

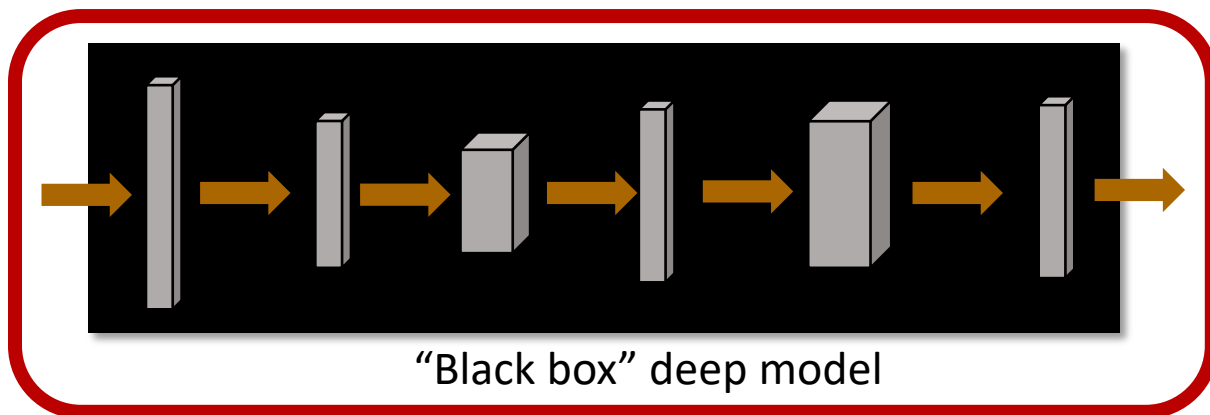
Deep Density Destructors (from a biased viewpoint)

David I. Inouye

Electrical and Computer Engineering
Purdue University

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

End-to-end learning

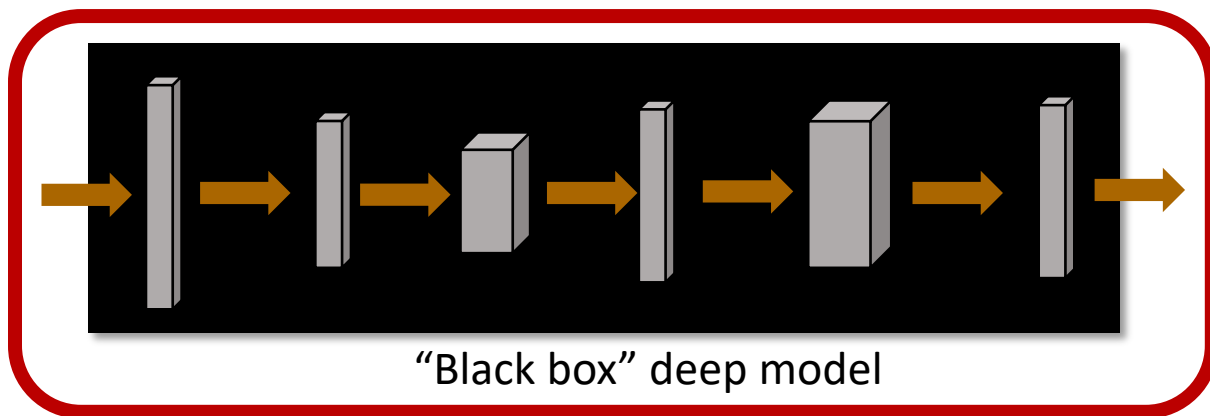


- ▶ Real NVP
- ▶ MAF
- ▶ GLOW
- ▶ Etc.

"Gray box" deep model

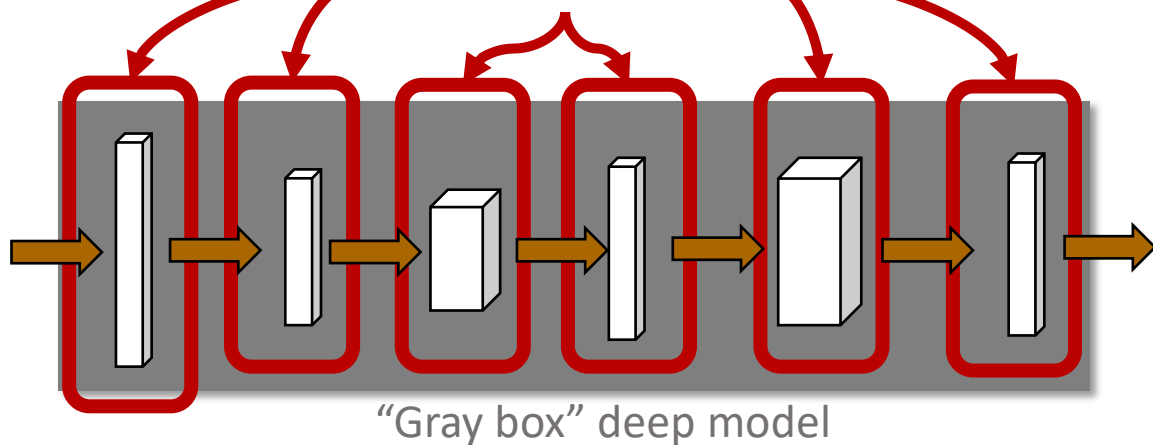
Modular deep learning would allow *local* learning within each component

