# Neural Manifold Geometry Encodes Feature Fields

**Julian Yocum**, Cam Allen, Bruno Olshausen, Stuart Russell

CHAI — Center for Human-Compatible Artificial Intelligence
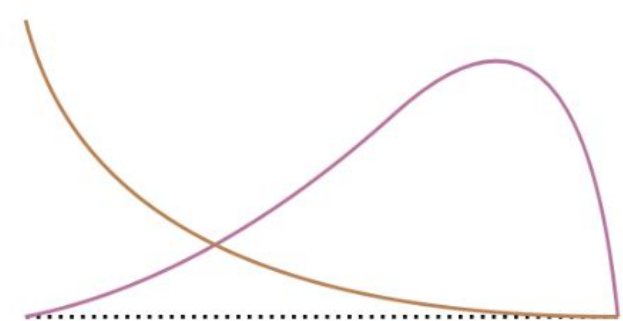
Berkeley — UNIVERSITY OF CALIFORNIA

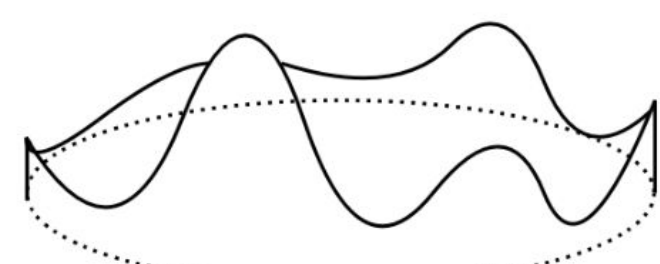**1** How do neural networks represent **functions** over **spaces**?
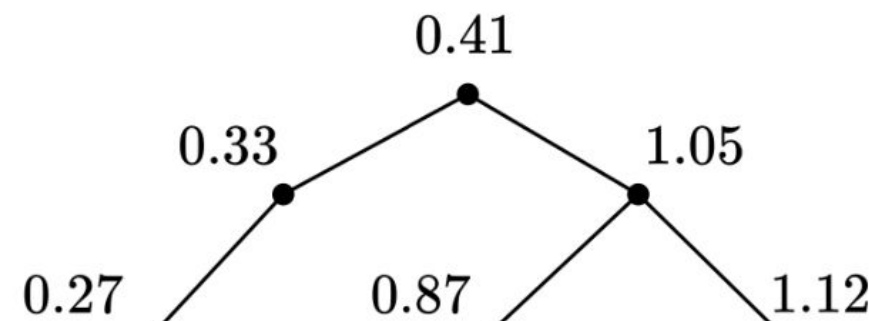
$f_a(x_1, I)$    $f_a(x_2, I)$    $f_b(x_3, S^1)$    $f_c(x_4, G)$

0.41
0.33    1.05
0.27    0.87    1.12

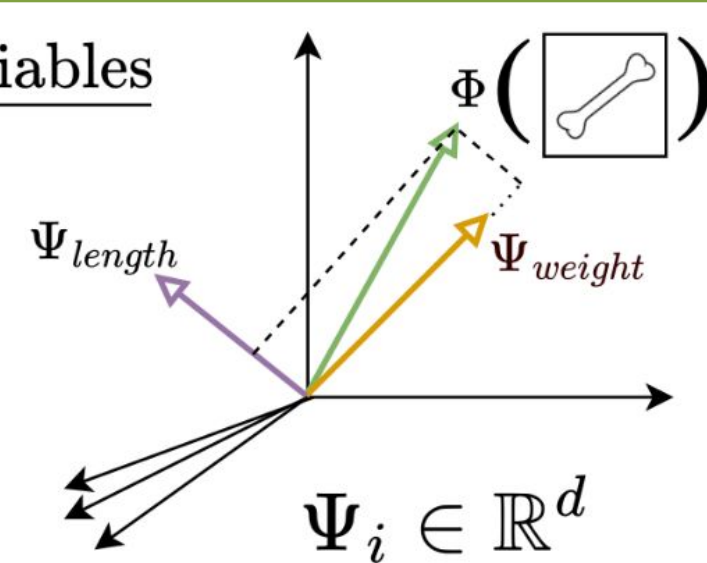$\mathcal{Z} = I = [0, 1]$        $\mathcal{Z} = S^1$        $\mathcal{Z} = G$

Given a **domain space Z**, a **feature field** is a distribution of functions over **Z**, e.g. value function, belief distribution.

**2**

While **linear probes** recover **scalars** from activations, we use **linear field probes** to recover **functions**.

