FAU

# Advanced Deep Learning
## Energy-based Models

**Florian Kordon[1], Katharina Breininger[2], Vincent Christlein[1]**

1 Pattern Recognition Lab, FAU
2 Artificial Intelligence in Medical Imaging, FAU

June, 7th 2023

# Intuition: What is Energy-based Modeling?

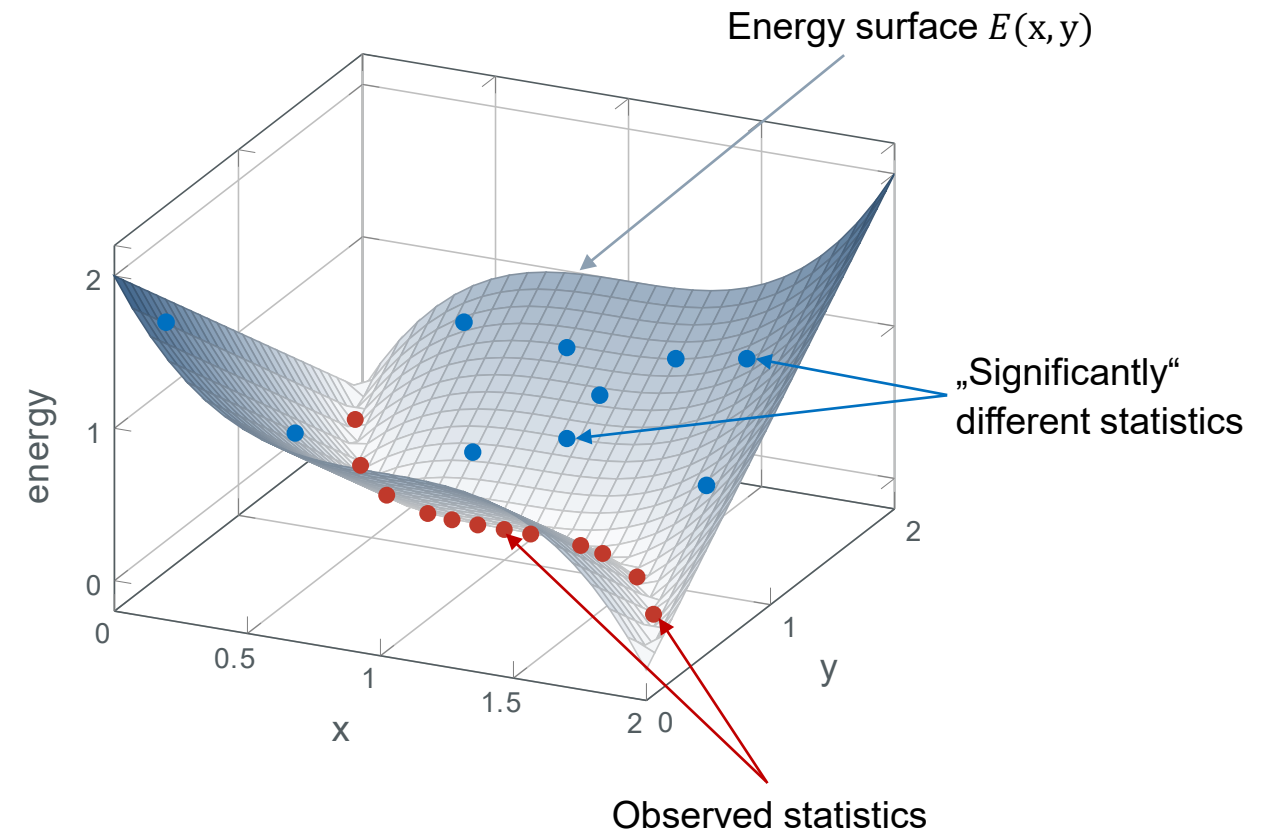# Today's lecture

What is energy-based modeling?

## Intuition

Assume some observations (red data points) that share a common set of characteristics.

These observations lie on a high-dimensional manifold $\mathcal{M}$.

→ Find an embedding function $\phi$ that maps $\mathcal{M}$ to a structure-preserving (potentially lower-dim.) representation $\mathcal{N}$.

$$\phi \colon \mathcal{M} \mapsto \mathcal{N}$$
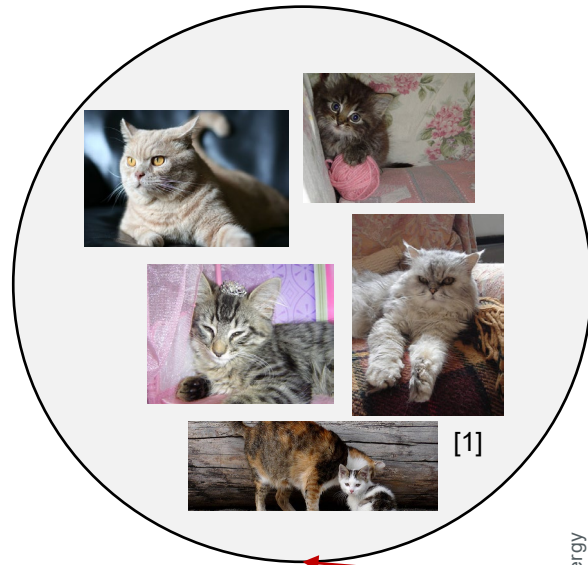
An energy-based model (EBM) performs such mapping by minimizing the value of the energy function $E$ for the observed statistics (red data points), while increasing it for all other statistics (blue data points).



Energy surface $E(\mathrm{x}, y)$

"Significantly" different statistics

Observed statistics

[1] Chris G. Willcocks. Deep Learning Lecture 7: Energy-based models. https://cwkx.github.io/data/teaching/dl-and-rl/dl-lecture7.pdf
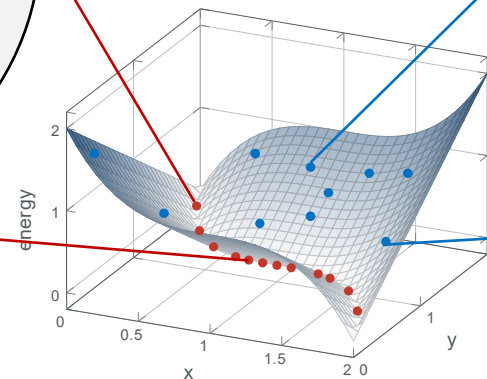
# Energy-based models

Illustrative example

We define an energy function $E(\mathbf{x})$ that assigns a small scalar value to images of different breeds of cats.

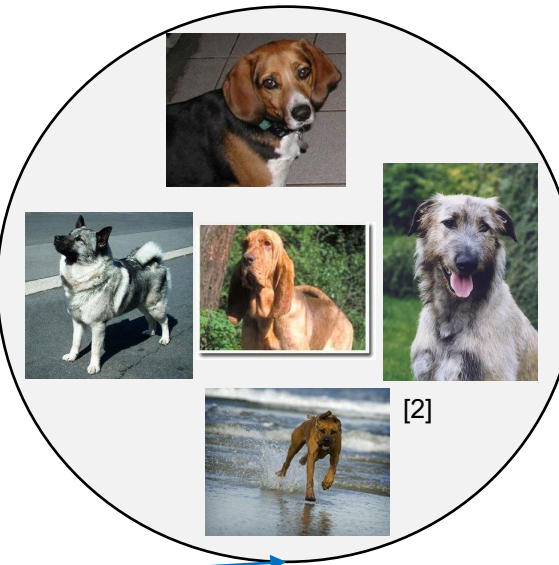**Observed cat statistics**

$$E(\mathbf{x}) \approx 0$$

**Unknown dog statistics**

$$E(\mathbf{x}') > 0$$



[1]

[2]

→ For different measurements (e.g., images of dogs), the function output should be significantly larger.

[1] https://www.kaggle.com/datasets/crawford/cat-dataset
[2] https://www.kaggle.com/datasets/jessicali9530/stanford-dogs-dataset

# Background and Fundamentals

# Energy-based models

Energy function $E$ as the backbone of energy-based models

FAU

## Fundamentals

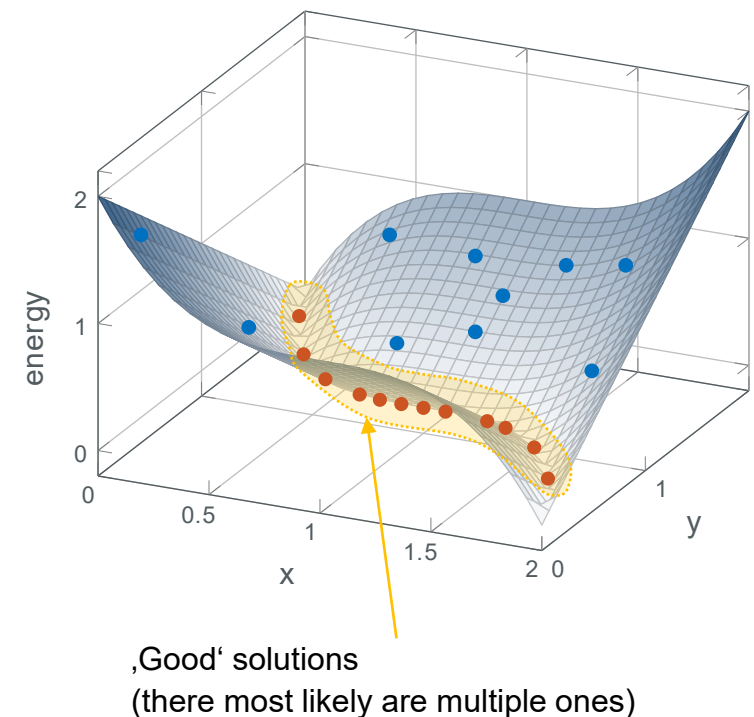The energy function $E(\mathrm{x}, \mathrm{y})$ measures the compatibility between an input data point $\mathrm{x} \in \mathbb{R}^{D_x}$ and a target variable $\mathrm{y} \in \mathbb{R}^{D_y}$ by assigning a scalar value (=energy).

$$E(\mathrm{x}, \mathrm{y}) \colon \mathbb{R}^{D_x} \times \mathbb{R}^{D_y} \mapsto \mathbb{R}$$

→ Takes small values for compatible $\mathrm{x}$ and $\mathrm{y}$, high values for less compatible $\mathrm{x}$ and $\mathrm{y}$.

