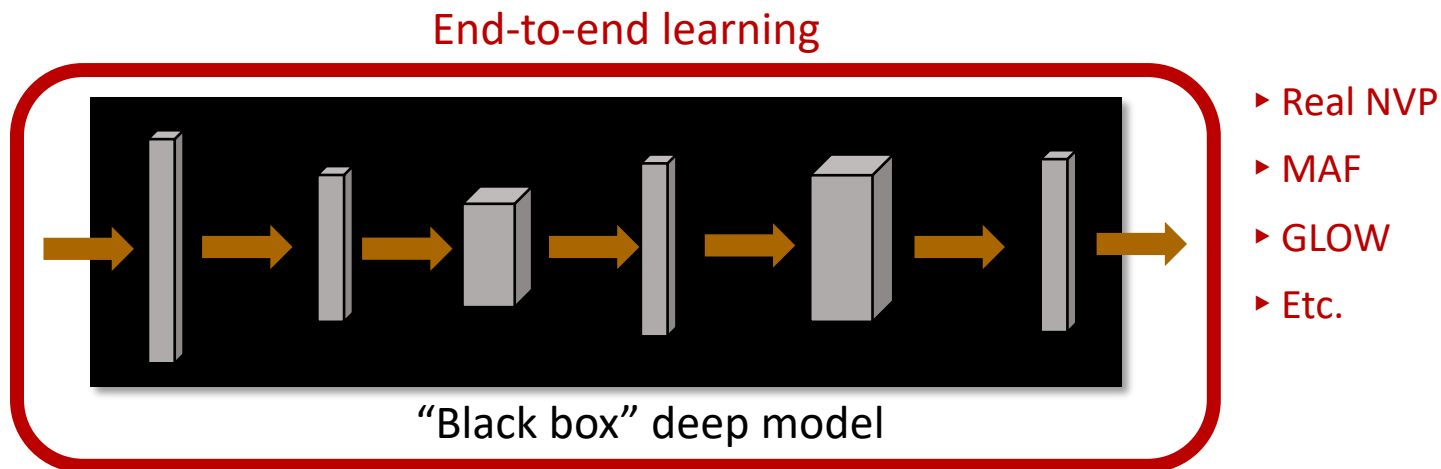


# Deep Density Destructors (from a biased viewpoint)

**David I. Inouye**

Electrical and Computer Engineering  
Purdue University

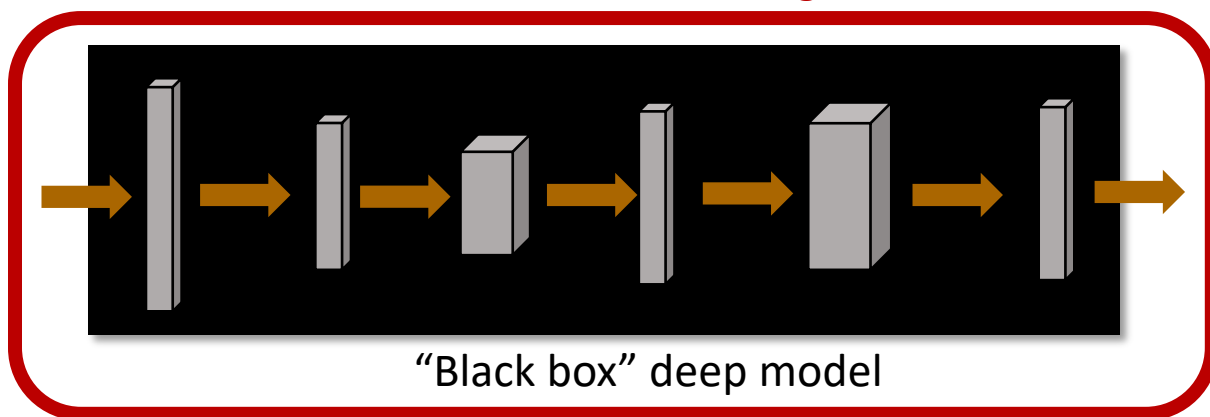
Previous deep normalizing flows are trained end-to-end where all components are optimized simultaneously



"Gray box" deep model

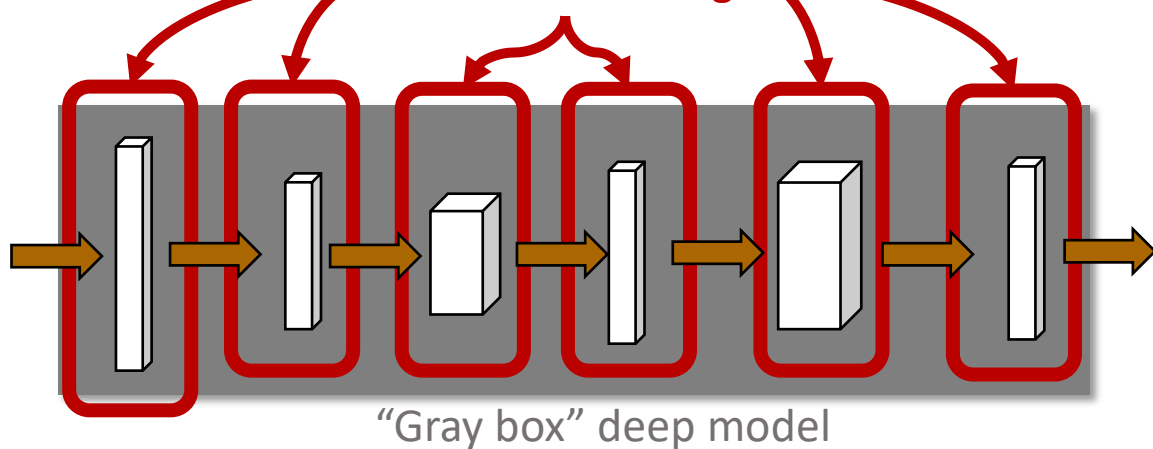
# Modular deep learning would allow *local* learning within each component

## End-to-end learning



- Real NVP
- MAF
- GLOW
- Etc.

## Modular learning



- Density destructors
- Each weak/shallow learning algorithm is independent
- Learning algorithms could be heterogeneous (e.g., SGD and decision trees)

