# Approximation Theory of Group Invariant Neural Networks

Nadav Dym

Technion

# Is this your first time in Banff?

In July 2003 (age 16) I attended:

Mathematical Biology: From molecultes to ecosystems: the legacy of Lee Segel

`While he liked talking about his work, he had the rare quality of actually being interested in hearing about other people's work (Daniel Segel, free translation)'

# Approximation Theory of Group Invariant Neural Networks

Supervised Machine Learning: Learn $f$ from examples

$x_i$



$f(x_i)$     `cat'          `dog'          `cat'          'dog'

$$\min_{h \in \mathrm{H}} \sum_{i=1}^{N} |f(x_i) - h(x_i)|^2$$

<span style="color:red">Neural Networks</span>

Activation function: $\sigma \colon \mathbb{R} \to \mathbb{R}$

Induces $\sigma \colon \mathbb{R}^d \to \mathbb{R}^d$     $\sigma(x_1, \ldots, x_d) = (\sigma(x_1), \ldots, \sigma(x_d))$

Affine functions $h^{(i)}(x) = A^{(i)}x + b^{(i)}$ where $h^{(i)} \colon \mathbb{R}^{w_{i-1}} \to \mathbb{R}^{w_i}$

Definition: We say that $\mathcal{N} \colon \mathbb{R}^d \to \mathbb{R}^m$ is a **fully connected neural network** if

$$\mathcal{N}(x) = h_{L+1} \circ \sigma \circ h_L \circ \sigma \circ \cdots \circ \sigma \circ h_0(x)$$

**Depth of $\mathcal{N}$** $:= L$

**Width of $\mathcal{N}$** $:=$ Maximal dimension $\max\limits_{1 \le i \le L} w_i$

# Approximation Theory of Group Invariant Neural Networks

Universality Theorem [Cybenko 1989, Pinkus 1999,many others in between]

If the activation function: $\sigma: \mathbb{R} \to \mathbb{R}$ is continuous and not polynomial

then for every compact $K \subseteq \mathbb{R}^d$, continuous $f: K \to \mathbb{R}$ and $\epsilon > 0$,

There exists a **fully connected neural network** $\mathcal{N}: \mathbb{R}^d \to \mathbb{R}$ **of depth L=1** (and arbitrarily large width)

$$\mathcal{N}(x) = h_1 \circ \sigma \circ h_0(x)$$

Such that

$$|f(x) - \mathcal{N}(x)| < \epsilon, \qquad \forall x \in K$$

Universality- provides justification for choosing neural networks as a function space for any continuous learning task.

# Approximation Theory of Group Invariant Neural Networks

Beyond universality- rates of approximation (More recent research)

Given $f : K \to \mathbb{R}$ which is Lispschitz/smooth/fractal and $\epsilon$ what width $W(\epsilon)$ and depth $L(\epsilon)$

 are necessary to achieve an $\epsilon$ approximation?

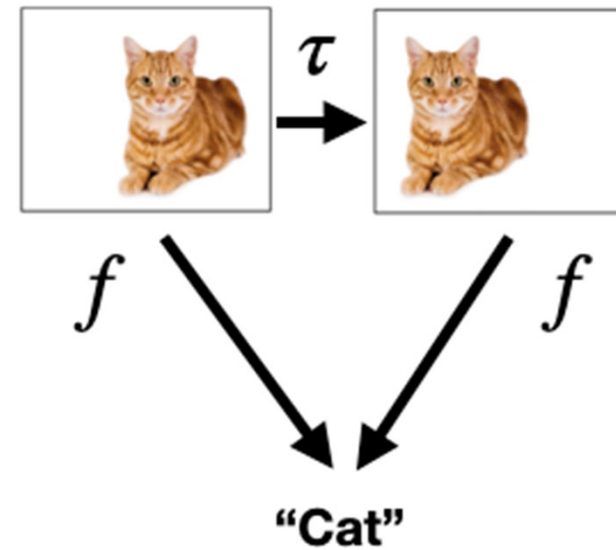# Approximation Theory of Group Invariant Neural Networks

$$\min_{h \in H} \sum_i |f(x_i) - h(x_i)|^2$$

Invariant networks:

Construct $H = H_{inv}$ so that all $h \in H_{inv}$ are invariant to the symmetries of $f$

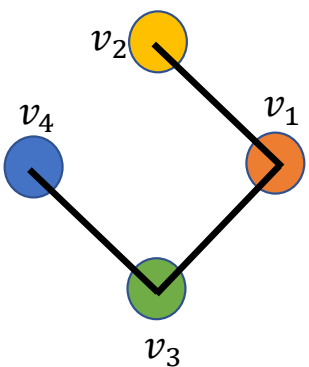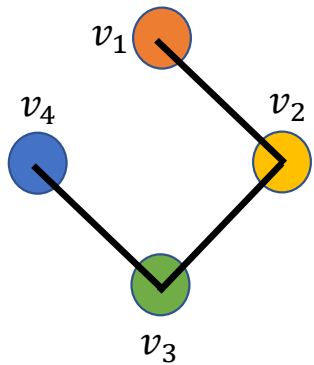(e.g., Convolutional Neural Networks for translation invariance)
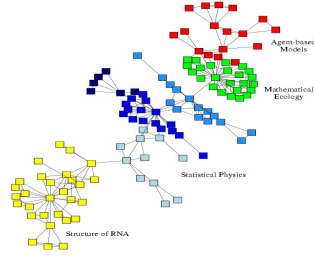
Many other examples..



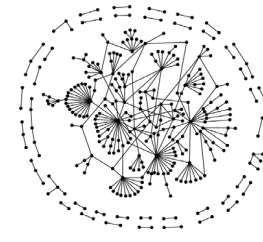Popular model class: Convolutional Neural Networks

# Invariant networks example 2: Learning on Graphs



Social networks

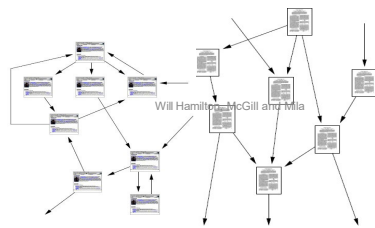Economic networks

Biomedical networks

Information networks:
Web & citations

Internet

Networks of neurons

9

Graph Neural Networks : Graph valued functions typically invariant to node relabeling

# Main example for today: point sets



$3 \times n$ points

$$X = \{x_1, x_2, \ldots, x_n\}$$

$$X = (x_1, x_2, \ldots, x_n) \sim \sigma_* X = (x_2, x_1, \ldots x_n)$$

$$\sigma \in S_n = permutations$$

# Orthogonal invariance



$$X = (x_1, x_2, \ldots, x_n) \sim R_* X = (R x_1, R x_2, \ldots R x_n)$$

$$R \in O(d) = \{ R \in \mathbb{R}^{d \times d} \mid R R^T = I_d \}$$

# Special Orthogonal=Rotation invariance



$$X = (x_1, x_2, \ldots, x_n) \sim R_* X = (R x_1, R x_2, \ldots R x_n)$$

$$R \in SO(d) = \{R \in \mathbb{R}^{d \times d} \mid RR^T = I_d, \det(R) = 1\}$$

# Rotation+Permutation invariance



**Point set symmetries:**

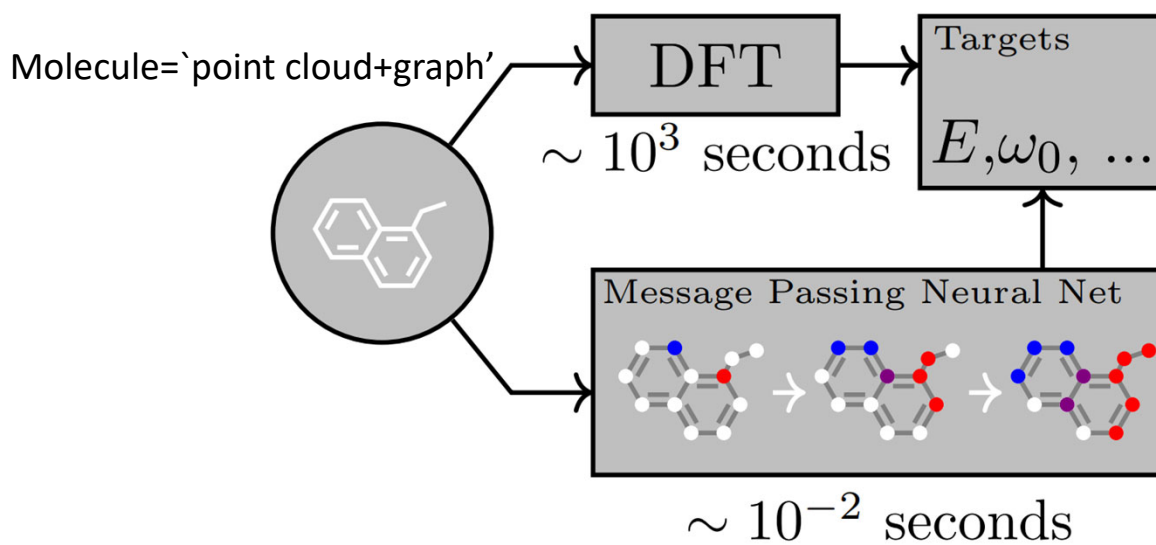Permutation $S_n$

Orthogonal $O(d)$

Rotation $SO(d)$

Orthogonal+Permutation

Rotation+Permutation

$$X = (x_1, x_2, \ldots, x_n) \sim (R, \sigma)_*(X) = (Rx_2, Rx_1, \ldots Rx_n)$$