→ $\mathrm{y}$ can represent, e.g., a binary or categorial assignment, a target image, or an object's coordinates.

During inference, we search for configurations of $\mathrm{x}$ for which $E(\mathrm{x}, \mathrm{y})$ is small.

$$\hat{\mathrm{y}} = \operatorname{argmin}_{\mathrm{y}} E(\mathrm{x}, \mathrm{y})$$



,Good' solutions
(there most likely are multiple ones)

[1] Chris G. Willcocks. Deep Learning Lecture 7: Energy-based models. https://cwkx.github.io/data/teaching/dl-and-rl/dl-lecture7.pdf
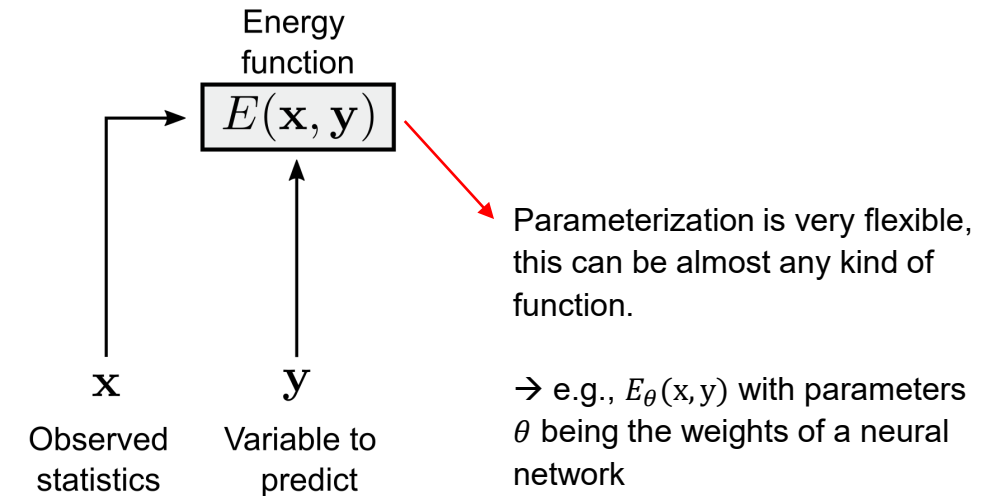
# Energy-based models

EBMs as implicit functions

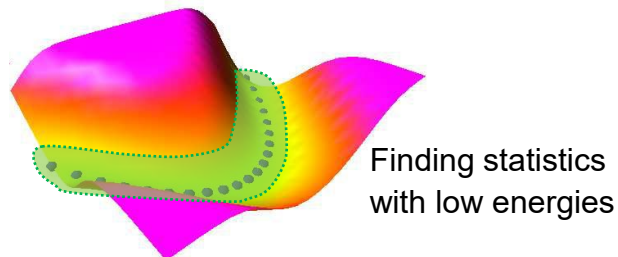→ A feed-forward model represents an explicit function that predicts $\hat{y}$ based off $x$

→ An energy-based model (EBM) is an implicit function that models the relation and compatability between $x$ and $y$

Divergence measure

$$\hat{\mathbf{y}} \rightarrow \boxed{c(\mathbf{y}, \hat{\mathbf{y}})}$$

$$g_\theta(\mathbf{x})$$

$\mathbf{x}$
Observed statistics

$\mathbf{y}$
Variable to be predicted

Energy function

$$\boxed{E(\mathbf{x}, \mathbf{y})}$$

$\mathbf{x}$
Observed statistics

$\mathbf{y}$
Variable to predict

Parameterization is very flexible, this can be almost any kind of function.

→ e.g., $E_\theta(x, y)$ with parameters $\theta$ being the weights of a neural network

[1] Yann LeCun et al. A Tutorial on Energy-Based Learning. Predicting structured data. 2006

# Energy-based models

Inference and training

## A) Inference procedure

During inference, we search for configurations of x for which the energy function $E_\theta(x, y)$ is small.

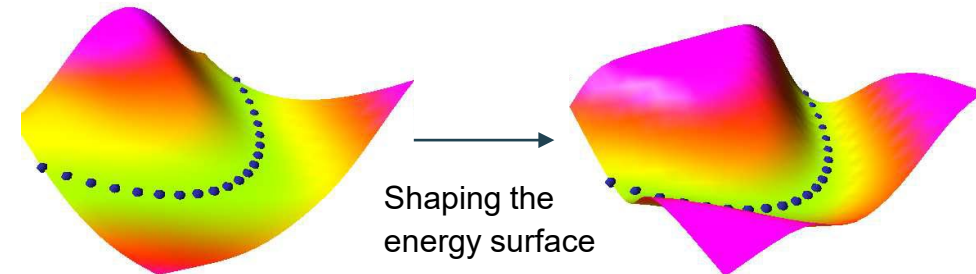$$\hat{y} = \text{argmin}_y \, E_\theta(x, y)$$



Finding statistics with low energies

## B) Training procedure

During training, we search for an energy function $E_\theta(x, y)$ from the space of functions $\mathcal{E}$ that yields the optimal y for any x.

This typically involves searching for the best set of parameters $\theta^*$ for a particular function $E$ using a loss functional $\mathcal{L}$
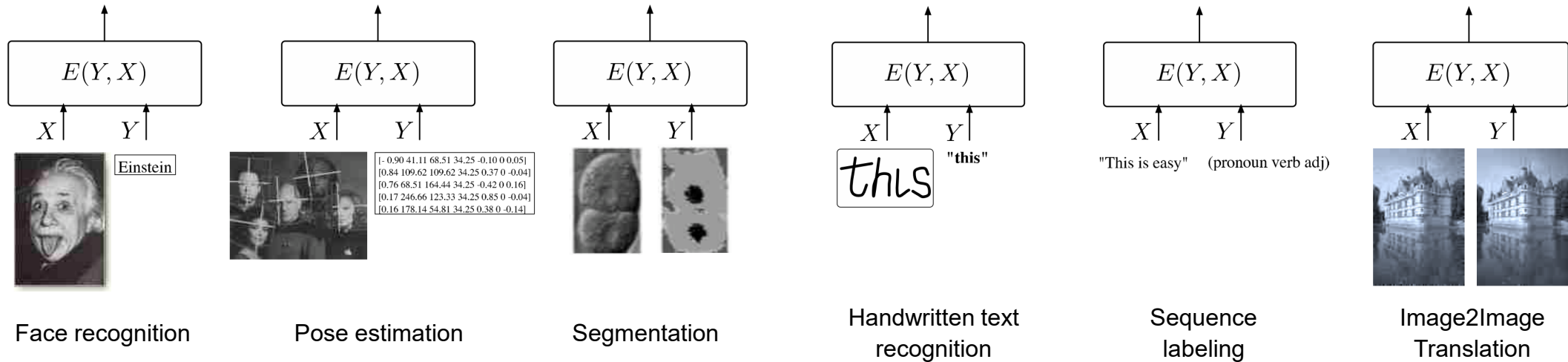
$$\theta^* = \text{argmin}_\theta \, \mathcal{L}(\theta, S)$$

, where $S$ is the set of training statistics $S = \{x_i, y_i\}$.



Shaping the energy surface

[1] Yann LeCun et al. A Tutorial on Energy-Based Learning. Predicting structured data. 2006

# Energy-based models

A) Inference strategies



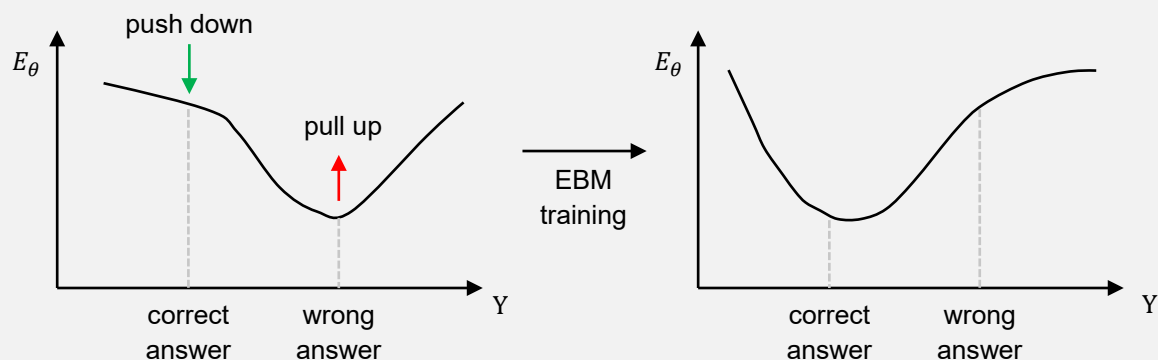| Face recognition | Pose estimation | Segmentation | Handwritten text recognition | Sequence labeling | Image2Image Translation |

Inference can be quite difficult for high-cardinality solution spaces $\mathcal{Y}$ (exhaustive search over all states is not feasible).

Suitable inference strategies include:

- gradient-decent (if $y$ is continuous and the energy function is smooth and differentiable)
- belief propagation
- min-sum (factor graphs)
- dynamic programming (Viterbi, A$*$).
- variational approaches with surrogate energies
- …

[1] Yann LeCun et al. A Tutorial on Energy-Based Learning. Predicting structured data. 2006

# Energy-based models

B) Training strategies

## Intuition on the training procedure (per-sample basis)



→ **Push down** the energy of the correct answer

→ **Pull up** the energies of the incorrect answers

## Standard loss functions (amongst many)

### Energy loss

$$\mathcal{L}_{\text{energy}}(E_\theta(\text{x}, \mathcal{Y}), \text{y}) = E_\theta(\text{x}, \text{y})$$

### Perceptron loss

$$\mathcal{L}_{\text{perceptron}}(E_\theta(\text{x}, \mathcal{Y}), \text{y}) = E_\theta(\text{x}, \text{y}) \min_{\text{y}' \in \mathcal{Y}} E(\text{x}, \text{y}')$$

### Hinge loss

$$\mathcal{L}_{\text{hinge}}(E_\theta(\text{x}, \mathcal{Y}), \text{y}) = \max(0, m + E_\theta(\text{x}, \text{y}) - E_\theta(\text{x}, \bar{\text{y}})$$

positive margin     most offending incorrect answer

→ See [1] for a broad overview of standard loss functions and their pros and cons

[1] Yann LeCun et al. A Tutorial on Energy-Based Learning. Predicting structured data. 2006
[2] Stefano Ermon, Yang Song. Energy-Based Models. Lecture Stanford University

# From energies to probability

## Limitations

The energies are uncalibrated (=un-normalized log-probability)!