# *Destructive learning* enables modular deep learning via “reverse engineering” data

## Reverse engineering phone

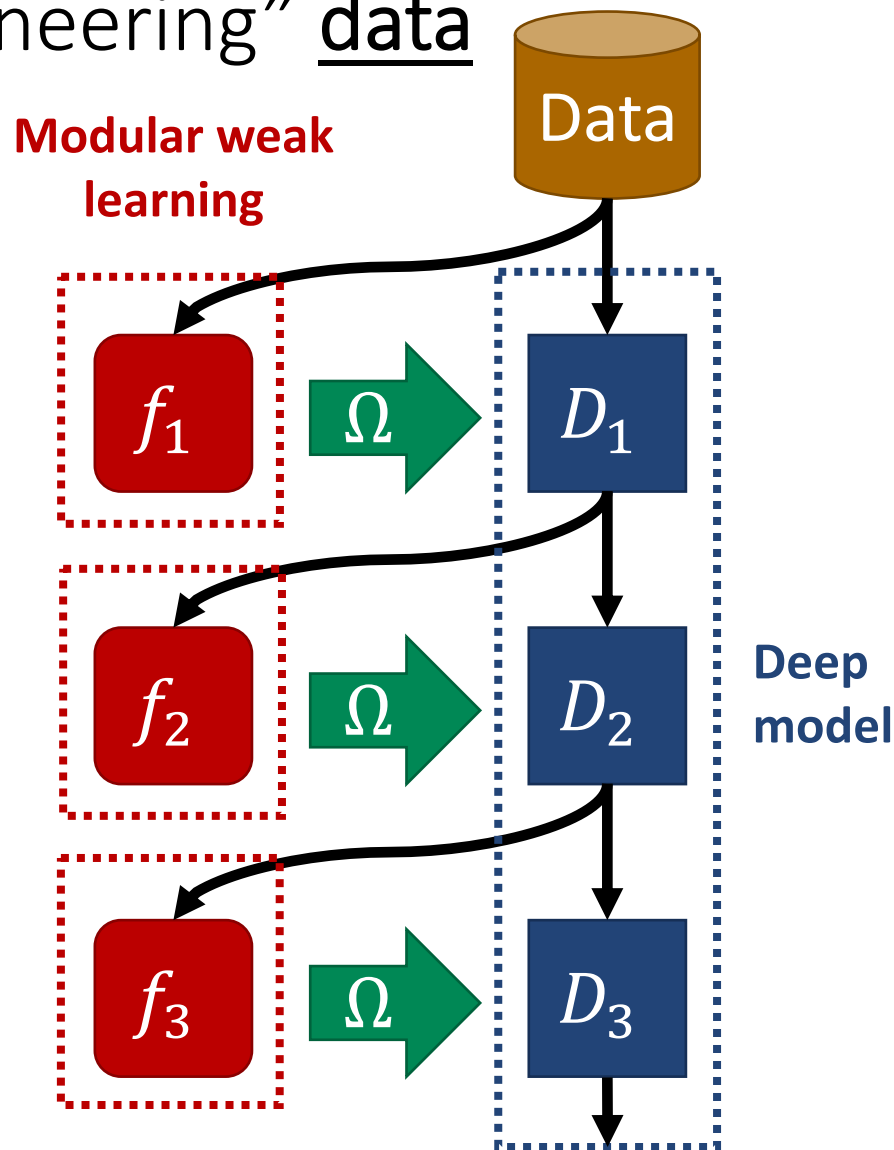
1. **Find** *part* to take off using **understanding and expertise**
2. **Determine how** to take off part in a reversible way (e.g., unscrewing bolts)
3. **Remove part**
4. **Repeat**

## Reverse engineering data

1. **Find** *patterns* in data via **shallow/weak learning**
2. **Map model** to destructive but invertible transformation
3. **Destroy the patterns** via transformation
4. **Repeat**

*Destructive learning* enables modular deep learning via “reverse engineering” data

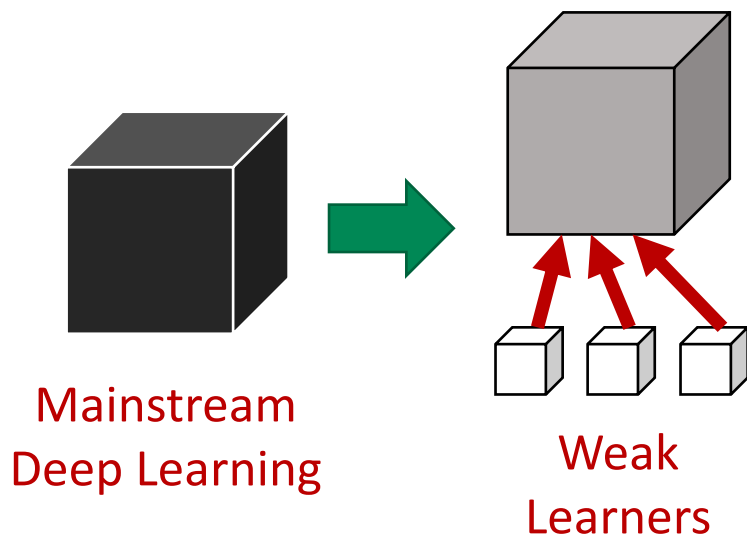
1. **Find** patterns in data via **shallow/weak learning**
2. **Map model** to destructive transformation
3. **Destroy the patterns** via transformation
4. **Repeat**



# Why use modular weak learning for deep models?

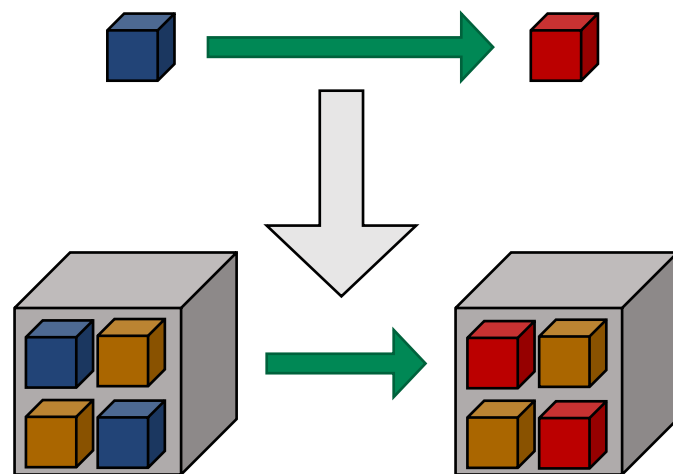
## Reuse

The **algorithms, insights and intuitions** of shallow learning can be lifted into the deep context



## Decoupling

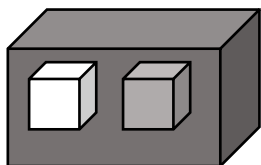
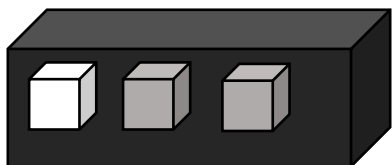
Components can be **debugged, tested and improved** separate from the system



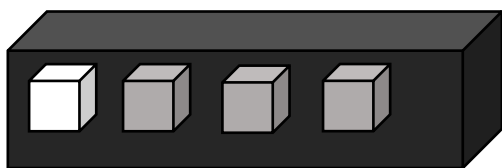
# Why use modular weak learning for deep models?

## Algorithmic Interpretability

Increasing or decreasing model complexity is straightforward



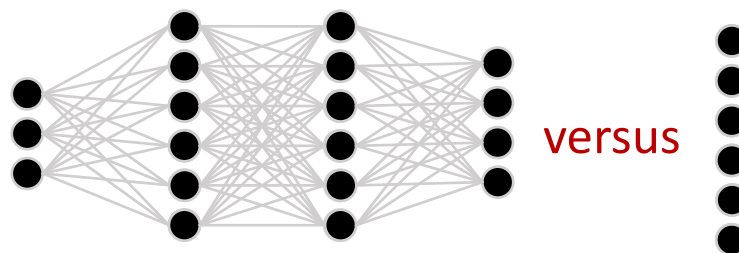
Shrink model  
if problem



Grow if  
more data

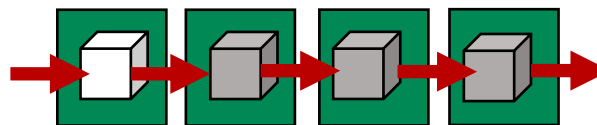
## Resource Constraints

Layer-wise training  
(memory bottleneck)



