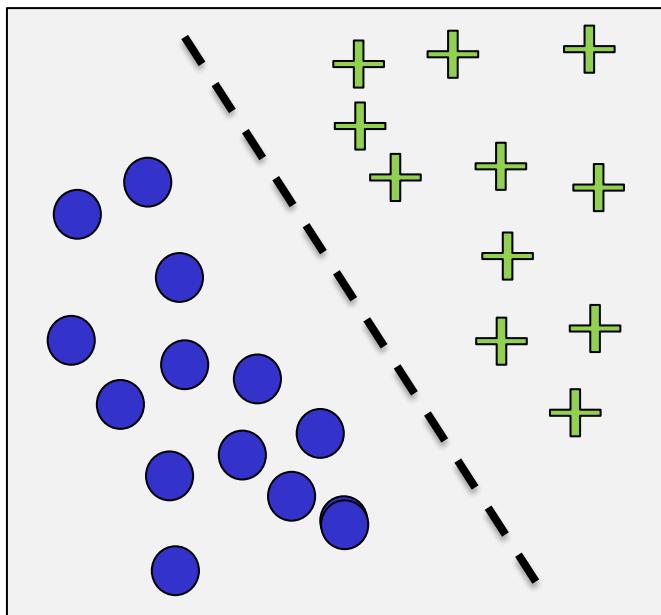


Réseaux de neurones  
IFT 780

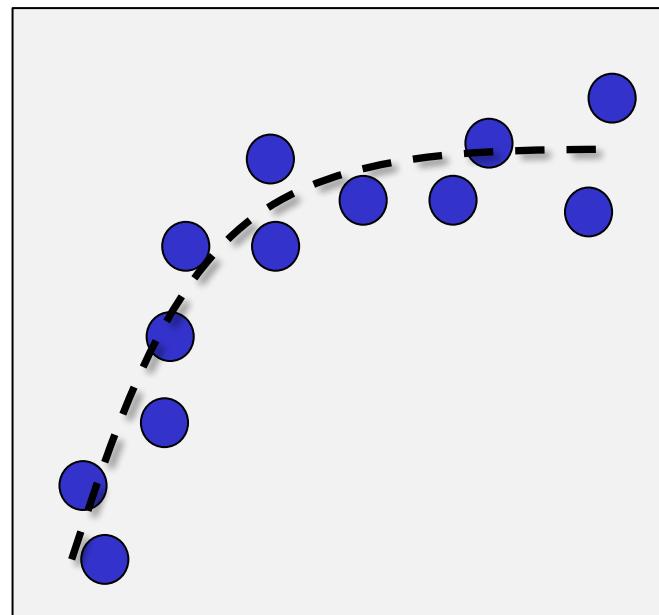
Modèles génératifs  
Par  
Pierre-Marc Jodoin, Antoine Théberge

# Jusqu'à présent : apprentissage supervisé

Classification



Régression



# Jusqu'à présent : apprentissage supervisé

Segmentation



Description



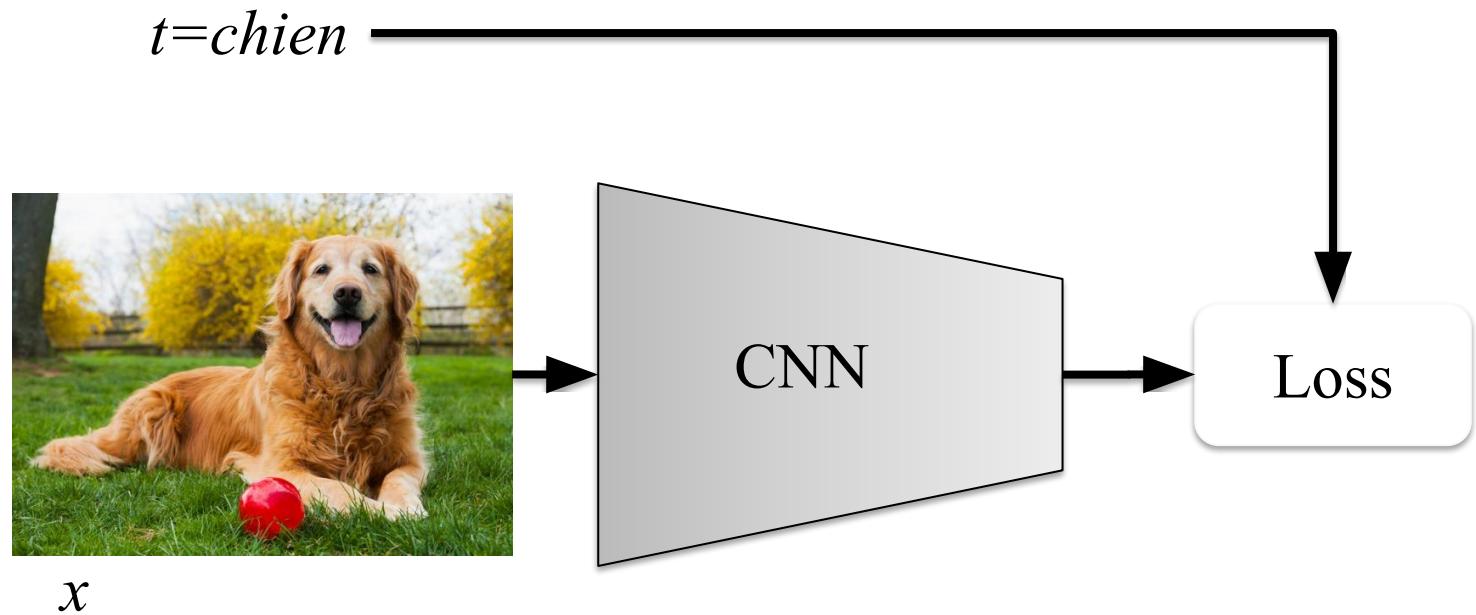
an elephant standing in a grassy field with  
trees in the background

# Apprentissage supervisé

Deux grandes familles d'applications

- **Classification** : la cible est un indice de classe  $t \in \{1, \dots, K\}$ 
  - Exemple : reconnaissance de caractères
    - ✓  $x$  : vecteur des intensités de tous les pixels de l'image
    - ✓  $t$  : identité du caractère
- **Régression** : la cible est un nombre réel  $t \in \mathbb{R}$ 
  - Exemple : prédiction de la valeur d'une action à la bourse
    - ✓  $x$  : vecteur contenant l'information sur l'activité économique de la journée
    - ✓  $t$  : valeur d'une action à la bourse le lendemain

# Apprentissage supervisé avec CNN



# Modèles discriminatifs vs. génératifs

$p(t|x)$

probabilité de la classe  $t$   
en sachant  $x$

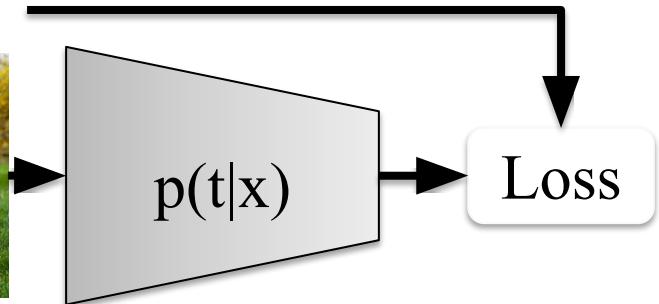
$p(x)$

probabilité de  $x$

$t=chien$

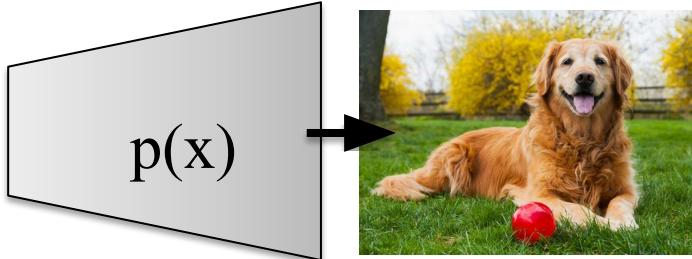


$x$



$p(t|x)$

Loss



$p(x)$

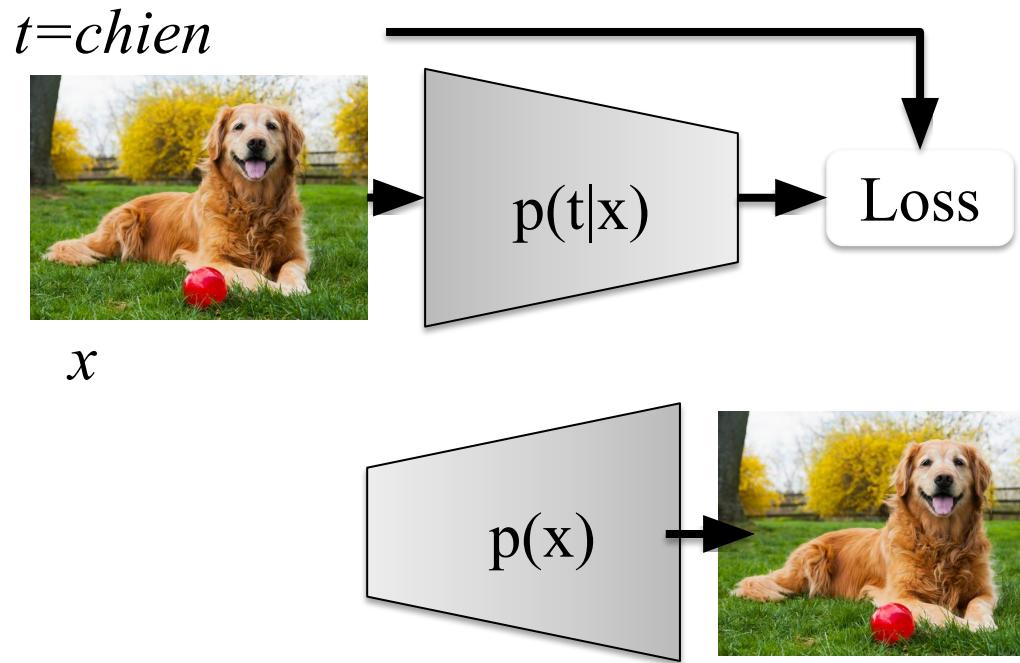
# Modèles discriminatifs vs. génératifs

$p(t|x) \rightarrow p(\text{chien} | \text{})$

probabilité de la classe t  
en sachant x

$p(x) \rightarrow p(\text{})$

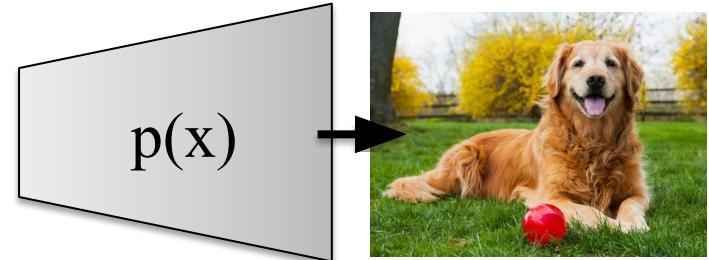
probabilité de x



# Modèles génératifs

Comment apprendre  $p(x)$  ?

$p(x) \rightarrow p(\text{dog})$   
probabilité de  $x$



Nous allons formuler  $p(x)$  par un réseau de neurones  $p(\vec{x}) = p_\theta(\vec{x})$

Selon un jeu de données  $D = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N\}$

$$\theta^* = \operatorname{argmax}_{\theta} \prod_{i=1}^N p_{\theta}(\vec{x}_i)$$

$$= \operatorname{argmax}_{\theta} \sum_{i=1}^N \log p_{\theta}(\vec{x}_i)$$

# Supervisé vs non supervisé

**Apprentissage supervisé** : il y a une cible

$$D = \{(x_1, t_1), (x_2, t_2), \dots, (x_N, t_N)\}$$

**Apprentissage non-supervisé** : la cible n'est pas fournie

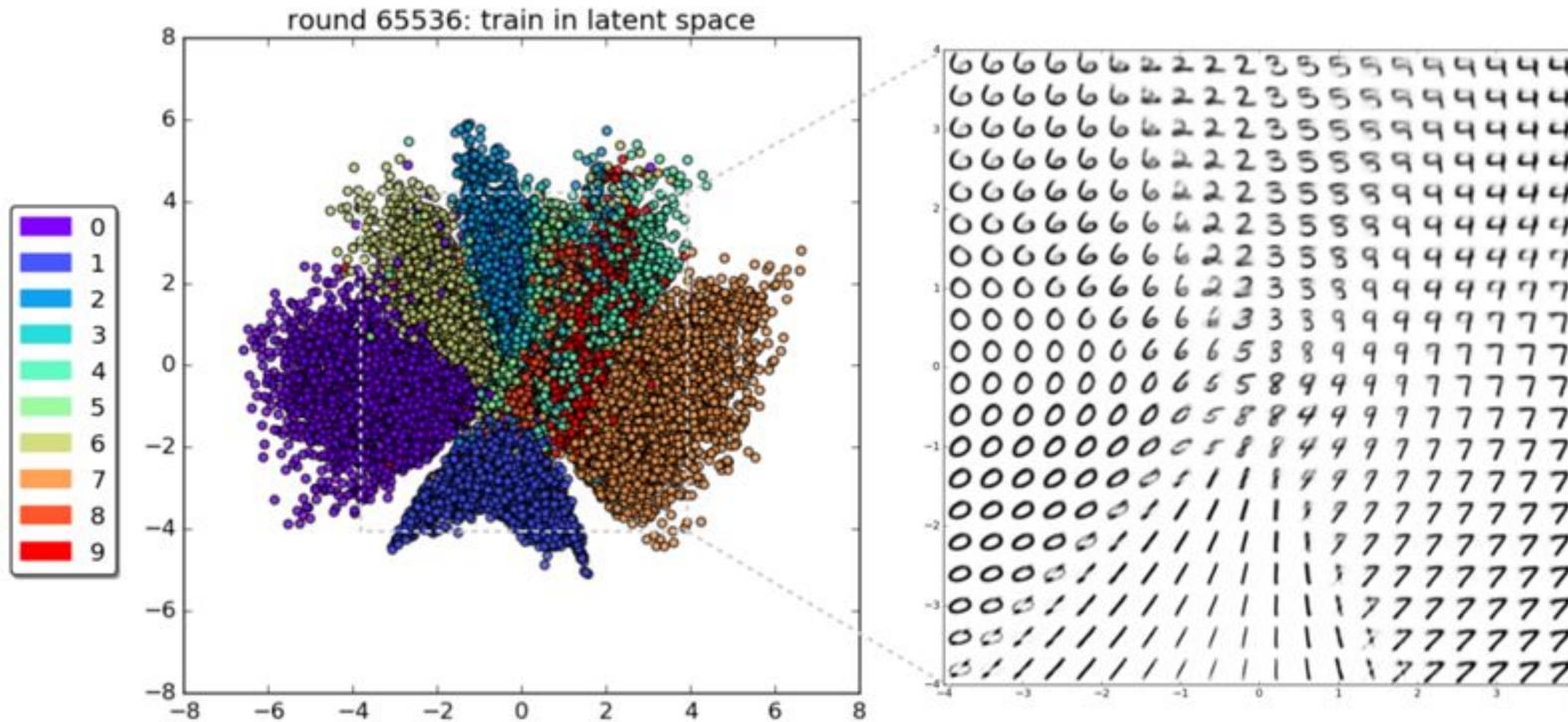
$$D = \{x_1, x_2, \dots, x_N\}$$

# Apprentissage non supervisé

Comprendre la distribution sous-jacente de données **non-étiquetées**

**Applications** : clustering (p.e. k-means), visualization, comprehension, etc.

**Exemple** : visualization de la distribution des images MNIST

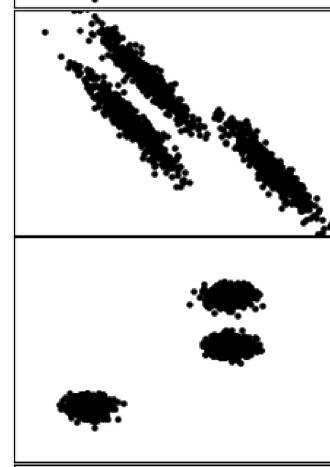
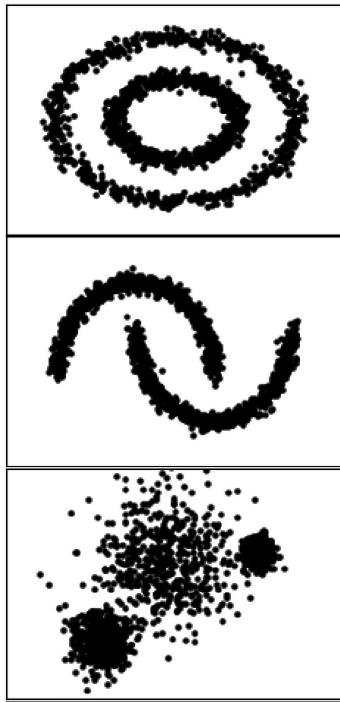


# Apprentissage non supervisé

Souvent, l'apprentissage non-supervisé inclut un (ou des) **variables latentes**.

**Variable latente:** variable aléatoire non observée mais sous-jacente à la distribution des données

**Ex:** clustering = retrouver la variable latente “cluster”

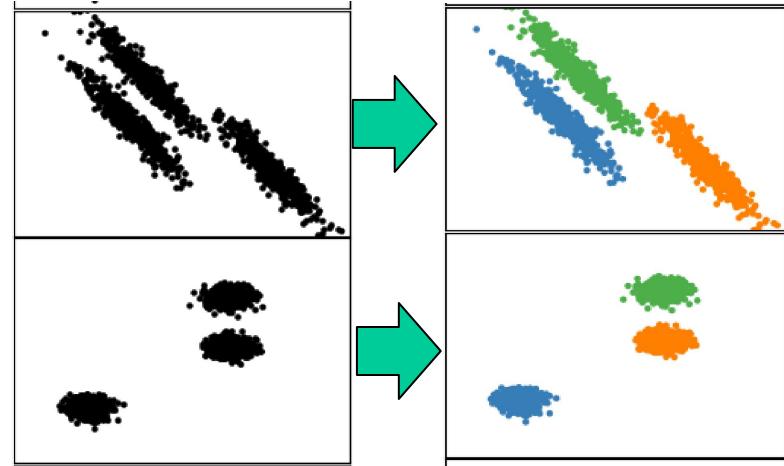
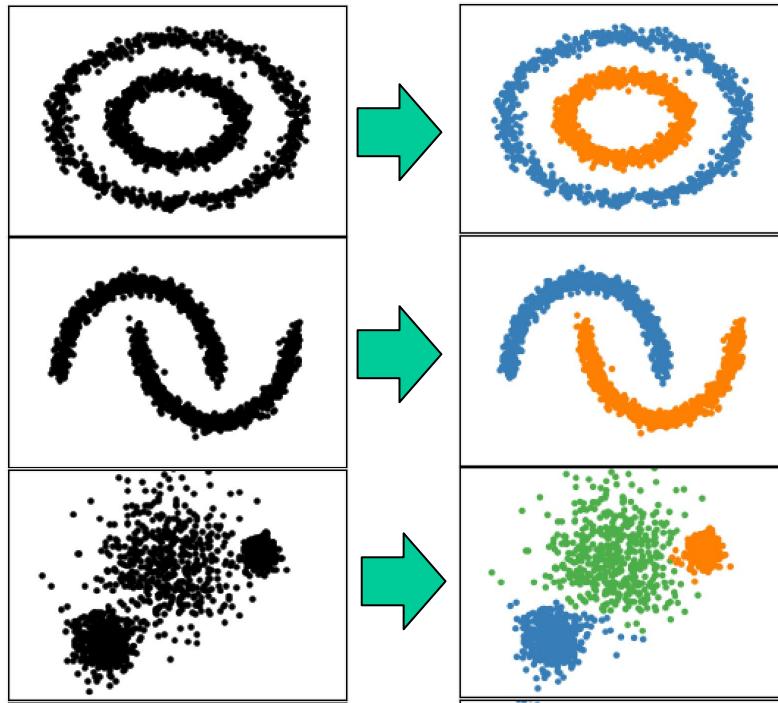


# Apprentissage non supervisé

Souvent, l'apprentissage non-supervisé inclut un (ou des) **variables latentes**.