What does that mean?

- Not a particular problem for low-cardinality decision making scenarios (just take the solution with the lowest energy)

- But: Two EBMs trained independently most likely have different energies scales. This renders model combination or comparison almost impossible.

- But: Without normalizing the energies, the likelihood of the observed statistics can change with different energy scales

→ We need to find a consistent way to embed EBMs in a framework of common energy units.

## Straight-forward solution?*

We embed the energies in a probability distribution $p(\mathrm{x})$ with properties:

1. non-negative variables: $p(\mathrm{x}) \geq 0$ ← easy to achieve

2. integrates to 1: $\int p(\mathrm{x})\, d\mathrm{x} = 1$ ← essential, but hard for non-trivial problems

*We continue with a simpler notation of the energy function $E_\theta(\mathrm{x})$ that models distribution membership. The conditional $p(\mathrm{y}|\mathrm{x})$ can thus be simplified to $p(\mathrm{x})$

[1] Yann LeCun et al. A Tutorial on Energy-Based Learning. Predicting structured data. 2006
[2] Stefano Ermon, Yang Song. Energy-Based Models. Lecture Stanford University

# From energies to probability

## Boltzmann-Gibbs distribution

$$p(\mathrm{x}) = \frac{1}{Z} \exp\left(-\frac{E(\mathrm{x})}{T}\right)$$

$$Z = \int \exp\left(-\frac{E(\mathrm{x})}{T}\right) d\mathrm{x}$$

- $x$: system state

- $E(x)$: system energy at state x

- $T$: system temperature

- $Z$: normalizing constant/partition function

## Parameterized EBM

$$q_\theta(\mathrm{x}) = \frac{1}{Z_\theta} \exp(-E_\theta(\mathrm{x}))$$

$$\boxed{Z_\theta = \int \exp(-E_\theta(\mathrm{x}))\, d\mathrm{x}}$$

- $x$: image, text, etc.

- $E_\theta(x)/T$: energy function (parameterized by a neural network)

**What can we do to compute the partition function?**

[1] Jianwen Xie, Ying Nian Wu. Theory and Applications of Energy-Based Generative Models. ICCV, 2021.

# From energies to probability

**Approach 1**: Is the task based on pair-wise comparison (e.g. denoising)?

For two data points $x, x'$, calculating $q_\theta(x)$ and $q_\theta(x')$ requires us to know $Z_\theta$. If we can use their ratio, however, we can avoid computating $Z_\theta$ entirely: $\frac{q_\theta(x)}{q_\theta(x')} = \exp(-E_\theta(x) + E_\theta(x'))$.

**Approach 2**: Choose energy function $E_\theta(x)$ such that we can compute $Z_\theta(x)$ analytically.

This approach is viable for trivial choices of the energy function. However, we can also build more complex methods by using products of individually normalized functions (autoregressive, product of experts) or mixtures of normalized objects (latent variables).

**Approach 3**: Approximate the partition function $Z_\theta(x)$ using a Monte Carlo (MC) estimate.

[1] Jianwen Xie, Ying Nian Wu. Theory and Applications of Energy-Based Generative Models. ICCV, 2021.
[2] Stefano Ermon, Yang Song. Energy-Based Models. Lecture Stanford University

# Maximum Likelihood Estimation for EBMs

MCMC for model synthesis

## Prerequisites

$n$ observed data points:

$$\{x_1, \dots, x_n\} \sim p_{\text{train}}(x)$$

Model:

$$q_\theta(x) = \frac{1}{Z_\theta} \exp(-E_\theta(x))$$

MLE objective:

$$\mathcal{L}(\theta; p) = \frac{1}{n} \sum_{i=1}^{n} \log q_\theta(x_i) \doteq \mathbb{E}_{p(x)}[\log q_\theta(x)]$$

### Derivative of negative log-likelihood (NLL)

$$\frac{\partial \mathcal{L}(\theta; p)}{\partial \theta} = \mathbb{E}_{p(x)}\left[\frac{\partial E_\theta(x)}{\partial \theta}\right] - \frac{\partial \log Z_\theta}{\partial \theta}$$

$$\overset{*}{=} \mathbb{E}_{p(x)}\left[\frac{\partial E_\theta(x)}{\partial \theta}\right] - \mathbb{E}_{q_\theta(x)}\left[\frac{\partial E_\theta(x)}{\partial \theta}\right]$$

$$\approx \frac{1}{n} \sum_{i=1}^{n} \frac{\partial E_\theta(x_i)}{\partial \theta} - \frac{1}{\tilde{n}} \sum_{i=1}^{\tilde{n}} \frac{\partial E_\theta(\tilde{x}_i)}{\partial \theta}$$

This is analytically intractable:

Consider grayscale image $x \in \mathbb{R}^{128 \times 128}$ (uint)

Solution space is $256^{16,384}$!

Approximate by Markov Chain Monte Carlo (MCMC)
$$\{\tilde{x}_1, \dots, \tilde{x}_n\} \sim q_\theta(x)$$

* [1] Oliver Woodford. Notes on Contrastive Divergence. Department of Engineering Science, University of Oxford, Tech. Rep 4 (2006).

# Maximum Likelihood Estimation for EBMs

## Contrastive Optimization

Recall the gradient of the NLL objective:

**Contrastive Approximation**

$$\frac{\partial \mathcal{L}(\theta; p)}{\partial \theta} \approx \frac{1}{n}\sum_{i=1}^{n}\frac{\partial E_\theta(\mathrm{x}_i)}{\partial \theta} - \frac{1}{\tilde{n}}\sum_{i=1}^{\tilde{n}}\frac{\partial E_\theta(\tilde{\mathrm{x}}_i)}{\partial \theta}$$
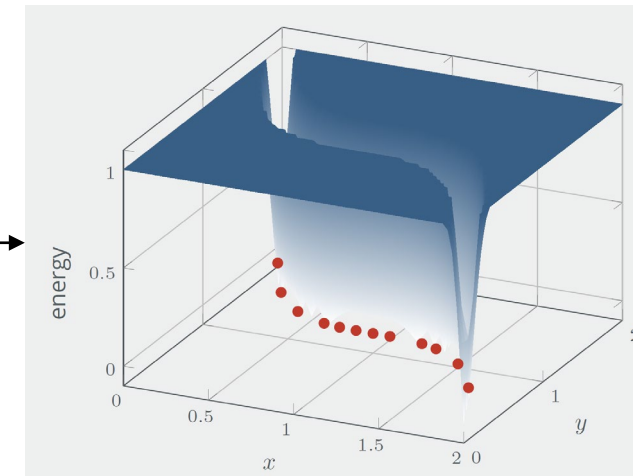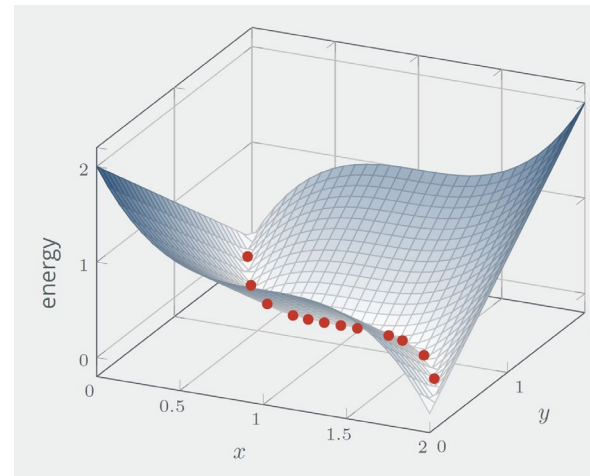
Minimize energies for observed statistics

Maximize energies for synthesized statistics from our current model

Optimum at equilibrium $q_\theta(\mathrm{x}) \simeq p(\mathrm{x})$

"The data samples from the training distribution should be more likely than a sample from our model."

**Caution**: The objective wants to strongly maximize the energy off the manifold, leading to non-smooth energies. → Consider loss regularization!



[1] Chris G. Willcocks. Deep Learning Lecture 7: Energy-based models. https://cwkx.github.io/data/teaching/dl-and-rl/dl-lecture7.pdf
[2] Yann LeCun et al. A Tutorial on Energy-Based Learning. Predicting structured data. 2006

# Gradient-based MCMC w. Langevin Dynamics

Synthesizing samples from $q_\theta(\mathrm{x})$ via vanilla MCMC works in theory but shows extremely long mixing times.

→ For distributions $q_\theta(\mathrm{x})$ that are continuous, we can exploit a score function $\frac{\partial \log q_\theta(\mathrm{x})}{\partial \mathrm{x}}$ for gradient-based synthesis.

→ Widely used approaches are, e.g., Langevin Dynamics (SGLD) or Hamiltonian Monte Carlo (HMC).

## MCMC with Stochastic Gradient Langevin Dynamics

$$\mathrm{x}^0 \sim \pi(\mathrm{x})$$

Stochastic gradient ascent*        Brownian motion

$$\mathrm{x}^{k+1} = \mathrm{x}^k \boxed{+ \frac{\eta}{2}\frac{\partial \log q_\theta(\mathrm{x}^k)}{\partial \mathrm{x}^k}} + \epsilon, \qquad \epsilon \sim \mathcal{N}(0, \eta)$$

$$= \mathrm{x}^k - \frac{\eta}{2}\frac{\partial E_\theta(\mathrm{x}^k)}{\partial \mathrm{x}^k} - \underbrace{\frac{\partial \log Z_\theta}{\partial \mathrm{x}^k}}_{=0} + \epsilon$$

$$= \mathrm{x}^k - \frac{\eta}{2}\frac{\partial E_\theta(\mathrm{x}^k)}{\partial \mathrm{x}^k} + \epsilon$$

*Recall that $f_\theta(\mathrm{x}) = -E_\theta(\mathrm{x})$ (per convention) denotes the negative energy function.