# Scientific applications (Chemistry, Physics)

Molecule=`point cloud+graph'



$\sim 10^3$ seconds

$\sim 10^{-2}$ seconds

[Neural Message Passing for Quantum Chemistry Gilmer et al. 2017]

# Symmetry preserving architectures for point sets

**Point set networks (permutation invariant)**
[PointNet: *Deep Learning on Point Sets for 3D Classification and Segmentation,* Qi et al. 2016]
[Deep sets, Zaheer et al. 2017]
[Set Transformer: A Framework for Attention-based Permutation-Invariant Neural Networks, Lee et al. 2019]
**...**

**PointNet/DeepSets** On $\{x_1, \dots, x_n\}$ consider **permutation invariant** functions of the form

$$\{x_1, \dots, x_n\} \mapsto \mathcal{N}^{(2)}\left(\sum_{i=1}^{n} \mathcal{N}^{(1)}(x_i)\right)$$

Or $\{x_1, \dots, x_n\} \mapsto \mathcal{N}^{(2)}\left(\max_i\{\mathcal{N}^{(1)}(x_i)| i = 1, \dots n\}\right)$

**Useful principle:** Invariance cannot be `ruined' by composition (by $\mathcal{N}^{(2)}$ in this example**)**

# Symmetry preserving architectures for point sets 2

**Point set networks (rotation invariant)**
Not so much…

**Point set networks (rotations+permutation invariant)**
[Tensor field networks: Rotation- and translation-equivariant neural networks for 3D point clouds, Thomas et al. 2018]
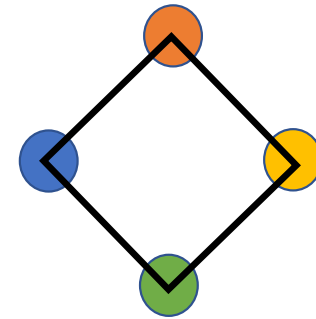[E(n) Equivariant Graph Neural Networks, Satorras et al. 2021]
[Directional Message Passing for Molecular Graphs, Gasteiger et al. 2020]
...

# Approximation Theory of Group Invariant Neural Networks

# Approximation Theory of Group Invariant Neural Networks

# Universality of invariant machine learning

Continuous functions $=\overline{Fully\ connected\ NNs}$

Continuous invariant functions $=\overline{H_{inv}}$

# Example: Universality for permutation invariant point set functions

**Question:** Can any continuous permutation invariant $f : \mathbb{R}^{d \times n} \to \mathbb{R}$

$$f(x_1, \ldots, x_n) = f\left(x_{\tau(1)}, \ldots, x_{\tau(n)}\right) \text{ for every permutation } \tau$$

Be approximated by functions of the form

$$\{x_1, \ldots, x_n\} \mapsto \mathcal{N}^{(2)}\left(\sum_{i=1}^{n} \mathcal{N}^{(1)}(x_i)\right)$$

# Throughout we will assume…

$(G, V)$ are **nice**, meaning

- $V$ is a real finite dimensional vector space

e.g., $V = \mathbb{R}^{d \times n}$

- $G$ is a compact matrix group defined by polynomial equations

e.g., $O(d) = \{R \in \mathbb{R}^{d \times d} | RR^T = I_d\}$

- The map $(g, v) \mapsto gv$ is polynomial

e.g., $(R, X) \mapsto RX$

# Standard approach: Invariant Universality via generators of the invariant ring

[Universal Approximations of Invariant Maps by Neural Networks, Yarotsky 2022]

Theorem [Hilbert, 1890]

Let $(V, G)$ be **nice**, then there exist a finite number of invariant polynomials $F_1, \dots, F_N: V \to \mathbb{R}$ such that all invariant polynomials are of the form

$$q(v) = p\big(F_1(v), \dots, F_N(v)\big), \text{ for some } p: \mathbb{R}^N \to \mathbb{R}$$

Remark

$F_1, \dots, F_N$ are called the **generators** of the ring

$$R(V, G) = \{F: V \to \mathbb{R} \text{ are } G \text{ invariant polynomials}\}$$

# Universality of invariant machine learning via generators of the invariant ring

[Universal Approximations of Invariant Maps by Neural Networks, Yarotsky 2022]

Corollary

Let $(V, G)$ be **nice**, and $F_1, \dots F_N$ be generators of the invariant ring. Then any continuous invariant function $f : V \to \mathbb{R}$

can be approximated on compact subsets of $V$ to arbitrary accuracy by

$$\mathcal{N}\big(F_1(v), \dots, F_N(v)\big), \text{ for some neural network } \mathcal{N} : \mathbb{R}^N \to \mathbb{R}$$

# Universality of invariant machine learning via generators of the invariant ring

[Universal Approximations of Invariant Maps by Neural Networks, Yarotsky 2022]

---

Issues

- Can we explicitly compute the generators $F_1, \ldots, F_N$?

  (often yes. In invariant theory this will be called `the first fundamental theorem for $(V, G)'$)

- How does $N$ depend on $\dim(V)$?

  (**often this is very bad**... we will see examples)

- Do we want to use polynomials for approximation?

  (let's ignore this for now)

---

# Point set `Orthogonal Universality via generators'

**Group**: $O(d) = \{R \in \mathbb{R}^{d \times d} | RR^T = I_d\}$

**Action:** $R_*(x_1, \dots, x_n) = (Rx_1, \dots, Rx_n)$

**$\sim n^2$ Generators:**

$$\langle x_i, x_j \rangle \quad 1 \leq i < j \leq n$$

**Universality:** All continuous $O(d)$ invariant functions $f$ can be approximated by functions of the form

$$\mathcal{N}(\langle x_1, x_1 \rangle, \langle x_1, x_2 \rangle, \dots, \langle x_n, x_n \rangle)$$

Where $\mathcal{N}$ is a (fully connected) neural network

# Point set `Special Orthogonal Universality via generators'

**Group**: $SO(d) = \{R \in \mathbb{R}^{d \times d} | RR^T = I_d \ and \ \det(R) = 1\}$

**Action:** $R_*(x_1, \ldots, x_n) = (Rx_1, \ldots, Rx_n)$

$\sim \binom{n}{d}$ **Generators:**

$$\langle x_i, x_j \rangle, \quad 1 \leq i \leq j \leq n \ \text{and} \ \det(x_{i_1}, \ldots, x_{i_d}) \ i_1 < i_2 < \cdots < i_d$$

**Universality:** All continuous $SO(d)$ invariant functions $f$ can be approximated by functions of the form

$$\mathcal{N}(\langle x_1, x_1 \rangle, \langle x_1, x_2 \rangle, \ldots, \langle x_n, x_n \rangle, \det(x_1, \ldots, x_d), \ldots \det(x_{n-d+1}, \ldots, x_n))$$

Where $\mathcal{N}$ is a (fully connected) neural network

# Point set `Permutation Universality via generators'

**Group**: $S_n = \{permutations \ \tau: \{1, \dots, n\} \to \{1, \dots, n\}\}$

**Action:** $\tau_*(x_1, \dots, x_n) = \left(x_{\tau^{-1}(1)}, \dots, x_{\tau^{-1}(n)}\right)$

$\boldsymbol{m(n, d)} = \binom{n+d}{d}$ **Generators:**

$(x_1, \dots, x_n) \mapsto \sum_{i=1}^{n} p_j(x_i)$ where $p_1, \dots, p_m$ form a basis for the space of polynomials of degree $\leq n$ in $d$ variables

**Universality:** All continuous $S_n$ invariant functions can be approximated by

$$\mathcal{N}\left(\sum_{i=1}^{n} p_1(x_i), \sum_{i=1}^{n} p_2(x_i), \dots \sum_{i=1}^{n} p_m(x_n)\right)$$

$$\text{Or } \mathcal{N}^{(2)}\left(\sum_{i=1}^{n} \mathcal{N}^{(1)}(x_i)\right)$$

# Number of generators for point set actions

| Group action on $\mathbb{R}^{d \times n}$ | Num of generators |
|---|---|
| $O(d)$ | $\sim n^2$ |
| $SO(d)$ | $\sim \binom{n}{d}$ |
| $S_n$ | $\binom{n+d}{d}$ |

