# On the Universality of Invertible Neural Networks

**Takeshi Teshima**[1][2]**, Isao Ishikawa**[3][2]

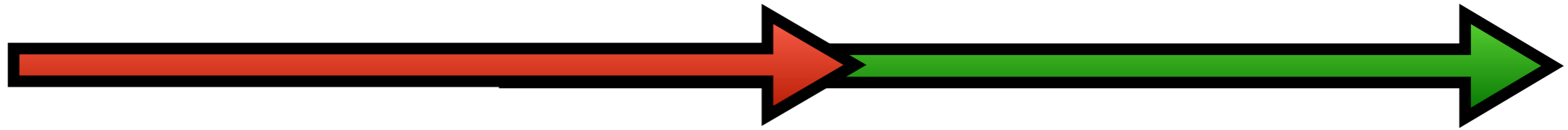[1]**The University of Tokyo, Japan**   [2]**RIKEN, Japan**   [3] **Ehime University, Japan**

**Joint work with  Koichi Tojo, Kenta Oono, Masahiro Ikeda, and Masashi Sugiyama.**

# Today's talk structure

## Part 1

Introduction.

Overview of what we did and why it's important.

## Part 2

Details of the theory.

Theoretical preliminaries and proof machinery.

**Takeshi Teshima** (https://takeshi-teshima.github.io)

Ph.D. candidate @ UTokyo
(advisor: Prof. Sugiyama)

Supported by:
RIKEN JRA Program
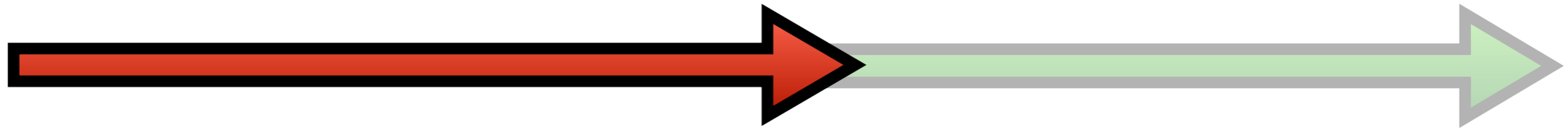and Masason Foundation.

**Recent Research Interests:**
General methodology of machine learning.
In particular: "Causality for machine learning"

- Causal mechanism transfer ← I used INNs in this work

- Causal data augmentation

# Today's talk structure

# Part 1

Introduction.

Overview of what we did and why it's important.

# Part 2

Details of the theory.

Theoretical preliminaries and proof machinery.

# Invertible Neural Networks (INNs)

## Goal

Understand theoretical props of **invertible neural networks (INNs)**.

## Invertible Neural Networks (INNs) generated by $\mathcal{G}$

Compositions of **flow maps/layers** $\mathcal{G}$ and **affine transforms** Aff.

$$f = g_1 \circ W_1 \circ \cdots \circ g_k \circ W_k \quad (g_i \in \mathcal{G}, W_i \in \text{Aff})$$

$\mathcal{G}$ is parametrized ("trainable") but **designed to be invertible**.

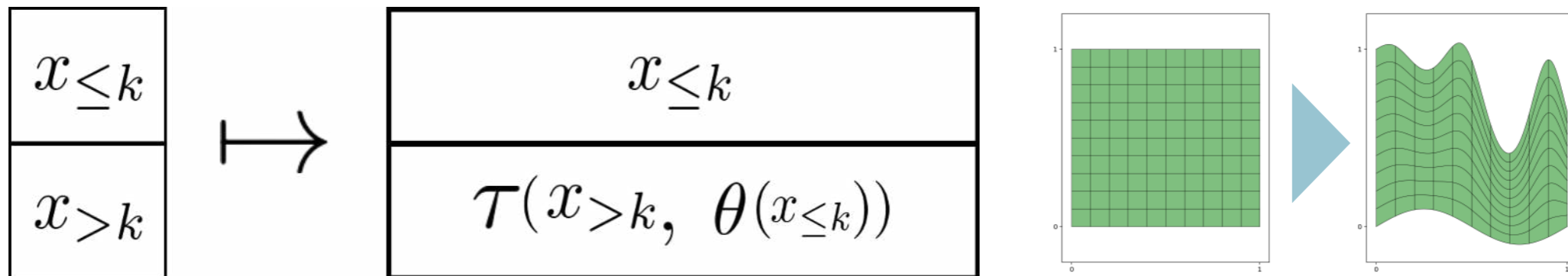($\mathcal{G}$ is often rather simple → Composed to model complex $f$)

## Example (Designs of flow layers $\mathcal{G}$)

- Coupling-based flow layers  [DKB14, PNRML19, KPB19]

- Neural ordinary differential equations [CRBD18]

**Coupling flows (CFs)**   [DKB14, PNRML19, KPB19]

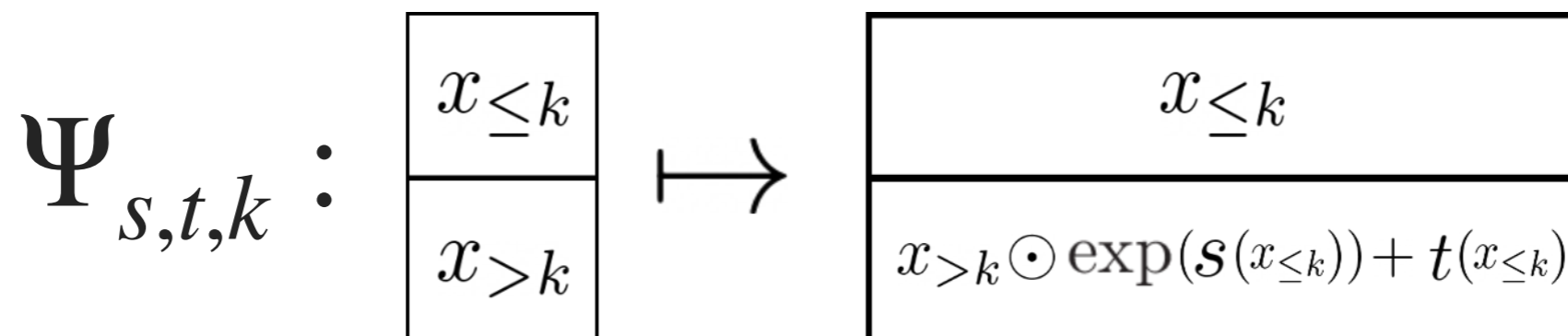$$\frac{x_{\leq k}}{x_{>k}} \longmapsto \frac{x_{\leq k}}{\tau(x_{>k}, \ \theta(x_{\leq k}))}$$