First order Euler discretion of a stochastic differential equation [2]

→ Gradient term (approximated on mini-batch) enforces dynamics to focus on regions of high probability

→ Brownian motion imposes noisy trajectory so that the dynamics explore the complete parameter space

[1] Jianwen Xie, Ying Nian Wu. Theory and Applications of Energy-Based Generative Models. ICCV, 2021.
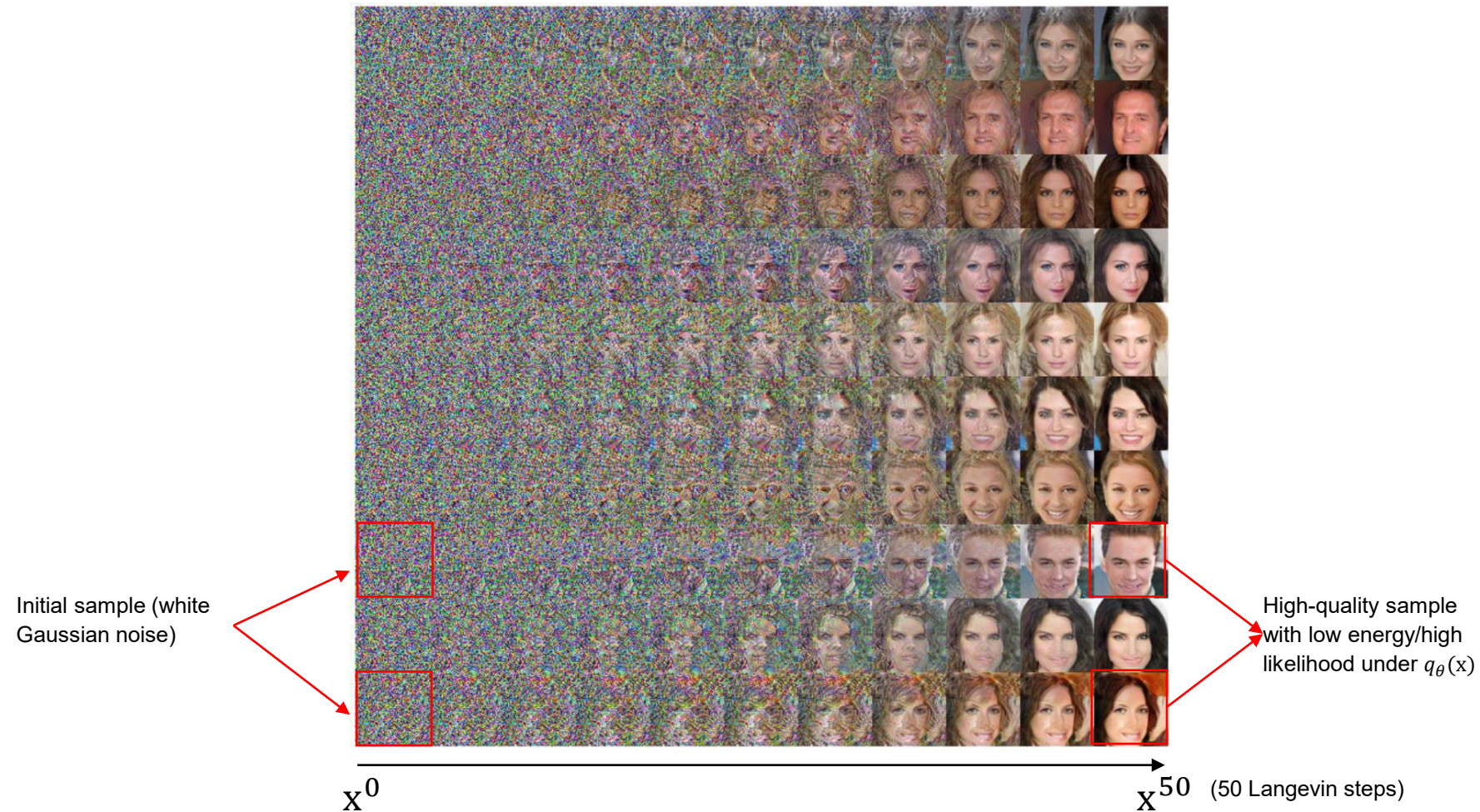[2] Max Welling, Yee Whye The. Bayesian Learning via Stochastic Gradient Langevin Dynamics. ICML 2011. https://www.stats.ox.ac.uk/~teh/research/compstats/sgld.pdf

# Gradient-based MCMC w. Langevin Dynamics

There exist different approaches for selecting the starting point of SGLD-based synthesis:

1. Contrastive Divergence: $x^0 \sim p_{\text{train}}(x)$ → Run finite MCMC from observed samples from the training dataset.

2. Persistent Chain: $x^0 \sim q_\theta^{\text{current epoch}-1}(x)$ → Run finite MCMC from synthesized examples from previous epoch.

3. Non-persistent Short-run MCMC [2]: $x^0 \sim \pi(x)$ → Run finite MCMC from Gaussian noise (see previous slide).

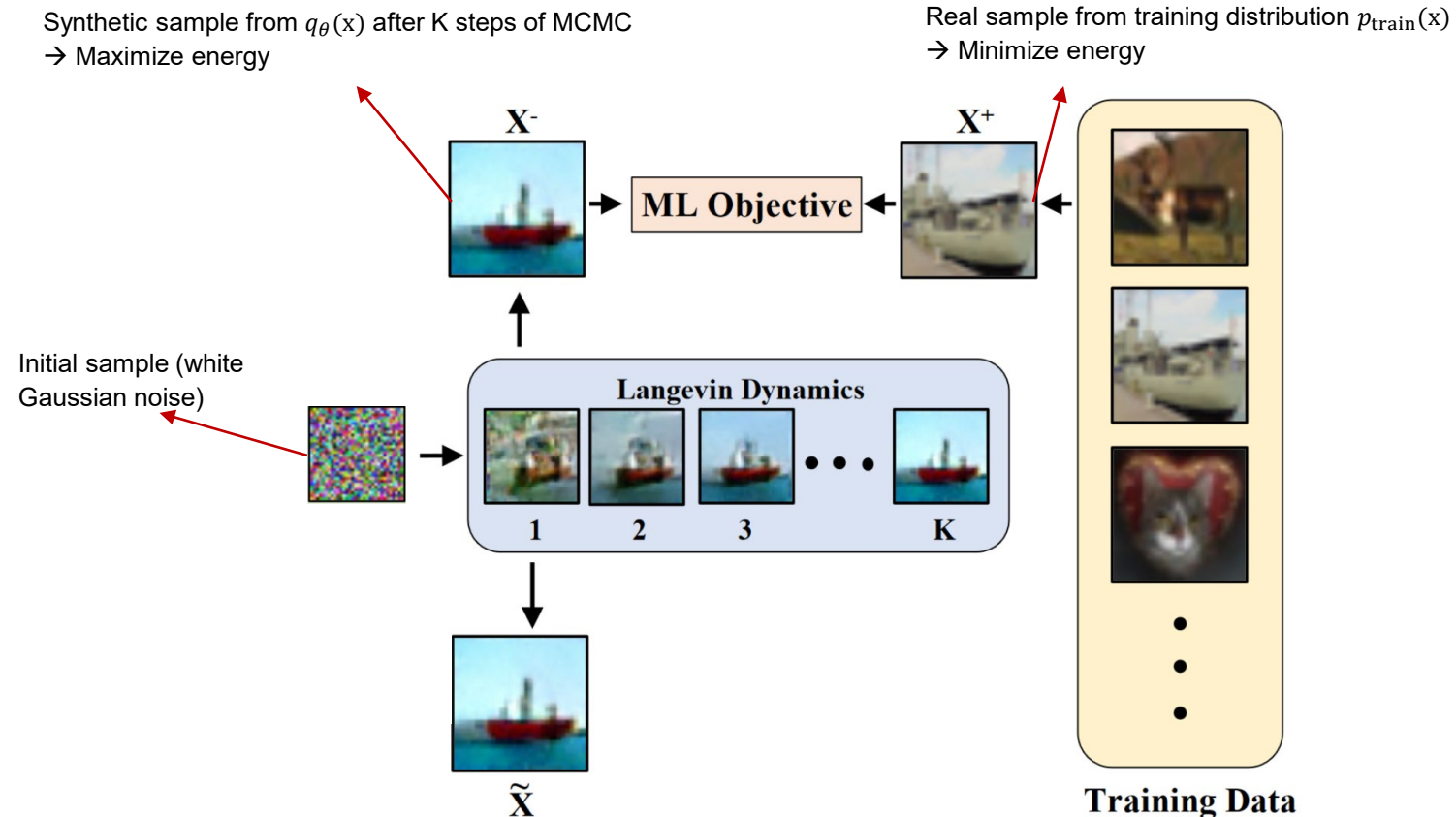[1] Jianwen Xie, Ying Nian Wu. Theory and Applications of Energy-Based Generative Models. ICCV, 2021.
[2] Erik Nijkamp, Mitch Hill, Song-Chun Zhu, Ying NianWu. On learning non-convergent non-persistent short-run MCMC toward energy-based model. NeurIPS, 2019.

# Gradient-based MCMC w. Langevin Dynamics

Initial sample (white Gaussian noise)

High-quality sample with low energy/high likelihood under $q_\theta(x)$

$x^0$      $x^{50}$   (50 Langevin steps)

[1] Yang Zhao, Jianwen Xie, Ping Li. Learning Energy-Based Generative Models via Coarse-to-Fine Expanding and Sampling. ICLR, 2021.

