

ELV: Embodied Learning and Vision

Mengye Ren

NYU

Spring 2026

elvcourse.org



Module 1: Deep Learning for Structured Outputs

Deep Learning

- Over decades, optimizing deep neural networks was not trivial.

Deep Learning

- Over decades, optimizing deep neural networks was not trivial.

- Progress came from (taken for granted nowadays):

- Initialization ~ weights

- Optimizers (Adam, gradient conditioning)

- Normalization (BN, LN, GN, etc.)

- Skip connection (recurrent net, residual net, preserving gradients)

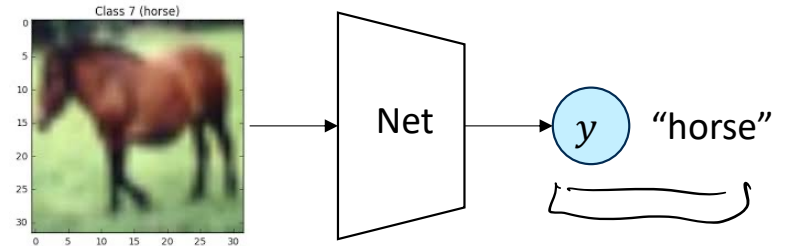
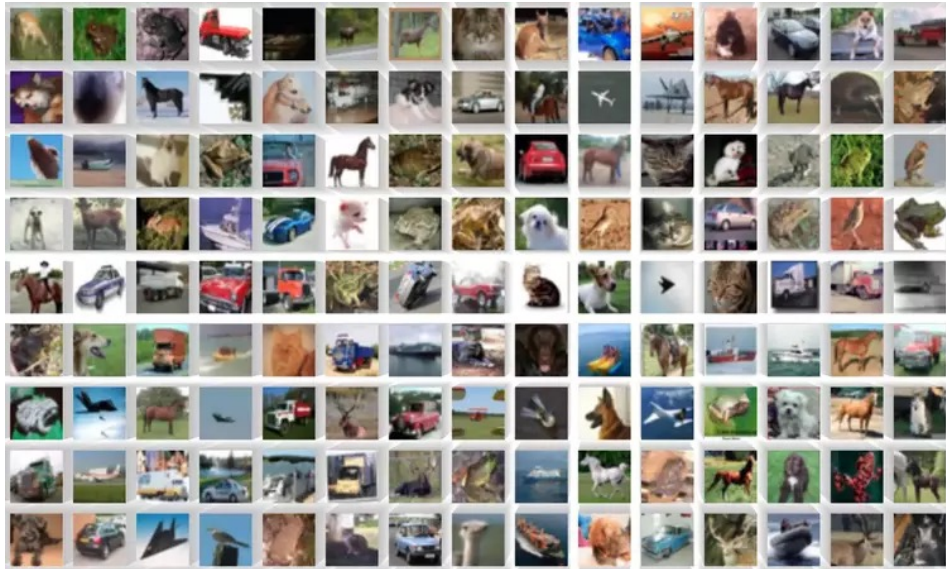
- Regularization (dropout, noise, augmentation)

- Attention (generalization, parallelization)



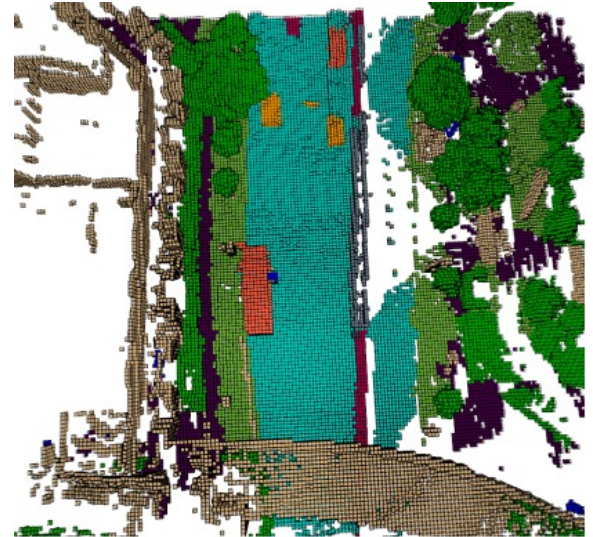
Classification

- To test how we can fit a deep neural network well, people have relied on simple benchmarks, such as image classification.



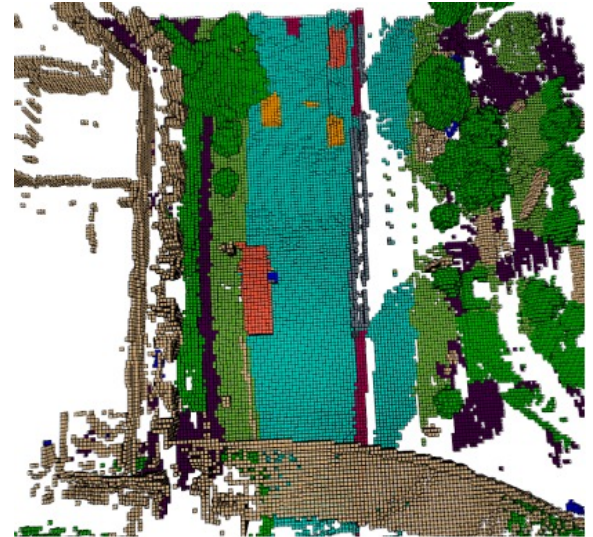
What are Structured Outputs?

- An embodied agent needs to have a structured output space.
 - Object localization, tracking, motion, spatial segmentation, 3d occupancy
 - Trajectory, forecasting, planning



What are Structured Outputs?

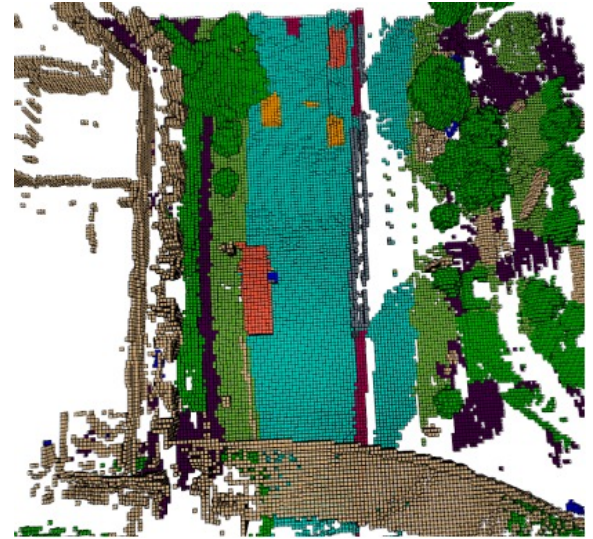
- An embodied agent needs to have a structured output space.
 - Object localization, tracking, motion, spatial segmentation, 3d occupancy
 - Trajectory, forecasting, planning
- A naïve solution is to simply have multiple output dimensions.



What are Structured Outputs?

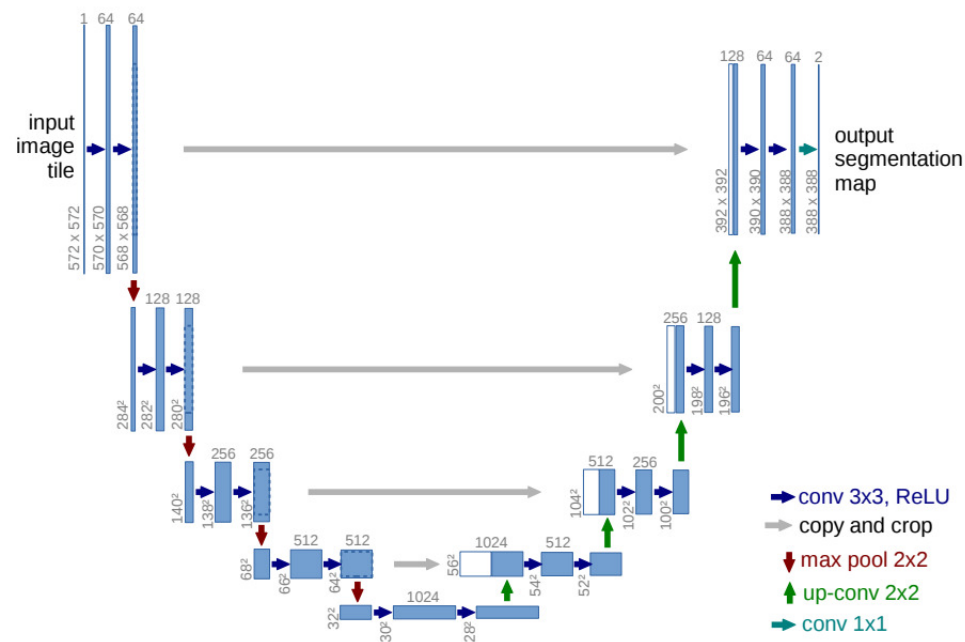
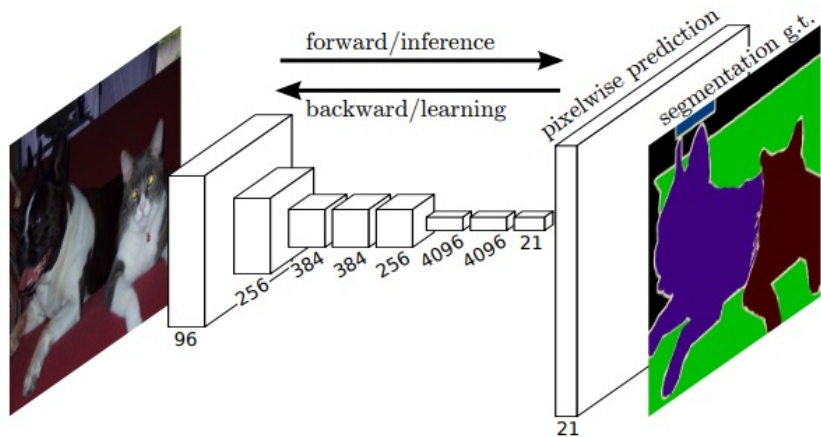
- An embodied agent needs to have a structured output space.
 - Object localization, tracking, motion, spatial segmentation, 3d occupancy
 - Trajectory, forecasting, planning
- A naive solution is to simply have multiple output dimensions.
- It often does not reason the joint probability

$$y = Wx \quad \textcircled{y_1} \textcircled{y_2} \textcircled{y_3} \dots \textcircled{x}$$
$$y_i = w_i^T x$$



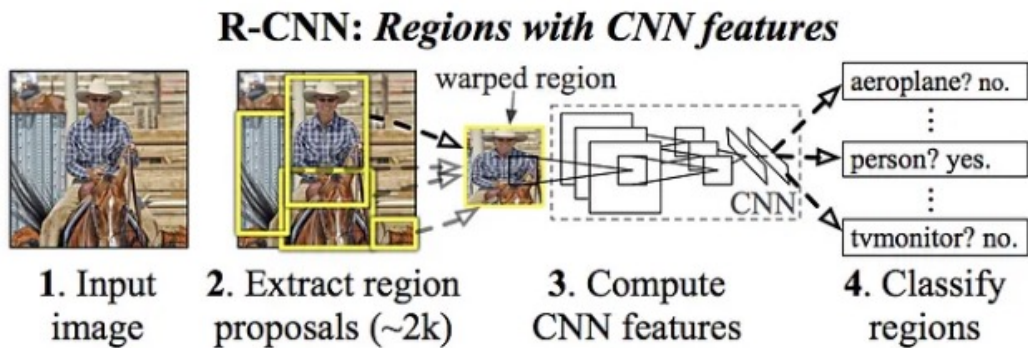
Network Architecture for Structured Outputs

- Segmentation
- Spatial, high-resolution

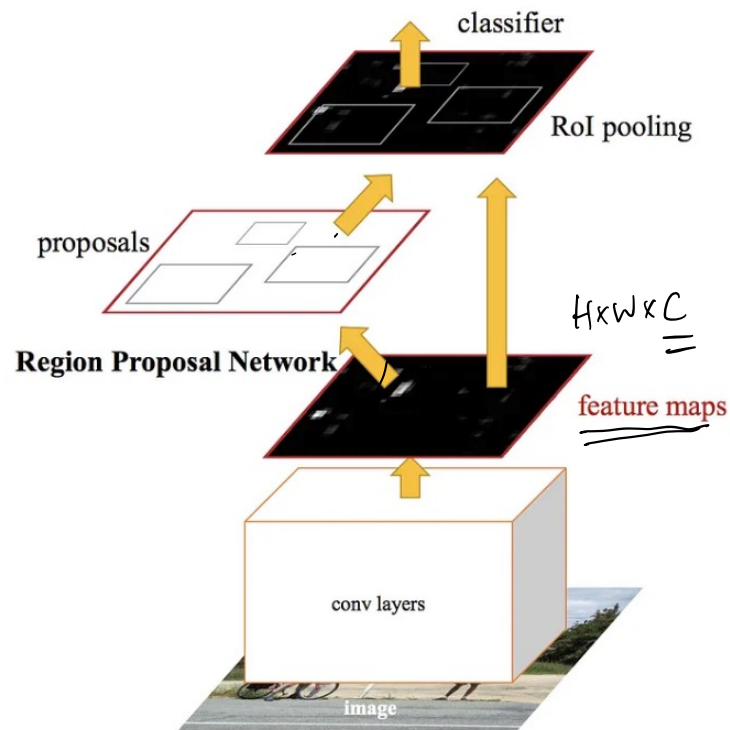


Network Architecture for Structured Outputs

- Object Detection



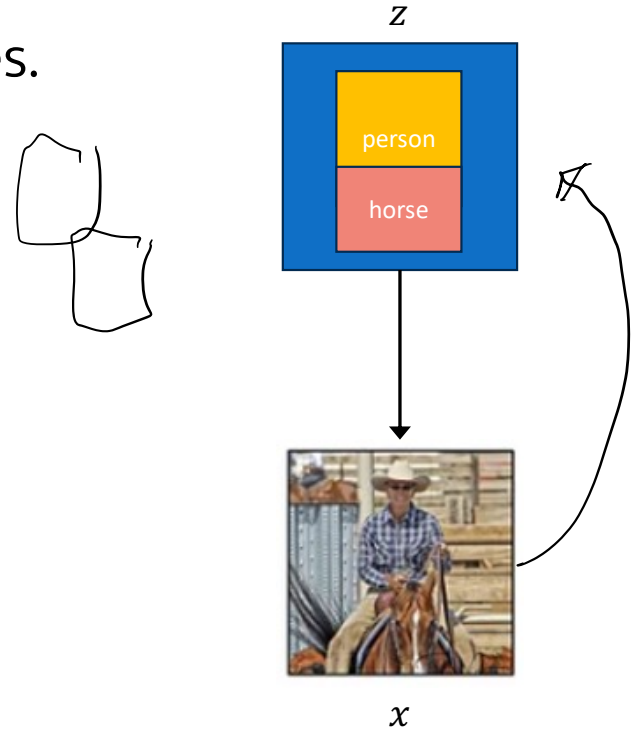
Girshick et al., 2013



Ren et al., 2015

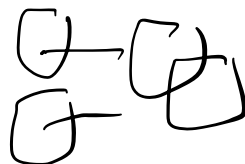
Object Detection as Inference

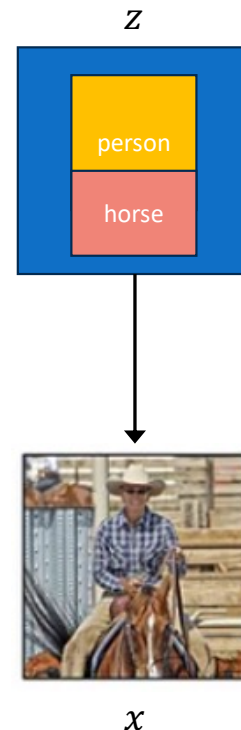
- Bounding boxes are structured latent variables.



