# Normalising Flow Models (Part 2)

Hao Dong

Peking University

# Normalising Flow Models (Part 2)

Hao Dong

Peking University

## So far

- Learning via **maximum likelihood** over the dataset D

$$\max_{\theta} \log p(D; \theta) = \sum_{x \in D} \log \pi \left( G_{\theta}^{-1}(x) \right) + \log \left| \det \left( \frac{\partial G_{\theta}^{-1}(\boldsymbol{x})}{\partial \boldsymbol{x}} \right) \right|$$

<span style="color:red">inverted function</span>     <span style="color:red">determinant of Jacobian</span>

- What we need?

  - 1）Prior $z \sim \pi(z)$ easy to sample

  - 2）Invertible transformations

  - 3）Determinants of Jacobian Efficient to compute

# Reference slides

- Hung-yi Li. Flow-based Generative Model
- Stanford "Deep Generative Models". Normalising Flow Models

- Coupling layer based normalising flow models
  - Coupling layer
  - NICE
  - Real NVP
  - Glow
- Autoregressive models as flow models
  - MAF
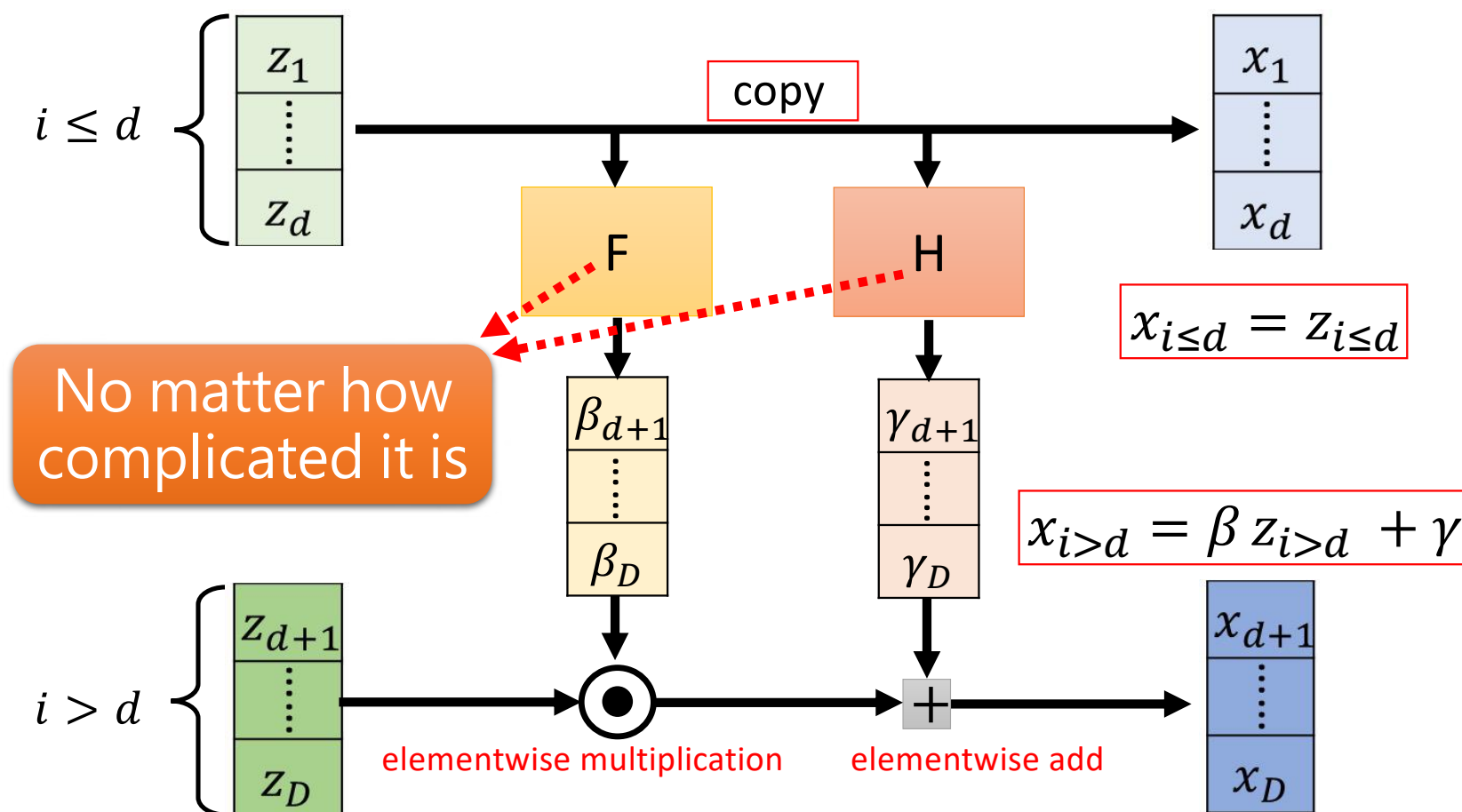  - IAF
  - Parallel Wavenet

