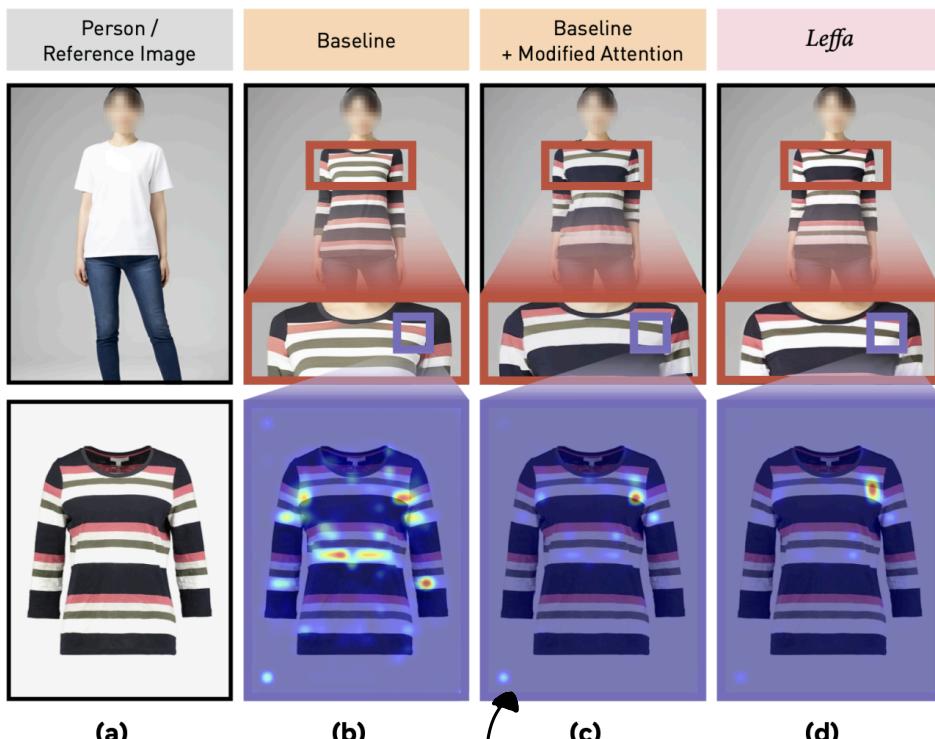


Leffy → generate a person image conditioned on reference images

↳ Learning flow fields in attention

Attention mechanism that tells the model which parts of the "reference" image should map the target



is not manually changed

→ specifically a loss function

Leffy is a new training method

↳ plugged into any diffusion-based approach

Steps: → VAE (Variational autoencoder) to compress

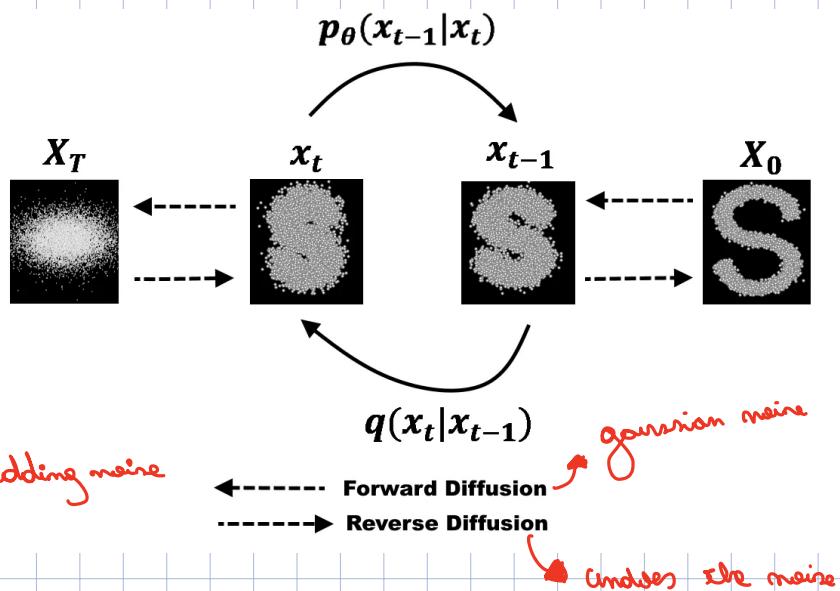
reduce computational cost

I → Encoder E → pixel space into latent space (\mathbb{Z}_0)

II

pick a timestep $t \in \{1, \dots, T\}$

Gaussian noise $\varepsilon \sim N(0, 1)$ with a mean of zero



- diffusion models learn to reverse the process of mixing an image
- the model practices "denoising"
- This process is like molecules spreading out in a gas
 - ↳ in physics, **diffusion** spreads particles from high density to low-density