Pipelined training  
(computation bottleneck)

Shallow/weak online learners



Distributed on different  
processors or devices

# Overview of iterative destructive learning

Motivation and intuition for modular destructive learning

Density destructors objective function

Modular and greedy deep destructive algorithm

- Simple density destructors
- Deep density destructors

Theory about algorithm: Monotonic decrease of objective

Density destructor results

Limitations and open problems

Iterative alignment and translation



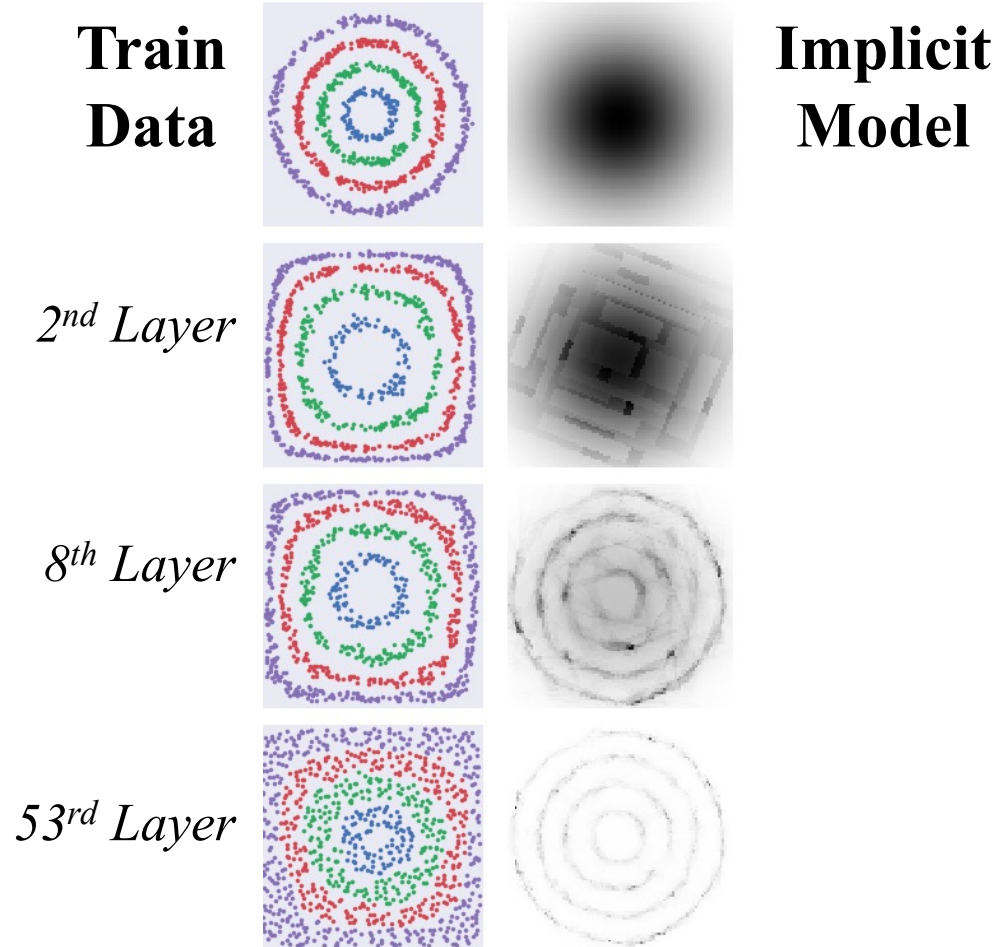
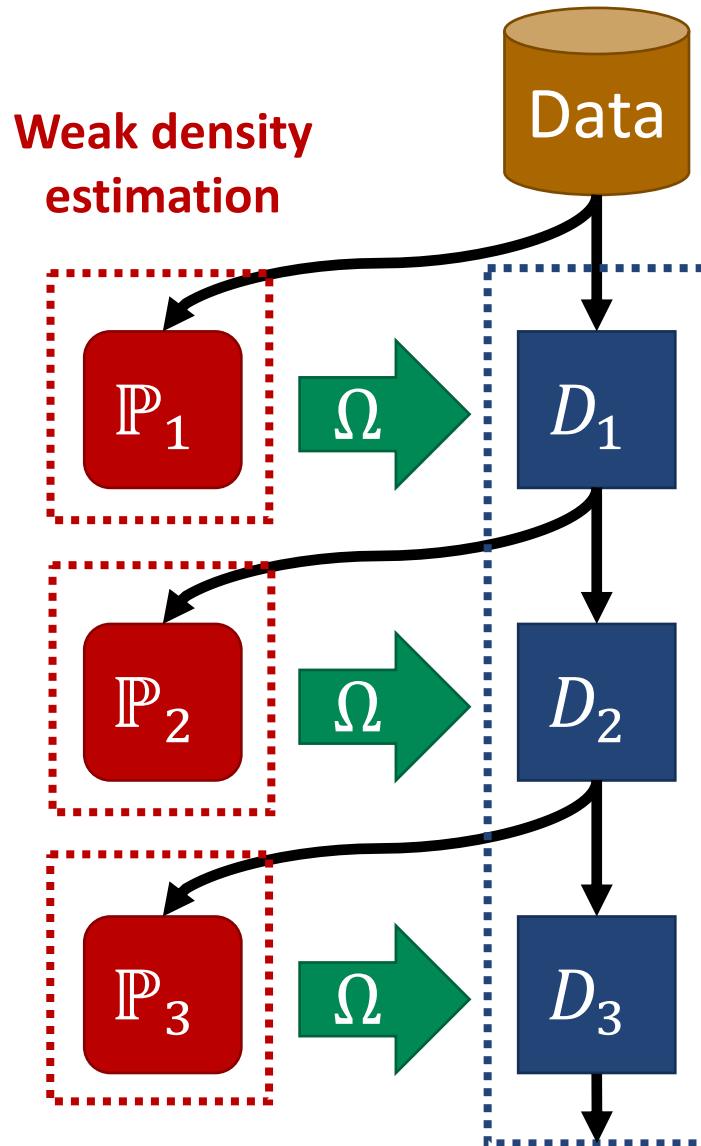
# Background for objective: KL equivalence lemma

- ▶ KL equivalence lemma: If  $z = D(x)$  for invertible  $D$ , then
$$KL(P_x(x), Q_x(x)) = KL(P_z(z), Q_z(z))$$
- ▶  $KL(P_x(x), Q_x(x))$
- ▶  $= E_{P_x} \left[ \log \frac{P_x(x)}{Q_x(x)} \right]$
- ▶  $= E_{P_x} \left[ \log \frac{P_z(D(x)) |J_D(x)|}{Q_z(D(x)) |J_D(x)|} \right]$  (Change of variables formula)
- ▶  $= E_{P_x} \left[ \log \frac{P_z(D(x))}{Q_z(D(x))} \right]$
- ▶  $= E_{P_z} \left[ \log \frac{P_z(D(D^{-1}(z)))}{Q_z(D(D^{-1}(z)))} \right]$  (Expectation change of variables LOTUS)
- ▶  $= E_{P_z} \left[ \log \frac{P_z(z)}{Q_z(z)} \right]$
- ▶  $= KL(P_z(z), Q_z(z))$