# Universality of invariant machine learning via ~~generating~~ *separating* invariants

Advocates: [Complete set of translation invariant measurements with Lipschitz bounds, Cahill et al. 2020]
[Group invariant max-filtering, Cahill et al. 2022]
[Low Dimensional Invariant Embeddings for Universal Geometric Learning, **Dym** and Gortler 2022]

**Definition (Separating invariants)**

Let $G$ be a group acting on $V$. We say that $H_1, \ldots, H_m : V \to \mathbb{R}$ are $(V, G)$ separating invariants if

- **Invariant: if** $u =_G v$ **then** $H_i(u) = H_i(v), \forall i = 1, \ldots, m$

- **Separating: if** $H_i(v) = H_i(u), \forall i = 1, \ldots, m$ **then** $v =_G u$

**Invariance** means that $V/_G \ni [v] \mapsto (H_1(v), \ldots, H_m(v))$ is well defined

**Separating** means that it is injective on $V/_G$

Example: $G = O(2)$ acts on $V = \mathbb{R}^{2 \times 2}$ via $R_*(x_1, x_2) = (Rx_1, Rx_2)$
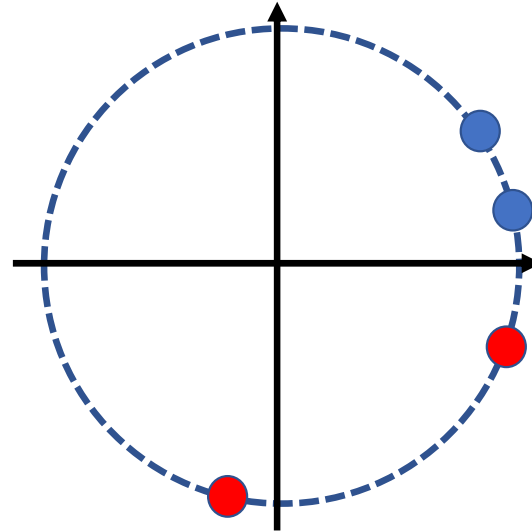
What invariants can we suggest? Are they separating?

How about:

$H_1(x_1, x_2) = \|x_1\|$ and $H_2(x_1, x_2) = \|x_2\|$?

We get separation by adding

$H_3(x_1, x_2) = \|x_1 - x_2\|$

# Separation vs generation: sufficiency for universality

We saw

*and $H_1, \ldots, H_m$ be continuous separating invariants*

Let $(V, G)$ be **nice**, and ~~$F_1, \ldots F_N$ be generators of the invariant ring~~. Then any continuous invariant function $f: V \to \mathbb{R}$

can be approximated on compact subsets of $V$ to arbitrary accuracy by

~~$\mathcal{N}\big(F_1(v), \ldots, F_N(v)\big)$~~, for some neural network $\mathcal{N}: \mathbb{R}^N \to \mathbb{R}$

$\mathcal{N}(H_1(v), \ldots H_m(v))$

**Remark:** This in fact implies the generator-based theorem, since **generators are always separators**

# Separation vs generation: cardinality
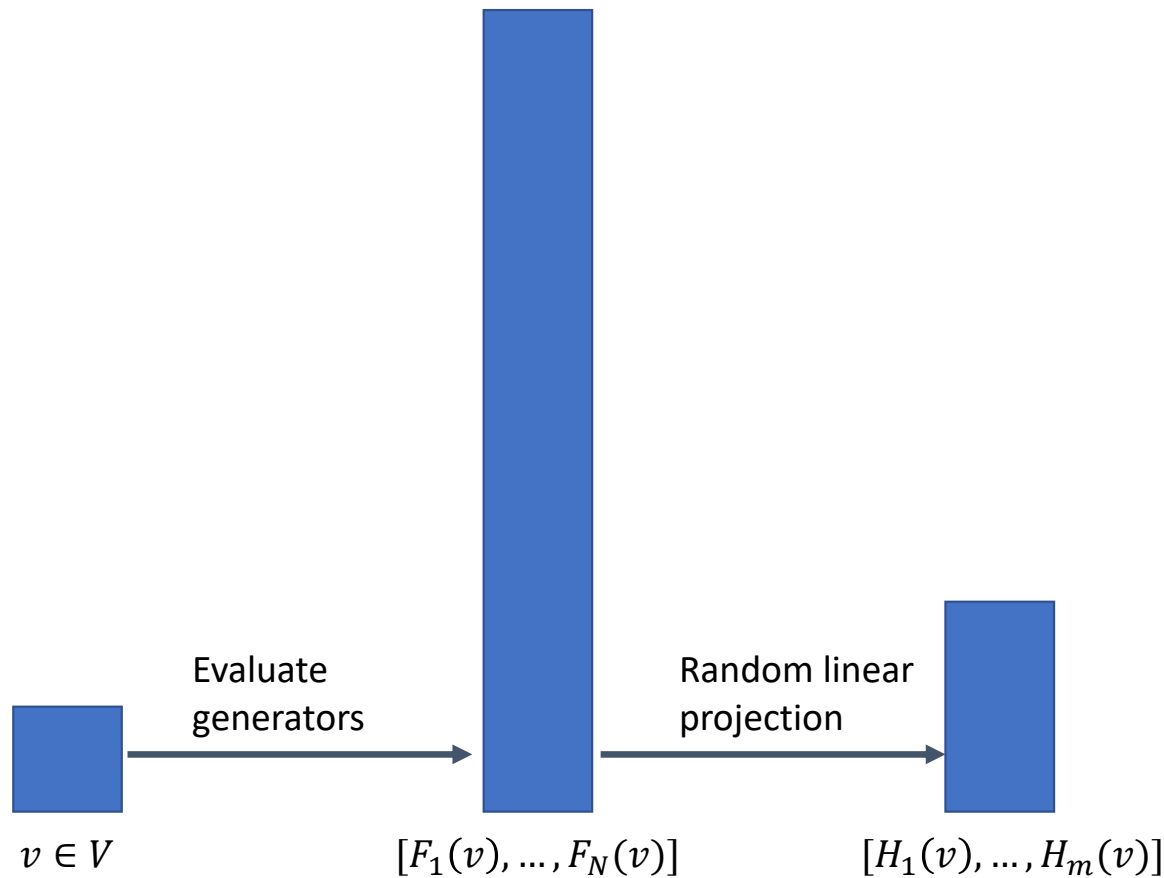
Theorem [E. S. Dufresne 2008]

If $(V, G)$ are **nice**, then there always exist **polynomial** separating invariants $H_1, \ldots, H_m : V \to \mathbb{R}$ of cardinality

$$m = 2 \dim(V) + 1$$

# Partial solution: low dimensional-separation via generation+`linear compression'

Evaluate generators

Random linear projection

$v \in V$   $[F_1(v), \dots, F_N(v)]$   $[H_1(v), \dots, H_m(v)]$

# Intermediate conclusions

| Group action on $\mathbb{R}^{d \times n}$ | Num of generators |
|---|---|
| $O(d)$ | $n^2$ |
| $SO(d)$ | $n^2 + \binom{n}{d}$ |
| $S_n$ | $\binom{n+d}{d}$ |

<span style="color:red">Let's start here</span>

Can we do better?   <span style="color:red">Yes</span>

# Efficient invariants: $SO(d)$

Example: $R \in SO(d)$ acts on $\mathbf{X} = (x_1, \dots, x_n) \in \mathbb{R}^{d \times n}$ $(d < n)$

$$R_*(x_1, \dots x_n) = (Rx_1, \dots, Rx_n)$$

Generators: $\sim \binom{n}{d}$

$$\left| x_i - x_j \right|^2 \ and \ \left| x_j \right|^2 \ and \ \det\left( x_{i_1}, \dots x_{i_d} \right)$$

Continuous family of separating invariants:

$$H(x_1, \dots, x_n; w, W) = |w_1 x_1 + \dots + w_n x_n|^2 + \det(XW)$$

Random separators: For almost all $w^{(1)}, W^{(1)} \dots, w^{(m)}, W^{(m)}, m = 2nd + 1$

$H\left( x_1, \dots, x_n; w^{(i)}, W^{(i)} \right)$ are invariant and separating!!!