Forward process

- each step only depends on the immediately previous step
- A Markov chain

$$q(x_{0:T} | x_0) = q(x_0) \prod_{t=1}^T q(x_t | x_{t-1})$$

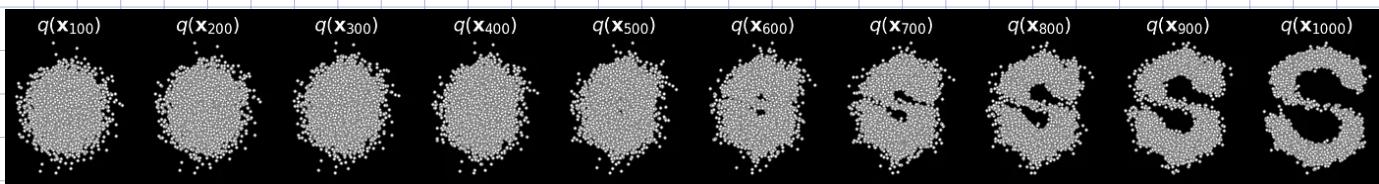
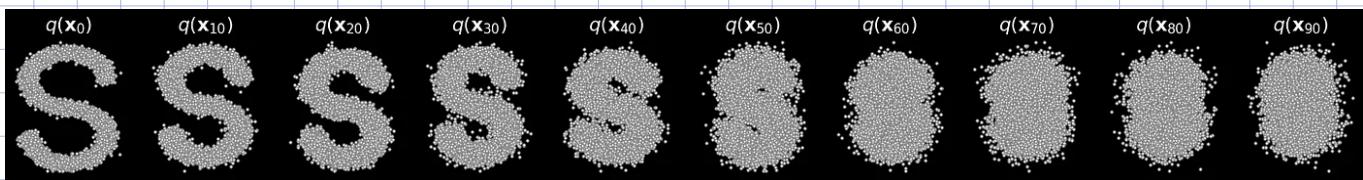
→ apply Gaussian noise in each step

- Long-term Dynamics → sampling from a distribution by using the grad
 - ↳ "tiny random walks"

Reconstruction

↪ we know how to go from $x_{t-1} \xrightarrow{t \rightarrow} x_t$
 $x_t \xrightarrow{t \rightarrow} x_{t-1}$ (train a Neural Net) U-Net
 net predicts how to
 remove noise from x_t

The neural network answers "Given the noisy image x_t , which pixels (and how) should we remove
 to get a cleaner version x_{t-1} "
 ↪ outputs the mean and variance



↑ reverse process

II

- Visual sign-on process

• Gromore mask → binary mask indicating which pixels belong to the gromore

- eliminate the gromore

- Add noise

Done pose image is w/o body layout

- Pose transfer → "These are the person's appearance details (face, cloth) to preserve"