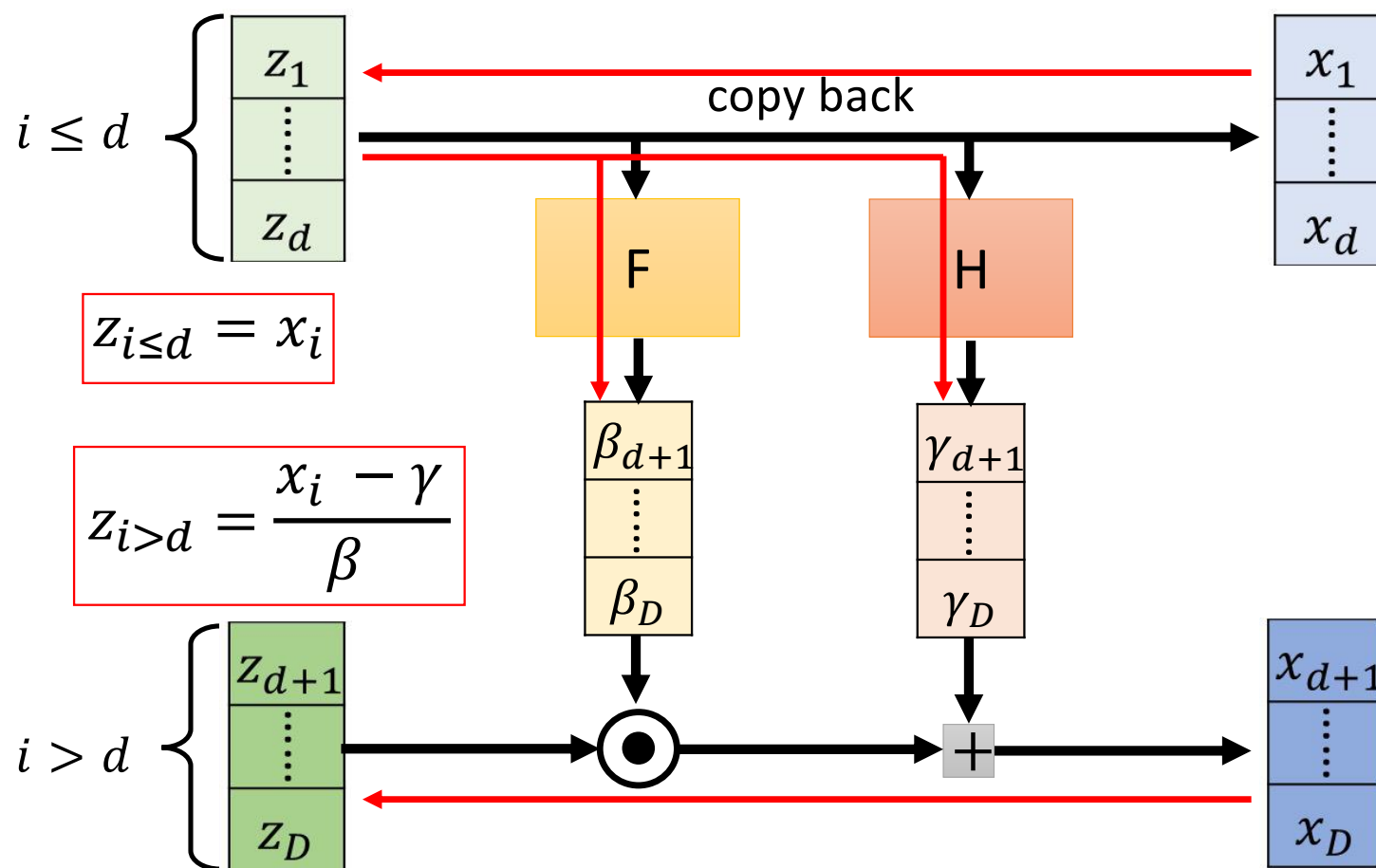- Coupling layer based normalising flow models
  - Coupling layer
  - NICE
  - Real NVP
  - Glow
- Autoregressive models as flow models
  - MAF
  - IAF
  - Parallel Wavenet

# Coupling Layer

NICE
https://arxiv.org/abs/1410.8516
Real NVP
https://arxiv.org/abs/1605.08803



$i \leq d$ : $z_1, \dots, z_d$

copy

$x_1, \dots, x_d$

F

H

No matter how complicated it is

$x_{i \leq d} = z_{i \leq d}$

$\beta_{d+1}, \dots, \beta_D$

$\gamma_{d+1}, \dots, \gamma_D$

$x_{i > d} = \beta \, z_{i > d} + \gamma$

$i > d$ : $z_{d+1}, \dots, z_D$

$\odot$

elementwise multiplication

$+$

elementwise add

$x_{d+1}, \dots, x_D$

# Coupling Layer

$i \le d$

$z_1$

$\vdots$

$z_d$

copy back

$x_1$

$\vdots$

$x_d$

F

H

$z_{i \le d} = x_i$

$z_{i > d} = \dfrac{x_i - \gamma}{\beta}$

$\beta_{d+1}$

$\vdots$

$\beta_D$

$\gamma_{d+1}$

$\vdots$

$\gamma_D$

$i > d$

$z_{d+1}$

$\vdots$
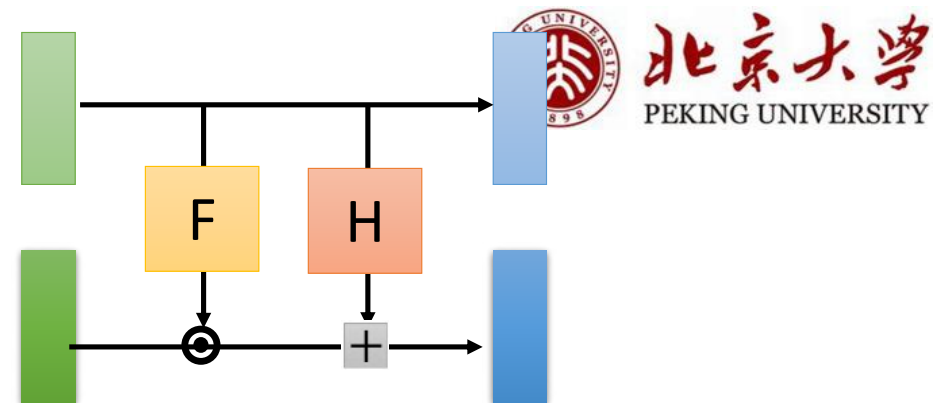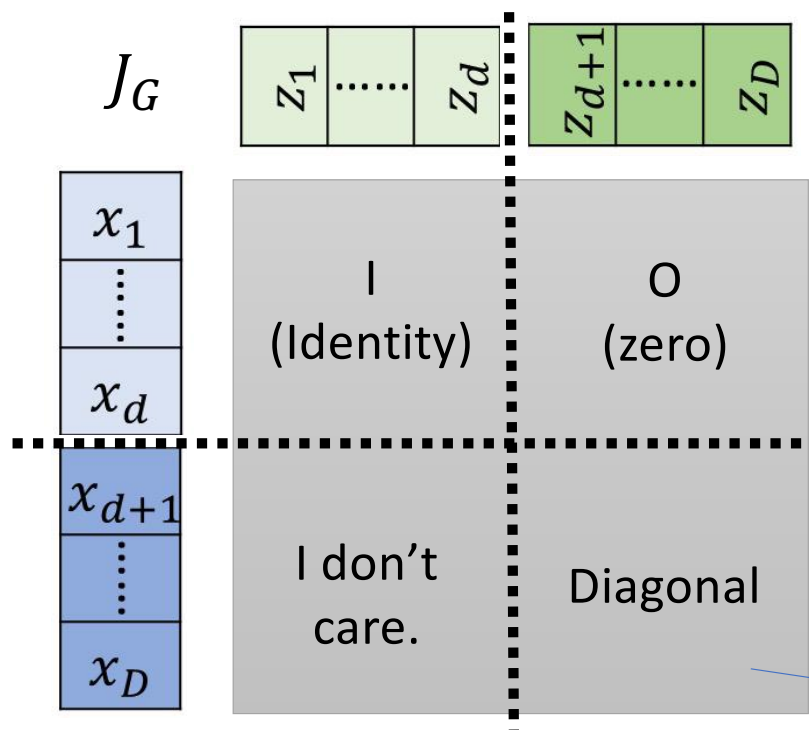
$z_D$