Object Detection as Inference

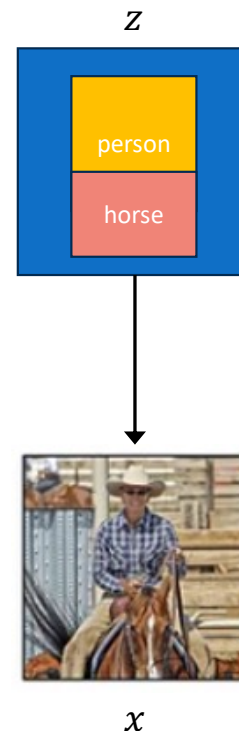
- Bounding boxes are structured latent variables.
- Occupancy as physical constraints.
 - One spatial 3D location can only present a single physical object.

$$y = W^T x.$$




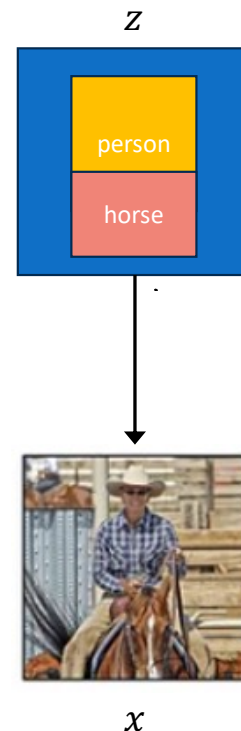
Object Detection as Inference

- Bounding boxes are structured latent variables.
- Occupancy as physical constraints.
 - One spatial 3D location can only present a single physical object.
- Object co-occurrence in the scene
 - Context



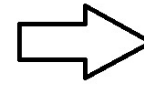
Object Detection as Inference

- Bounding boxes are structured latent variables.
- Occupancy as physical constraints.
 - One spatial 3D location can only present a single physical object.
- Object co-occurrence in the scene
 - Context
- The role of the network is to perform “inference” on the latent variables.



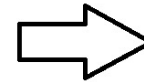
Non-Maximal Suppression

- Box Proposals: Samples of boxes that may come from a single object (latent)



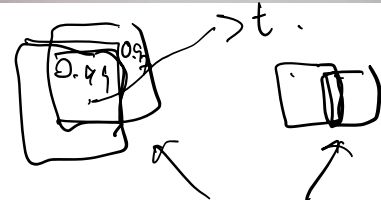
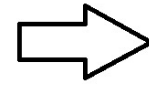
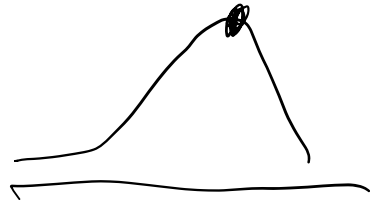
Non-Maximal Suppression

- Box Proposals: Samples of boxes that may come from a single object (latent)
- Each with a confidence score, density representation of the inferred latent distribution $q(z|x)$.

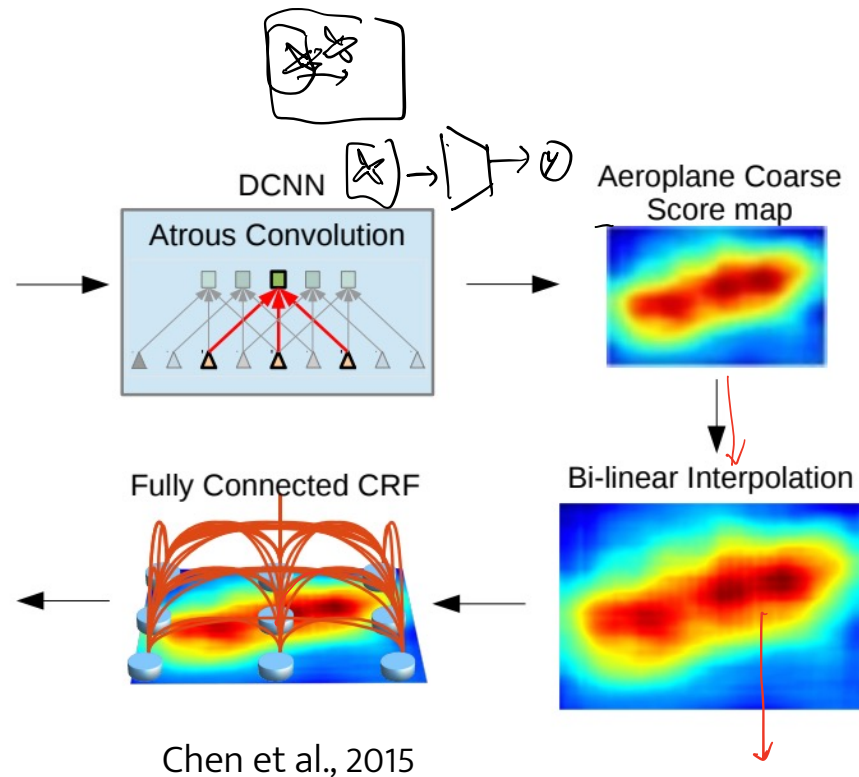
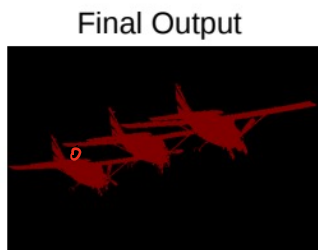
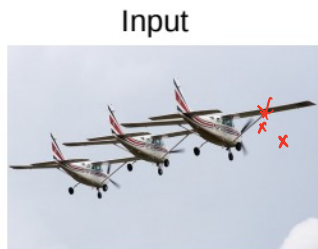
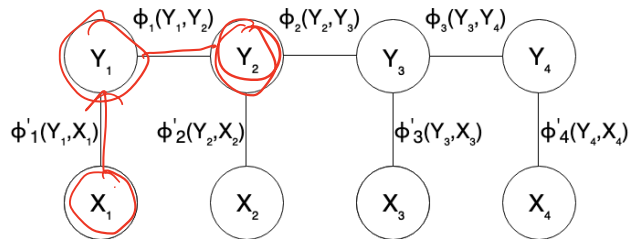


Non-Maximal Suppression

- Box Proposals: Samples of boxes that may come from a single object (latent)
- Each with a confidence score, density representation of the inferred latent distribution $q(z|x)$.
- MAP: take the mode of the distribution



Segmentation as CRF Inference



Inference Problem

- Knows $p(x|z)$.

Inference Problem

- Knows $p(x|z)$.
- Wants to know $p(z|x)$. Bayes rule.

$$\underline{p(z|x; \theta)} = \frac{p(x, z; \theta)}{\int_z p(x, z; \theta)} \cdot$$



Inference Problem

- Knows $p(x|z)$.
- Wants to know $p(z|x)$. Bayes rule.

$$p(z|x; \theta) = \frac{p(x, z; \theta)}{\int_z p(x, z; \theta)}$$

- Brute force

all possible latent state.

Inference Problem

- Knows $p(x|z)$.
- Wants to know $p(z|x)$. Bayes rule.

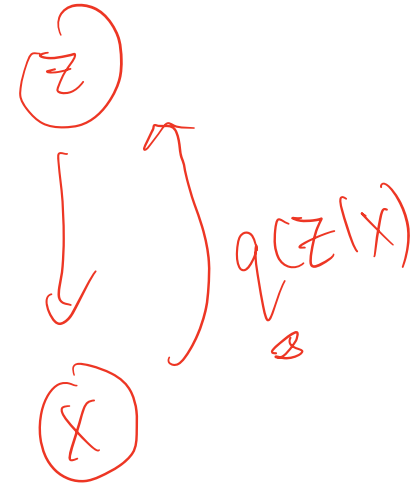
$$p(z|x; \theta) = \frac{p(x, z; \theta)}{\int_z p(x, z; \theta)}.$$

- Brute force
- Message passing, sum-product, belief propagation (DP)

Inference Problem

- Knows $p(x|z)$.
- Wants to know $p(z|x)$. Bayes rule.

$$p(z|x; \theta) = \frac{p(x, z; \theta)}{\int_z p(x, z; \theta)}.$$



- Brute force
- Message passing, sum-product, belief propagation (DP)
- Variational inference, mean field

Inference Problem

- Knows $p(x|z)$.
- Wants to know $p(z|x)$. Bayes rule.

$$p(z|x; \theta) = \frac{p(x, z; \theta)}{\int_z p(x, z; \theta)} \cdot \mathbb{E}$$

- Brute force
- Message passing, sum-product, belief propagation (DP)
- Variational inference, mean field
- Stochastic sampling, MCMC

Variational Inference

- Jensen's inequality to get ELBO.

$$\begin{aligned}\log p(x) &= \log \int_z p(x, z) \\ &= \log \int_z p(x, z) \frac{q(z)}{q(z)}\end{aligned}$$

Variational Inference

- Jensen's inequality to get ELBO.

$$\begin{aligned}\log p(x) &= \log \int_z p(x, z) \\ &= \log \int_z p(x, z) \frac{q(z)}{q(z)} \\ &= \log \left(\mathbb{E}_q \left[\frac{p(x, z)}{q(z)} \right] \right)\end{aligned}$$

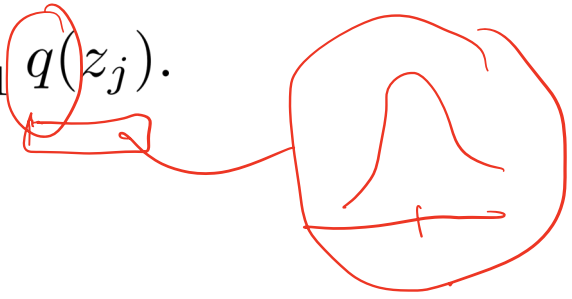
Variational Inference

- Jensen's inequality to get ELBO.

$$\begin{aligned} \log p(x) &= \log \int_z p(x, z) \\ &= \log \int_z p(x, z) \frac{q(z)}{q(z)} \\ &= \log \left(\mathbb{E}_q \frac{p(x, z)}{q(z)} \right) \\ &\geq \mathbb{E}_q \log \frac{p(x, z)}{q(z)} \\ &= \mathbb{E}_q \log p(x, z) - \mathbb{E}_q \log q(z) = \mathcal{L}. \end{aligned}$$

Mean-Field Inference

- If there are many latent variables, we can assume factorization (local variational approximation):

$$\underline{q(z_1, \dots, z_m)} = \prod_{j=1}^m q(z_j).$$


z_1 z_2, z_3



Mean-Field Inference

- If there are many latent variables, we can assume factorization (local variational approximation):

$$q(z_1, \dots, z_m) = \prod_{j=1}^m q(z_j).$$


$$\mathcal{L} = \log p(x) + \underbrace{\sum_{j=1}^m \mathbb{E}_{q(z_j)} \log p(z_j | z_{-j}, x)}_{\text{local}} - \underbrace{\mathbb{E}_{q(z_j)} \log(q(z_j))}_{\text{entropy}}.$$

Inference Operations

- CRF with pairwise energy. Use x as labels.

$$E(\mathbf{x}) = \sum_i \psi_u(x_i) + \sum_{i < j} \psi_p(x_i, x_j),$$

↓ unary. ↓ binary.



[Krähenbühl & Koltun, 2012]

Inference Operations

- CRF with pairwise energy. Use x as labels.

$$E(\mathbf{x}) = \sum_i \psi_u(x_i) + \sum_{i < j} \psi_p(x_i, x_j),$$

$$\psi_p(x_i, x_j) = \mu(x_i, x_j) \underbrace{\sum_{m=1}^K w^{(m)} k^{(m)}(\mathbf{f}_i, \mathbf{f}_j)}_{k(\mathbf{f}_i, \mathbf{f}_j)}.$$

[Krähenbühl & Koltun, 2012]

Inference Operations

- CRF with pairwise energy. Use x as labels.

$$E(\mathbf{x}) = \sum_i \psi_u(x_i) + \sum_{i < j} \psi_p(x_i, x_j),$$

$$\psi_p(x_i, x_j) = \mu(x_i, x_j) \underbrace{\sum_{m=1}^K w^{(m)} k^{(m)}(\mathbf{f}_i, \mathbf{f}_j)}_{k(\mathbf{f}_i, \mathbf{f}_j)}$$

$$k^{(m)}(\mathbf{f}_i, \mathbf{f}_j) = \exp\left(-\frac{1}{2}(\mathbf{f}_i - \mathbf{f}_j)^\top \Lambda^{(m)}(\mathbf{f}_i - \mathbf{f}_j)\right).$$
$$\underbrace{w^{(1)} \exp\left(-\frac{|p_i - p_j|^2}{2\theta_\alpha^2} - \frac{|I_i - I_j|^2}{2\theta_\beta^2}\right)}_{\text{appearance kernel}} + \underbrace{w^{(2)} \exp\left(-\frac{|p_i - p_j|^2}{2\theta_\gamma^2}\right)}_{\text{smoothness kernel}}$$

$$\mu(x_i, x_j) = [x_i \neq x_j]$$

[Krähenbühl & Koltun, 2012]

Inference in Fully Connected CRF

- Iterative mean-field inference.

