

Anonymous Author(s)

Affiliation

Address

email

Abstract

This position paper explores the often overlooked geometric implications of current functional forms used in deep learning and proposes a new paradigm to be termed *Isotropic Deep Learning*. It has been demonstrated that the functional forms of current deep learning influence the activation distributions. Through training, broken symmetries in functional forms can induce broken symmetries in embedded representations. Thus producing a geometric artefact in representations which is not task-necessitated, and solely due to human-imposed choices of functional forms. There appears to be no strong a priori justification for why such a representation or functional form is desirable, while this paper proposes several detrimental effects of the current formulation. As a result, a modified framework for functional forms will be explored with the goal of unconstraining representations by elevating a rotational symmetry-inductive bias throughout the network. This framework is encouraged to be adopted as a new default. This direction is proposed to improve network performance. Preliminary functions are proposed, including activation functions. Since this overhauls almost all functional forms characterising modern deep learning, it is suggested that this shift may constitute a new branch of deep learning.

1 Introduction

Current deep learning models typically employ an elementwise functional form. This is particularly evident in activation functions, sometimes called ridged activation functions. These functions are often displayed univariately as shown in Eqn. 1, with σ being a specific activation function implemented, e.g. ReLU, Tanh, etc.

$$f : \mathbb{R} \rightarrow \mathbb{R}, \quad x \mapsto f(x) = \sigma(x) \quad (1)$$

However, this display choice obfuscates a crucial (standard) basis dependence. This is explicitly displayed in, what should be considered a more implementation-correct multivariate form, of Eqn. 2. This reveals the functional form's usually hidden \hat{e}_i basis dependence. The multivariate form is depicted for an n neuron layer, with activation vector $\vec{x} \in \mathbb{R}^n$. This standard basis dependence is arbitrary and appears as a historical precedent, rather than appropriate inductive bias, discussed further in App A.

$$\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad \vec{x} \mapsto \mathbf{f}(\vec{x}) = \sum_{i=1}^n \sigma(\vec{x} \cdot \hat{e}_i) \hat{e}_i \quad (2)$$

Due to this basis dependence, non-linear transformations differ angularly in effect Bird [2025]. Therefore, this will be termed an *anisotropic function*, indicating this rotational asymmetry. Due to the pervasive use of these functional forms, including optimisers, normalisers, regularisers, activation functions etc., current deep learning as a paradigm may consequently be termed *anisotropic deep*

learning. Despite its implications, this *choice* of basis-dependent anisotropy appears unappreciated and incidental in the development of most contemporary models, with anisotropic forms are almost treated as axiomatic to deep learning rather than a considered choice. Hence, re-evaluating and systematically reformulating this foundational aspect of modern deep learning, with such wide-reaching consequences, is felt to constitute a new branch: *Isotropic Deep Learning*.

This asymmetry in the non-linear transform is particularly about the standard (Kronecker) basis vectors and their negative $\{+\hat{e}_i, -\hat{e}_i\}_{\forall i}$, due to the elementwise nature. Therefore, it can be said to distinguish the standard basis - a '*distinguished basis*'¹. For example, the standard basis is visible in Fig. 1 showing the mapping of elementwise-tanh on a variety of test shapes. Most generally, this



Figure 1: [Insert caption here]

choice of functional forms can be considered to break *continuous* rotational symmetry, and reduce it to a *discrete* rotational symmetry — a permutation symmetry of the standard basis. In effect, if the function is treated in its multivariate form, it is equivariant to a permutation of the components of its vector decomposed in the standard basis. For an element of the permutation group $\mathbf{P} \in \mathcal{S}_n$, the following equivariance relation holds $f(\mathbf{P}\vec{x}) = \mathbf{P}f(\vec{x})$.

Non-linearities are usually pivotal to the network’s ability to achieve a desired computation, as seen through the universal approximation theorem’s Cybenko [1989] explicit dependence on the form of the activation function Hornik [1991]. The non-linearities produce differing local transformations, such as stretching, compressing, and generally reshaping a manifold. Consequently, the network may be expected to adapt by moving representations to geometries about these distinguished directions, to use specific local transforms to achieve the desired computation. Hence, an anisotropy about a distinguished basis is induced into the activation distribution shown by.

This has been empirically demonstrated by Bird [2025]: training results in the broken symmetry of the functional forms, inducing a broken symmetry in the activations. Since these non-linear zones are centred around the distinguished bases, the embedded representations are expected to move to beneficial angular arrangements about the arbitrarily imposed privileged basis’s geometry. For example, they appear to move towards the non-linearities’ extremums, aligned, anti-aligned or other geometries, through training. This may correspond to a local, dense or sparse codings and superposition, respectively.

Therefore, the network has adapted its representations through training, for probably various optimisation reasons, due to these functional form choices. The causal hypothesis aids in explaining the observed tendency of privileged-basis alignment. This is the causal hypothesis underlying the author’s position: **functional forms should be a deliberate and considered decision, with a suitably optimal, and minimally harmful, default**. Currently, anisotropy results in a human-caused representational collapse onto the privileged basis, and it is suggested that this is frequently not a task-necessitated collapse. There appears to be little justification as to why this is universally

¹This is suggested generalisation from a '*privileged basis*' discussed in Elhage et al. [2022]. The change to '*distinguished basis*' reflects that the basis may be more-or-less aligned to the representation; whereas '*privileged basis*' is felt to suggest a basis which is more aligned to the embedded activations. The term 'basis' will be retained even though the set of '*distinguished vectors*' may also be under-/over complete for spanning the whole space, as demonstrated in Bird [2025]. There may be multiple distinguished bases, such as aligned and anti-aligned with the standard basis for the model.

desirable, with several key negative implications discussed in *Sec. 2*. Without a priori justification, this inductive bias may be detrimental to computation, so unconstraining the activation appears generally preferable.

Overall, historic and frequently overlooked functional forms for modern deep learning directly influence the models’ activations and therefore behaviour. There exists a functional form basis dependence which appears entirely arbitrary, obfuscated and hence neglected. Yet, a causal link between this arbitrary basis and activations has been empirically demonstrated, and hence, a resultant effect on the final performance of the model is hypothesised. These are often underappreciated choices which have consequences and should be well-justified and studied. This is the position of the authors.

Throughout the rest of this position paper, **it is argued that a departure from this anisotropic functional form paradigm towards the isotropic paradigm is generally preferable as an inductive bias**, unless otherwise justified. It encourages the reader to be conscious of these choices when designing a model, as well as the usual architectural tool kit. Particularly, isotropic choices, equivalent to basis independence, may be thought to unconstrain the representations into more optimal arrangements for a task. The tenets of this paradigm are suggested for all architectures on general tasks.

2 Hypothesised Problems of Anisotropy

This section argues that current functional forms impose unintended anisotropic performance detriments, indicating that *Isotropic Deep Learning*, once substantially developed, is proposed as the default inductive bias unless an alternative is task-necessitated.