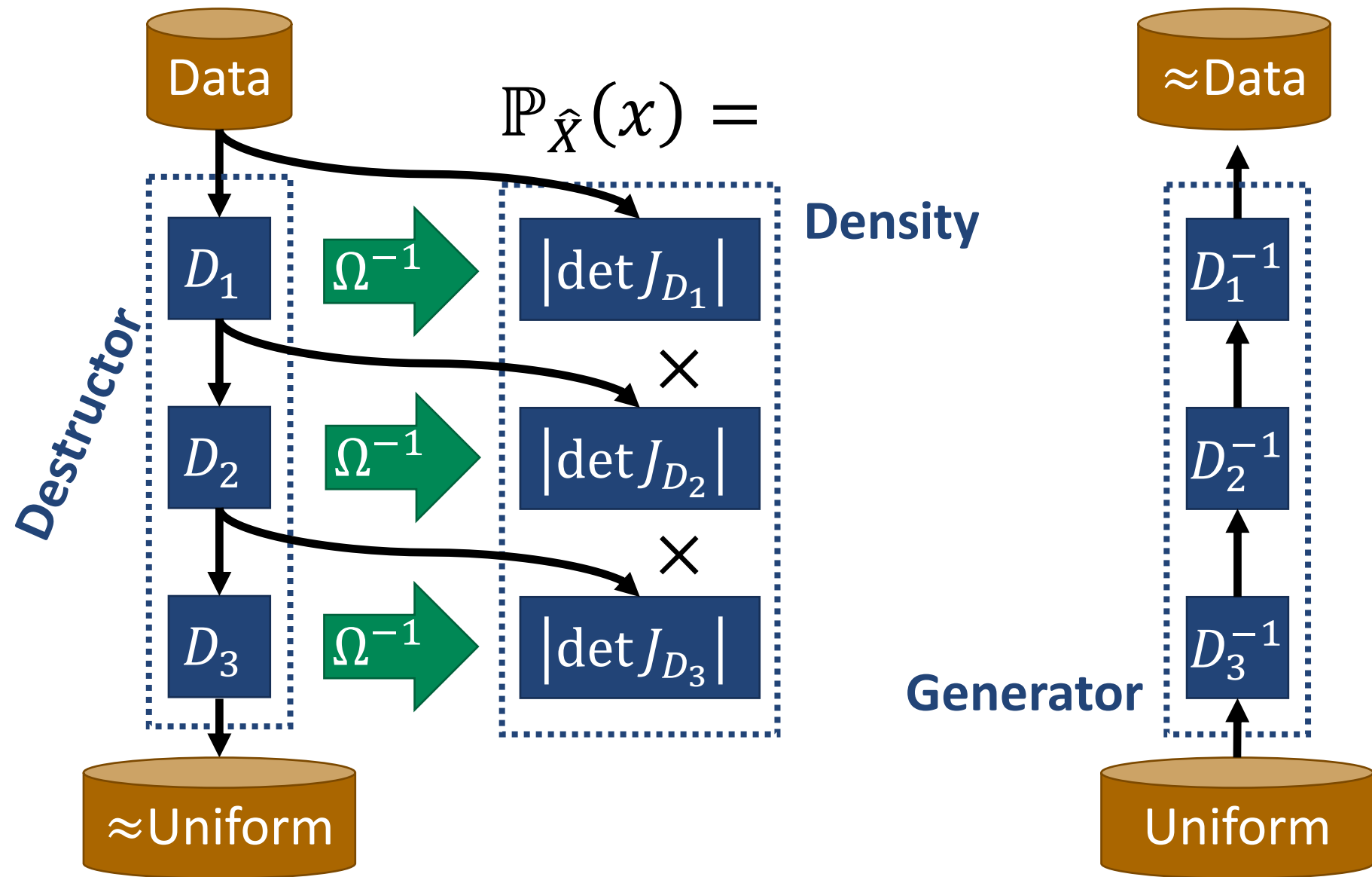
# Destructive learning objective is equivalent to MLE

- ▶ The destructive learning objective, where  $z = D(x)$ , and  $U_z(z)$  is the uniform density function  
$$\arg \min_D KL(P_z(z; D), U(z))$$
- ▶ Simple corollary is that objective above is MLE:
- ▶  $KL(P_z(z; D), U(z))$
- ▶  $= KL(P_x(x), Q_x(x; D))$  (KL equivalence, MLE objective)
- ▶  $= KL(P_x(x), |J_D(x)|U(D(x)))$  (In terms of  $D$ )
- ▶  $= KL(P_x(x), |J_D(x)|)$  ( $U(z) = 1$ )

# Algorithm: Deep density destructors via sequence of *weak* destructors

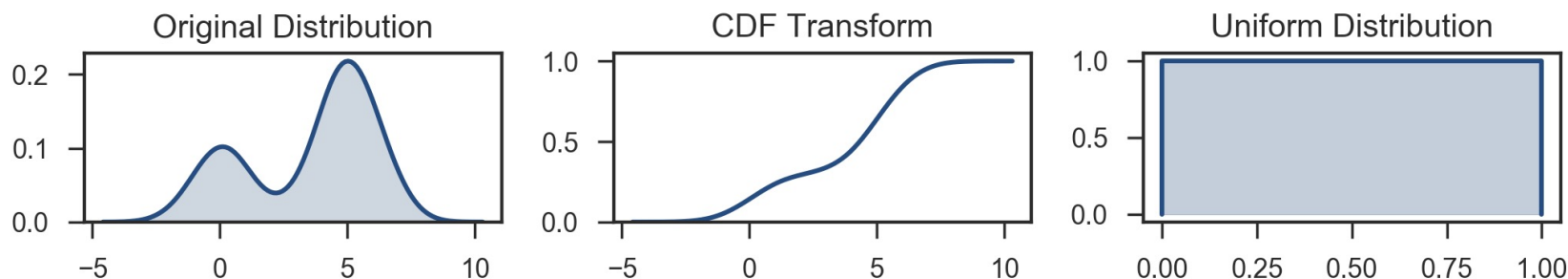


# Density computation and sample generation



# Definition: *Density destructors* generalize the univariate CDF transformation

## ► Univariate: CDF transformation



## ► The map $\Omega(\mathbb{P}) = D$ should:

1. Encode the density  $\mathbb{P}$  into  $D$ , i.e.  $\exists \Omega^{-1}$ .
2. Ensure  $D$  destroys all patterns in  $\mathbb{P}$  when applied to the random variable, i.e. the distribution of  $D_X(X)$  is **maximum entropy**.