**Idea**: Keep some dimensions unchanged. (Strong constraint!)

**CF-INN** = Coupling-flow based INN.

**Affine-coupling flows (ACFs)**   [DKB14,DSB17,KD18]

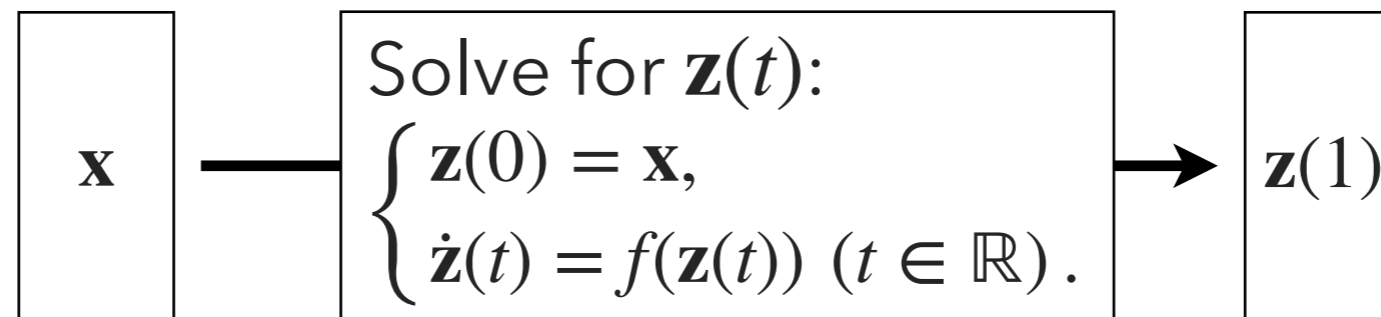One of the simplest CFs using **coordinate-wise affine transformation**:

$$\Psi_{s,t,k} : \frac{x_{\leq k}}{x_{>k}} \longmapsto \frac{x_{\leq k}}{x_{>k} \odot \exp(s(x_{\leq k})) + t(x_{\leq k})}$$

**NODE layer** $\qquad \mathrm{Lip}(\mathbb{R}^d) := \left\{ f \colon \mathbb{R}^d \to \mathbb{R}^d \mid f \text{ \textbf{is Lipschitz}} \right\}$

For each $f \in \mathrm{Lip}(\mathbb{R}^d)$, we define an invertible map $\mathbf{x} \mapsto \mathbf{z}(1)$
via an initial value problem [DJ76]

$$
\mathbf{x} \longrightarrow
\boxed{
\begin{aligned}
&\text{Solve for } \mathbf{z}(t)\text{:} \\
&\begin{cases} \mathbf{z}(0) = \mathbf{x}, \\ \dot{\mathbf{z}}(t) = f(\mathbf{z}(t)) \ (t \in \mathbb{R}). \end{cases}
\end{aligned}
}
\longrightarrow \boxed{\mathbf{z}(1)}
$$

**NODE layers** [CRBD18]

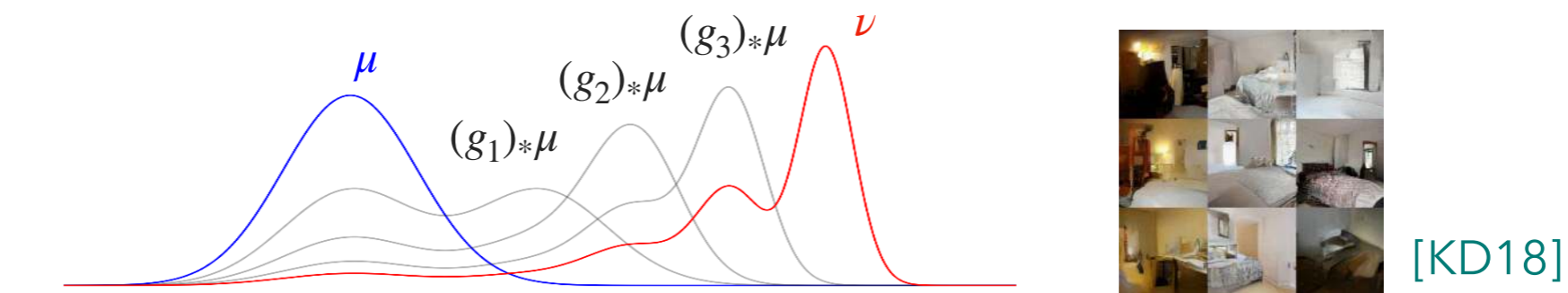Then, for $\mathscr{H} \subset \mathrm{Lip}(\mathbb{R}^d)$, consider the set of NODEs:

$$\mathrm{NODEs}(\mathscr{H}) := \{ \mathbf{x} \mapsto \mathbf{z}(1) \mid f \in \mathscr{H} \}$$

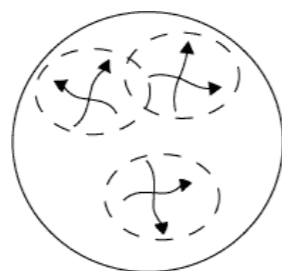**Useful properties of INNs (for nicely designed $\mathcal{G}$)**

✓ **Explicit and efficient invertibility**.

✓ **Tractability** of Jacobian determinant (for nicely designed $\mathcal{G}$).

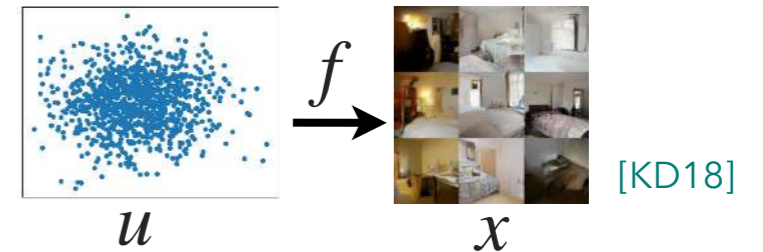**Usages of INNs**

- Approximate distributions (normalizing flows).



[KD18]

- Approximate invertible maps (feature extraction & manipulation).



[DSB17]

$u$    $f$    $x$   [KD18]

## Normalizing Flows

Express $x$ as a transformation $f$ of a real vector $u$ sampled from $p_u$:

$$x = f(u) \text{ where } u \sim p_u$$

## Examples

- Generative modeling [DSB17,KD18,OLB+18,KLSKY19,ZMWN19]
- Probabilistic inference [BM19,WSB19,LW17,AKRK19]
- Semi-supervised learning [IKFW20]

## Training by Maximum Likelihood (Invertibility+Tractable Jacobian!)

By change of variables formula:

$\downarrow$ easily invertible

$$\log p_x(x) = \log p_u(f^{-1}(x)) + \log \left| \det J_{f^{-1}}(x) \right|$$    ($J_{f^{-1}}$: Jacobian of $f^{-1}$)

$\uparrow$ known      $\uparrow$ tractable

## Feature Extraction & Manipulation



$f^{-1}$

$f$

[DSB17]

$u$

$x$

1. Extract latent representation $u$ from $x$ by $f$.

2. Modify $u$ in the latent space (e.g., interpolation).
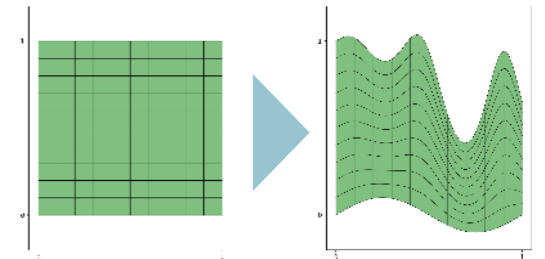
3. Map back to the data space by $f^{-1}$.

## Examples

- Generative modeling [DSB17,KD18,OLB+18,KLSKY19,ZMWN19]
- Semi-supervised learning [IKFW20]
- Transfer learning [TSS20]

INN $f$ is used for **distribution modeling** (application 1) and **invertible function modeling** (application 2).

**BUT...**

$\mathscr{G}$ relies on special designs to maintain good properties. (e.g., CF layers keep some dimensions unchanged)



**Complications**

- The layers have clever specific designs (e.g., ACFs).
- Function composition is the only way to build complex models. (Operations such as addition or multiplications are not allowed.)

**Research question**

**Can these INNs have sufficient representation power?**

(Restricted function form → restricted representation power?)