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

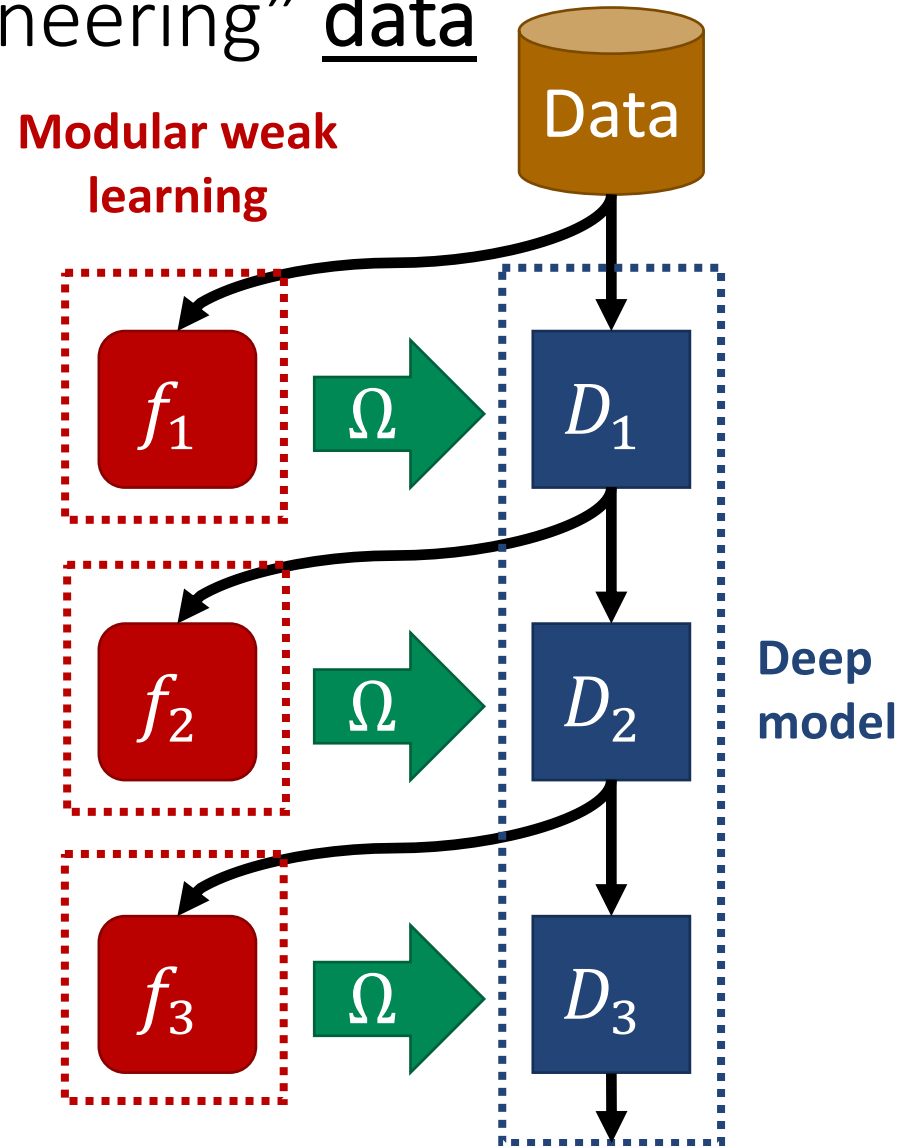
1. **Find *part*** to take off using **understanding and expertise**
2. **Determine how** to take off part in a *reversible* way (e.g., unscrewing bolts)
3. **Remove part**
4. **Repeat**

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

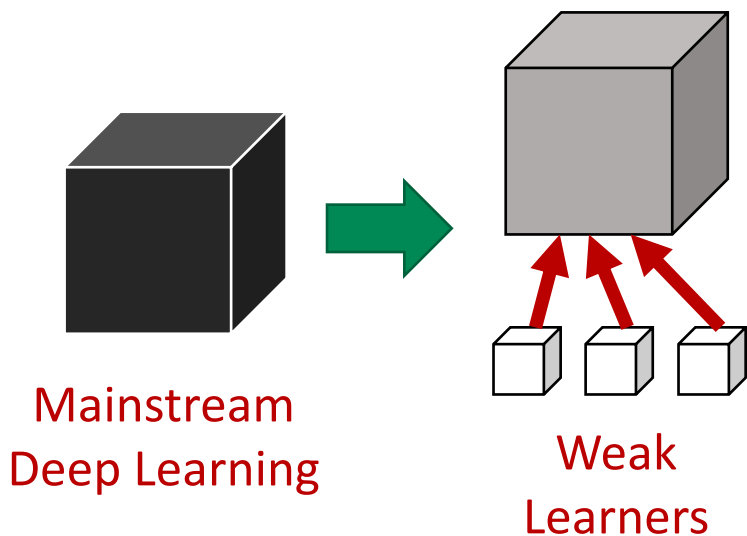
1. **Find** patterns in data via **shallow/weak learning**
2. **Map model** to destructive transformation
3. **Destroy the patterns** via transformation
4. **Repeat**



Why use modular weak learning for deep models?