## ► A *density destructor* is an **invertible** transformation such that

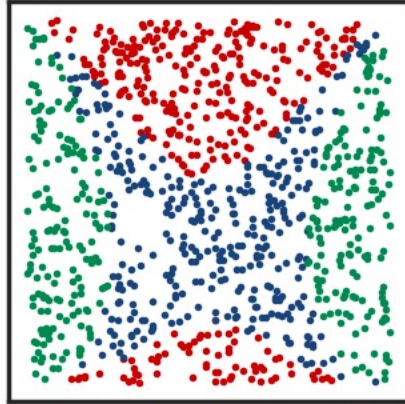
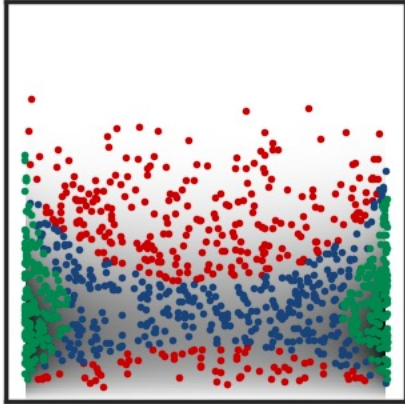
$$X \sim \mathbb{P}_X$$

$$D_X(X) \sim \text{Uniform}([0, 1]^d)$$

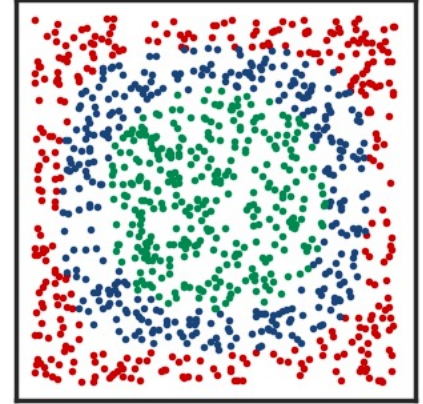
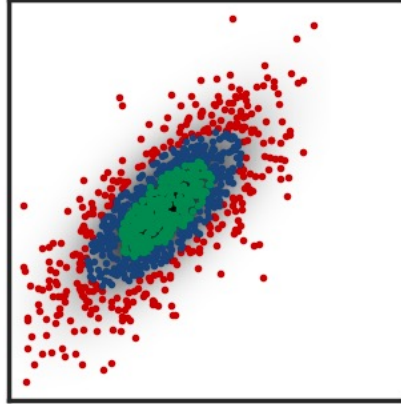
- $\Omega^{-1}(D_X) = |\det J_{D_X}| = \mathbb{P}_X \leftarrow \text{Closed-form density!}$
- Different from multivariate CDF function:  $F(x): \mathbb{R}^d \rightarrow [0, 1]$

# Many shallow densities can be mapped to destructors

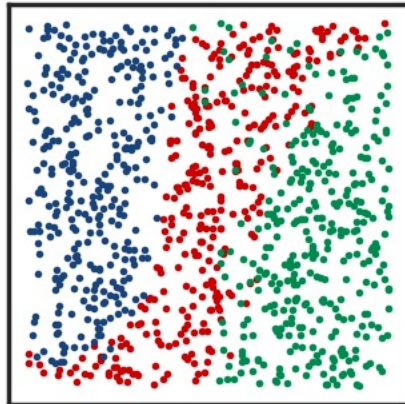
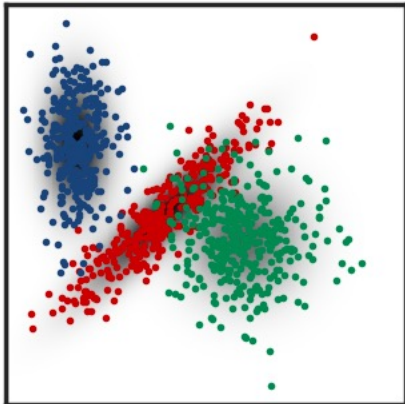
Independent (Beta distributions)



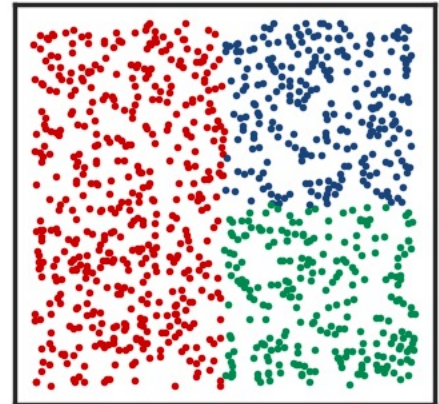
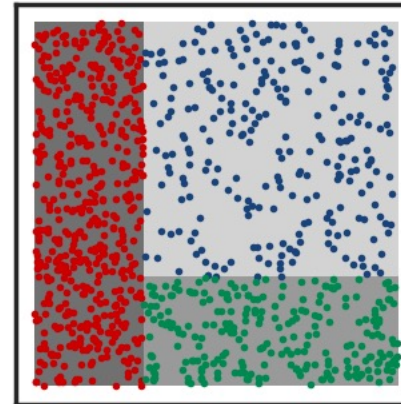
Multivariate Gaussian



Gaussian Mixture



Decision Tree Density



Data before (left) and after (right) transformation via corresponding density destructor.  
Note: Color is just to show correspondence between areas before and after transformation.



# Examples of simple closed-form destructors

Description	Density	Transformation
<b>Autoregressive Density</b>	$\prod_{s=1}^d \mathbb{P}_s(x_s   \mathbf{x}_{1:s-1})$	$[F_1(x_1), F_2(x_2   x_1), \dots, F_d(x_d   \mathbf{x}_{1:s-1})]$
Mixture of Gaussians Conditionals (e.g. MADE, MAF)	$\prod_{s=1}^d \left[ \sum_{t=1}^m \pi_t(\mathbf{x}_{1:s-1}) \times \mathbb{P}_{\mathcal{N}}(x_s   \mu_{st}(\mathbf{x}_{1:s-1}), \sigma_{st}^2(\mathbf{x}_{1:s-1})) \right]$	$\left[ F_1(x_1), F_2(x_2   x_1), \dots, F_d(x_d   x_1, \dots, x_{s-1}) \right]$
Block Gaussian Conditionals (e.g. Real NVP, NICE)	$\mathbb{P}_{\mathcal{N}}(\mathbf{x}_{1:t}   0, \mathbf{I}) \times \mathbb{P}_{\mathcal{N}}(\mathbf{x}_{t+1:d}   \boldsymbol{\mu}(\mathbf{x}_{1:t}), \boldsymbol{\sigma}^2(\mathbf{x}_{1:t}))$	$\left[ \Phi(\mathbf{x}_{1:t}), \Phi\left(\frac{x_{t+1} - \mu_{t+1}(\mathbf{x}_{1:t})}{\sigma_{t+1}(\mathbf{x}_{1:t})}\right), \dots, \Phi\left(\frac{x_d - \mu_d(\mathbf{x}_{1:t})}{\sigma_d(\mathbf{x}_{1:t})}\right) \right]$
<b>Linear Projection Density</b>	$\mathbb{P}_{\psi}(W\mathbf{x})$	$D_{\theta}(W\mathbf{x})$
Independent Components (e.g. Gaussianization via ICA)	$\prod_{s=1}^d \mathbb{P}_s(\mathbf{w}_s^T \mathbf{x})$	$\mathbf{F}(W\mathbf{x})$
Gaussian (e.g. via PCA)	$\mathbb{P}_{\mathcal{N}}(\mathbf{x}   \boldsymbol{\mu}, \Sigma)$	$\Phi(\Sigma^{-\frac{1}{2}}(\mathbf{x} - \boldsymbol{\mu}))$
<b>Copula-based Density</b>	$\mathbb{P}^{\text{cop}}(\mathbf{F}(\mathbf{x})) \prod_{s=1}^d \mathbb{P}_s(x_s)$	$D_{\theta}(\mathbf{F}(\mathbf{x}))$
Gaussian Copula	$\mathbb{P}_R^{\mathcal{N}\text{-cop}}(\mathbf{F}(\mathbf{x})) \prod_{s=1}^d \mathbb{P}_s(x_s)$	$\Phi(R^{-\frac{1}{2}}\Phi^{-1}(\mathbf{F}(\mathbf{x})))$
<b>Gaussian Mixture</b> (note that $F_s(x_s   \mathbf{x}_{-s})$ is computable)	$\sum_{t=1}^m \pi_t \mathbb{P}_{\mathcal{N}}(\mathbf{x})$	$[F_1(x_1), F_2(x_2   x_1), \dots, F_d(x_d   x_1, \dots, x_{s-1})]$
Examples of new destructors enabled by our unified destructor framework		
<b>Piecewise Density (or Tree Density)</b>	$\{\mathbb{P}_{\psi_{\ell}}(\mathbf{x}), \text{ if } \mathbf{x} \in \mathcal{L}_{\ell}\},$ where $\mathcal{L}_{\ell}$ are the disjoint subspaces of the leaves.	$\{D_{\theta_{\ell}}(\mathbf{x}), \text{ if } \mathbf{x} \in \mathcal{L}_{\ell}\}$
Piecewise Uniform (e.g. DET)	$\{c_{\ell}, \text{ if } \mathbf{x} \in \mathcal{L}_{\ell}\}$	$\{\text{diag}(\mathbf{a}_{\ell})\mathbf{x} + \mathbf{b}_{\ell}, \text{ if } \mathbf{x} \in \mathcal{L}_{\ell}\}$
<b>Image-Specific Feature Pairs</b>	$\prod_{P \in \mathcal{P}} \mathbb{P}_P(x_{P(1)}, x_{P(2)}),$ where feature pairs $\mathcal{P}$ are based on pixel locality.	$\{D_P(x_{P(1)}, x_{P(2)}), \forall P \in \mathcal{P}\}$