# Training and Sampling Algorithm

## Overview

Synthetic sample from $q_\theta(x)$ after K steps of MCMC
$\rightarrow$ Maximize energy

Real sample from training distribution $p_{\text{train}}(x)$
$\rightarrow$ Minimize energy

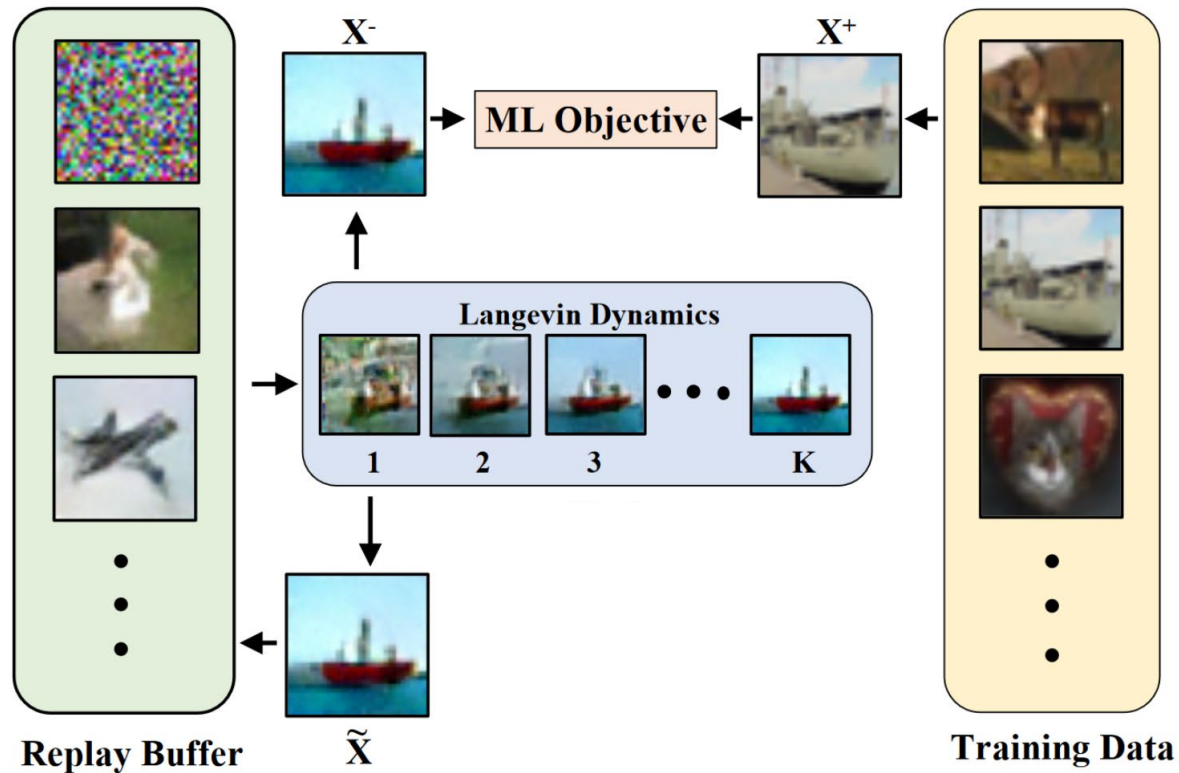

Initial sample (white Gaussian noise)

We can use the trained EBM to synthesize new (and quite realistic) samples. (Interpretation as a latent generative model).

[1] Yilun Du, Igor Mordatch. Implicit Generation and Generalization in Energy-Based Models. NeurIPS, 2019.

# Training and Sampling Algorithm

## General approach

**Replay Buffer**   $\tilde{\mathbf{X}}$   **Training Data**

Leads to mixture of a) non-persistent short-run MCMC and b) persistent MCMC

### Standard training and sampling algorithm for EBMs

**Input** data distribution $p(x)$, MCMC mixing steps $K$

1: Initialize synthesis buffer $\mathcal{B} \leftarrow \emptyset$

2: **while** not converged **do:**

3:  sample from dataset: $x_i^+ \sim p_{\text{train}}$

3:  initial synthesis: $x_i^0 \sim \mathcal{B}$ with 95% probability, $x_i^0 \sim \mathcal{U}$ or $x_i^0 \sim \mathcal{N}$ otherwise

4:  **for** sample step $k = 1$ to K **do:**

5:
$$\tilde{x}^k = \tilde{x}^{k-1} - \frac{\eta}{2}\frac{\partial E_\theta(\tilde{x}^{k-1})}{\partial \tilde{x}^{k-1}} + \epsilon, \;\; \epsilon \sim \mathcal{N}(0, \eta)$$

6:  **end for**

7:  $x_i^- \leftarrow \Omega(\tilde{x}_i^k)$

8:  Contrastive Divergence: $\mathcal{L}_{\text{CD}} = \frac{1}{N}\sum_i E_\theta(x_i^+) - E_\theta(x_i^-)$

9:  L2 regularization: $\mathcal{L}_{\text{L2}} = \frac{1}{N}\sum_i E_\theta(x_i^+)^2 + E_\theta(x_i^-)^2$

10:  Optimization with SGD/Adam/…: $\frac{\partial}{\partial \theta}(\mathcal{L}_{\text{CD}} + \lambda\mathcal{L}_{\text{L2}})$

11:  add to buffer: $\mathcal{B} \leftarrow \mathcal{B} \cup x_i^-$

12: **end while**

→ Alternating between sampling $q_\theta(\mathrm{x})$ and updating parameters $\theta$ of $q_\theta(\mathrm{x})$

[1] Yilun Du, Igor Mordatch. Implicit Generation and Generalization in Energy-Based Models. NeurIPS, 2019.
[2] Phillip Lippe. Tutorial 8: Deep Energy-Based Generative Models. https://uvadlc-notebooks.readthedocs.io/en/latest/tutorial_notebooks/tutorial8/Deep_Energy_Models.html#Training-algorithm

# Adversarial Interpretation of Training and Sampling

Recall the contrastive algorithm where we synthesize data points using MCMC.

$$\frac{\partial \mathcal{L}(\theta; p)}{\partial \theta} \approx \frac{1}{n} \sum_{i=1}^{n} \frac{\partial E_\theta(x_i)}{\partial \theta} - \frac{1}{\tilde{n}} \sum_{i=1}^{n} \frac{\partial E_\theta(\tilde{x}_i)}{\partial \theta}$$

$$= \frac{\partial}{\partial \theta} \left[ \frac{1}{n} \sum_{i=1}^{n} E_\theta(x_i) - \frac{1}{\tilde{n}} \sum_{i=1}^{n} E_\theta(\tilde{x}_i) \right]$$

→ Define value function $V(\{\tilde{x}_i\}, \theta) = \frac{1}{n} \sum_{i=1}^{n} E_\theta(x_i) - \frac{1}{\tilde{n}} \sum_{i=1}^{n} E_\theta(\tilde{x}_i)$

→ The alternating i) Learning and ii) Sampling procedures play an adversarial minmax game:

$$\min_{\{\tilde{x}_i\}} \max_{\theta} V(\{\tilde{x}_i\}, \theta)$$

[1] Jianwen Xie, Ying Nian Wu. Theory and Applications of Energy-Based Generative Models. ICCV, 2021.

# Non-persistent Short-Run MCMC

**Problem** If $q_\theta(\mathrm{x})$ is multi-modal, different MC chains might get trapped in distinct local modes → no mixing.

**Approach** Nijkamp et al. [2] propose to run a non-convergent, non-mixing, and non-persistent chain that starts from always the same initial noise distribution $q_0$ (uniform, Gaussian) and run it for a small fixed number of steps $K$ towards $q_\theta$. → Short-run MCMC

## Short-run MCMC

$M_\theta$: $\quad$ $K$-step MCMC transition kernel

$z \sim q_0$: $\quad$ starting point (i.e., latent variables)

$$\gamma_\theta(\mathrm{x}) = (M_\theta q_0)(\mathrm{z}) = \int q_0(\mathrm{z}) M_\theta(\mathrm{x}|\mathrm{z}) dz$$

Marginal distribution of sample x after K-step MCMC from $q_0$

$$\mathrm{x} = M_\theta(\mathrm{z}, \epsilon)$$

If training converges, EBMs tends to have low entropy and Langevin dynamics behaves like GD with disabled noise term $\epsilon$ → x = $M_\theta(\mathrm{z})$

→ No longer maximum-likelihood estimation, but moment-matching estimation (MME) with the following relation to solve:

$$\mathbb{E}_{p(\mathrm{x})}\left[\frac{\partial E_\theta(\mathrm{x})}{\partial \theta}\right] = \mathbb{E}_{\gamma_\theta(\mathrm{x})}\left[\frac{\partial E_\theta(\mathrm{x})}{\partial \theta}\right]$$

[1] Jianwen Xie, Ying Nian Wu. Theory and Applications of Energy-Based Generative Models. ICCV, 2021.
[2] Erik Nijkamp, Mitch Hill, Song-Chun Zhu, Ying NianWu. On learning non-convergent non-persistent short-run MCMC toward energy-based model. NeurIPS, 2019.

# Non-persistent Short-Run MCMC

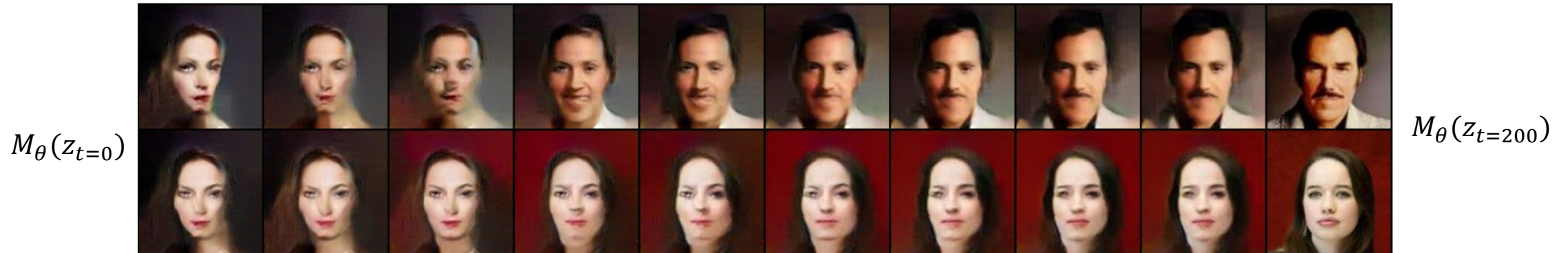**Image synthesis** using Short-run MCMC ($K = 100$, $q_0 \sim \mathcal{U}[-1,1]$) on CelebA (64x64)



**Image synthesis** using Short-run MCMC ($K = 100$, $q_0 \sim \mathcal{U}[-1,1]$) on CelebA (128x128)

[1] Erik Nijkamp, Mitch Hill, Song-Chun Zhu, Ying NianWu. On learning non-convergent non-persistent short-run MCMC toward energy-based model. NeurIPS, 2019.

