**PDF = "probability per unit volume"**

- 1D : RV w/ PDF $p_X(x)$ $\Rightarrow$ $\Pr(X \in [x, x+\Delta x]) \approx p_X(x)\,\Delta x$

- nD : " $\Rightarrow$ $\Pr(X \in \text{tiny box around } x) \approx p_X(x) \cdot (\text{volume of that box})$

**Jacobian determinant = local volume scaling**

- For a smooth map $f: \mathbb{R}^d \to \mathbb{R}^d$, near a point $x$, the map is approximately linear:
  $\quad\hookrightarrow f(x+\delta) \approx f(x) + J_f(x)\,\delta \qquad \Rightarrow \qquad \underline{(\text{new volume}) \approx |\det J_f(x)| \cdot (\text{old volume})}$ <span style="color:purple">$\Delta$ of variables formula</span>

**change-of-var. gives the flow likelihood :** Define $z = f_\theta(q; h)$, choose simple PDF for $z$ : $p_z(z) = \mathcal{N}(0, I)$

$\hookrightarrow$ Take a tiny region $A$ around $q$ $\Rightarrow$ its img in $z$-space $= B = f(A)$

$\quad\hookrightarrow$ $\because$ this is reparameterizing the same outcomes : $\underline{\Pr(q \in A \mid h) = \Pr(z \in B \mid h)}$

Approximating both sides
using $\Pr(X \in \text{tiny box around } x) \approx p_X(x) \cdot (\text{volume of that box})$

$: p_\theta(q \mid h)\,\text{Vol}(A) \approx p_z(z \mid h)\,\text{Vol}(B)$

since $z = f(q; h)$ : $\text{Vol}(B) \approx |\det J_f(q; h)|\,\text{Vol}(A)$

$\downarrow$

$\therefore p_\theta(q \mid h)\,\text{Vol}(A) \approx p_z(f(q; h))\,|\det J_f(q; h)|\,\text{Vol}(A)$

$\quad\hookrightarrow$ Taking logs (base e) & using Gaussian $p_z$ : <span style="color:purple">$\overset{\text{Flow training objective}}{}$</span> $-\log p_\theta(q \mid h) = \underbrace{\frac{1}{2}\|z\|^2}_{\substack{\text{make } z \\ \text{look Gaussian}}} - \underbrace{\log|\det J_f(q; h)|}_{\text{volume correction}} + \text{const.}$

⚠ **Generic neural nets can't model $f_\theta(q; h)$**

$\hookrightarrow$ Problems : 1. Computing $\det(J_f(q; h))$ is inefficient
$\qquad\qquad\quad$ 2. $f$ is not invertible $\qquad$ <span style="color:purple">e.g. if using MLP</span>

$\downarrow$

Need special architecture for flow models $\Rightarrow$ e.g. coupling layers!

**Coupling layers :**



a) Forward mapping
$h_1$ — $\phi[h_1]$ — $h_1'$
$h_2$ — $g[h_2, \phi[h_1]]$ — $h_2'$
Input — Output

b) Inverse mapping
$h_1$ — $\phi[h_1]$ — $h_1'$
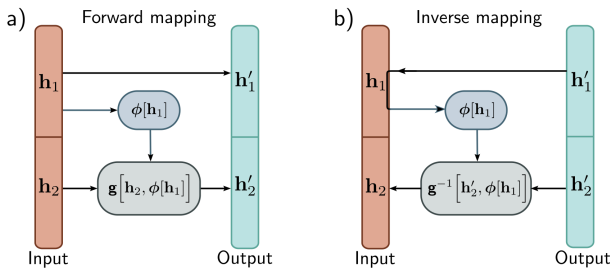$h_2$ — $g^{-1}[h_2', \phi[h_1]]$ — $h_2'$
Input — Output

**Figure 16.6** Coupling flows. a) The input (orange vector) is divided into $h_1$ and $h_2$. The first part $h_1'$ of the output (cyan vector) is a copy of $h_1$. The output $h_2'$ is created by applying an invertible transformation $g[\bullet, \phi]$ to $h_2$, where the parameters $\phi$ are themselves a (not necessarily invertible) function of $h_1$. b) In the inverse mapping, $h_1 = h_1'$. This allows us to calculate the parameters $\phi[h_1]$ and then apply the inverse $g^{-1}[h_2', \phi]$ to retrieve $h_2$.

- Additive coupling is invertible & $\det(J)$ efficient to compute, but it can't expand / contract volume $\rightarrow$ $\hookrightarrow \therefore$ <u>limited expressiveness</u>

$\Rightarrow$ e.g. Let $x \in \mathbb{R}^d$ $\Rightarrow$ split into $x = (x_a, x_b)$, $x_a \in \mathbb{R}^{d_a}$, $x_b \in \mathbb{R}^{d_b}$

- Additive coupling : Define $y = g(x; c)$ $\begin{cases} y_a = x_a \\ y_b = x_b + t_\theta(x_a, c) \end{cases}$
  $\hookrightarrow$ where $t_\theta$ = any NN outputting a vector $\in \mathbb{R}^{d_b}$

  $>$ Invertible : ✓ can recover $x$ from $y$
  $\quad x_a = y_a$
  $\quad x_b = y_b - t_\theta(y_a, c)$ $\Big\}$ $g^{-1}$ exists & is easy to compute

  $>$ Jacobian w/ easy det : ✓
  $\quad J = \dfrac{\partial y}{\partial x} = \begin{bmatrix} \partial y_a/\partial x_a & \partial y_a/\partial x_b \\ \partial y_b/\partial x_a & \partial y_b/\partial x_b \end{bmatrix} = \begin{bmatrix} I & 0 \\ \partial t/\partial x_a & I \end{bmatrix}$