**Paper 1: Coupling-based invertible neural networks are universal diffeomorphism approximators (NeurIPS 2020)**　[TIT+20]　NEURAL INFORMATION PROCESSING SYSTEMS　**Oral paper!**

- Proposed **a general theoretical framework** to analyze the representation power (universalities) of invertible models.

- Analyzed **CF-INNs** (**ACFs** and more advanced ones).

**Paper 2: Universal Approximation Property of Neural Ordinary Differential Equations (NeurIPS 2020 DiffGeo4DL Workshop)** [TTI+20]　NEURAL INFORMATION PROCESSING SYSTEMS

- Analyzed **NODEs**, building on the general framework.

- (with minor modification to the general framework)

Here,

**"Representation power" = Universal approximation property.**

> **Definition** (informal) [C89,HSW89]

$\sup$- **($L^p$-) universal approximator**:
the model can approximate any target function
w.r.t. $\sup$- ($L^p$-) norm on a compact set.



> **Definition** (informal)

A model is a **distributional universal approximator** if it can transform one distribution arbitrarily close to any distribution.

$$(g_n)_*\mu \xrightarrow[n\to\infty]{} \nu$$

(**weak** convergence).

**Definition (Approximation target $\mathscr{D}^2$)**

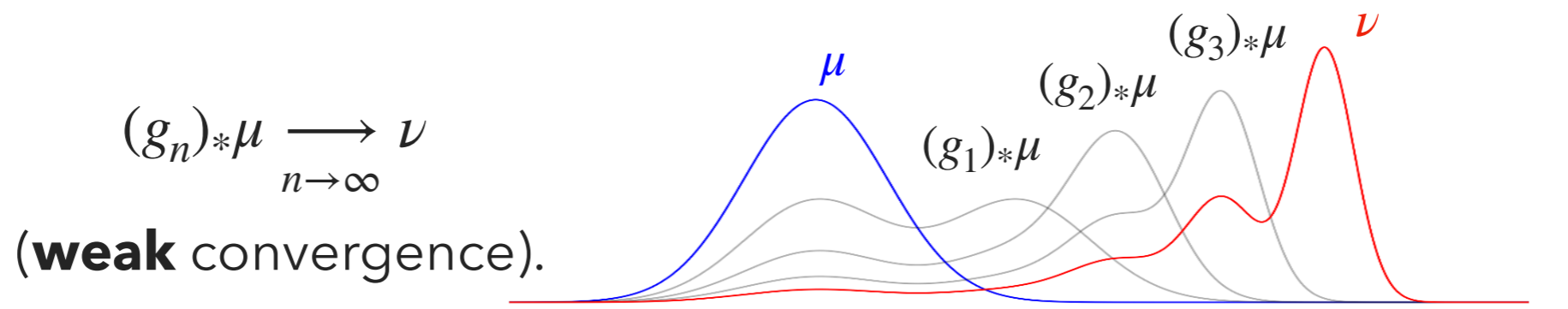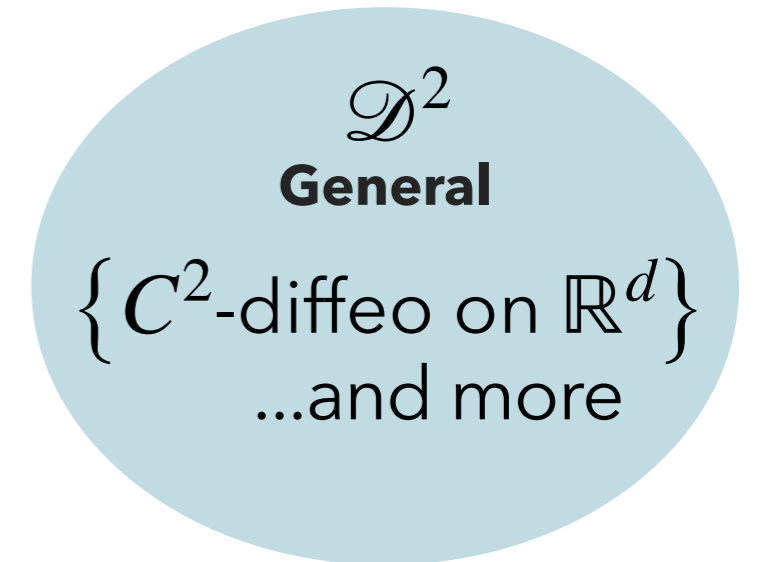Fairly **large set** of smooth invertible maps.

$$\mathscr{D}^2 := \left\{ C^2\text{-diffeo of the form } f : U_f \to f(U_f) \right\}$$
$$(U_f \subset \mathbb{R}^d : \text{open } C^2\text{-diffeo to } \mathbb{R}^d)$$

$$\mathscr{D}^2$$
**General**

$$\left\{ C^2\text{-diffeo on } \mathbb{R}^d \right\}$$
...and more

**Paper 1 Result (Theoretical Framework)** (under mild regularity conditions)

| $L^p$-/sup-univ. for $\mathscr{D}^2$ | $\Leftrightarrow$ | $L^p$-/sup-univ. for $\Xi$ | $\Leftrightarrow$ | $L^p$-/sup-univ. for $\mathcal{S}_c^\infty$ | $\Rightarrow$ | Distributional-univ. |

$\Xi$ : "flow endpoints"

**Application of a structure theorem in differential geometry**

$$\mathscr{D}^2$$
**General**

$$\Xi \qquad \mathcal{S}_c^\infty$$
**Specific**

$$\boxed{\begin{array}{l} L^p\text{-/sup-univ.} \\ \text{for } \mathscr{D}^2 \end{array}} \Leftrightarrow \boxed{\begin{array}{l} L^p\text{-/sup-univ.} \\ \text{for } \Xi \end{array}} \Leftrightarrow \boxed{\begin{array}{l} L^p\text{-/sup-univ.} \\ \text{for } \mathscr{S}_c^\infty \end{array}} \Rightarrow \boxed{\begin{array}{l} \text{Distributional-} \\ \text{univ.} \end{array}}$$

## Paper 1 Result (Examples of Universal Coupling Flows)

- **Sum-of-squares polynomial flow** (SoS-flow)   [JSY19]

- **Deep sigmoidal flow** (DSF; aka. NAF)   [HKLC18]

$\mathscr{D}^2$
**General**

$\mathscr{S}_c^\infty$

**Specific**

yield sup-**univ. INNs for** $\mathscr{S}_c^\infty$ **(and hence for** $\mathscr{D}^2$**, and also Dist-univ.).**

(stronger than in  [JSY19, HKLC18]).

## Paper 1 Result (Affine Coupling Flows yield universal INNs)

Affine Coupling Flows yield $L^p$-univ. INNs for $\mathscr{S}_c^\infty$
(and hence for $\mathscr{D}^2$, and also Dist-univ.).

$L^p$-/sup-univ. for $\mathscr{D}^2$ $\Leftrightarrow$ $L^p$-/sup-univ. for $\Xi$ $\Leftrightarrow$ $L^p$-/sup-univ. for $\mathscr{S}_c^\infty$ $\Rightarrow$ Distributional-univ.

**Paper 2 Result (Analysis of NODEs)**

NODEs yield **sup**-univ. INNs for $\Xi$

(and hence **sup**-univ. for $\mathscr{D}^2$. Also Dist-univ.).