$|x_i - x_j|^2$ *and* $|x_j|^2$ and $\det(x_{i_1}, \ldots x_{i_d})$

| Group action on $\mathbb{R}^{d \times n}$ | Num of generators | Num of separators | Complexity per separator? | |
|---|---|---|---|---|
| $O(d)$ | $n^2$ | $2n \cdot d + 1$ | | |
| $SO(d)$ | $n^2 + \binom{n}{d}$ | $2n \cdot d + 1$ | $nd^2$ | $\lvert w_1 x_1 + \cdots + w_n x_n \rvert^2 + \det(XW)$ |
| $S_n$ | $\binom{n+d}{d}$ | $2n \cdot d + 1$ | | |

# Efficient Invariants: SO(d) and beyond

For the action of $SO(d)$ on $\mathbb{R}^{d \times n}$, the following is a *continuous family of separating invariants*

$$H(\mathbf{x_1}, \dots, \mathbf{x_n}; \mathbf{w}, \mathbf{W}) = |\mathbf{w_1}\mathbf{x_1} + \cdots + \mathbf{w_n}\mathbf{x_n}|^2 + \det(\mathbf{XW})$$

Definition: Let $(V, G)$ be **nice**. We sat that a function $H: V \times \mathbb{R}^{d_W} \to \mathbb{R}$ is a *continuous family of separating invariants* if it satisfies the following conditions:

- **Invariance:** If $v =_G v'$ then $H(v; w) = H(v'; w)$ for all $w \in \mathbb{R}^{d_W}$

- **Separation:** If $v \neq_G v'$ then there exists $w \in \mathbb{R}^{d_W}$ such that $H(v; w) \neq H(v'; w)$

# Finite Witness Theorem

Finite Witness Theorem [Dym and Gortler 2022] (weakened version):

Let $(V, G)$ be **nice**. Let $H: V \times \mathbb{R}^{d_w} \to \mathbb{R}$ be a family of separating **polynomial** invariants.

Set $m = 2 \dim(V) + 1$. Then for Lebesgue almost every $w^{(1)}, \dots, w^{(m)} \in \mathbb{R}^{d_w}$, the functions $H_1, \dots H_m$ defined by

$$H_i(v) = H(v; w^{(i)})$$

are separating invariants.

Remarks

- Cardinality is often not optimal

- Proof idea comes from [On signal reconstruction without phase, Balan, Casazza and Edidin 2006] relies on Real Algebraic Geometry

Finite Witness Theorem [Dym and Gortler 2022] (weakened version):

Let $(V, G)$ be **nice**. Let $H: V \times \mathbb{R}^{d_w} \to \mathbb{R}$ be a family of separating **polynomial** invariants.

Set $m = 2 \dim(V) + 1$. Then for Lebesgue almost every $w^{(1)}, \dots, w^{(m)} \in \mathbb{R}^{d_w}$, the functions $H_1, \dots H_m$ defined by

$$H_i(v) = H(v; w^{(i)})$$

are separating invariants.

---

Proof idea:

- Consider the `**lifted bad set'**

$$B = \left\{ (v, v', w^{(1)}, \dots w^{(m)}) \in V \times V \times \mathbb{R}^{d_w \times m} \,\middle|\, v \neq_G v' \ \textbf{but} \ H(v; w^{(i)}) = H(v'; w^{(i)}), \forall i = 1 \dots m \right\}$$

- This set is a subset of a $\textbf{2} \dim(V) + md_w$ dimensional vector space defined by $m$ equations

  "therefore" $\dim(B) = md_w + (2\dim(V) - m) = md_w - 1$

- The dimension of the `projected bad set' is no larger

$$B_{proj} = \left\{ (w^{(1)}, \dots w^{(m)}) \in \mathbb{R}^{d_w \times m} \,\middle|\, \exists (v, v') \ s.t. \ v \neq_G v' \ \textbf{but} \ H(v; w^{(i)}) = H(v'; w^{(i)}), \forall i = 1 \dots m \right\}$$

- $dim(B_{proj}) = md_w - 1 < dim(\mathbb{R}^{d_w \times m})$

- Most $(w^{(1)}, \dots w^{(m)})$ are not in $B_{proj}$, and so are separating

Finite Witness Theorem [Dym and Gortler 2022] (weakened version):

Let $(V, G)$ be **nice**. Let $H: V \times \mathbb{R}^{d_W} \to \mathbb{R}$ be a continuous family of separating polynomial invariants.

Set $m = 2\dim(V) + 1$. Then for Lebesgue almost every $w^{(1)}, \ldots, w^{(m)} \in \mathbb{R}^{d_W}$, the functions $H_1, \ldots H_m$ defined by

$$H_i(v) = H(v; w^{(i)})$$

are separating invariants.

<span style="color:red">Real algebraic geometry</span>

<span style="color:red">Full Proof</span>

Proof ~~idea~~ (inspired by phase retrieval paper):

- Consider the `**lifted bad set'**

$$\boldsymbol{B} = \left\{ \left(\boldsymbol{v}, \boldsymbol{v}', \boldsymbol{w}^{(1)}, \ldots \boldsymbol{w}^{(m)}\right) \in \boldsymbol{V} \times \boldsymbol{V} \times \mathbb{R}^{d_w \times m} \middle| \boldsymbol{v} \neq_{\boldsymbol{G}} \boldsymbol{v}' \text{ } \boldsymbol{but} \text{ } \boldsymbol{H}(\boldsymbol{v}; \boldsymbol{w}^{(i)}) = \boldsymbol{H}(\boldsymbol{v}'; \boldsymbol{w}^{(i)}), \forall i = 1 \ldots \boldsymbol{m} \right\}$$

- This set is a subset of a $\boldsymbol{2\dim(V) + m d_w}$ dimensional vector space defined by $\boldsymbol{m}$ equations

<span style="color:red">"therefore"</span> $\boldsymbol{\dim(B) = m d_w + (2\dim(V) - m) = m d_w - 1}$

- The dimension of the `projected bad set' is no larger

$$\boldsymbol{B_{proj}} = \left\{ \left(\boldsymbol{w}^{(1)}, \ldots \boldsymbol{w}^{(m)}\right) \in \mathbb{R}^{d_w \times m} \middle| \exists (\boldsymbol{v}, \boldsymbol{v}') \text{ } \boldsymbol{s.t.} \text{ } \boldsymbol{v} \neq_{\boldsymbol{G}} \boldsymbol{v}' \text{ } \boldsymbol{but} \text{ } \boldsymbol{H}(\boldsymbol{v}; \boldsymbol{w}^{(i)}) = \boldsymbol{H}(\boldsymbol{v}'; \boldsymbol{w}^{(i)}), \forall i = 1 \ldots \boldsymbol{m} \right\}$$

- $\boldsymbol{dim(B_{proj}) = m d_w - 1 < dim(\mathbb{R}^{d_w \times m})}$

- Most $\left(\boldsymbol{w}^{(1)}, \ldots \boldsymbol{w}^{(m)}\right)$ are not in $\boldsymbol{B_{proj}}$, and so are separating

# Finite Witness Theorem-Applications

| Group action on $\mathbb{R}^{d \times n}$ | Num of generators | Num of separators | Complexity per separator? |
|---|---|---|---|
| $O(d)$ | $n^2$ | $2n \cdot d + 1$ | $n \cdot d$ |
| $SO(d)$ | $n^2 + \binom{n}{d}$ | $2n \cdot d + 1$ | $n \cdot d^2$ |
| $S_n$ | $\binom{n+d}{d}$ | $2n \cdot d + 1$ | $n \cdot log(n)$ |

# Recent work- Analytic Finite Witness Theorem

Analtyic Finite Witness Theorem [Amir, Gortler, Avni, Ravina, Dym  2023] (weakened version): **Analytic**

Let $(V, G)$ be **nice**. Let $H: V \times \mathbb{R}^{d_w} \to \mathbb{R}$ be a continuous family of separating **polynomial** invariants.

Set $m = 2 \dim(V) + 1$. Then for Lebesgue almost every $w^{(1)}, \dots, w^{(m)} \in \mathbb{R}^{d_w}$, the functions $H_1, \dots H_m$

defined by

$$H_i(v) = H(v; w^{(i)})$$

are separating invariants.

`Proof'

 Real Algebraic Geometry $\mapsto$ Real analytic geometry, o-minimal systems and related concepts

# Application: Permutation invariant networks (with analytic activations)

Theorem [Amir, Gortler, Avni, Ravina, Dym 2023]

Let $d, n$ be natural numbers and set $m = 2nd + 1$.

If $\sigma: \mathbb{R} \mapsto \mathbb{R}$ is **analytic** and not polynomial, then for Lebesgue almost every $A \in \mathbb{R}^{m \times d}$ and $b \in \mathbb{R}^m$

the permutation invariant function

$$\mathbb{R}^{d \times n} \ni (x_1, \ldots, x_n) \mapsto \sum_{i=1}^{n} \sigma(Ax_i + b)$$

is separating

# Finite Witness Theorem

Analtyic Finite Witness Theorem-stronger (but not strongest) version

Let $(V, G)$ be **nice**. Let $H: V \times \mathbb{R}^{d_w} \to \mathbb{R}$ be a continuous family of separating **analytic** invariants. Set $m = 2 \dim(V) + 1$. Then for Lebesgue almost every $w^{(1)}, \ldots, w^{(m)} \in \mathbb{R}^{d_w}$, the functions $H_1, \ldots H_m$ defined by

$$H_i(v) = H\big(v; w^{(i)}\big)$$

are separating invariants.