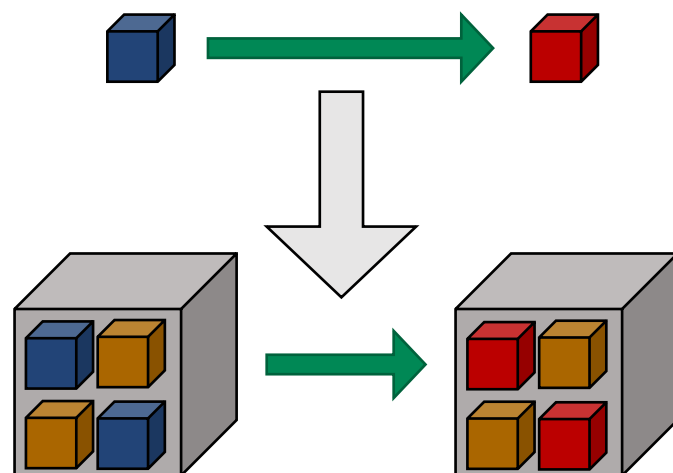
Reuse

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



Decoupling

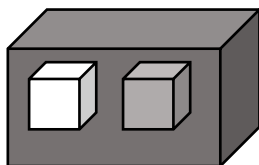
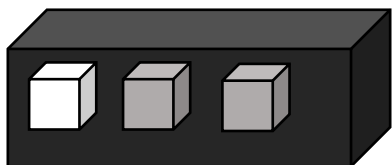
Components can be **debugged, tested and improved** separate from the system



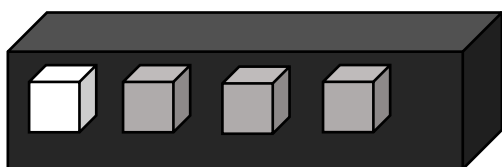
Why use modular weak learning for deep models?

Algorithmic Interpretability

Increasing or decreasing model complexity is straightforward



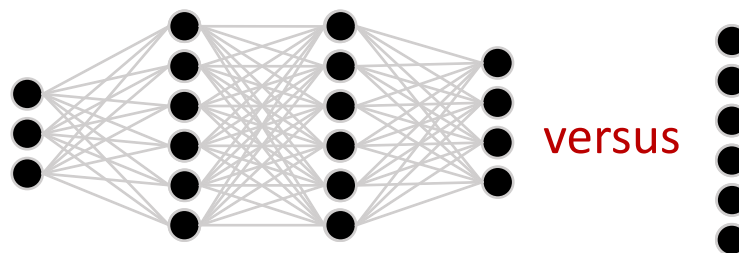
Shrink model if problem



Grow if more data

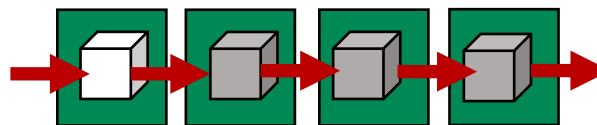
Resource Constraints

Layer-wise training (memory bottleneck)



Pipelined training (computation bottleneck)

Shallow/weak online learners



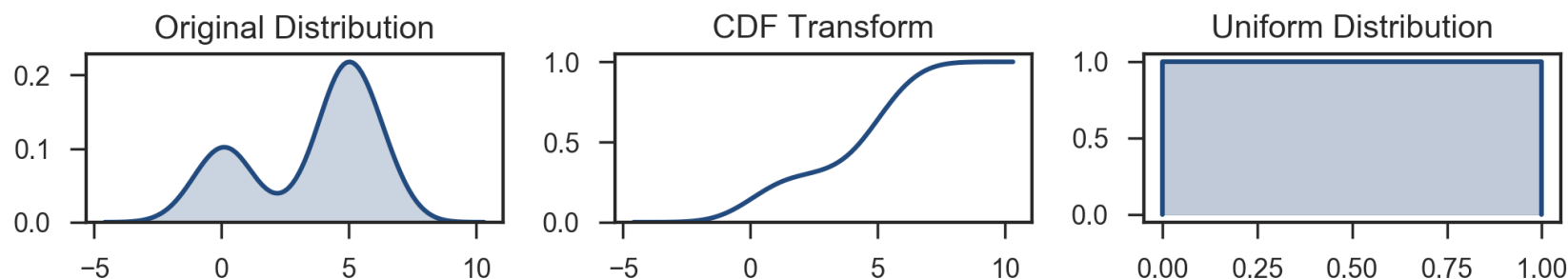
Distributed on different processors or devices

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 Ω , which can be challenging
- ▶ Often requires more layers to achieve similar result because of optimization

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

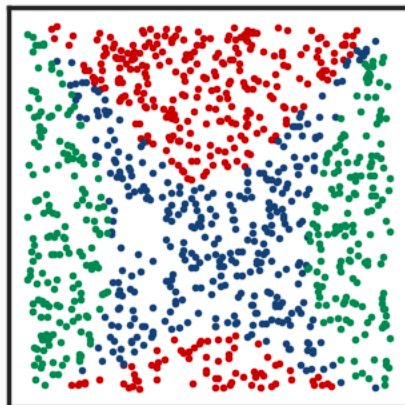
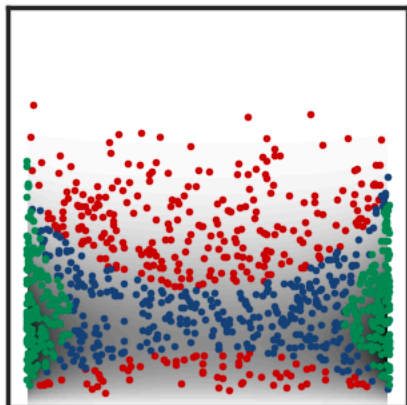
$$X \sim \mathbb{P}_X$$