This section lays out a non-exhaustive set of arguments discussing some implications that anisotropic functional forms may cause. These mainly centre on the role of the activation functions, since this is the area the author has primarily explored in their PhD thus far. To the author’s knowledge, some of these failure modes are newly characterised phenomena, such as the so-called ‘*neural refractive problem*’.

Furthermore, there are some specific instances when anisotropy may be detrimental to performance, one of those is in the self-attention step of transformers, speculated in *App. E*.

2.1 The Neural Refractive Problem

The ‘*neural refractive problem*’ describes how linear trajectories of activations may converge or diverge from their initial consistent path after an activation function is applied. This is analogous to a light ray refracting through an optically varying medium or boundary.

It appears to be a phenomenon in all anisotropic activation functions to date. The ‘refraction’ typically occurs more significantly at larger magnitudes — potentially a failure mode under network extrapolation. Mathematically, this has several representations, a magnitude-varying ‘dynamic refraction’ shown in *Eqn. 3* or differentially in *Eqn. 4*. Also defined is a ‘static refraction’ definition shown in *Eqn. 5*. These are described for a multivariate activation function \mathbf{f} and vector $\vec{x} = \alpha \hat{x}$ where \hat{x} is a unit vector. This relation may be satisfied for a single direction, a subset of the space or all directions $\hat{x} \in \mathcal{X} \subseteq S^n$. The relations generally show how the activation function alters the direction of its input vector in an anisotropic manner.

$$\exists \hat{x} \in S^{n-1}, \exists \alpha_1 \neq \alpha_2 > 0 : \frac{\mathbf{f}(\alpha_1 \hat{x})}{\|\mathbf{f}(\alpha_1 \hat{x})\|} \neq \frac{\mathbf{f}(\alpha_2 \hat{x})}{\|\mathbf{f}(\alpha_2 \hat{x})\|} \quad (3)$$

$$\exists \hat{x} \in S^{n-1}, \exists \alpha_0 : \left. \frac{\partial}{\partial \alpha} \frac{\mathbf{f}(\alpha \hat{x})}{\|\mathbf{f}(\alpha \hat{x})\|} \right|_{\alpha_0} \neq \vec{0} \quad (4)$$

$$\exists \hat{x} \in S^{n-1}, \exists \alpha : \frac{\mathbf{f}(\alpha \hat{x})}{\|\mathbf{f}(\alpha \hat{x})\|} \neq \hat{x} \quad (5)$$

It can be seen that along a straight-line trajectory in direction \hat{x} , the result of the activation function is a curved line if dynamically refracted. Therefore, if the linear feature hypothesis is followed, *every* linear feature, in these directions, becomes curved following the activation function. The network may

exploit some of this curvature to construct new linear features in the subsequent layers; however, there may be many instances where this curvature is detrimental to established semantics. The network may lose semantic separability, produce magnitude-based semantic inconsistency or compensatory maladaptations in later layers, due to this refraction. This may hinder network performance and is resolved by isotropic choices.

Particularly detrimental, in both refraction cases, can be the loss of semantic separability. If two distinct trajectories, representing different semantics, are transformed into curves which intersect or converge, then the separability of these concepts is lost. For example, suppose one direction is a linear feature for the presence of a dog in an image, whilst the other is for a horse. In that case, if these activations are of a magnitude where the activation function causes convergence, the identity of the activation’s meaning can be misrepresented.

This may be particularly consequential for functions such as Sigmoid and Tanh, since large magnitude inputs end up at particular limit points (discussed as trivial representational alignments in Bird [2025]). For example, Tanh produces the limit points shown in Eqn. 6 when $\hat{x} \cdot \hat{e}_i \neq 0$ for all i . If there exists an i , s.t. $\hat{x} \cdot \hat{e}_i = 0$, then the corresponding index also has a 0 as a limit point. A fully-connected layer can only effectively separate two such converging directions at a time, which are then further curved by a subsequent activation function.

$$\lim_{\alpha \rightarrow \infty} \mathbf{f}(\alpha \hat{x}) = \sum_{i=1}^N \tanh(\alpha \hat{x} \cdot \hat{e}_i) \hat{e}_i \approx \sum_{i=1}^N \pm \hat{e}_i = (\pm 1, \dots, \pm 1)^T \quad (6)$$

Consequently, semantic separability is lost except for 3^n discrete limit points for Tanh and Sigmoid. Therefore, embedded activations may be expected to align with these limit points, an empirically observed tendency Bird [2025]. Similarly, ReLU has one distinct limit point, $\vec{0}$, but otherwise an orthant unaffected by neural refraction. The authors speculate whether this is an additional reason for the success of ReLU, due to only a subset of directions experiencing the neural refraction phenomena. Furthermore, this would suggest an advantage of Leaky-ReLU: despite featuring static-refraction, directions do not become overlapped, so semantic separability is retained. Otherwise, the network may expend training time on producing robust semantic separability, a needless compensatory adaptation, which may lower representational capacity or extend training as a result.

More generally, dynamic deflection of trajectories may cause semantic ambiguity for the network, where only samples interpolable from training samples are reliably semantically identifiable. Particularly, *the larger the deflection, the greater the semantic ambiguity expected*. More considerable deflections typically occur at larger magnitudes in many current functions. Therefore, a magnitude-dependent semantic inconsistency may arise due to such deflections. A deflection function can be a trivial diagnostic measure defined by Eqn. 7 for a particular activation function.

$$\theta(\alpha; \hat{x}, \mathbf{f}) = \arccos \left(\frac{\mathbf{f}(\alpha \hat{x}) \cdot \hat{x}}{\|\mathbf{f}(\alpha \hat{x})\|} \right) \quad (7)$$

This may explain why the network may perform excessively poorly on out-of-training-distribution samples. For example, suppose a linear feature roughly represents the quantity of cows in a field. In that case, the network may fail to extrapolate its function when an anomalous amount of cows are present, as this would be a very large magnitude of the linear feature. Therefore, the deflection is unprecedented and becomes uninterpretable. The activation function would result in a loss of semantic consistency. Consequently, a network seeking to preserve linear features may constrain activation magnitudes, through training, to regions where the non-linear response is approximately predictable and stable to avoid the damaging consequences of neural refractions.

Angular anisotropies fundamentally cause the refraction phenomenon. If compression and rarefaction of certain angular regions occur, linear features will be deflected in various ways. A fix for this is introducing isotropy — the initial motivation for developing the paradigm. This does not prevent compression and rarefaction of activation distributions in general, as a bias can be added to reintroduce these useful phenomena predictably. It is argued that these are only an issue when they affect linear, not affine, features in a potentially unpredictable and thus semantically uninterpretable way.