# Overview and Recap

**What did we do?** 🔍 **Theoretically investigated:**
**Are our INNs expressive enough?**

INNs = Invertible neural networks

**Why important?** 🏠 **Models without a representation**
**power guarantee are hard to rely on.**

**What is the result?** 🍎 **"Coupling-based INNs (CF-INNs)" and**
**"NODE-based INNs (NODE-INNs)" are**
**"universal function approximators"**
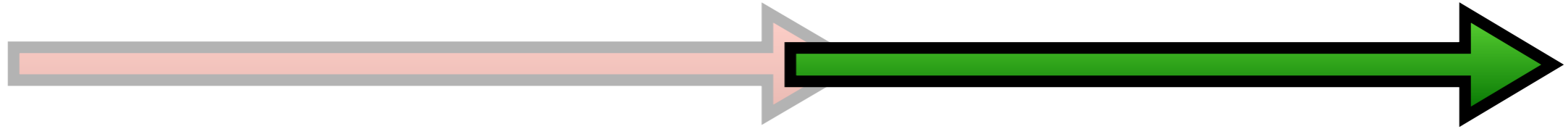**despite their special architectures.**

**Message**

**CF-INNs and NODE-INNs can be relied on in modeling**
**invertible functions and probability distributions.**

# Today's talk structure

## Part 1

Introduction.

Overview of what we did and why it's important.

## Part 2

Details of the theory.

Theoretical preliminaries and proof machinery.

# Self-introduction

## Isao Ishikawa

Assistant professor
@Center for Data Science, Ehime University

Supported by CREST JPMJCR1913

**Recent Research Interests:**

Mathematical analysis of theoretical backgrounds of machine learning and data analysis

- Analysis of representation power of neural networks

- Data analysis via Koopman operator

# Contents of Part 2

1. Idea of proof

2. Notion of universalities

3. Machinery for proof

    i) Compatibility of approximation and composition

    ii) Structure theorem of diffeomorphism group

4. Proof outline of universality of NODE

5. Proof of results in paper 1

## Difficulty

- We cannot use **techniques of functional analysis**!

    - INNs and $\mathscr{D}^2$ are **not** linear spaces

        Recall : $\mathscr{D}^2 := \left\{ C^2\text{-diffeo of the form } f : U_f \to f(U_f) \right\}$

        ($U_f \subset \mathbb{R}^d$ : open $C^2$-diffeo to $\mathbb{R}^d$)

    - Existing methods do not work....(e.g. Hahn-Banach theorem, Fourier transform, Stone-Weirestrass theorem, e.t.c)

## Idea

- Utilize a concrete structure of the **diffeomorphism group** !

$\mathcal{M}$: model, set of measurable bijection from $\mathbb{R}^d$ to $\mathbb{R}^d$ (e.g. INNs)

$\mathcal{F}$: target functions $f : U_f \to f(U_f)$ (e.g. $\mathcal{D}^2$)

$\mathcal{M}$ is an $L^p$-**universal approximator** for $\mathcal{F}$ **if**

$\forall f \in \mathcal{F}$, $\forall \varepsilon > 0$, $\forall K \subset U_f$: compact , $\exists g \in \mathcal{M}$
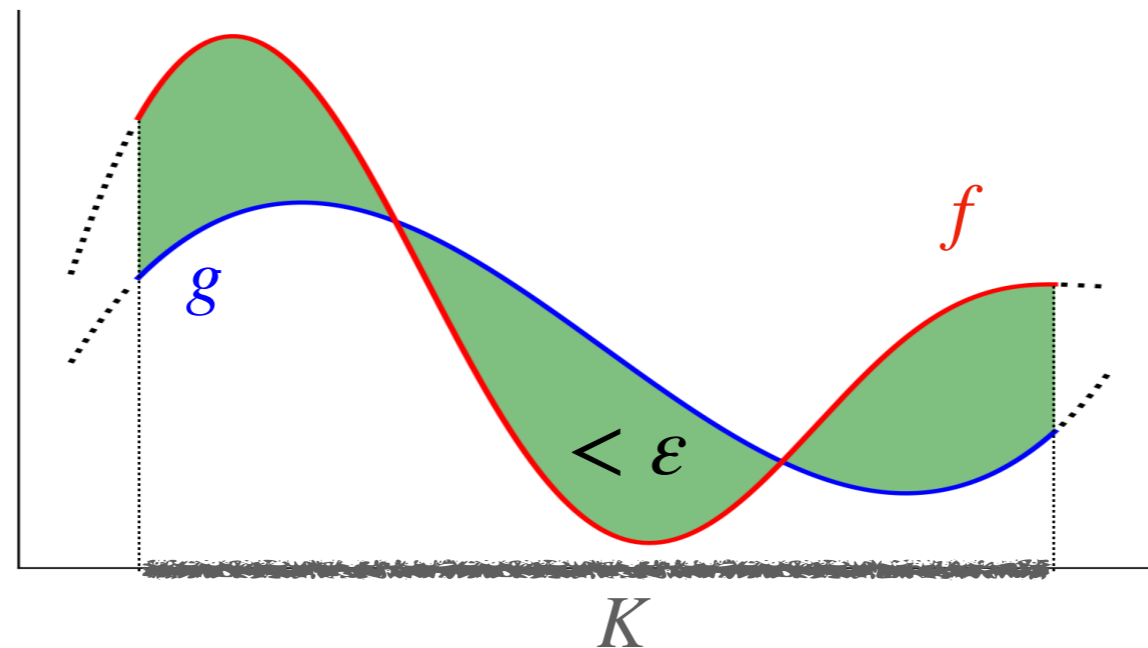
$$\int_K |f(x) - g(x)|^p\, dx < \varepsilon$$

$\mathcal{M}$ : model, set of measurable bijection from $\mathbb{R}^d$ to $\mathbb{R}^d$ (e.g. INNs)

$\mathcal{F}$ : target functions $f : U_f \to f(U_f)$ (e.g. $\mathcal{D}^2$)

$\mathcal{M}$ is an sup-**universal approximator** for $\mathcal{F}$ **if**

$\forall f \in \mathcal{F}, \ \forall \varepsilon > 0, \ \forall K \subset U_f$ : compact , $\exists g \in \mathcal{M}$

$$\sup_{x \in K} |f(x) - g(x)| < \varepsilon$$

**Proposition**

A model $\mathscr{M}$ is a sup-universal approximator for a target $\mathscr{F}$

$$\Downarrow$$

A model $\mathscr{M}$ is an $L^p$-universal approximator a target $\mathscr{F}$

- Is a composition of approximations an approximation of the composition ?

- We may reduce the problem to approximations of small constituents

## Proposition

$\mathcal{M}$ : a set of piecewise $C^1$-diffeomorphisms

$F_1, \ldots, F_r$ : **linearly increasing** piecewise $C^1$-diffeomorphims

Assume $\exists H_i \in \mathcal{M}$ such that

$$H_i \approx F_i \ (L^p\text{-approximation on any compact sets})$$

Then, for compact set $K \subset \mathbb{R}^d$, there exist $G_1, \ldots, G_r \in \mathcal{M}$ such that

$$G_r \circ \cdots \circ G_1 \approx F_r \circ \cdots \circ F_1 \ (L^p\text{-approximation on } K)$$

## Remark

If $\mathcal{M}$ is composed of **locally bounded** maps and $F_i$'s are **continuous**,

we have a similar proposition for sup-universal approximators.

