

# **MoGlow: Probabilistic and controllable motion synthesis using normalising flows (2019)**

Gustav Eje Henter\* Simon Alexanderson\* Jonas Beskow

Division of Speech, Music and Hearing, KTH Royal Institute of Technology, Stockholm, Sweden

- Prediction models suffer from regression to mean pose (“mean collapse”) and artefacts like foot sliding.
- Normalising Flow (Glow) to map pose to latent space
- Sample a natural motion given a condition (e.g., walk along a trajectory)
- LSTM to generate motion sequence of arbitrary length and ensure temporal consistency
- Autoregression to ensure temporal consistency

# Background: Normalising Flow

- Is a generative model like VAE, GAN etc.
- Estimates data distribution  $P_{\text{data}}(\mathbf{x})$
- Allows exact probability computation unlike VAEs
- Easier to train than VAEs (may suffer “posterior collapse”, i.e., strong decoders yield models where  $z$  has little impact), and GANs

### Assumptions and definitions:

- probability of observing data point  $\mathbf{x}$  (unknown):  $P_{\text{data}}(\mathbf{x})$
- samples drawn from  $p_{\text{data}}$  are independent (like images, CIFAR, MNIST)

$$P_{\text{data}}(\mathbf{X}) = P_{\text{data}}(\mathbf{x}_1) P_{\text{data}}(\mathbf{x}_1) \dots P_{\text{data}}(\mathbf{x}_N) = \prod_i P_{\text{data}}(\mathbf{x}_i)$$

- we want to create a model from which we can sample:  $P_{\text{model}}(\mathbf{x})$
- in order to build as much “realistic” model as we can we want to have:  $P_{\text{model}}(\mathbf{x}) = P_{\text{data}}(\mathbf{x})$
- this allows us to define target metric which tells us how far from true distribution we are:

$$D(P_{\text{data}}(\mathbf{x}), P_{\text{model}}(\mathbf{x})) = 0 \text{ if } P_{\text{model}}(\mathbf{x}) = P_{\text{data}}(\mathbf{x})$$

## Assumptions and definitions:

- The **natural choice** for **D** is Kullback Leibler divergence (*this is a critical point for all derivations below*):

$$\text{KL} (P_{\text{data}} (\mathbf{x}) , P_{\text{model}} (\mathbf{x})) = \sum_i P_{\text{data}} (\mathbf{x}) \log \left( \frac{P_{\text{data}} (\mathbf{x})}{P_{\text{model}} (\mathbf{x})} \right)$$

- Our generator  $P_{\text{model}}$  will be defined by some neural network, let's put explicitly its dependence on some parameters (neural network weights)

$$P_{\text{model}} (\mathbf{x}) = P_{\text{model}} (\mathbf{x}; \theta)$$

- We want to minimize KL w.r.t model parameters, hence we need to compute the gradients

$$\nabla_{\theta} \text{KL} = \nabla_{\theta} \sum_i P_{\text{data}} (\mathbf{x}) \log \left( \frac{P_{\text{data}} (\mathbf{x})}{P_{\text{model}} (\mathbf{x}; \theta)} \right) = -\nabla_{\theta} \sum_i P_{\text{data}} (\mathbf{x}) \log (P_{\text{model}} (\mathbf{x}; \theta))$$

- The samples are explicitly sampled from data distribution (N goes to infinity):

$$\nabla_{\theta} \text{KL} = -\frac{1}{N} \nabla_{\theta} \sum_{\mathbf{x}_i \sim P_{\text{data}} (\mathbf{x})} \log (P_{\text{model}} (\mathbf{x}_i; \theta)) = -\frac{1}{N} \nabla_{\theta} \log \prod_{\mathbf{x}_i \sim P_{\text{data}} (\mathbf{x})} P_{\text{model}} (\mathbf{x}_i; \theta)$$

- The samples are explicitly sampled from data distribution (N goes to infinity):

$$\nabla_{\theta} \text{KL} = -\frac{1}{N} \nabla_{\theta} \sum_{\mathbf{x}_i \sim P_{\text{data}}(\mathbf{x})} \log (P_{\text{model}}(\mathbf{x}_i; \theta)) = -\frac{1}{N} \nabla_{\theta} \log \prod_{\mathbf{x}_i \sim P_{\text{data}}(\mathbf{x})} P_{\text{model}}(\mathbf{x}_i; \theta)$$

- **Recall assumption that:**  $P_{\text{data}}(\mathbf{X}) = P_{\text{data}}(\mathbf{x}_1) P_{\text{data}}(\mathbf{x}_1) \dots P_{\text{data}}(\mathbf{x}_N) = \prod_i P_{\text{data}}(\mathbf{x}_i)$
- This allows us to write similar expression for model (but X are drawn from “real” data)

$$\nabla_{\theta} \text{KL} = -\frac{1}{N} \nabla_{\theta} \log P_{\text{model}}(\mathbf{X}; \theta)$$