# Density destructor algorithm performs greedy layer-wise construction of deep destructor

1. Simple density estimation (GMM, Gaussian, tree density, etc.)

$$Q^t \leftarrow \arg \min_{Q \in \mathcal{Q}} KL(P(x^{t-1}), Q(x^{t-1}))$$

2. Map density to simple destructor layer

$$d^t = \Omega(Q^t)$$

3. Transform data for next layer

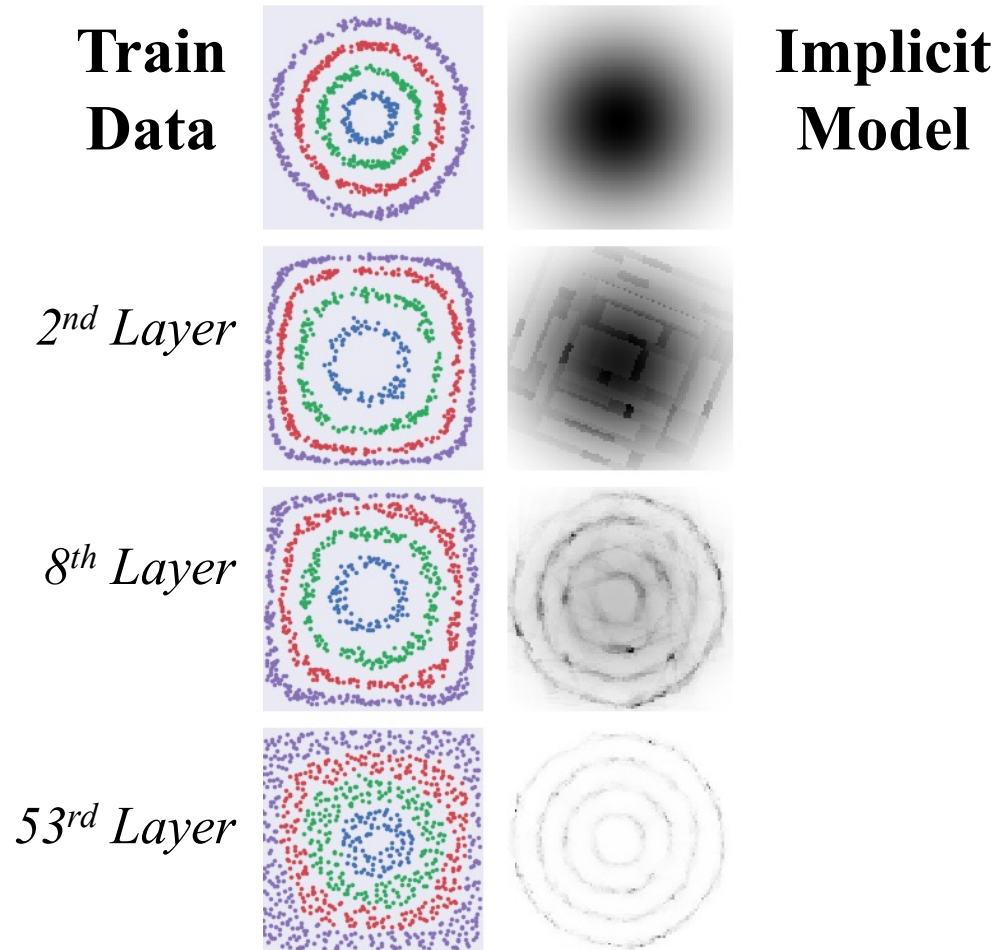
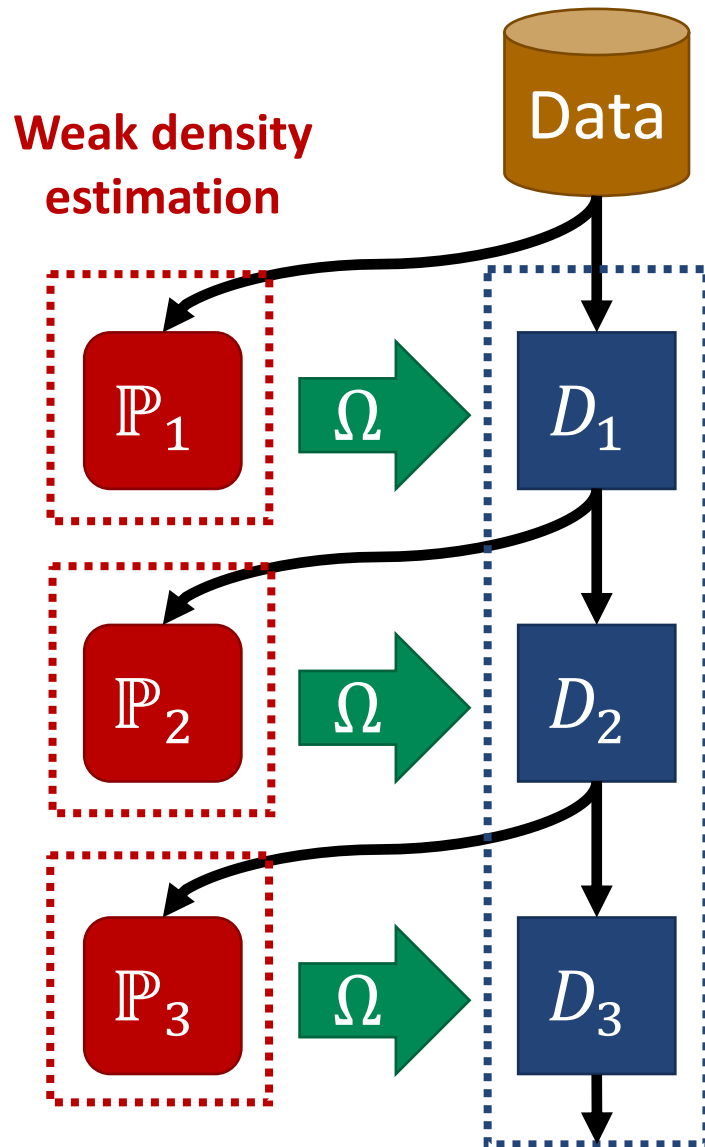
$$x^t = d^t(x^{t-1})$$