# Non-persistent Short-Run MCMC

$M_\theta(z_1)$                                                                 $M_\theta(z_2)$

**Interpolation** using Short-run MCMC as generative latent model. Transitions depits $M_\theta(z_p)$ with interpolated noise $z_p = p z_1 + \sqrt{1 - p^2 z_2}$ where $p \in [0,1]$ on CelebA (64x64).
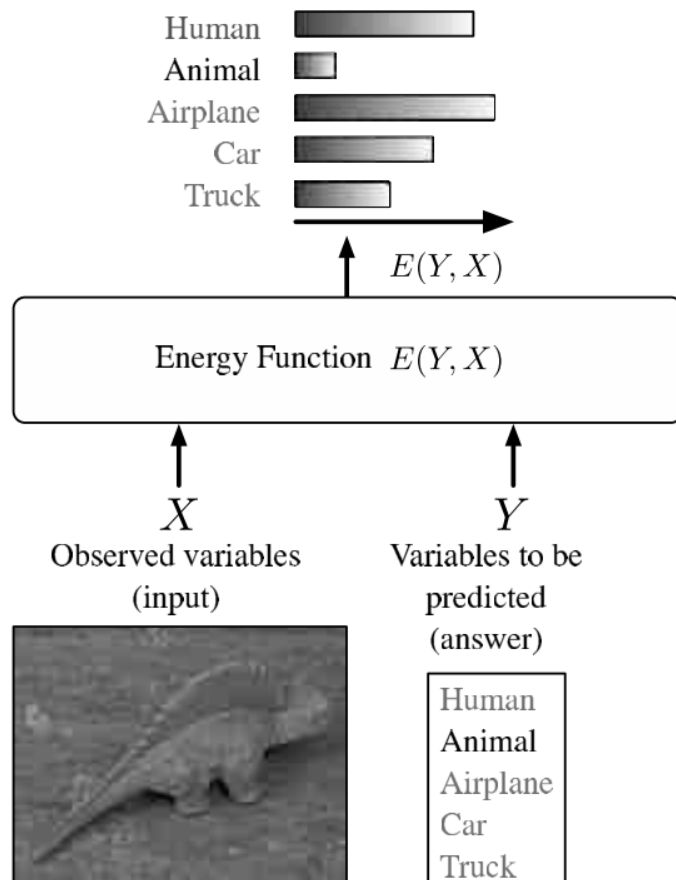


$M_\theta(z_{t=0})$                                                             $M_\theta(z_{t=200})$

**Reconstruction** using Short-run MCMC as generative latent model. Transitions depits $M_\theta(z_t)$ over time $t$ from random noise at $t = 0$ to reconstruction of observed example at $t = 200$ on CelebA (64x64).
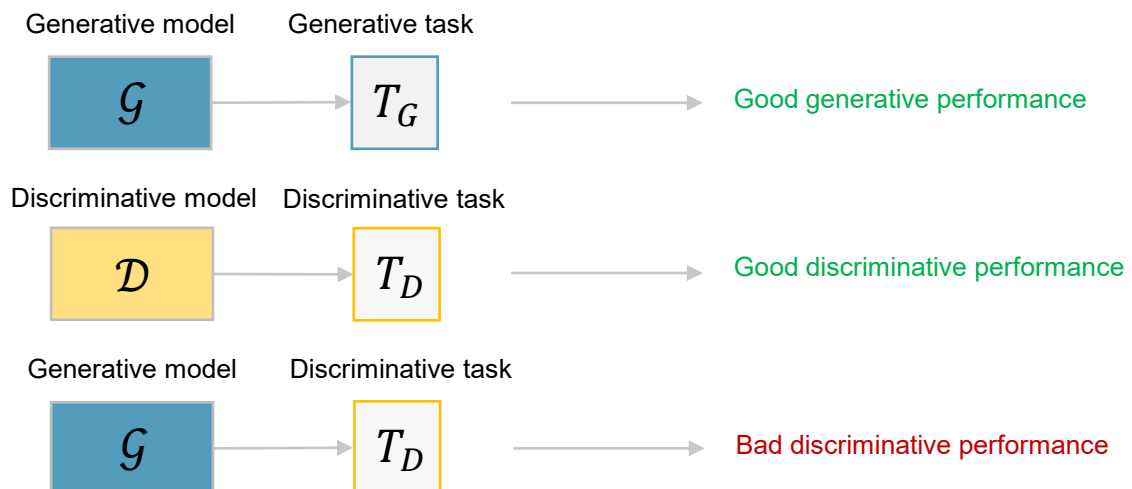
[1] Erik Nijkamp, Mitch Hill, Song-Chun Zhu, Ying NianWu. On learning non-convergent non-persistent short-run MCMC toward energy-based model. NeurIPS, 2019.

# Joint Energy-based Models (JEM)

# Joint Energy-based Models (JEM)

Human
Animal
Airplane
Car
Truck

$E(Y,X)$

Energy Function $E(Y,X)$

$X$
Observed variables
(input)

$Y$
Variables to be
predicted
(answer)

Human
Animal
Airplane
Car
Truck

We can use the joint probability density $p(x,y)$ for classification and recognition tasks: $\hat{y} = \operatorname{argmin}_y E(x,y)$

**General observation**

Discriminative capabilities of generative models are inferior to discriminative models that specialize on classification tasks!



Generative model    Generative task
$\mathcal{G}$    $T_G$    Good generative performance

Discriminative model    Discriminative task
$\mathcal{D}$    $T_D$    Good discriminative performance

Generative model    Discriminative task
$\mathcal{G}$    $T_D$    Bad discriminative performance

[1] Yann LeCun et al. A Tutorial on Energy-Based Learning. Predicting structured data. 2006

# Joint Energy-based Models (JEM)

**Gratwohl et al. made a fascinating observation**

(1) Categorical distribution obtained by classifier:

$$q_\theta(y|\mathbf{x}) = \frac{\exp(\overbrace{f_\theta(\mathbf{x})[y]}^{\text{logits for target class label}})}{\sum_{y'} \exp(f_\theta(\mathbf{x})[y'])}$$

Pseudo-probabilities via softmax over all y

(2) We can re-interpret the logits as unnormalized densities of the joint distribution $q_\theta(\mathrm{x}, y)$!

$$q_\theta(\mathbf{x}, y) = \frac{\exp(\overbrace{f_\theta(\mathbf{x})[y]}^{E_\theta(\mathrm{x}, y) = -f_\theta(\mathrm{x})[y]})}{\underbrace{Z(\theta)}_{\text{Partition function}}}$$

Marginalization over y $\longrightarrow$

$$q_\theta(\mathbf{x}) = \sum_y q_\theta(\mathbf{x}, y) = \frac{\sum_y \exp(f_\theta(\mathbf{x})[y])}{\underbrace{Z(\theta)}_{\text{Density model for x}}}$$

logits for target class label

(3) Energy function at point x:

$$E_\theta(\mathbf{x}) = -\mathrm{LogSumExp}_y(f_\theta(\mathbf{x})[y]) = -\log \sum_y \exp(f_\theta(\mathbf{x})[y])$$

[1] Gratwohl et al. Your classifier is secretely an Energy-based model and you should treat it like one, ICLR 2020

# Joint Energy-based Models (JEM)

(4) Efficient training with factorization

→ Recall that the posterior is given by $p(y|\mathbf{x}) = \frac{p(\mathbf{x},y)}{p(\mathbf{x})}$

→ Factor the likelihood to facilitate training a hybrid model (generative and discriminative aspects)

$$\log q_\theta(\mathbf{x}, y) = \underbrace{\log q_\theta(\mathbf{x})}_{\substack{\text{optimize density model} \\ \text{with NLL + SGLD}}} + \underbrace{\log q_\theta(y|\mathbf{x})}_{\substack{\text{optimize with standard} \\ \text{cross-entropy}}}$$



[1] Gratwohl et al. Your classifier is secretely an Energy-based model and you should treat it like one, ICLR 2020

# Application Examples

# Application examples

Energy-based inpainting

Can we recover missing information from incomplete training data in an unsupervised fashion?



**Proposed solution:** Learning + synthesis of new examples + recovery of incomplete training samples

→ Combination of two Langevin dynamics:

1.  Start from white noise and synthesize new example $a_i$

2.  Start from incomplete data and synthesize recovered data $b_i$

→ Update rule: $\theta_{t+1} = \theta_t - \eta_t \left[ \frac{1}{n}\sum_{i=1}^{n} \frac{\partial E_\theta(a_i)}{\partial \theta} - \frac{1}{n}\sum_{i=1}^{n} \frac{\partial E_\theta(b_i)}{\partial \theta} \right]$

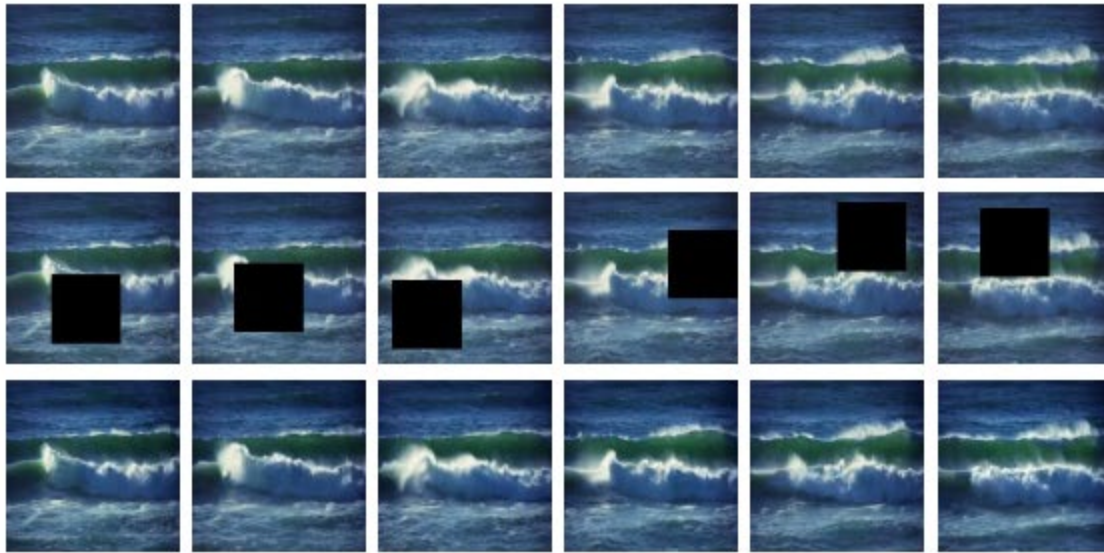[1] Jianwen Xie, Ying Nian Wu. Theory and Applications of Energy-Based Generative Models. ICCV, 2021.
[2] Jianwen Xie, Song-Chun Zhu, Ying Nian Wu. Synthesizing Dynamic Pattern by Spatial-Temporal Generative ConvNet. CVPR 2017.
[3] Jianwen Xie, Song-Chun Zhu, Ying Nian Wu. Learning Energy-based Spatial-Temporal Generative ConvNetfor Dynamic Patterns. PAMI 2019.