Feature Variables

$\Phi(\phantom{bone})$
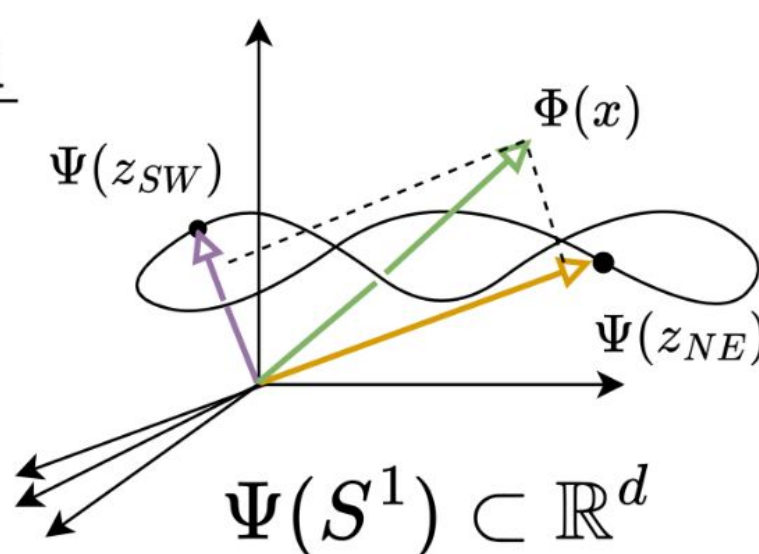$\Psi_{length}$    $\Psi_{weight}$
$\Psi_i \in \mathbb{R}^d$

$f_i(x) = \Phi(x) \cdot \Psi_i$

$\Phi(\phantom{bone}) \cdot \Psi_{weight} = 1.03$    $f_{weight}(\cdot) = $ How heavy (kg)?

$\Phi(\phantom{bone}) \cdot \Psi_{length} = 0.26$    $f_{length}(\cdot) = $ How long (m)?

Feature Field

$\Psi(z_{SW})$    $\Phi(x)$
$\Psi(z_{NE})$
$\Psi(S^1) \subset \mathbb{R}^d$

$f(x, z) = \Phi(x) \cdot \Psi(z)$
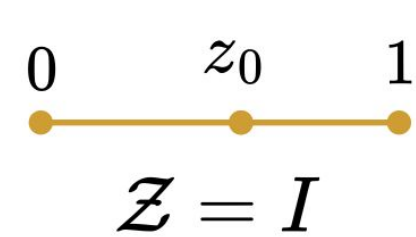
$f(z_{NE})$    $f(z_{SW})$

$f(\cdot, z) = Pr(\text{bone in direction } z)$

$f(x, S^1)$

**3** Neural networks represent a feature field by embedding its domain space into activation space.
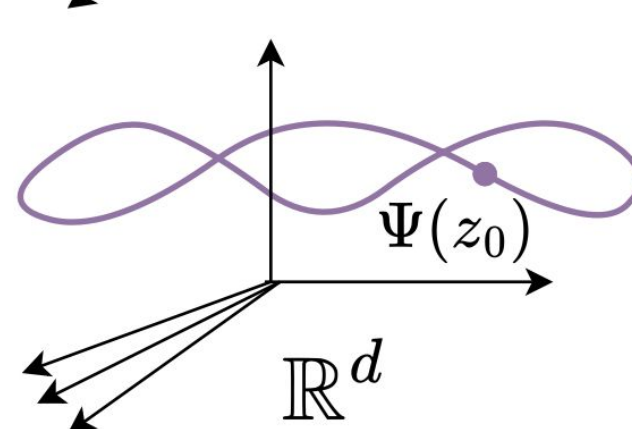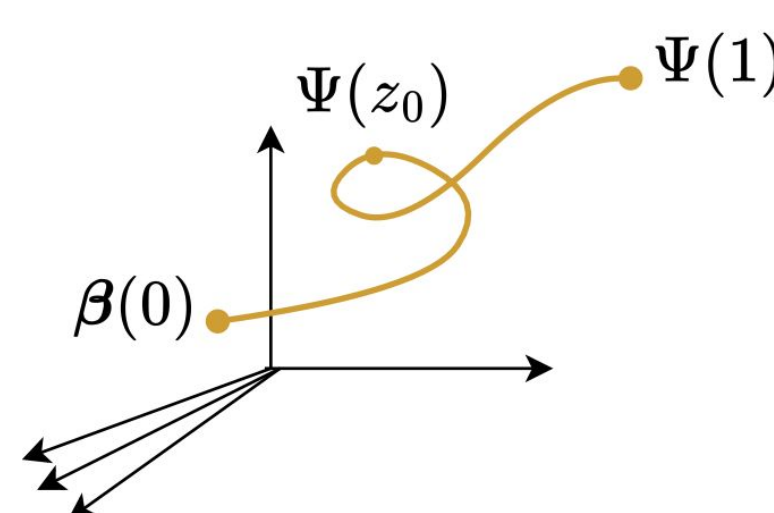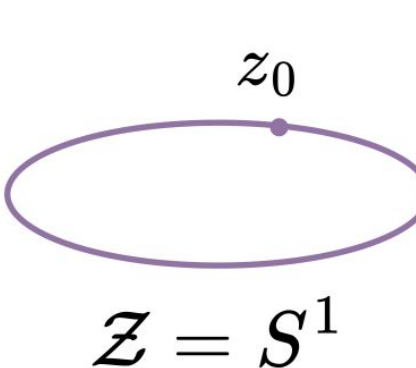
Domain Space $\mathcal{Z}$