- Hence by minimizing **KL divergence** we **maximize log-likelihood** of observed data, both can be used to obtain same result. However with one it maybe do in easier way.

$$\text{minimize} \quad \text{KL} \qquad \text{maximize} \quad \log P_{\text{model}}(\mathbf{X}; \theta)$$

Ignoring the gradient operator, want to minimise KL == minimise model likelihood (log is preferred due to numerical stability)

# Background: Normalising Flow

- To estimate  $P_{\text{model}}$ , we start from a “simple” latent distribution (characterised by latent random variable  $\mathbf{Z}$ ), and find a series of invertible functions that will “recover”  $P_{\text{model}}$  (characterized by latent RV  $\mathbf{X}$ ), i.e. (we choose a latent distribution of uniform Gaussian in this example):

$$\mathbf{Z} \sim \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$$

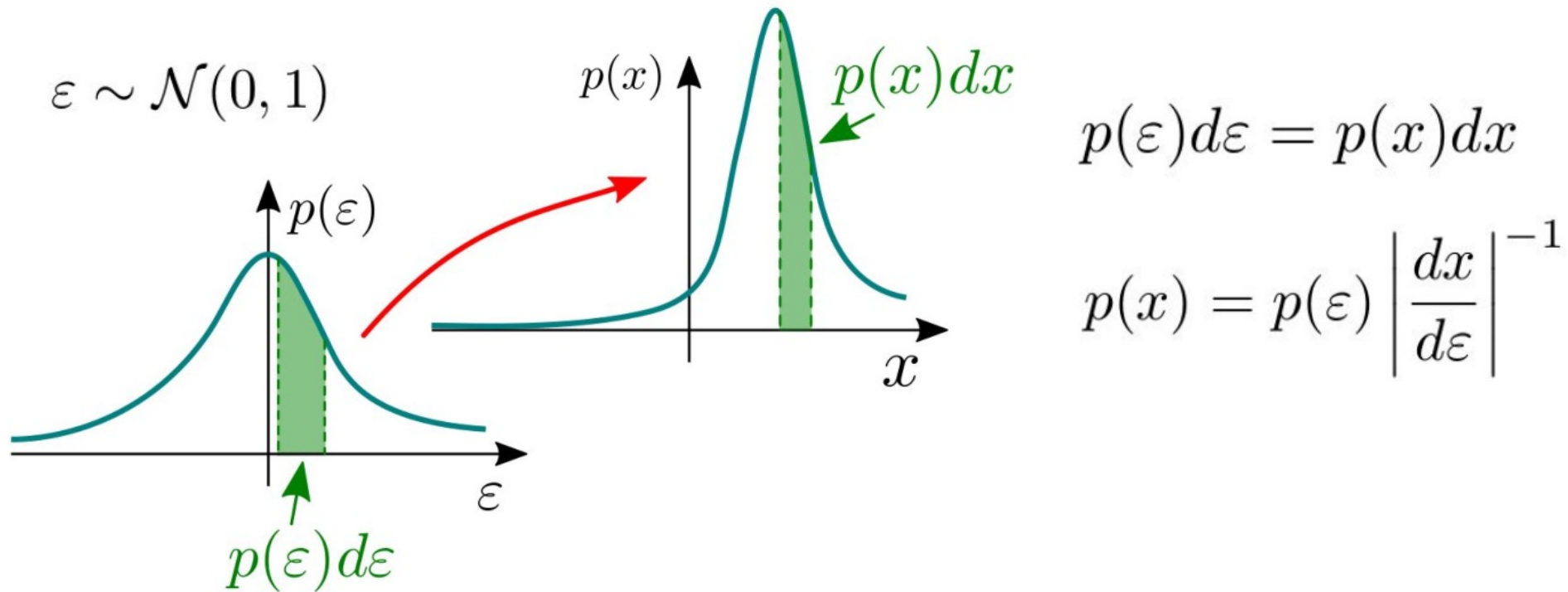
$$\mathbf{x} = \mathbf{f}(\mathbf{z}) = \mathbf{f}_1(\mathbf{f}_2(\dots \mathbf{f}_N(\mathbf{z})))$$

$$\mathbf{z} = \mathbf{z}_N \xrightarrow{\mathbf{f}_N} \mathbf{z}_{N-1} \xrightarrow{\mathbf{f}_{N-1}} \dots \xrightarrow{\mathbf{f}_2} \mathbf{z}_1 \xrightarrow{\mathbf{f}_1} \mathbf{z}_0 = \mathbf{x}$$

$$\mathbf{z}_n(\mathbf{x}) = \mathbf{f}_n^{-1} \circ \dots \circ \mathbf{f}_1^{-1}(\mathbf{x}).$$

# Background: Normalising Flow

- When we transform from  $\mathbf{z}_i$  to  $\mathbf{z}_{i+1}$  (assume  $\mathbf{z}_0 = \mathbf{x}$  and  $\mathbf{z}_n = \mathbf{z}$ ), we want the densities preserved. 1D example:

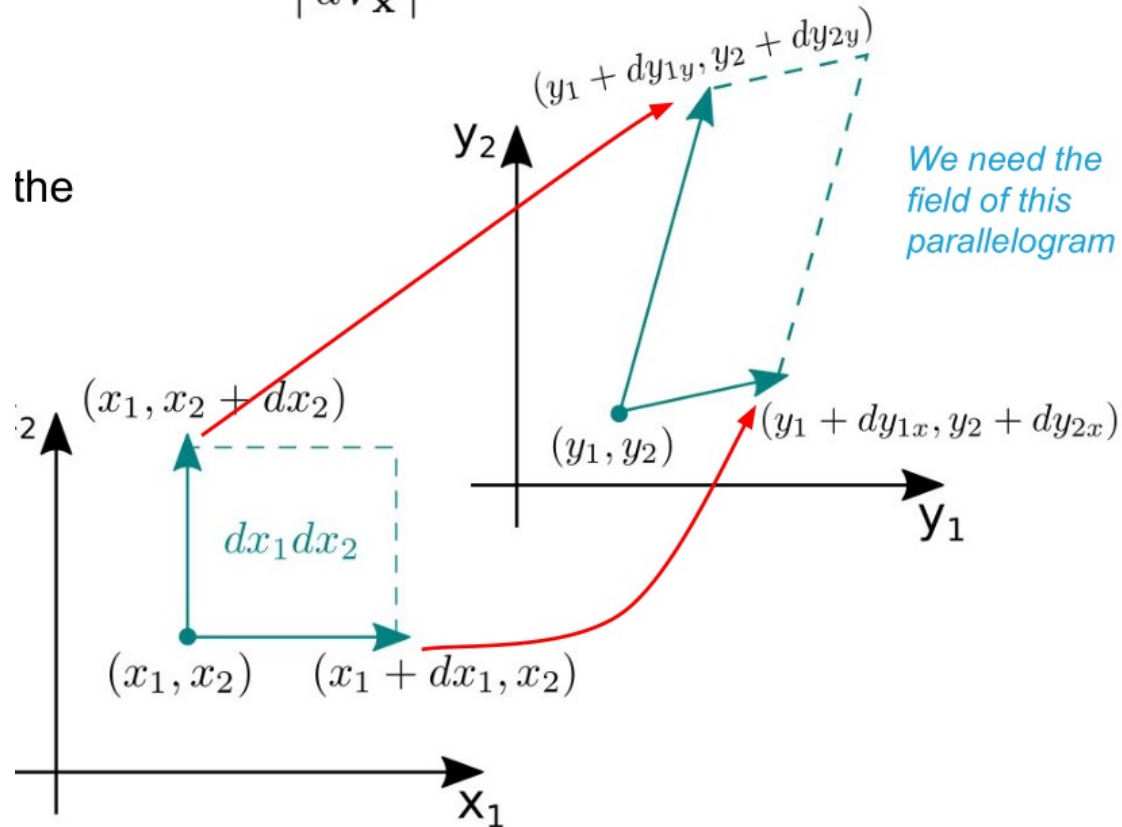




# Background: Normalising Flow

- For  $\text{dim} > 1$  we want volume preservation
- 2D case:

$$p(\mathbf{y}) = p(\mathbf{x}) \left| \frac{dV_{\mathbf{y}}}{dV_{\mathbf{x}}} \right|^{-1}$$



We derive:

$$p(\mathbf{y}) = p(\mathbf{x}) |\det \mathbf{J}|^{-1}$$

# Background: Normalising Flow

- General case for 1 step of transformation:

$$q'(\mathbf{z}') = q(\mathbf{z}) \left| \det \frac{\partial f(\mathbf{z}^{(n)})}{\partial \mathbf{z}^{(n)}} \right|^{-1}$$

- Through chaining:

$$\begin{array}{l} \mathbf{z}^{(0)} \sim p(\mathbf{z}) \\ \mathbf{z}^{(1)} = f^{(1)}(\mathbf{z}^{(0)}) \\ \mathbf{z}^{(2)} = f^{(2)}(\mathbf{z}^{(1)}) \\ \dots \\ \mathbf{z}^{(n)} = f^{(n)}(\mathbf{z}^{(n-1)}) \end{array} \quad \longrightarrow \quad q^{(n)}(\mathbf{z}^{(n)}) = q(\mathbf{z}^{(0)}) \prod_{i=1}^n \left| \det \frac{\partial \mathbf{z}^{(i)}}{\partial \mathbf{z}^{(i-1)}} \right|^{-1}$$