4. Update deep destructor

$$D^t = d^t \circ D^{t-1}$$



# Deep density destructors via sequence of *weak* destructors



# Destructor algorithm can be shown to monotonically decrease the negative log likelihood after every iteration/layer

- ▶ The destructive learning objective, where  $z = D(x)$ , and  $U(z)$  is the uniform density function

$$\arg \min_D KL(P_z(z), U(z))$$

- ▶ Want: Every iteration decreases objective:

$$KL(P_{d^t \circ D^{t-1}(x)}, U) \leq KL(P_{D^{t-1}(x)}, U)$$

- ▶ Let  $x = D^{t-1}(x^{(0)})$  and  $z = d^t(x) = d^t(D^{t-1}(x^{(0)}))$

- ▶  $KL(P_z(z; D), U(z))$

- ▶  $= KL(P_x(x), Q_x(x; D))$

(KL equivalence lemma)

- ▶  $\leq KL(P_x(x), Q_x(x; D = Id))$

(minimization is better than one particular)

- ▶  $= KL(P_x(x), |J_D(x)|U(D(x)))$  (Expand in terms of  $D$ )

- ▶  $= KL(P_x(x), U(x))$

# Reuse results: Deep density destructors can be built from simple and well-understood components

- ▶ MNIST  $d = 784$
- ▶ CIFAR-10  $d = 3072$
- ▶ Autoregressive flow baselines (**DNN-based**)
  - ▶ MADE [Germain et al., 2015]
  - ▶ Real NVP [Dinh, et al. 2017]
  - ▶ MAF [Papamakarios et al. 2017]
- ▶ Our deep copula method
  - ▶ **PCA + histograms**

	MNIST			CIFAR-10		
	LL	D	T	LL	D	T
<i>Models from MAF paper computed on Titan X GPU</i>						
Gaussian	-1367	1	0.0	2367	1	0.0
MADE	-1385	1	0.0	448	1	0.2
MADE MoG	-1042	1	0.1	-53	1	0.3
Real NVP	-1329	5	0.2	2600	5	1.4
Real NVP	-1765	10	0.2	2469	10	1.0
MAF	-1300	5	0.1	2941	5	3.7
MAF	-1314	10	0.2	<b>3054</b>	10	7.5
MAF MoG	-1100	5	0.2	2822	5	3.9

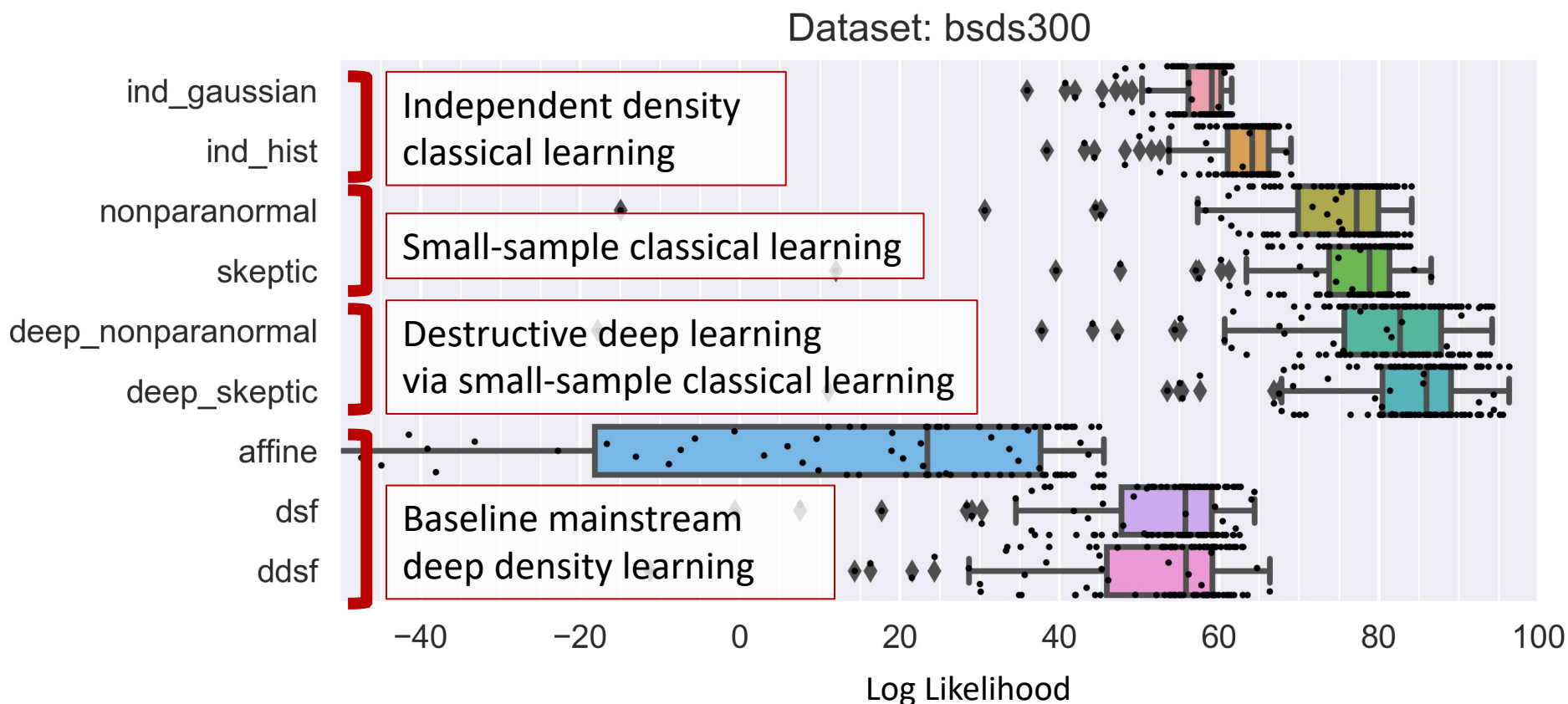
*Our proposed destructors computed on 10 CPUs*

Copula	-1028	5	0.2	2626	17	10.1
--------	-------	---	-----	------	----	------

LL = Log Likelihood (higher is better)

D = # of layers, T = Time

# Modularity enables classical learning improvements to carry over to deep learning



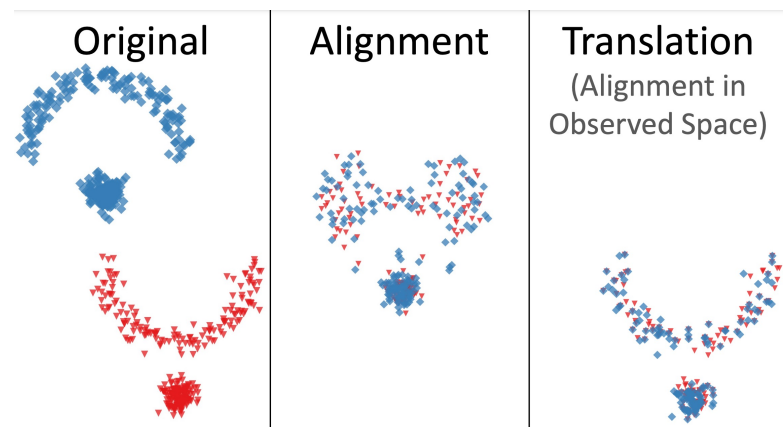
Small-sample experiment where number of dimensions is 63 and number of training samples is 30. Notice how mainstream deep learning fails in this setting.