Can be a low dimensional subset of some higher dimensional vector space, providing it is `reasonable'
e.g., a countable union of sets defined by polynomial and analytic equalities and inequalities
Or image of these sets under an analytic functions

# Adding to the table…

| Group action on $\mathbb{R}^{d \times n}$ | Num of generators | Num of separators | Complexity per separator? |
|---|---|---|---|
| $O(d)$ | $n^2$ | $2n \cdot d + 1$ | $n \cdot d$ |
| $SO(d)$ | $n^2 + \binom{n}{d}$ | $2n \cdot d + 1$ | $n \cdot d^2$ |
| $S_n$ | $\binom{n+d}{d}$ | $2n \cdot d + 1$ | $n(d + \log(n))$ |
| $O(d) \times S_n$ | ? | $2n \cdot d + 1$ | $\mathbf{n^d}$ |
| $SO(d) \times S_n$ | ? | $2n \cdot d + 1$ | $\mathbf{n^d}$ |

# Parting questions

- Separating invariants are injective mappings $f: V/G \to R^m$. Do they preserve distances?

- Separating invariants for surfaces? (one example: conformal welding)

# Collaborators



Steven J. Gortler
Harvard

Tal Amir

Snir Hordan

Technion

Ilai Avni

Ravina Ravina

[Neural Injective Functions for Multisets, Measures and Graphs
via a Finite Witness Theorem.
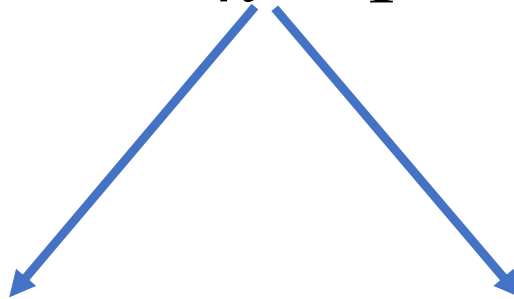Amir, Gortler, Avni, Ravina and Dym 2023]

# Thank you!

[Low Dimensional Invariant Embeddings for Universal
Geometric Learning
Dym and  Gortler 2022]

[Complete Neural Networks for Euclidean Graphs
Hordan, Amir, Gortler, and Dym 2023]

$$\text{O(d)} \times S_n \text{ separation}$$

Generically separating
Geometric Message Passing

Fully separating
$(d-1)$ order geometric message
passing

# Geometric message passing

e.g. EGNN [E(n) equivariant graph neural networks, Sattoras et al. 2021]

Input: $x_1, \ldots, x_n \in \mathbb{R}^d$

set $h_1^{(0)}, \ldots, h_n^{(0)} = 0$

$h_i^{(t)} = f_{agg}\left(h_i^{(t-1)}, \left\{h_j^{(t-1)}, |x_i - x_j|, j = 1, \ldots n\right\}\right)$ (repeat $T$ times)

$h_{global}(x_1, \ldots, x_n) = f_{readout}(\{h_1^{(T)}, \ldots, h_n^{(T)}\})$

$h_{global}(x_1, \ldots, x_n)$ is $O(d) \times S_n$ invariant

# Geometric message passing-separation

e.g. EGNN [E(n) equivariant graph neural networks, Sattoras et al. 2021]

Input: $x_1, \ldots, x_n \in \mathbb{R}^d$

set $h_1^{(0)}, \ldots, h_n^{(0)} = 0$

$h_i^{(t)} = f_{agg}\left(h_i^{(t-1)}, \left\{h_j^{(t-1)}, |x_i - x_j|, j = 1, \ldots n\right\}\right)$ (repeat $T$ times)

$h_{global} = f_{readout}(\{h_1^{(T)}, \ldots, h_n^{(T)}\})$

Permutation invariant and separating

# `Hard' to separate

[Incompleteness of Atomic Structure Representations. Physical Review Letters Pozdynakov et al. 2020]

| 0 | 1 | 2 | -2 | -1 |
|---|---|---|----|----|
| 1 | 1 | 0 | 0 | -1 |
| 1 | 0 | 2 | -2 | 0 |

$$\neq_G$$

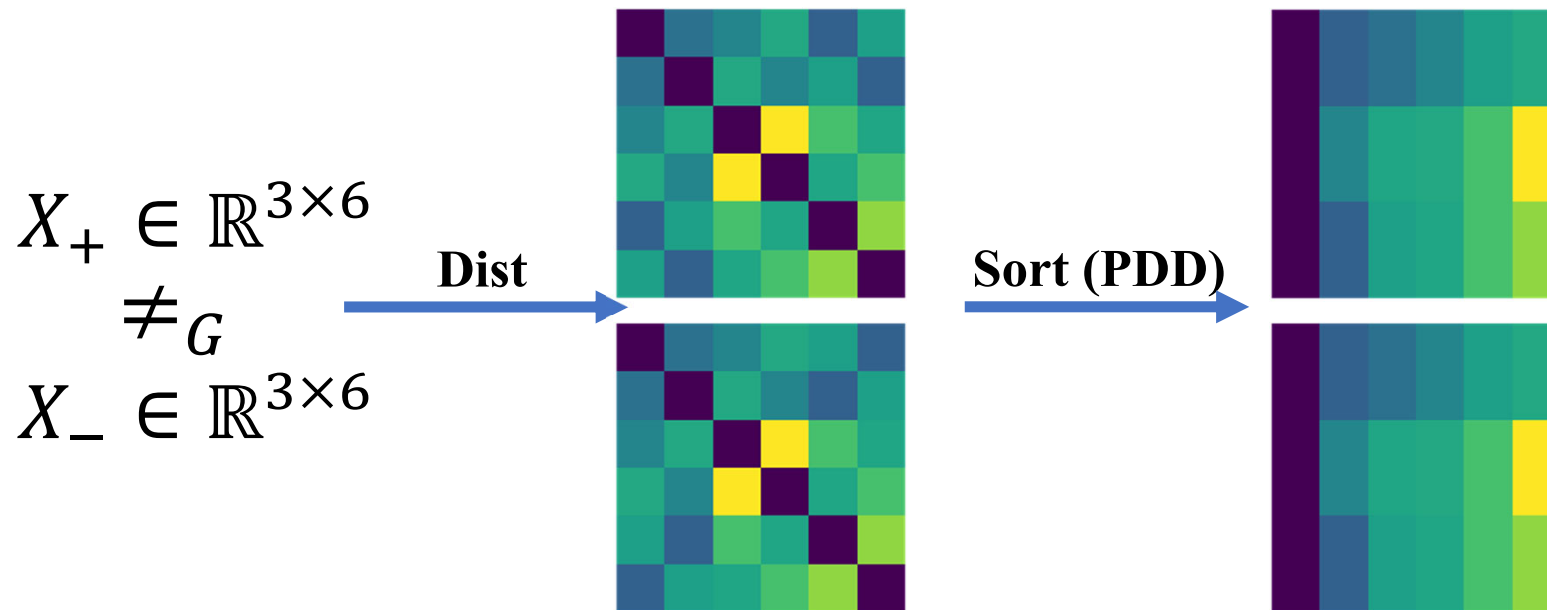| 0 | 1 | 2 | -2 | -1 |
|---|---|---|----|----|
| 1 | 1 | 0 | 0 | -1 |
| -1 | 0 | 2 | -2 | 0 |

**Dist** →

**Sort** →



- **Cannot** be separated by MPNN with $T = 1$

- **Can** be separated by MPNN with $T \geq 2$

`Harder'

[Incompleteness of graph neural networks
for points clouds in three dimensions, Pozdnyakov and Ceriotti 2022]

$$X_+ \in \mathbb{R}^{3 \times 6}$$
$$\neq_G$$
$$X_- \in \mathbb{R}^{3 \times 6}$$

**Dist** →

**Sort (PDD)** →

Cannot be separated by MPNN for any $T$

# Geometric K-order message passing

[Sign and Basis Invariant Networks for Spectral Graph Representation Learning, Lim et al. 2022]
[Is distance matrix enough for geometric deep learning, Li et al. 2023]

Assume K= 3 for notation simplicity

$$h^{(0)}(i,j,k)(X) = \begin{pmatrix} \langle x_i, x_i \rangle & \langle x_i, x_j \rangle & \langle x_i, x_k \rangle \\ \langle x_j, x_i \rangle & \langle x_j, x_j \rangle & \langle x_j, x_k \rangle \\ \langle x_k, x_i \rangle & \langle x_k, x_j \rangle & \langle x_k, x_k \rangle \end{pmatrix}$$

$$h^{(t)}(i,j,k)(X) = f_{agg}\left( h^{(t-1)}(i,j,k), \left\{ \begin{pmatrix} h^{(t-1)}(s,j,k) \\ h^{(t-1)}(i,s,k) \\ h^{(t-1)}(i,j,s) \end{pmatrix}, \quad s = 1, \dots, n \right\} \right)$$

$$h_{global} = f_{readout}\left\{ h^{(T)}(i,j,k) \big| (i,j,k) \in [n]^3 \right\}$$

Permutation invariant
+separating

<u>Theorem [Hordan, Amir, Gortler, Dym, 2023]</u>

For every $X, Y \in \mathbb{R}^{d \times n}$ we have that the d-order message passing with $T = 1$ is separating:

It gives the same output $h_{global}(X) = h_{global}(Y)$ if and only if $X, Y$ are related by a permutation

and orthogonal transformation.