**Variable latente:** variable aléatoire non observée mais sous-jacente à la distribution des données

Ex: clustering = retrouver la variable latente “cluster”



# Pourquoi une variable latente?

Plus facile de représenter  $p(x, y)$ ,  $p(x | y)$ ,  $p(y)$   
que  $p(x)$

# Apprentissage non supervisé

**Variable latente:** variable aléatoire non observée mais sous-jacente à la distribution des données

**Rappel:**

- si  $x_1$  et  $x_2$  sont des variables aléatoires indépendantes:  
 $p(x_1, x_2) = p(x_1)(x_2)$  et  $p(x_1|x_2) = p(x_1)$

# Apprentissage non supervisé

**Variable latente:** variable aléatoire non observée mais sous-jacente à la distribution des données

**Rappel:**

- si  $x_1$  et  $x_2$  sont des variables aléatoires indépendantes:  
 $p(x_1, x_2) = p(x_1)(x_2)$  et  $p(x_1|x_2) = p(x_1)$
- si  $x_1$  est dépendant de  $x_2$ , mais pas l'inverse:  
 $p(x_1|x_2) \neq p(x_1)$  MAIS  $p(x_2|x_1) = p(x_2)$   
 $p(x_1, x_2) = p(x_1|x_2)p(x_2)$

# Apprentissage non supervisé

**Variable latente:** variable aléatoire non observée mais sous-jacente à la distribution des données

**Rappel:**

- si  $x_1$  et  $x_2$  sont des variables aléatoires indépendantes:  
 $p(x_1, x_2) = p(x_1)(x_2)$  et  $p(x_1|x_2) = p(x_1)$
- si  $x_1$  est dépendant de  $x_2$ , mais pas l'inverse:  
 $p(x_1|x_2) \neq p(x_1)$  MAIS  $p(x_2|x_1) = p(x_2)$   
 $p(x_1, x_2) = p(x_1|x_2)p(x_2)$

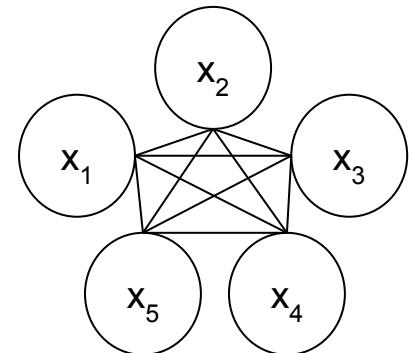
Exemple:  $x_1$  = lire dans le jardin  
 $x_2$  = il pleut dehors

# Apprentissage non supervisé

**Variable latente:** variable aléatoire non observée mais sous-jacente à la distribution des données

**Rappel:**

si  $\vec{x} = (x_1, x_2, x_3, x_4, x_5)$  sont des variables aléatoires dépendantes

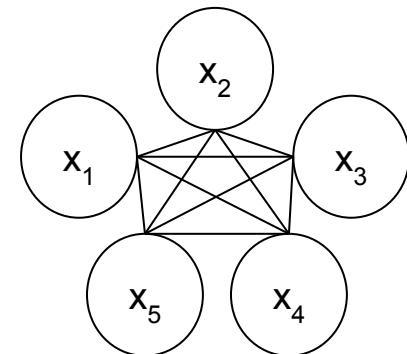


# Apprentissage non supervisé

**Variable latente:** variable aléatoire non observée mais sous-jacente à la distribution des données

**Rappel:**

si  $\vec{x} = (x_1, x_2, x_3, x_4, x_5)$  sont des variables aléatoires dépendantes



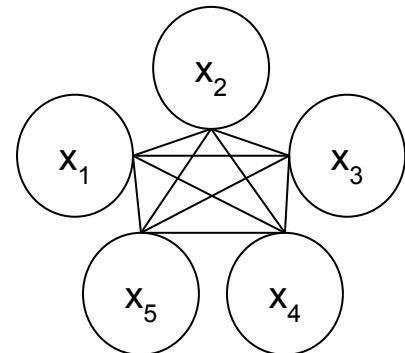
$p(\vec{x}) = p(x_1, x_2, x_3, x_4, x_5)$ , compliqué à calculer car il faut considérer chaque combinaison

# Apprentissage non supervisé

**Variable latente:** variable aléatoire non observée mais sous-jacente à la distribution des données

**Rappel:**

si  $\vec{x} = (x_1, x_2, x_3, x_4, x_5)$  sont des variables aléatoires dépendantes



$p(\vec{x}) = p(x_1, x_2, x_3, x_4, x_5)$ , compliqué à calculer car il faut considérer chaque combinaison

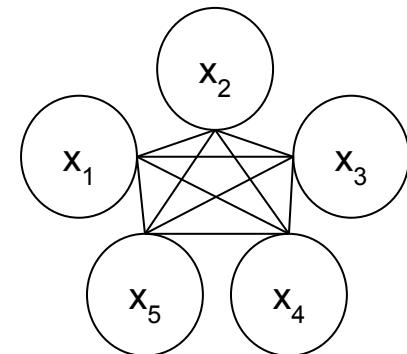
Règle en chaîne des probabilités:  $p(x_1, x_2, \dots, x_n) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1})$

# Apprentissage non supervisé

**Variable latente:** variable aléatoire non observée mais sous-jacente à la distribution des données

**Rappel:**

si  $\vec{x} = (x_1, x_2, x_3, x_4, x_5)$  sont des variables aléatoires dépendantes



$p(\vec{x}) = p(x_1, x_2, x_3, x_4, x_5)$ , compliqué à calculer car il faut considérer chaque combinaison

Règle en chaîne des probabilités:  $p(x_1, x_2, \dots, x_n) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1})$

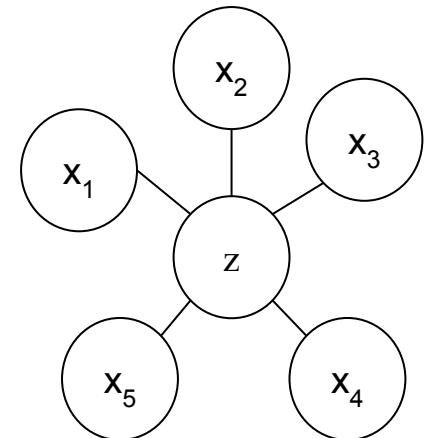
Pas si mal pour 5 variables, mais dans une image, chaque pixel est une variable !

# Apprentissage non supervisé

**Variable latente:** variable aléatoire non observée mais sous-jacente à la distribution des données

**Pour simplifier le problème** on présume que les variables  $x_i$  dépendent d'une variable latente  $z$  et sont indépendantes entre elles

$$p(\vec{x}, z) = p(x_1, x_2, x_3, x_4, x_5, z)$$

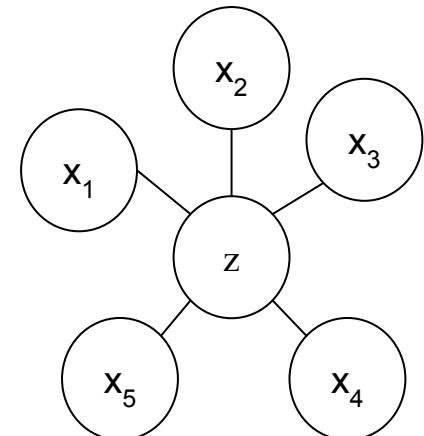


# Apprentissage non supervisé

**Variable latente:** variable aléatoire non observée mais sous-jacente à la distribution des données

Pour simplifier le problème on présuppose que les variables  $x_i$  dépendent d'une variable latente  $z$  et sont indépendantes entre elles

$$\begin{aligned} p(\vec{x}, z) &= p(x_1, x_2, x_3, x_4, x_5, z) \\ &= p(z)p(x_1|z)\dots p(x_5|z) \end{aligned}$$

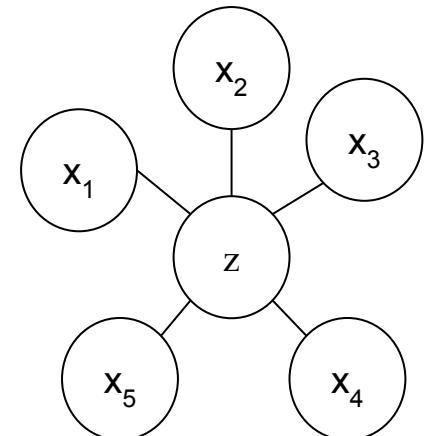


# Apprentissage non supervisé

**Variable latente:** variable aléatoire non observée mais sous-jacente à la distribution des données

Pour simplifier le problème on présuppose que les variables  $x_i$  dépendent d'une variable latente  $z$  et sont indépendantes entre elles

$$\begin{aligned} p(\vec{x}, z) &= p(x_1, x_2, x_3, x_4, x_5, z) \\ &= p(z)p(x_1|z)\dots p(x_5|z) \\ &= p(z) \prod_{i=1}^5 p(x_i|z) \end{aligned}$$

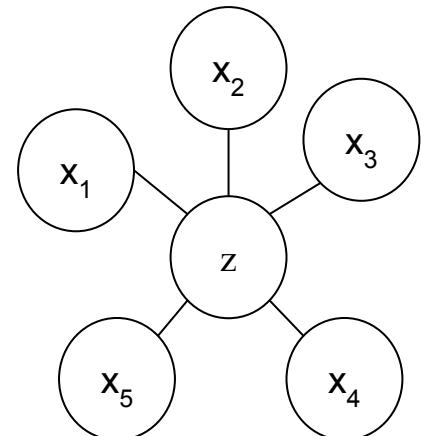


# Apprentissage non supervisé

**Variable latente:** variable aléatoire non observée mais sous-jacente à la distribution des données

Pour simplifier le problème on présuppose que les variables  $x_i$  dépendent d'une variable latente  $z$  et sont indépendantes entre elles

$$\begin{aligned} p(\vec{x}, z) &= p(x_1, x_2, x_3, x_4, x_5, z) \\ &= p(z)p(x_1|z)\dots p(x_5|z) \\ &= p(z) \prod_{i=1}^5 p(x_i|z) \\ &= p(z)p(\vec{x}|z) \end{aligned}$$

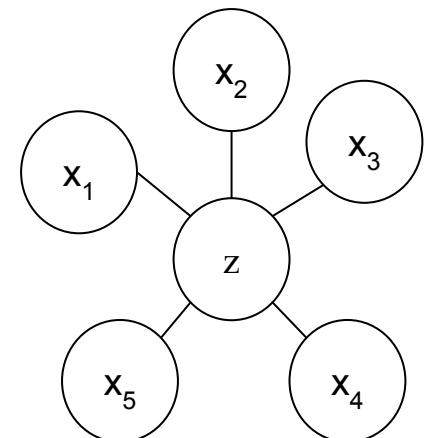


# Apprentissage non supervisé

**Variable latente:** variable aléatoire non observée mais sous-jacente à la distribution des données

Pour simplifier le problème on présuppose que les variables  $x_i$  dépendent d'une variable latente  $z$  et sont indépendantes entre elles

$$\begin{aligned} p(\vec{x}, z) &= p(x_1, x_2, x_3, x_4, x_5, z) \\ &= p(z)p(x_1|z)\dots p(x_5|z) \\ &= p(z) \prod_{i=1}^5 p(x_i|z) \\ &= p(z)p(\vec{x}|z) \end{aligned}$$



On peut toujours présumer que  $z$  existe !

$$\begin{aligned} p(\vec{x}) &= \int p(\vec{x}, z) dz \quad \text{=< marginalisation} \\ &= \int p(\vec{x}|z)p(z)dz \end{aligned}$$

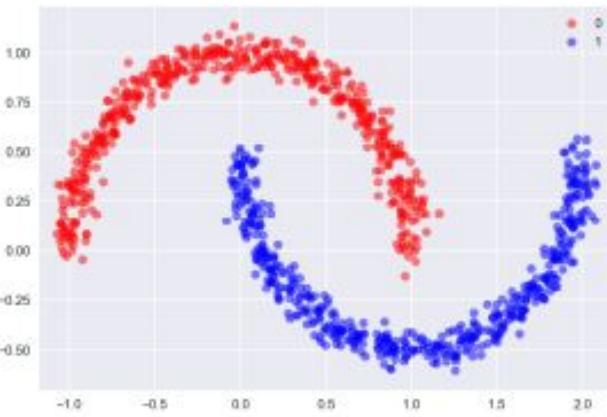
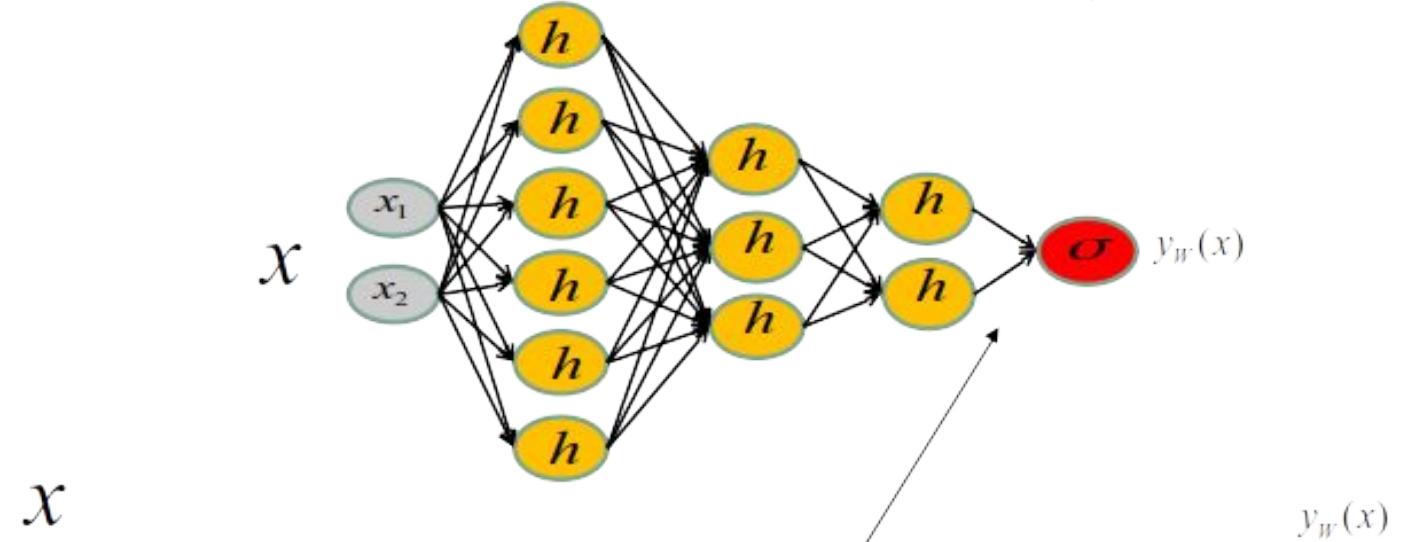
En apprentissage non-supervisé

nous nous appuierons sur

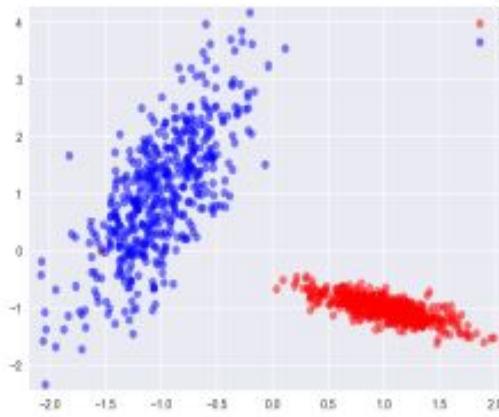
**2 propriétés** des réseaux de neurones

# Propriété 1

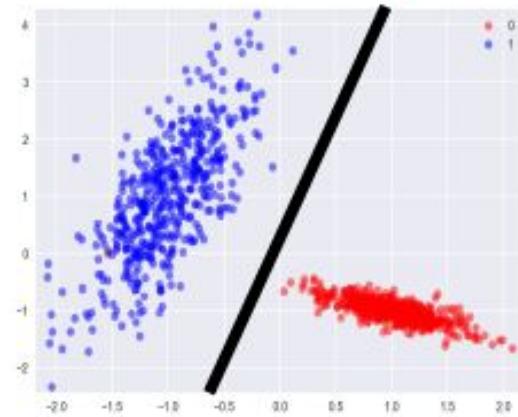
Les réseaux de neurones sont d'excellentes machines pour projeter des données brutes vers un **espace dimensionnel plus faible** dont les propriétés dépendent de la *loss* (ici **espace linéaire car la sortie du réseau est un classifieur linéaire**).



Données en entrée



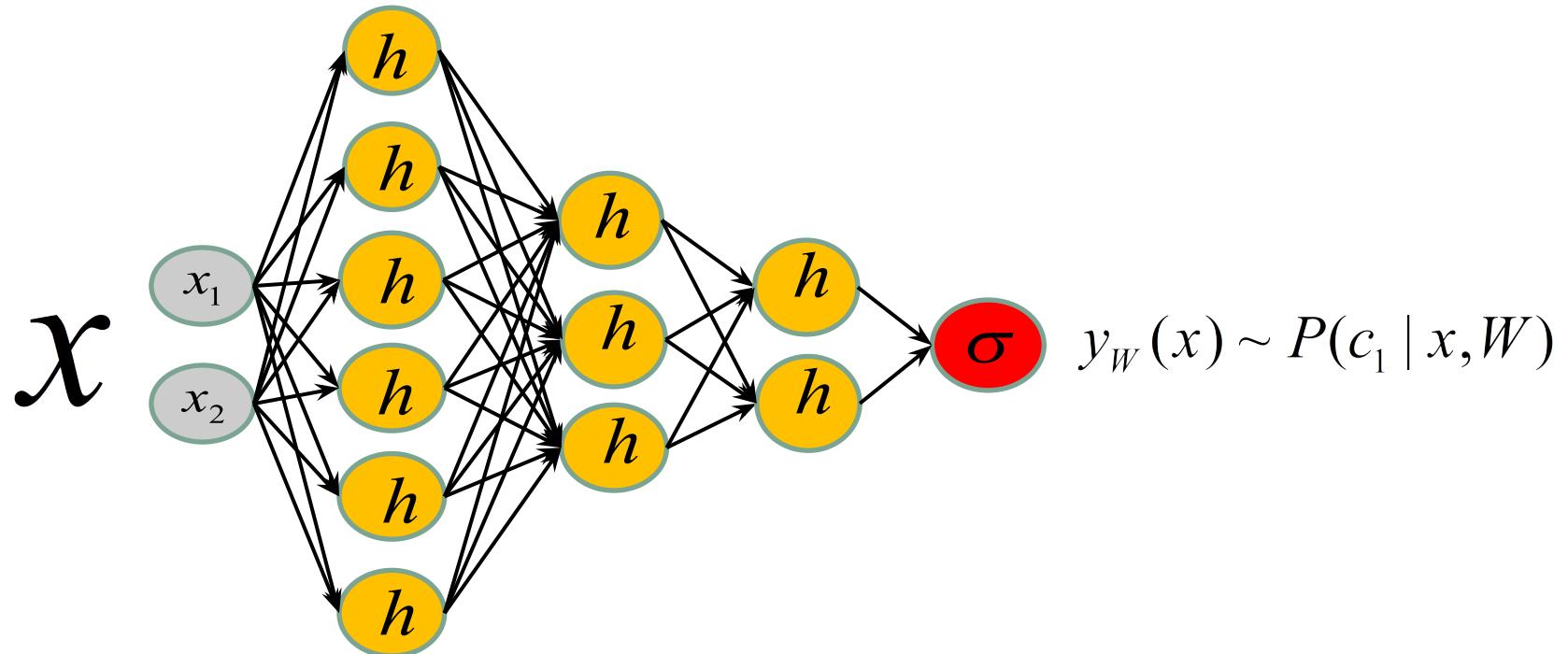
Sortie de la dernière couche



Sortie du réseau

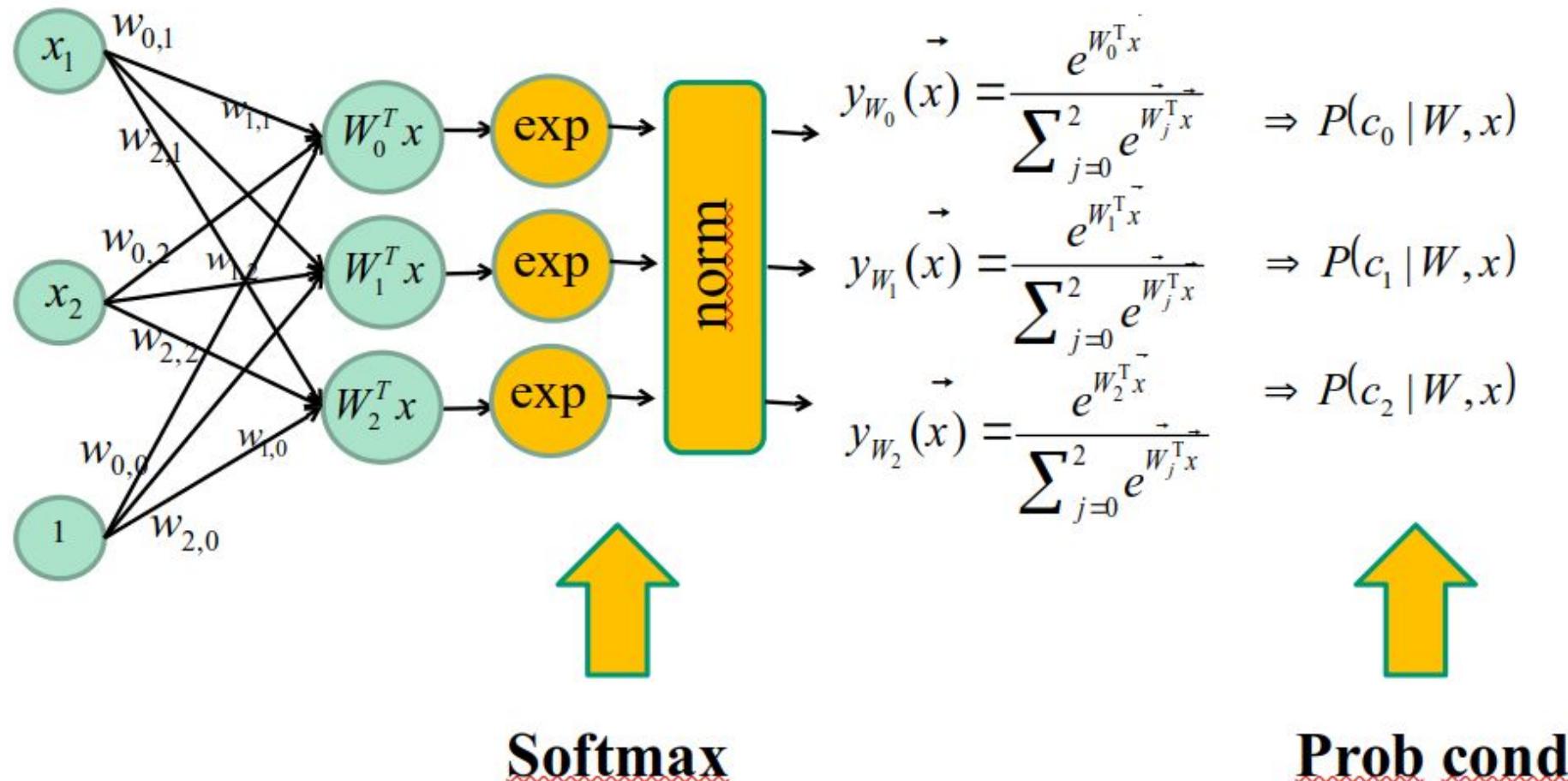
# Propriété 2

Les réseaux de neurones sont d'excellentes machines pour estimer des **probabilités conditionnelles**.



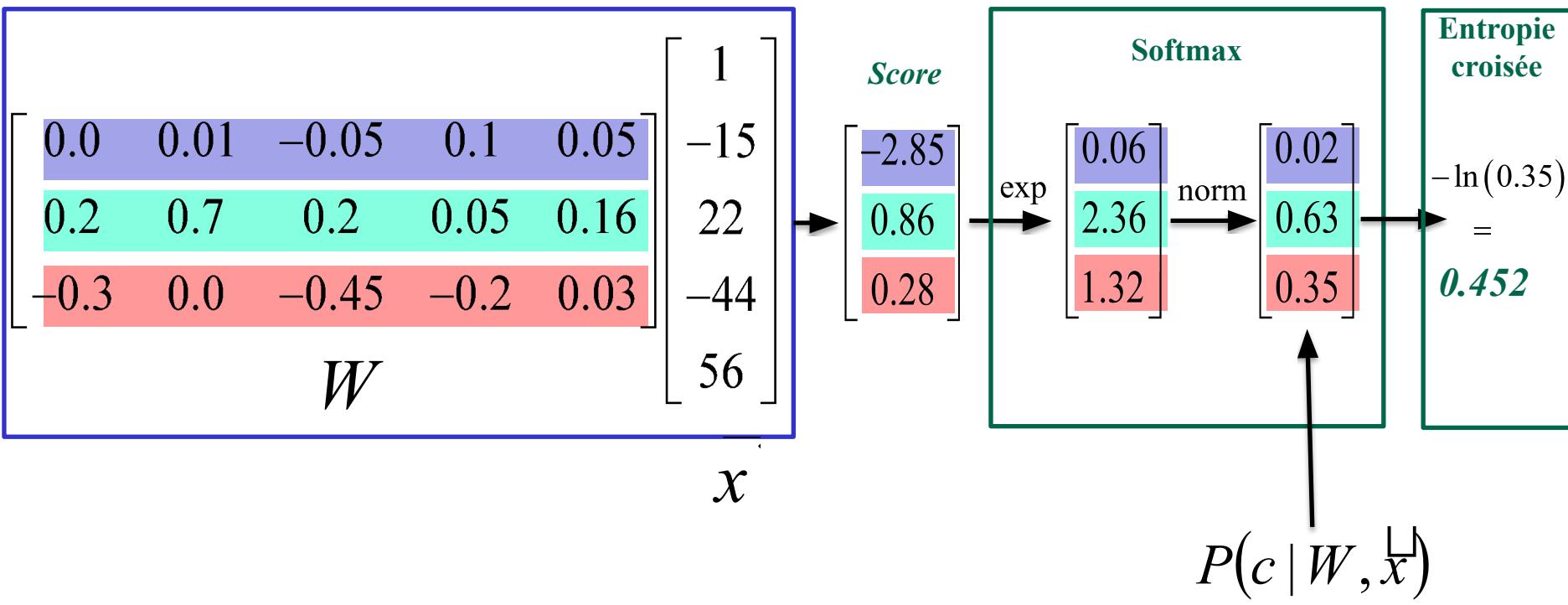
# Propriété 2

Les réseaux de neurones sont d'excellentes machines pour estimer des **probabilités conditionnelles**.