# Limitations of destructive modular learning

- ▶ Unlikely to perform as well as joint learning
  - ▶ Greedy vs joint optimization
  - ▶ Local vs global optimization
- ▶ Must create destructor mapping  $\Omega$ , which can be challenging
- ▶ Often requires more layers to achieve similar result because of optimization
- ▶ Normalizing flows transform to simple known distribution
  - ▶ What about transforming between any two distributions?

# Translating or aligning two arbitrary distributions is a more general task

- ▶ Flow-based methods for alignment/translation
  - ▶ Likelihood-based: AlignFlow<sup>[1]</sup>, LRMF<sup>[2]</sup>
  - ▶ Optimal Transport: Iterative Alignment Flow



- ▶ Iterative Alignment Flow:  
(Greedy) fast and easy alignment  
of multiple distributions!

[1] Aditya Grover, Christopher Chute, Rui Shu, Zhangjie Cao, Stefano Ermon: AlignFlow: Cycle Consistent Learning from Multiple Domains via Normalizing Flows. AAAI 2020: 4028-4035

[2] Ben Usman, Avneesh Sud, Nick Dufour, Kate Saenko: Log-Likelihood Ratio Minimizing Flows: Towards Robust and Quantifiable Neural Distribution Alignment. NeurIPS 2020

# Iterative Alignment Flow: Extending iterative approach to translation/alignment task

- ▶ Motivation: Hard to directly align two or more distributions in a high dimensional space
- ▶ Proposed solution: Decompose the high dimensional problem into a series of **simple** 1D problems
- ▶ Background
  - ▶ Optimal transport in 1D
  - ▶ Max-sliced Wasserstein Distance (max-SWD)
- ▶ Objective and iterative algorithm
  - ▶ Min-max optimization (different from GAN)
  - ▶ Iterative algorithm

## Background: Optimal transport problems in 1D are known in closed-form!

- ▶ Computing the Wasserstein distance in 1D is known in **closed-form**

$$W_1(\hat{p}_x, \hat{q}_y) = \sum_{i=1}^n |x_i - y_i|$$

- ▶ Where the data is **sorted**  $x_1 \leq x_2 \leq \dots \leq x_n$  and  $y_1 \leq y_2 \leq \dots \leq y_n$
- ▶ The optimal Monge map is known in **closed-form**

$$T^*(x) = F_Y^{-1}(F_X(x))$$



## Background: Max sliced Wasserstein distance simplifies to a “worst case” 1D Wasserstein problem

- ▶ Recall Wasserstein-1 Distance

$$W_1(p_X, p_Y) = \left( \min_T \mathbb{E}_{p_X}[\|x - T(x)\|_1] \right. \\ \left. \text{s. t. } p_{T(X)} = p_Y \right)$$

- ▶ Max sliced Wasserstein-1 divergence

$$\max\_SW_1(p_X, p_Y) = \max_{\theta: \|\theta\|_2=1} W_1(p_{X^T\theta}, p_{Y^T\theta})$$

- ▶ Where  $\theta$  defines the direction of the projection/slice
  - ▶ Note that this finds the **largest** difference along a **1D** projection/slice
- ▶ The max-SW is 0 if and only if the distributions are aligned

# Iterative alignment flow objective is a different type of min-max optimization

- ▶ The objective becomes a min-max optimization:

$$\min_{T_X, T_Y} \max_{-} SW_1(p_{T_X(X)}, p_{T_Y(Y)})$$
$$= \min_{T_X, T_Y} \max_{\theta: \|\theta\|_2=1} W_1(p_{T_X(X)}^T \theta, p_{T_Y(Y)}^T \theta)$$

- ▶ The inner optimization tries to find the “worst-case” projection that maximizes the divergence.
- ▶ The outer optimization tries to find the transformation that minimizes this “worst-case” divergence.
- ▶ Theoretical minimum occurs when distributions are aligned per the property of max-SW

# Algorithm: Iteratively solve min and max problems

- Find structure  $\theta$  by solving max problem

$$\theta^{(i)} = \arg \max_{\theta: \|\theta\|_2=1} W_1(p_{X^T \theta}, p_{Y^T \theta})$$

- Fix  $\theta^{(i)}$  and solve for simple OT Monge maps via **closed-form solutions** (i.e.,  $\Omega$ )

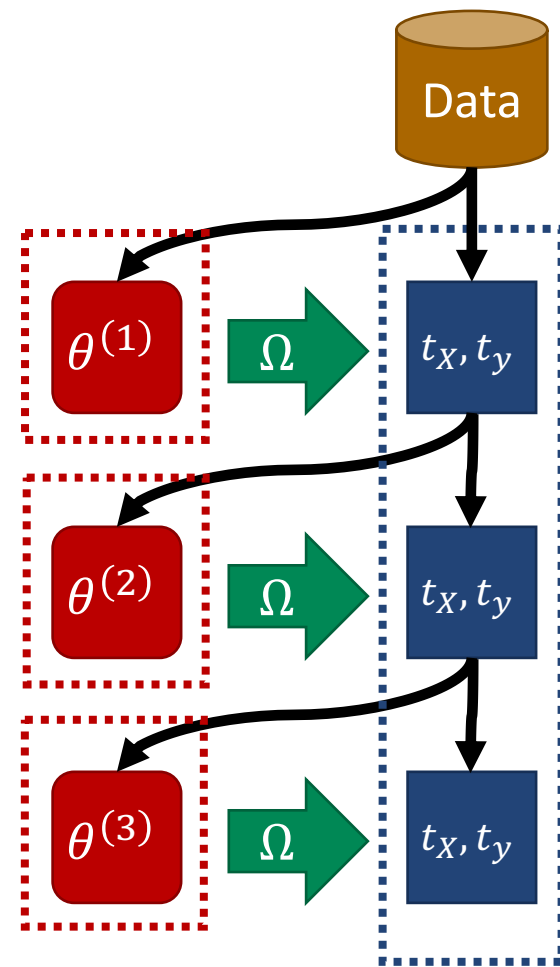
$$t_X^{(i)}, t_Y^{(i)} = \arg \min_{t_X, t_Y} W_1(p_{t_X(X)^T \theta^{(i)}}, p_{t_Y(Y)^T \theta^{(i)}})$$

- Transform data

$$X \leftarrow t_X^{(i)}(X), Y \leftarrow t_Y^{(i)}(Y)$$

- Update deep transforms

$$T_X^{(i)} \leftarrow t_x^{(i)} \circ T_X^{(i-1)}$$



Iterative alignment flows can be used to translate between multiple domains via this simple iterative algorithm