Algorithm 1 Mean field in fully connected CRFs

Initialize Q

while not converged **do**

$\tilde{Q}_i^{(m)}(l) \leftarrow \sum_{j \neq i} k^{(m)}(\mathbf{f}_i, \mathbf{f}_j) Q_j(l)$ for all m

$\hat{Q}_i(x_i) \leftarrow \sum_{l \in \mathcal{L}} \mu^{(m)}(x_i, l) \sum_m w^{(m)} \tilde{Q}_i^{(m)}(l)$

$Q_i(x_i) \leftarrow \exp\{-\psi_u(x_i) - \hat{Q}_i(x_i)\}$

normalize $Q_i(x_i)$

end while

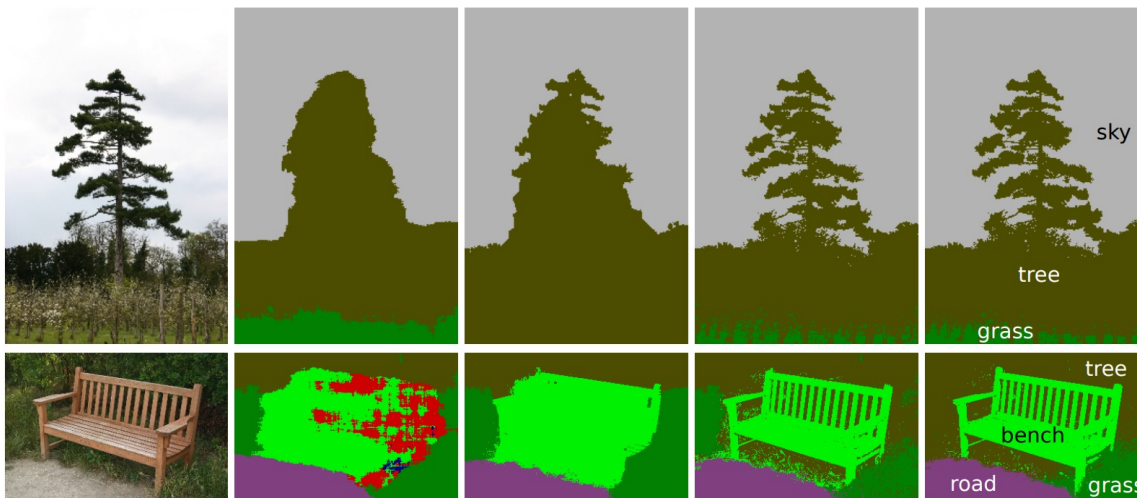
▷ $Q_i(x_i) \leftarrow \frac{1}{Z_i} \exp\{-\phi_u(x_i)\}$

▷ See Section 6 for convergence analysis

▷ **Message passing** from all X_j to all X_i

▷ **Compatibility transform**

▷ **Local update**



(a) Image

(b) Unary classifiers

(c) Robust P^n CRF

(d) Fully connected CRF, MCMC inference, 36 hrs

(e) Fully connected CRF, our approach, 0.2 seconds

[Krähenbühl & Koltun, 2012]

CRFs as RNNs

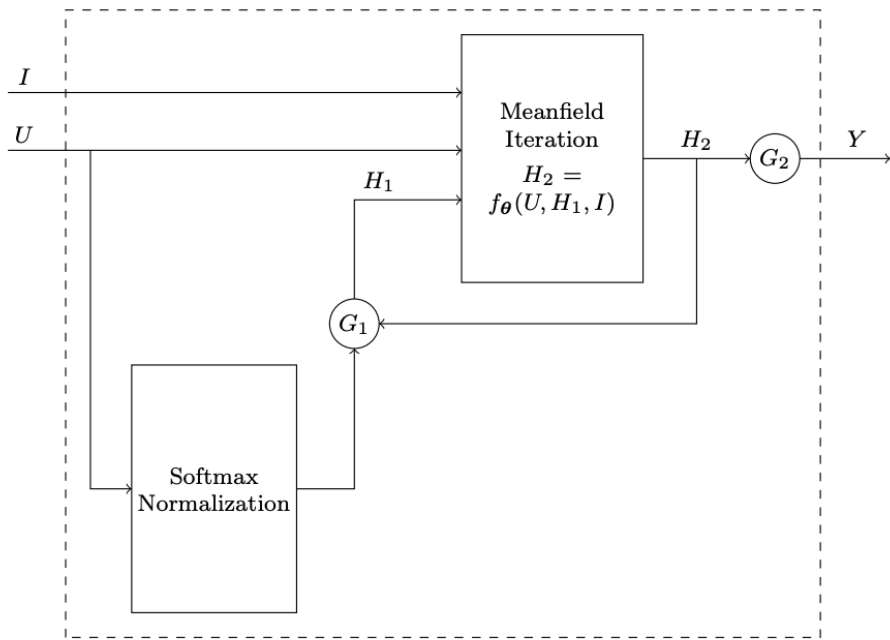
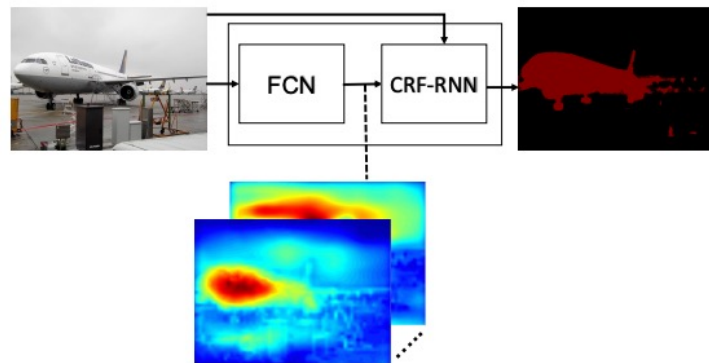


Figure 2. **The CRF-RNN Network.** We formulate the iterative mean-field algorithm as a Recurrent Neural Network (RNN). Generating functions G_1 and G_2 are fixed as described in the text.

$Q_i(l) \leftarrow \frac{1}{Z_i} \exp(U_i(l))$ for all i ▷ Initialization
while not converged do
 $\tilde{Q}_i^{(m)}(l) \leftarrow \sum_{j \neq i} k^{(m)}(\mathbf{f}_i, \mathbf{f}_j) Q_j(l)$ for all m ▷ Message Passing
 $\check{Q}_i(l) \leftarrow \sum_m w^{(m)} \tilde{Q}_i^{(m)}(l)$ ▷ Weighting Filter Outputs
 $\hat{Q}_i(l) \leftarrow \sum_{l' \in \mathcal{L}} \mu(l, l') \check{Q}_i(l')$ ▷ Compatibility Transform
 $\check{\check{Q}}_i(l) \leftarrow U_i(l) - \hat{Q}_i(l)$ ▷ Adding Unary Potentials
 $Q_i \leftarrow \frac{1}{Z_i} \exp(\check{\check{Q}}_i(l))$ ▷ Normalizing
end while

Summary

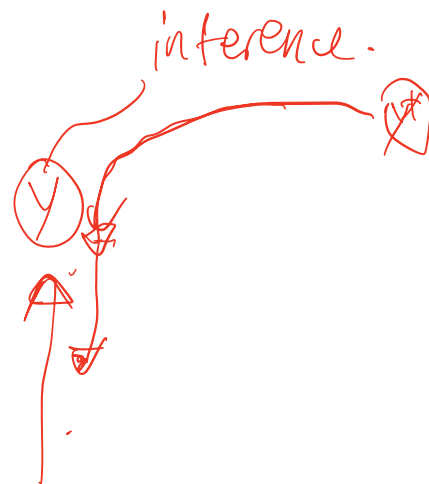
- Perception of high dimensional objects can be viewed as inferring latent variables with probabilistic distributions.

Summary

- Perception of high dimensional objects can be viewed as inferring latent variables with probabilistic distributions.
- We can impose structure.

Summary

- Perception of high dimensional objects can be viewed as inferring latent variables with probabilistic distributions.
- We can impose structure.
- We can learn through the inference process.
 - Taking the inference process into account.
 - Learning representations that matter.



Some Nuances

- If the process is deterministic or unimodal, standard deep networks may work.

Some Nuances

- If the process is deterministic or unimodal, standard deep networks may work.
- Network forward propagation vs. relaxed probabilistic inference.

Some Nuances

- If the process is deterministic or unimodal, standard deep networks may work.
- Network forward propagation vs. relaxed probabilistic inference.
- Having a stronger prior has the potential to be more data efficient.

Some Nuances

- If the process is deterministic or unimodal, standard deep networks may work.
- Network forward propagation vs. relaxed probabilistic inference.
- Having a stronger prior has the potential to be more data efficient.
- And you will need structured / generative learning when there are multiple modes.
 - E.g. Motion generation and planning: there can be multiple future trajectories



Deep Generative Models

- Autoregressive models

Deep Generative Models

- Autoregressive models
- Energy-based models
 - CRF
 - Max-Margin

Deep Generative Models

- Autoregressive models
- Energy-based models
 - CRF
 - Max-Margin
- Normalizing flow

Deep Generative Models

- Autoregressive models
- Energy-based models
 - CRF
 - Max-Margin
- Normalizing flow
- Diffusion models

Autoregressive Modeling

- Another type of output is autoregressive modeling.

Autoregressive Modeling

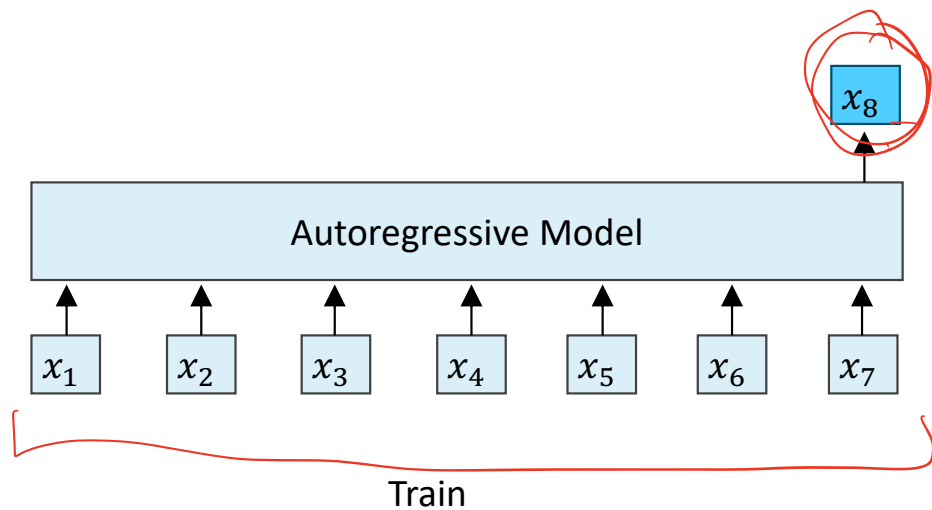
- Another type of output is autoregressive modeling.
- Example: Object detection/segmentation.

Autoregressive Modeling

- Another type of output is autoregressive modeling.
- Example: Object detection/segmentation.
- Intuition: Our visual attention focus on one object at a time.

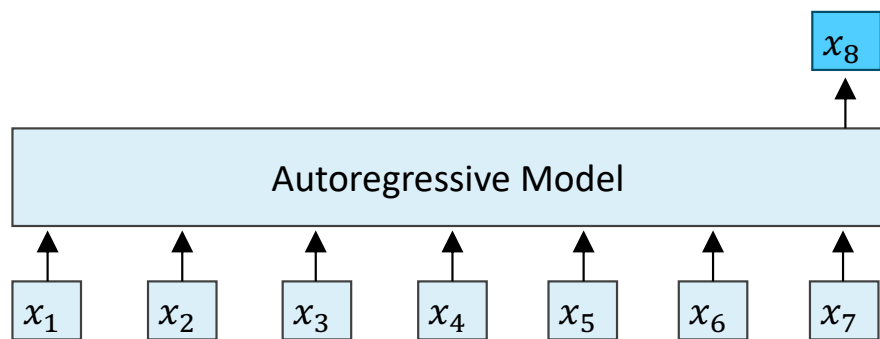
Teacher Forcing

- Teacher Forcing: Pretend that you know the whole sequence $\{y_1 \dots y_t\}$ at training time, and train for y_{t+1} .

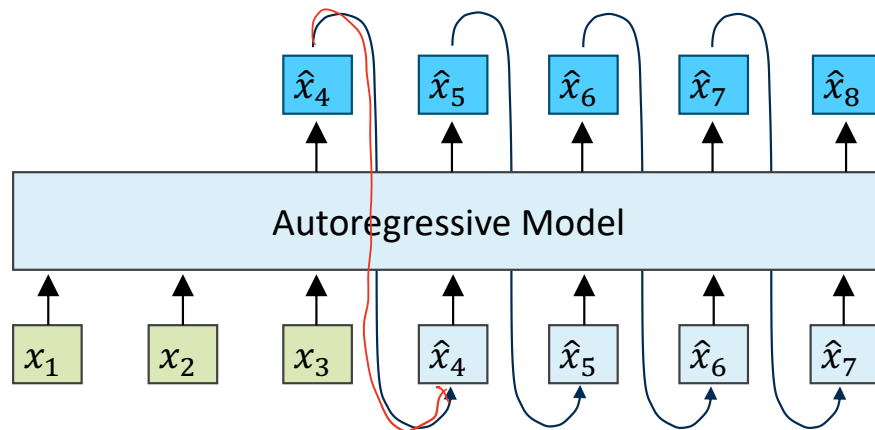


Teacher Forcing

- Teacher Forcing: Pretend that you know the whole sequence $\{y_1 \dots y_t\}$ at training time, and train for y_{t+1} .



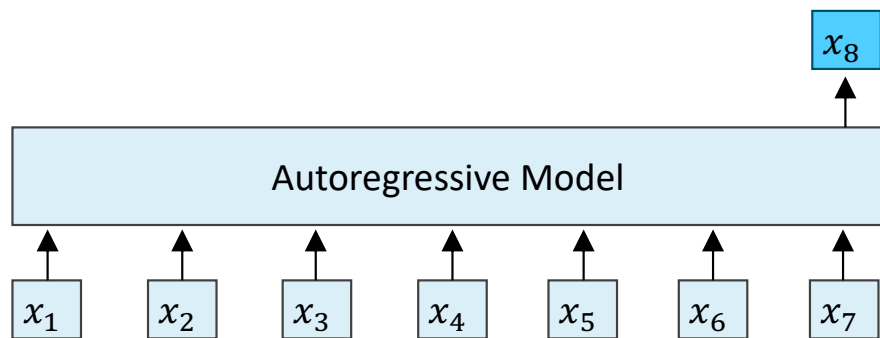
Train



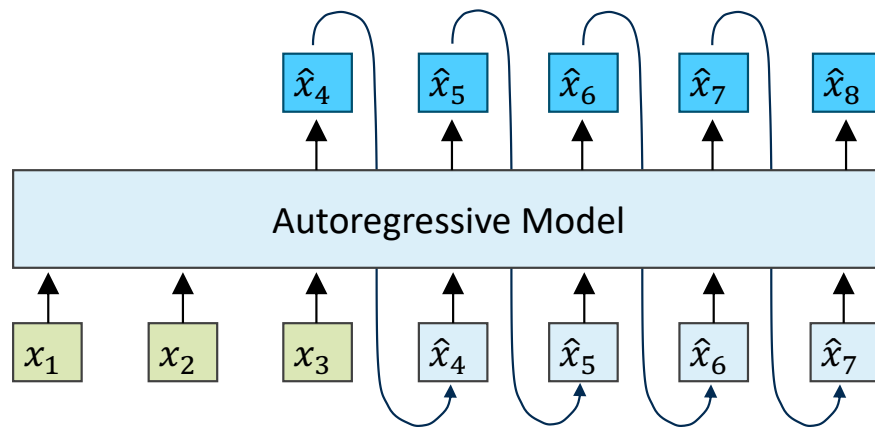
Test

Teacher Forcing

- Teacher Forcing: Pretend that you know the whole sequence $\{y_1 \dots y_t\}$ at training time, and train for y_{t+1} .
- Problem: You have to know the ordering.



Train



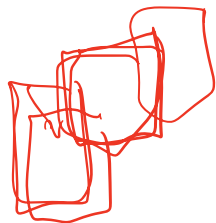
Test

More on Ordering

- There are many set to set problems.