$\odot$

$+$

$x_{d+1}$

$\vdots$

$x_D$

8

# Coupling Layer

- Learning via **maximum likelihood** over the dataset D

$$\max_{\theta} log\, p(D; \theta) = \sum_{x \in D} log\, \pi\left(G_{\theta}^{-1}(x)\right) + log\left|det\left(\frac{\partial G_{\theta}^{-1}(x)}{\partial x}\right)\right|$$

<span style="color:red">Jacobian</span>
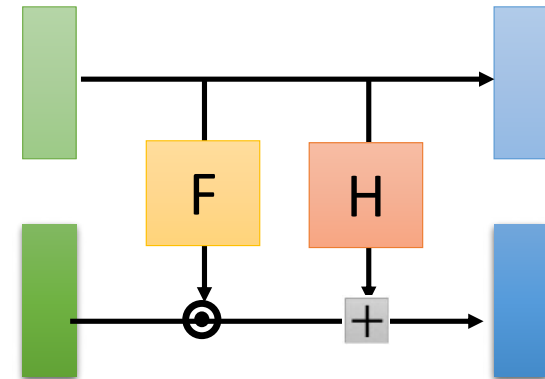
# Coupling Layer



$J_G$



$$det(J_G)$$

$$= \frac{\partial x_{d+1}}{\partial z_{d+1}} \frac{\partial x_{d+2}}{\partial z_{d+2}} \cdots \frac{\partial x_D}{\partial z_D}$$