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

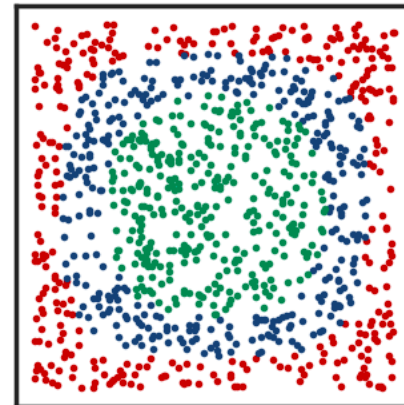
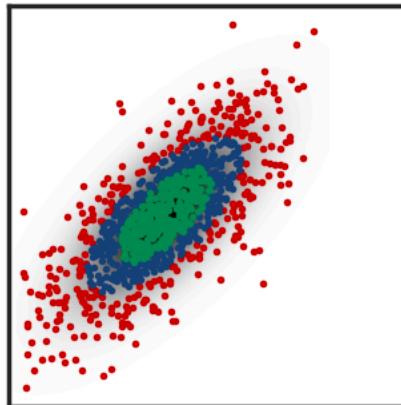
- $\Omega^{-1}(D_X) = |\det J_{D_X}| = \mathbb{P}_X \leftarrow$ **Closed-form density!**
- Different from multivariate CDF function: $F(x): \mathbb{R}^d \rightarrow [0, 1]$

Many shallow densities can be mapped to destructors

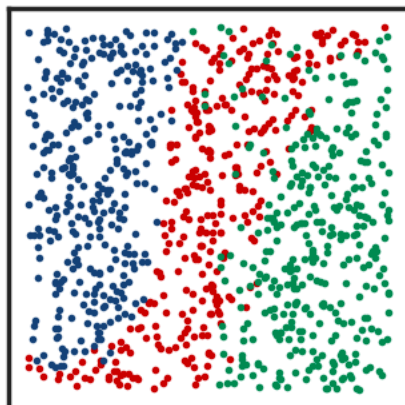
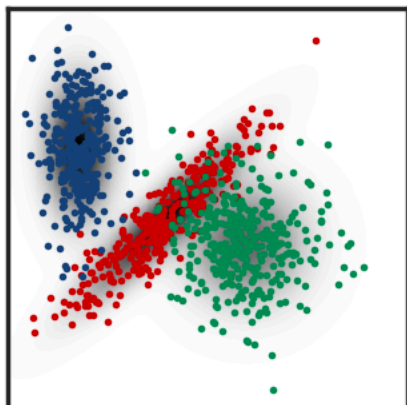
Independent (Beta distributions)



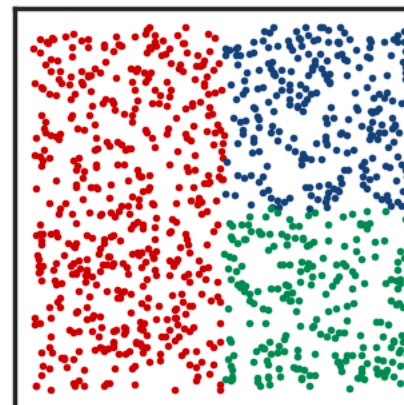
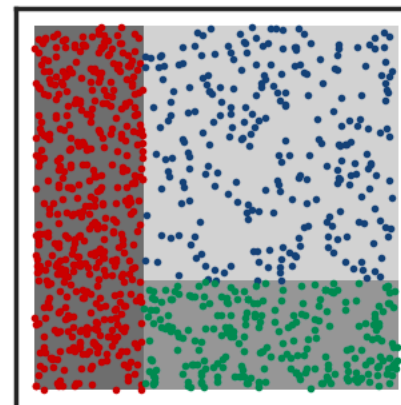
Multivariate Gaussian