windmill

fountain

50% salt and pepper masking

[1] Jianwen Xie, Ying Nian Wu. Theory and Applications of Energy-Based Generative Models. ICCV, 2021.
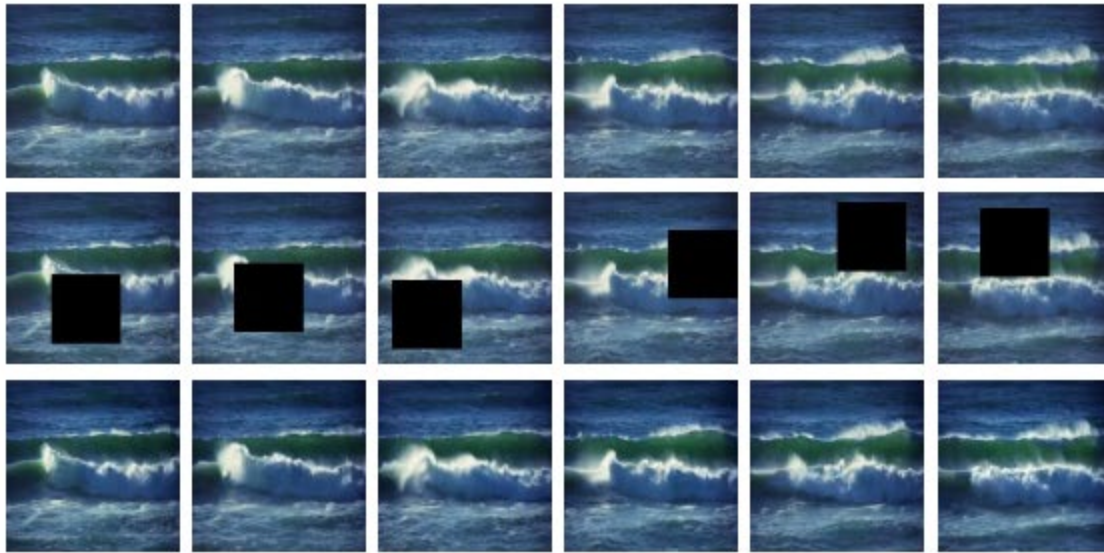[2] Jianwen Xie, Song-Chun Zhu, Ying Nian Wu. Synthesizing Dynamic Pattern by Spatial-Temporal Generative ConvNet. CVPR 2017.
[3] Jianwen Xie, Song-Chun Zhu, Ying Nian Wu. Learning Energy-based Spatial-Temporal Generative ConvNetfor Dynamic Patterns. PAMI 2019.

# Application examples

Energy-based inpainting



ocean

flag

50% salt and pepper masking

[1] Jianwen Xie, Ying Nian Wu. Theory and Applications of Energy-Based Generative Models. ICCV, 2021.
[2] Jianwen Xie, Song-Chun Zhu, Ying Nian Wu. Synthesizing Dynamic Pattern by Spatial-Temporal Generative ConvNet. CVPR 2017.
[3] Jianwen Xie, Song-Chun Zhu, Ying Nian Wu. Learning Energy-based Spatial-Temporal Generative ConvNetfor Dynamic Patterns. PAMI 2019.
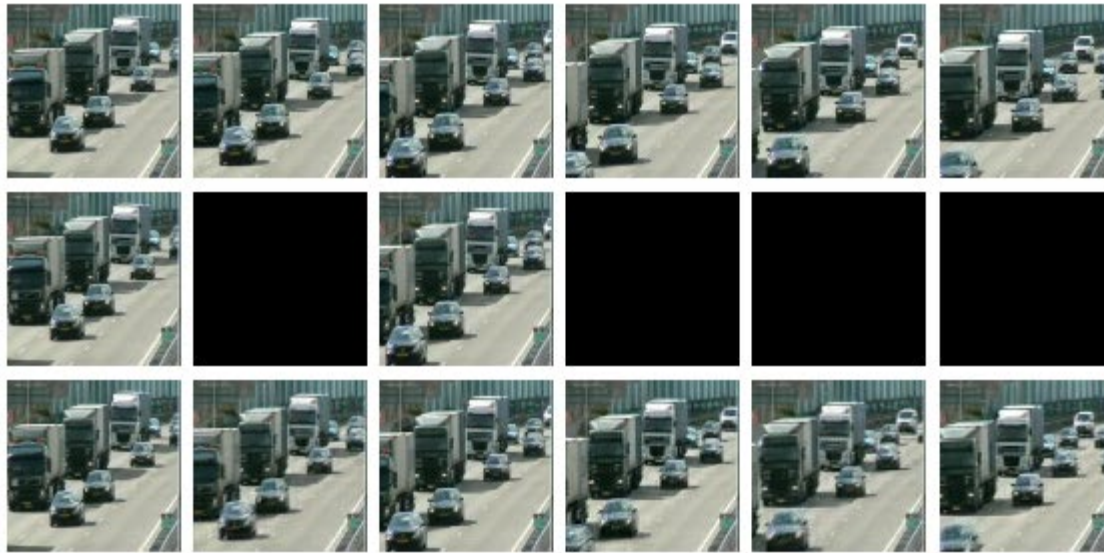
ocean

flag

single region masks

[1] Jianwen Xie, Ying Nian Wu. Theory and Applications of Energy-Based Generative Models. ICCV, 2021.
[2] Jianwen Xie, Song-Chun Zhu, Ying Nian Wu. Synthesizing Dynamic Pattern by Spatial-Temporal Generative ConvNet. CVPR 2017.
[3] Jianwen Xie, Song-Chun Zhu, Ying Nian Wu. Learning Energy-based Spatial-Temporal Generative ConvNetfor Dynamic Patterns. PAMI 2019.

# Application examples

## Energy-based inpainting

traffic

playing

50% missing frames

[1] Jianwen Xie, Ying Nian Wu. Theory and Applications of Energy-Based Generative Models. ICCV, 2021.
[2] Jianwen Xie, Song-Chun Zhu, Ying Nian Wu. Synthesizing Dynamic Pattern by Spatial-Temporal Generative ConvNet. CVPR 2017.
[3] Jianwen Xie, Song-Chun Zhu, Ying Nian Wu. Learning Energy-based Spatial-Temporal Generative ConvNet for Dynamic Patterns. PAMI 2019.

# Application examples

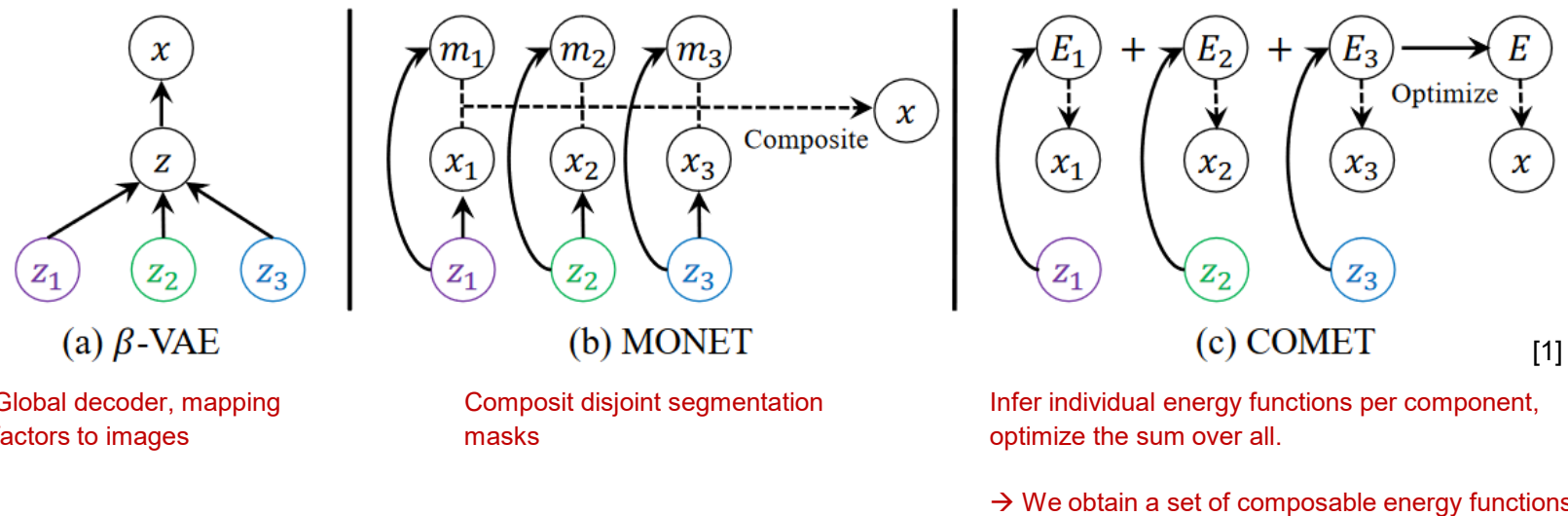Energy-based inpainting



(a) removing a moving boat in the lake

(b) removing a walking person in front of fountain

Background inpainting

[1] Jianwen Xie, Ying Nian Wu. Theory and Applications of Energy-Based Generative Models. ICCV, 2021.
[2] Jianwen Xie, Song-Chun Zhu, Ying Nian Wu. Synthesizing Dynamic Pattern by Spatial-Temporal Generative ConvNet. CVPR 2017.
[3] Jianwen Xie, Song-Chun Zhu, Ying Nian Wu. Learning Energy-based Spatial-Temporal Generative ConvNet for Dynamic Patterns. PAMI 2019.