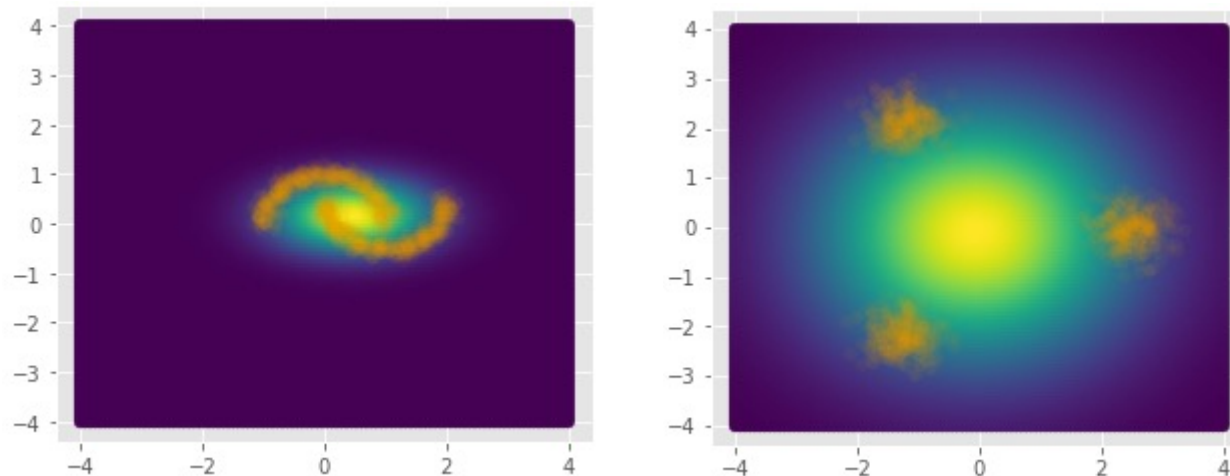
- Using log probability (or log likelihood if summed over minibatch):

$$\ln p_{\theta}(\mathbf{x}) = \ln p_{\mathcal{N}}(\mathbf{z}_N(\mathbf{x})) + \sum_{n=1}^N \ln \left| \det \frac{\partial \mathbf{z}_n(\mathbf{x})}{\partial \mathbf{z}_{n-1}} \right|$$

- For the choice of invertible function  $f_i$ , we may use a simple affine transformation e.g.,

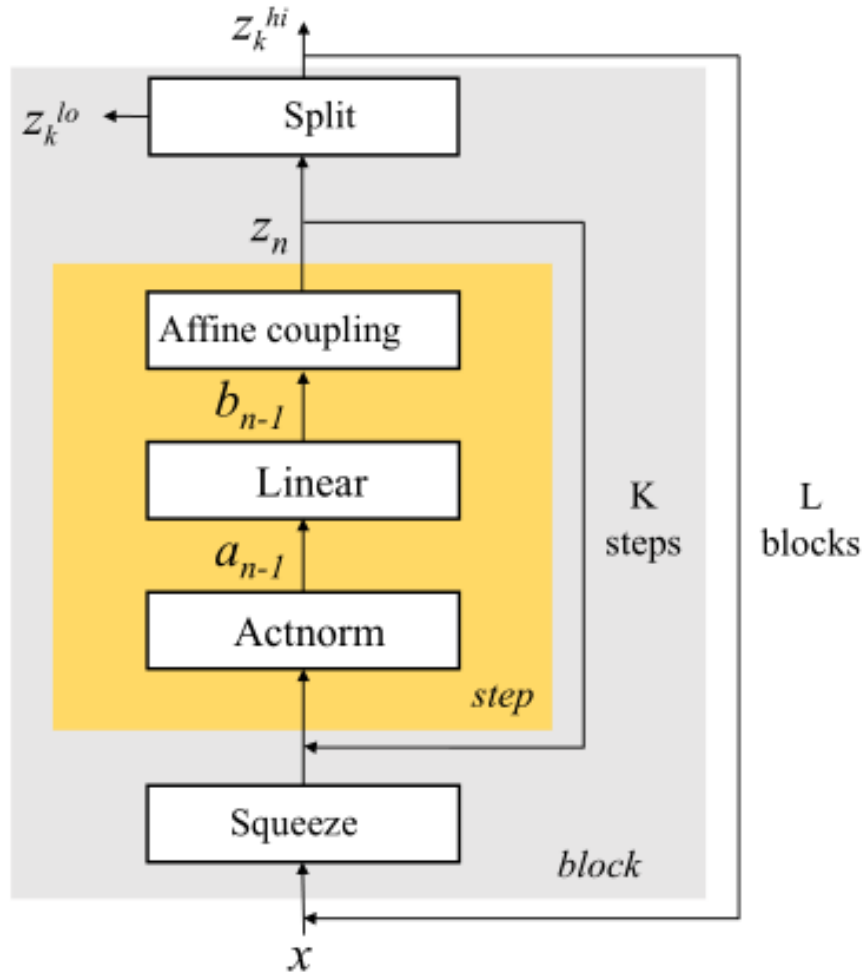
$$x = \mu + \varepsilon\sigma$$
$$\varepsilon = (x - \mu) / \sigma$$

- But then it struggles to fit (i.e., transform from the unit Gaussian to) a complex, nonlinear data distribution:



- Several good candidates proposed: MADE, RealNVP, Glow
- Glow is used in this project.