$$= \beta_{d+1} \beta_{d+2} \cdots \beta_D$$

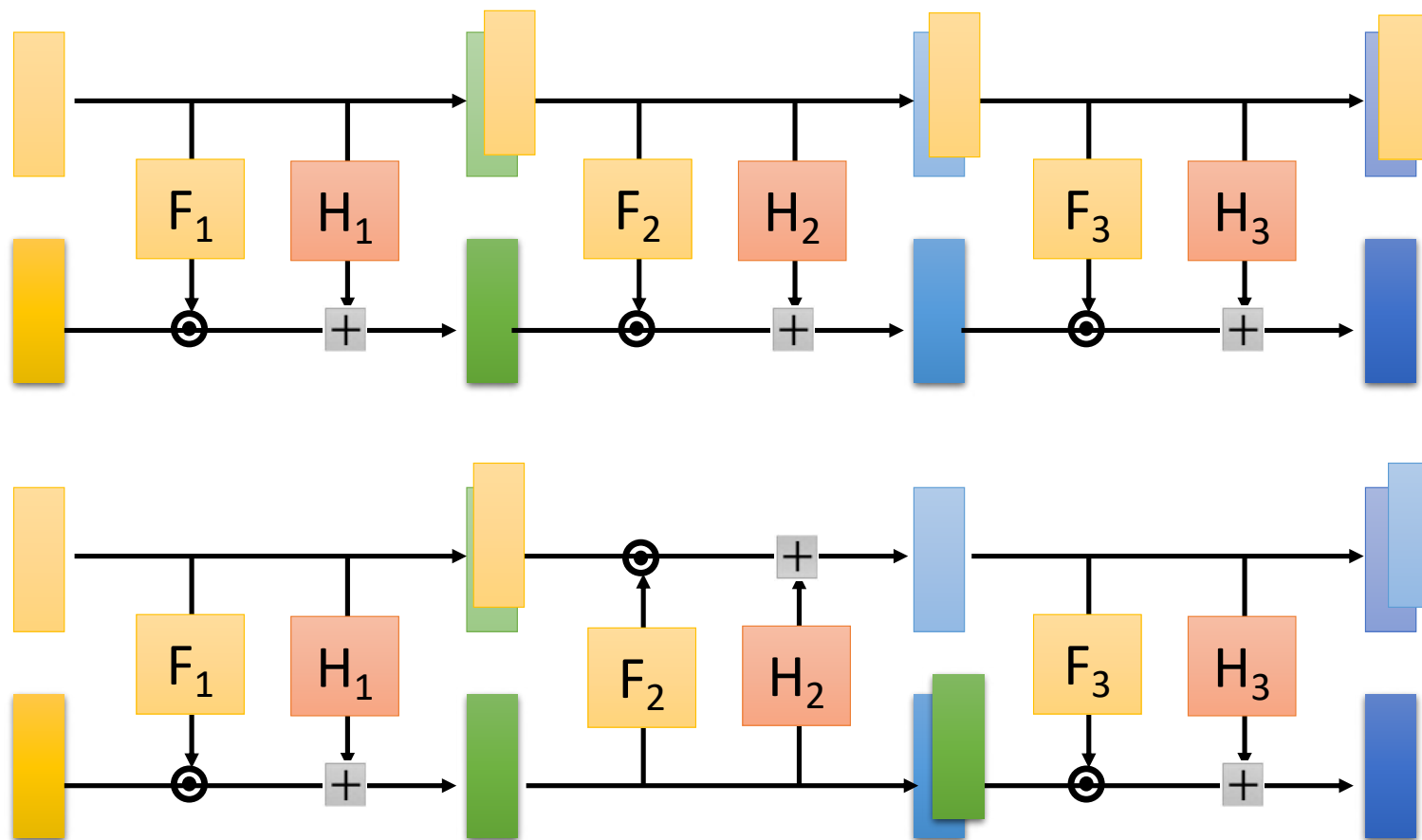$$x_{i>d} = \beta z_{i>d} + \gamma$$

# Coupling Layer

- We can use coupling layer to design invertible function and calculate the determinant of Jacobian efficiently!
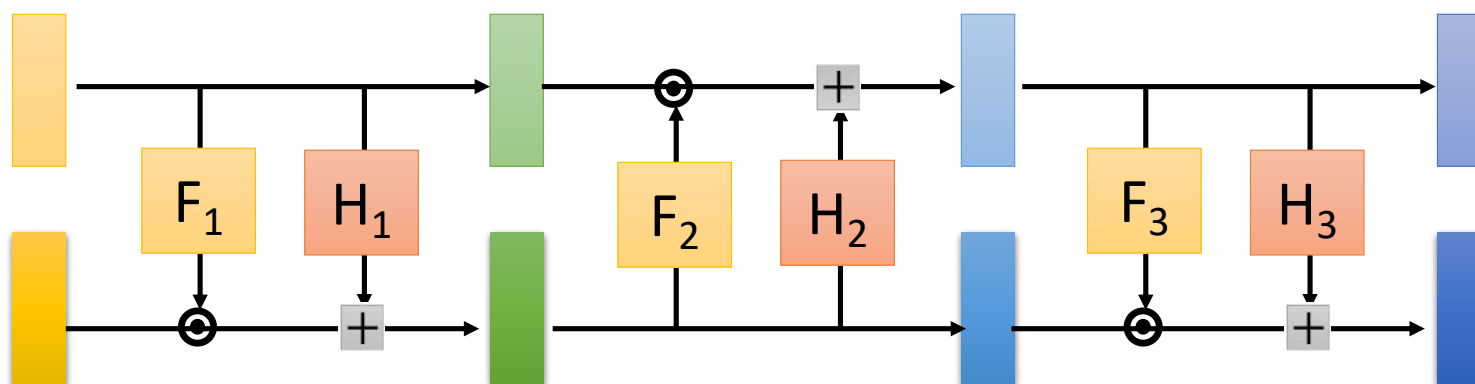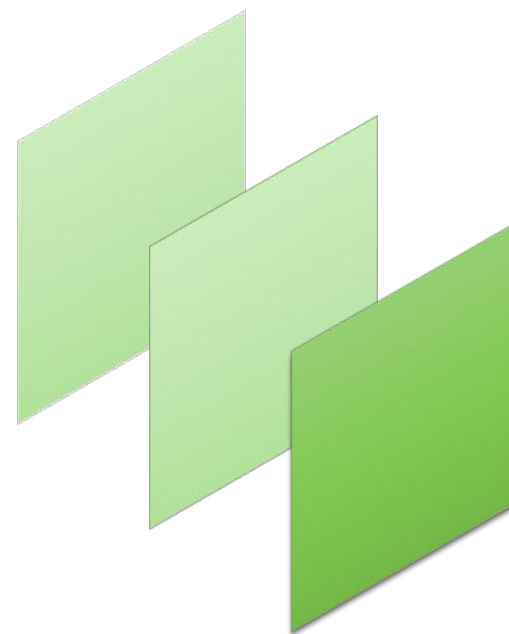
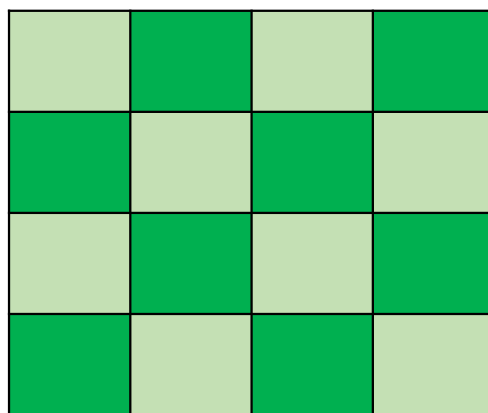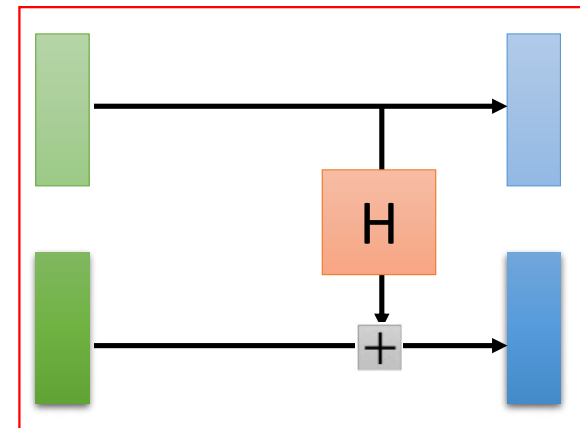# Coupling Layer - Stacking

# Coupling Layer

- Coupling layer based normalising flow models
  - Coupling layer
  - NICE
  - Real NVP
  - Glow
- Autoregressive models as flow models
  - MAF
  - IAF
  - Parallel Wavenet

# NICE: Nonlinear Independent Components Estimation

- **Additive** coupling layers
  - Partition the variables **z** into two disjoint subsets
  - $x_{1:d} = z_{1:d}$
  - $x_{d+1:n} = z_{d+1:n} + H(z_{1:d})$
  - **Volume preserving transformation** since determinant is 1.