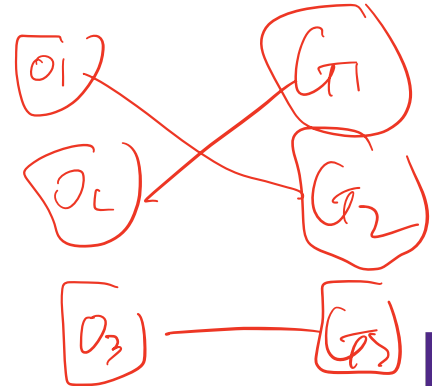
More on Ordering

- There are many set to set problems.
- E.g. Detection, segmentation, generating multiple objects, clustering



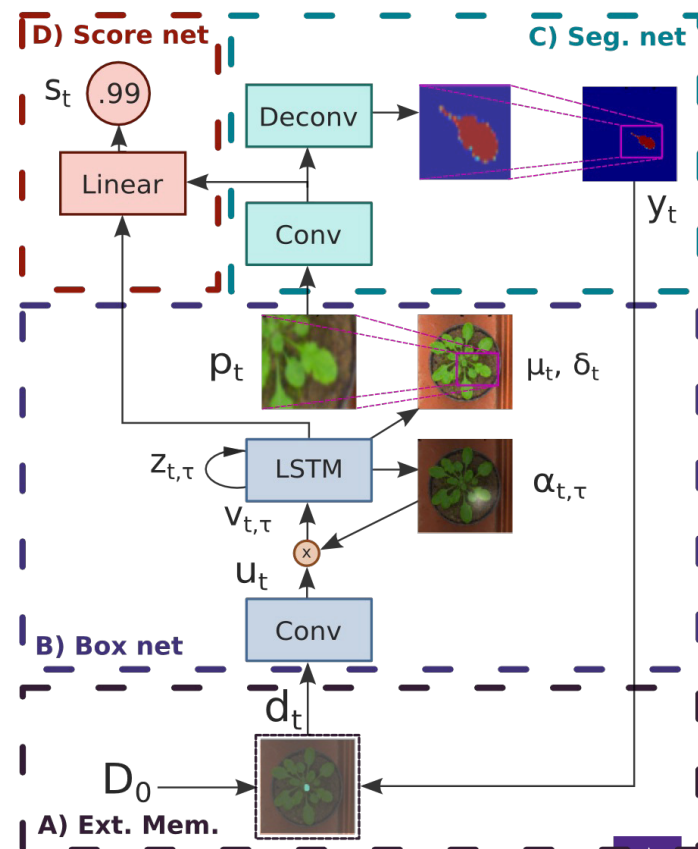
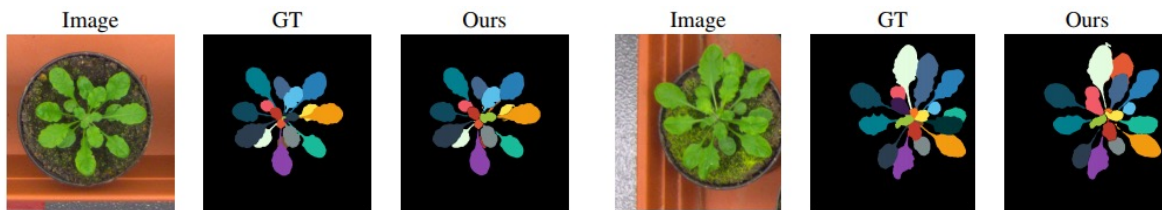
More on Ordering

- There are many set to set problems.
- E.g. Detection, segmentation, generating multiple objects, clustering
- Input does not follow a particular order. Loss function should not favor a particular order.
 - Attention: The attention operation is order-invariant.
 - Matching: The teacher can match the next input based on the closest matching ground truth instances.



Instance Segmentation using Attention

- Attending to one region at a time.
- Zooming in for segmentation.
- End-to-end differentiable box proposals and external memory.
- State-of-the-art on leaf segmentation problems for many years.



Summary

- Autoregressive models
 - Exact likelihood (instead of a lower bound or approximation)
 - Conditional probability may be easier to model
 - Need an ordering
 - Can take a long time to decode

Summary

- Autoregressive models
 - Exact likelihood (instead of a lower bound or approximation)
 - Conditional probability may be easier to model
 - Need an ordering
 - Can take a long time to decode
- Autoregressive models aren't just for language.
 - Planning
 - Set prediction problems
 - Perception

Summary

- Autoregressive models
 - Exact likelihood (instead of a lower bound or approximation)
 - Conditional probability may be easier to model
 - Need an ordering
 - Can take a long time to decode
- Autoregressive models aren't just for language.
 - Planning
 - Set prediction problems
 - Perception
- Generating the next variable is a different objective from learning a good representation of the whole sequence.

Summary

- Autoregressive models
 - Exact likelihood (instead of a lower bound or approximation)
 - Conditional probability may be easier to model
 - Need an ordering
 - Can take a long time to decode
- Autoregressive models aren't just for language.
 - Planning
 - Set prediction problems
 - Perception
- Generating the next variable is a different objective from learning a good representation of the whole sequence.
- No hierarchical planning.

Energy-Based Modeling

- The goal is to learn a distribution $p(x)$ to model the set of examples.

Energy-Based Modeling

- The goal is to learn a distribution $p(x)$ to model the set of examples.
- Assuming Boltzmann distribution:

$$p(x; \theta) = \frac{1}{Z(\theta)} \exp(-E(x; \theta))$$

$$Z(\theta) = \int_{\mathcal{X}} \exp(-E(x; \theta))$$

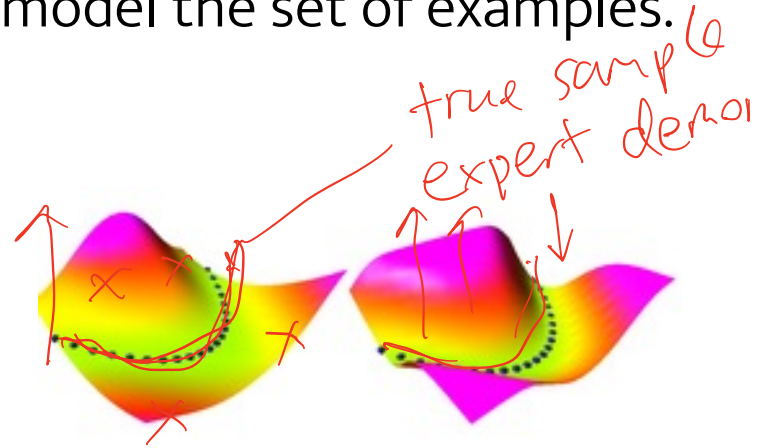
→ network.

Energy-Based Modeling

- The goal is to learn a distribution $p(x)$ to model the set of examples.
- Assuming Boltzmann distribution:

$$p(x; \theta) = \frac{1}{Z(\theta)} \exp(-E(x; \theta))$$

$$Z(\theta) = \int_X \exp(-E(x; \theta))$$



Picture from LeCun

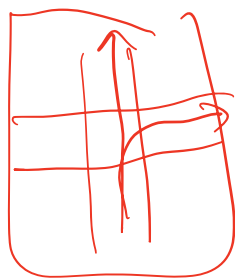
- Maximizing the likelihood:

$$\frac{\partial}{\partial \theta} \log p(x; \theta) = \mathbb{E}_{x' \sim p(x)} \left[\frac{\partial E(x'; \theta)}{\partial \theta} \right] - \frac{\partial E(x; \theta)}{\partial \theta}.$$

Handwritten red annotations: A red box surrounds the left side of the equation. Red circles highlight $\mathbb{E}_{x' \sim p(x)}$ and $\frac{\partial E(x'; \theta)}{\partial \theta}$. A red arrow points from the expectation term to the left plot. A red arrow points from the second term to the right plot. Handwritten red text 'real data.' is on the right.

Structured Prediction

- EBM can be easily adapted to model the joint distribution of x and y .
- Requires optimization of y at inference time.



traj: y .

visual scene: x .

arg min $E(x, y)$

y

Optimization

- Computing $\mathbb{E}_{x' \sim p(x)} \left[\frac{\partial E(x'; \theta)}{\partial \theta} \right]$ is non-trivial. We need to sample from $p(x)$.

Optimization

- Computing $\mathbb{E}_{x' \sim p(x)} \left[\frac{\partial E(x'; \theta)}{\partial \theta} \right]$ is non-trivial. We need to sample from $p(x)$.
- But we only know $E(x)$!

Optimization

- Computing $\mathbb{E}_{x' \sim p(x)} \left[\frac{\partial E(x'; \theta)}{\partial \theta} \right]$ is non-trivial. We need to sample from $p(x)$.
- But we only know $E(x)$!
- Through MCMC samplers: MH, Langevin, HMC, Gibbs.

Optimization

- Computing $\mathbb{E}_{x' \sim p(x)} \left[\frac{\partial E(x'; \theta)}{\partial \theta} \right]$ is non-trivial. We need to sample from $p(x)$.
- But we only know $E(x)$!
- Through MCMC samplers: MH, Langevin, HMC, Gibbs.
- Approximations: Using truncated steps (Contrastive Divergence)



Optimization

- Computing $\mathbb{E}_{x' \sim p(x)} \left[\frac{\partial E(x'; \theta)}{\partial \theta} \right]$ is non-trivial. We need to sample from $p(x)$.
- But we only know $E(x)$!
- Through MCMC samplers: MH, Langevin, HMC, Gibbs.
- Approximations: Using truncated steps (Contrastive Divergence)
- Score Matching: Tries to model the score function $\nabla_x \log p_\theta(x)$
 - If we know the gradient, we can improve the samples.
 - Want the model gradient to match with that from the data distribution
 - Closely related to diffusion models

Max-Entropy Inverse RL

- Want to model the reward function of trajectories.
- Given a reward function, the policy distribution can be modeled as a Boltzmann distribution.
- Given expert trajectories, we can maximize the likelihood under such distribution induced by the learned reward

$$p(\tau | \mathcal{O}_{1:T}, \psi) \propto \cancel{p(\tau)} \exp \left(\sum_t r_\psi(\mathbf{s}_t, \mathbf{a}_t) \right)$$

can ignore (independent of ψ)

maximum likelihood learning: $\max_{\psi} \frac{1}{N} \sum_{i=1}^N \log p(\tau_i | \mathcal{O}_{1:T}, \psi) = \max_{\psi} \frac{1}{N} \sum_{i=1}^N r_\psi(\tau_i) - \log Z$

Slide credit:
Sergey Levine



Max-Entropy Inverse RL

$$\max_{\psi} \frac{1}{N} \sum_{i=1}^N r_{\psi}(\tau_i) - \log Z$$

$$Z = \int p(\tau) \exp(r_{\psi}(\tau)) d\tau$$

$$\nabla_{\psi} \mathcal{L} = \frac{1}{N} \sum_{i=1}^N \nabla_{\psi} r_{\psi}(\tau_i) - \underbrace{\frac{1}{Z} \int p(\tau) \exp(r_{\psi}(\tau)) \nabla_{\psi} r_{\psi}(\tau) d\tau}_{p(\tau | \mathcal{O}_{1:T}, \psi)}$$

current policy.

$$\nabla_{\psi} \mathcal{L} = \underbrace{E_{\tau \sim \pi^*(\tau)} [\nabla_{\psi} r_{\psi}(\tau_i)]}_{\text{estimate with expert samples}} - \underbrace{E_{\tau \sim p(\tau | \mathcal{O}_{1:T}, \psi)} [\nabla_{\psi} r_{\psi}(\tau)]}_{\text{soft optimal policy under current reward}}$$

estimate with expert samples

soft optimal policy under current reward

Slide credit:
Sergey Levine



Max-Margin Learning

- In addition to probabilistic learning, we can also apply the max-margin framework.

Max-Margin Learning

- In addition to probabilistic learning, we can also apply the max-margin framework.
- Review: SSVM vs. CRF.

Max-Margin Learning

- In addition to probabilistic learning, we can also apply the max-margin framework.
- Review: SSVM vs. CRF.

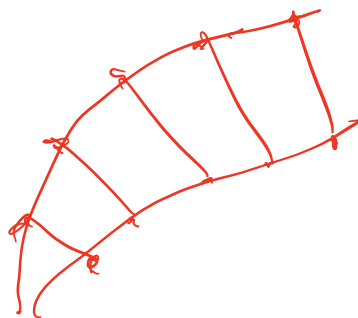
$$\arg \min_{\theta} \sum_{n=1}^N \max[0, m + \underbrace{E(x_n; \theta)}_{\substack{\text{n}^{\text{th}} \text{ demo}}} - \underbrace{E(x^*; \theta)}_{\substack{\text{best output} \\ \text{from model}}}]] + \lambda \|\theta\|_2^2.$$

Max-Margin Learning

- In addition to probabilistic learning, we can also apply the max-margin framework.
- Review: SSVM vs. CRF.

$$\arg \min_{\theta} \sum_{n=1}^N \max[0, m + E(x_n; \theta) - E(x^*; \theta)] + \lambda \|\theta\|_2^2.$$

- Margin can be difference in trajectories.



Max-Margin Learning

- In addition to probabilistic learning, we can also apply the max-margin framework.
- Review: SSVM vs. CRF.

$$\arg \min_{\theta} \sum_{n=1}^N \max[0, m + E(x_n; \theta) - E(x^*; \theta)] + \lambda \|\theta\|_2^2.$$

- Margin can be difference in trajectories.
- Non-probabilistic

Max-Margin Learning

- In addition to probabilistic learning, we can also apply the max-margin framework.
- Review: SSVM vs. CRF.

$$\arg \min_{\theta} \sum_{n=1}^N \max[0, m + E(x_n; \theta) - E(x^*; \theta)] + \lambda \|\theta\|_2^2.$$

- Margin can be difference in trajectories.
- Non-probabilistic
- Still need to run optimization to find the best x^*

samples!

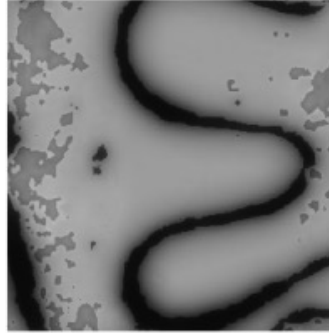


Max-Margin Planning

mode 1 - training



mode 1 - learned cost map over novel region



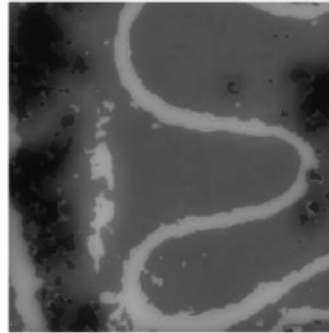
mode 1 - learned path over novel region



mode 2 - training



mode 2 - learned cost map over novel region

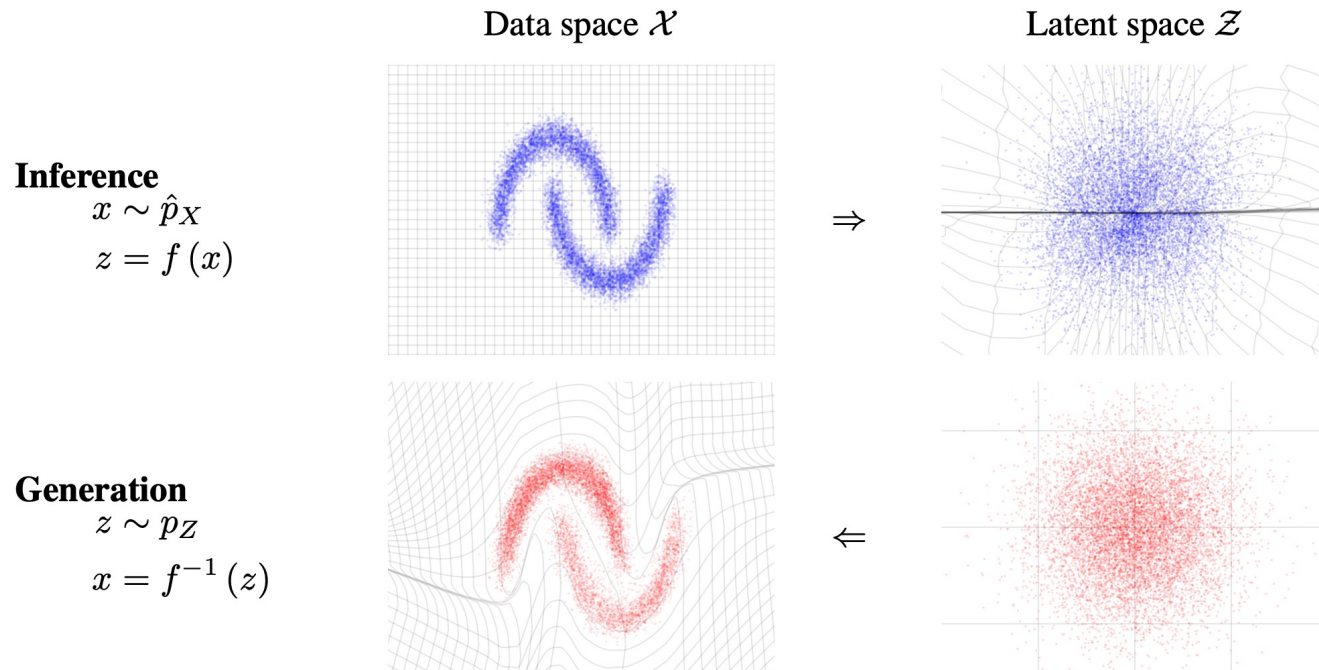


mode 2 - learned path over novel region



Learning Latent Representations

- z may contain useful high-level compressed information of the input.