The phenomenon is eliminated from networks by rearranging Eqn. 5 shown in Eqn. 8, then applying the simplification $\|\mathbf{f}(\alpha \hat{x})\| = \sigma(\alpha)$ in Eqn. 9.

$$\mathbf{f}(\alpha \hat{x}) = \|\mathbf{f}(\alpha \hat{x})\| \hat{x}' \quad (8)$$

161

$$\mathbf{f}(\alpha \hat{x}) = \sigma(\alpha) \hat{x}' \quad (9)$$

162 Finally choosing $\hat{x}' = \mathbf{R}\hat{x}$ for isotropy and $\mathbf{R}\hat{x} = \mathbf{I}_n \hat{x} = \hat{x}$ for simplicity, shown in Eqn. 10.

$$\mathbf{f}(\alpha \hat{x}) = \sigma(\alpha) \hat{x} \quad (10)$$

163 In standard notation, Eqn. 10 can be rewritten into the *final functional form for isotropic activation*
164 *functions* shown in Eqn. 11.

$$\mathbf{f}(\vec{x}) = \sigma(\|\vec{x}\|) \hat{x} \quad (11)$$

165 This can be generalised as a result of rotational equivariance of the function. This can be expressed
166 as a condition in Eqn. 12, which uses a commutator bracket for convenience, with $\forall \mathbf{R} \in \text{SO}(n)$.
167 This bracket can be used to similarly define the current permutation-anisotropic paradigm, by using
168 the transform $\forall \mathbf{P} \in \mathcal{S}_n$ instead of the rotation. This may be recognised as superficially similar to
169 equivariant neural networks, due to an analogous equivariance relation; however, the differences in
170 both implementation and motivations are discussed in App ??.

$$[\mathbf{R}, \mathbf{f}] = (\mathbf{R}\mathbf{f} - \mathbf{R}\mathbf{f}) = \vec{0} \quad (12)$$

171 The relation may be more familiar as $\mathbf{f}(\mathbf{R}\vec{x}) = \mathbf{R}\mathbf{f}(\vec{x})$. This relation only applies to single-
172 argument functions and requires generalising to more circumstances. A preliminary condition may
173 be $\mathbf{f}(\mathbf{R}\vec{a}_1, \dots, \mathbf{R}\vec{a}_N) = \mathbf{R}\mathbf{f}(\vec{a}_1, \dots, \vec{a}_N)$ for $\mathbf{f} : \bigotimes_N \mathbb{R}^n \rightarrow \mathbb{R}^n$.

174 This introduces the general isotropic functional form for activation functions given in Eqn 13. This
175 should be a piecewise function, defined using the identity at $\vec{x} = \vec{0}$, but this is suppressed for
176 simplicity. This functional form is $\mathcal{O}(n)$ time for \mathbb{R}^n . Future work is establishing a universal
177 approximation theorem for this functional form, as this is ongoing research for the author's PhD.

$$\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad \vec{x} \mapsto \mathbf{f}(\vec{x}) = \sigma(\|\vec{x}\|) \hat{x} \quad (13)$$

178 This is not to be confused with the radial-basis functional form displayed in Eqn. 14.

$$\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad \vec{x} \mapsto \mathbf{f}(\vec{x}) = \sum_{i=1}^N \sigma(\|\vec{x} - \vec{c}_i\|) \hat{e}_i \quad (14)$$

179 The ‘neural refractive problem’ outlines how semantic meanings may become intertwined or am-
180 biguous due to current functional forms consistently skewing linear features in undesirable ways. A
181 hypothesis is developed that this may be especially detrimental for out-of-distribution activations,
182 which are likely to be most deflected and hence most semantically corrupted. Thus, the network’s
183 generalisation then fails excessively. Finally, it may be expected that the network produces com-
184 pensatory adaptations for the phenomena. Since neural refraction is a non-linear and anisotropic
185 phenomenon, it cannot be inverted by a single subsequent layer, potentially wasting training time on
186 corrections to the activation distributions due to unintended refraction. These refraction corrections
187 may be limited in scope and produce maladaptive behaviour outside of the original training data.

188 2.2 Weight Locking, Optimisation Barriers and Disconnected Basins

189 ‘Weight locking’ is a term to describe how particularly the weight parameter² may suffer from being
190 stuck in local-minima found further into loss valleys encountered after a sufficient amount of training.
191 This arises only due to the anisotropic functional form’s *discrete* permutation symmetry.

192 Qualitatively, this is because the semantically meaningful linear features tend to become aligned
193 with geometric positions about the distinguished bases. Any small perturbation to a parameter
194 may misalign activations to the network’s existing ‘understood’ semantics. In effect, further small
195 perturbations to the parameters may move activations from a semantically-aligned to a semantically-
196 dislocated state, thus the activation’s meaning becomes ambiguous, forfeiting performance and
197 resulting in a ‘false’ local-minima created only by the *discrete* nature of the permutation symmetry.
198 It may also suggest that the optimisation barrier may be some function of the angular separation of
199 the semantic directions. Consequently, creating a plethora of architectural local minima in the space.
200 Only sufficiently large perturbations to a parameter may move activations between two differing
201 semantically aligned directions.

²Though similarly applies to a ‘locking’ of the bias to $\vec{0}$.

202 This semantic dislocation is an emergent consequence of breaking the continuous symmetric forms,
 203 as without continuous rotational symmetry, it results in the dual phenomena of connectivity of many
 204 minima basins being lost. In effect, enforcing the isotropy constraints results in sets of continuously
 205 connected local minima which can be smoothly transformed into one another, by corresponding
 206 parameter rotations shown in Eqn. 15, a consequence of Eqn. 11. If this is downgraded to discrete
 207 rotational symmetry (i.e. permutation symmetry), then artificial optimisation barriers may reemerge
 208 in these basins. In this case, only a sufficiently large perturbation to the parameters may dislodge the
 209 network into a more optimal minima, while more minor perturbations are insufficient. Effectively, the
 210 discrete permutation symmetry may result in a discretised lattice solution for the parameters, much
 211 like how it breaks the symmetry of activations through training too Bird [2025].

$$\forall \mathbf{R} \in \text{SO}(n) : \underbrace{\mathbf{W}^l \mathbf{R}^\top}_{\mathbf{W}^{l'}} \mathbf{f} \left(\underbrace{\mathbf{R} \mathbf{W}^{l-1}}_{\mathbf{W}^{l-1}} \vec{x} + \underbrace{\mathbf{R} \vec{b}}_{\vec{b}'} \right) = \mathbf{W}^l \mathbf{f} \left(\mathbf{W}^{l-1} \vec{x} + \vec{b} \right) \quad (15)$$

212 This exists as a qualitative intuition since, until robust methods to determine semantically meaningful
 213 directions are produced, this hypothesis remains difficult to verify. Nevertheless, steps can be taken
 214 immediately to counteract the problem, and this is to introduce isotropy to connect these minima.

215 2.3 Emergence of Linear Features and Semantic Interpolatability

216 As previously mentioned, symmetry-broken functional forms induce symmetry-broken representa-
 217 tions. Thus, *approximately* discrete embedding directions are tended towards. It may be conjectured
 218 that semantically meaningful linear directions may also be encouraged to discretise, aligning with
 219 these anisotropic embeddings. This generally appears to be the case, with notable counterexamples.
 220 Moreover, these counterexamples are for networks which *do not feature anisotropic functional forms*.