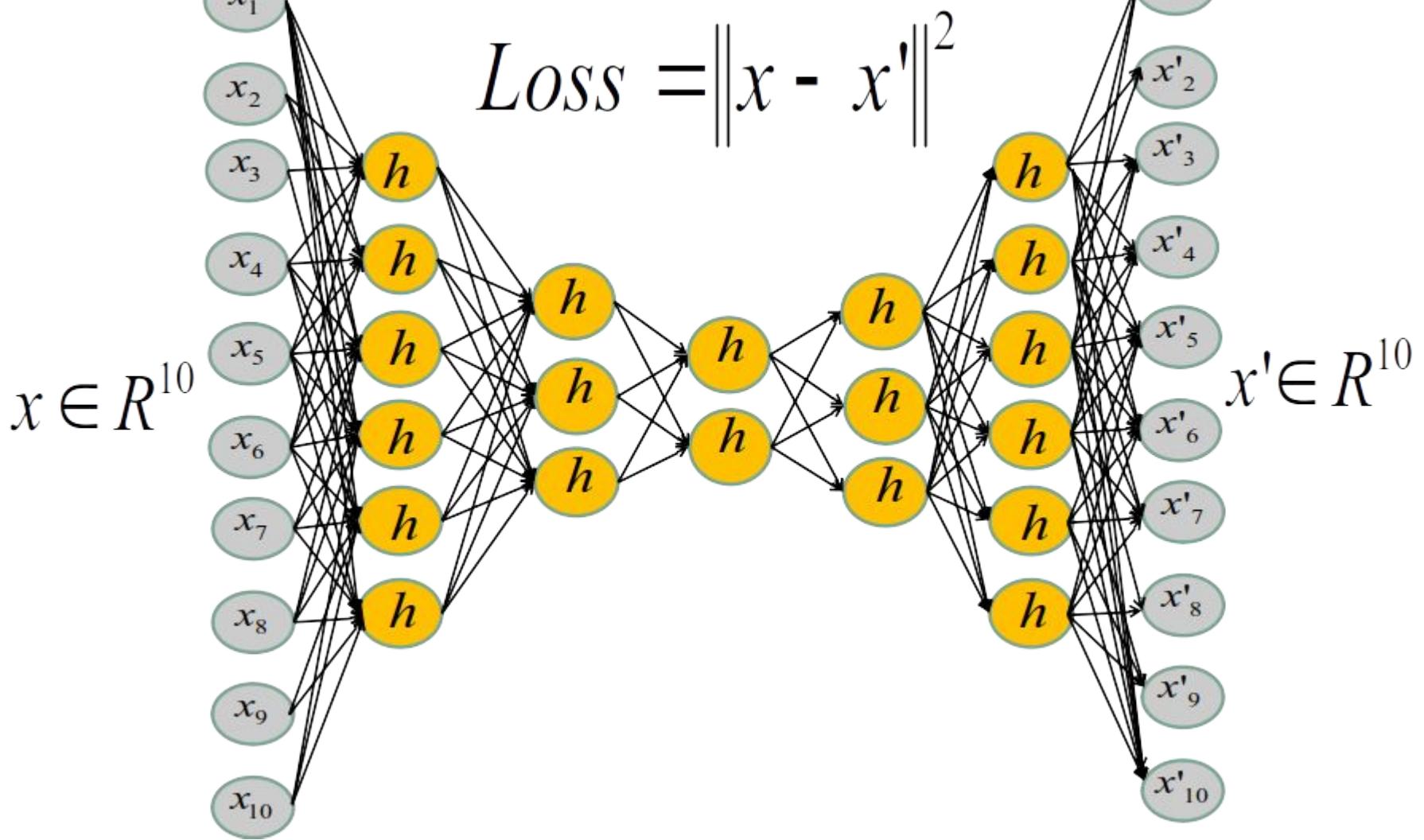
$$\vec{x} = \begin{bmatrix} -15 \\ 22 \\ -44 \\ 56 \end{bmatrix}, t = 2$$

# Rappel

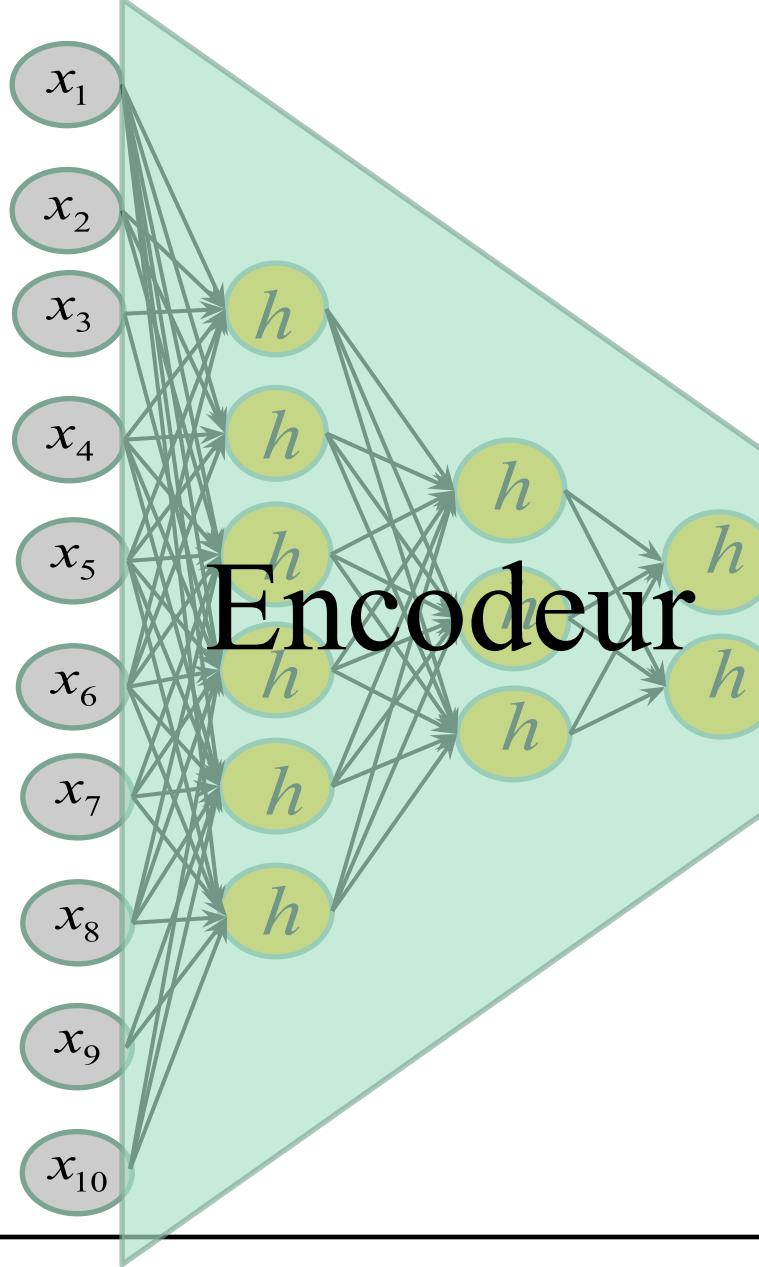


Comment utiliser un réseau de neurones pour apprendre  
la **configuration sous-jacente** de données non étiquetés?