# Application examples

Unsupervised Learning of Compositional Energy Concepts

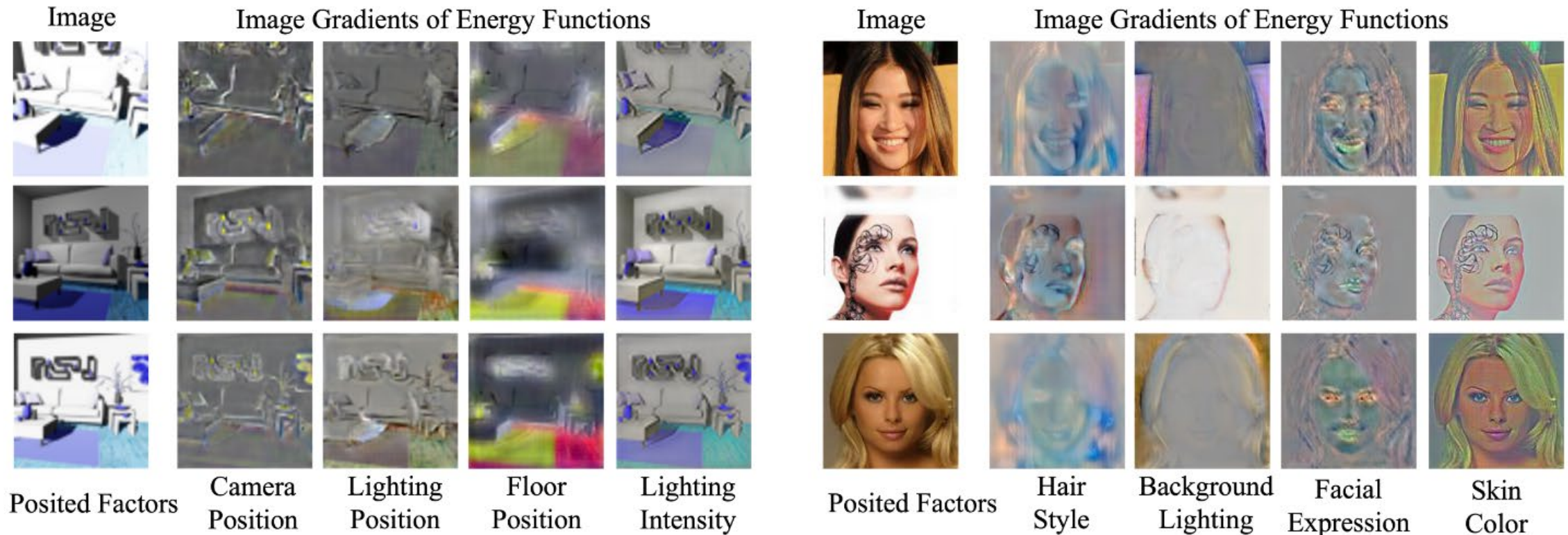Can we identify and discover visual concepts (local or global factor of variation) using EBMs?

**Proposed solution:** COMET[1] infers multiple energy functions individually with separate minimal energy states, each capturing a distinct factor of variation.



(a) $\beta$-VAE       (b) MONET       (c) COMET       [1]

Global decoder, mapping factors to images

Composit disjoint segmentation masks

Infer individual energy functions per component, optimize the sum over all.

→ We obtain a set of composable energy functions!

[1] Yilun Du, Shuang Li, Yash Sharma, Joshua B. Tenenbaum, Igor Mordatch. Unsupervised Learning of Compositional Energy Concepts. NeurIPS, 2021.

# Application examples
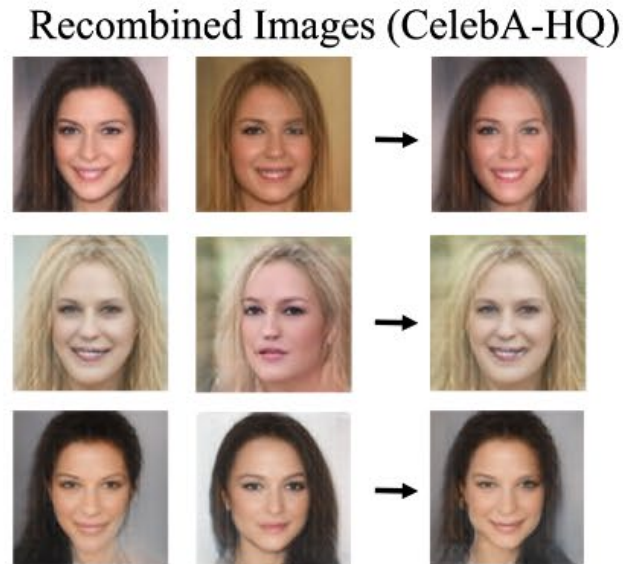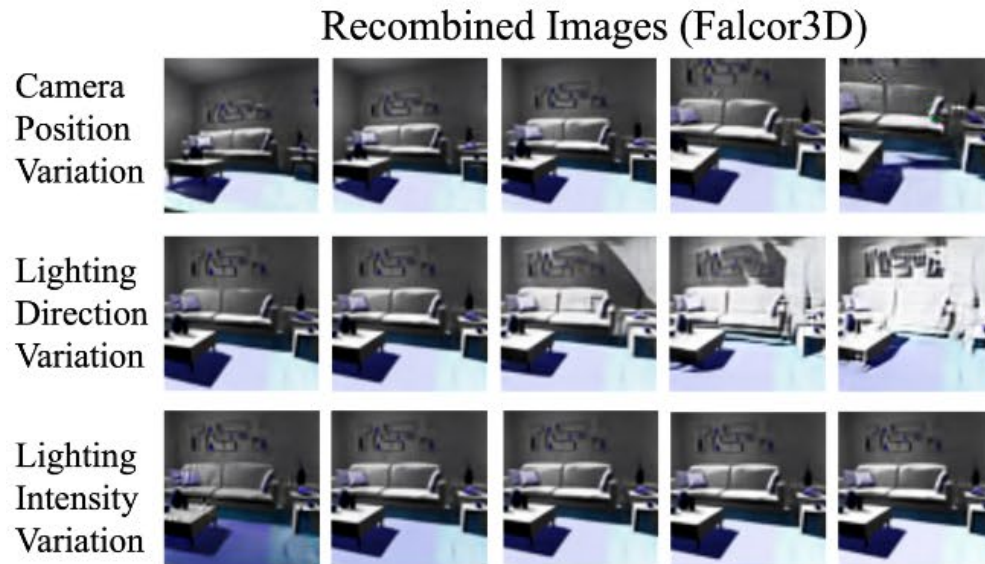
Unsupervised Learning of Compositional Energy Concepts



Gradients of the individual energy functions (k=4) wr.t. an image. The gradient images indicate those aspects of an image the respective energy function attends to. (Factor labels are assigned by visual inspection.)

[1] Yilun Du, Shuang Li, Yash Sharma, Joshua B. Tenenbaum, Igor Mordatch. Unsupervised Learning of Compositional Energy Concepts. NeurIPS, 2021.

## Unsupervised Learning of Compositional Energy Concepts



Recombination of individual energy functions.

[1] Yilun Du, Shuang Li, Yash Sharma, Joshua B. Tenenbaum, Igor Mordatch. Unsupervised Learning of Compositional Energy Concepts. NeurIPS, 2021.

## Unsupervised Learning of Compositional Energy Concepts



Decompose and recombine energy functions representing both local and global factors

[1] Yilun Du, Shuang Li, Yash Sharma, Joshua B. Tenenbaum, Igor Mordatch. Unsupervised Learning of Compositional Energy Concepts. NeurIPS, 2021.

# Application examples
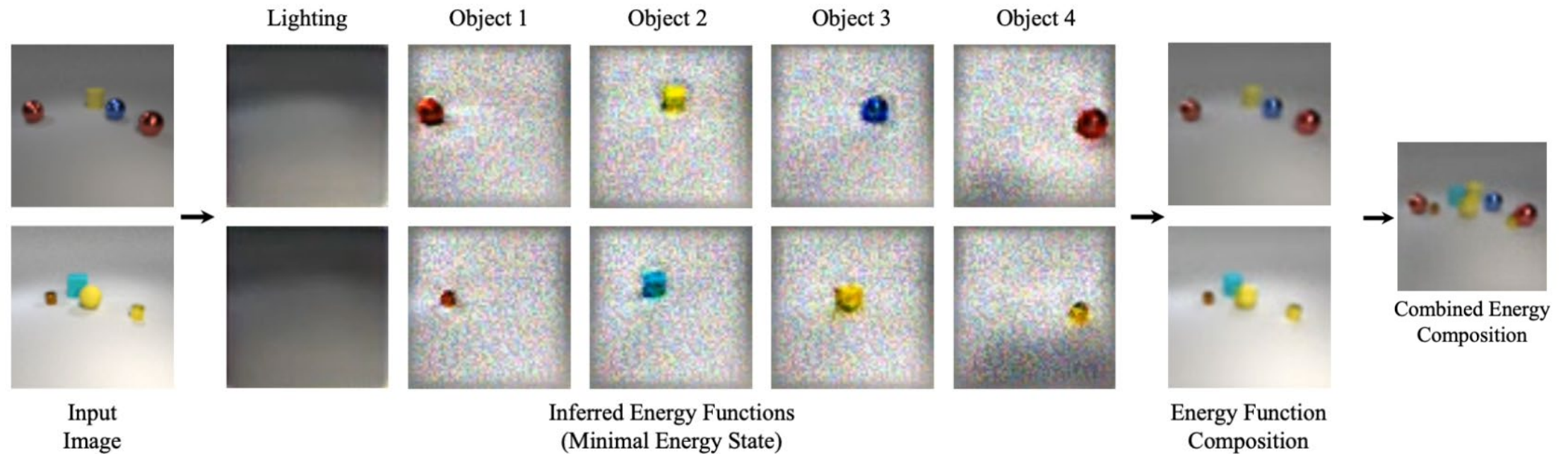
Unsupervised Learning of Compositional Energy Concepts



Decompose and recombine energy functions representing both local and global factors

[1] Yilun Du, Shuang Li, Yash Sharma, Joshua B. Tenenbaum, Igor Mordatch. Unsupervised Learning of Compositional Energy Concepts. NeurIPS, 2021.

# Conclusion and Take-Away Points

# Conclusion and Take-Away Points

$$q_\theta(\mathrm{x}) = \frac{1}{\int \exp(-E_\theta(\mathrm{x}))\, dx} \exp(-E_\theta(\mathrm{x})) = \frac{1}{Z_\theta} \exp(-E_\theta(\mathrm{x}))$$

→ Energy-based Models (EBMs) are powerful generative methods for capturing characteristics, regularities, and constraints of (high-dim.) data distributions

→ The energy function can be parameterized by almost any kind of model architecture!

→ Composition of energies is mostly straight-forward (product of experts, etc.)

→ There exist lots of fundamental connections to other model families (Diffusion Modeling, GANs, Normalizing Flows, Graphical Models, etc.)

**But**:

→ Likelihood-based learning is not trivial and oftentimes computationally expensive

→ Sampling is hard and training is oftentimes instable