221 However, many real-life semantics are continuums: colours, positions of objects, broad morphology,
 222 even within a single species. Representational collapse onto a single discrete semantic may lose this
 223 important nuance. It appears a poor inductive bias to have functional forms encourage discretised
 224 representations. Isotropic functions do not prevent discrete semantics, but they don't encourage them
 225 either, enabling continuous ones since they generally unconstrain the representations. Therefore,
 226 moving towards isotropy is hypothesised to encourage embeddings to be more smoothly distributed,
 227 taking on intermediate values between typically discrete linear features and substantially enlarging
 228 the expressivity and representation capacity of networks — only limited by concept interferences.

229 In this case, the discrete concept of representation capacity may become irrelevant; each layer may
 230 express different continuous arrangements, where differing concepts are angularly suppressed and
 231 expressed in analogy to the linear features hypothesis. Instead, the '*magnitude-direction hypothesis*' is
 232 proposed as a continuous extension, magnitudes indicating the amount of stimulus present, direction
 233 indicating the particular concept. Activations then populate this more continuous manifold.

234 This may also produce a better organised semantic map at each layer of the network since intermediate
 235 representations may now relate otherwise discrete features, and therefore, this connectivity can bring
 236 them into continuous proximity (which 'weight locking' may typically prevent). This may additionally
 237 aid researchers in comparing representational alignment between models and biology³.

238 Therefore, in general settings, the inductive bias of isotropy appears more appropriate as a default
 239 unless justification for anisotropy is present. It is an inductive bias that meaningful semantics are
 240 often continuous and interpolatable, while retaining discrete semantics when task necessitated as
 241 opposed to human imposition. Hence, it generalises the discrete linear features paradigm into a more
 242 continuous setting. It suggests that a continuous rotational symmetry allows a continuous embedding,
 243 since functional forms do not induce direction-based symmetry breaking. It is hoped this will enable
 244 networks to acquire a more optimal and natural representation embedding for the given task.

³Though there is no guarantee that *all* basins are connected, so therefore would not necessarily be alignable through continuous rotation transforms. Furthermore, the success of ensemble methods with diverse constituent models would suggest that diverse disconnected basins remain, complicating representational alignment even under isotropy conditions, though isotropy may help somewhat.

3 Alternative View: Embedding Folding

4 Isotropic Implementations

This position paper argues for the implementation of isotropic functional forms into neural networks as a default inductive bias. Near-term adoption may be rate-limited due to the development of suitable functions, particularly since *anisotropic deep learning* has a substantial head start and analogues to existing functions are not so trivial to produce. Despite this, several preliminary implementations are outlined in this section as a starting point; however, these functions are far from optimal and substantial research and development are required to bring isotropic deep learning into practicality.

Nevertheless, below is a non-exhaustive list of activation functions and some of their consequences. A brief summary of other functions, including optimisers, regularisers, and normalisers, is also discussed. It is hoped that a directed search for new functions may occur.

4.1 Activation Functions

As stated, the isotropic functional form for activation functions is given in *Eqn. 13*. In *Tab. 4.1*, it is compared with the other common functional forms. It is clear in this comparison that the isotropic functional form is basis-independent and relatively simple. Further criteria, in addition to isotropy, are also a performance necessity, but will be outlined in future work.

Beginning from this functional form, familiar analogous to elementwise functions can be developed: isotropic-Tanh, isotropic-Relu, and isotropic-Leaky-Relu. However, it is hoped that the development of the paradigm will produce further activation functions that are not just analogues of existing activation functions but exploit the novel properties of isotropy for optimal performance.

Radial Basis Form	Elementwise Form	Isotropic Form
$\mathbf{f}(\vec{x}) = \sum_{i=1}^N \sigma(\ \vec{x} - \vec{c}_i\) \hat{e}_i$	$\mathbf{f}(\vec{x}) = \sum_{i=1}^n \sigma(\vec{x} \cdot \hat{e}_i) \hat{e}_i$	$\mathbf{f}(\vec{x}) = \sigma(\ \vec{x}\) \hat{x}$

Isotropic-Tanh is described in *Eqn. 16*. In basis directions, \hat{e}_i , it is equal in function to standard elementwise-tanh, as indicated by its name. It is bounded, up to a norm of one, but does not angularly saturate like standard tanh, allowing activation to continue semantically shifting.

$$\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad \vec{x} \mapsto \mathbf{f}(\vec{x}) = \tanh(\|\vec{x}\|) \hat{x} \quad (16)$$

This function is reasonably cheap, computation of $r = \|\vec{x}\|$, $\tanh(r)$ and $\text{sech}^2(r)$, need only be computed once (including for backward-pass) rather than per-component like the anisotropic functional forms. The vector norms are naturally constrained to $[0, 1)$ acting as an implicit normaliser. Around the origin, the transform is approximately the identity: $\lim_{r \rightarrow 0} \mathbf{J}(r\hat{x}) = \mathbf{I}_n$, justifying $\mathbf{f}(\vec{0}) = \vec{0}$, to preserve a smooth gradient. It is also globally 1-Lipschitz.

Isotropic-ReLU is shown in *Eqn. 17*, an analogue to its traditional implementation.

$$\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad \vec{x} \mapsto \mathbf{f}(\vec{x}; R_0) = \max(\|\vec{x}\| - R_0, 0) \hat{x} \quad (17)$$

In effect, all activations are reduced by a threshold magnitude, R_0 , with negative resultant magnitudes set to zero. Variations can be made to this activation function as shown in *Eqns. 18* and *19*, which include a maximum magnitude, R_∞ or do not reduce magnitudes except for below R_0 , respectively.

$$\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad \vec{x} \mapsto \mathbf{f}(\vec{x}; R_0, R_\infty) = \min(\max(\|\vec{x}\| - R_0, 0), R_\infty) \hat{x} \quad (18)$$

$$\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad \vec{x} \mapsto \mathbf{f}(\vec{x}; R_0) = \begin{cases} \vec{0} & : \|\vec{x}\| < R_0 \\ \vec{x} & : \|\vec{x}\| \geq R_0 \end{cases} \quad (19)$$

These activation functions continue to use $\mathbf{f}(\vec{0}) = \vec{0}$ property.