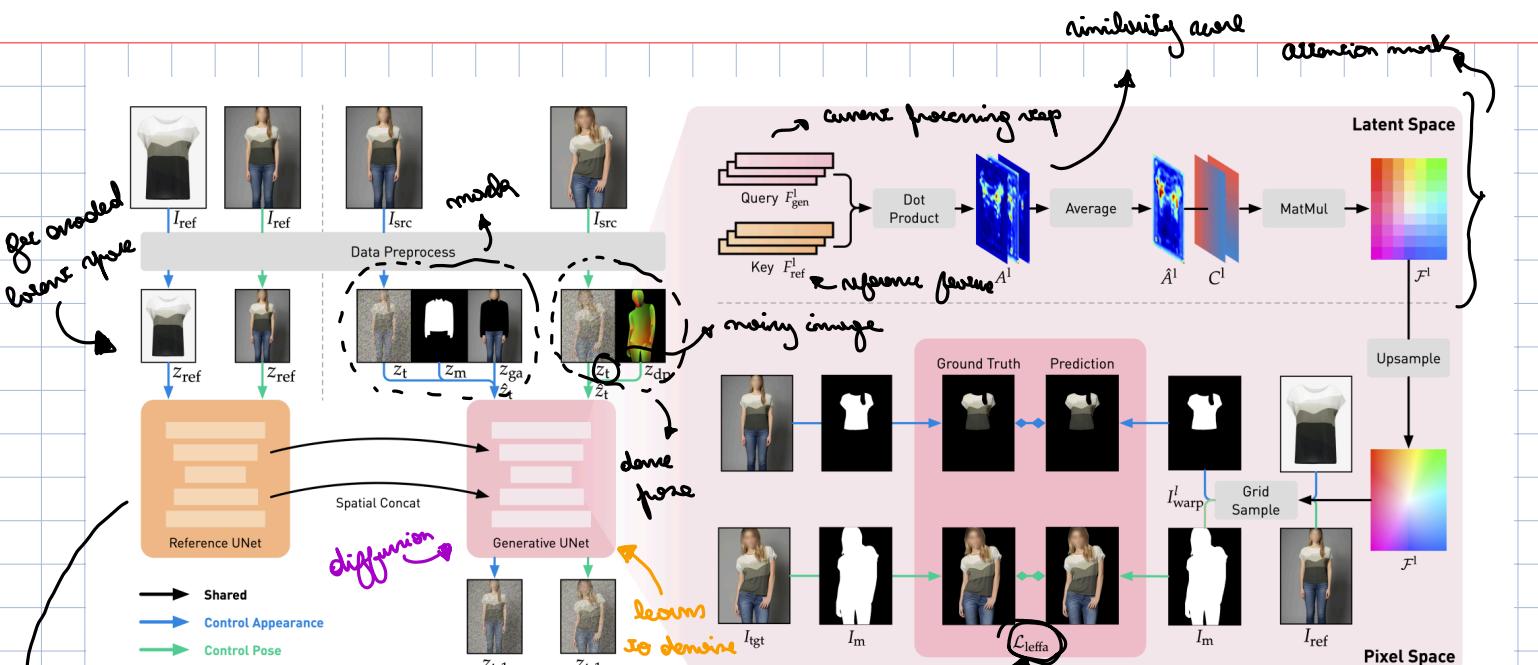


Figure 3 An overview of our *Leffa* training pipeline for controllable person image generation. The left is our diffusion-based baseline; the right is our *Leffa* loss. Note that I_{src} and I_{tgt} are the same image during training.

network's internal "description" of the image content
↳ capture "feature map"

specialized loss function in charge of formalizing patterns

$$Q \in \mathbb{R}^{m \times m} \quad (\text{the current feature map})$$

$$K \in \mathbb{R}^{m \times m} \quad (\text{reference map})$$

- $(Q K)^T \rightarrow$ similarity score for "which reference apes are similar to my current latent"

• Softmax (prob)

$$z_t$$

1. You have a noisy latent z_t that you feed into the Generative UNet.

2. Within the UNet, each layer outputs a feature map—call the l -th layer's output F^l_{gen} .

3. Meanwhile, the Reference UNet produces its own feature map F^l_{ref} describing the garment/person in the reference image.

4. Cross-attention at layer l takes:

- F^l_{gen} as "Query,"
- F^l_{ref} as "Key/Value."
- It aligns them to inject the correct garment/pose details into the Generative UNet's feature map.

At the end of the Generative UNet, you get a predicted denoised latent z_{t-1} . That z_{t-1} is still in the VAE's latent space; it's just less noisy than z_t . When you finally reach z_0 (fully denoised), you decode it back to an RGB image.

a less quality refinement

Leffa loss \rightarrow allocates detail distortion through learning flow

fields in attention

distortion arises when the target query in all layers fails to attend to the corresponding region

$\xrightarrow{\text{target}}$ $\xrightarrow{\text{Self-Att. loss}}$ guides the target query to specially focus on the correct region
↳ based on the attention map

(!!) The problem now is that the attention layer may be looking at the wrong place

flow field \Rightarrow maps an output field to a reference pixel

↳ compare to ground truth \rightarrow learn more to diverse

$$L_{\text{finegr}} = L_{\text{diffusion}} + \underbrace{L_{\text{Self-Att.}}}_{\xrightarrow{\text{less weight of } I_{\text{left}}}}$$

Self-Consistency for Human Parsing

image pixel from human body \longrightarrow semantic category

target \rightarrow improve the model performance and generalization by progressively refining the many levels during training

SCHP \rightarrow purification strategy
↳ self-consistency procedure
↳ self-consistency mechanism
↳ cover the many training level via a model and label mutually promoting process.

Methodology

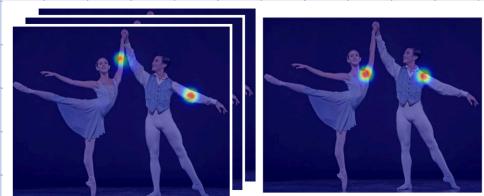
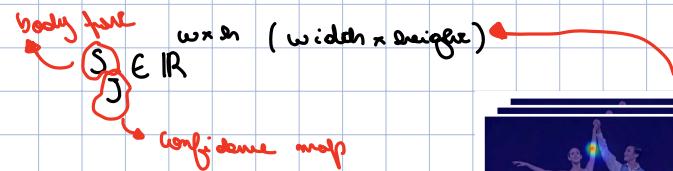
CE2P $\xrightarrow{\text{boundary maps}}$
↳ forming branch
↳ edge branch
↳ fusion branch
↳ combined to produce a refined human parsing pred

Openpose (Merced)

1. 2D confidence maps (S) \wedge 2D Vector fields L (PAFs)

then ...