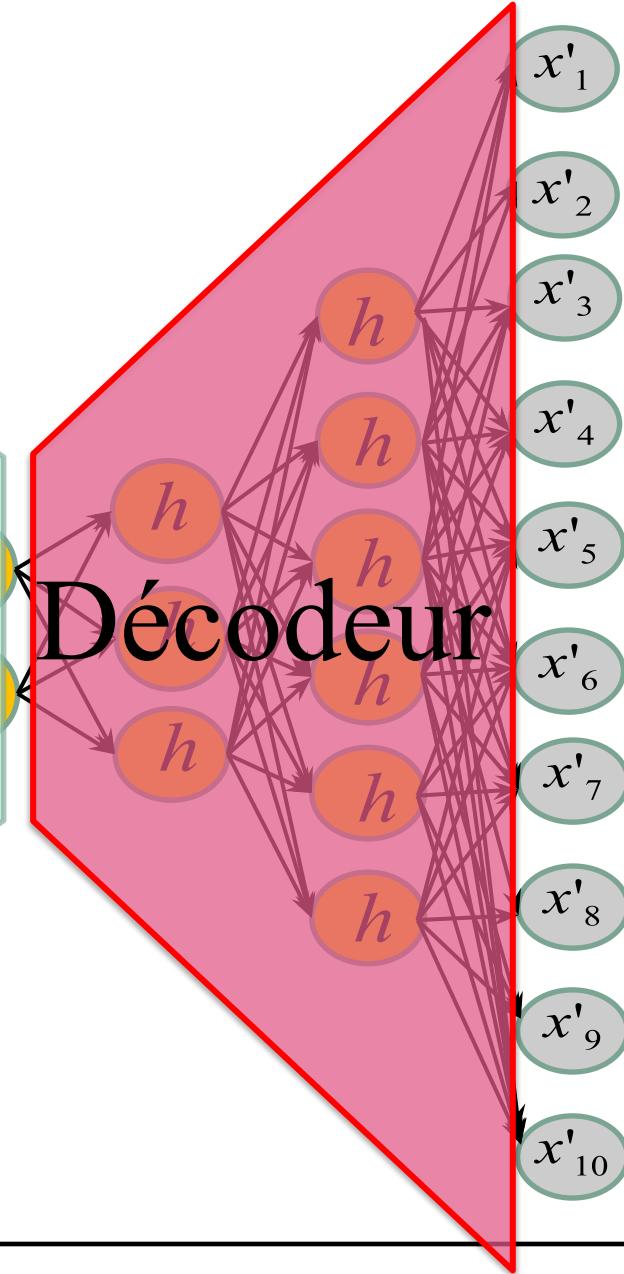
$x$

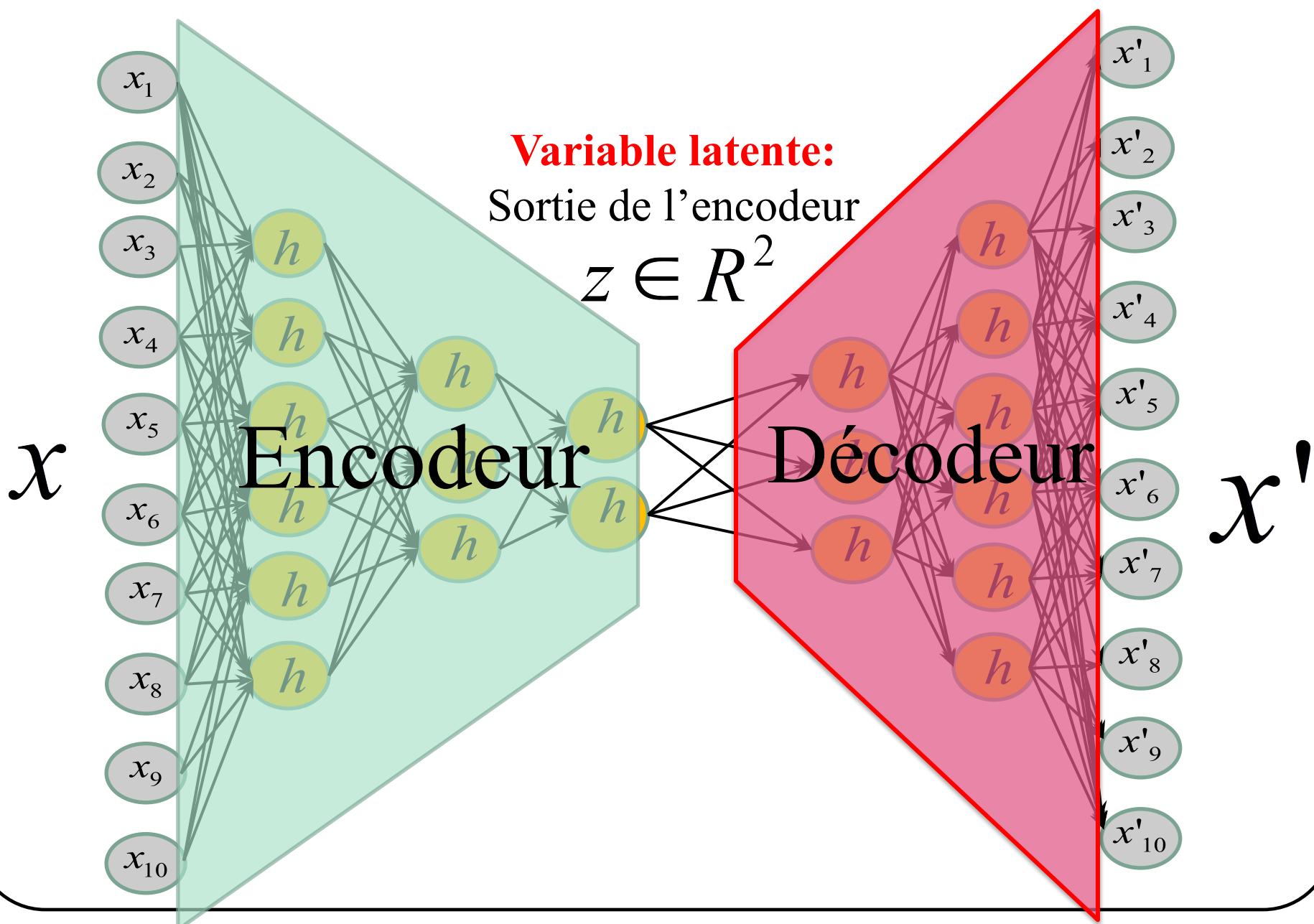


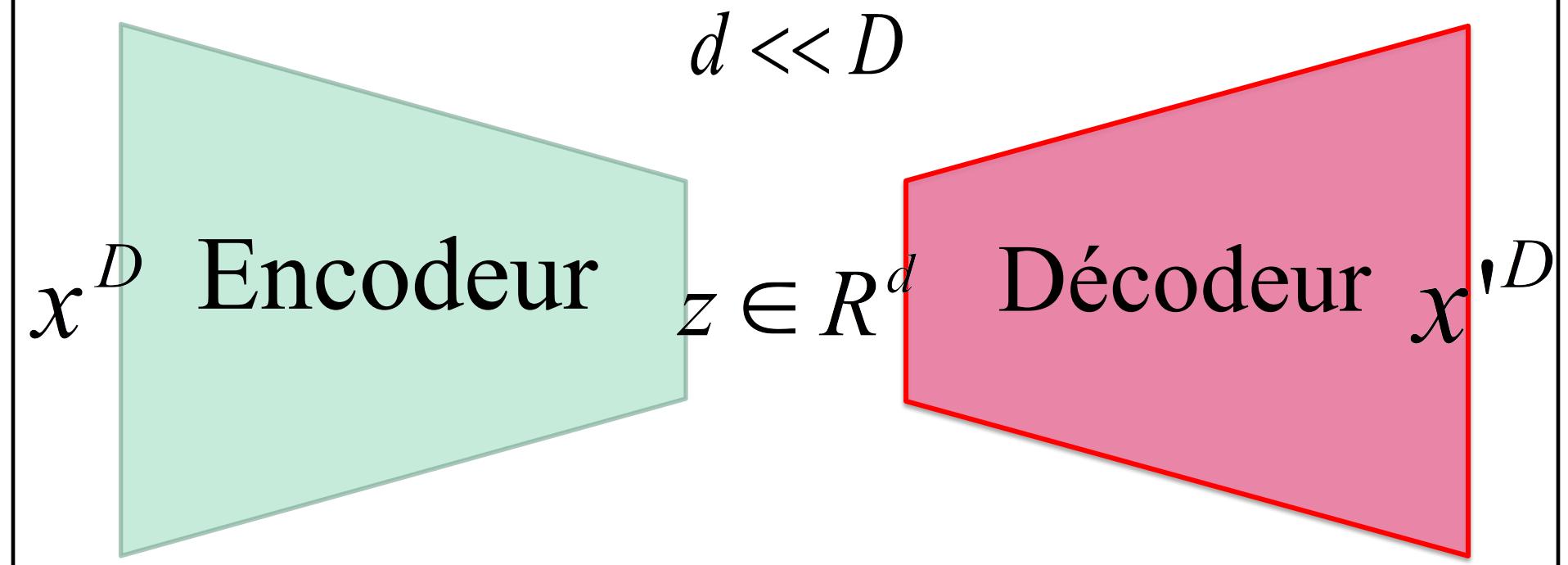
Encodeur

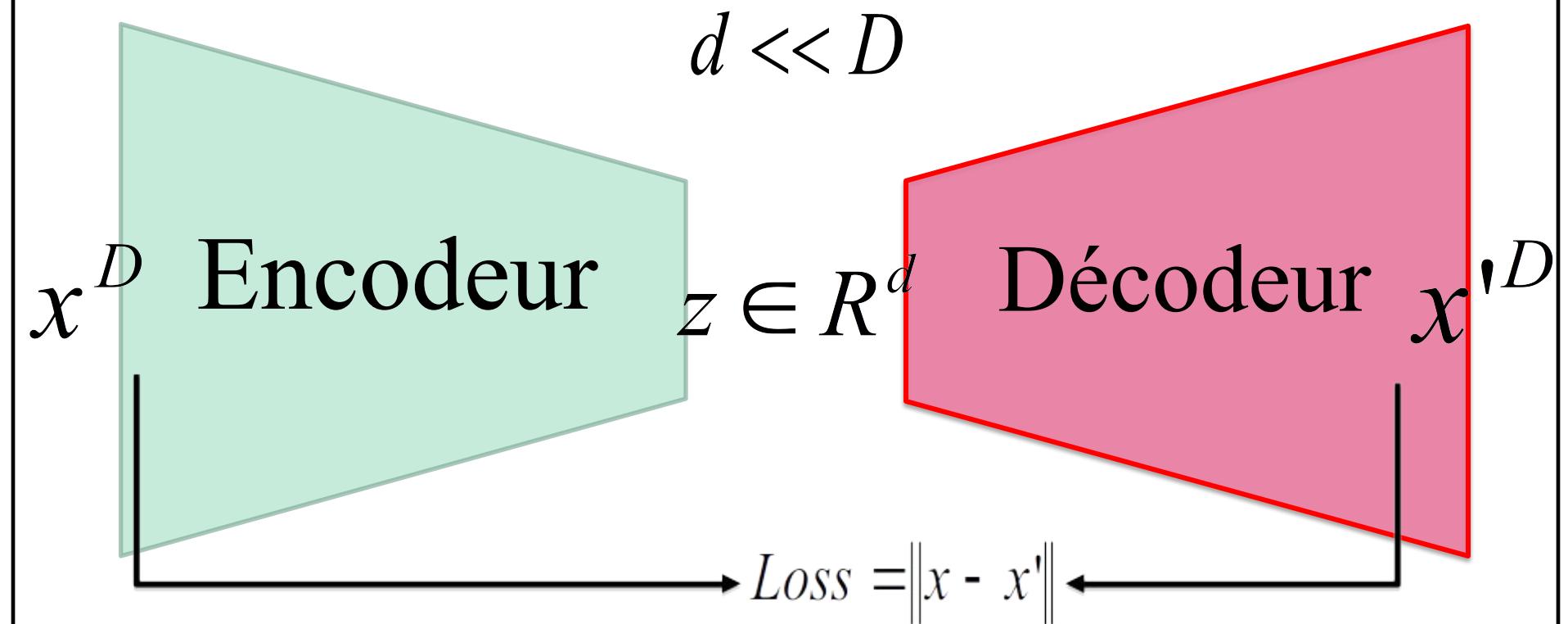
Décodeur

$x'$

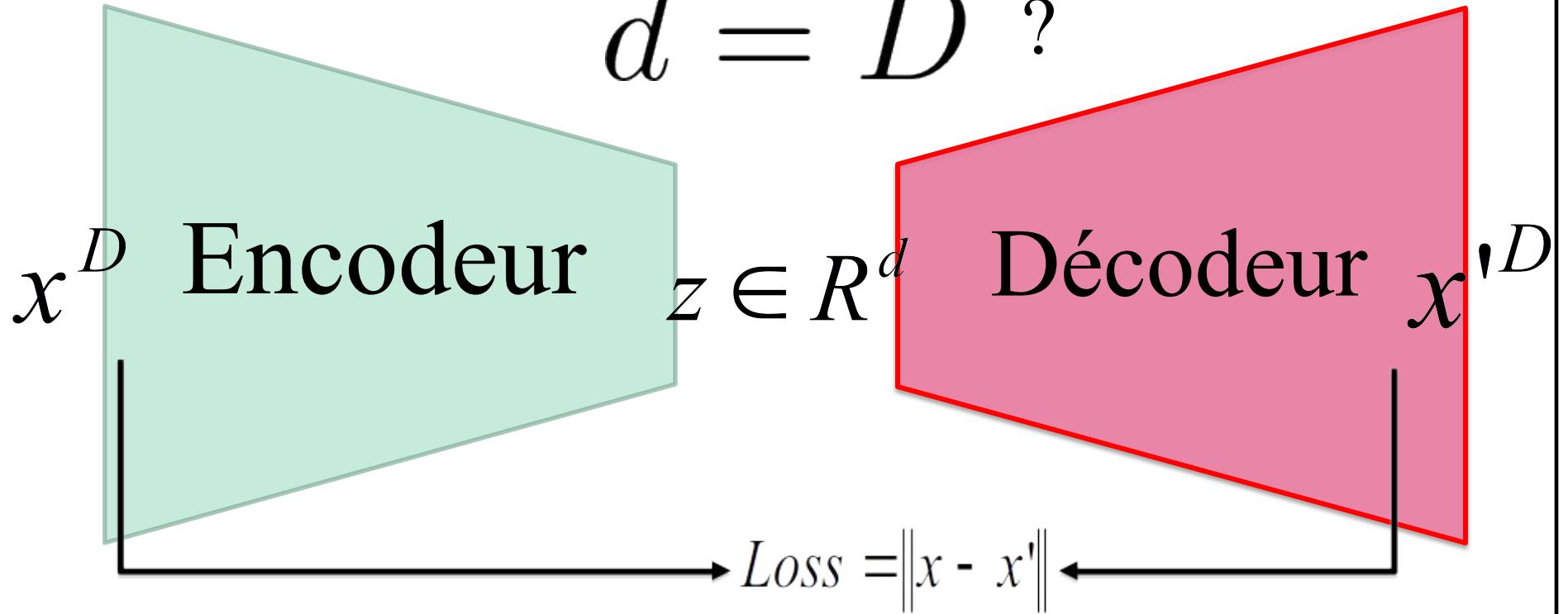


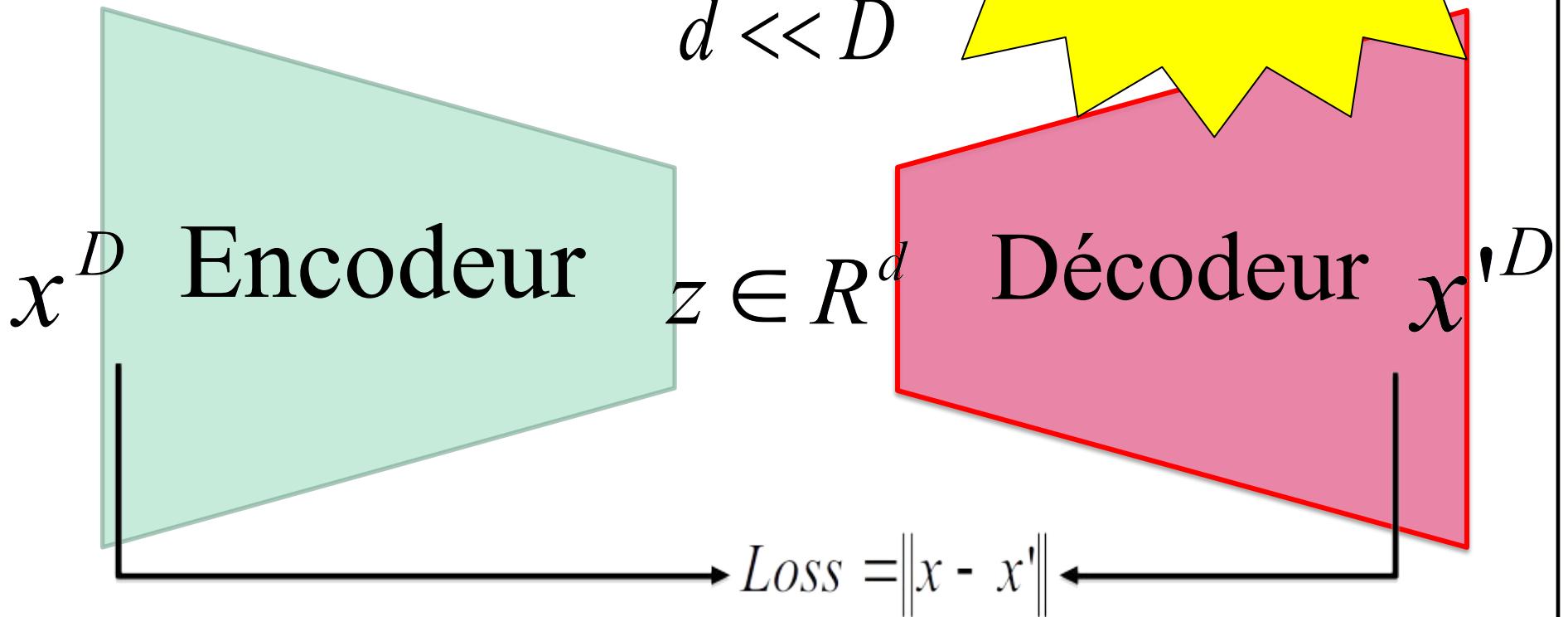




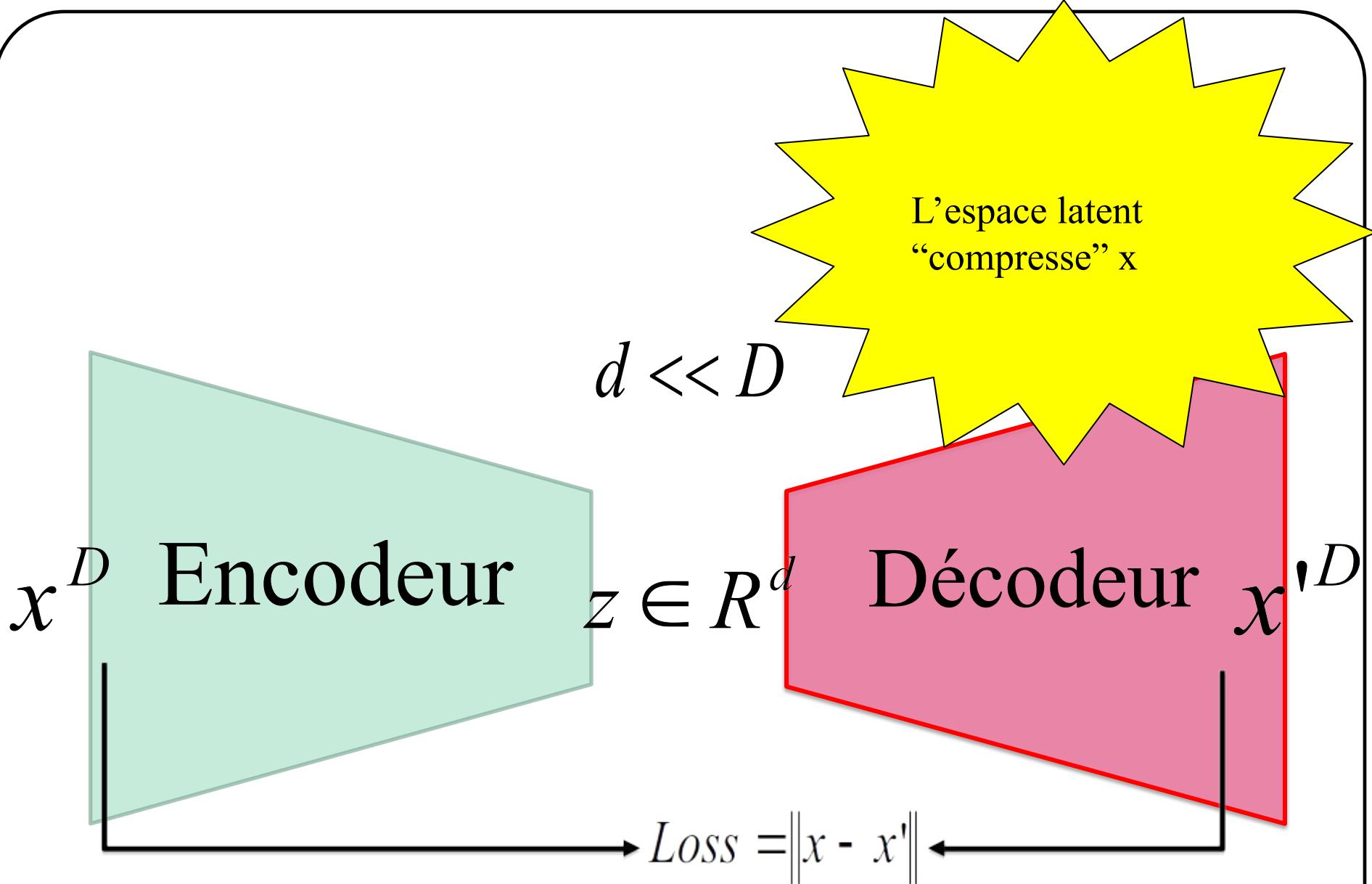


$$d = D ?$$



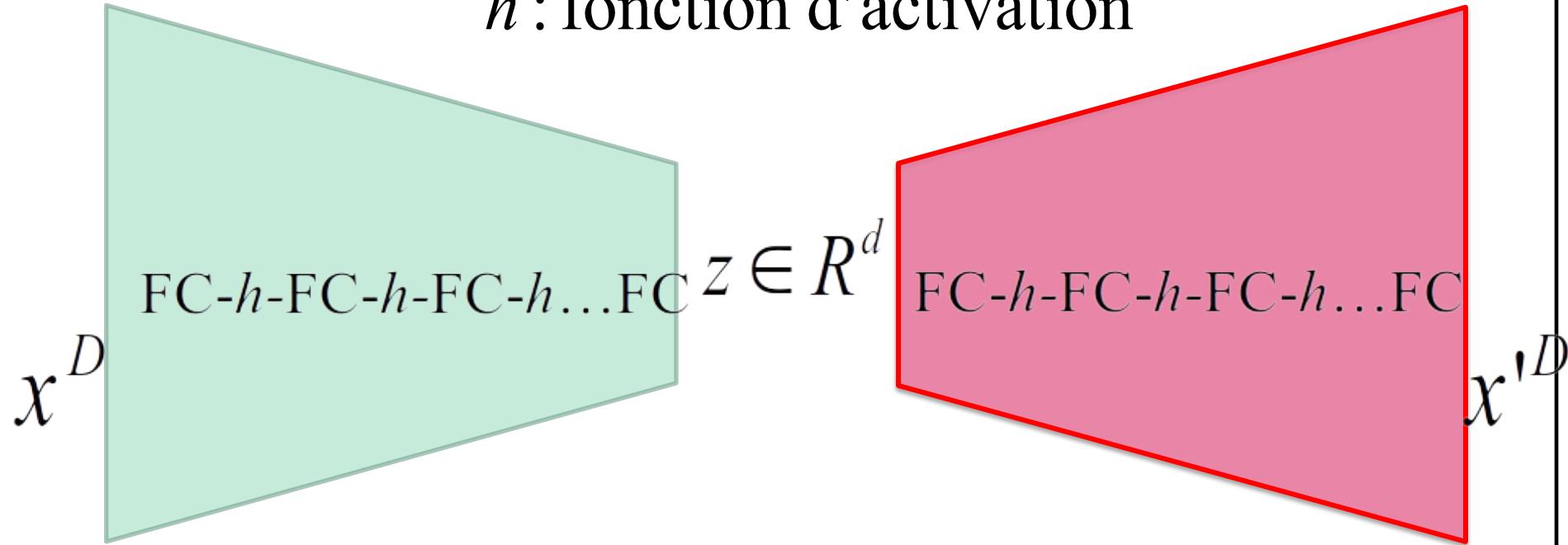


N'apprend pas la  
fonction identité !

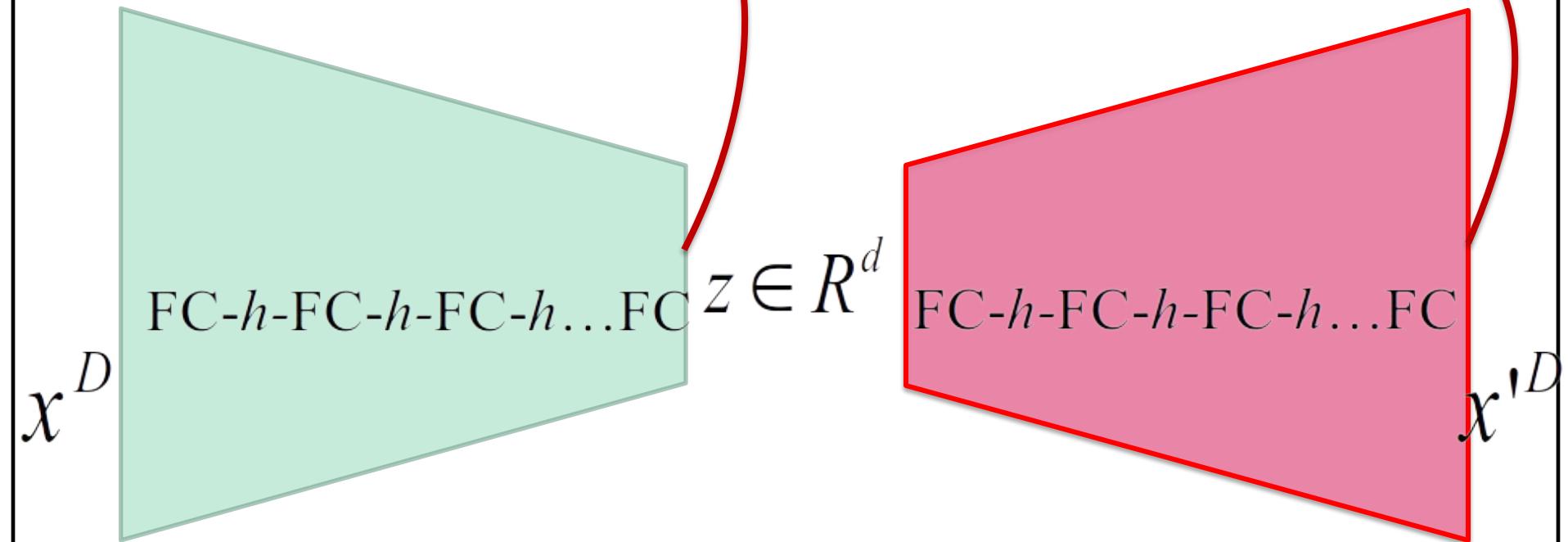


# Couches pleinement connectées

$h$  : fonction d'activation

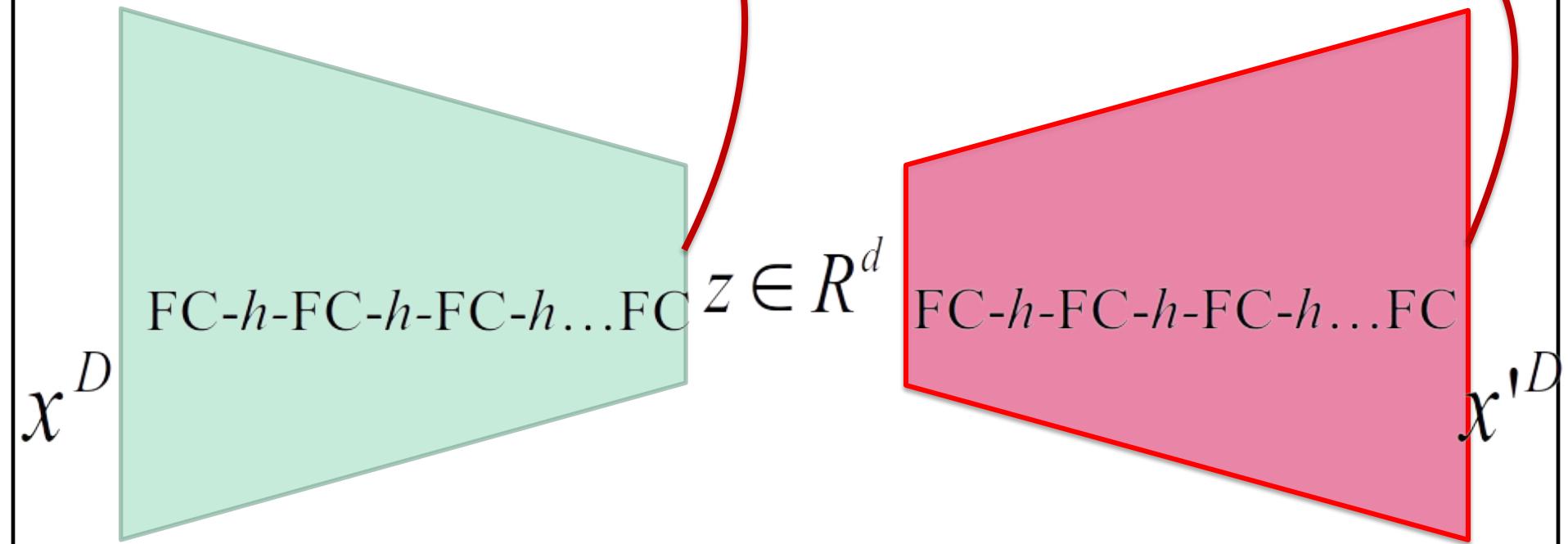


En général...  
pas de fonction d'activation à la sortie  
de l'encodeur et du décodeur



Selon vous, **pourquoi?**

En général...  
pas de fonction d'activation à la sortie  
de l'encodeur et du décodeur

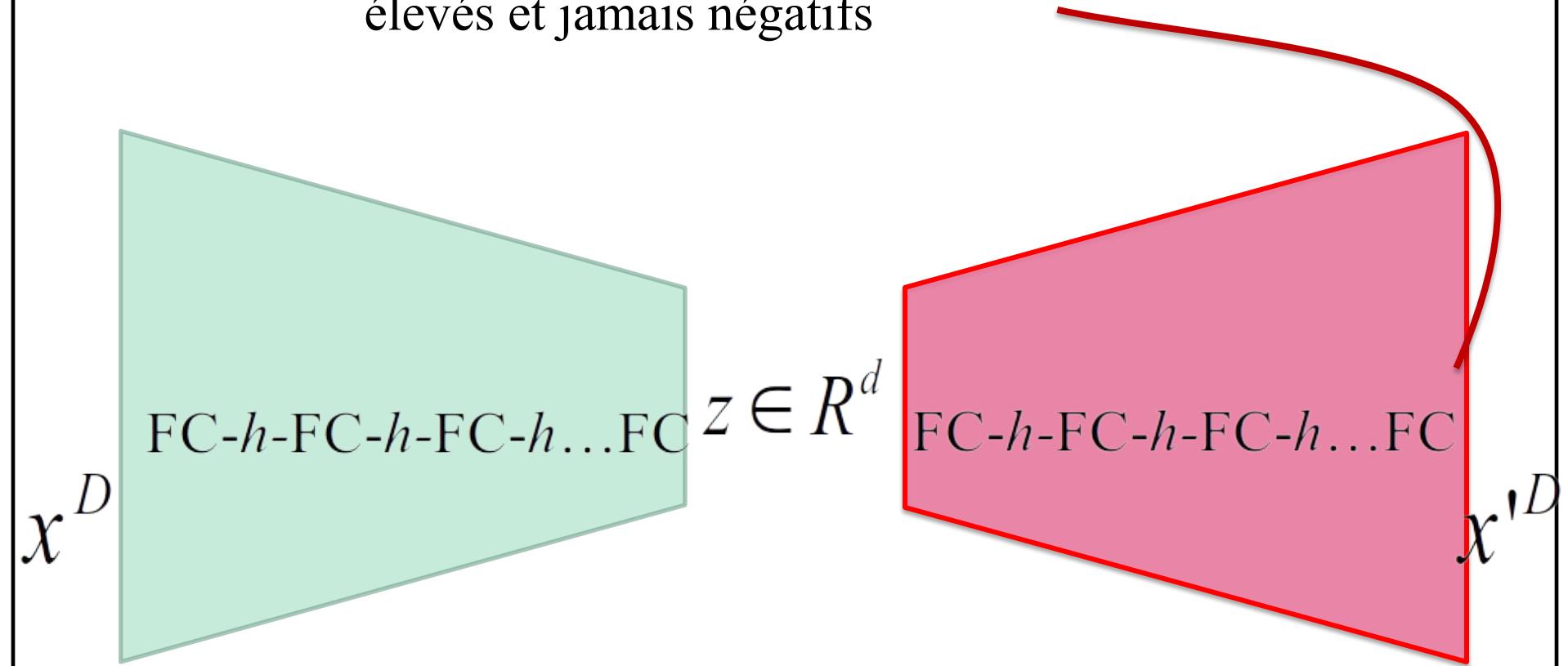


Selon vous, **pourquoi?**

Parce qu'on veut que le  $x$  encodé soit réellement un point  
dans l'espace latent

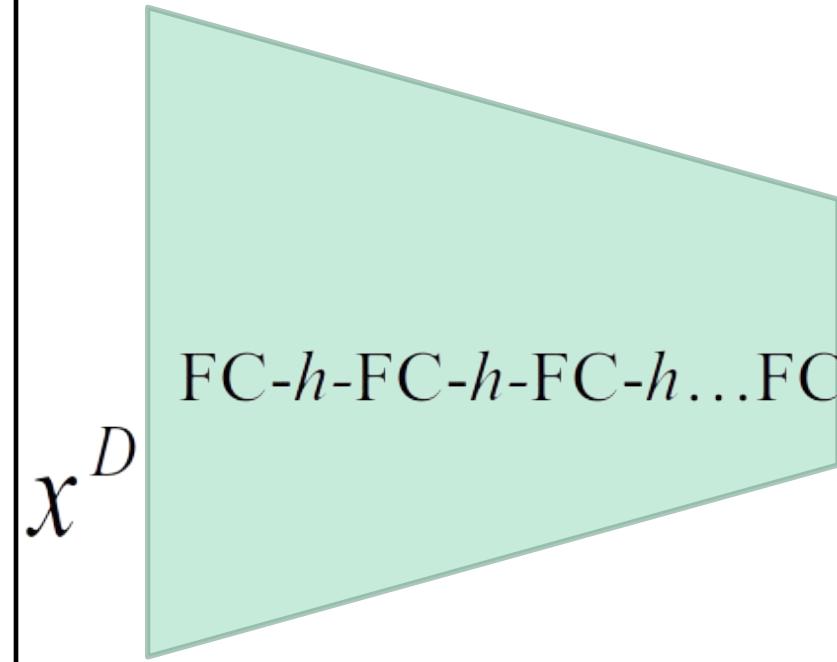
Parfois sigmoïde en sortie lorsque les pixels ont des niveaux de gris entre 0 et 1.

ou ReLU lorsque les niveaux de gris peuvent être élevés et jamais négatifs



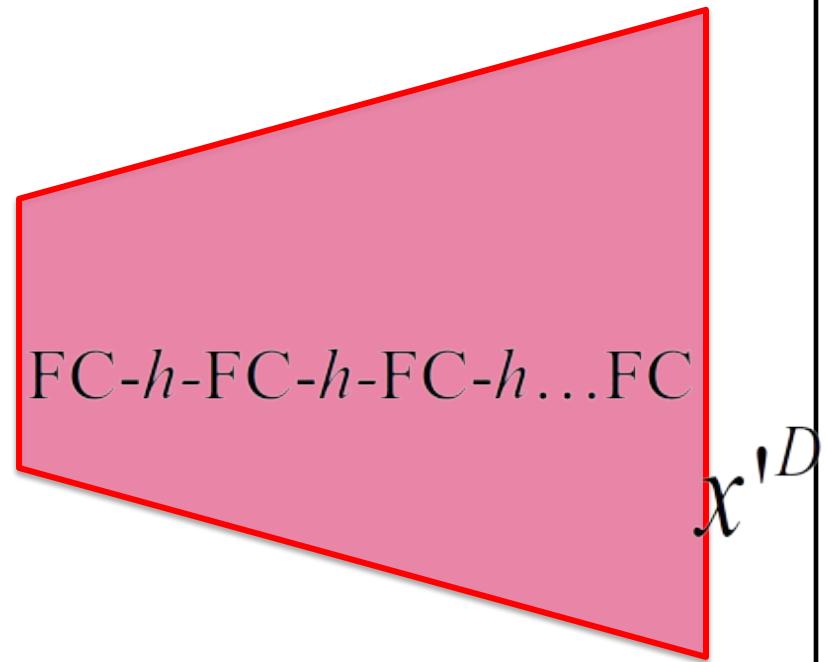
Le nombre de neurones

**Décroît ou se maintient**  
d'une couche à l'autre



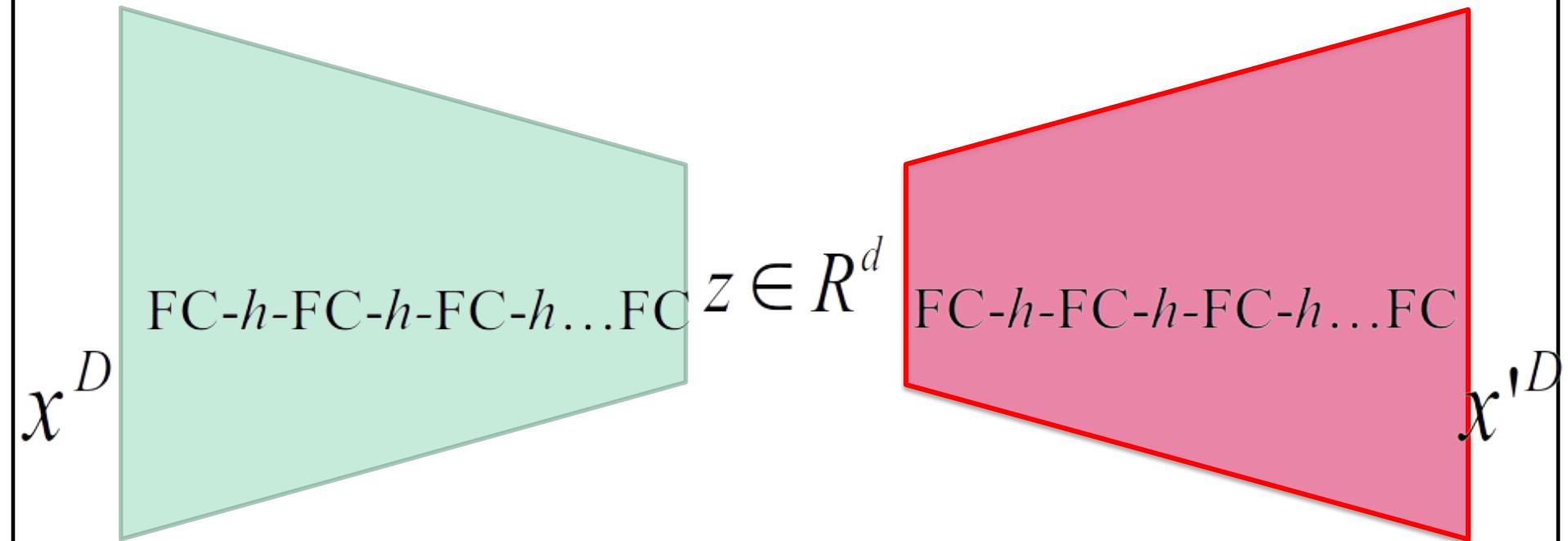
Le nombre de neurones

**Augmente ou se maintient**  
d'une couche à l'autre



Très souvent...