279 **Isotropic-Leaky-ReLU** follows a similar form to ReLU; however, it linearly rescales the magni-
 280 tudes below the threshold, forming a ball of smaller rescaled magnitudes. It is displayed in *Eqn. 20*,
 281 with a small value $0 < \alpha \ll 1$.

$$\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad \vec{x} \mapsto \mathbf{f}(\vec{x}; \alpha, R_0) = \begin{cases} \alpha \vec{x} & : \|\vec{x}\| < R_0 \\ \vec{x} - (1 - \alpha) R_0 \hat{x} & : \|\vec{x}\| \geq R_0 \end{cases} \quad (20)$$

282 **Isotropic-Soft-ReLU**, $\alpha = 0$, and ‘Isotropic-SoftLeaky-ReLU’, $\alpha \in (0, 1)$, are left in the deriva-
 283 tive form of the radial part $\sigma'(r)$ in *Eqn. 21*, where $\phi(r)$ is a monotonically increasing function .
 284 There are very many suitable candidates fulfilling this ϕ and one may be selected which has a suitable
 285 balance between performance, computation-cost and desirable properties. Imposed is $0 < \delta < R_0$,
 286 where $\delta < R_0$ is the centre-point of the interpolation window and 2δ is the width of this window.
 287 Consequently, the function blends smoothly between two linear regions of differing scaling.

$$\sigma : \mathbb{R} \rightarrow \mathbb{R}, \quad r \mapsto \sigma'(r; R_0, \delta, \alpha, \phi) = \begin{cases} \alpha & : \|\vec{x}\| \leq R_0 - \delta \\ \frac{\phi\left(\frac{r - R_0 + \delta}{2\delta}\right)}{\phi\left(\frac{r - R_0 + \delta}{2\delta}\right) + \phi\left(\frac{R_0 + \delta - r}{2\delta}\right)} & : R_0 - \delta < \|\vec{x}\| < R_0 + \delta \\ 1 & : \|\vec{x}\| \geq R_0 + \delta \end{cases} \quad (21)$$

288 **Isotropic-Sinusoids**, is a possibility which does not have an analogue within the current anisotropic
 289 paradigm — demonstrating a broad array of new possibilities. With the aforementioned isotropic-
 290 Relu-like functions, the networks may normalise magnitudes such that the distributions ‘escape’ the
 291 non-linear regions of the function, due to utilising the non-linearity constructively, which can initially
 292 be an unpredictable learning hurdle. Therefore, a function that introduces non-linearity throughout
 293 the space may be desirable. Proposed is ‘isotropic-Sinusoids’, which allows for distributions to be
 294 compressed, rarefied and folded (for $|\lambda_m| > 1$) in a predictable manner. It is hoped the network can
 295 utilise this for effective computation. This activation function is demonstrated in *Eqn. 22*. It includes
 296 a monotonicity-violating parameter $\lambda_m \in \mathbb{R}$, which may be useful.

$$\mathbf{f}(\vec{x}) = \vec{x} + \lambda_m \sin(\|\vec{x}\|) \hat{x} \quad (22)$$

297 5 Argument Against Isotropy (and Quasi-Isotropic Activation Functions)

298 It is argued that anisotropies result in an activation distribution shift, which may be detrimental to the
 299 network’s performance. This is because this inductive bias is typically introduced universally, and if
 300 no justification exists for this particular distribution, then it may be a suboptimal imposition by the
 301 network designer. However, it could also be argued that some symmetry-breaking anisotropy may be
 302 beneficial, by clustering parts of the activation distribution, leading to classifications that develop
 303 more quickly. In classification, one of the most common applications of deep learning, this clustering
 304 may be a suitable a priori justification for anisotropy. Therefore, introducing isotropy may limit the
 305 network’s performance in this case.

306 Despite this, current activation functions produce anisotropies along a Cartesian grid, due to their
 307 standard basis dependence. This particular arrangement does not seem justified through classification.
 308 Anisotropies can be introduced in a more general arrangement, enabling a more uniform distribution
 309 of directions from which semantics could develop.

310 **A middleground** may be to relax the hard isotropy condition and introduce slight symmetry
 311 breaking in many directions. Then the network has many distinguished vectors, a subset of which it
 312 may align its representations too in a task-dependent manner. Therefore, it does not favour a particular
 313 basis, but still introducing some desirable consequences of anisotropy. If one further restricts the
 314 functions to not feature dynamic refraction, then it limits detrimental anisotropic effects.

315 This approximately basis-free anisotropy appears preferable since it does not constrain representations
 316 to the arbitrary standard basis. It can be achieved by introducing many small perturbations to the
 317 direction unit-vector only, producing a softer symmetry breaking.

One method is to apply a non-linearity based on rounding the vector’s directions. This is shown in Eqn. 23, where $\lceil \cdot \rceil$ indicates the rounding operation and $\phi(\vec{x}) \neq \vec{x}$. In fact, the anisotropic perturbation may be implemented as simply as: $\phi(\vec{x}) = \beta \vec{x}$ for $\beta \neq 1$. The overall angular term is approximately unit-normed, but can be trivially modified to be exactly norm-1.

$$\Phi(\hat{x}; \alpha) = \frac{\lceil \alpha \hat{x} \rceil}{\alpha} + \phi\left(\hat{x} - \frac{\lceil \alpha \hat{x} \rceil}{\alpha}\right) \approx \hat{x} \quad (23)$$

This produces a quasi-isotropic functional form shown in Eqn. 24, with an isotropy-breaking parameter α . Slight anisotropic refraction is added, independent of magnitude, such that it is predictable and thus extrapolatable to the network. Due to the angular rarefaction and compression by the proposed non-linearity, representation over- and underdensities may then occur, where semanticity may begin to be assigned. However, for $\alpha \rightarrow \infty$, isotropy is continuously reintroduced and could be an optimisable parameter.

$$\mathbf{f}(\vec{x}) = \sigma(\|\vec{x}\|) \Phi(\hat{x}) \quad (24)$$

5.1 Real-Time Dynamical Network Topology

An appealing feature of isotropic deep learning is the relation displayed in Eqn. 15, showing that due to rotational equivariance a rotation to one weight matrix can be counteracted with the inverse-rotation of another, preserving the network’s function. Consequently, a particular gauge can be chosen which expresses the weights in a beneficial basis.

One such basis may be a magnitude ordering of the singular values for the matrices. One could then set a threshold for the singular value to determine if each corresponding direction in such a matrix has a meaningful contribution to the overall functionality. If it is deemed to have negligible value, it can be pruned with little adverse effect on the network.

Moreover, ζ latent neurons can be included, with zero-initialised singular values fully connected to existing neurons. Since the Jacobians of the isotropic activation functions are not strictly diagonal, these latent neurons may be rapidly trained if required. Therefore, the otherwise static fully-connected network is now dynamic, growing and shrinking with task-necessitated demand, with minimal impact to performance with these actions. This is only enabled through an isotropic functional form, due to the continuous rotational symmetry available. It poses an interesting research direction, where transfer learning and task-swapping may become more straightforward. Output and input neurons could also be appended and removed in such a way, allowing for real-time changes to a dataset, or even training on multiple datasets. Such a procedure could be trivially extended to convolutional networks, allowing a dynamic number of kernels.

This could offer substantial insight into how parameters may be shared between tasks in real-time. For example, we may postulate that if a new dataset is introduced partway through training on a different dataset, there might be a short-term parameter increase until the network parameter-sharing begins, followed by a phase of pruning until a more compact architecture is reached. These network dynamics may be incredibly insightful.