**Definition** (compactly supported diffeomorphisms)

$\mathrm{Diff}_c^2$: the set of $C^2$-diffeomorphisms $f : \mathbb{R}^d \to \mathbb{R}^d$ such that $f(x) = x$ outside a compact subset ($U_f = \mathbb{R}^d$).

$$\mathscr{D}^2$$

$$\mathrm{Diff}_c^2$$

**Theorem** (Herman, Thurston, Epstein, and Mather)

$\mathrm{Diff}_c^2$ is a **simple group** (does not have any proper normal subgroup except $\{\mathrm{Id}\}$)

## Proposition

For $f \in \mathscr{D}^2$ $(f : U_f \to \mathbb{R}^d$ $)$ and compact subset $K \subset U_f$, there
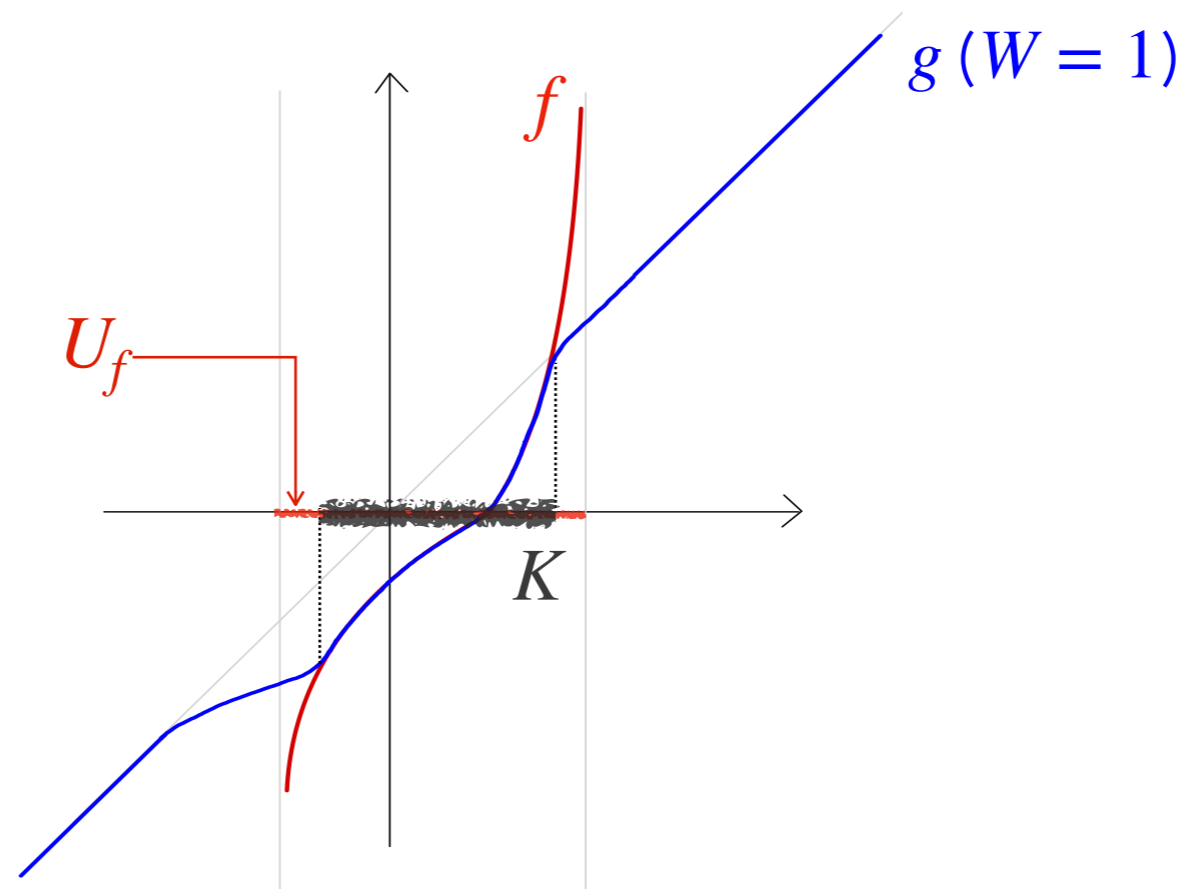
exist an affine transform $W \in \mathrm{Aff}$ and $g \in \mathrm{Diff}^2_c$ such that

$$f|_K = W \circ g|_K.$$

**Definition** (flow endpoints $\Xi$)

$g \in \mathrm{Diff}_c^2$: **flow endpoint** **if** there exists a **continuous** and **"additive"** map $\phi : [0,1] \to \mathrm{Diff}_c^2$ such that $\phi(0) = \mathrm{Id}$ and $\phi(1) = g$

$\Xi := \{\text{flow endpoints}\}$

*w.r.t. Whitney topology.*

$\forall s, t \, . \, s+t \in [0,1] \, , \quad \phi(s+t) = \phi(s) \circ \phi(t)$

**Proposition**

The set of finite compositions of flow endpoints (the group generated by $\Xi$) is a **nontrivial normal subgroup** of $\mathrm{Diff}_c^2$.

**Corollary**

For $g \in \mathrm{Diff}_c^2$, there exist **finite** flow endpoints $g_1, \ldots, g_m \in \Xi$ such that

$$g = g_1 \circ \cdots \circ g_m.$$

In particular,

| $L^p$-/sup-univ. for $\mathscr{D}^2$ | $\iff$ | $L^p$-/sup-univ. for $\Xi$ |

$f \in \mathscr{D}^2$: target, $K \subset U_f$: compact

$$f|_K$$
$$\parallel \quad \ll \text{ Extend } f|_K$$
$$\exists W \circ h \text{ (Aff \& compactly supported } C^2\text{-diffeomorphism)}$$
$$\parallel \quad \ll \textbf{structure theorem of diffeomorphism group}$$
$$\exists h_1 \circ h_2 \circ \cdots \text{ (\textbf{flow endpoints})}$$
$$\wr\wr$$

element of $\text{NODEs}(\mathscr{H})$   $\text{NODEs}(\mathscr{H}) := \{ \mathbf{x} \mapsto \mathbf{z}(1) \mid f \in \mathscr{H} \}$

**Paper 2 Result (Analysis of NODEs)**

NODEs yield $\sup$-univ. INNs for $\Xi$

(and hence $\sup$-univ. for $\mathscr{D}^2$. Also Dist-univ.).

$L^p$-/sup-univ. for $\mathscr{D}^2$

$\Longleftrightarrow$

$L^p$-/sup-univ. for $\mathscr{S}_c^\infty$

$$\mathscr{S}_c^\infty := \left\{ \tau : \text{compactly supported } \tau(\mathbf{x}, y) = (\mathbf{x}, u(\mathbf{x}, y)) \right\} \subset \mathrm{Diff}_c^2$$

$$u : \mathbb{R}^{d-1} \to \mathbb{R}, \ (\mathbf{x}, y) \in \mathbb{R}^{d-1} \times \mathbb{R}$$

$f|_K$

$f \in \mathscr{D}^2$: target, $K \subset U_f$: compact

$\parallel$ ≪ Extend $f|_K$

$\exists W \circ h$ (Aff & compactly supported $C^2$-diffeomorphism)

$\parallel$ ≪ **structure theorem of diffeomorphism group**

$\exists h_1 \circ h_2 \circ \cdots$ (**flow endpoints** $\Xi$)

$\parallel$

$\exists g_1 \circ g_2 \circ \cdots$ (nearly Ids)

$\parallel$

$\sigma_1 \circ \tau_1 \circ \cdots$ (**permutations** & $\mathscr{S}_c^\infty$)

**Decompose** $f|_K$ **into simpler mappings**

**Definition** (nearly-Id elements)

$g \in \mathrm{Diff}_c^2$: **nearly-Id element if** $\|dg(x) - I\| < 1$ for $x \in \mathbb{R}^d$