  $\det(\text{triangular mtx}) = \prod \text{diag. entries}$
  $\therefore \det(J) = \det(I) \cdot \det(I) = 1$

<span style="color:red">★ FrEIA AllInOneBlock combines affine coupling, permutation, ... etc. common normalizing flow ops.</span>

- **Affine coupling:** $y = g(x; c)$ $\begin{cases} y_a = x_a \\ y_b = x_b \odot \exp(S_\theta(x_a, c)) + t_\theta(x_a, c) \end{cases}$ — *elementwise mult.*  $(y = ax + b \Rightarrow$ "affine")

  $\hookrightarrow S_\theta = $ scale network $\Rightarrow S_\theta(x_a, c) = (s_1, \ldots, s_{d_b})$

  - **Invertible:** ✓  $x_a = y_a$

    $x_b = (y_b - t_\theta(x_a, c)) \odot \frac{1}{\exp(S_\theta(x_a, c))} = (y_b - t_\theta(x_a, c)) \odot \exp(-S_\theta(x_a, c))$

  - **Jacobian w/ easy det:**  $(y_b)_K = (x_b)_K \cdot \exp(S_K(x_a, c)) + t_K(x_a, c)$

    when differentiating w.r.t. $x_b$: $\frac{\partial (y_b)_K}{\partial (x_b)_K} = \exp(S_K(x_a, c))$, for $j \neq K$: $\frac{\partial (y_b)_K}{\partial (x_b)_j} = 0$

    $\therefore \frac{\partial y_b}{\partial x_b} = \text{diag}\left(\exp(s_1), \ldots, \exp(s_{d_b})\right)$

    $\Rightarrow J = \begin{bmatrix} I & 0 \\ * & \text{diag}(\exp(s)) \end{bmatrix}$ $\Rightarrow \det(J) = \det(I) \cdot \det(\text{diag}(\exp(s))) = \prod_{K=1}^{d_b} \exp(s_K) = \exp\left(\sum_{K=1}^{d_b} s_K\right)$

    $\Rightarrow \log|\det J| = \sum_{K=1}^{d_b} s_K$  ✓ *easy to compute* ☺

⚠ **But half the vars don't change?!** ↰

↓

In one coupling layer, $x_a$ stays exactly the same $(y_a = x_a)$ $\Rightarrow$ only $x_b$ changes

↳ **Fix:** after each coupling layer, <u>permute</u> s.t. different subset of $x$ becomes $x_a$ next layer

$\Rightarrow$ over multiple layers, every coord. gets to be in the "changed group" multiple times

**Conditioning preserves invertibility:**  [Recall: coupling layer requires generating $s$ & $t$ from NN]

$\hookrightarrow y_b = x_b \odot \exp(s) + t$

- **Unconditional case:** $(s, t) = NN_\theta(x_a)$

- **Conditional case:** $(s, t) = NN_\theta(x_a, c)$

$\Rightarrow$ $c$ is like a hyperparam: we only need $x \mapsto y = g_c(x)$ to be invertible for each fixed $c$

$\hookrightarrow$ $c$ is given @ both forward & reverse time

↓

e.g. affine coupling w/ conditioning: $y_a = x_a$, $(s, t) = NN_\theta(x_a, c)$, $y_b = x_b \odot \exp(s) + t$

↓

Showing this is invertible:  1. $x_a = y_a$
2. $(s, t) = NN_\theta(y_a, c)$
3. Solve for $x_b$: $y_b = x_b \odot \exp(s) + t$

**Intuition for applying flow models to our reachability problem:**

$\hookrightarrow$ Goal post-training: for given $H = c$, sampling different $z \sim \mathcal{N}(0, I)$ gives different valid $Q$

↓

Flow model learns a geometry-aware warp s.t. the solution distr. becomes simple in latent space

**Concrete cINN example w/ our setup:**  e.g. choose a split for one layer: $x_a = (x, y)$  $x_b = (\cos\theta, \sin\theta)$

$\hookrightarrow$ e.g. affine coupling layer: $y_a = x_a$, $(s, t) = NN_\theta(x_a, H)$  where $s, t \in \mathbb{R}^2$, $y_b = x_b \odot \exp(s) + t$

aka. given position + target, you rescale/shift angle features

$\hookrightarrow$ next layer can permute to $x_a = (\cos\theta, \sin\theta)$, $x_b = (x, y)$ $\Rightarrow$ angle features control how
you rescale/shift position features

⇓

stack enough of these & we get a rich coupling btwn parts of $Q$ relative to $H$

clamp exists in FrEIA ∵ affine coupling uses $\exp(s) \rightarrow$ if $s$ gets big, $\exp(s)$ blows up

    ↳ Implementations squash $s$ into a bounded range

## why cINN is potentially good for our problem:

- Multi-modal : different modes correspond to different $z$-values
  
      ↓
  
  Training is max. likelihood : penalizes missing prob. mass (dropping modes)

## Dimensions in our simplified setup : (maps to code implementation)

- Input to flow : $x = Q_{feat} \in \mathbb{R}^4$

- condition : $c = H \in \mathbb{R}^2$

- Output (forward) : $z \in \mathbb{R}^4$

- sampling (reverse) : sample $z \sim \mathcal{N}(0, I) \Rightarrow$ compute $x = f^{-1}(z; c)$

- convert $x = (x, y, \underline{\cos\theta, \sin\theta})$ back to $Q = (x, y, \theta)$ w/ $\theta = \text{atan2}(\sin\theta, \cos\theta)$

                          | normalize to unit len. before atan2 for stability