It appears it may side-step the lottery ticket hypothesis in choosing optimal network size prior to training. Due to the computational cost, this does not need to be computed at every step, only periodically, and can be performed layerwise.

6 Conclusion

In this position paper, the isotropic functional form is proposed as a better default inductive bias for deep learning. Current forms have been demonstrated in the literature to produce task-unmotivated representational artifacts which may limit semantic expressibility of the networks. It is further argued that the current anisotropic functional forms may have detrimental effects in performance and learning, through the ‘neural refractive problem’, ‘weight locking’ and ‘discrete semantics’. Hence, removing such constraints from the model is argued to also unconstrain the representations from any particular basis. It is then expected that the network can produce a more natural activation representation based upon task necessities, rather than human imposed by functional forms.

The adjustment to functional forms is through promoting the existing discrete rotational symmetry (permutation symmetry) of modern deep learning to a continuous special-orthogonal symmetry. This

366 has substantial consequences for the form of almost every function in the modern day deep learning:
367 activation functions, regularisers, normalisers, optimisers and more. It is proposed that the breadth of
368 this reformulation, is best represented as a novel branch of deep learning: *Isotropic deep learning* to
369 distinguish it from existing paradigms.

370 In this position paper, several example isotropic functions are showcased as a starting point, yet these
371 are analogues and not expected to be inherently optimal or better since they display superficial their
372 similarity to existing functions. Instead, it is the authors position that this change to isotropic deep
373 learning is generally advantageous, but may need substantial time to development as a paradigm
374 such that better optimised implementations are discovered which suitably leverage the new properties
375 available from isotropy. Therefore, empirical work will be presented in following papers as to not
376 distract from the primary arguments motivating this shift to Isotropic deep learning.

377 It is hoped that the proposed ideas may stimulate the communities' interest into beginning a search
378 for such Isotropic functions, which will hopefully bring the paradigm into widespread applicability
379 and adoption.

References

- George Bird. The spotlight resonance method: Resolving the alignment of embedded activations. In *Second Workshop on Representational Alignment at ICLR 2025*, 2025. URL <https://openreview.net/forum?id=alxPpqVRzX>.
- Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposition, 2022. URL <https://arxiv.org/abs/2209.10652>.
- George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.
- Kurt Hornik. Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4(2):251–257, 1991. ISSN 0893-6080. doi: [https://doi.org/10.1016/0893-6080\(91\)90009-T](https://doi.org/10.1016/0893-6080(91)90009-T). URL <https://www.sciencedirect.com/science/article/pii/089360809190009T>.
- R. Quian Quiroga, L. Reddy, G. Kreiman, C. Koch, and I. Fried. Invariant visual representation by single neurons in the human brain. *Nature*, 435(7045):1102–1107, Jun 2005. ISSN 1476-4687. doi: 10.1038/nature03687. URL <https://doi.org/10.1038/nature03687>.
- Katharine M Cammack, Thomas R Reppert, and Denise R Cook-Snyder. The simpsons neuron: A case study exploring neuronal coding and the scientific method for introductory and advanced neuroscience courses. *J Undergrad Neurosci Educ*, 20(1):C1–C10, December 2021.
- G Kreiman, C Koch, and I Fried. Category-specific visual responses of single neurons in the human medial temporal lobe. *Nat Neurosci*, 3(9):946–953, September 2000.
- Charles G. Gross. Genealogy of the “grandmother cell”. *The Neuroscientist*, 8(5):512–518, 2002. doi: 10.1177/107385802237175. URL <https://doi.org/10.1177/107385802237175>. PMID: 12374433.
- Charles E Connor. Neuroscience: friends and grandmothers. *Nature*, 435(7045):1036–1037, June 2005.
- Jerzy Konorski. Learning, perception, and the brain: Integrative activity of the brain. an interdisciplinary approach. *Science*, 160(3828):652–653, 1968. doi: 10.1126/science.160.3828.652. URL <https://www.science.org/doi/abs/10.1126/science.160.3828.652>.
- H B Barlow. Summation and inhibition in the frog’s retina. *J Physiol*, 119(1):69–88, January 1953.
- Daniel J Graham and David J Field. Sparse coding in the neocortex. *Evolution of nervous systems*, 3: 181–187, 2006.
- J. Y. Lettvin, H. R. Maturana, W. S. McCulloch, and W. H. Pitts. What the frog’s eye tells the frog’s brain. *Proceedings of the IRE*, 47(11):1940–1951, 1959. doi: 10.1109/JRPROC.1959.287207.
- H. K. Hartline. The response of single optic nerve fibers of the vertebrate eye to illumination of the retina. *American Journal of Physiology-Legacy Content*, 121(2):400–415, 1938. doi: 10.1152/ajplegacy.1938.121.2.400. URL <https://doi.org/10.1152/ajplegacy.1938.121.2.400>.
- Warren S. McCulloch and Walter Pitts. A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5(4):115–133, Dec 1943. ISSN 1522-9602. doi: 10.1007/BF02478259. URL <https://doi.org/10.1007/BF02478259>.
- Frank Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65 6:386–408, 1958. URL <https://api.semanticscholar.org/CorpusID:12781225>.
- Andrew L Maas, Awni Y Hannun, Andrew Y Ng, et al. Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml*, volume 30, page 3. Atlanta, GA, 2013.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9:1735–80, 12 1997. doi: 10.1162/neco.1997.9.8.1735.

- 427 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolu-
428 tional neural networks. *Advances in neural information processing systems*, 25, 2012.
- 429 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz
430 Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1706.03762)
431 [1706.03762](https://arxiv.org/abs/1706.03762).
- 432 Kai Hu and Barnabas Poczos. Rotationout as a regularization method for neural network, 2020. URL
433 <https://openreview.net/forum?id=r1e7M6VYwH>.
- 434 Wallace Givens. Computation of plain unitary rotations transforming a general matrix to triangular
435 form. *Journal of the Society for Industrial and Applied Mathematics*, 6(1):26–50, 1958. doi:
436 [10.1137/0106004](https://doi.org/10.1137/0106004). URL <https://doi.org/10.1137/0106004>.
- 437 Adelaide P Yiu, Valentina Mercaldo, Chen Yan, Blake Richards, Asim J Rashid, Hwa-Lin Liz Hsiang,
438 Jessica Pressey, Vivek Mahadevan, Matthew M Tran, Steven A Kushner, Melanie A Woodin,
439 Paul W Frankland, and Sheena A Josselyn. Neurons are recruited to a memory trace based on
440 relative neuronal excitability immediately before training. *Neuron*, 83(3):722–735, August 2014.
- 441 Lingxuan Chen, Kirstie A Cummings, William Mau, Yosif Zaki, Zhe Dong, Sima Rabinowitz,
442 Roger L Clem, Tristan Shuman, and Denise J Cai. The role of intrinsic excitability in the evolution
443 of memory: Significance in memory allocation, consolidation, and updating. *Neurobiol. Learn.*
444 *Mem.*, 173(107266):107266, September 2020.