A modified $d - 1$ message passing algorithm is also separating
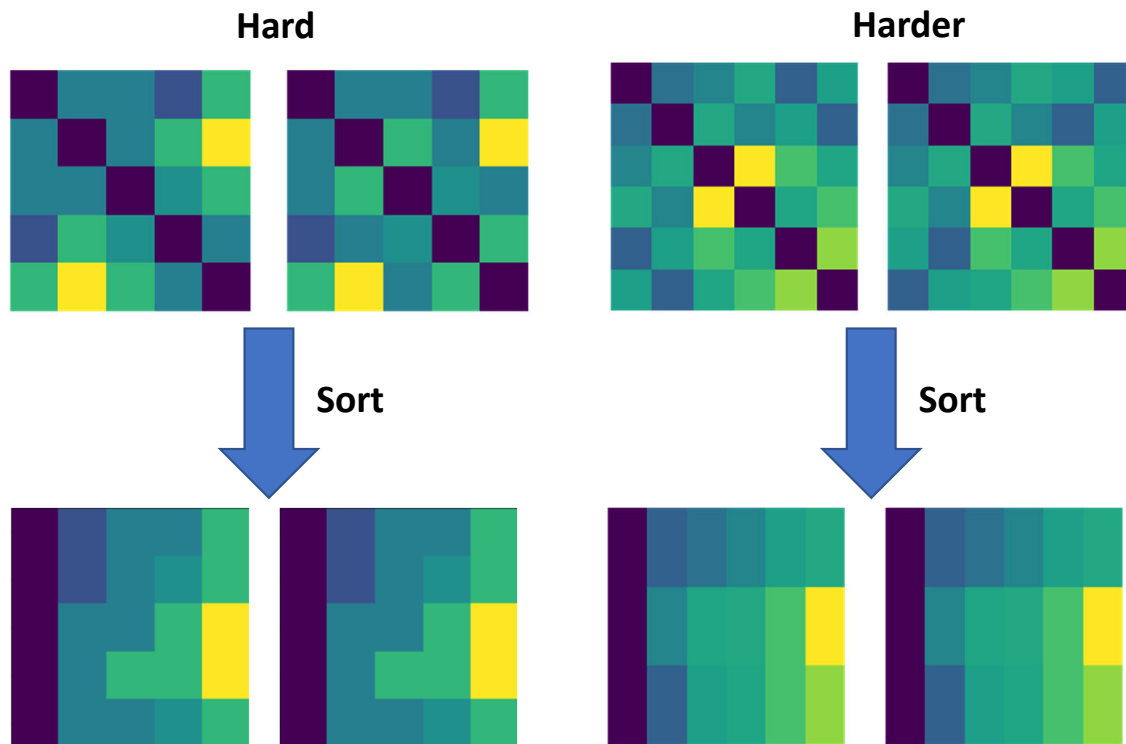
<u>Theorem [Rose et al. 2023]</u>

The original $d - 1$ message passing algorithm is also separating

# Complexity

- Full $O(d) \times S_n$ separation with $(d-1) - WL$ requires computing $2nd + 1$ invariants with computational complexity of $n^d$ each, using our permutation invariant separating functions

- This also uses the dependence of the theorem on intrinsic dimension. Considering extrinsic dimension only would lead to exponential blowup

# Separation experiment



| | Hard | Harder |
|---:|:---|:---|
| MPNN | Yes | No |
| $(d-1)$ MPNN | Yes | Yes |

# Separation of existing invariant architectures:

$O(d) \times S_n$ Invariant architectures

$(d-1)MPNN$     $MPNN$

| Point Clouds | GramNet | GeoEGNN | EGNN | LinearEGNN | MACE | TFN | DimeNet | GVPGNN |
|---|---|---|---|---|---|---|---|---|
| Hard1[2] | 1.0 | 0.998 | 0.5 | 1.0 | 1.0 | 0.5 | 1.0 | 1.0 |
| Hard2 [2] | 1.0 | 0.97 | 0.5 | 1.0 | 1.0 | 0.5 | 1.0 | 1.0 |
| Hard3 [2] | 1.0 | 0.85 | 0.5 | 1.0 | 1.0 | 0.55 | 1.0 | 1.0 |
| Harder [1] | 1.0 | 0.899 | 0.5 | 0.5 | 1.0 | 0.5 | 1.0 | 1.0 |
| Cholesky dim=6 | 1.0 | Irrelevant | 0.5 | 0.5 | 1.0 | Irrelevant | Irrelevant | Irrelevant |
| Cholesky dim=8 | 1.0 | Irrelevant | 0.5 | 0.5 | 1.0 | Irrelevant | Irrelevant | Irrelevant |
| Cholesky dim=12 | N/A | Irrelevant | 0.5 | 0.5 | 0.5 | Irrelevant | Irrelevant | Irrelevant |

Dataset composed of two point clouds which are hard to separate+rotations+permutations+noise

# We didn't discuss…

- **Generic separation:** Separation up to a set of measure zero. Need only $\dim(V) + 1$ invariants

- **Stability:** Invariant and separating $H: V \to \mathbb{R}^m$ can be identified with $H: {}^V/_G \to \mathbb{R}^m$ injective.

  Is $H$ bi-Lipschitz with respect to

$$d([v], [v']) = \min_{g \in G} ||gv - v'||$$

[Permutation invariant representations with applications to graph deep learning, Balan Haghani Singh 2022]
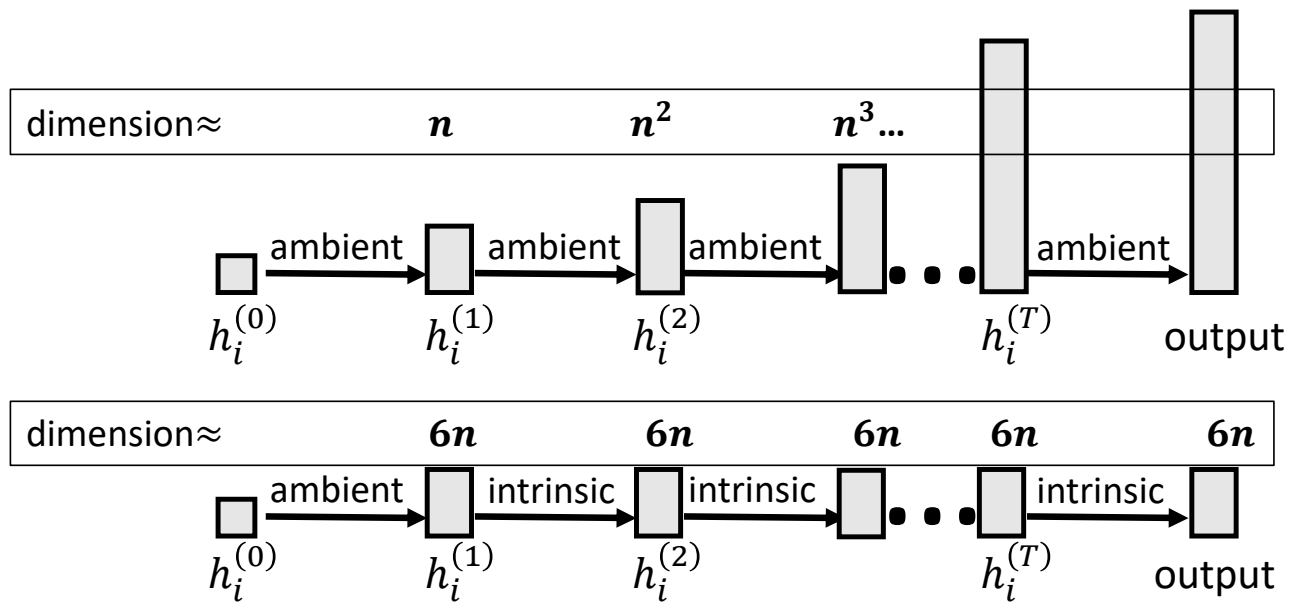[Group-invariant max flitering, Cahill Iverson Dixon and Packer]

TODO

$$h_i^{(t)} = f_{agg}\left(h_i^{(t-1)}, \left\{h_j^{(t-1)}, |x_i - x_j|, j = 1, \ldots n\right\}\right) \text{ (repeat } T \text{ times)}$$