- Additive coupling layers are composed together (with arbitrary partitions of variables in each layer)

- Final layer of NICE applies a rescaling transformation

Dinh et al., 2014. Nonlinear Independent Components Estimation

# NICE - Rescaling layers

- **Rescaling** layers
  - Forward:
    - $x_i = beta_i z_i$, where $s_i > 0$ is the scaling factor for the i-th dimension.
  - Inverse:
    - $z_i = x_i / beta_i$
  - Jacobian:
    - $J = diag(\boldsymbol{beta})$

# Samples generated via NICE



(a) Model trained on MNIST

(b) Model trained on TFD

# Samples generated via NICE
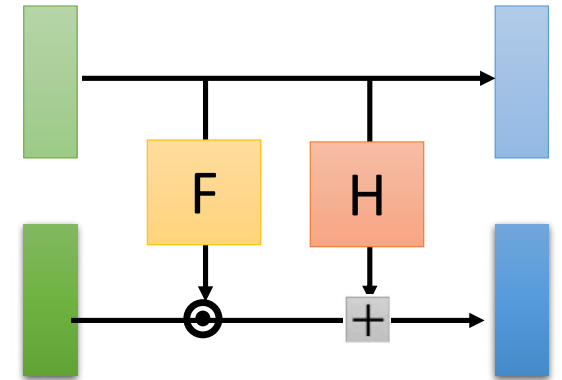


(c) Model trained on SVHN

(d) Model trained on CIFAR-10

- Coupling layer based normalising flow models
  - Coupling layer
  - NICE
  - **Real NVP**
  - Glow
- Autoregressive models as flow models
  - MAF
  - IAF
  - Parallel Wavenet

# Real NVP



- Coupling layers
  - Partition the variables **z** into two disjoint subsets
  - $x_{1:d} = z_{1:d}$
  - $x_{d+1:n} = z_{d+1:n} \odot F(z_{1:d}) + H(z_{1:d})$
  - **Non-volume preserving transformation** in general since determinant can be less than or greater than 1

- Coupling layers are composed together (with arbitrary partitions of variables in each layer)

Dinh et al., 2017. Density estimation using Real NVP

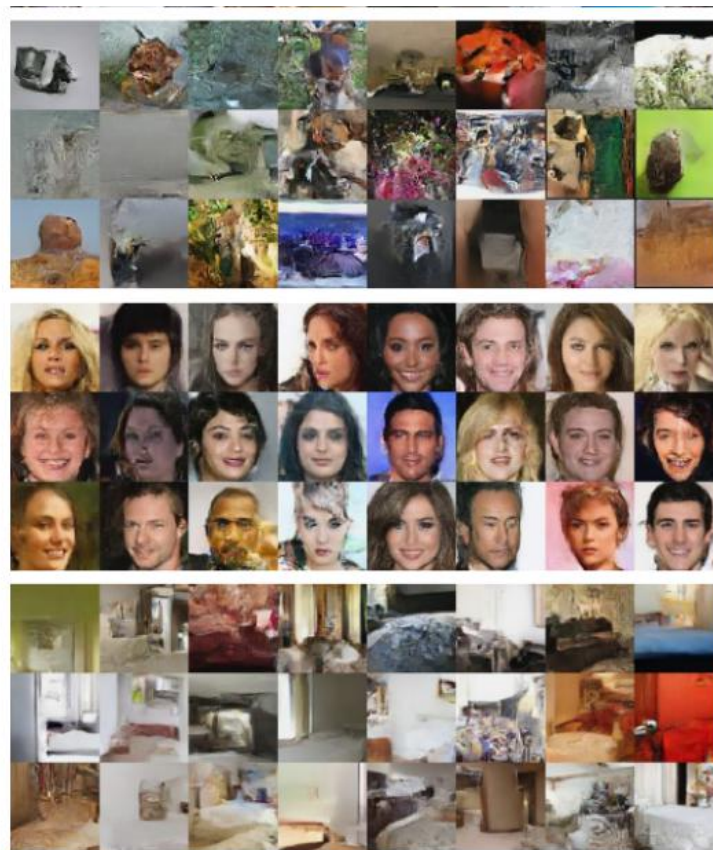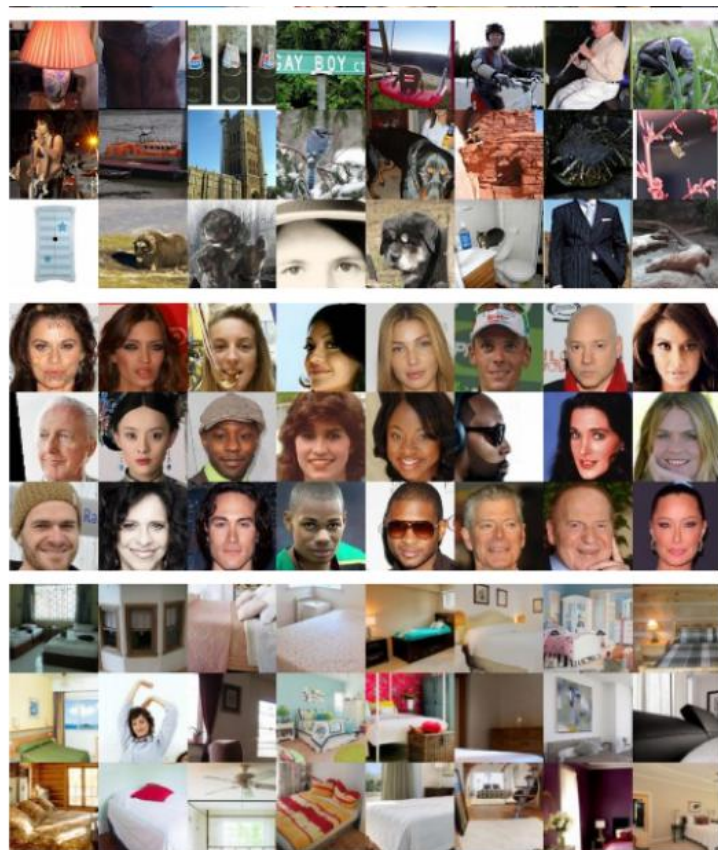# Samples generated via Real-NVP

- Coupling layer based normalising flow models
  - Coupling layer
  - NICE
  - Real NVP
  - Glow
- Autoregressive models as flow models
  - MAF
  - IAF
  - Parallel Wavenet

# Glow: Generative Flow with Invertible 1×1 Convolutions



Kingma et al. Glow: Generative Flow with Invertible 1x1 Convolutions

# 1x1 Convolution

$W$ can shuffle the channels.

If $W$ is invertible, it is easy to compute $W^{-1}$.

| 3 | | 0 | 0 | 1 | | 1 |
|---|---|---|---|---|---|---|
| 1 | = | 1 | 0 | 0 | | 2 |
| 2 | | 0 | 1 | 0 | | 3 |

# 1x1 Convolution

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix}$$



W

3 x 3

$z$

$x$

$$x = f(z) = Wz$$

$$J_f = \begin{bmatrix} \partial x_1/\partial z_1 & \partial x_1/\partial z_2 & \partial x_1/\partial z_3 \\ \partial x_2/\partial z_1 & \partial x_2/\partial z_2 & \partial x_2/\partial z_3 \\ \partial x_3/\partial z_1 & \partial x_3/\partial z_2 & \partial x_3/\partial z_3 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix} = W$$

# 1x1 Convolution

$$\left(det(W)\right)^{d \times d}$$

If $W$ is 3x3, computing $det(W)$ is easy.