0    $z_0$    1
$\mathcal{Z} = I$

$\Psi : I \to \mathbb{R}^d$

Domain Embedding $\Psi(\mathcal{Z})$

$\Psi(z_0)$    $\Psi(1)$
$\beta(0)$

$z_0$
$\mathcal{Z} = S^1$

$\Psi : S^1 \to \mathbb{R}^d$

$\Psi(z_0)$
$\mathbb{R}^d$
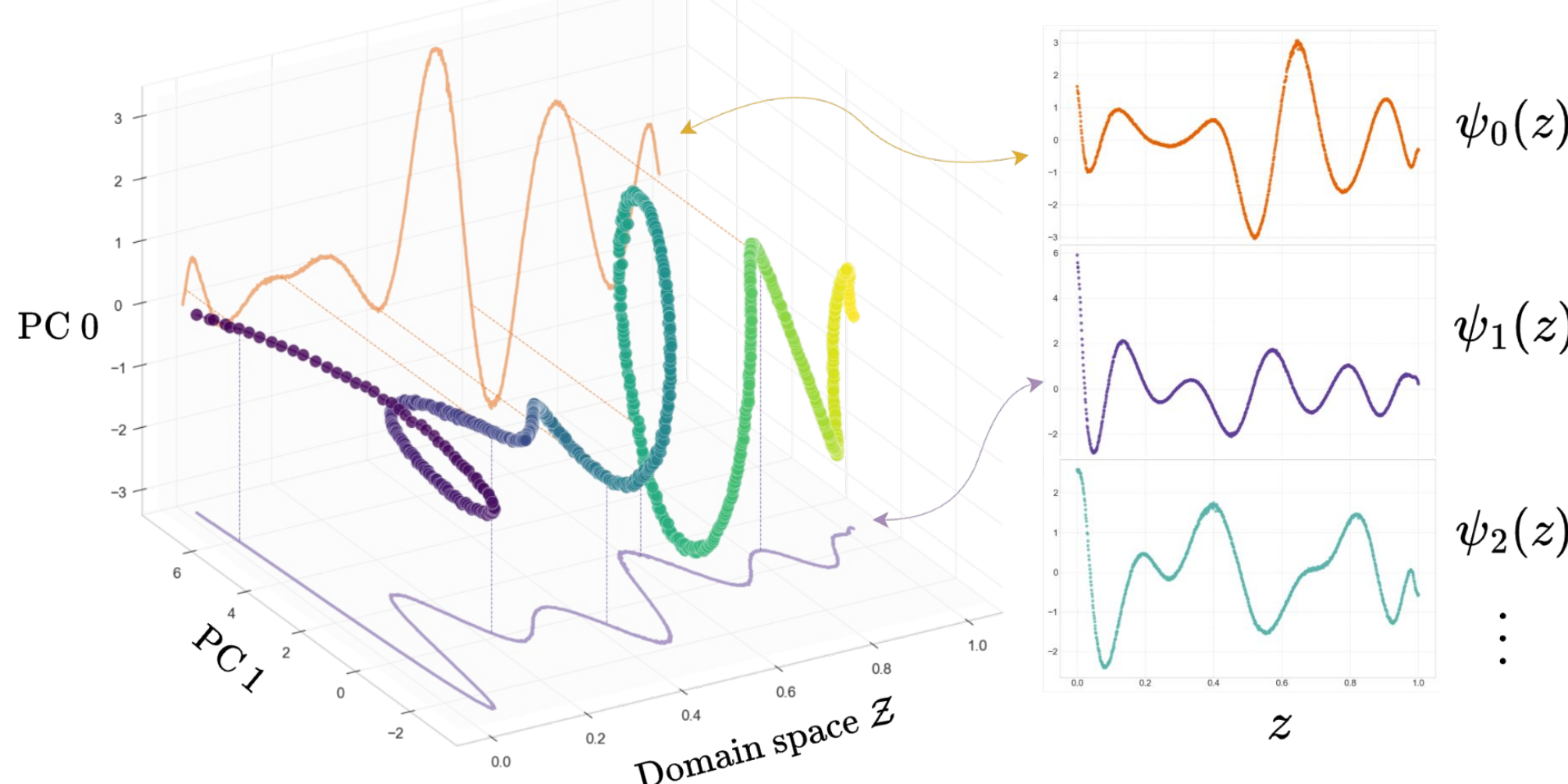
This **domain embedding** preserves the topology of the latent domain space for continuous feature fields.

**4** The geometry of the **domain embedding encodes basis functions** over **Z**.

Domain Embedding Geometry    ⟷    Basis Functions

PC 0
PC 1
Domain space $\mathcal{Z}$

$\psi_0(z)$
$\psi_1(z)$
$\psi_2(z)$
$\vdots$
$z$

Neural networks represent functions over **Z** as linear combinations of these basis functions.

$$f(\phi, x) = \sum_{i=1}^{d} a_i(x) \psi_i(\phi)$$

$f(x, z) = a_0(x) \cdot \phantom{} + a_1(x) \cdot \phantom{} + a_2(x) \cdot \phantom{} + \cdots$