La structure de l'encodeur est le dual de celle du décodeur



# Autoencodeur jouet de MNIST

```
class autoencoder(nn.Module):

    def __init__(self):
        super(autoencoder, self).__init__()
        self.encoder = nn.Sequential(
            nn.Linear(28 * 28, 128), nn.ReLU(True),
            nn.Linear(128, 64), nn.ReLU(True),
            nn.Linear(64, 12), nn.ReLU(True),
            nn.Linear(12, 2))

        self.decoder = nn.Sequential(
            nn.Linear(2, 12), nn.ReLU(True),
            nn.Linear(12, 64), nn.ReLU(True),
            nn.Linear(64, 128), nn.ReLU(True),
            nn.Linear(128, 28 * 28))

    def forward(self, x):
        z = self.encoder(x)
        x_prime = self.decoder(z)
        return x_prime
```

Espace latent 2D

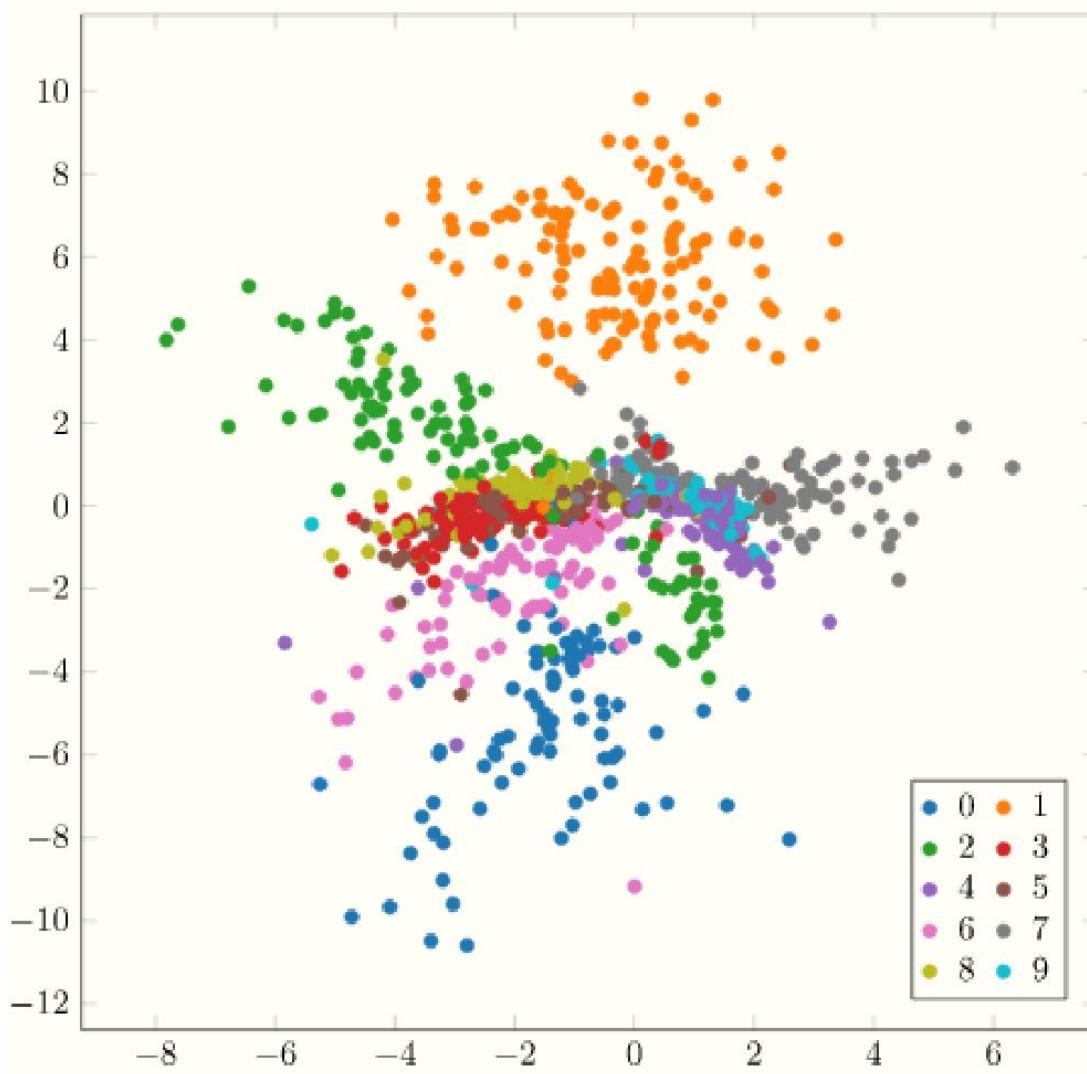
# Autoencodeur jouet de MNIST

```
class autoencoder(nn.Module):  
  
    def __init__(self):  
        super(autoencoder, self).__init__()  
  
        self.encoder = nn.Sequential(  
            nn.Linear(28 * 28, 128), nn.ReLU(True),  
            nn.Linear(128, 64), nn.ReLU(True),  
            nn.Linear(64, 12), nn.ReLU(True),  
            nn.Linear(12, 2))  
  
        self.decoder = nn.Sequential(  
            nn.Linear(2, 12), nn.ReLU(True),  
            nn.Linear(12, 64), nn.ReLU(True),  
            nn.Linear(64, 128), nn.ReLU(True),  
            nn.Linear(128, 28 * 28))  
  
    def forward(self, x):  
        z = self.encoder(x)  
        x_prime = self.decoder(z)  
  
        return x_prime
```

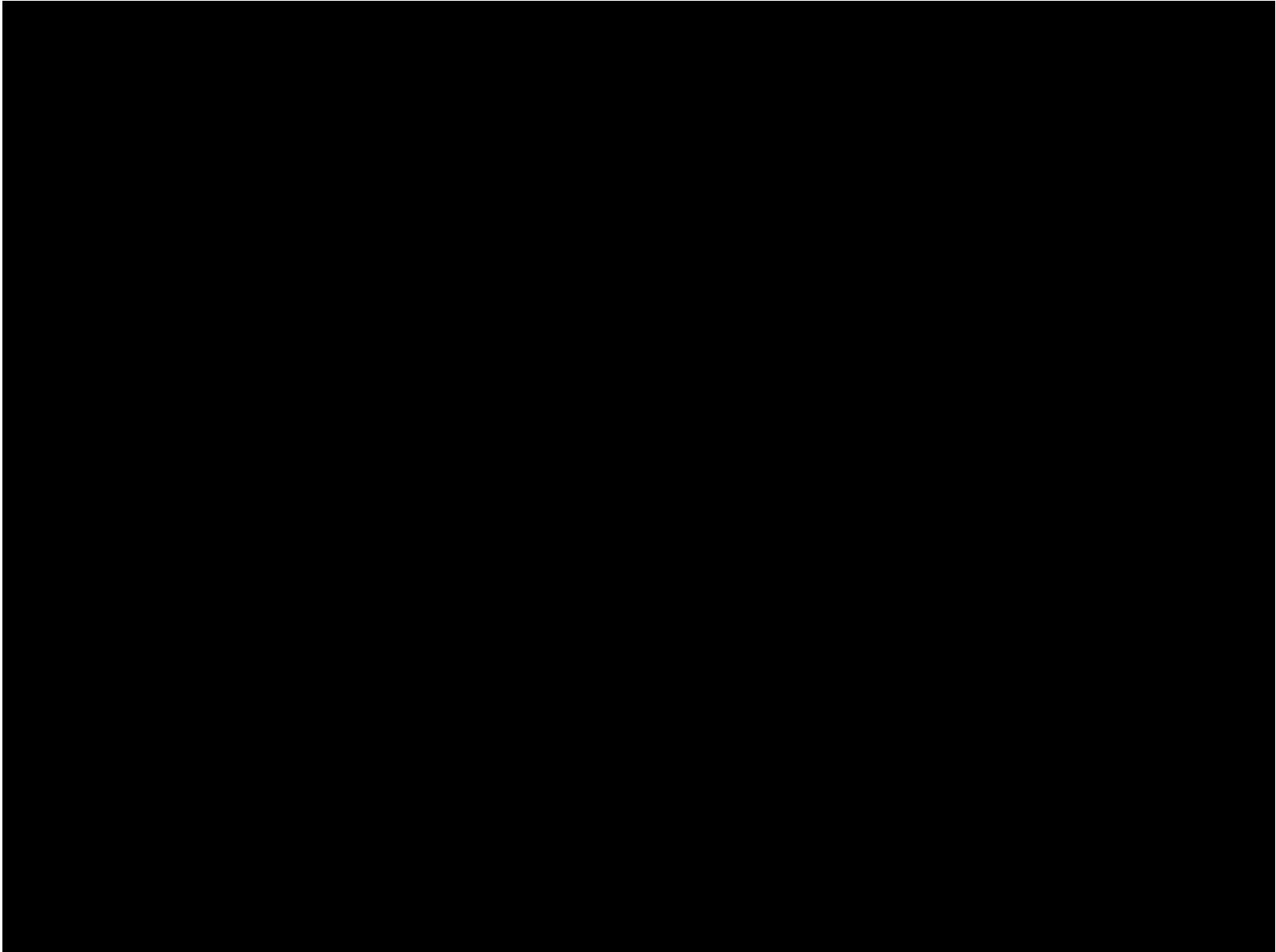
symétrie

# Visualisation de l'espace latent pour 1000 images MNIST

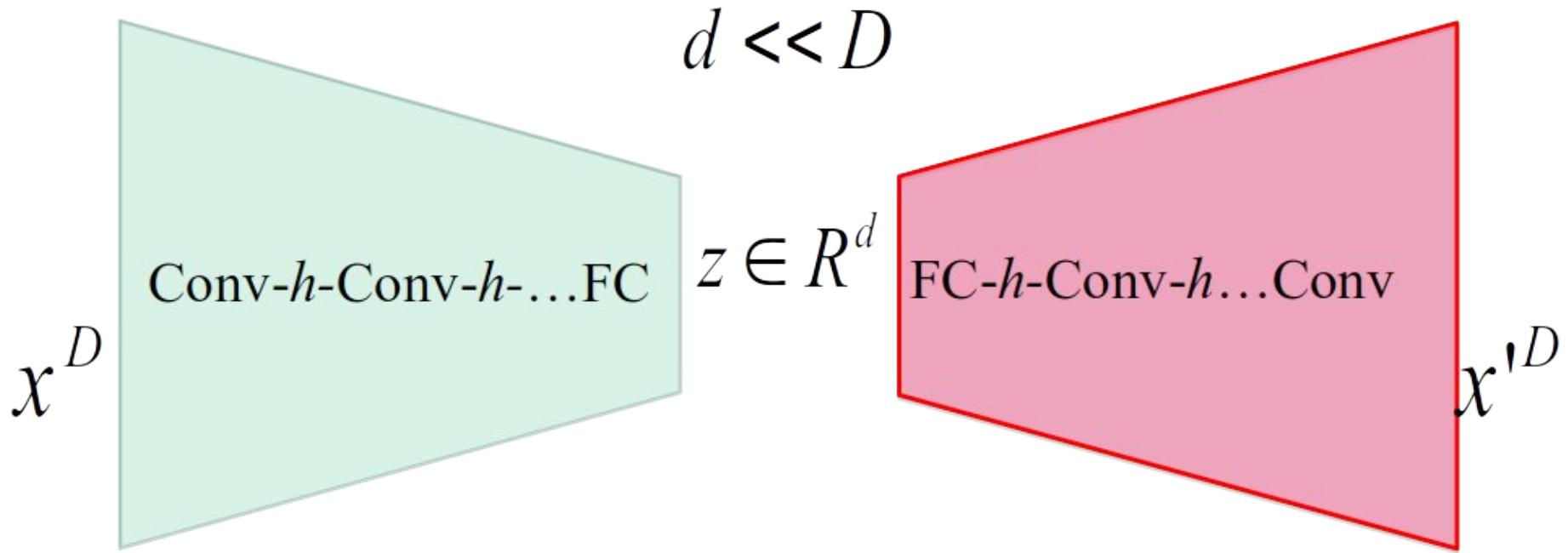
chaque image correspond à 1 point 2D



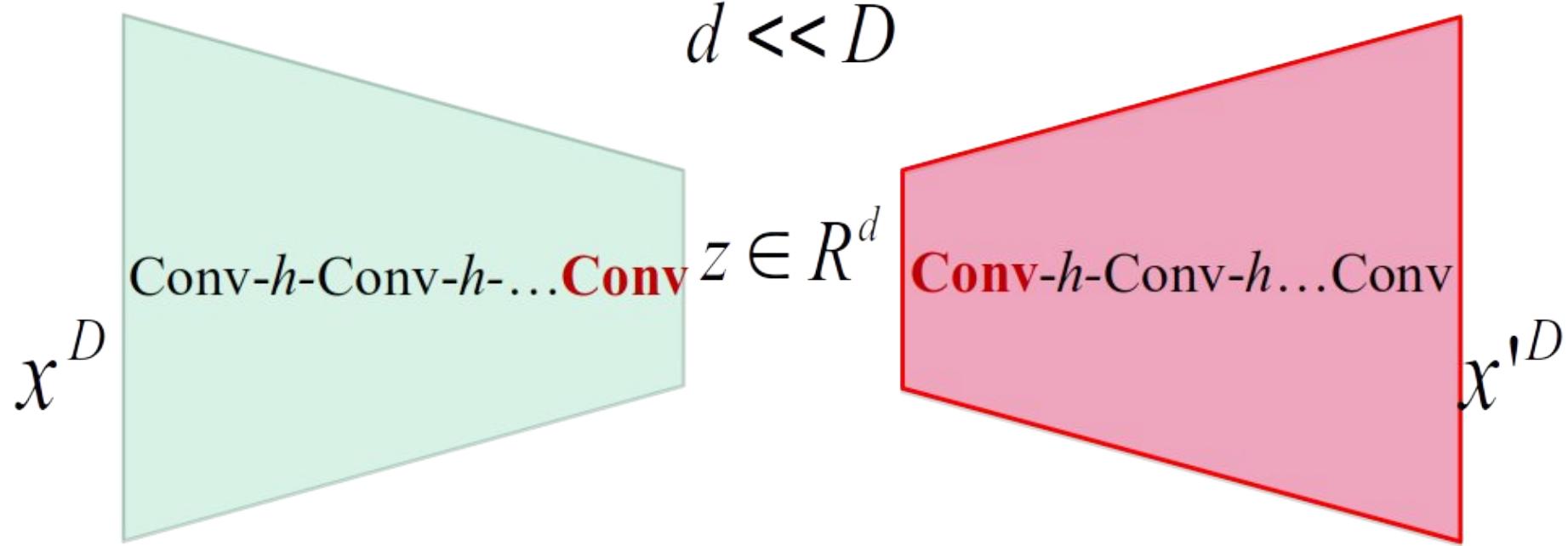
Générer des images avec le décodeur.



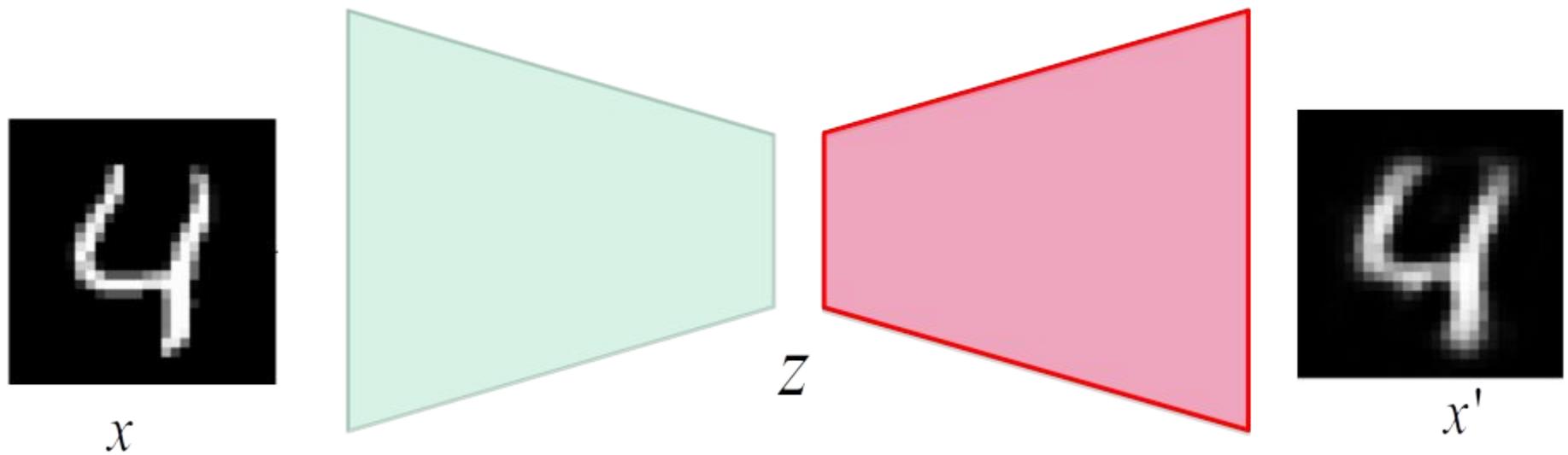
# Couches convolutives



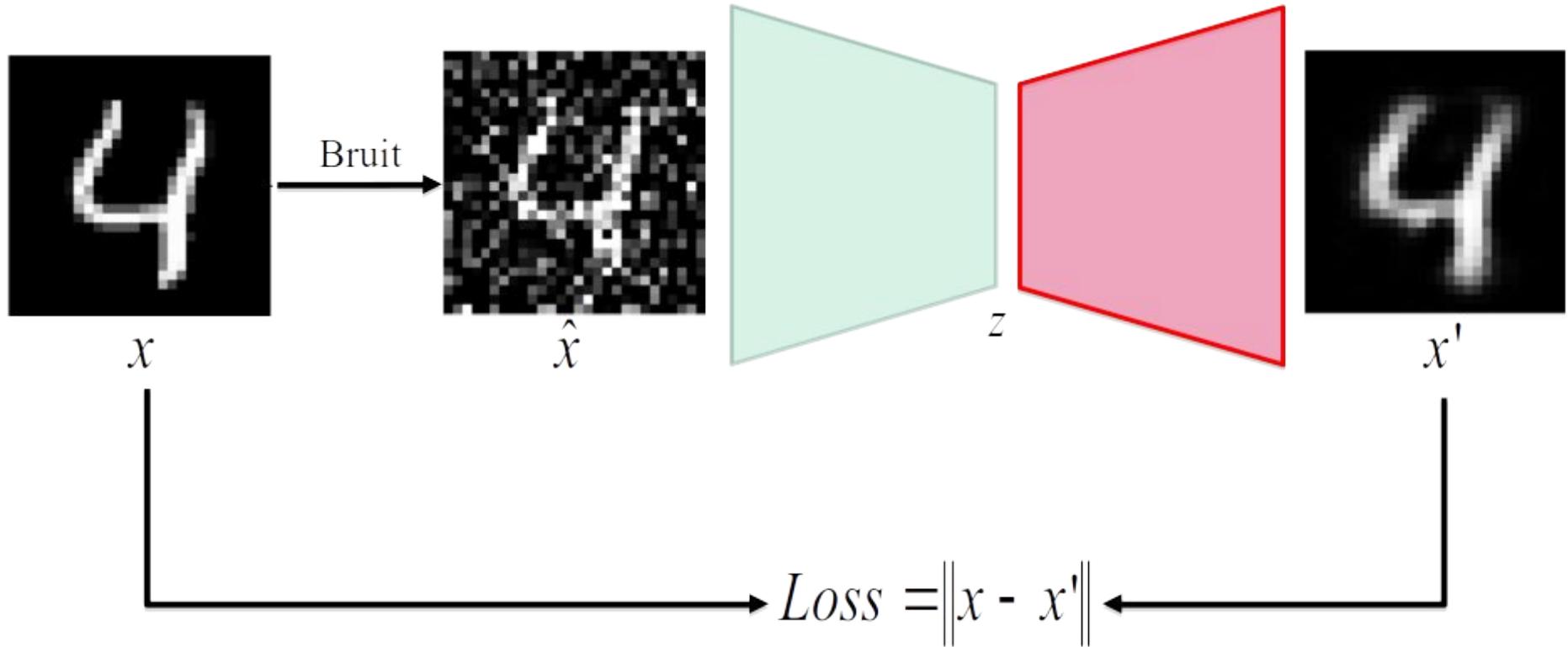
# Autoencodeur pleinement convolutif



## Autoencodeur de base

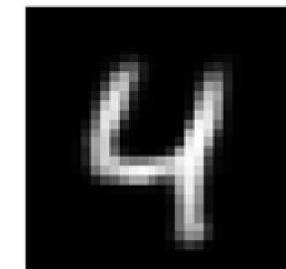
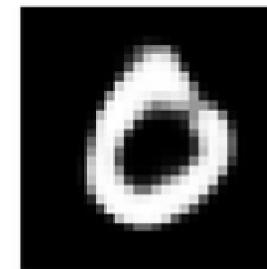
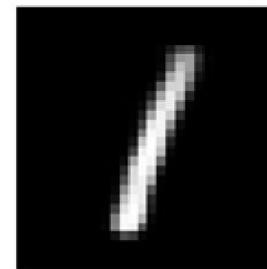
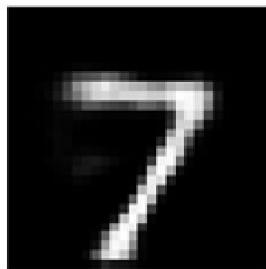
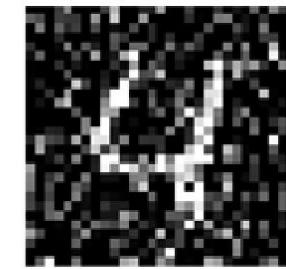
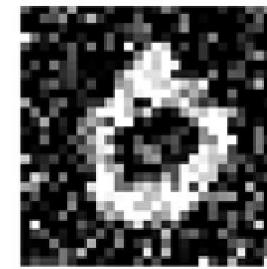
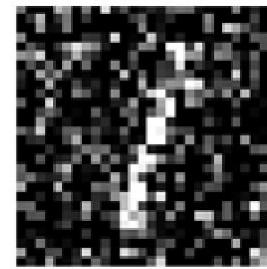
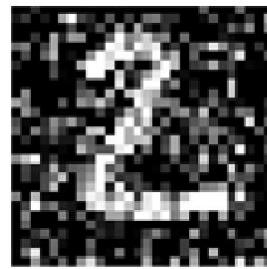
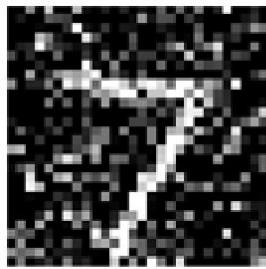