Gaussian Mixture



Decision Tree Density

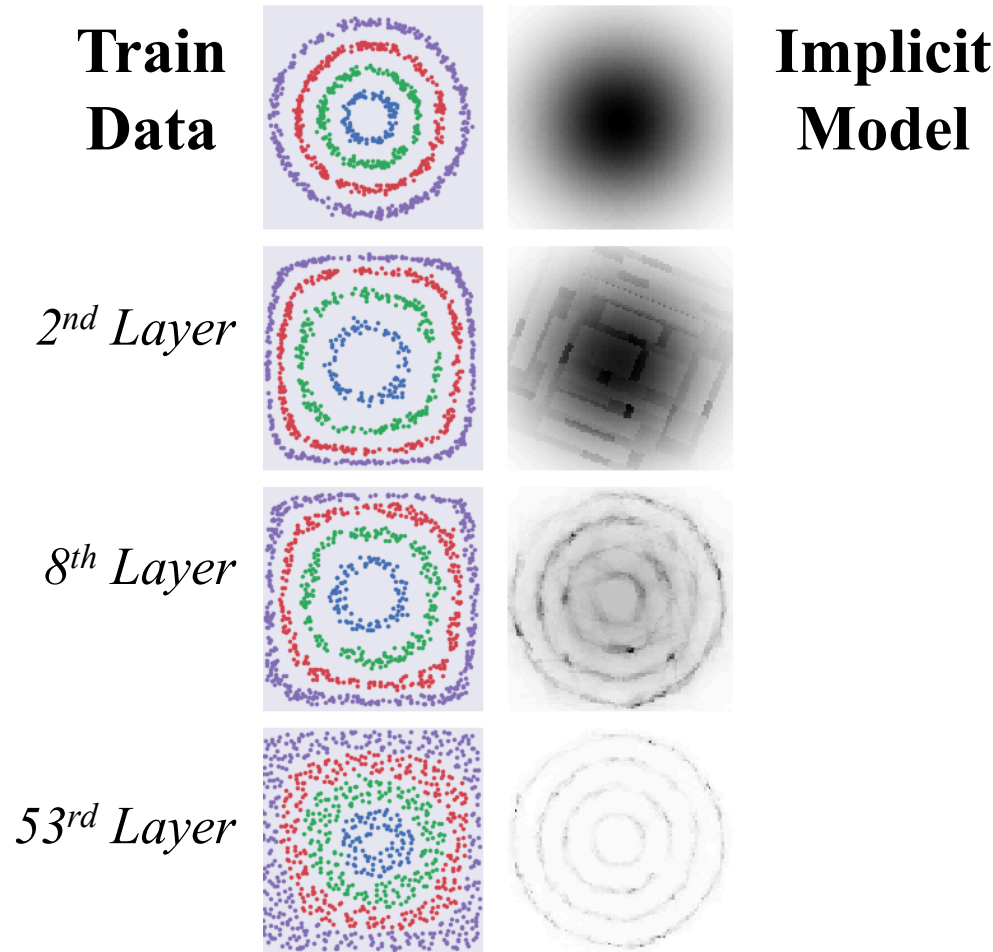
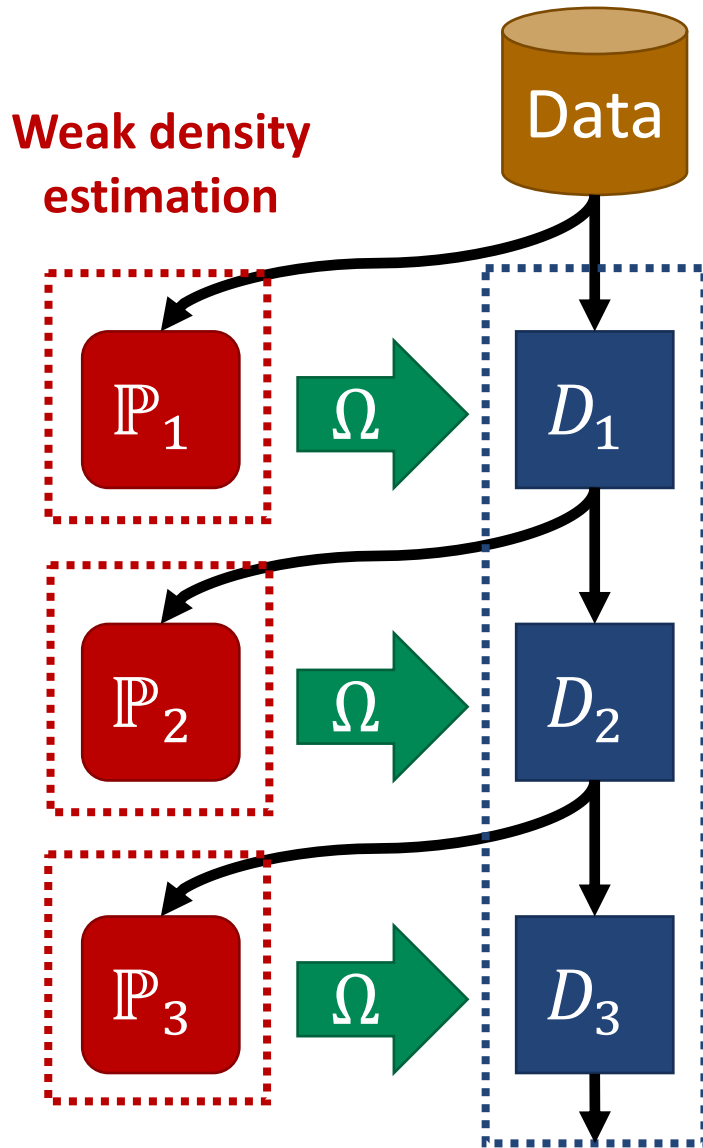


Data before (left) and after (right) transformation via corresponding density destructor.
Note: Color is just to show correspondence between areas before and after transformation.