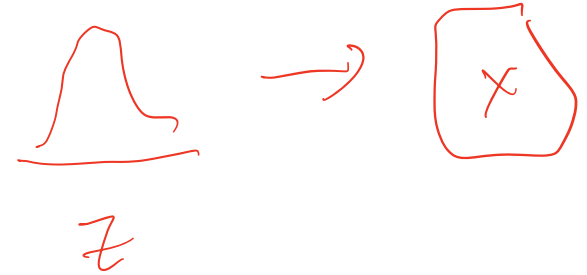
Variational Autoencoders (VAE)

- \mathbf{z} from a simple prior distribution.
- Optimize ELBO.

$$\log p_{\theta}(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

Variational Autoencoders (VAE)

- \mathbf{z} from a simple prior distribution.
- Optimize ELBO.



$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &= \log \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \\ &= \log \int \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}\end{aligned}$$

Variational Autoencoders (VAE)

- \mathbf{z} from a simple prior distribution.
- Optimize ELBO.

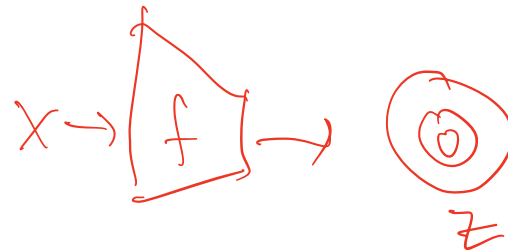
$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &= \log \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \\ &= \log \int \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \\ &\geq -\mathbb{D}_{\text{KL}}[q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})] + \mathbb{E}_q[\log p_{\theta}(\mathbf{x}|\mathbf{z})]\end{aligned}$$

VAE Problems

- Hard to transform the inputs to the prior distribution with a network.
- The coarser the approximation of q , the worse the lower bound.

$$\begin{aligned} ELBO &= \int_z q(z) \log \frac{p(x, z)}{q(z)} \\ &= \mathbb{E}_{q(z)} \log \frac{p(z|x)p(x)}{q(z)} \end{aligned}$$

VAE Problems



- Hard to transform the inputs to the prior distribution with a network.
- The coarser the approximation of q , the worse the lower bound.

$$\begin{aligned} ELBO &= \int_z q(z) \log \frac{p(x, z)}{q(z)} \\ &= \mathbb{E}_{q(z)} \log \frac{p(z|x)p(x)}{q(z)} \\ &= -\mathbb{E}_{q(z)} \log \frac{q(z)}{p(z|x)} + \mathbb{E}_{q(z)} \log p(x) \\ &= \underbrace{-KL(q(z) || p(z|x))}_{\downarrow} + \underbrace{\log p(x)}_{\text{true posterior.}} \end{aligned}$$

Normalizing Flows

- Requires a bijection between the input and output.

Normalizing Flows

- Requires a bijection between the input and output.
- Needs to maintain dimensionality.

Normalizing Flows

- Requires a bijection between the input and output.
- Needs to maintain dimensionality.

$$p(x)$$

$$p(z) = p(x) \cdot |\det J|$$

$$\frac{1}{2} \|z\|^2$$

$$\max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} \log p_{\text{NF}}(x; \theta) = \log p_0(\overset{z}{f_{\theta}(x)}; \theta) + \log \left(\left| \det \left(\frac{\partial f_{\theta}(x)}{\partial x} \right) \right| \right)$$

↓
gaussian.

Normalizing Flows

- Requires a bijection between the input and output.
- Needs to maintain dimensionality.
- Needs an efficient way to estimate determinant of Jacobian (typically $O(n^3)$).

$$\max_{\theta} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log p_{\text{NF}}(\mathbf{x}; \theta) = \log p_0(f_{\theta}(\mathbf{x}); \theta) + \log \left(\left| \det \left(\frac{\partial f_{\theta}(\mathbf{x})}{\partial \mathbf{x}} \right) \right| \right)$$

Adding Constraints in Transformation

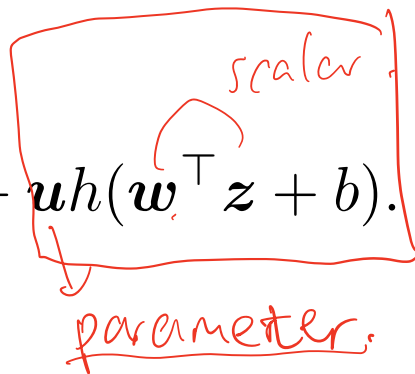
- Rezende & Mohamed (2015):

$$f(\mathbf{z}) = \mathbf{z} + \mathbf{u}h(\mathbf{w}^\top \mathbf{z} + b).$$

Adding Constraints in Transformation

- Rezende & Mohamed (2015):

$$f(\mathbf{z}) = \mathbf{z} + \mathbf{u}h(\mathbf{w}^\top \mathbf{z} + b).$$



- Determinant of Jacobian:

$$\psi(\mathbf{z}) = h'(\mathbf{w}^\top \mathbf{z} + b)\mathbf{w}.$$

Adding Constraints in Transformation

- Rezende & Mohamed (2015):

$$f(\mathbf{z}) = \mathbf{z} + \mathbf{u}h(\mathbf{w}^\top \mathbf{z} + b).$$

- Determinant of Jacobian:

$$\psi(\mathbf{z}) = h'(\mathbf{w}^\top \mathbf{z} + b)\mathbf{w}.$$

$$\left| \det \frac{\partial f}{\partial \mathbf{z}} \right| = \left| \det(\mathbf{I} + \mathbf{u}\psi(\mathbf{z})^\top) \right| = |1 + \mathbf{u}^\top \psi(\mathbf{z})|.$$

Masked Autoregressive Flow (MAF)

- Another way to make Jacobian determinant computable -> Lower triangular matrix.

Masked Autoregressive Flow (MAF)

- Another way to make Jacobian determinant computable -> Lower triangular matrix.
- Consider the autoregression of generating one dimension at a time.

$$x_i = \mu_i(x_{1:i-1}) + \sigma_i(x_{1:i-1})z_i, \quad z_i = \frac{x_i - \mu_i(x_{1:i-1})}{\sigma_i(x_{1:i-1})}.$$

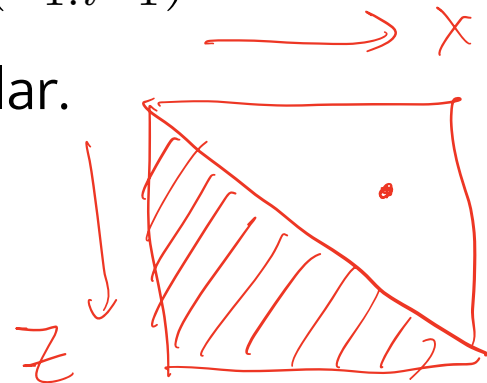
Masked Autoregressive Flow (MAF)

- Another way to make Jacobian determinant computable -> Lower triangular matrix.
- Consider the autoregression of generating one dimension at a time.

$$x_i = \mu_i(x_{1:i-1}) + \sigma_i(x_{1:i-1})z_i, \quad z_i = \frac{x_i - \mu_i(x_{1:i-1})}{\sigma_i(x_{1:i-1})}.$$

- z_i does not depend on x_j for any $j > i$. Lower triangular.

z x



Masked Autoregressive Flow (MAF)

- Another way to make Jacobian determinant computable -> Lower triangular matrix.
- Consider the autoregression of generating one dimension at a time.

$$x_i = \mu_i(x_{1:i-1}) + \sigma_i(x_{1:i-1})z_i, \quad z_i = \frac{x_i - \mu_i(x_{1:i-1})}{\sigma_i(x_{1:i-1})}.$$

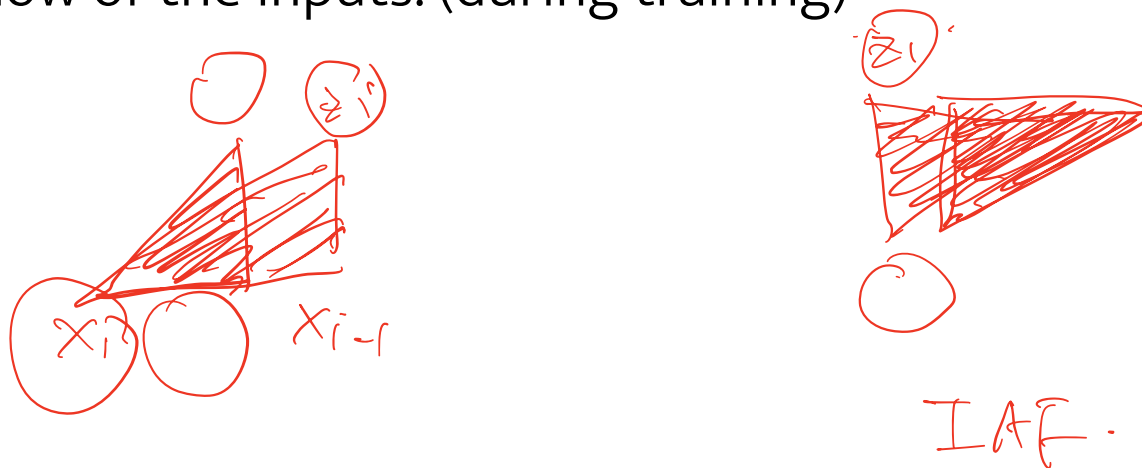
- z_i does not depend on x_j for any $j > i$. Lower triangular.

$$\log |\det J| = \sum_{i=1}^D \log \left| \frac{1}{\sigma_i(x_{1:i-1})} \right| = - \sum_{i=1}^D \log \sigma_i(x_{1:i-1})$$

$$L = \sum_i \left(\frac{1}{2} \|z_i\|_2^2 + \log \sigma_i \right)$$

MAF vs IAF

- MAF inference can be parallelized by directly passing different window of the inputs. (during training)



MAF vs IAF

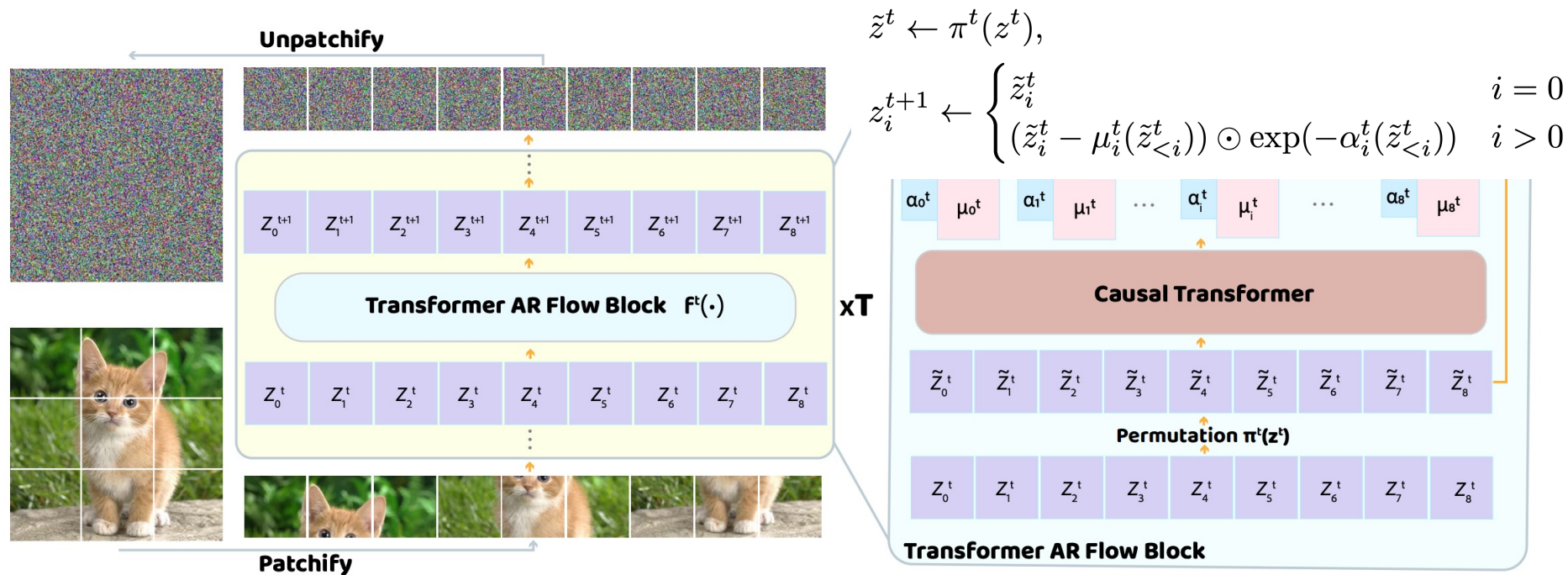
- MAF inference can be parallelized by directly passing different window of the inputs. (during training)
- Sampling needs to be sequential: One input dimension at a time.

MAF vs IAF

- MAF inference can be parallelized by directly passing different window of the inputs. (during training)
- Sampling needs to be sequential: One input dimension at a time.
- Inverse Autoregressive Flow (IAF) does the reverse. Parallel sample, sequential inference (slow training).

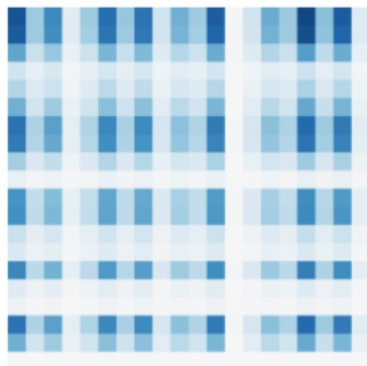
Block Autoregressive Flows

- Using blocks, permutations, and element-wise operations



Full-Rank Jacobian

- Is it possible to make the architecture more flexible?