**Proposition**

For a flow endpoint $g \in \mathrm{Diff}_c^2$, there exist nearly-Id elements $g_1, \ldots, g_m \in \mathrm{Diff}_c^2$ such that

$$g = g_1 \circ \cdots \circ g_m.$$

$\because$ $g = \phi(1)$ ($\phi : [0,1] \to \mathrm{Diff}_c^2$ : "additive" and continuous)

Then, $g = \phi(1/m)^m$ and $\phi(1/m) \to \mathrm{Id}$ as $m \to \infty$

Thus, we define $g_1 = g_2 = \ldots = g_m = \phi(1/m)$ for sufficiently large $m$ ∎

## Proposition

For a nearly-$\mathrm{Id}$ element $g \in \mathrm{Diff}_c^2$, there exist $\tau_1, \ldots, \tau_d \in \mathcal{S}_c^2$ and $\sigma_1, \ldots, \sigma_d \in \mathfrak{S}_d$ such that

$$g = \sigma_1 \circ \tau_1 \circ \cdots \circ \sigma_m \circ \tau_m.$$

## Lemma for this proposition

For $g = (g_i)_{i=1}^d \in \mathrm{Diff}_c^2$, if for any $k = 1, \ldots, d$, the submatrix of its jacobian

$$\left( \frac{\partial g_{i+k-1}}{\partial x_{j+k-1}}(x) \right)_{i,j=1,\ldots,d-k-1} \qquad dg = \begin{pmatrix} \ddots & & \nearrow^{invertible.} \\ & \boxed{\phantom{x}} & \\ & & \end{pmatrix}$$

is invertible for all $x$, then there exit $\tau_1, \ldots, \tau_d \in \mathcal{S}_c^2$ and $\sigma_1, \ldots, \sigma_d \in \mathfrak{S}_d$ such that

$$g = \sigma_1 \circ \tau_1 \circ \cdots \circ \sigma_m \circ \tau_m.$$

$L^p$-/sup-univ. for $\mathscr{D}^2$

$\Longleftrightarrow$

$L^p$-/sup-univ. for $\mathscr{S}_c^\infty$

$$\mathscr{S}_c^\infty := \left\{ \tau : \text{compactly supported } \tau(\mathbf{x}, y) = (\mathbf{x}, u(\mathbf{x}, y)) \right\} \subset \mathrm{Diff}_c^2$$

$$u : \mathbb{R}^{d-1} \to \mathbb{R}, \ (\mathbf{x}, y) \in \mathbb{R}^{d-1} \times \mathbb{R}$$

$f|_K \qquad\qquad f \in \mathscr{D}^2 : \text{ target, } K \subset U_f : \text{ compact}$

$\parallel \ll$ Extend $f|_K$

$\exists W \circ h$ (Aff & compactly supported $C^2$-diffeomorphism)

$\parallel$

$\ll$ **structure theorem of diffeomorphism group**

$\exists h_1 \circ h_2 \circ \cdots$ (**flow endpoints** $\Xi$)

$\parallel$

$\exists g_1 \circ g_2 \circ \cdots$ (nearly Ids)

$\parallel$

$\sigma_1 \circ \tau_1 \circ \cdots$ (**permutations &** $\mathscr{S}_c^\infty$)

**Decompose** $f|_K$ **into simpler mappings**

**You show**          **You get**

$$\text{sup-univ. for } \mathscr{S}_c^\infty \quad \Longrightarrow \quad \text{sup-univ. for } \mathscr{D}^2$$

$$\Downarrow \qquad\qquad\qquad \Downarrow$$

$$L^p\text{-univ. for } \mathscr{S}_c^\infty \quad \Longrightarrow \quad L^p\text{-univ. for } \mathscr{D}^2$$

$$\mathscr{D}^2$$
**General**

$$\mathscr{S}_c^\infty$$
**Specific**

Regrading guarantees for existing INN architectures:

- **Sum-of-squares polynomial flow** (SoS-flow)

- **Deep sigmoidal flow** (DSF; aka. NAF)

Previously known/claimed [JSY19, HKLC18]:

$$\text{sup-universality for } \mathscr{S}_c^\infty$$

$$\Downarrow$$

$$\text{sup-universality for } \mathscr{D}^2$$
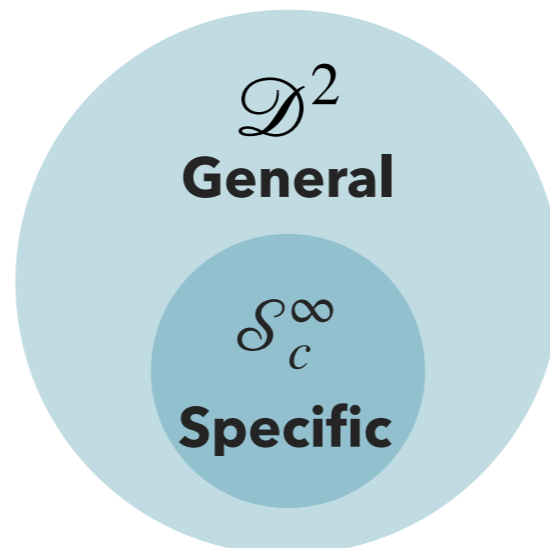
$$\mathscr{D}^2$$
**General**

$$\mathscr{S}_c^\infty$$
**Specific**

**Definition** (distributional universal approximator)

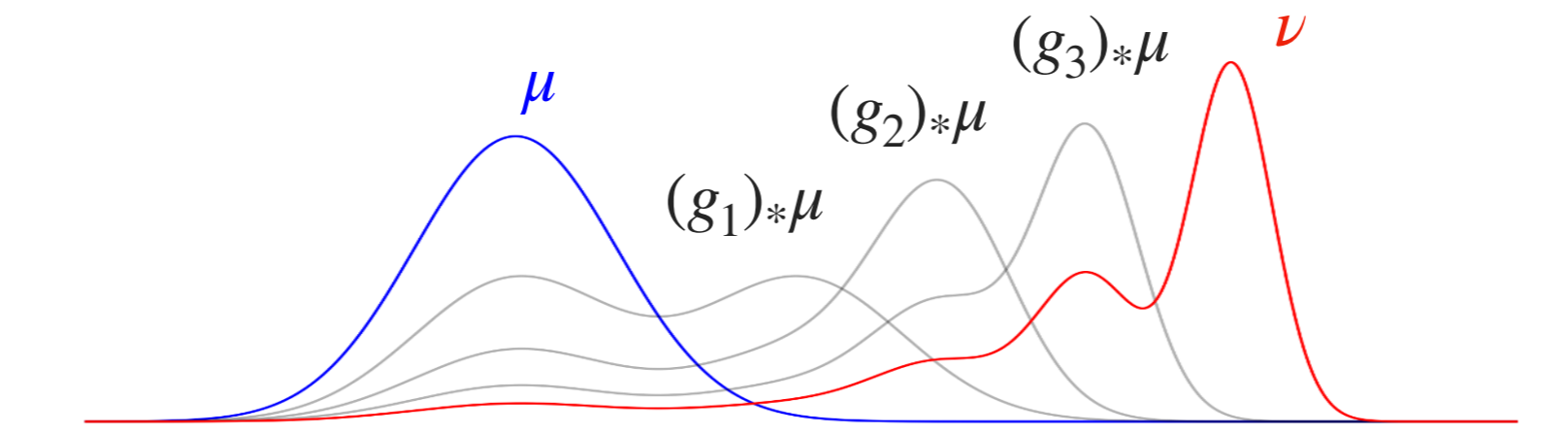$\mathcal{M}$: model, set of measurable bijection from $\mathbb{R}^d$ to $\mathbb{R}^d$ (e.g. INNs)