A Historical Precedent

Neural networks, of all forms, represent their information through a neural code which remains under debate Quiroga et al. [2005], Cammack et al. [2021], Kreiman et al. [2000]. Several codings have been proposed and it appears that initially anisotropic functions may have indirectly arisen from one of these, local coding, which was relatively more common around the advent of artificial neural networks. Local coding is a neuropsychological hypothesis of one neuron’s activation representing the presence of one real-world stimulus. This is also commonly known as the grandmother neuron interpretation Gross [2002], Connor [2005], and less so as gnostic neurons and gnostic fields Konorski [1968]. This hypothesis was developed Barlow [1953], Konorski [1968] and debated from the early 1950s through to the 1970s Gross [2002], with continuing discussion to the present day Quiroga et al. [2005], Connor [2005], Graham and Field [2006].

Crucially, McCulloch and Pitts co-authored a paper supporting the presence of several differing pattern feature detectors in frog’s retinas Lettvin et al. [1959] expanding on Hartline’s earlier work of similar findings Hartline [1938]. Thus, it was shown that retinal ganglions can respond to distinct real-world patterns, so-called ‘bug detectors’ for the frog, acting much like the description of grandmother neurons. These are the same McCulloch and Pitts who are earlier credited as the inventors of the binary threshold network McCulloch and Pitts [1943] which led to the first perceptron neural network being developed. In the very first sentence of McCulloch and Pitts [1943] work, it states:

Because of the "all-or-none" character of nervous activity, neural events and the relations among them can be treated by means of propositional logic.

—A Logical Calculus of the Ideas Immanent in Nervous Activity McCulloch and Pitts [1943], p. 115.

Hence this premises artificial neural networks on a binary logic, inducing this early basis-dependent nature. This appears to have implicitly encouraged the adoption of the heaviside step function in early deep learning, which has a continuous lineage to the vast array of elementwise functions utilised today. This lineage is connected through a series incremental modifications from discrete heaviside step functions, such as differentiable sigmoid-based functions, to non-saturating ReLU through to its variants. Which in-turn may have influenced the wider array of functions, like elementwise dropout or approximating Hessians as a diagonal in adaptive optimisers.

Therefore, I argue that the pervasiveness of anisotropic deep learning can be traced back to a trajectory set in the early developments in the field — an underappreciated *choice* of functional form influenced, by the discoveries of their time.

A further indication is in one of McCulloch and Pitts [1943]’s concluding statements:

[...] pushed to ultimate psychic units or “psychons,” for a psychon can be no less than the activity of a single neuron. Since that activity is inherently propositional, all psychic events have an intentional, or “semiotic,” character. The “all-or-none” law of these activities, and the conformity of their relations to those of the logic of propositions, insure that the relations of psychons are those of the two-valued logic of propositions. Thus in psychology, introspective, behavioristic or physiological, the fundamental relations are those of two-valued logic.

—A Logical Calculus of the Ideas Immanent in Nervous Activity McCulloch and Pitts [1943], p. 131.

This “psychon”, with semiotic qualities and two-valued logic, is highly suggestive of a local coding approach — aligning with the paradigm of vector decomposition in modern neural networks through generalising the two-value logic. This local coding perspective is argued to have led to a preference for elementwise logic and functional forms. Hence, the activation vectors within neural networks are treated as just an array of numbers, as opposed to a direction and magnitude, which is the defining feature of the elementwise anisotropic forms.

As a result, the roots of anisotropies can be seen in neural networks since their first conception. At this time, the non-linear activation function utilised was the heaviside step function and this persisted into the multilayer perceptron Rosenblatt [1958], but later generalised whilst keeping an elementwise form Maas et al. [2013]. These continuous-generalised perceptron layers continue to be prevalent becoming a pervasive backbone of many current models Hochreiter and Schmidhuber

498 [1997], Krizhevsky et al. [2012], Vaswani et al. [2023], as well as being adapted into other structures
499 such as convolutional neural networks.
500 [unfinished]

501 **B Distinction from Equivariant Networks of Geometric Deep Learning**

502 Both isotropic deep learning and equivariant networks share a similar equivariant relation, which
503 may make them appear superficially similar. However, they differ substantially in how this relation
504 is implemented, motivated and consequences. This section will briefly outline those differences.
505 [unfinished]

506 B.1 Normalisers, Regularisers and Optimisers

507 Do not wish to discourage anisotropic distributions if they are beneficial, afterall, the goal is to
 508 unconstrain the network and regularisers should reflect this. The scale factor α controls how many
 509 anisotropic-centres there are
 510 gradient clipping Can imagine a dense thick layer hybridising the a
 511 [unfinished]

512 C Stochastic Isotropy — Producing Immediate Anisotropic Analogs

513 One method to approximate isotropy with current functions is by stochastically choosing a basis for
 514 the anisotropic function to operate on. This enables anisotropic functions to be used in an isotropic
 515 network, without inducing a representational alignment to an arbitrary basis.

516 For example current (anisotropic-)dropout would appear to privilege the basis anti-aligned with
 517 the standard basis, to maximally preserve information when a direction of the standard basis is
 518 collapsed. So it is expected to incur an arbitrary basis dependence onto the representations. However,
 519 anisotropic-dropout can be applied to a stochastically chosen basis. This randomness is hypothesised
 520 to prevent a representational-anisotropy induced by an arbitrarily chosen basis.

521 This can be achieved by producing a basis, uniformly drawn from the layer’s special-orthogonal
 522 symmetry: $\mathbf{B} \sim \text{SO}(n)$. There are many such methods to produce such a uniform random ma-
 523 trix, with varying computational costs, such as via exponentiation of Lie generator scaled by an
 524 appropriately drawn random variable, the gram-schmidt procedure and many others. The existing
 525 matrix-multiplication procedure is computationally cumbersome, so simpler formulations may be
 526 desirable. The following proposed forms are only a starting point for converting anisotropic functions
 527 directly to stochastically-isotropic forms. In practice, isotropic functions should be constructed from
 528 the ground-up rather than merely analogous functions converted from existing anisotropic ones.

529 For the example of standard dropout, shown in Eqn. 25, it can be made stochastically-isotropic by
 530 including the basis-transform shown in Eqn. 26. Where \vec{x} is the activation vector, with normalisation
 531 factor S_a and S_{si} , dropout-mask $M_i = \vec{M} \cdot \hat{e}_i$ and standard-basis vectors \hat{e}_i .

$$\vec{x}' = S_a \sum_{i=1}^N M_i (\vec{x} \cdot \hat{e}_i) \hat{e}_i \quad (25)$$

532

$$\vec{x}' = S_{si} \sum_{i=1}^N (\mathbf{B} \vec{M} \cdot \hat{e}_i) (\vec{x} \cdot \hat{e}_i) \hat{e}_i \quad (26)$$

533 Similar formulations have been demonstrated, such as RotationOut by Hu and Póczos [2020], which
 534 showed generally improved performance when the basis is effectively stochastically rotated. However,
 535 this method remains stochastically anisotropic as the rotations are generated through Given’s rotations
 536 Givens [1958], which is not uniform over the space of special-orthogonal matrices — a necessity
 537 for full stochastic-isotropy. Nevertheless, the implementation by Hu and Póczos [2020] is somewhat
 538 encouraging.