Low-rank



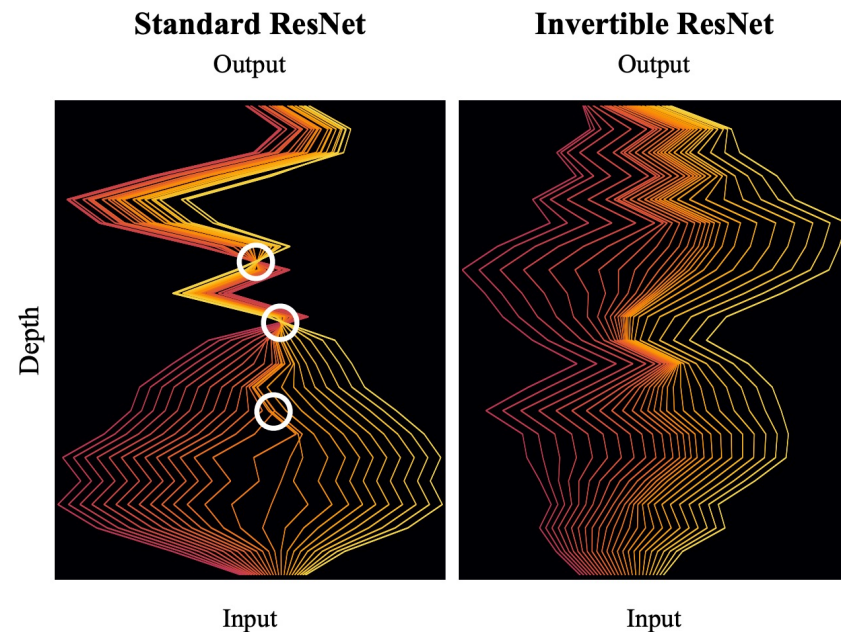
Auto-regressive



Full-rank

Invertibility of ResNet

$$x_{t+1} \leftarrow x_t + g_{\theta_t}(x_t)$$



Invertibility of ResNet

$$x_{t+1} \leftarrow x_t + g_{\theta_t}(x_t)$$

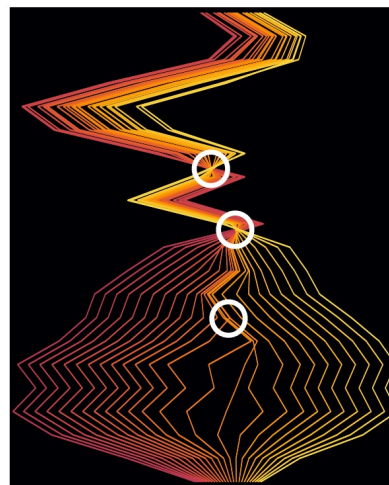
$\hat{x}_t \in x_{t+1} - g_{\theta_t}(x_{t+1})$

- The inverse can be recovered from fixed point iteration if $\text{Lip}(g) < 1$.

$$x_t^0 := x_{t+1}$$

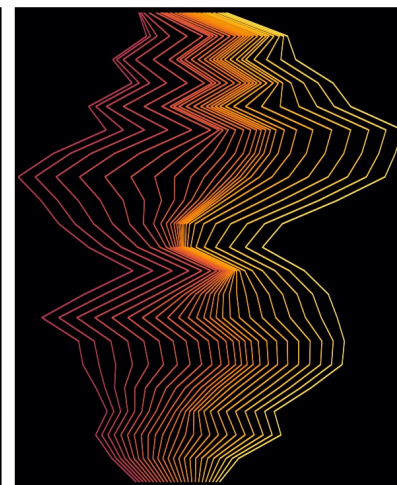
$$x_t^{k+1} := x_{t+1} - g_{\theta_t}(x_t^k)$$

Standard ResNet
Output



Input

Invertible ResNet
Output



Input

Invertibility of ResNet

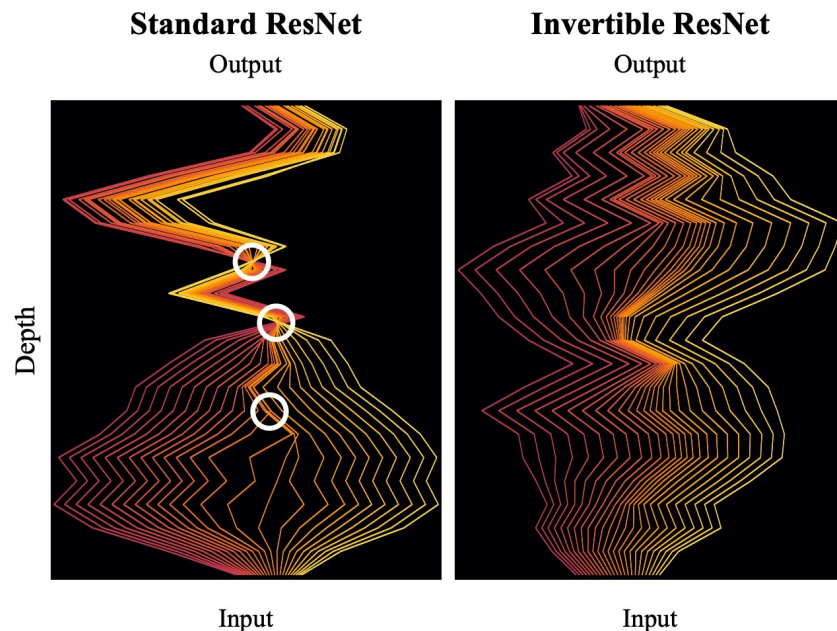
$$x_{t+1} \leftarrow x_t + g_{\theta_t}(x_t)$$

- The inverse can be recovered from fixed point iteration if $\text{Lip}(g) < 1$.

$$x_t^0 := x_{t+1}$$

$$x_t^{k+1} := x_{t+1} - g_{\theta_t}(x_t^k)$$

- MAF uses lower triangular Jacobian structure while iResNet uses approximation.



Estimating Jacobian Determinant

Estimating Jacobian Determinant

- Using power series to estimate $\det J$.

$$\log |\det J_F| = \text{tr}(\log J_F).$$

$$\text{tr}(\log(I + J_g)) = \sum_{k=1}^{\infty} (-1)^{k+1} \frac{\text{tr}(J_g^k)}{k}.$$

Estimating Jacobian Determinant

- Using power series to estimate $\det J$.

$$\log |\det J_F| = \text{tr}(\log J_F).$$

$$\text{tr}(\log(I + J_g)) = \sum_{k=1}^{\infty} (-1)^{k+1} \frac{\text{tr}(J_g^k)}{k}.$$

Handwritten red annotations showing the trace of a matrix product. The top part shows $v^T J v$ with a horizontal line underneath, and the bottom part shows $v^T J$.

- To compute the trace without the full Jacobian. Use a stochastic estimate:

$$\text{tr}(A) = \mathbb{E}_v[v^T A v].$$

Estimating Jacobian Determinant

- Using power series to estimate $\det J$.

$$\log |\det J_F| = \text{tr}(\log J_F).$$

$$\text{tr}(\log(I + J_g)) = \sum_{k=1}^{\infty} (-1)^{k+1} \frac{\text{tr}(J_g^k)}{k}.$$

- To compute the trace without the full Jacobian. Use a stochastic estimate:

$$\text{tr}(A) = \mathbb{E}_v[v^\top A v].$$

Draw v from $\mathcal{N}(\check{0}, I)$

$$w^T := v^T$$

$$\ln \det := 0$$

for $k = 1$ **to** n **do**

$$w^T := w^T J_g \text{ (vector-Jacobian product)}$$

$$\ln \det := \ln \det + (-1)^{k+1} w^T v / k$$

end for

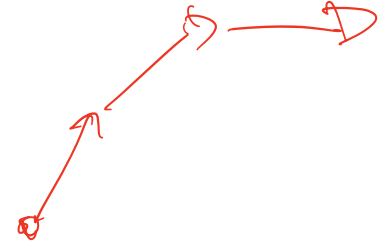
Summary: Normalizing Flow

- Requires invertibility and tractability on Jacobian determinant
- Constraint on the functional form
- Or estimation of Jacobian / iterative inverse (numerical approximation).

power series
+ stochastic
estimate.

fixed.

Inference as Integration



- Residual Flow: $\phi(x) = x + \delta u(x)$

$$\frac{\phi(x) - x}{\delta} = u(x)$$

Inference as Integration

- Residual Flow: $\phi(x) = x + \delta u(x)$

$$\frac{\phi(x) - x}{\delta} = u(x)$$

CNF
↑

- Can be written in continuous time (continuous normalizing flow).

$$x_t = \phi_t(x_0) = x_0 + \int_0^t \delta u_s(x_s) ds.$$

Inference as Integration

- Residual Flow: $\phi(x) = x + \delta u(x)$

$$\frac{\phi(x) - x}{\delta} = u(x)$$

- Can be written in continuous time (continuous normalizing flow).

$$x_t = \phi_t(x_0) = x_0 + \int_0^t \delta u_s(x_s) ds.$$

- Integration: Euler, or directly run an ODE solver

Inference as Integration

- Residual Flow: $\phi(x) = x + \delta u(x)$

$$\frac{\phi(x) - x}{\delta} = u(x)$$

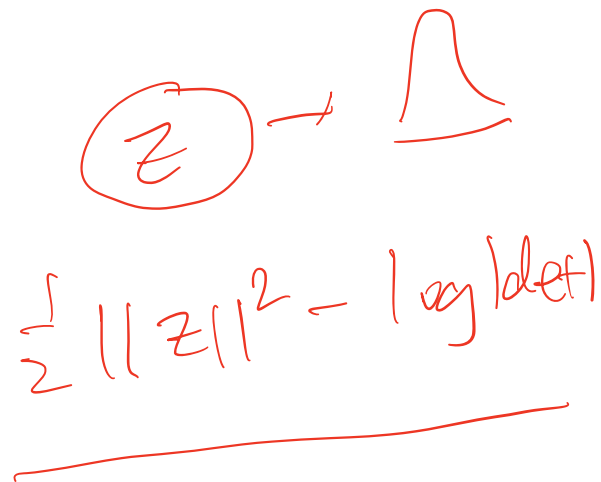
- Can be written in continuous time (continuous normalizing flow).

$$x_t = \phi_t(x_0) = x_0 + \int_0^t \delta u_s(x_s) ds.$$

- Integration: Euler, or directly run an ODE solver
- Backprop solving an ODE without storing intermediate states

Flow Matching

- Taking integrals during training is cumbersome, and the learning signal from the top latent is weak.

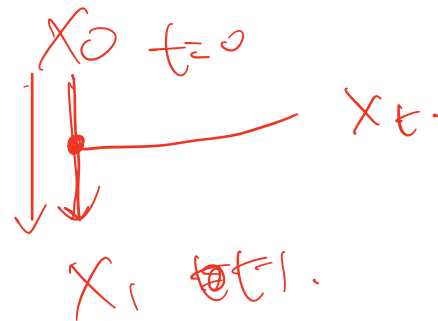


Handwritten diagram illustrating the loss function for a latent variable z . A circle containing z is connected by an arrow to a bell-shaped curve. Below this, the equation $\int \frac{1}{2} \|z\|^2 - \log |\det|$ is written and underlined.

Flow Matching

- Taking integrals during training is cumbersome, and the learning signal from the top latent is weak.

- Given an original input x_1 and a noise x_0 , $x_t = (1 - t)x_0 + tx_1$



Flow Matching

- Taking integrals during training is cumbersome, and the learning signal from the top latent is weak.
- Given an original input x_1 and a noise x_0 , $x_t = (1 - t)x_0 + tx_1$
- Want to fit a velocity $v_t(x)$ to match $u_t = x_1 - x_0$



Flow Matching

- Taking integrals during training is cumbersome, and the learning signal from the top latent is weak.
- Given an original input x_1 and a noise x_0 , $x_t = (1 - t)x_0 + tx_1$
- Want to fit a velocity $v_t(x)$ to match $u_t = x_1 - x_0$
- Flow Matching learning objective:

$$\mathcal{L}_{FM}(\theta) = \mathbb{E}_{t, x_0, x_1} [\|v(x_t, t) - u_t(x_0, x_1)\|^2]$$

Diffusion Models

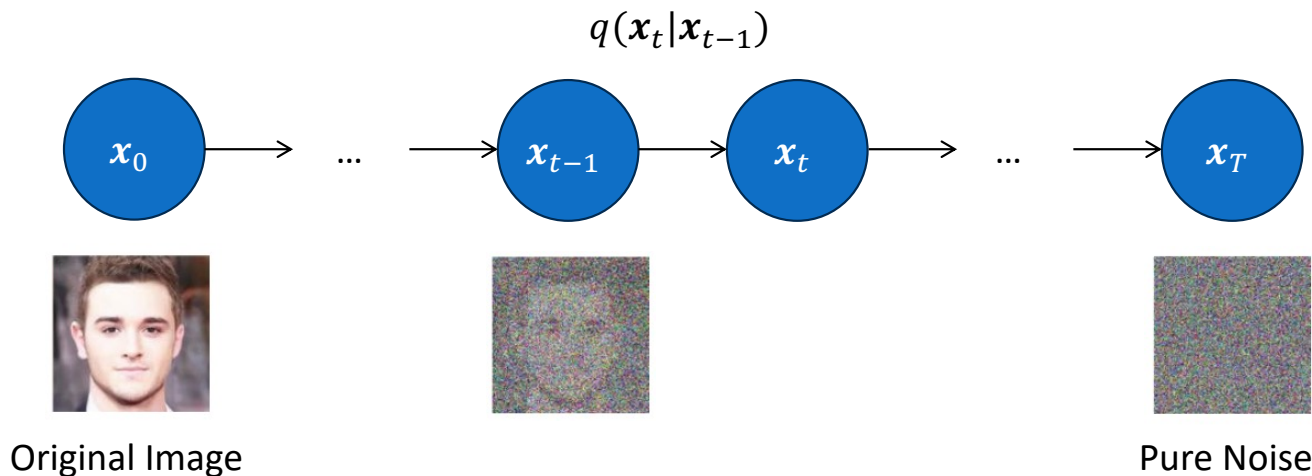
- A popular generative model formulation.

Diffusion Models

- A popular generative model formulation.
- The intuition is to iteratively denoise from Gaussian random noises into an image.

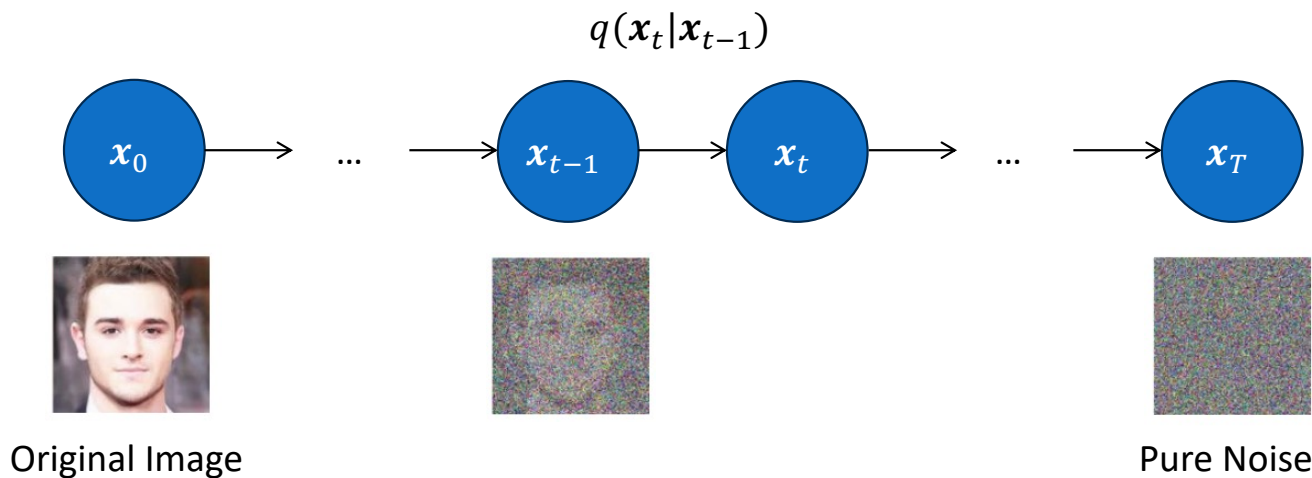
Diffusion Models

- A popular generative model formulation.
- The intuition is to iteratively denoise from Gaussian random noises into an image.