$d \times d$ positions (pixels)

R G B

$W$

$W$

0

0

$W$

# Image results: Glow



Figure 5: Linear interpolation in latent space between real images
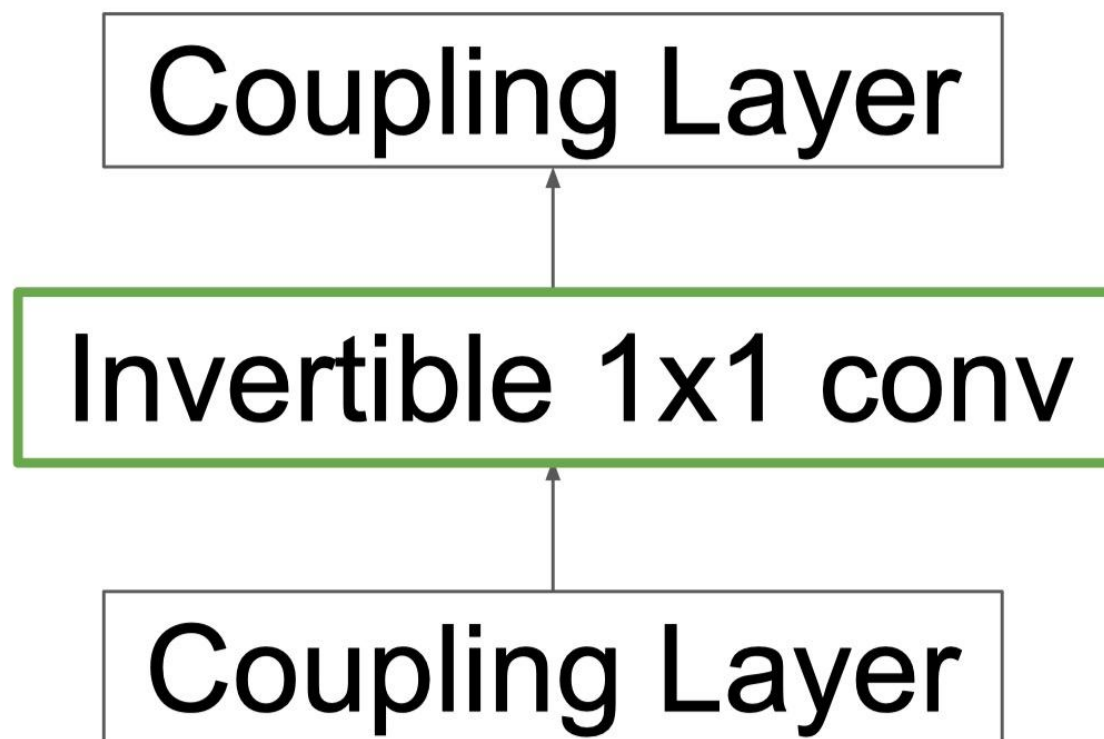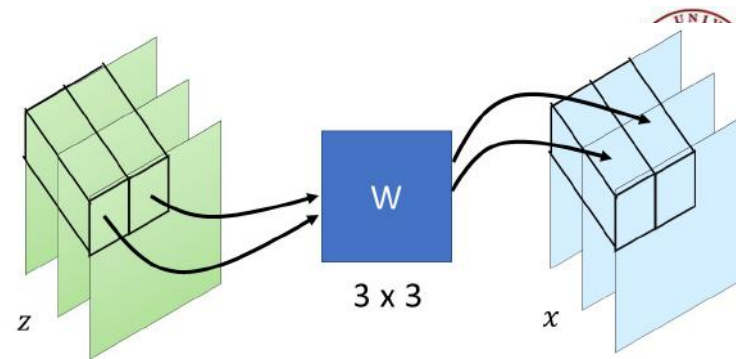
Coupling Layer

Invertible 1x1 conv

Coupling Layer

- Coupling layer based normalising flow models
  - Coupling layer
  - NICE
  - Real NVP
  - Glow

- Autoregressive models as flow models
  - MAF
  - IAF
  - Parallel Wavenet

# Autoregressive models as flow models

- Consider a Gaussian autoregressive model:
  - $p(\boldsymbol{x}) = \prod_{i=1}^{n} p(x_i | \boldsymbol{x}_{<i})$
  - Such that $p(x_i | \boldsymbol{x}_{<i}) = N(\mu_i(x_1, \ldots, x_{i-1}), \exp(\alpha_i(x_1, \ldots, x_{i-1}))^2)$, $\mu_i, \alpha_i$ are neural networks.
- Sampler for this model:
  - Sample $z_i \sim N(0,1)$
  - Let $x_i = \exp(\alpha_i) z_i + \mu_i$ ← look like coupling layer ~~
- **Flow interpretation:** transform **z** to **x** via invertible transformation (parameterised by $\mu_i, \alpha_i$)

# Masked Autoregressive Flow (MAF)



$$x_i = z_i \cdot \exp(\alpha_i) + \mu_i \quad \forall i \in \{1 \ldots n\}$$

- Forward: **(z to x)**
  - $x_i = z_i \exp(\alpha_i) + \mu_i$
  - Then calculate $\alpha_{i+1}, \mu_{i+1}$
- Sampling is sequential and slow (like autoregressive)

Figure adapted from Eric Jang's blog          Papamakarios et al. Masked Autoregressive Flow for Density Estimation 30

# Masked Autoregressive Flow (MAF)



Figure adapted from Eric Jang's blog

- Inverse **(x** to **z)**
  - $z_i = (x_i - \mu_i)\exp(-\alpha_i)$
- can be done in parallel.

- Jacobian is lower diagonal; hence determinant can be computed efficiently
- Likelihood evaluation is easy and parallelisable

- $\max\limits_{\theta} logp(D;\theta) = \sum_{x \in D} log\pi\left(G_\theta^{-1}(x)\right) + log\left|det\left(\frac{\partial G_\theta^{-1}(x)}{\partial \boldsymbol{x}}\right)\right|$

- MAF can calculate $G_\theta^{-1}(x)$ parallel.