$\mathcal{P}$:  absolutely continuous probability measures

$\mathcal{M}$ is a **distributional universal approximator** **if**

$\forall \mu, \nu \in \mathcal{P}, \ \exists \{g_n\}_{n=1}^{\infty} \subset \mathcal{M}$

$$(g_n)_* \mu \xrightarrow[n \to \infty]{} \nu \ \ (\textbf{weak} \text{ convergence}).$$

## **Proposition**

A model $\mathscr{M}$ is a $L^p$-universal approximator for a target $\mathscr{D}^2$

$$\Downarrow$$

A model $\mathscr{M}$ is a distributional universal approximator

In summary, we obtain

| $L^p$-/sup-univ. for $\mathscr{D}^2$ | $\Leftrightarrow$ | $L^p$-/sup-univ. for $\Xi$ | $\Leftrightarrow$ | $L^p$-/sup-univ. for $\mathscr{S}_c^\infty$ | $\Rightarrow$ | Distributional-univ. |

$\mathcal{H}$ : functions on $\mathbb{R}^{d-1}$ (e,g, MLPs)

$\mathrm{INN}_{\mathcal{H}\text{-ACF}}$ is an INN with the flow layers composed of

$$\Psi_{d-1,s,t}(\mathbf{x}, y) := \left(\mathbf{x}, e^{s(\mathbf{x})}y + t(\mathbf{x})\right)$$

$$(\mathbf{x}, y) \in \mathbb{R}^{d-1} \times \mathbb{R}, s, t \in \mathcal{H}$$
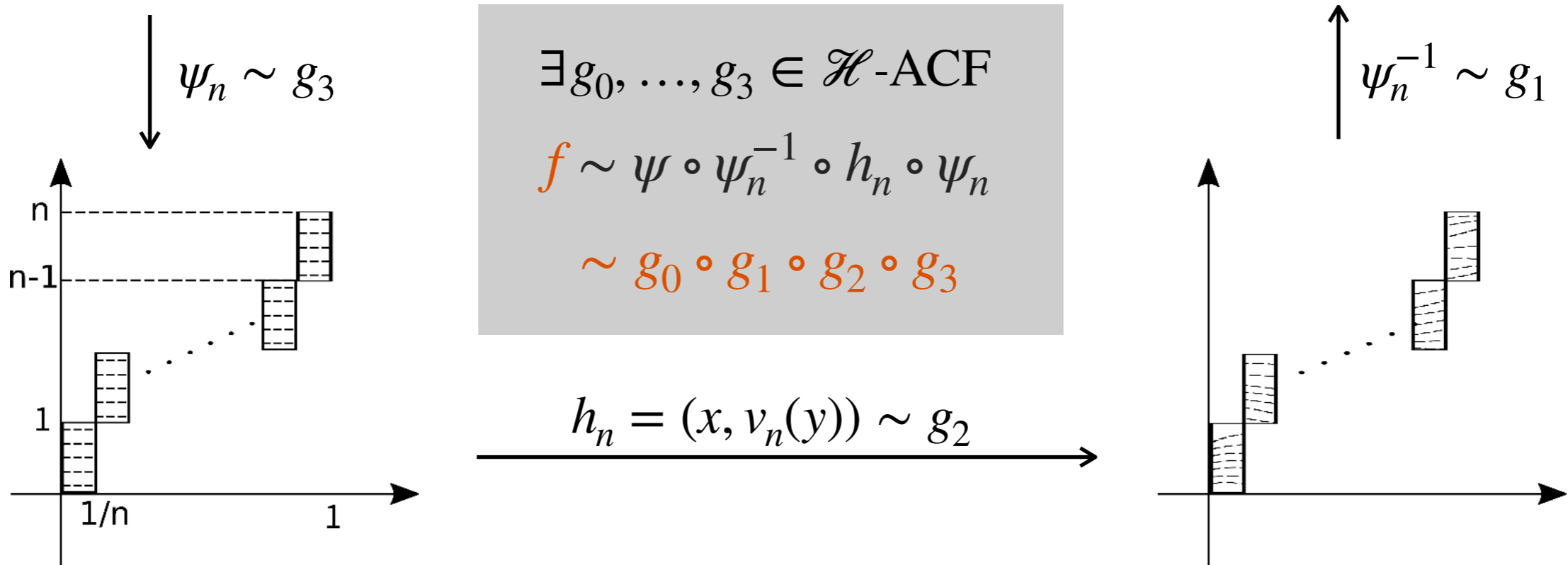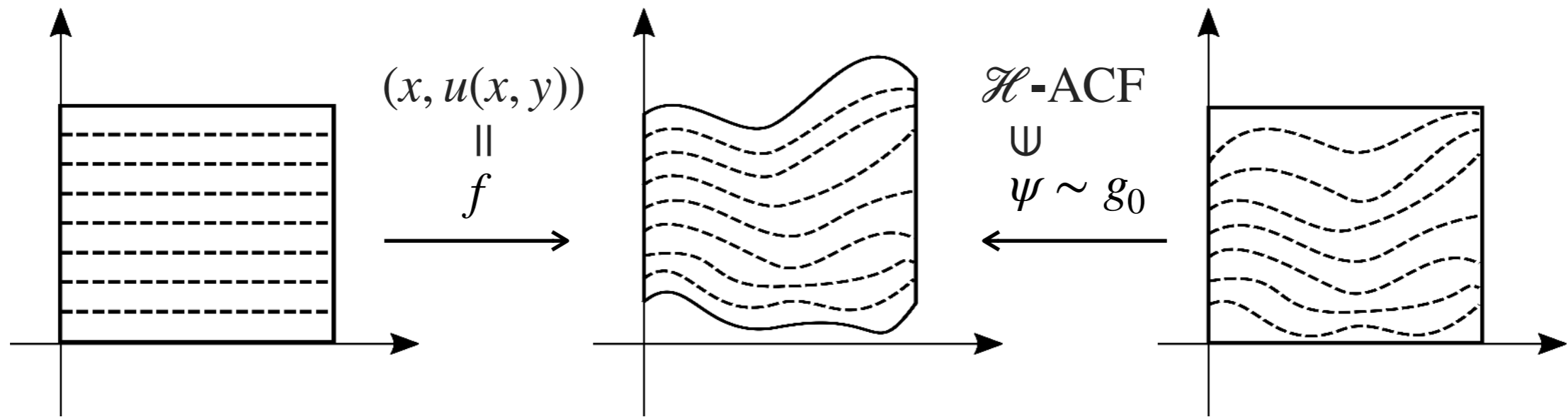
One of the simplest CF-INN

**Theorem**

Assume $\mathcal{H}$ arbitrarily approximates any element in $C_c^\infty(\mathbb{R}^{d-1})$,

and is composed of piecewise $C^1$-functions (e.g. MLPs with

ReLU activation, RKHS with Gaussian kernel, e.t.c).