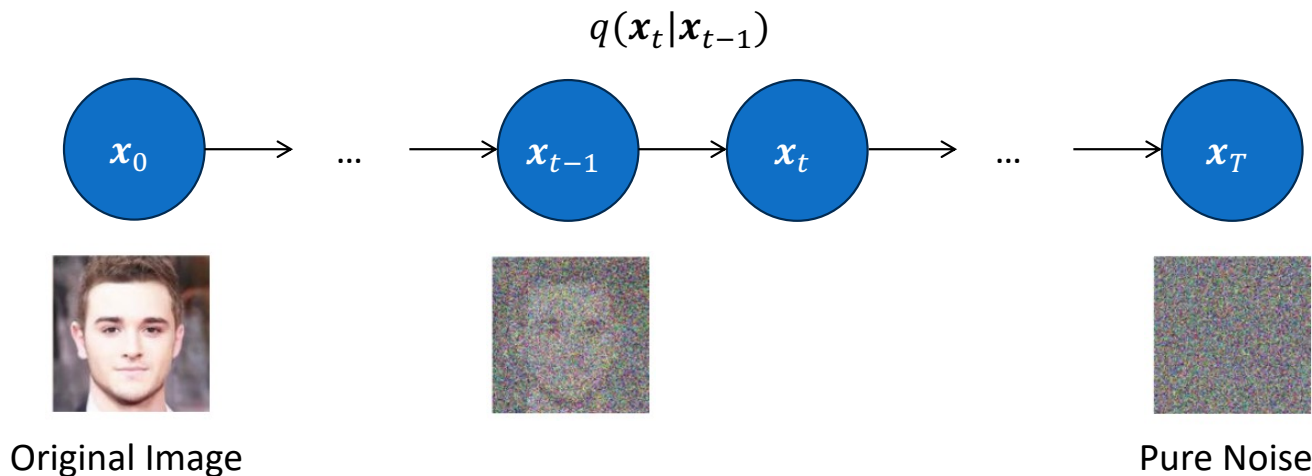
Diffusion Models

- Forward process: $q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I)$.



Diffusion Models

- Forward process: $q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I)$.
- You can also write: $x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon_t$, $\epsilon_t \sim \mathcal{N}(0, I)$.



Properties of the Forward Process

- Forward process: $x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon_t$

Properties of the Forward Process

- Forward process: $x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon_t$
- Write x_t as a function of x_0 with larger noises:

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I).$$

Properties of the Forward Process

- Forward process: $x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon_t$
- Write x_t as a function of x_0 with larger noises:

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I).$$

- Cumulative schedule: $\alpha_t = 1 - \beta_t$. $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$.

Properties of the Forward Process

- Forward process: $x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon_t$
- Write x_t as a function of x_0 with larger noises:

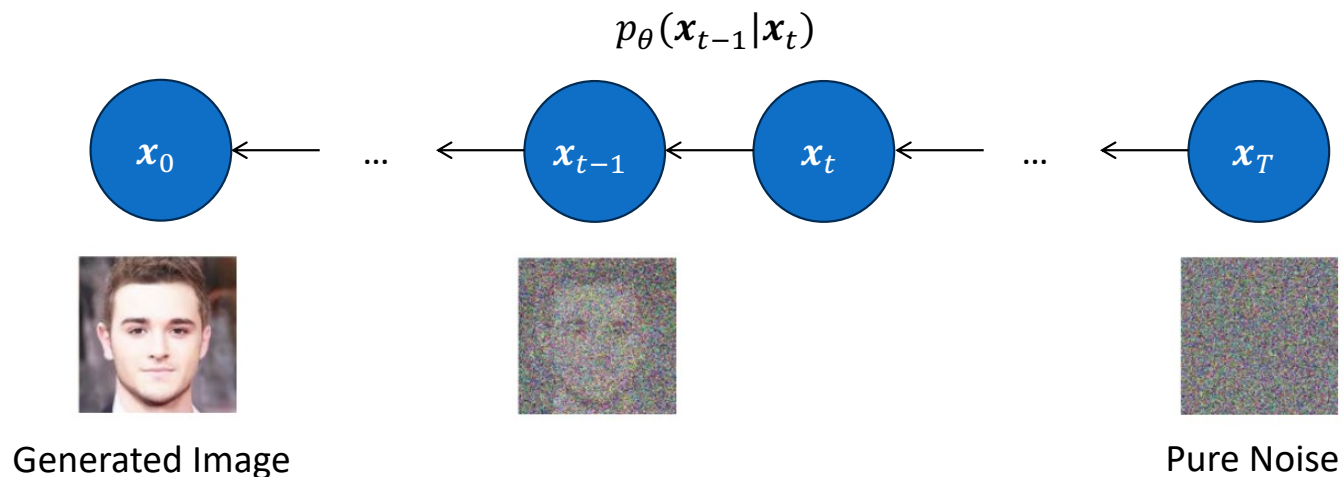
$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I).$$

- Cumulative schedule: $\alpha_t = 1 - \beta_t$. $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$.
- $x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$.

Reverse Process

- So, we need to learn a “model”:

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_{\theta}(x_t, t)).$$



Reverse Process

- Compute μ_θ ? Derive $p(x_{t-1}|x_t)$.

Reverse Process

- Compute μ_θ ? Derive $p(x_{t-1}|x_t)$.

- Bayes rule:

$$q(x_{t-1}|x_t) = \frac{q(x_t|x_{t-1})q(x_{t-1})}{q(x_t)}.$$

Reverse Process

- Compute μ_θ ? Derive $p(x_{t-1}|x_t)$.

- Bayes rule:
$$q(x_{t-1}|x_t) = \frac{q(x_t|x_{t-1})q(x_{t-1})}{q(x_t)}.$$

- But we don't know the marginal distribution $q(x_{t-1})$. We only know q_T and $q(x_t|x_{t-1})$.

Reverse Process

- Compute μ_θ ? Derive $p(x_{t-1}|x_t)$.

- Bayes rule:
$$q(x_{t-1}|x_t) = \frac{q(x_t|x_{t-1})q(x_{t-1})}{q(x_t)}.$$

- But we don't know the marginal distribution $q(x_{t-1})$. We only know q_T and $q(x_t|x_{t-1})$.
- Solution: Condition on the original input x_0 :

$$q(x_{t-1}|x_t, x_0) = \frac{q(x_t|x_{t-1})q(x_{t-1}|x_0)}{q(x_t|x_0)}.$$

Reverse Process

$$q(x_{t-1}|x_t, x_0) = \frac{q(x_t|x_{t-1})q(x_{t-1}|x_0)}{q(x_t|x_0)}.$$

Reverse Process

$$q(x_{t-1}|x_t, x_0) = \frac{q(x_t|x_{t-1})q(x_{t-1}|x_0)}{q(x_t|x_0)}.$$

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}_t, \tilde{\beta}I).$$

Reverse Process

$$q(x_{t-1}|x_t, x_0) = \frac{q(x_t|x_{t-1})q(x_{t-1}|x_0)}{q(x_t|x_0)}.$$

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}_t, \tilde{\beta}I).$$

$$\tilde{\mu}_t = \frac{\sqrt{\alpha_t}\bar{\beta}_{t-1}}{\bar{\beta}_t}x_t + \frac{\sqrt{\alpha_{t-1}}\beta_t}{\bar{\beta}_t}x_0 = \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon\right).$$

Reverse Process

$$q(x_{t-1}|x_t, x_0) = \frac{q(x_t|x_{t-1})q(x_{t-1}|x_0)}{q(x_t|x_0)}.$$

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}_t, \tilde{\beta}I).$$

$$\tilde{\mu}_t = \frac{\sqrt{\alpha_t}\bar{\beta}_{t-1}}{\bar{\beta}_t}x_t + \frac{\sqrt{\alpha_{t-1}}\beta_t}{\bar{\beta}_t}x_0 = \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon\right).$$

- Want: train up a μ_θ to match with $\tilde{\mu}_t$.

Training

- Sometimes it is more common to predict the denoising vector ϵ instead of μ .

$$\tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right),$$

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t) \right).$$

Training

- Sometimes it is more common to predict the denoising vector ϵ instead of μ .
- Randomly pick at a time step and predict the difference between the noisy and the original.

$$\tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right),$$

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t) \right).$$

Training

- Sometimes it is more common to predict the denoising vector ϵ instead of μ .
- Randomly pick at a time step and predict the difference between the noisy and the original.

$$\tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right),$$

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t) \right).$$

Algorithm 1 Training

- 1: **repeat**
 - 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
 - 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
 - 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 5: Take gradient descent step on
 $\nabla_\theta \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2$
 - 6: **until** converged
-

Sampling

- How do we sample an image?

Algorithm 2 Sampling

- 1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 2: **for** $t = T, \dots, 1$ **do**
 - 3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$
 - 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
 - 5: **end for**
 - 6: **return** \mathbf{x}_0
-

Sampling

- How do we sample an image?
- We know μ_θ which will help us transition from x_t to x_{t-1} .

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t) \right).$$

Algorithm 2 Sampling

- 1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 2: **for** $t = T, \dots, 1$ **do**
 - 3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$
 - 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
 - 5: **end for**
 - 6: **return** \mathbf{x}_0
-

Sampling

- How do we sample an image?
- We know μ_θ which will help us transition from x_t to x_{t-1} .

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t) \right).$$

- Sample from $\mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \sigma_t^2)$
. σ_t can either be β_t or $\tilde{\beta}_t$ derived from the posterior.

Algorithm 2 Sampling

- 1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 2: **for** $t = T, \dots, 1$ **do**
 - 3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$
 - 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
 - 5: **end for**
 - 6: **return** \mathbf{x}_0
-

Deep Generative Models Summary

Guided Diffusion

- We can add guidance on the diffusion updates at inference time.

Classifier Guidance / External Score Model

$$\hat{\epsilon} \leftarrow \epsilon_{\theta}(x_t) - \sqrt{1 - \bar{\alpha}_t} \nabla_{x_t} \log p_{\phi}(y|x_t)$$
$$\underline{x}_{t-1} \leftarrow \text{sample from } \mathcal{N}(\mu + s\Sigma \nabla_{x_t} \log p_{\phi}(y|x_t), \Sigma) \quad \underline{x}_{t-1} \leftarrow \sqrt{\bar{\alpha}_{t-1}} \left(\frac{x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon}$$

Guided Diffusion

- We can add guidance on the diffusion updates at inference time.

Classifier Guidance / External Score Model

$$\hat{\epsilon} \leftarrow \epsilon_{\theta}(x_t) - \sqrt{1 - \bar{\alpha}_t} \nabla_{x_t} \log p_{\phi}(y|x_t)$$
$$x_{t-1} \leftarrow \text{sample from } \mathcal{N}(\mu + s\Sigma \nabla_{x_t} \log p_{\phi}(y|x_t), \Sigma) \quad x_{t-1} \leftarrow \sqrt{\bar{\alpha}_{t-1}} \left(\frac{x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon}$$

- We also can train a conditional diffusion model.

repeat

$$(\mathbf{x}, \mathbf{c}) \sim p(\mathbf{x}, \mathbf{c})$$

$\mathbf{c} \leftarrow \emptyset$ with probability p_{uncond} \triangleright Sample data with conditioning from the dataset
 \triangleright Randomly discard conditioning to train unconditionally

$$\lambda \sim p(\lambda)$$

\triangleright Sample log SNR value

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\mathbf{z}_{\lambda} = \alpha_{\lambda} \mathbf{x} + \sigma_{\lambda} \epsilon$$

\triangleright Corrupt data to the sampled log SNR value

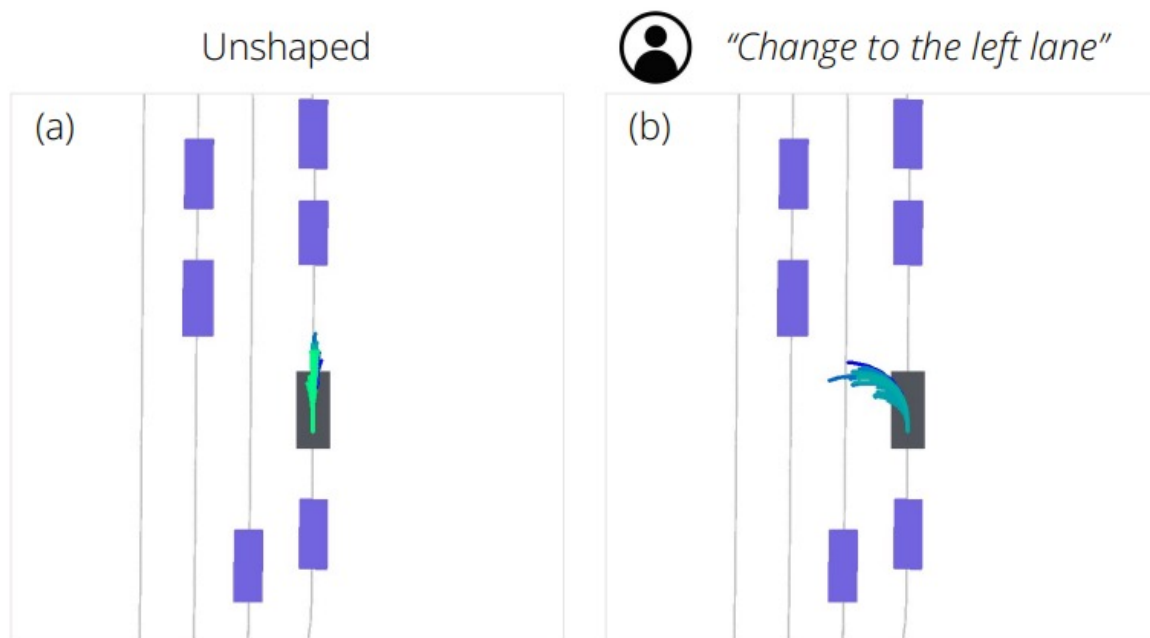
Take gradient step on $\nabla_{\theta} \|\epsilon_{\theta}(\mathbf{z}_{\lambda}, \mathbf{c}) - \epsilon\|^2$

\triangleright Optimization of denoising model

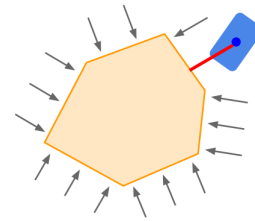
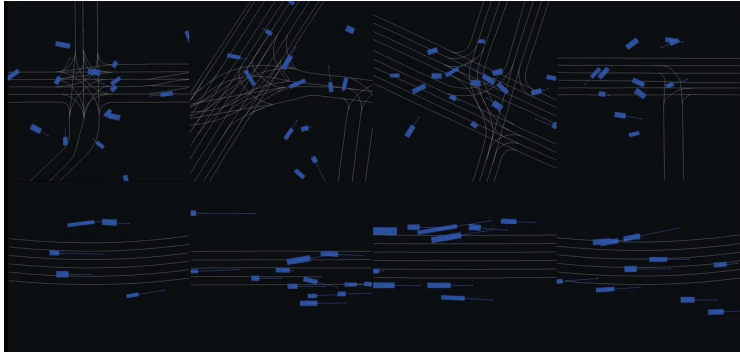
until converged

Test-Time Guidance / Adaptation

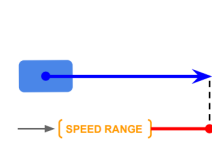
- Diffusion can be combined / guided at test time.