Permutation invariant and separating

| dimension≈ | | $\boldsymbol{n}$ | | $\boldsymbol{n^2}$ | | $\boldsymbol{n^3}\ldots$ | | |
|---|---|---|---|---|---|---|---|---|

ambient → ambient → ambient → $\cdots$ → ambient →

$h_i^{(0)}$  $h_i^{(1)}$  $h_i^{(2)}$  $h_i^{(T)}$  output

| dimension≈ | | $\boldsymbol{6n}$ | $\boldsymbol{6n}$ | $\boldsymbol{6n}$ | $\boldsymbol{6n}$ | $\boldsymbol{6n}$ |
|---|---|---|---|---|---|---|

ambient → intrinsic → intrinsic → $\cdots$ → intrinsic →

$h_i^{(0)}$  $h_i^{(1)}$  $h_i^{(2)}$  $h_i^{(T)}$  output

# Separation of existing architectures: (when) does it happen?

$O(d) \times S_n$ Invariant architectures

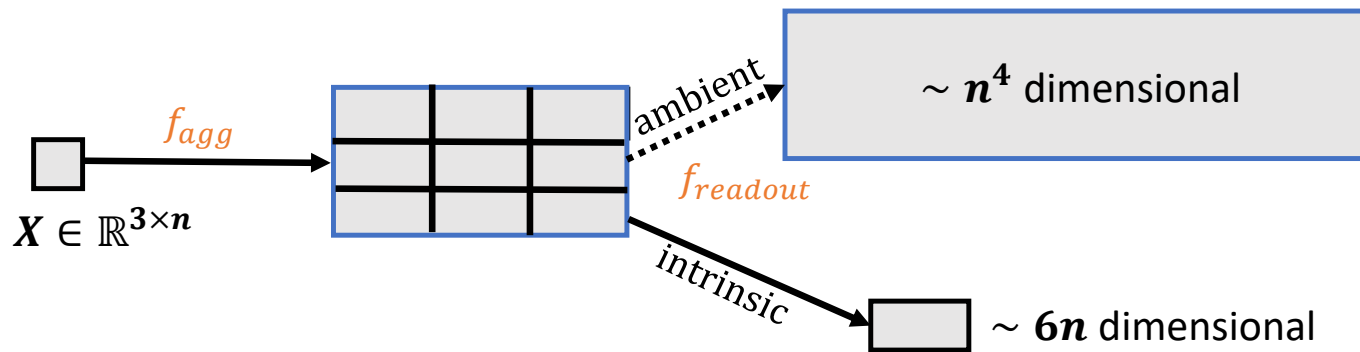| Theoretical separation | Yes (ours) | Yes (ours) | No | No | ? | Sort of | ? | ? |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Point Clouds | GramNet | GeoEGNN | EGNN | LinearEGNN | MACE | TFN | DimeNet | GVPGNN |
| Hard1[2] | 1.0 | 0.998 | 0.5 | 1.0 | 1.0 | 0.5 | 1.0 | 1.0 |
| Hard2 [2] | 1.0 | 0.97 | 0.5 | 1.0 | 1.0 | 0.5 | 1.0 | 1.0 |
| Hard3 [2] | 1.0 | 0.85 | 0.5 | 1.0 | 1.0 | 0.55 | 1.0 | 1.0 |
| Harder [1] | 1.0 | 0.899 | 0.5 | 0.5 | 1.0 | 0.5 | 1.0 | 1.0 |
| Cholesky dim=6 | 1.0 | Irrelevant | 0.5 | 0.5 | 1.0 | Irrelevant | Irrelevant | Irrelevant |
| Cholesky dim=8 | 1.0 | Irrelevant | 0.5 | 0.5 | 1.0 | Irrelevant | Irrelevant | Irrelevant |
| Cholesky dim=12 | N/A | Irrelevant | 0.5 | 0.5 | 0.5 | Irrelevant | Irrelevant | Irrelevant |

Dataset composed of two point clouds which are hard to separate+rotations+permutations+noise

# Proof of theorem (intuition)

# Full $O(d) \times S_n$ separation (1): Cardinality

$$h^{(1)}(i,j,k)(X) = f_{agg}\left(h^{(0)}(i,j,k), \left\{\begin{pmatrix} h^{(0)}(s,j,k) \\ h^{(0)}(i,s,k) \\ h^{(0)}(i,j,s) \end{pmatrix}, \quad s = 1, \ldots, n\right\}\right)$$

$$h_{global} = f_{readout}\left\{h^{(1)}(i,j,k) \big| (i,j,k) \in [n]^3\right\}$$

# Phase retrieval

# Better solution: imported from phase retrieval

**Phase retrieval:** we want to reconstruct a signal $z \in \mathbb{C}^n$ from phaseless linear measurements

$$H_i(z) = \left| \left\langle w^{(i)}, z \right\rangle \right|^2, i = 1, \ldots m$$

$S^1$ **invariance:** For all $\theta$ we have that $\left| H_i\left(e^{i\theta} z\right) \right| = |H_i(z)|$ so we can only hope for reconstruction up to a global phase factor, that is

$$H_i(z) = H_i(\hat{z}) \quad \Longleftrightarrow \quad z = e^{i\theta} \hat{z} \text{ for some } \theta$$

In other words, we would like $H_1, \ldots H_m$ to be separating

# Better solution: imported from phase retrieval

**Theorem** [On signal reconstruction without phase, Balan, Casazza and Edidin 2006]

If $m = 4n - 2$ then for Lebesgue almost all $w^{(1)}, \ldots, w^{(m)} \in \mathbb{R}^n$ the functions $H_1, \ldots H_m$ defined by

$$H_i(z) = \left| \left\langle w^{(i)}, z \right\rangle \right|^2, i = 1, \ldots m$$

are separating with respect to the action of $S^1$ on $\mathbb{C}^n$

**Remark:** In our context we think of $V = \mathbb{C}^n$ is a real vector space of dimension $2n$. So

$$m = 4n - 2 < 2 \dim(V) + 1 = 4n + 1$$

**Remark:** Note that all invariant are obtained by taking sample of $H(z; w)$ which is polynomial in both its argument $z$ and its parameters $w$

# Separation vs. generation for phase retrieval

**Theorem** [On signal reconstruction without phase, Balan, Casazza and Edidin 2006]

If $m = 4n - 2$ then for Lebesgue almost all $w^{(1)}, \ldots, w^{(m)} \in \mathbb{R}^n$ the functions $H_1, \ldots H_m$ defined by

$$H_i(\mathbf{z}) = \left| \left\langle \mathbf{w}^{(i)}, \mathbf{z} \right\rangle \right|^2, i = 1, \ldots m$$

are separating with respect to the action of $S^1$ on $\mathbb{C}^n$

In contrast, there are $\sim n^2$ generators for the ring of invariant polynomials:

$$H_{s,t}(z_1, \ldots, z_n) = z_s \overline{z_t}$$

# Invariant universality rephrased

Assume $G$ acts on $V$

**Orbit:** $$[v] = \{w \in V | \exists g \in G,\ w = gv\ \}$$

**Quotient space:** $$V/_G = \{[v] |\ v \in V\}$$

If $f: V \to Y$ is invariant then it induces a well-defined $\hat{f}: V/_G \to Y$ via

$$\hat{f}([v]) = f(v)$$

# Invariant universality via Invariant embeddings

If $\widehat{F}: {}^V\!/_G \to \mathbb{R}^m$ is invariant and *injective, then any $\hat{f}: {}^V\!/_G \to Y$* is of the form

$$\hat{f}([v]) = h \circ \widehat{F}([v]), \text{ for an appropriate } h: \mathbb{R}^m \to Y$$

On the image of $\widehat{F}$ we have $h = \hat{f} \circ \left(\widehat{F}\right)^{-1}$