- One glow step (inverse of fi)

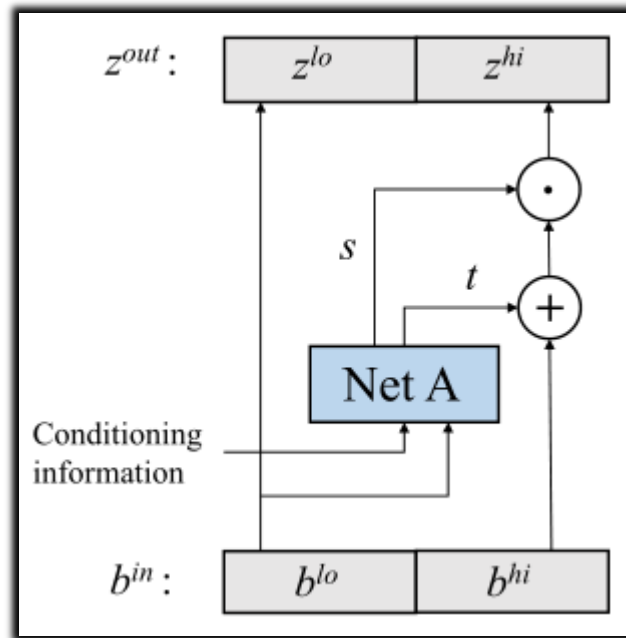
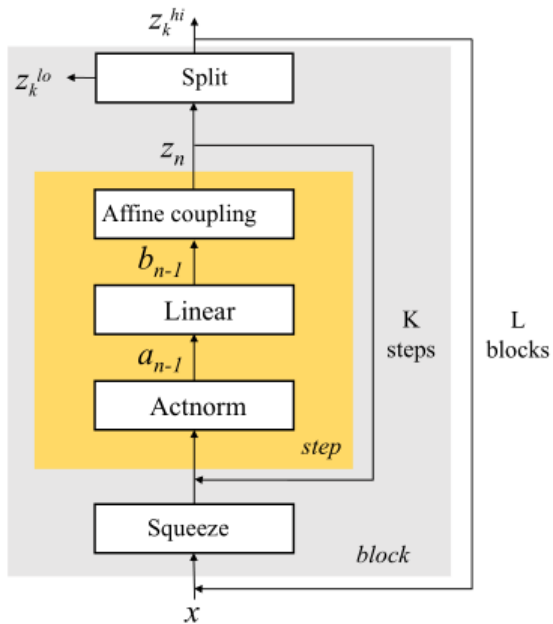


**Actnorm:**  $\mathbf{a}_{t,n} = \mathbf{s}_n \odot \mathbf{z}_{t,n} + \mathbf{t}_n$   
Kind of like a reversible batchnorm

**Linear transformation:**  $\mathbf{b}_{t,n} = \mathbf{W}_n \mathbf{a}_{t,n}$   
Weight matrix is LU-decomposed where one of the diagonals is constrained to contain just ones to ease computation of  $|\det J|$ .  
Note that by keeping the diagonals of the other matrix to have value  $\neq 0$ ,  $\mathbf{W}$  is kept invertible.  
Akin to permutation layer in RealNVP.

# Background: Normalising Flow – Glow

- One glow step (inverse of fi)