## Autoencodeur pour débruitage



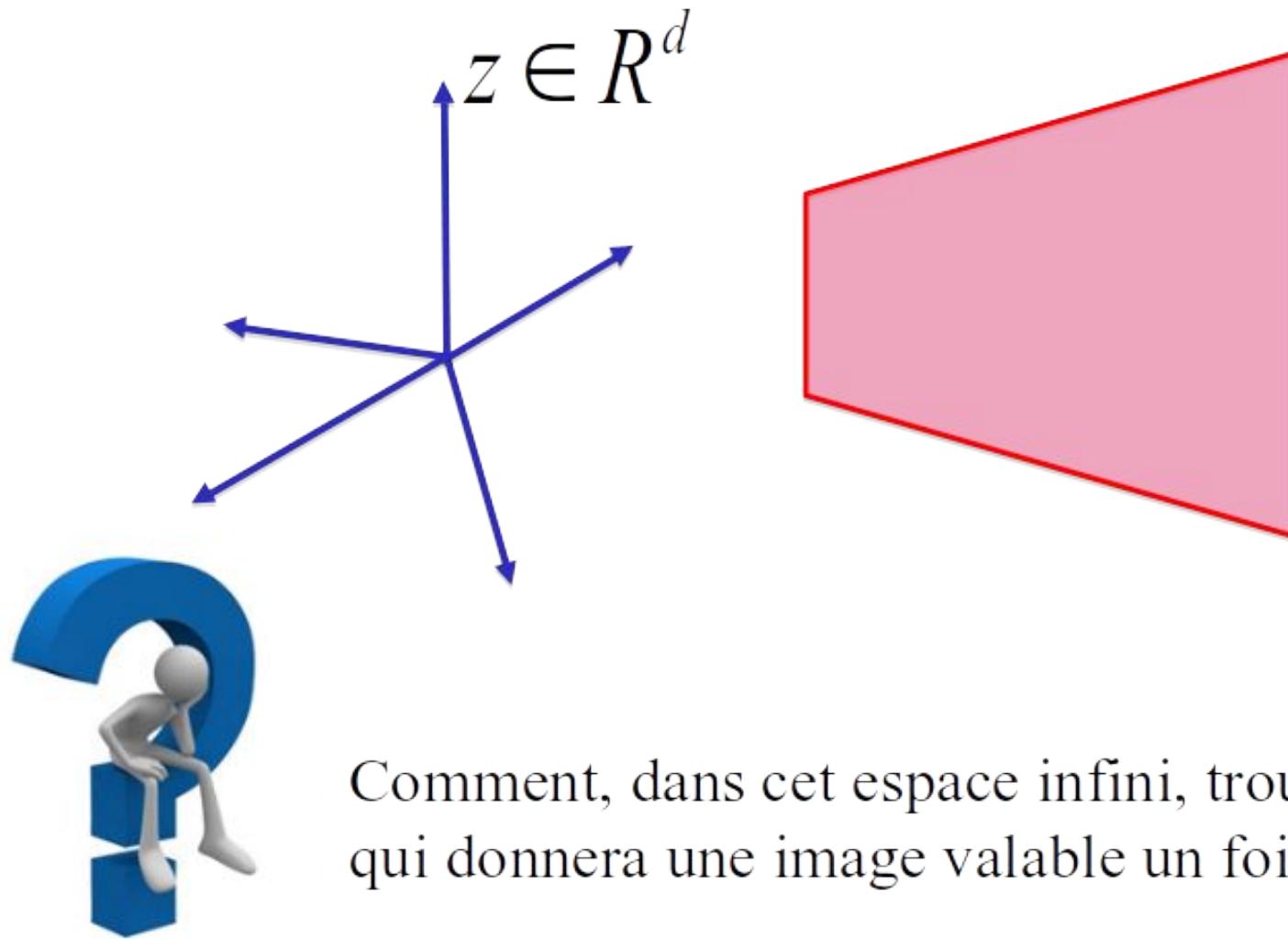
# Une fois entraîné...

Signal bruité en entrée



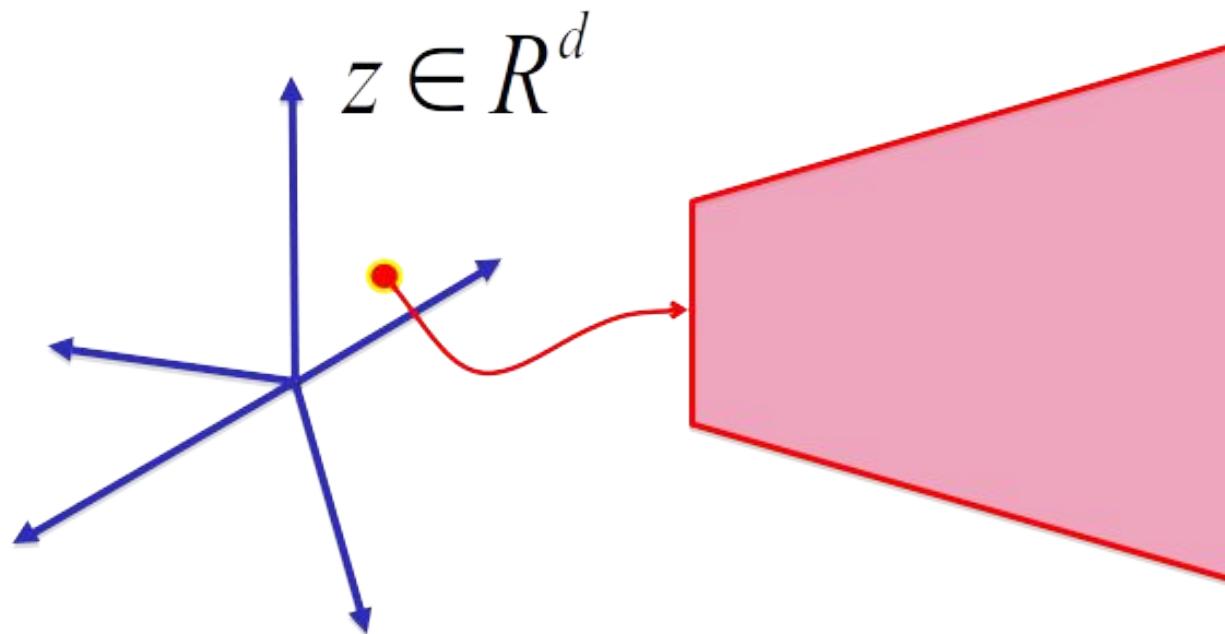
Signal reconstruit et débruité

En général, l'espace latent possède entre 16 et 128 dimensions.  
Ça peut être parfois plus, et parfois moins.

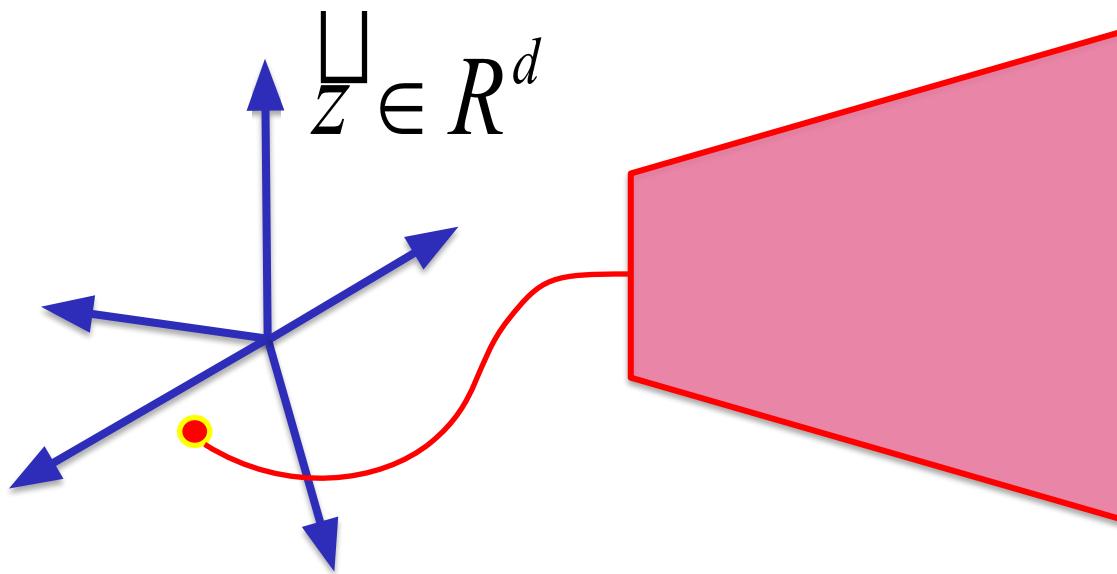


Comment, dans cet espace infini, trouver un point  $z$  qui donnera une image valable un fois décodée?

Avec de la chance, on peut sélectionner un point au hasard et reproduire une « bonne » image (ici une image « MNIST »)



Malheureusement, la vaste majorité du temps, on reproduira du bruit



Au lieu d'apprendre à **reproduire**  
**un signal d'entrée...**

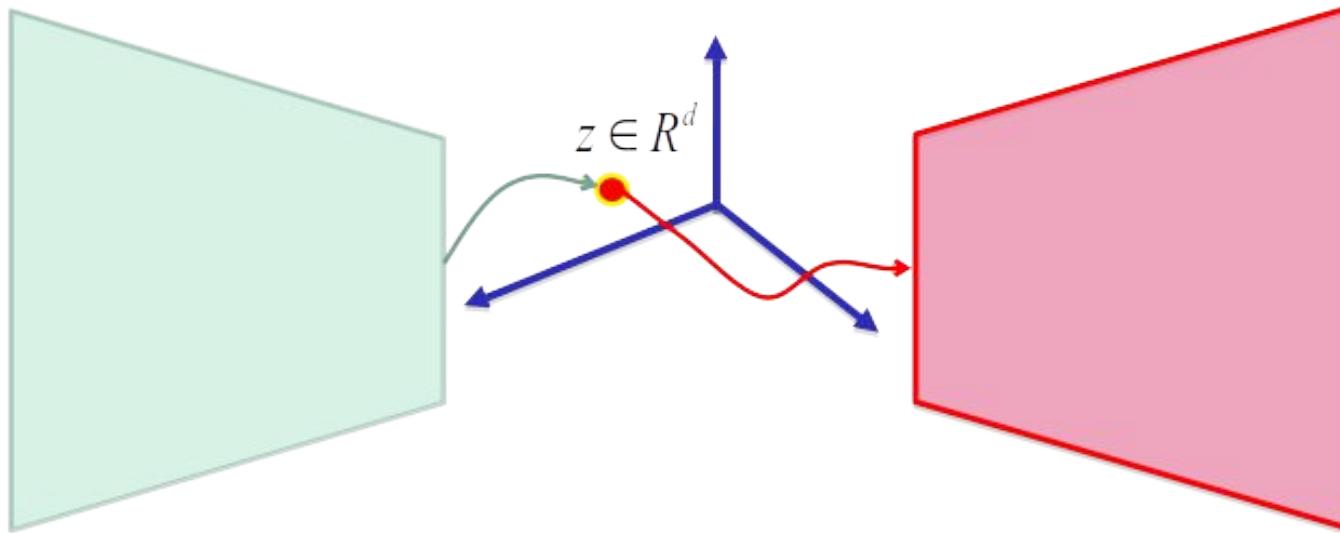


Apprendre à reproduire une **distribution**  $p(\mathbb{z})$   
**connue** de sorte qu'un **point échantillonné**  
**et décodé** de cette distribution correspond à  
un signal reconstruit valable

## Autoencodeur de base



$x$

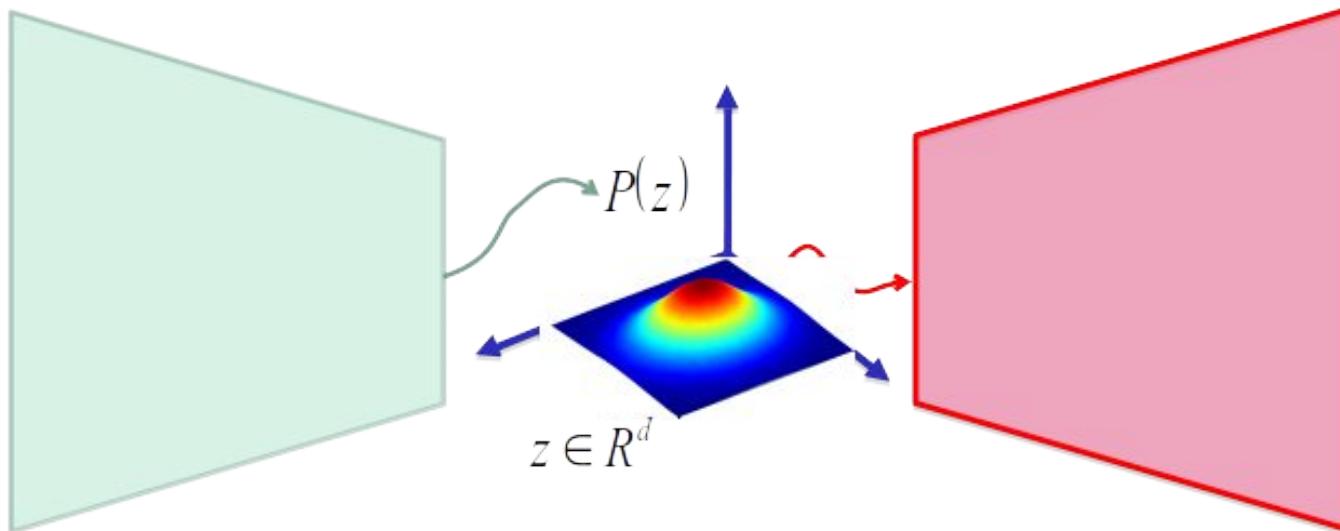


$x'$

## Autoencodeur variationnel



$x$

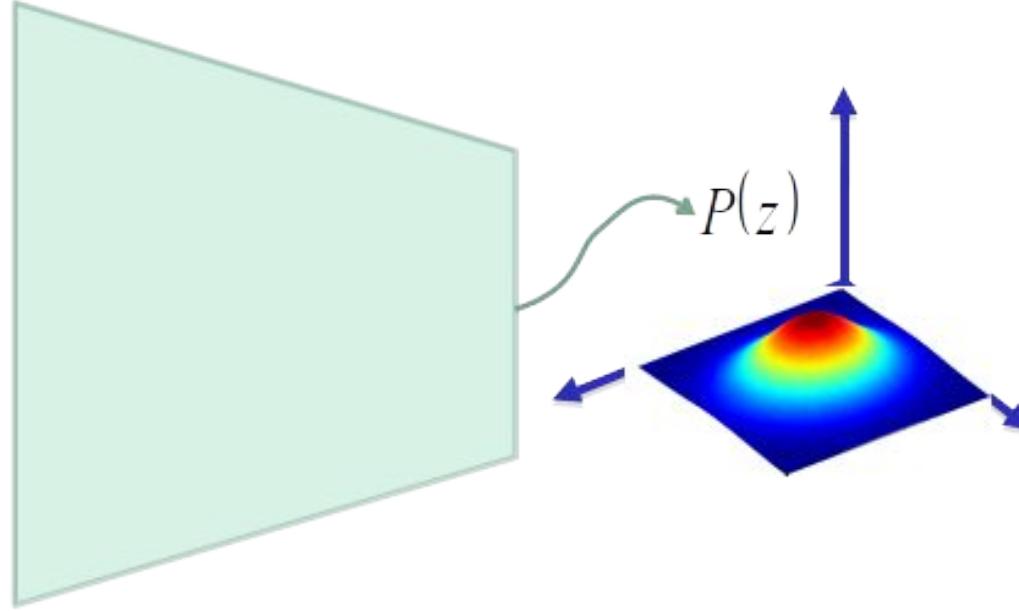


$x'$

# Autoencodeur variationnel



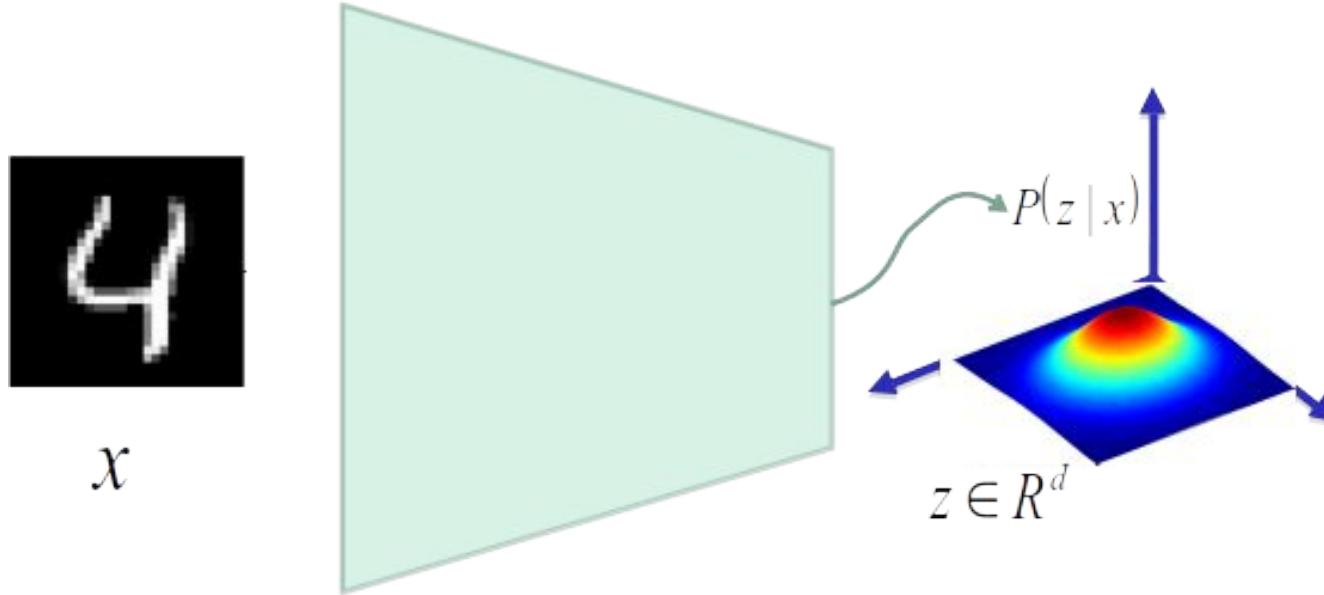
$x$



L'encodeur produit une distribution  $P(z)$   
et non juste un point  $z$

# Autoencodeur variationnel

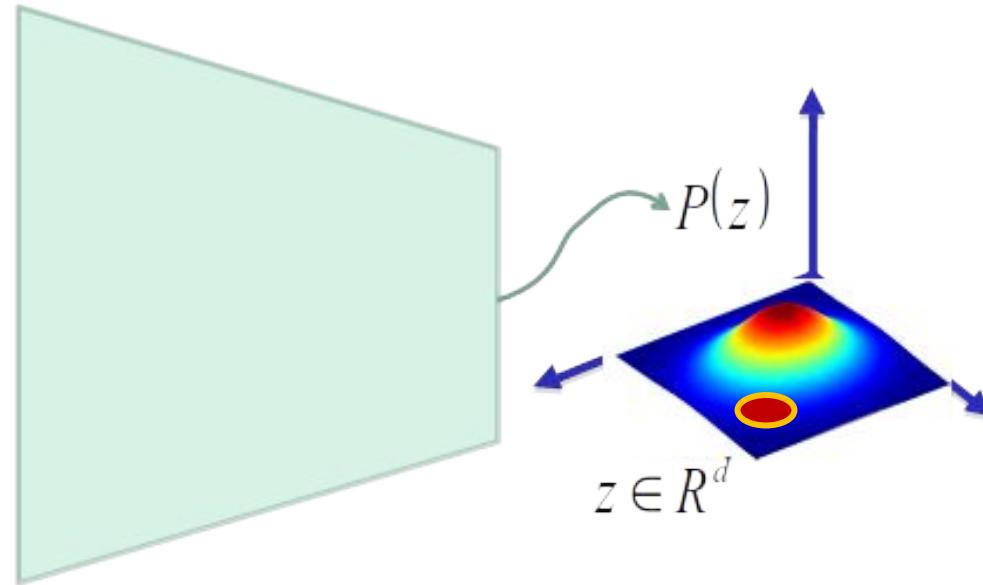
## Remarque 1



Puisque la distribution de  $z$  dépend de  $x$   
on dira que la distribution apprise est