539 Overall, this procedure can be generalised and applied to any existing anisotropic function, yet it is
 540 generally preferable to construct an isotropic function from first principles rather than relying on
 541 stochastic-isotropy which may be computationally costly.

542 C.1 Considering Correllating the Stochastic-Isotropy

543 A curious extention, particularly to stochastically-isotropic dropout, would be to correllate the
 544 random-bases in time. This may produce a time-like structure in a network’s embedded activation
 545 distribution.

546 If one imagines a random walk of the rotation matrices: $\text{SO}(n) \ni \mathbf{R}^{(t+dt)} = \mathbf{R}^{(t)} \delta \mathbf{R}$, with
 547 $\delta \mathbf{R} = e^{\tau \cdot \vec{n}}$, with τ being the corresponding (normalised) anti-symmetric generators for rotations and
 548 $\vec{n} \sim \mathcal{N}(\vec{0}, \sigma \mathbf{I}_n)$ with $0 < \sigma \ll 1$. This procedure results in a random walk of the rotation matrix at
 549 each time step.

550 Following this, a time-correlated Bernoulli distribution can be defined. Beginning with $\vec{D}^{(0)}$, divide up
 551 the layer of neurons into two sets: inactive $I_n = \left\{ i | \vec{D}_i^{(n-1)} = 0 \right\}$ and active $A_n = \left\{ i | \vec{D}_i^{(n-1)} = 1 \right\}$.
 552 Then we have two hyper-parameters: the standard dropout probability λ and an overlap probability Γ ,
 553 such that $|A_n| q + |I_n| \Gamma = (|A_n| + |I_n|) \lambda$ - where q is not a free parameter. If $|A_n| = 0$ or $|I_n| = 0$,
 554 then temporarily define $q = \Gamma = \lambda$. If not, then one needs to prevent unnormalised probabilities as
 555 shown in Eqn. 27.

$$\Gamma = \max \left(0, \max \left(\lambda + \frac{|A_n|}{|I_n|} (\lambda - 1), \min \left(1, \min \left(\lambda + \frac{|A_n|}{|I_n|} \lambda, \Gamma \right) \right) \right) \right) \quad (27)$$

556 Leading to $q = \lambda + \frac{|I_n|}{|A_n|} (\lambda - \Gamma)$. Then use one Bernoulli function across all active neurons
 557 using $\mathbb{R}^{|A_n|} \ni \vec{D}_{(A)}^{(n)} \sim \text{BernoulliDist.}^{|A_n|}(q)$ likewise for inactive neurons $\mathbb{R}^{|I_n|} \ni \vec{D}_{(I)}^{(n)} \sim$
 558 $\text{BernoulliDist.}^{|I_n|}(\Gamma)$. Therefore, correlating the inactive neurons across the time steps, whilst still
 559 introducing a degree of random dropout. Thus, the ‘basis of dropout’ undergoes a random walk at
 560 every time step, and neurons are randomly chosen to be dropped from the network, with a differing
 561 likelihood if they were just previously dropped. The coherence time can be adjusted through Γ , for
 562 the specific time-dependent task needed.

563 This creates a link between the stimulus’ presentation time to the network and the neurons its alters,
 564 such that stimuli presented in a smaller time window perturb similar subset of the network’s neurons.
 565 This may produce an encoding similar to that found in human cognition, where neurons are thought
 566 to go through excitability cycles of slightly differing frequencies and phases. When the excitability is
 567 higher, information (engrams) preferentially encodes upon those neurons Yiu et al. [2014], Chen et al.
 568 [2020]. As groups of neurons begin to decohere, there remains some overlap, such that memories
 569 are interlaced if they occur within a temporal window of coherence. *This potentially gives neural*
 570 *networks using isotropic dropout an advantage when it comes to time-series data.*

571 **D Taxonomy of Functional Forms**

572 [unfinished]

E Isotropy In Transformers

It is argued that isotropic deep learning may be a more appropriate inductive bias for deep learning. However, there may also be some architectures which especially benefit from its inclusion. One of these is the self-attention step of transformers Vaswani et al. [2023], where isotropic-tanh may be of particular benefit, in replacing the softmax operation .

Softmax is defined through elements being bounded between zero and one, $\mathbf{f}(\vec{x}) \cdot \hat{e}_i \in [0, 1]$ and summing to one. As a consequence it is non-negative and there are regimes where this may be limiting.

It has been shown that representations can exist in an antipodal superposition Elhage et al. [2022], particularly when stimuli do not coexist, so interference is minimised. Such a scenario may be a quantity which is continuous, but negative and positive values are mutually exclusive. Many of these semantics exist in the real-world: daytime-to-nighttime. These could be represented through a zero-to-one scale, but zero may be better represented as a neutral middle-point instead. In this case, the negative of a semantic direction may be equally meaningful. It may be expected that this is the case in the self-attention step.

Moreover, the sum-to-one case may not always be desirable: it always encourages a change to the semantics, when considering the residual-step-modification. **This may force an semantic correction to an activation in transformers even when it is inappropriate.**

The self-attention step comparese the pairwise similarities between several vectors grouped into the so-called ‘keys’ and ‘queries’. The degree-of-similarity then affects how much of another semantic is expressed: the ‘values’. However, the softmax layer prevents a negative expression of these value semantics.

A more suitable choice may be isotropic-tanh. In analogy of its sum-to-one constraint, its vector-magnitude is at maximum one, $0 \leq \|\mathbf{f}(\vec{x})\| \leq 1$, whilst elementwise its values are $-1 \leq \mathbf{f}(\vec{x}) \cdot \hat{e}_i \leq 1$. Hence, it can express a negative of the value semantic, or any scaling of it between -1 and 1 . This suggests that isotropic-tanh may be quite a appealing drop-in-replacement for softmax in the attention step, at least conceptually. Its continuous rotational symmetry may also offer advantage, since the underlying outer-product of self attention $QK^T = \vec{x}^T W_Q^T W_K \vec{x} \hat{=} \vec{x}^T W_{kq}^T \vec{x}$, is also isotropic in \vec{x} , then enabling more even interpolation between, and perturbation to, the values. Hence, an isotropic adaption to a self-attention may appear as shown in Eqn. 28, which will be explored in future work.

$$\text{Attention}(Q, K, V) = \text{Isotropic-Tanh} \left(\frac{QK^T}{\sqrt{d_k}} \right) V \quad (28)$$

However, this does not make transformers ‘isotropic’ as a whole, since further anisotropic steps exist.