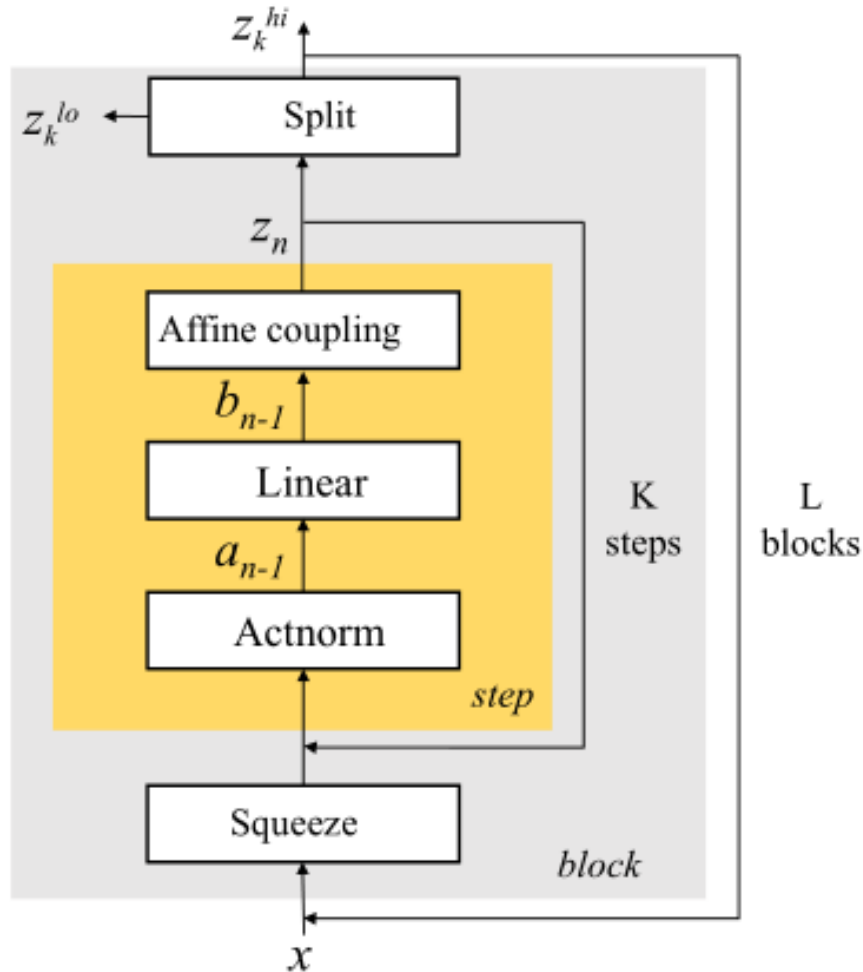


## Affine coupling layer:

Half of the layer goes through unchanged ( $b_{lo}$ )  
Other half ( $b_{hi}$ ) appended to the transformed  
(by Net A or  $A_n$ ; possibly non-invertible)  $b_{lo}$ .

Function	Reverse Function
$\mathbf{x}_a, \mathbf{x}_b = \text{split}(\mathbf{x})$ $(\log \mathbf{s}, \mathbf{t}) = \text{NN}(\mathbf{x}_b)$ $\mathbf{s} = \exp(\log \mathbf{s})$ $\mathbf{y}_a = \mathbf{s} \odot \mathbf{x}_a + \mathbf{t}$ $\mathbf{y}_b = \mathbf{x}_b$ $\mathbf{y} = \text{concat}(\mathbf{y}_a, \mathbf{y}_b)$	$\mathbf{y}_a, \mathbf{y}_b = \text{split}(\mathbf{y})$ $(\log \mathbf{s}, \mathbf{t}) = \text{NN}(\mathbf{y}_b)$ $\mathbf{s} = \exp(\log \mathbf{s})$ $\mathbf{x}_a = (\mathbf{y}_a - \mathbf{t}) / \mathbf{s}$ $\mathbf{x}_b = \mathbf{y}_b$ $\mathbf{x} = \text{concat}(\mathbf{x}_a, \mathbf{x}_b)$

- One glow step (inverse of fi)



### Hierarchical Decomposition:

Process iterated for  $L$  blocks,  
Which can increase the power to model long-range dependencies.

- Consider that in a vanilla NF model, one can sample a random point in  $\mathbf{z}$  and reconstruct a likely  $\mathbf{x}$  in the data distribution.
- But this is not very helpful if you want to generate a sequence of poses that is not totally random.
- That is why we condition the NF during training and testing using a control signal, and it is basically appended (author uses the term “squeezed”) in the affine coupling layers.
- First we present the fixed-length model which is a convolutional NF model that accepts a fixed-length sequence:  $\mathbf{x}_t \in \mathbb{R}^D, \mathbf{c}_t \in \mathbb{R}^C$
- $\mathbf{x}_t$  is pose at time  $t$ , and  $\mathbf{c}_t$  is the control signal at time  $t$  or the “global” control parameter.



# Model: Fixed-length models with control

$$\begin{aligned} [z_{t,n+1}^{\text{lo}}, z_{t,n+1}^{\text{hi}}] &= [b_{t,n}^{\text{lo}}, (b_{t,n}^{\text{hi}} + t'_{t,n}) \odot s'_{t,n}] \\ [s'_{t,n}, t'_{t,n}] &= A_n(b_{t,n}^{\text{lo}}) \end{aligned} \quad \longrightarrow \quad [s'_{t,n}, t'_{t,n}] = A_n(b_{t-1:t+1,n}^{\text{lo}}).$$

Add past frames info  
(convolutional fixed-length model)

Add conditioning

$$[s'_{t,n}, t'_{t,n}] = A_n(b_{t-1:t+1,n}^{\text{lo}}, c_{t-1:t+1,n})$$

- The official model that is proposed is, as opposed to the fixed-length sequence model from earlier, an autoregressive model that can generate sequence of arbitrary length.
- Because the fixed-length model contains a time-sequence (which may be long),  $L = 4$  blocks,  $K = 32$  each block.
- The autoregressive model only has  $L = 1$  block,  $K = 16$  steps each.
- LSTM used within affine coupling layer