Then, $\mathrm{INN}_{\mathcal{H}\text{-ACF}}$ is an $L^p$-universal approximator for $\mathcal{S}_c^\infty$

- We may assume $K = [0,1]^2$

$(x, u(x, y))$
$\|$
$f$

$\mathscr{H}$-ACF
$\cup$
$\psi \sim g_0$

$\psi_n \sim g_3$

$\psi_n^{-1} \sim g_1$

$\exists g_0, \ldots, g_3 \in \mathscr{H}\text{-ACF}$

$f \sim \psi \circ \psi_n^{-1} \circ h_n \circ \psi_n$

$\sim g_0 \circ g_1 \circ g_2 \circ g_3$

$h_n = (x, v_n(y)) \sim g_2$

$$\psi_n := \Psi_{d-1,1,t_n} \quad t_n := \sum_{k=0}^{n-1} k \mathbf{1}_{[k/n,(k+1)/n)} \quad v_n(y) = \begin{cases} u(k/n, y) + k & y \in [k, k+1), \\ y & \text{otherwise} . \end{cases}$$

| $L^p$-/sup-univ. for $\mathscr{D}^2$ | $\Leftrightarrow$ | $L^p$-/sup-univ. for $\Xi$ | $\Leftrightarrow$ | $L^p$-/sup-univ. for $\mathscr{S}_c^\infty$ | $\Rightarrow$ | Distributional-univ. |
|---|---|---|---|---|---|---|

## Paper 1 Result (Affine Coupling Flows yield universal INNs)

Affine Coupling Flows yield $L^p$-univ. INNs for $\mathscr{S}_c^\infty$
(and hence for $\mathscr{D}^2$, and also Dist-univ.).

## Remark

The representation power of invertible neural networks based on affine coupling flow is empirically known, and they were **conjectured** distributional universal approximator. We **affirmatively** answer this question.

## Conclusion

- Proposed a general theoretical framework to analyze the representation power (universalities) of invertible models.

- Guarantee the representation power of CF-INNs as an $L^p$ -universal approximator.

- Guarantee the representation power of NODE-INNs as a $\sup$ -universal approximator.

## Future work

- Quantitative analysis: Estimate the number of layers required for the approximation given the smoothness of the target.

Our papers are available at

[1] https://papers.nips.cc/paper/2020/hash/2290a7385ed77cc5592dc2153229f082-Abstract.html

[2] http://arxiv.org/abs/2012.02414

## Message

**CF-INNs and NODE-INNs can be relied on in modeling invertible functions and probability distributions.**

# Appendix

# References

[C89] Cybenko, G. (1989).
Approximation by superpositions of a sigmoidal function.
Mathematics of Control, Signals, and Systems, 2, 303–314.

[HSW89] Hornik, K., Stinchcombe, M., & White, H. (1989).
Multilayer feedforward networks are universal approximators.
Neural Networks, 2(5), 359–366.

[JSY19] Jaini, P., Selby, K. A., & Yu, Y. (2019).
Sum-of-squares polynomial flow.
Proceedings of the 36th International Conference on Machine Learning, 97, 3009–3018.

[HKLC18] Huang, C.-W., Krueger, D., Lacoste, A., & Courville, A. (2018).
Neural autoregressive flows.
Proceedings of the 35th International Conference on Machine Learning, 80, 2078–2087.

[KD18] Kingma, D. P., & Dhariwal, P. (2018).
Glow: Generative flow with invertible 1x1 convolutions.
In Advances in Neural Information Processing Systems 31 (pp. 10215–10224).

[PNRML19] Papamakarios, G., Nalisnick, E., Rezende, D. J., Mohamed, S., & Lakshminarayanan, B. (2019).
Normalizing flows for probabilistic modeling and inference.
ArXiv:1912.02762 [Cs, Stat].

[KPB19] Kobyzev, I., Prince, S., & Brubaker, M. A. (2019).
Normalizing flows: An introduction and review of current methods.
ArXiv:1908.09257 [Cs, Stat].

[DKB14]    Dinh, L., Krueger, D., & Bengio, Y. (2014).
           NICE: Non-linear independent components estimation.
           ArXiv:1410.8516 [Cs.LG].

[DSB17]    Dinh, L., Sohl-Dickstein, J., & Bengio, S. (2017).
           Density estimation using Real NVP.
           Fifth International Conference on Learning Representations (ICLR)

[AKRK19]   Ardizzone, L., Kruse, J., Rother, C., & Köthe, U. (2019).
           Analyzing inverse problems with invertible neural networks.
           7th International Conference on Learning Representations.

[BM19]     Bauer, M., & Mnih, A. (2019).
           Resampled priors for variational autoencoders.
           In Proceedings of machine learning research, 89, 66–75.

[LW17]     Louizos, C., & Welling, M. (2017).
           Multiplicative normalizing flows for variational Bayesian neural networks.
           In Proceedings of the 34th International Conference on Machine Learning,70, 2218-2227.

[NMT+19]   Nalisnick, E. T., Matsukawa, A., Teh, Y. W., Görür, D., & Lakshminarayanan, B. (2019).
           Hybrid models with deep and invertible features.
           In Proceedings of the 36th International Conference on Machine Learning, 97, 4723–4732.

[IKFW20]   Izmailov, P., Kirichenko, P., Finzi, M., & Wilson, A. G. (2020).
           Semi-supervised learning with normalizing flows.
           Proceedings of the 37th International Conference on Machine Learning.

[OLB+18] Oord, A., Li, Y., Babuschkin, I., Simonyan, K., Vinyals, O., Kavukcuoglu, K., Driessche, G., Lockhart, E., Cobo, L., Stimberg, F., Casagrande, N., Grewe, D., Noury, S., Dieleman, S., Elsen, E., Kalchbrenner, N., Zen, H., Graves, A., King, H., … Hassabis, D. (2018). Parallel WaveNet: Fast high-fidelity speech synthesis. Proceedings of the 35th International Conference on Machine Learning, 80, 3918–3926.

[TSS20] Teshima, T., Sato, I., & Sugiyama, M. (2020). Few-shot domain adaptation by causal mechanism transfer. Proceedings of the 37th International Conference on Machine Learning.

[KLSKY19] Kim, S., Lee, S.-G., Song, J., Kim, J., & Yoon, S. (2019). FloWaveNet: A generative flow for raw audio. In Proceedings of the 36th International Conference on Machine Learning, 97, 3370–3378.

[ZMWN19] Zhou, C., Ma, X., Wang, D., & Neubig, G. (2019). Density matching for bilingual word embedding. Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), 1588–1598.

[WSB19] Ward, P. N., Smofsky, A., & Bose, A. J. (2019). Improving exploration in soft-actor-critic with normalizing flows policies. ArXiv:1906.02771 [Cs, Stat].

[CRBD18]   R. T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. K. Duvenaud. (2018).
            Neural ordinary differential equations.
            Advances in Neural Information Processing Systems 31, 6571–6583.

[LLS20]    Q. Li, T. Lin, and Z. Shen. (2020).
            Deep learning via dynamical systems: an approximation perspective.
            arXiv:1912.10382 [cs, math, stat].

[DJ76]     W. Derrick and L. Janos. (1976).
            A global existence and uniqueness theorem for ordinary differentialequations.
            Canadian Mathematical Bulletin, 19(1), 105–107.

[LBH15]    Y. LeCun, Y. Bengio, and G. Hinton. (2015).
            Deep learning.
            Nature, 521(7553), 436–444.

[ALG19]    C. Anil, J. Lucas, and R. Grosse. (2019).
            Sorting out Lipschitz function approximation.
            Proceedings of the 36th International Conference on Machine Learning, PMLR 97, 291–301.

[PPMF20]   Pumarola, A., Popov, S., Moreno-Noguer, F., & Ferrari, V. (2020).
            C-Flow: Conditional Generative Flow Models for Images and 3D Point Clouds.
            2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 7946–7955.