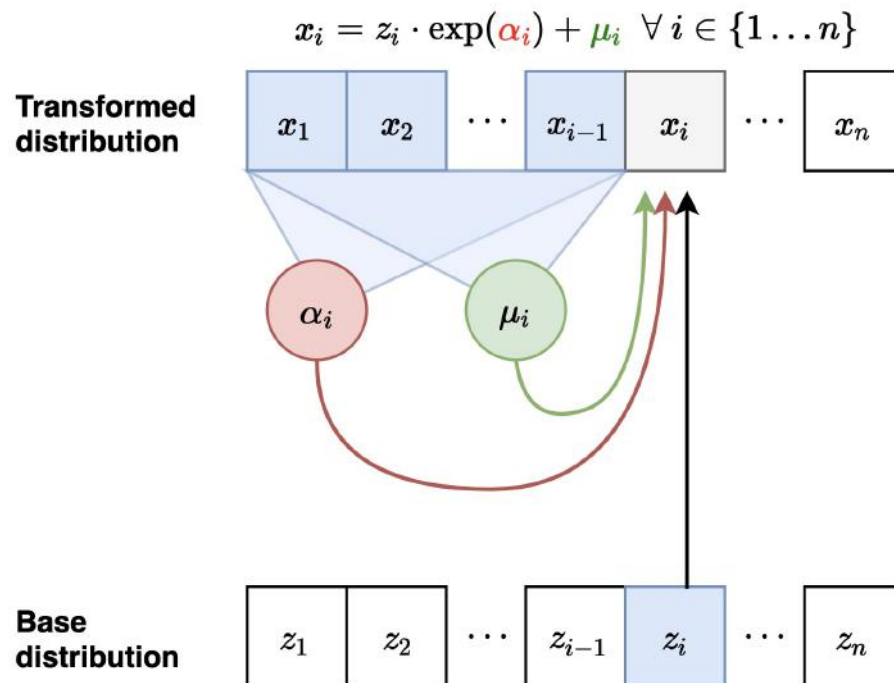- MAF: Fast likelihood evaluation (parallel), slow sampling (sequential)

- Coupling layer based normalising flow models
  - Coupling layer
  - NICE
  - Real NVP
  - Glow

- Autoregressive models as flow models
  - MAF

  - IAF

  - Parallel Wavenet

# Inverse Autoregressive Flow (IAF)
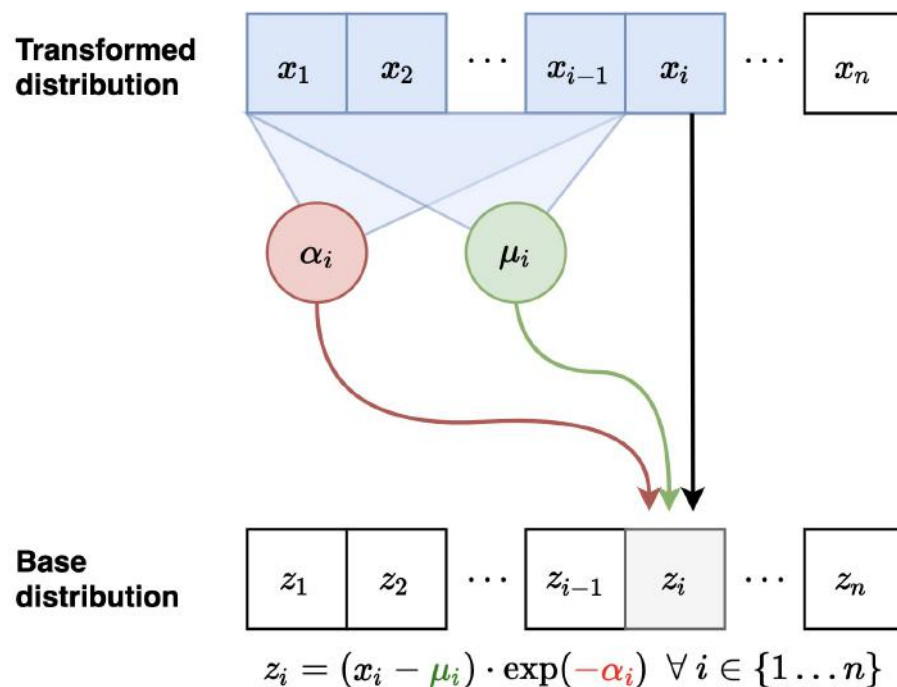


$$x_i = z_i \cdot \exp(\alpha_i) + \mu_i \quad \forall\, i \in \{1 \ldots n\}$$

**Transformed distribution**

$x_1$ | $x_2$ | $\cdots$ | $x_{i-1}$ | $x_i$ | $\cdots$ | $x_n$

$\alpha_i$    $\mu_i$

**Base distribution**

$z_1$ | $z_2$ | $\cdots$ | $z_{i-1}$ | $z_i$ | $\cdots$ | $z_n$

- Forward: **(z to x)**
  - $x_i = z_i \exp(\alpha_i) + \mu_i$
  - parallel
- Inverse **(x to z)**
  - $z_i = (x_i - \mu_i)\exp(-\alpha_i)$
  - Then compute $\mu_i, \alpha_i$
  - sequential

Figure adapted from Eric Jang's blog

Kingma et al. Improving Variational Inference with Inverse Autoregressive Flow

# Inverse Autoregressive Flow (IAF)

- Fast to sample (parallel)

- Slow to evaluate likelihoods of data points during training (sequential)

- Fast to evaluate likelihoods of a generated point (we only need to cache $z_1, z_2, \ldots, z_n$)

# IAF is inverse of MAF



Figure: Inverse pass of MAF (**left**) vs. Forward pass of IAF (**right**)

# IAF vs. MAF

- Computational tradeoffs
  - MAF: Fast likelihood evaluation, slow sampling
  - IAF: Fast sampling, slow likelihood evaluation
- MAF more suited for training based on MLE, density estimation
- IAF more suited for real-time generation
- Can we get the best of both worlds?