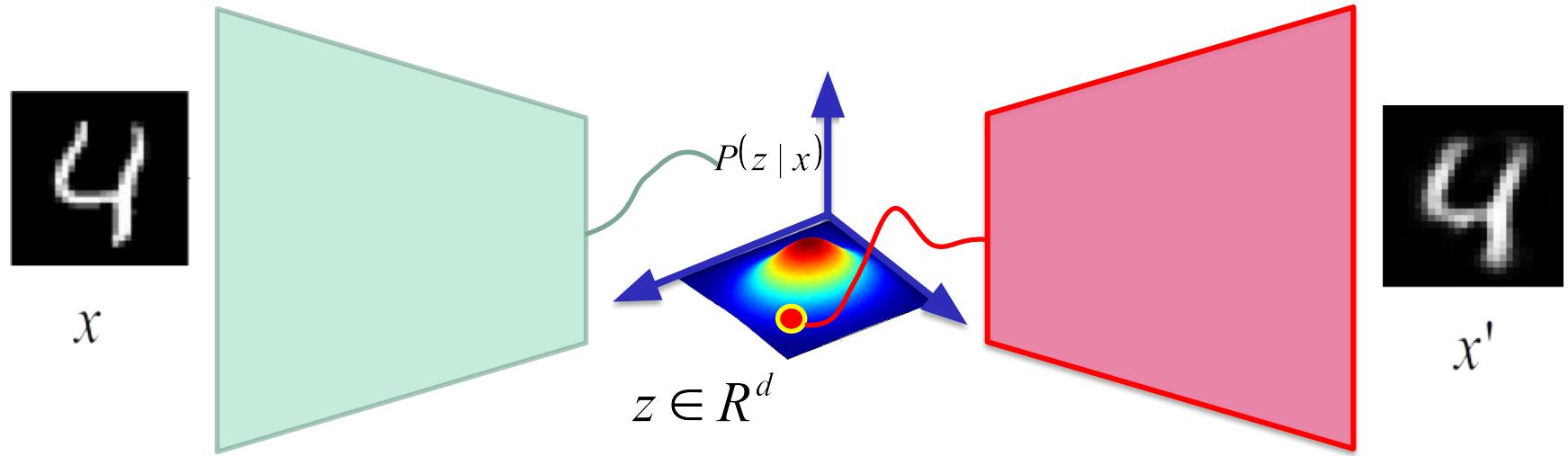
$$P(z | x)$$

# Autoencodeur variationnel



On échantillonne un  $z \sim p(z)$  au hasard

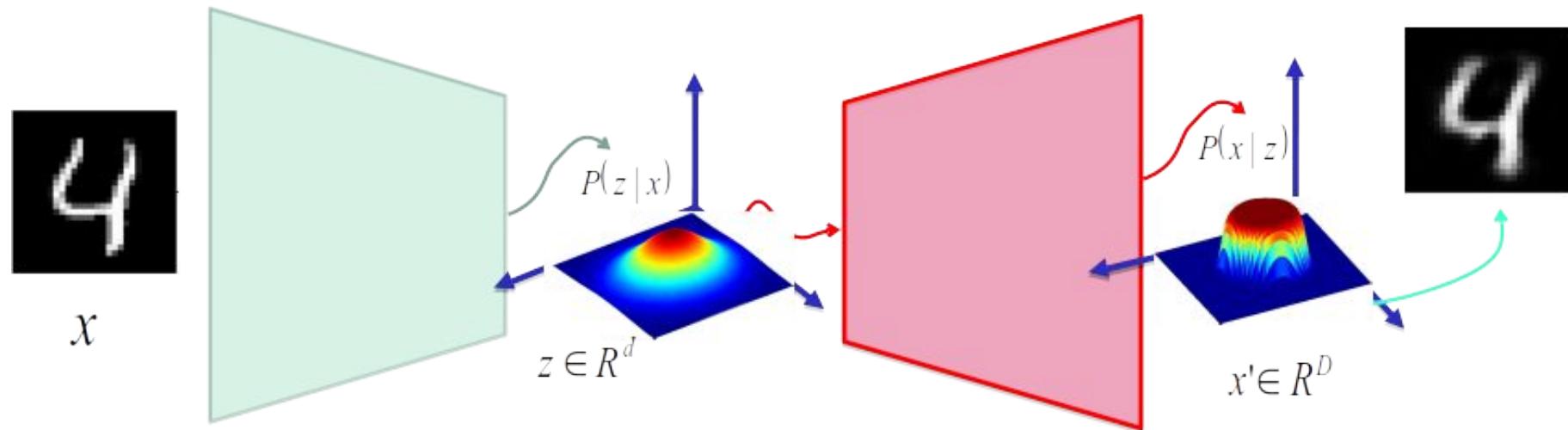
# Autoencodeur variationnel



On reconstruit  $x'$

## Autoencodeur variationnel

### Remarque 2



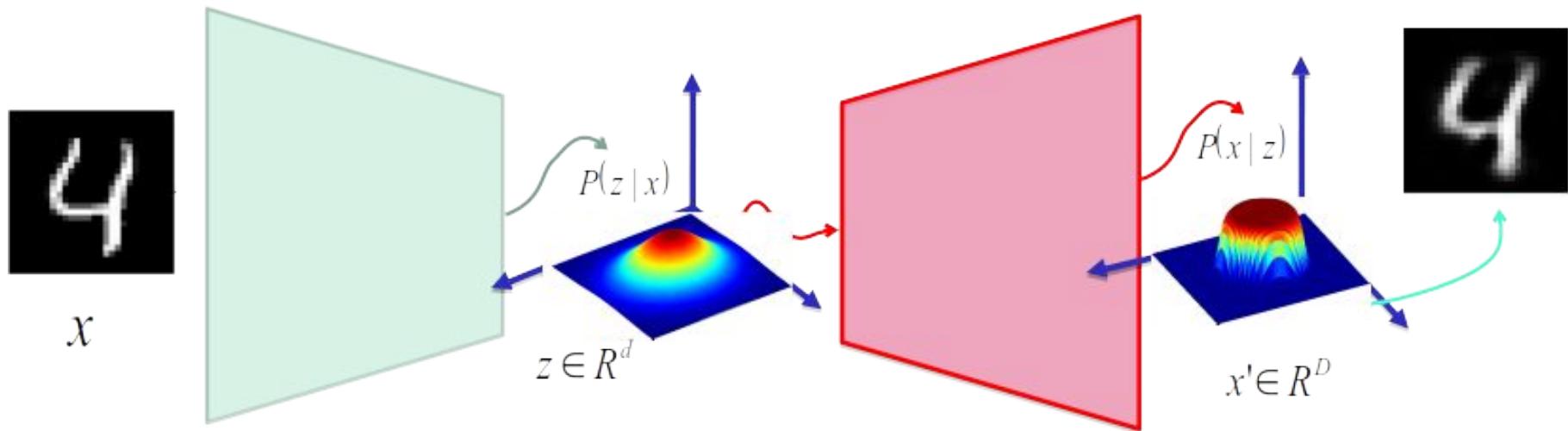
Le **décodeur** peut également produire une distribution de probabilités

$$P(x|z)$$

et  $x'$  est un point échantillonné au hasard de  $P(x|z)$

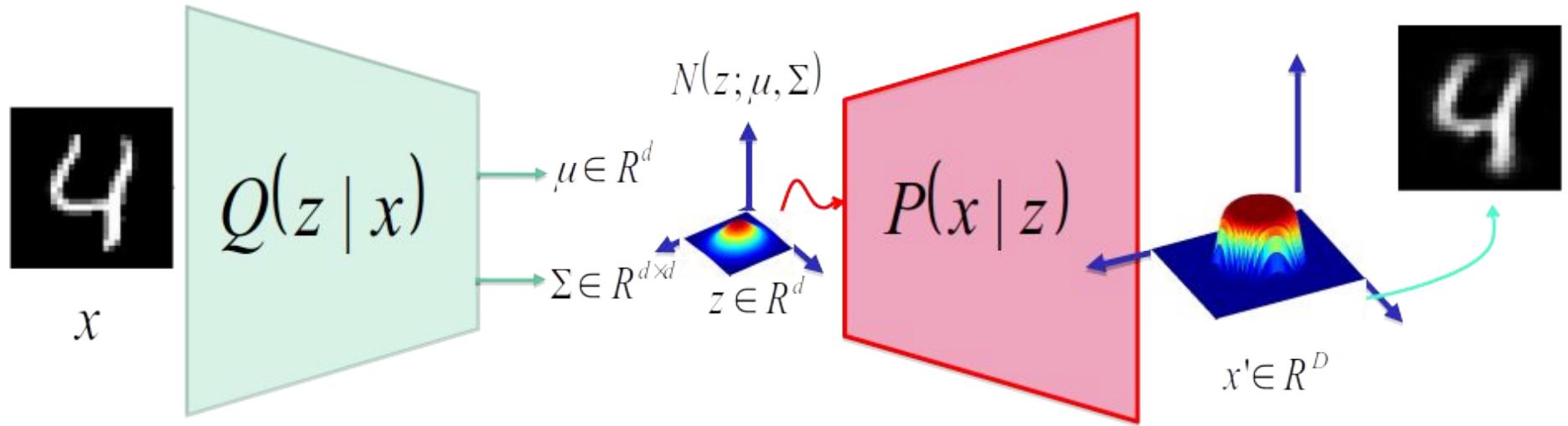
# Autoencodeur variationnel

## Remarque 3

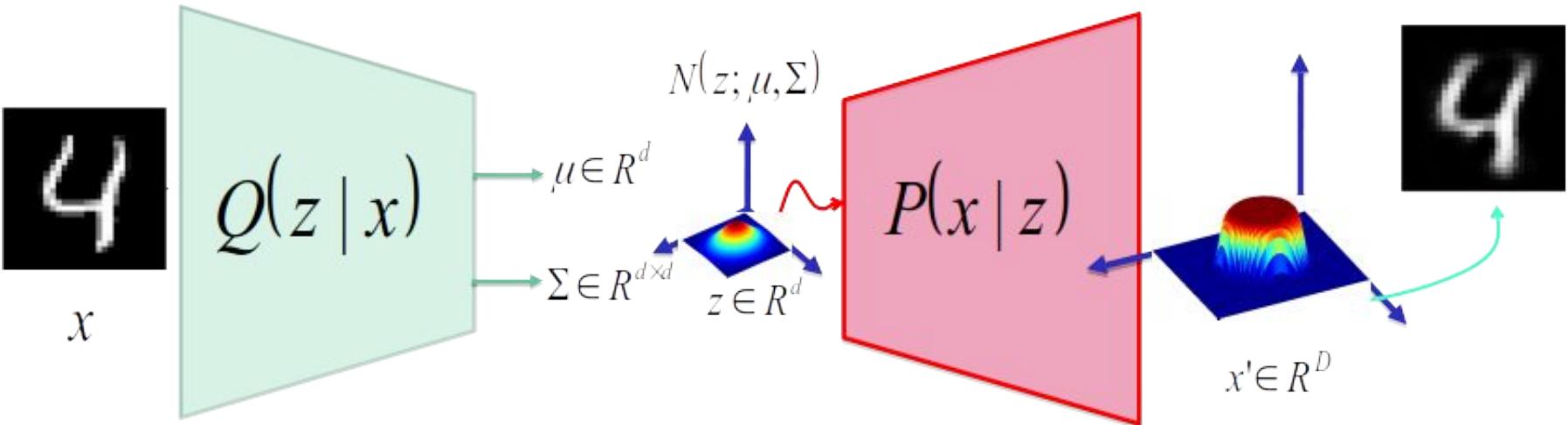


La distribution  $P(z | x)$  peut être très complexe et difficile à échantillonner, on va donc l'approximer par une distribution plus simple... une gaussienne

$$Q(z | x) \approx P(z | x)$$


$$Q(z | x) \sim \text{Gaussienne}$$

## Autoencodeur variationnel Remarque 4

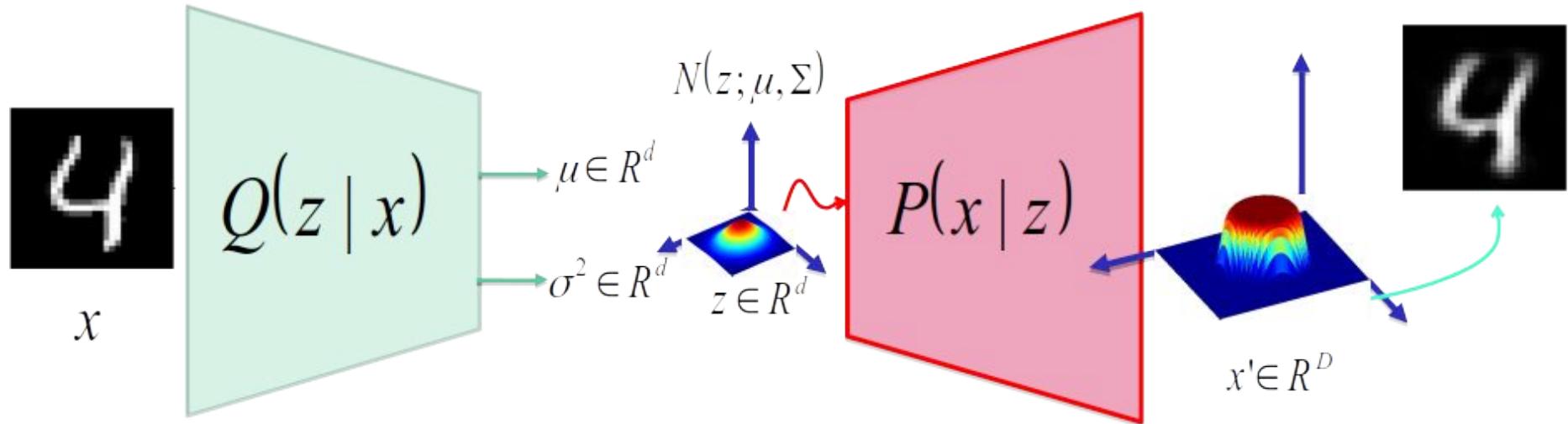


Pour simplifier les calculs, on va supposer que  $\Sigma$  est diagonale

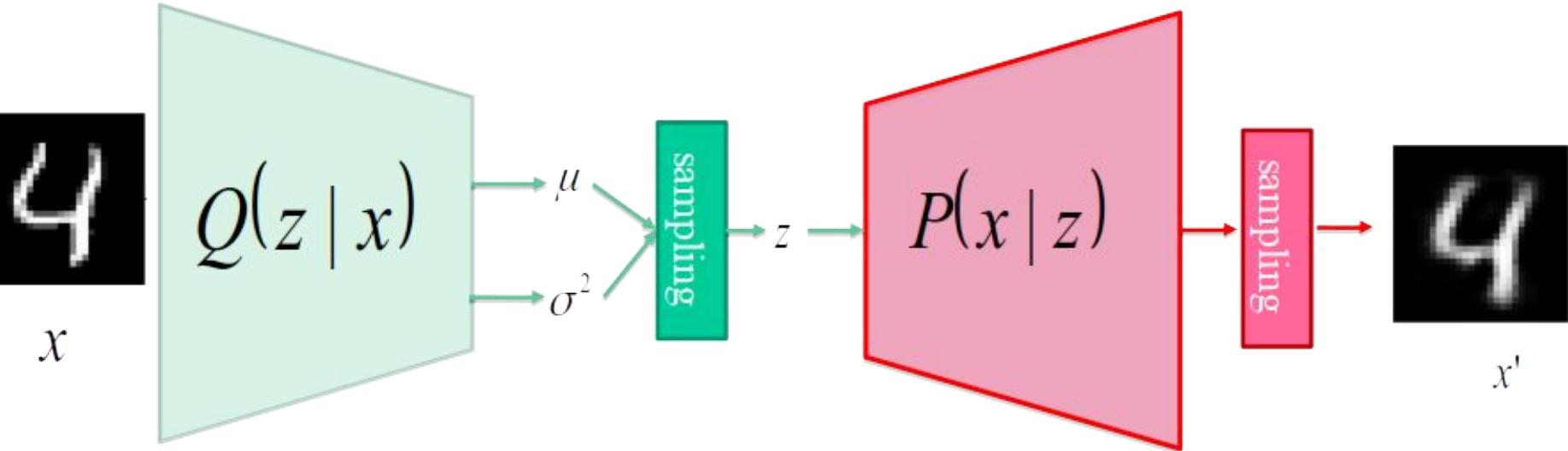
$$\Sigma = \begin{pmatrix} \sigma_1^2 & 0 & 0 & 0 & 0 \\ 0 & \sigma_2^2 & 0 & 0 & 0 \\ 0 & 0 & \sigma_3^2 & 0 & 0 \\ 0 & 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & 0 & \sigma_d^2 \end{pmatrix}$$

On va donc **prédirer un vecteur de variances et non une matrice**  
Présume que les pixels sont dépendants de  $z$ , mais indépendants entre eux