**Goal:** Find injective $\widehat{F}: {}^V\!/_G \to \mathbb{R}^m$

# Invariant embeddings and separating invariants

**Goal:** Find injective $\widehat{F}: V/_G \to \mathbb{R}^m$



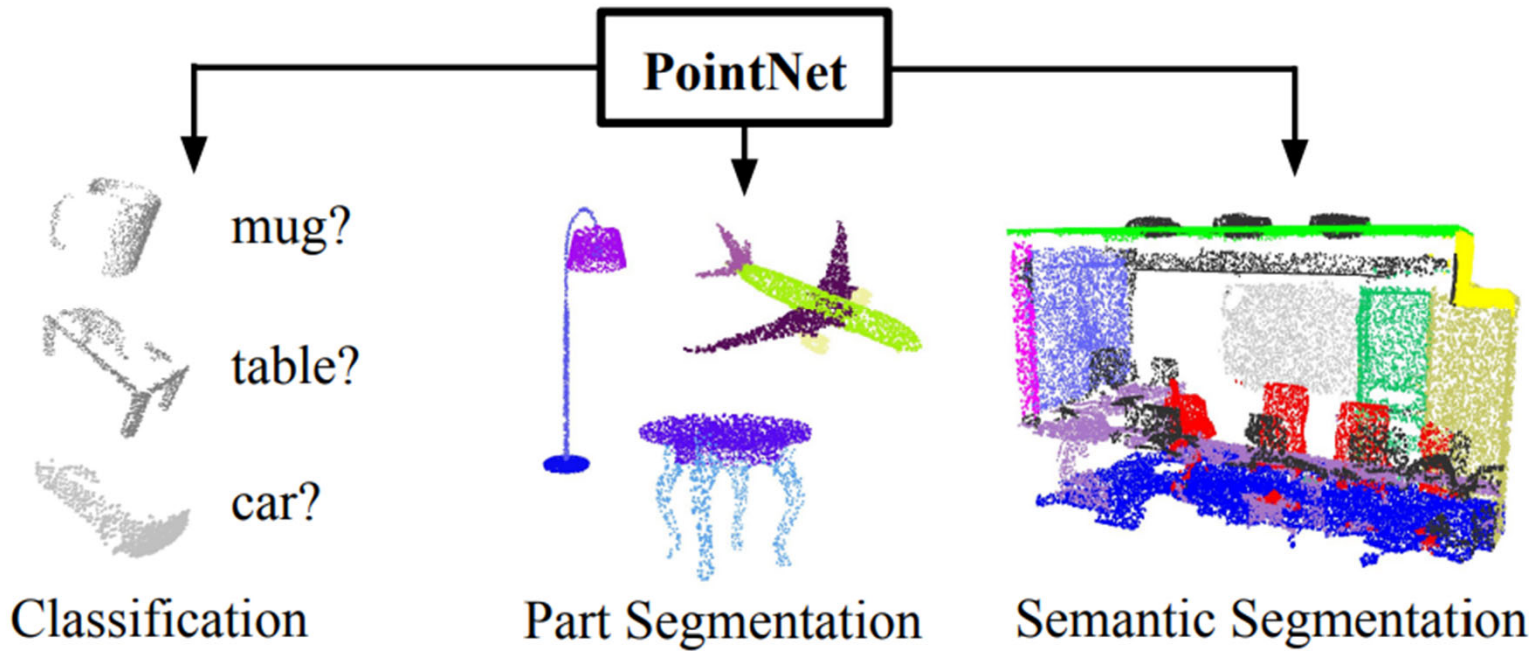**Goal:** Find invariant and separating $F: V \to \mathbb{R}^m$

- **Invariant:** if $[w] = [v]$ then $F(v) = F(w)$

- **Separating:** If $F(v) = F(w)$ then $[w] = [v]$

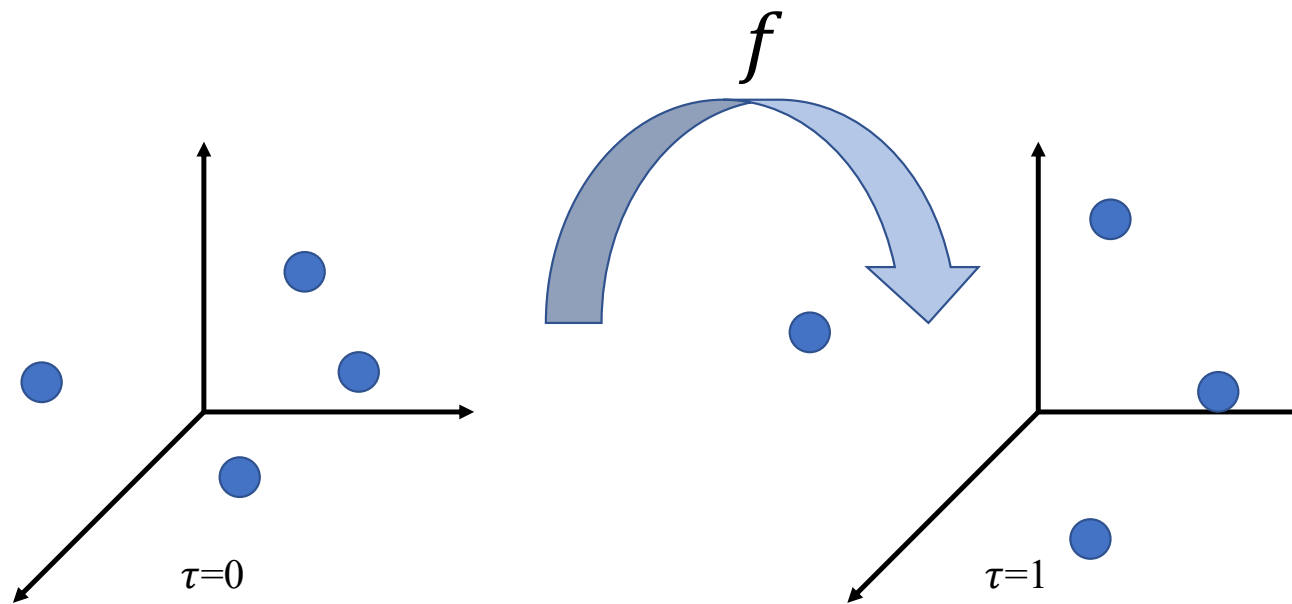# Conclusion: things we didn't discuss

- Stability
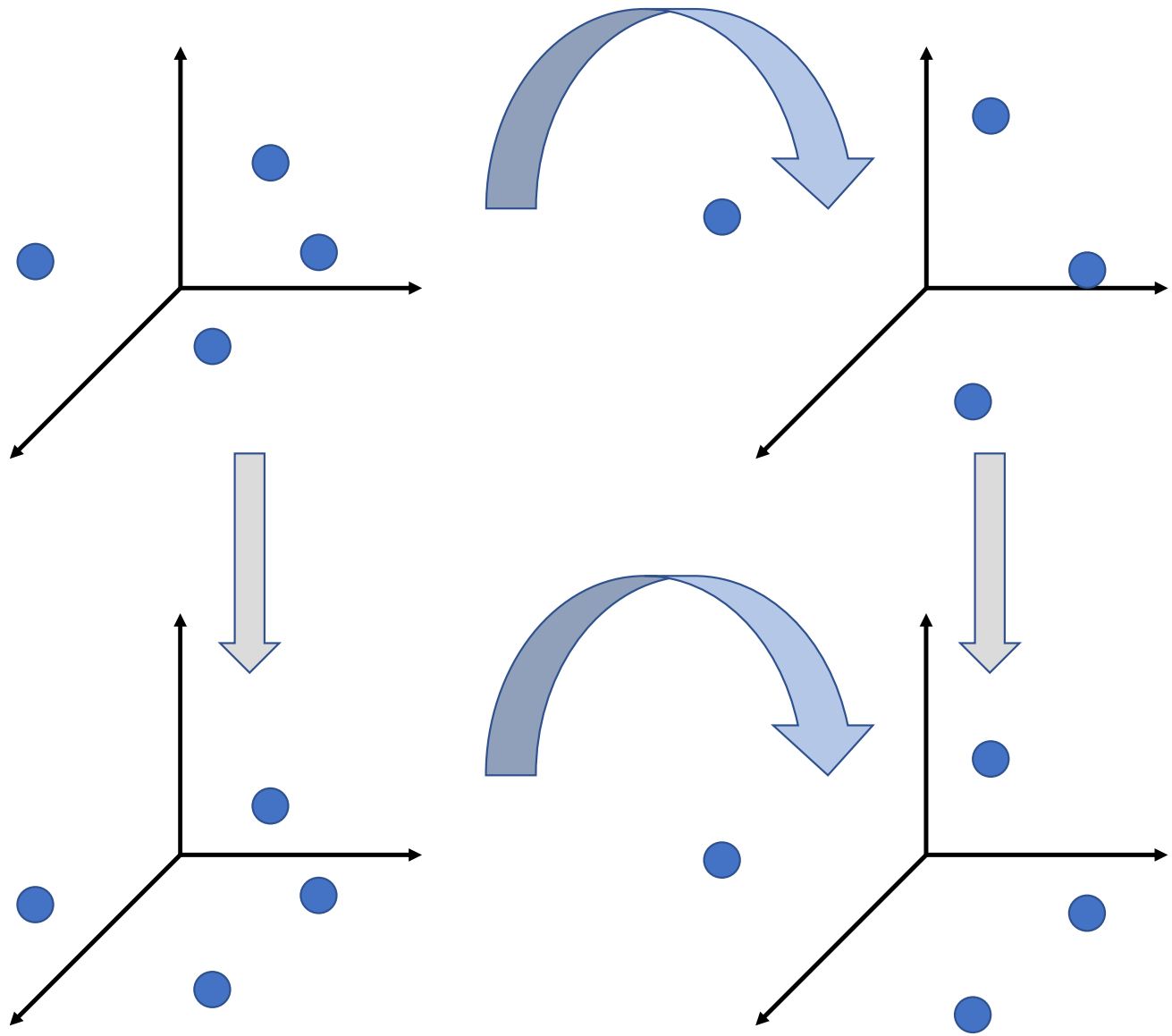
- Equivariance

- Performance

# Invariance vs. equivariance



PointNet

Classification — mug? table? car?

Part Segmentation

Semantic Segmentation

# Equivariance: For Physics simulation

# N-body problem

Equivariant to
- Permutation
- Translation
- Orthogonal
- Lorenz!

$f$

$\tau=0$

$\tau=1$