- Coupling layer based normalising flow models
  - Coupling layer
  - NICE
  - Real NVP
  - Glow

- Autoregressive models as flow models
  - MAF
  - IAF

  - Parallel Wavenet

# Parallel Wavenet

MAF: $x \mapsto z$ parallel
IAF: $z \mapsto x$ parallel

- Two part training with a teacher (MAF) and student model (IAF)

- Teacher can be efficiently trained via MLE.

- Once teacher is trained, initialise a student model parameterised by IAF. Student model cannot efficiently evaluate density for external data points but allows for efficient sampling

- **Key observation:** IAF can also efficiently evaluate densities of its own generations (via caching the noise variates $z_1, z_2, \dots, z_n$)

Oord et al. Parallel WaveNet: Fast High-Fidelity Speech Synthesis. 2017

# Parallel Wavenet

MAF: $x \mapsto z$ parallel
IAF: $z \mapsto x$ parallel



$z_i = (x_i - \mu_i) \cdot \exp(-\alpha_i) \ \forall i \in \{1 \dots n\}$

Inverse pass of MAP (training)

**WaveNet Teacher**

Linguistic features - - - - →

Teacher Output $P(x_i | x_{<i})$

Generated Samples $x_i = g(z_i | z_{<i})$

**WaveNet Student**

Linguistic features - - - - →

Student Output $P(x_i | z_{<i})$

Input noise $z_i$

forward pass of IAF (sampling)

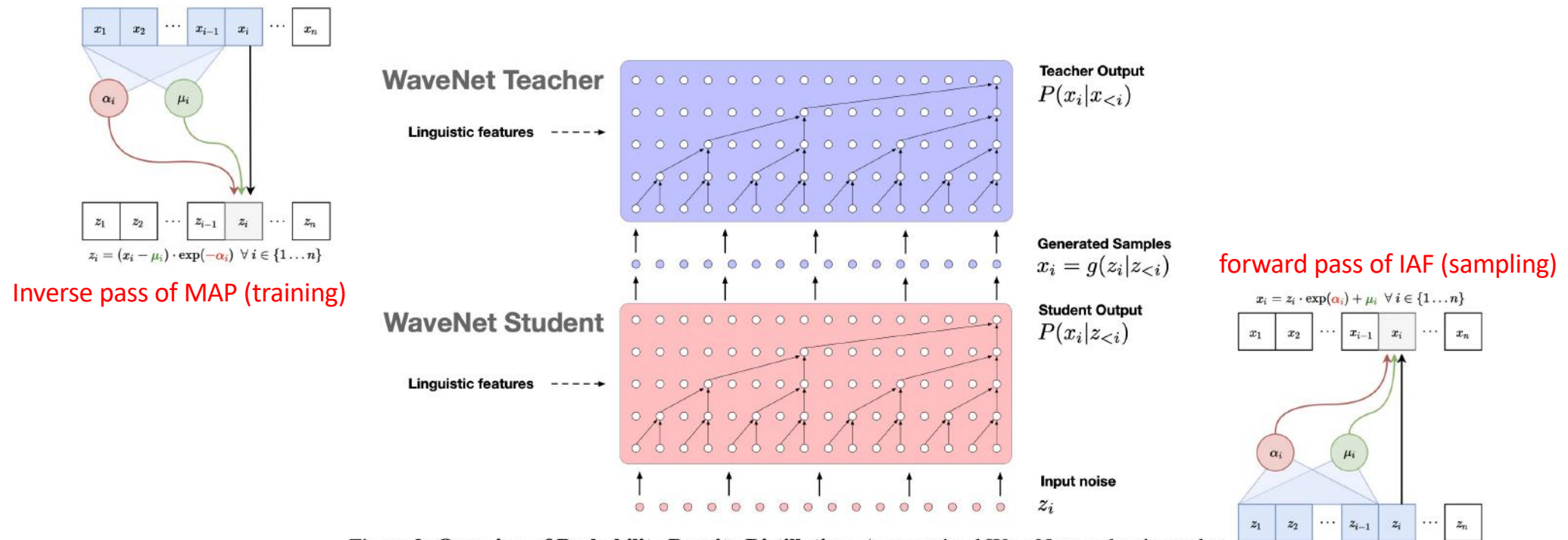$x_i = z_i \cdot \exp(\alpha_i) + \mu_i \ \forall i \in \{1 \dots n\}$

Figure 2: **Overview of Probability Density Distillation**. A pre-trained WaveNet teacher is used to score the samples $x$ output by the student. The student is trained to minimise the KL-divergence between its distribution and that of the teacher by maximising the log-likelihood of its samples under the teacher and maximising its own entropy at the same time.

Oord et al. Parallel WaveNet: Fast High-Fidelity Speech Synthesis. 2017

# Parallel Wavenet

MAF: $x \mapsto z$ parallel
IAF: $z \mapsto x$ parallel

- **Probability density distillation**: Student distribution is trained to minimise the KL divergence between student (s) and teacher (t)
$$D_{KL}(s,t) = E_{\boldsymbol{x} \sim s}[\log(s(\boldsymbol{x})) - \log(t(\boldsymbol{x}))]$$

- Evaluating and optimising Monte Carlo estimates of this objective requires:
  - Samples **x** from student model (IAF)
  - Density of **x** assigned by student model (IAF)
  - Density of **x** assigned by teacher model (MAF)

- All operations above can be implemented efficiently!

# Parallel Wavenet: Overall algorithm

- Training
  - Step 1: Train teacher model (MAF) via MLE
  - Step 2: Train student model (IAF) to minimize KL divergence with teacher
- Test-time: Use student model for testing
- Improves sampling efficiency over original Wavenet (vanilla autoregressive model) by 1000x!
- Useful in speech synthesis

- Coupling layer based normalising flow models
  - Coupling layer
  - NICE          add only
  - Real NVP      add+mul
  - Glow          conv 1x1
- Autoregressive models as flow models
  - MAF           fast train, slow test
  - IAF           fast test, slow train
  - Parallel Wavenet  fast train, fast test

# Summary of Normalising Flow Models

- Transform simple distributions into more complex distributions via change of variables

- Jacobian of transformations should have tractable determinant for efficient learning and density estimation

- Computational tradeoffs in evaluating forward and inverse transformations

# Thanks