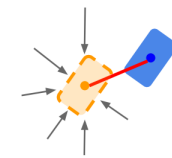
Generating Simulation Scenes



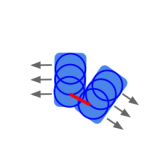
Spatial Region Constraint



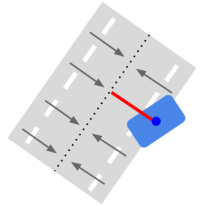
Actor Attribute Constraint



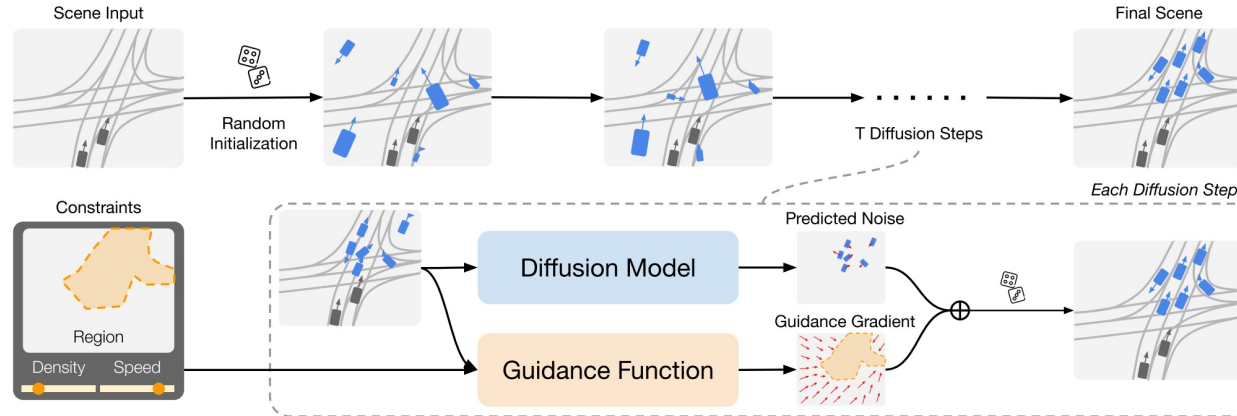
Initial Scene Constraint



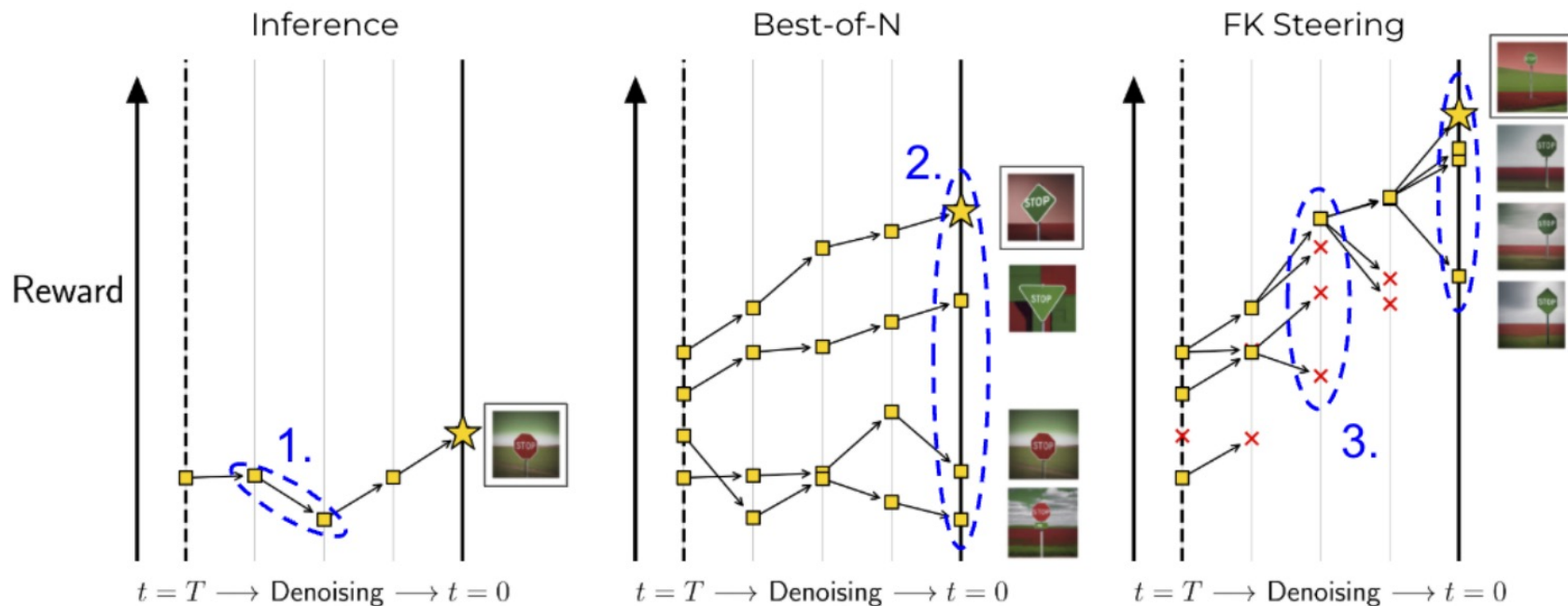
Collision Constraint



On-road Constraint



Non-Differentiable Rewards?



Prompt: “a green stop sign in a red field”

- (1) Iteratively de-noise $x_T \rightarrow x_{T-1} \rightarrow \dots \rightarrow x_0$.
- (2) Generate multiple samples (*particles*).
- (3) Resample promising particles at *intermediate* steps.

Diffusion Models for Detection

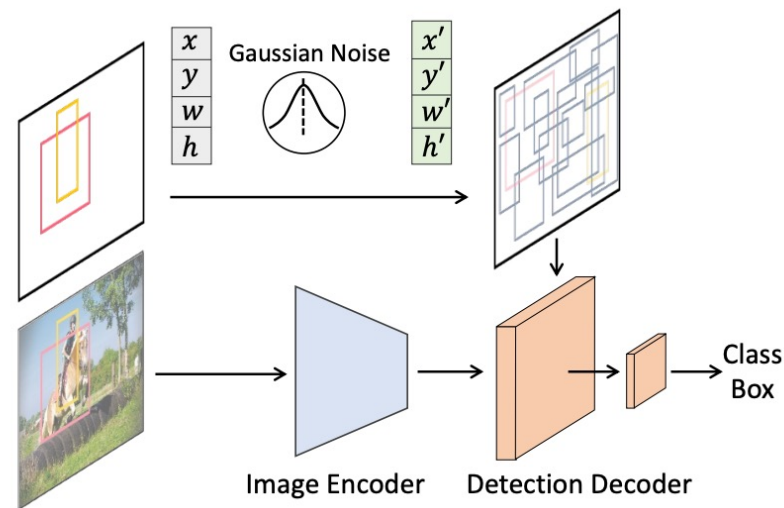
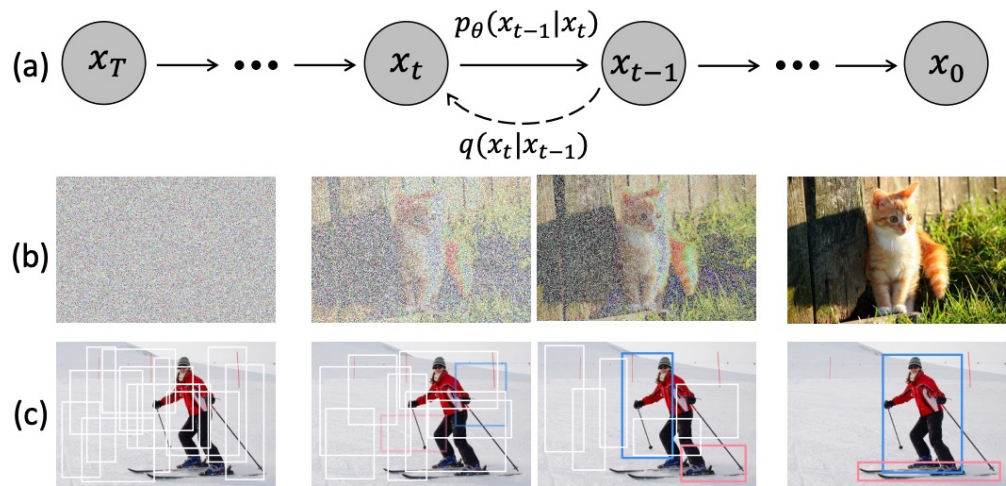
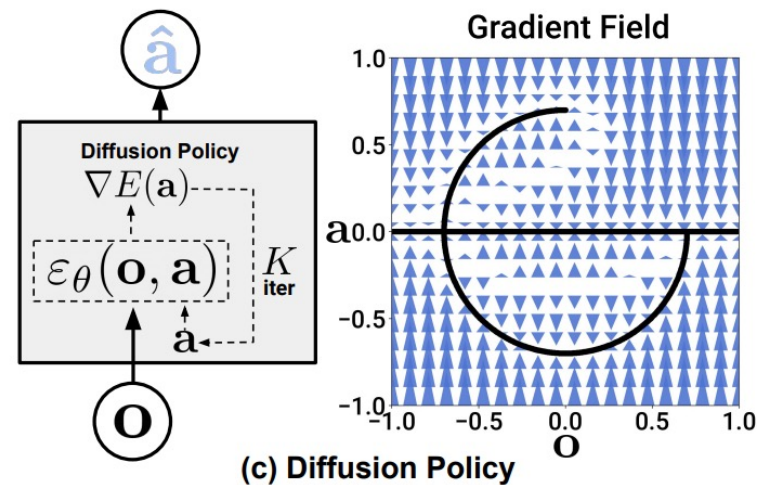
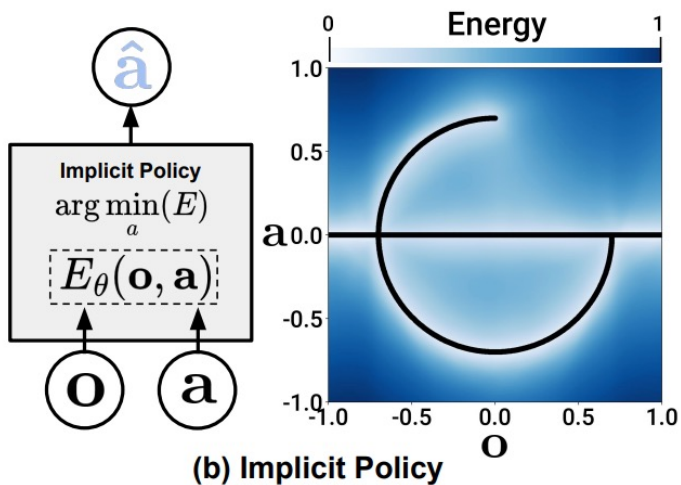
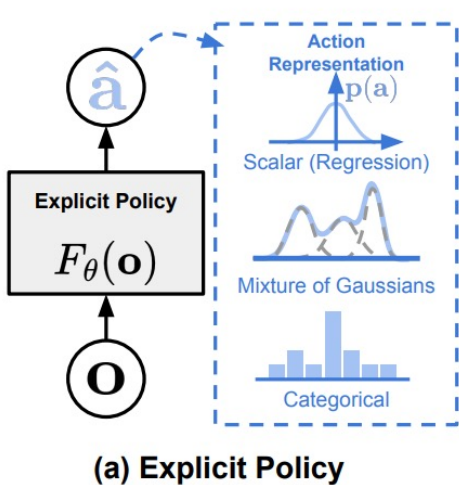
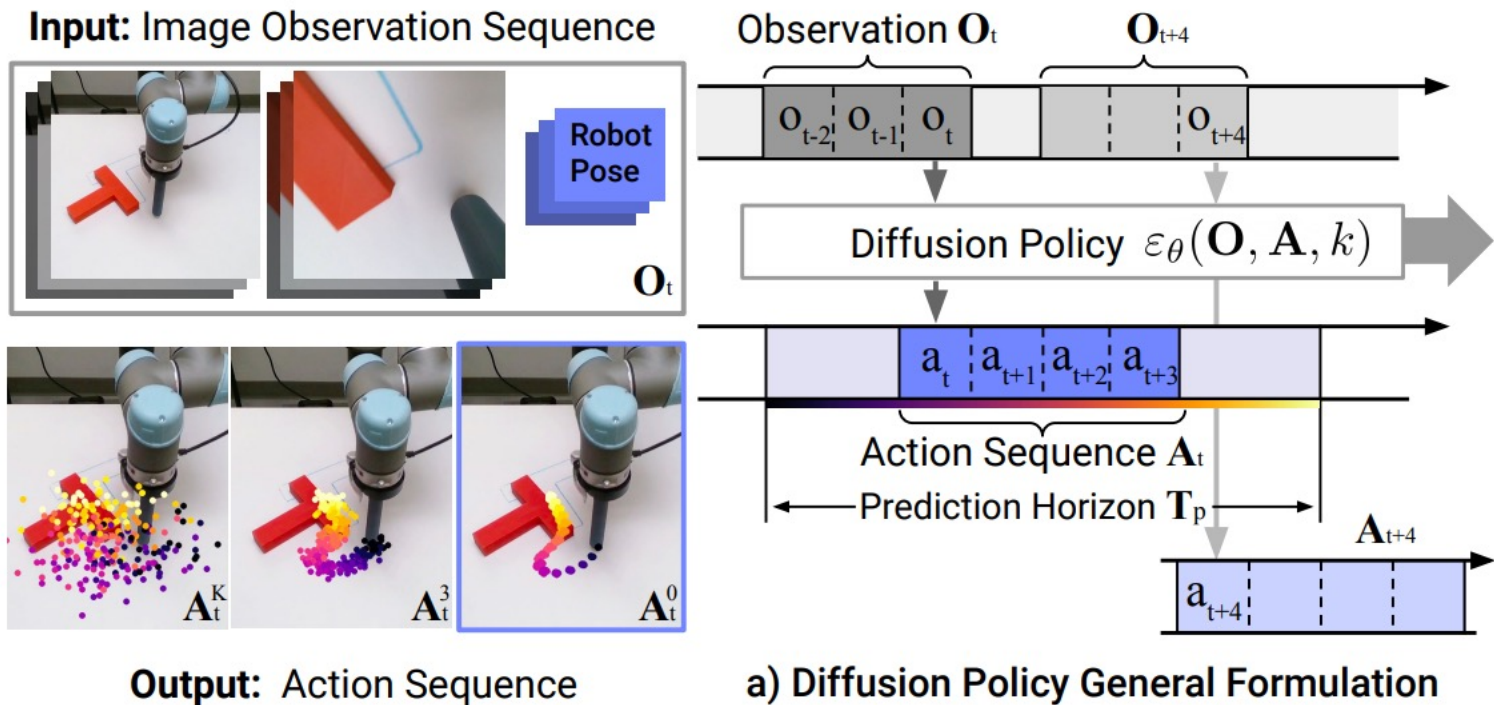


Figure 1. **Diffusion model for object detection.** (a) A diffusion model where q is the diffusion process and p_θ is the reverse process. (b) Diffusion model for image generation task. (c) We propose to formulate object detection as a denoising diffusion process from noisy boxes to object boxes.

Planning and Control



Diffusion for Planning and Control



Summary: DL for Structured Outputs

- Expanding the output dimension has limitations.

Summary: DL for Structured Outputs

- Expanding the output dimension has limitations.
- Requires us thinking about deep generative models.
 - Graphical models
 - Autoregressive
 - Energy-based
 - Normalizing flow
 - Diffusion

Summary: DL for Structured Outputs

- Expanding the output dimension has limitations.
- Requires us thinking about deep generative models.
 - Graphical models
 - Autoregressive
 - Energy-based
 - Normalizing flow
 - Diffusion
- Understand relations, pros and cons.

Summary: DL for Structured Outputs

- Expanding the output dimension has limitations.
- Requires us thinking about deep generative models.
 - Graphical models
 - Autoregressive
 - Energy-based
 - Normalizing flow
 - Diffusion
- Understand relations, pros and cons.
- Application in embodied environments.

What's Next

- Tutorial on HPC
- Next week: 3D vision, mapping