$$S = (s_1, s_2, \dots, s_J)$$

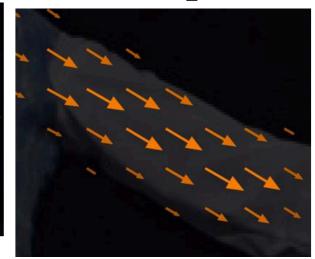


$$L_c = (L_{c1}, L_{c2}, \dots, L_{cJ})$$

c one per vector field

$$L_c \in \mathbb{R}^{w \times h \times J^2}$$

creates a 2D vector



finally

$$S_j \cap L_c$$

Join by greedy inference

→
queue



Network Architecture

- The iterative prediction architecture refines predictions

over stages $t \in \{1, \dots, T\}$ with intermediate supervision at each stage

Simultaneous Detection and Association

- feature maps of F $\xrightarrow{\text{process}} L^1 = \phi^1(\tilde{x})$ Stage 1 (PAFS)
- so at each stage $L^t = \phi^t(F, L^{t-1}), \forall (2 \leq t \leq T_p)$ total TAPs
 ↳ respect previous L^{t-1}

Same for confidence map

$$S^t = p^t(F, L^{t_p}, S^{t-1}) \quad \text{refers to the CNN for inference at stage } t \quad \forall t_p < t \leq t_p + t_c$$

→ refinement of detection over stages



Stage 1

Stage 2

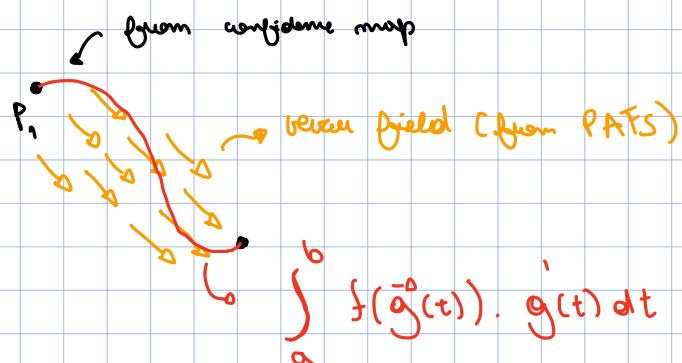
Stage 3

Confidence maps for Pose detection

- generate individual confidence maps $S_{j,k}^t$ for each person j,k .
- (II) Decide which parts belong together to form a person's pose
 PAFs → 2D vector field for each limb
- preserve both location and orientation information across the region of support of the limb
 ground truth



Pose Association and Pose Assembly



Sum up how well PAFs vectors align with the direction

DonePose → mapping all human pixels of RGB to the 3D-Surface
 ↳ coco dataset ground-truth colored Img 2 Surface

Diffusion model:

- slowly destroy structure in a data distribution
- then
 - learn a reverse diffusion process that restores structure in data

Forward diffusion process

- each step of the forward diffusion process is defined as

$$q(x_t | x_{t-1}) = N(x_t, \underbrace{\sqrt{1 - \beta_t} x_{t-1}}_u; \underbrace{\beta_t I}_{\sigma^2})$$

↳ output of forward process

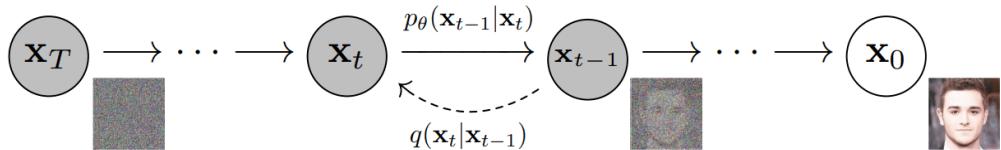
Schedule → range from 0 to 1

↳ OpenAI in 2021 proposed cosine schedule

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

Reverse diffusion process

↳ Convert noise into an image



- Diffusion model predicts the entire noise to be removed in a given timestep

- diffusion model $\xrightarrow{\text{predicts}}$ noise present in the current image

then ...

scheduler (reverse process) removes the predicted noise from the image

Architecture

↳ U-Net with embeddings

→ each step in reverse process adds embeddings with information

about current timestep and prompt

Resnet-block (residual network)

allow deeper networks to learn residual functions rather than direct mapping

learn the difference between input and desired output

If direct mapping were simply to keep input unchanged then

$$\text{residual } \underline{\underline{F(x) = 0}}$$

diffusion models: used to process and refine feature representations at various resolutions.

Downsample Block

→ reduces spatial resolution (height x width) of feature maps

+ receptive field \Rightarrow integrate context over larger areas of img

Diffusion models: reduce compute cost and refine feature representations

→ use of positional encoding

Self-Attention Block → allow position in the feature map