# Model: MoGlow

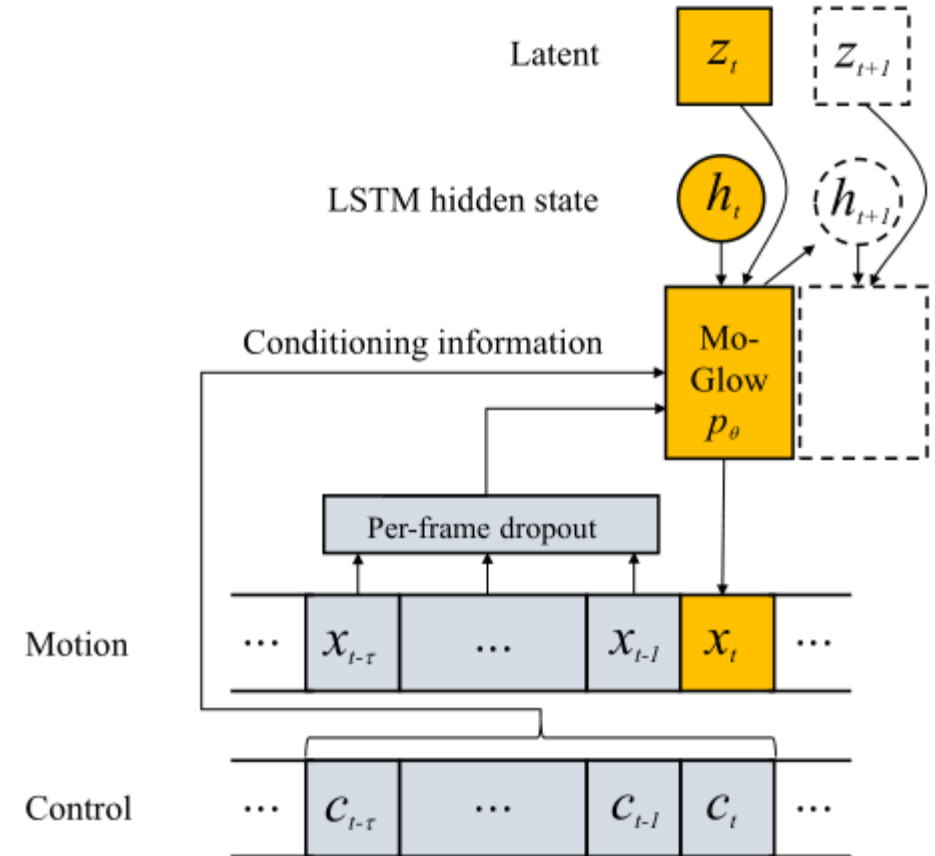
- Model is formulated as follows:

$$p_{\theta}(\mathbf{x} | \mathbf{c}) = p(\mathbf{x}_{1:\tau} | \mathbf{c}_{1:\tau}) \cdot \prod_{t=\tau+1}^T p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t-\tau:t-1}, \mathbf{c}_{t-\tau:t}, \mathbf{h}_t) \quad (16)$$

$$\mathbf{h}_{t+1} = \mathbf{g}_{\theta}(\mathbf{x}_{t-\tau:t-1}, \mathbf{c}_{t-\tau:t}, \mathbf{h}_t) \quad (17)$$

$$[s'_{t,n}, t'_{t,n}] = A_n(b_{t,n}^{\text{lo}}, \mathbf{x}_{t-\tau:t-1}, \mathbf{c}_{t-\tau:t,n}, \mathbf{h}_t), \quad (18)$$

- Notice that the hidden state  $\mathbf{h}_t$  of LSTM is also added as a conditioning var (along with control parameter), and allows for longer time dependence than  $\tau$ .