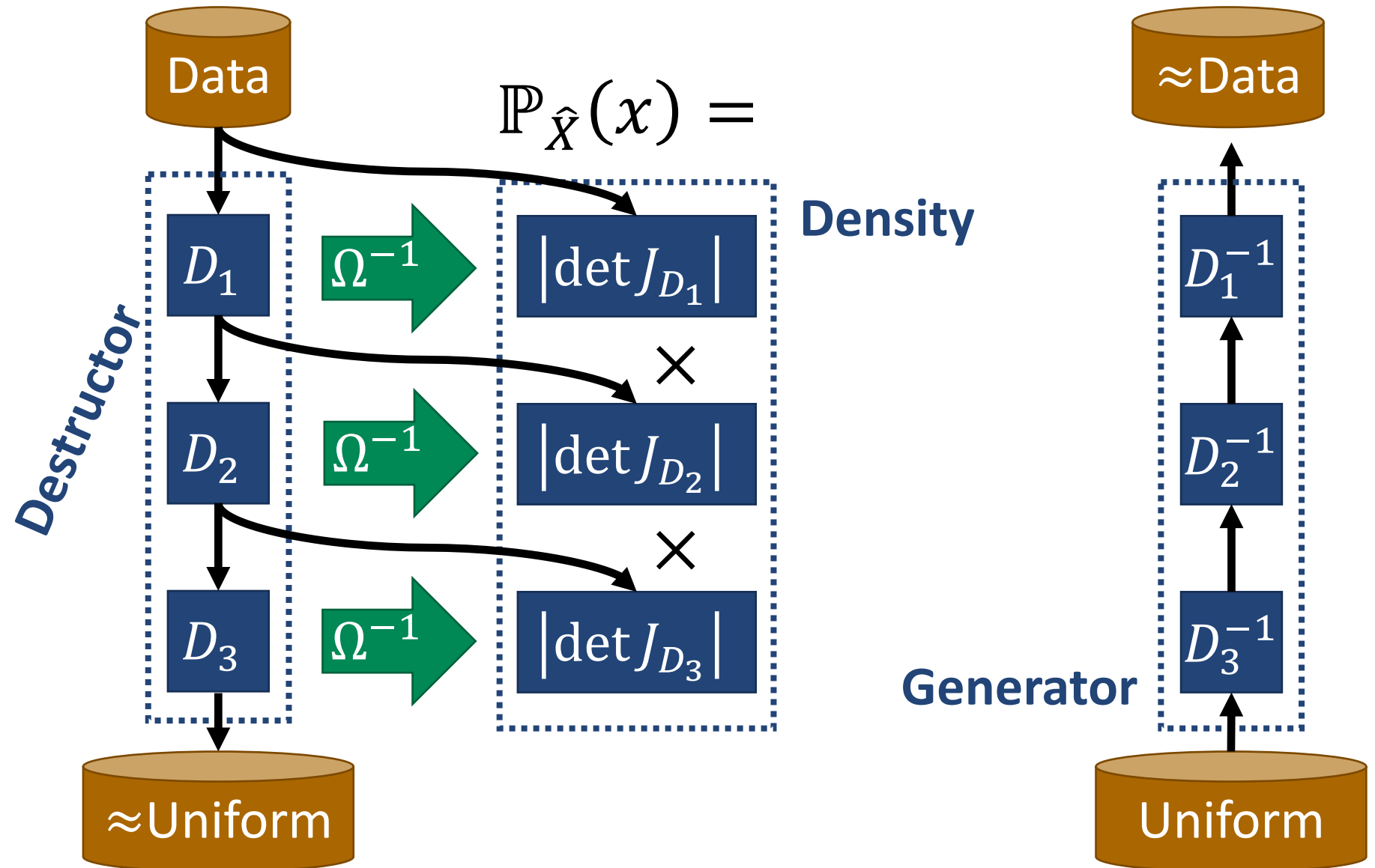
Example Destructors

| Description | Density | Transformation |
|--|--|--|
| Autoregressive Density | $\prod_{s=1}^d \mathbb{P}_s(x_s \mathbf{x}_{1:s-1})$ | $[F_1(x_1), F_2(x_2 x_1), \dots, F_d(x_d \mathbf{x}_{1:s-1})]$ |
| 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) | $\mathbb{P}_{\mathcal{N}}(\mathbf{x}_{1:t} 0, \mathbf{I}) \times \mathbb{P}_{\mathcal{N}}(\mathbf{x}_{t+1:d} \boldsymbol{\mu}(\mathbf{x}_{1:t}), \boldsymbol{\sigma}^2(\mathbf{x}_{1:t}))$ | $\left[\Phi(\mathbf{x}_{1:t}), \Phi\left(\frac{x_{t+1} - \mu_{t+1}(\mathbf{x}_{1:t})}{\sigma_{t+1}(\mathbf{x}_{1:t})}\right), \dots, \Phi\left(\frac{x_d - \mu_d(\mathbf{x}_{1:t})}{\sigma_d(\mathbf{x}_{1:t})}\right) \right]$ |
| Linear Projection Density | $\mathbb{P}_{\psi}(W\mathbf{x})$ | $D_{\theta}(W\mathbf{x})$ |
| 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) | $\mathbb{P}_{\mathcal{N}}(\mathbf{x} \boldsymbol{\mu}, \Sigma)$ | $\Phi(\Sigma^{-\frac{1}{2}}(\mathbf{x} - \boldsymbol{\mu}))$ |
| Copula-based Density | $\mathbb{P}^{\text{cop}}(\mathbf{F}(\mathbf{x})) \prod_{s=1}^d \mathbb{P}_s(x_s)$ | $D_{\theta}(\mathbf{F}(\mathbf{x}))$ |
| Gaussian Copula | $\mathbb{P}_R^{\mathcal{N}\text{-cop}}(\mathbf{F}(\mathbf{x})) \prod_{s=1}^d \mathbb{P}_s(x_s)$ | $\Phi(R^{-\frac{1}{2}}\Phi^{-1}(\mathbf{F}(\mathbf{x})))$ |
| Gaussian Mixture (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 | | |
| Piecewise Density (or Tree Density) | $\{\mathbb{P}_{\psi_{\ell}}(\mathbf{x}), \text{ if } \mathbf{x} \in \mathcal{L}_{\ell}\},$ where \mathcal{L}_{ℓ} 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) | $\{c_{\ell}, \text{ if } \mathbf{x} \in \mathcal{L}_{\ell}\}$ | $\{\text{diag}(\mathbf{a}_{\ell})\mathbf{x} + \mathbf{b}_{\ell}, \text{ if } \mathbf{x} \in \mathcal{L}_{\ell}\}$ |
| Image-Specific Feature Pairs | $\prod_{P \in \mathcal{P}} \mathbb{P}_P(x_{P(1)}, x_{P(2)}),$ where feature pairs \mathcal{P} are based on pixel locality. | $\{D_P(x_{P(1)}, x_{P(2)}), \forall P \in \mathcal{P}\}$ |