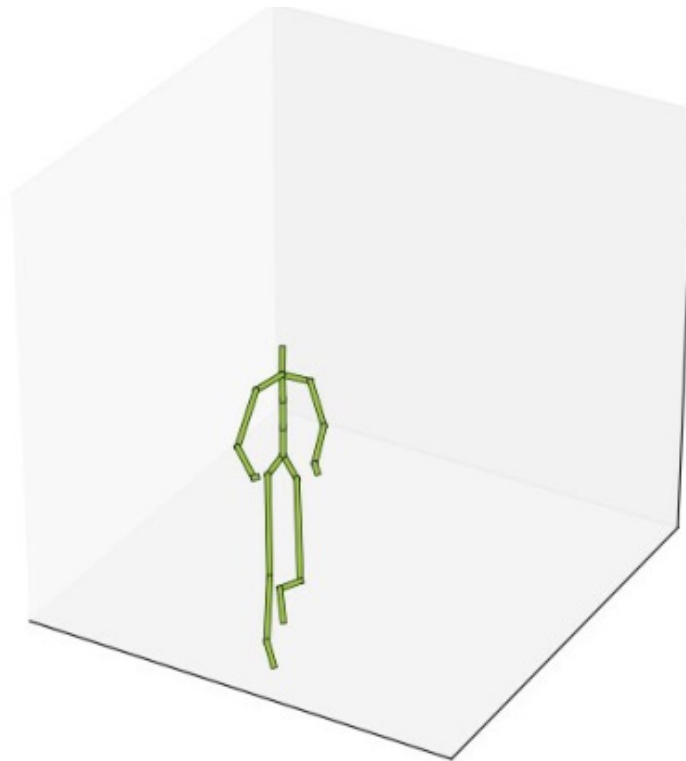


# Implementation: Locomotion Trials



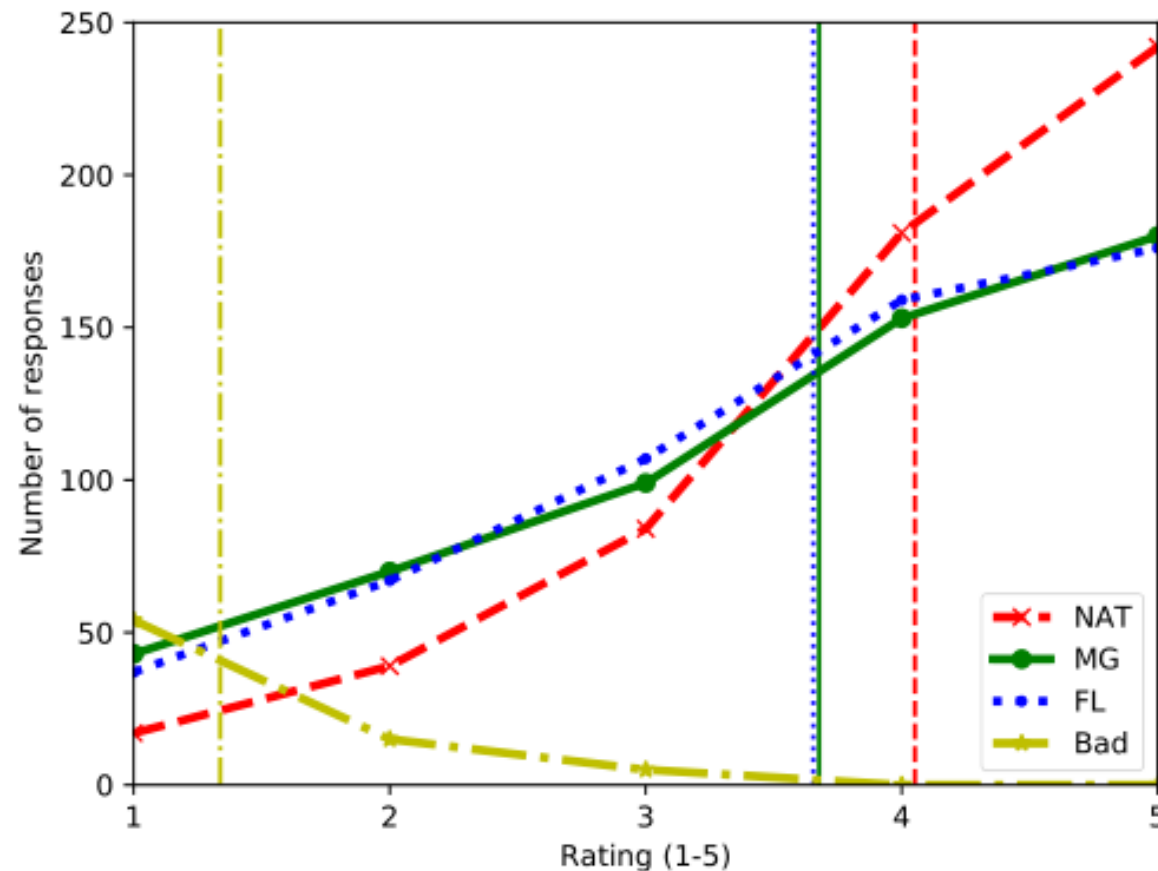
- Trained on **locomotion trials** from CMU and HDM05 motion capture data sets (part of AMASS).
- The input signals are local 3D positions of 21 joints, and the control signals are the displacement of the root at each frame.
- Data augmented by mirroring spatially and temporally.
- Fixed length model trained 20k steps, autoregressive 80k steps with Adam optimiser.

- During inference provided the path at each frame a random  $\mathbf{z}$  is sampled and with the conditionals (control signal & hidden states) a random walking motion is generated autoregressively.
- Synthetic motion examples: <https://youtu.be/IYhJnDBWyeo>



# Results

- For quantitative evaluation, participants were asked to give a subjective evaluation of how “natural” the generated/real motions were on a scale: 1-5.



- Official open-source implementation available:  
<https://github.com/simonalexanderson/StyleGestures>
- Updated research: *Style-Controllable Speech-Driven Gesture Synthesis Using Normalising Flows* came out in 2020