Deep density destructors via sequence of *weak* destructors



Density computation and sample generation



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

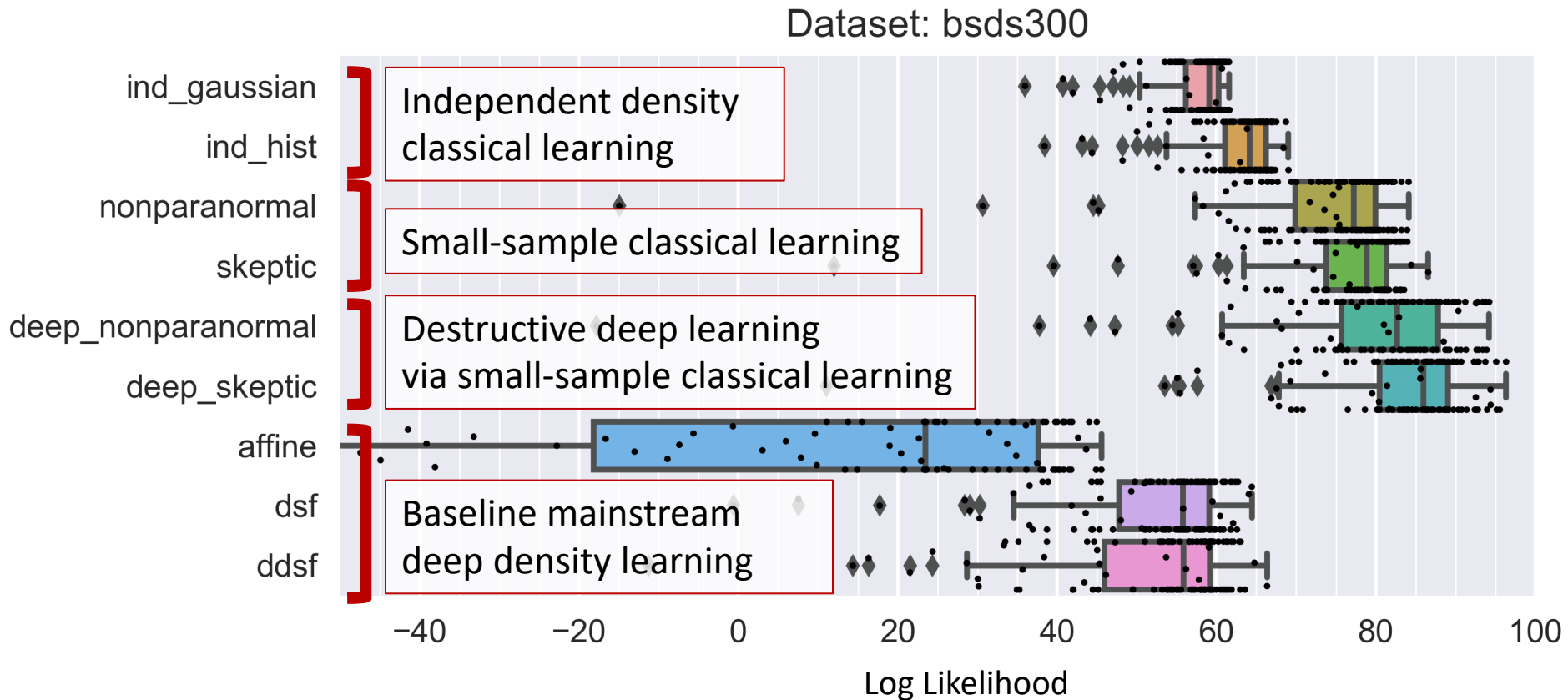
- ▶ MNIST $d = 784$
- ▶ CIFAR-10 $d = 3072$
- ▶ Autoregressive flow baselines (**DNN-based**)
 - ▶ MADE [Germain et al., 2015]
 - ▶ Real NVP [Dinh, et al. 2017]
 - ▶ MAF [Papamakarios et al. 2017]
- ▶ Our deep copula method
 - ▶ **PCA + histograms**

| | 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 | 3054 | 10 | 7.5 |
| MAF MoG | -1100 | 5 | 0.2 | 2822 | 5 | 3.9 |
| <i>Our proposed destructors computed on 10 CPUs</i> | | | | | | |
| 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.

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

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

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 $U_z(z)$ is the uniform density function

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

- ▶ KL equivalence lemma, let $z = D(x)$ for invertible D

$$KL(P_x(x), Q_x(x)) = KL(P_z(z), Q_z(z))$$

- ▶ Simple corollary is that objective above is MLE:

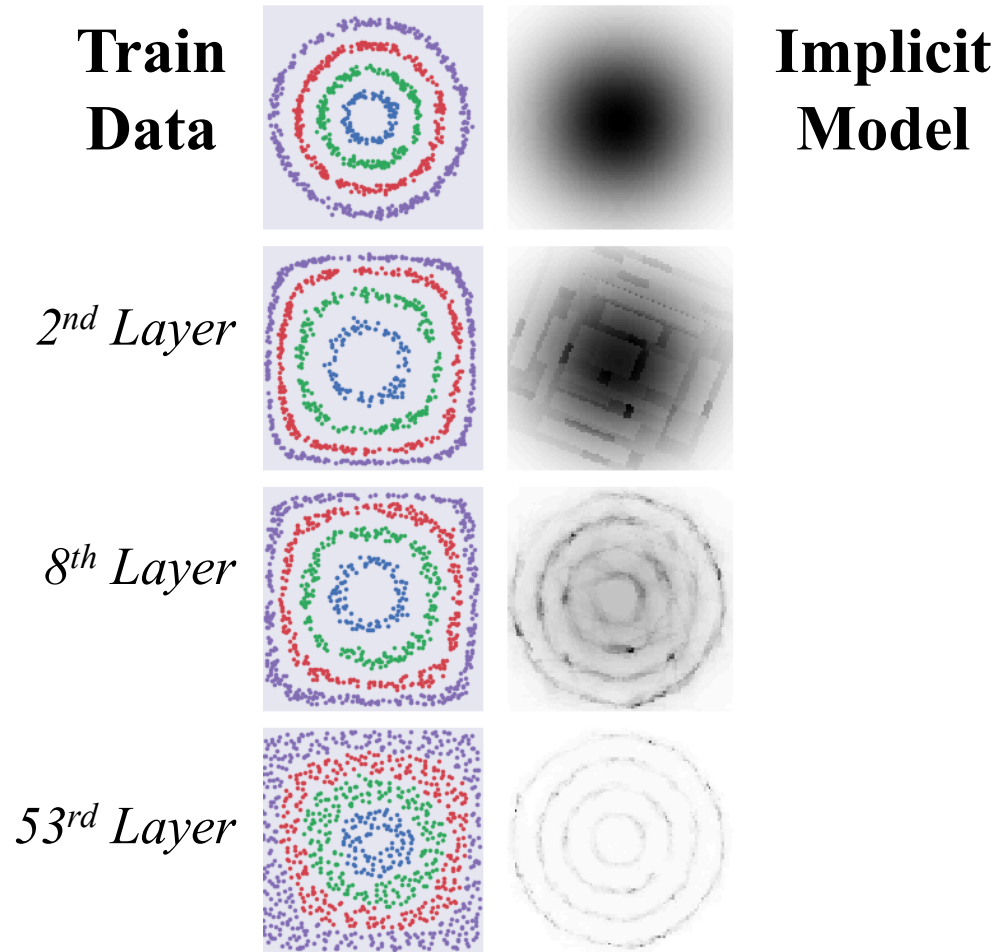
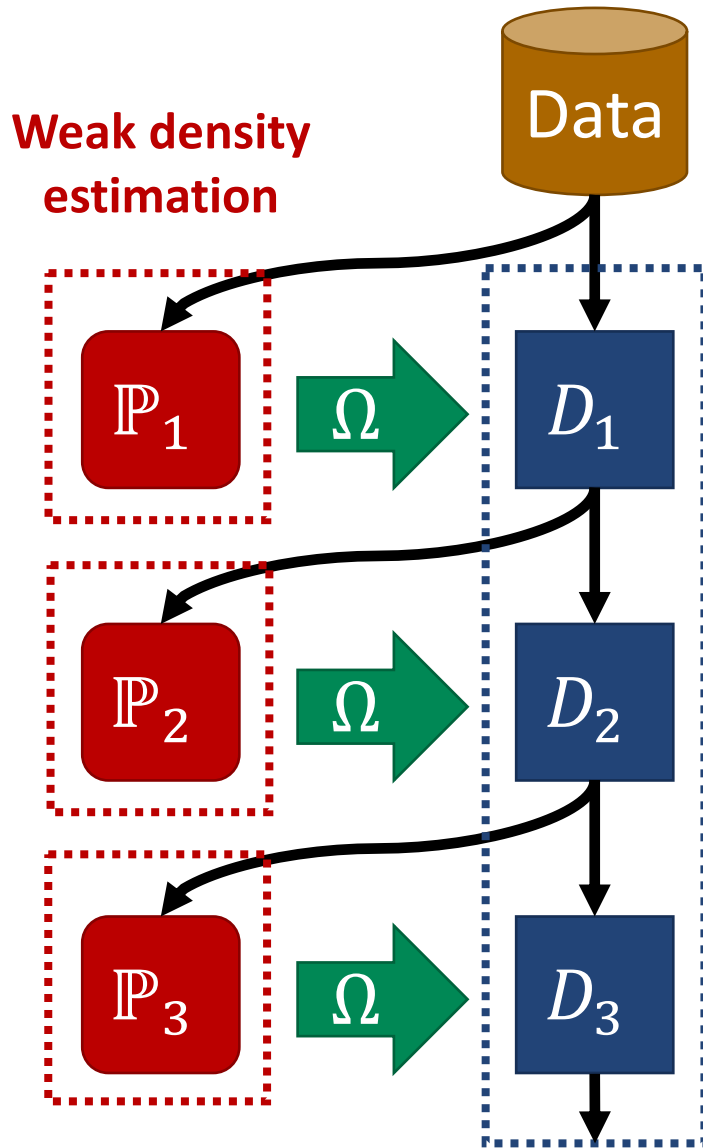
- ▶ $KL(P_z(z; D), U(z))$

- ▶ $= KL(P_x(x), Q_x(x; D))$ (KL equivalence, MLE objective)

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

- ▶ $= KL(P_x(x), |J_D(x)|)$ ($U(z) = 1$)

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:

$$KL(P_z(d^t(x)), U(d^t(x))) \leq KL(P_x(d^{t-1}(x)), U(d^{t-1}(x)))$$

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

- ▶ $KL(P_z(z; D), U(z))$

- ▶ $= KL(P_x(x), Q_x(x; D))$

(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)

- ▶ $= KL(P_x(x), U(x))$