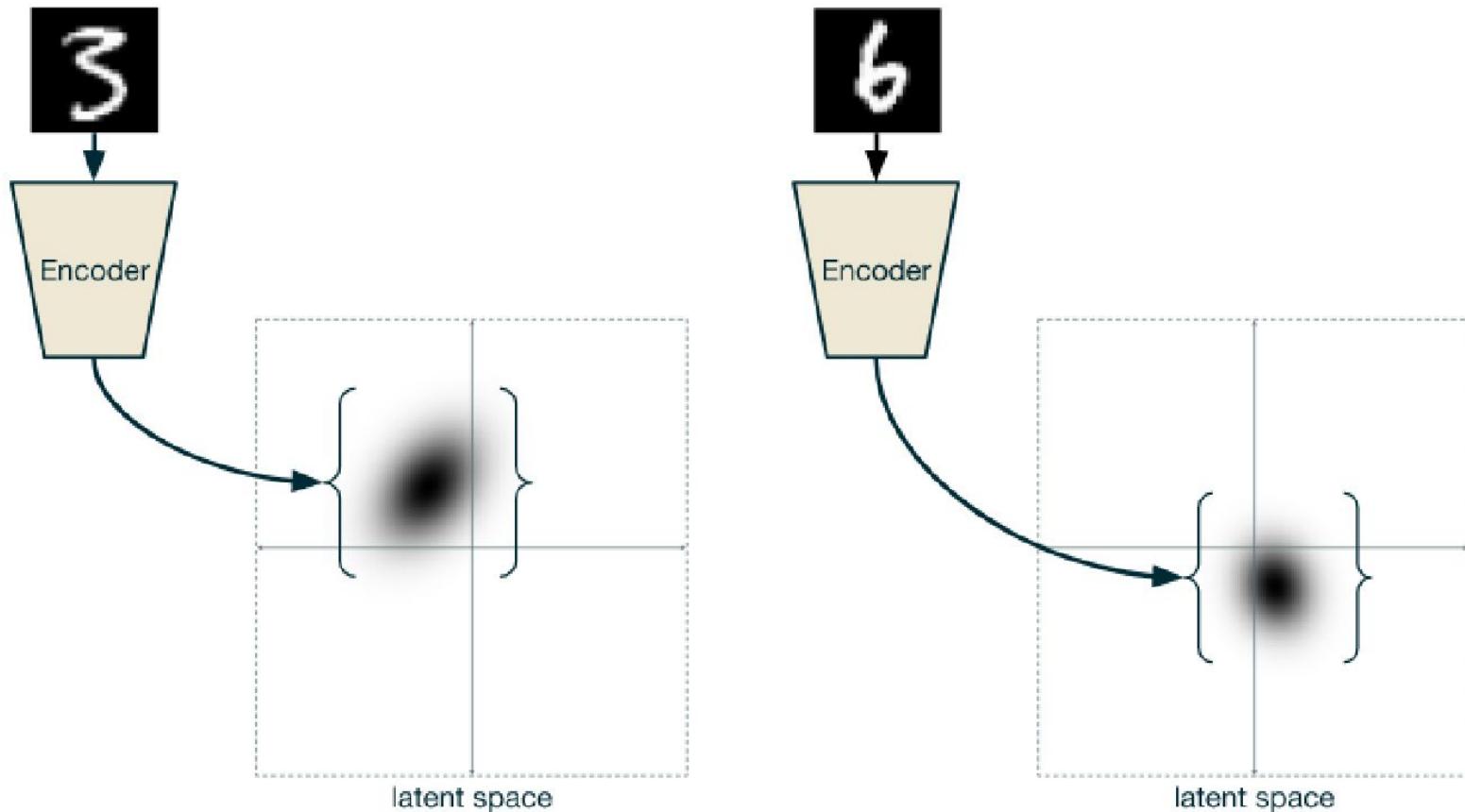
# Autoencodeur variationnel



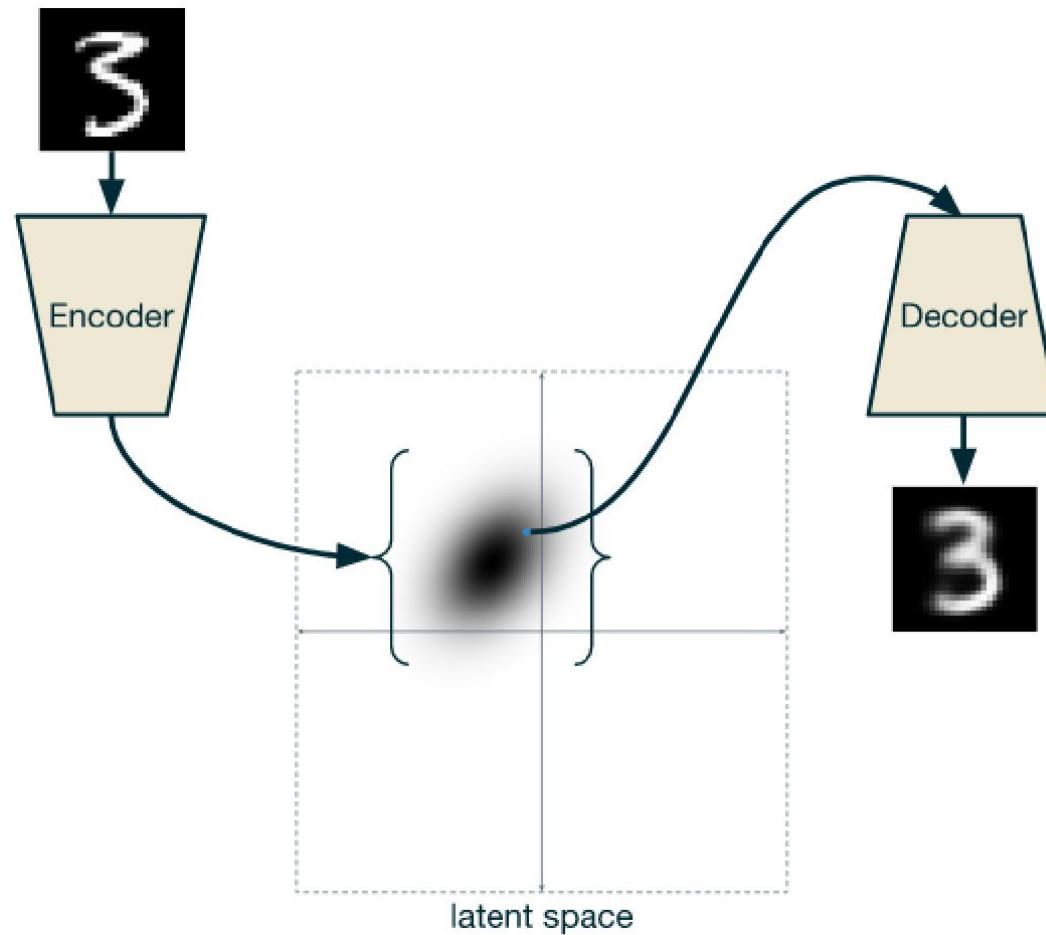
## Autoencodeur variationnel



# Autre façon de voir les choses...

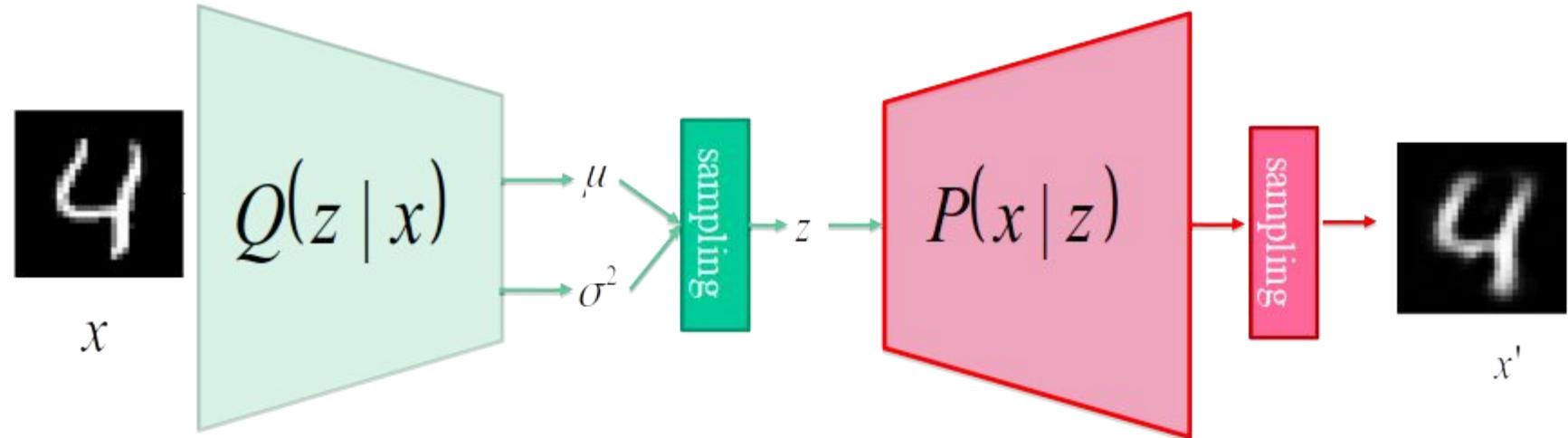


# Autre façon de voir les choses...



## Autoencodeur variationnel

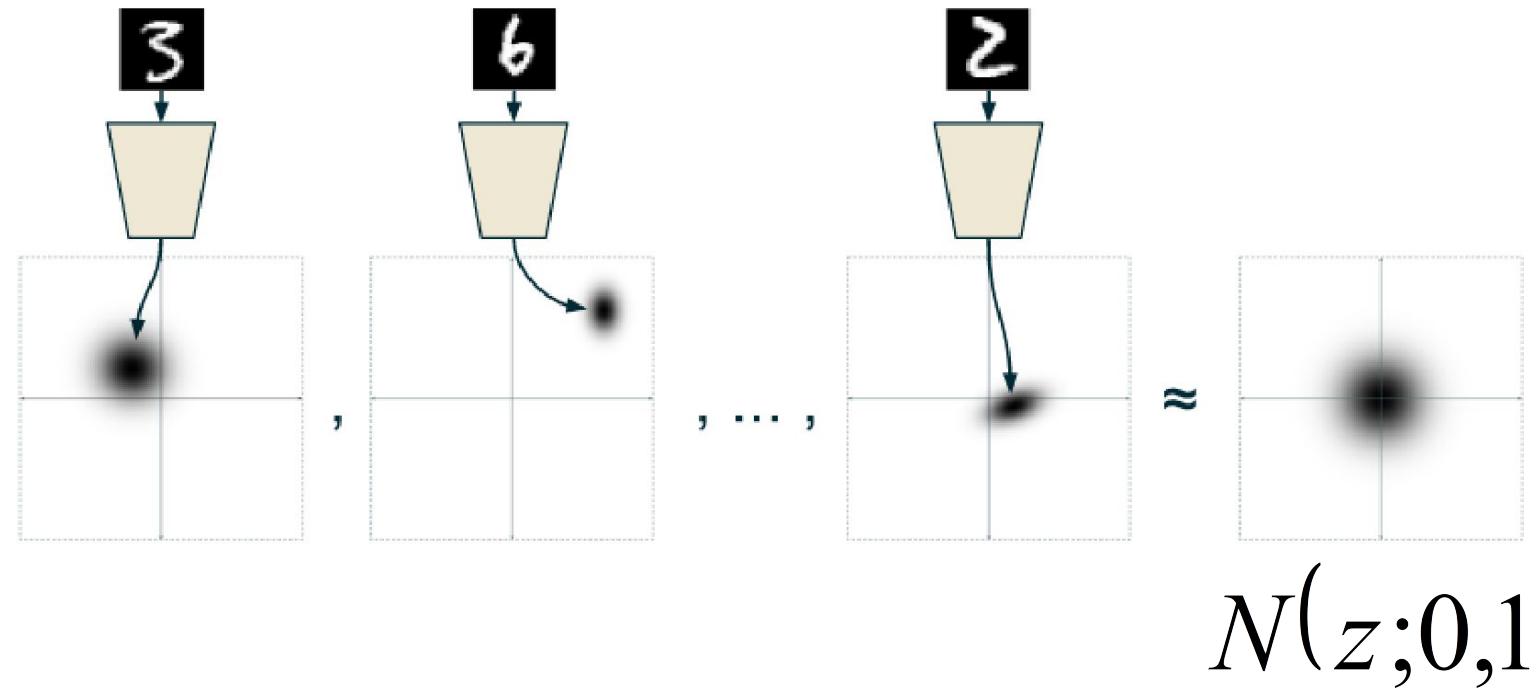
### Remarque 5



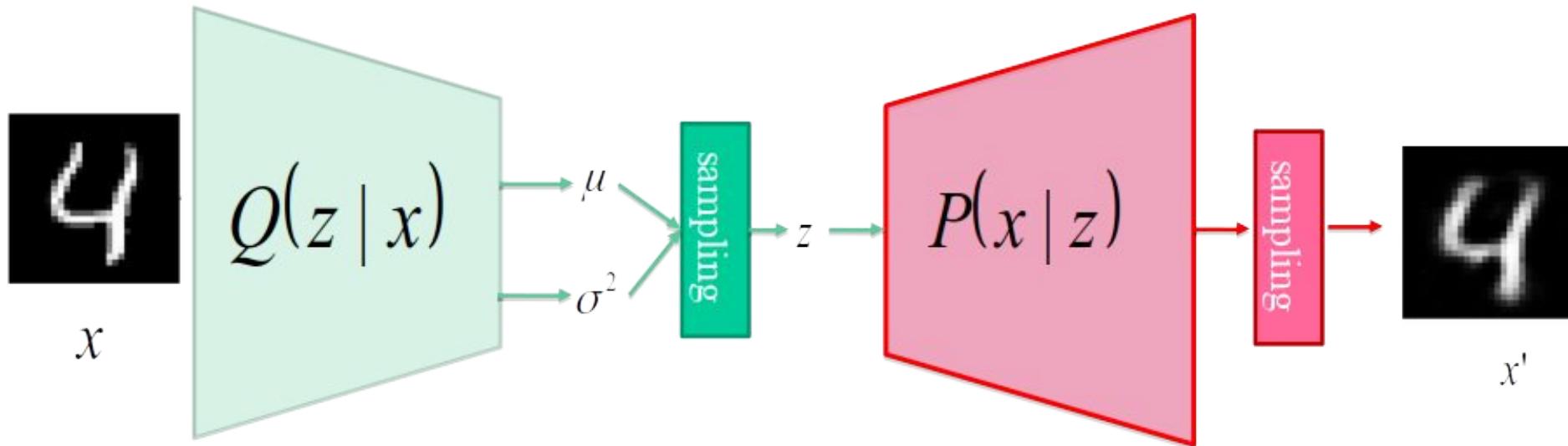
Distribution *a priori* de  $z$  : gaussienne centrée à 0 et de variance 1

$$P(z) = N(z; 0, 1)$$

# Autre façon de voir les choses...



## Autoencodeur variationnel



*ELBO loss : Evidence Lower Bound*

$$Loss = KL(N(z; 0, 1), N(z; \mu, \Sigma)) - \log(P(x | z))$$

Perte encodeur

Perte décodeur

## Autoencodeur variationnel

D.Kingma, M.Welling, Auto-Encoding Variational Bayes, arXiv:1312.6114v10 ([Annexe B](#))

*ELBO loss : Evidence Lower Bound*

$$Loss = \underbrace{KL(N(z;0,1), N(z; \mu, \Sigma))}_{\text{Perte encodeur}} - \underbrace{\log(P(x | z))}_{\text{Perte décodeur}}$$

Si on suppose que  $P(x|z)$  est gaussien

$$Loss = \frac{1}{2} \sum_{i=1}^d (1 + \log(\sigma_i^2) + \mu_i^2 - \sigma_i^2) - \lambda ||\vec{x} - \vec{x}'||^2$$


# Autoencodeur variationnel

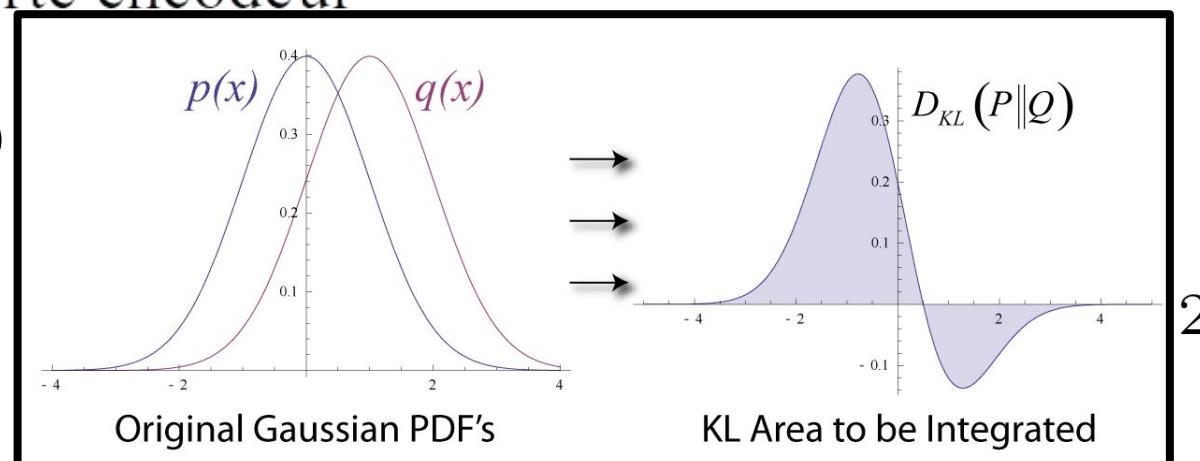
D.Kingma, M.Welling, **Auto-Encoding Variational Bayes**, arXiv:1312.6114v10 ([Annexe B](#))

*ELBO loss : Evidence Lower Bound*

$$Loss = \underbrace{KL(N(z;0,1), N(z; \mu, \Sigma))}_{\text{Perte encodeur}} - \underbrace{\log(P(x | z))}_{\text{Perte décodeur}}$$

Si on suppose que  $P(x|z)$

$$Loss = \frac{1}{2} \sum_{i=1}^d (1 +$$



$\underbrace{\text{Perte encodeur}}$

$\underbrace{\text{Perte décodeur}}$

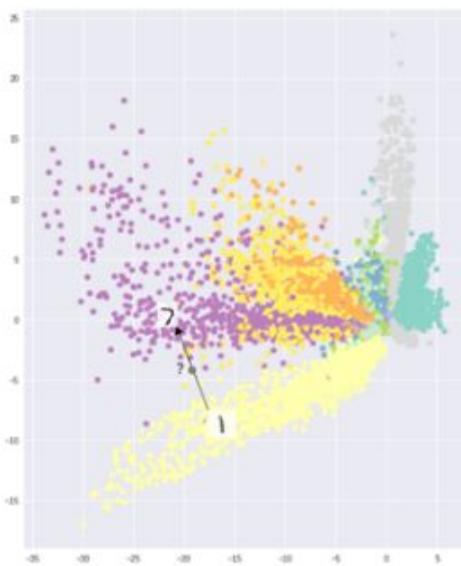
# Autoencodeur variationnel

D.Kingma, M.Welling, **Auto-Encoding Variational Bayes**, arXiv:1312.6114v10 (*Annexe B*)

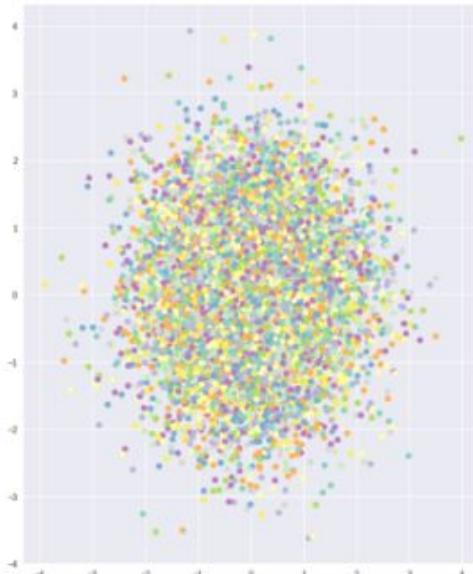
*Evidence Lower Bound*

$$\text{Loss} = \underbrace{\text{KL}(N(z;0,1), N(z;\mu,\Sigma))}_{\text{Perte encodeur}} - \underbrace{\log(P(x|z))}_{\text{Perte décodeur}}$$

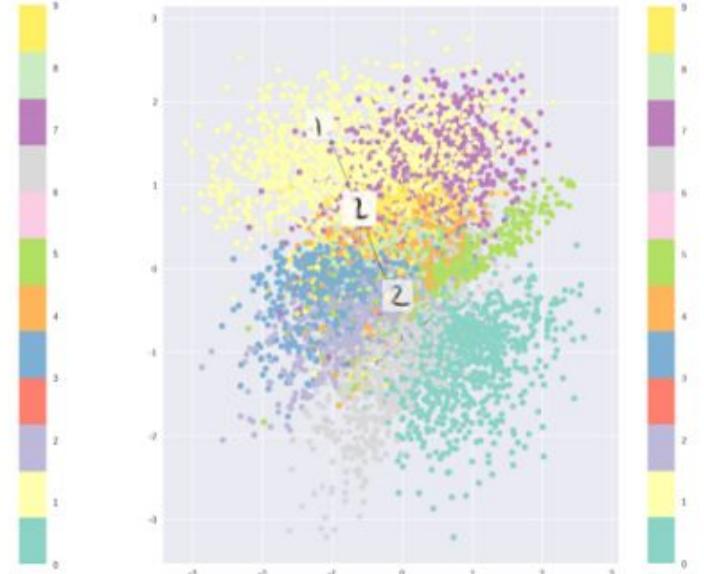
Only reconstruction loss



Only KL divergence



Combination

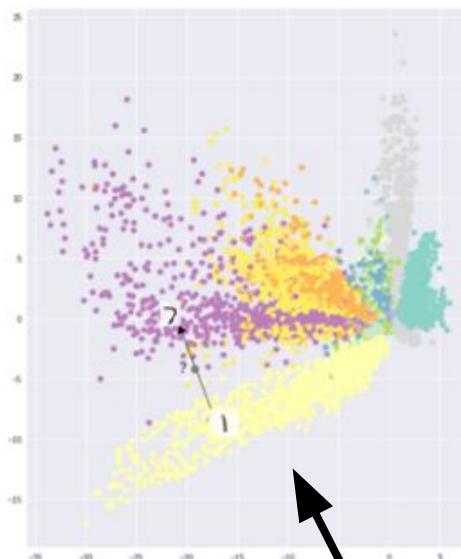


# Autoencodeur variationnel

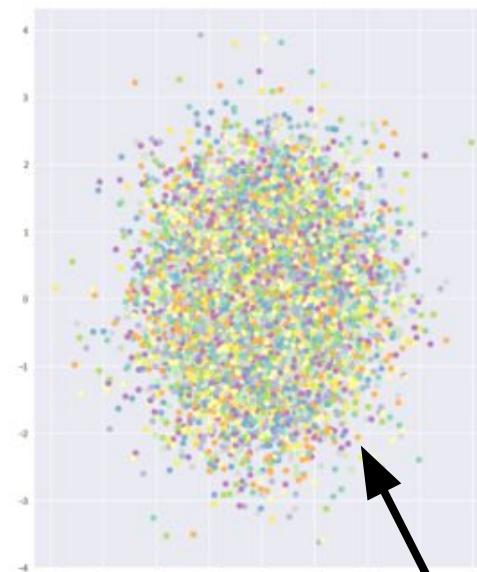
D.Kingma, M.Welling, **Auto-Encoding Variational Bayes**, arXiv:1312.6114v10 (*Annexe B*)

*ELBO loss : Evidence Lower Bound*

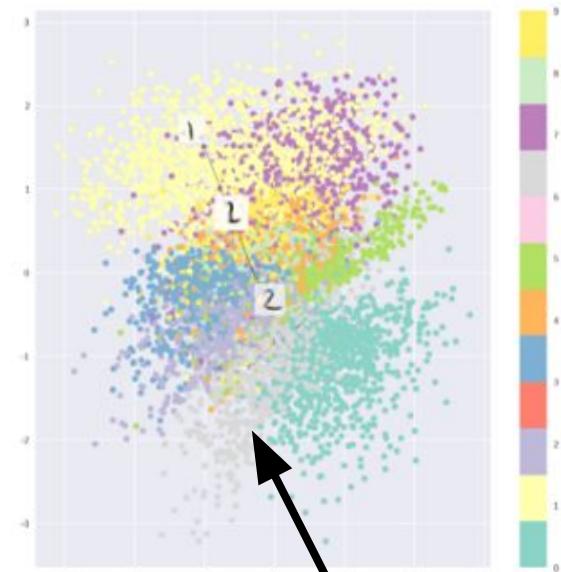
Only reconstruction loss



Only KL divergence



Combination



<https://www.jeremyjordan.me/variational-autoencoders/>

Pas gaussien, difficile à échantillonner

Trop gaussien, tout est indistinguorable

Ensemble de distributions séparées qui forment une gaussienne

## Autoencodeur variational

D.Kingma, M.Welling, Auto-Encoding Variational Bayes, and

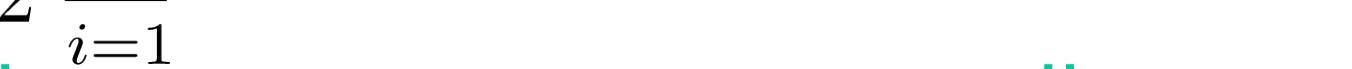
ELBO loss : Evidence Lower Bound

## À voir plus tard

$$Loss = KL(N(z; 0, 1), N(z; \mu, \Sigma))$$

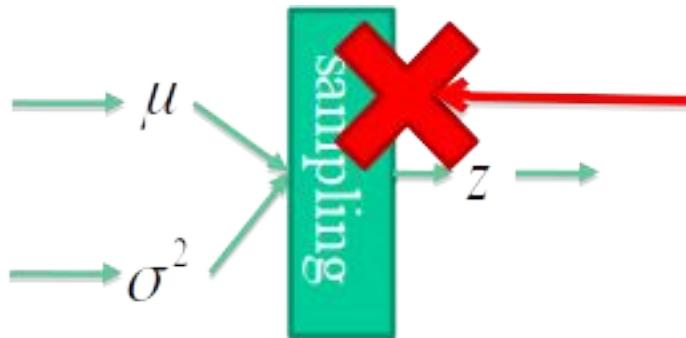
The diagram shows the total loss function as the sum of two components. A green bracket labeled "Perte encodeur" spans from the start of the equation to the term  $N(z; \mu, \Sigma)$ . A yellow bracket labeled "Perte décodeur" spans from  $\mu$  to the end of the equation  $N(\lambda|z)$ .

Si on suppose que  $P(x|z)$  est gaussien

$$Loss = \frac{1}{2} \sum_{i=1}^d (1 + \log(\sigma_i^2) + \mu_i^2 - \sigma_i^2) - \lambda ||\vec{x} - \vec{x}'||^2$$


# Autoencodeur variationnel

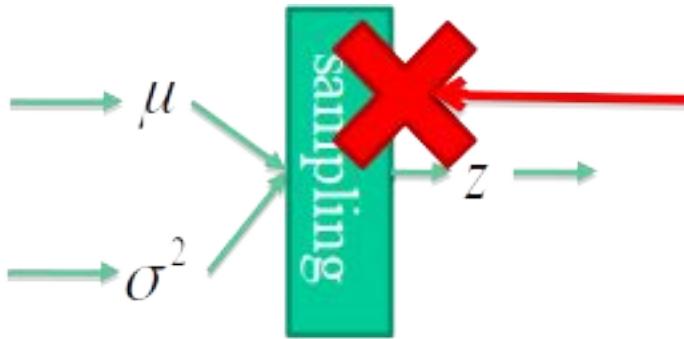
## Remarque 6



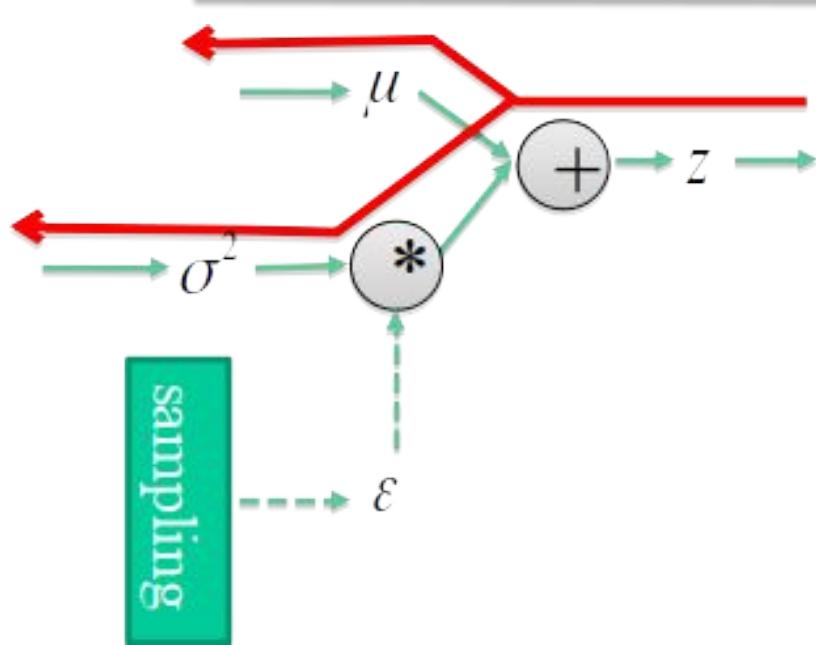
Pas de rétro-propagation à travers  
un processus d'échantillonage  
 $z \sim N(z; \mu, \Sigma)$

# Autoencodeur variationnel

## Remarque 6



Pas de rétro-propagation à travers  
un processus d'échantionnage  
 $z \sim N(z; \mu, \Sigma)$



*Reparameterization trick*

$$z = \mu + \varepsilon \sigma^2$$

# Autoencodeur variationnel jouet MNIST : d=32 dim

```
class VAE(nn.Module):
    def __init__(self):
        super(VAE, self).__init__()

        self.encoder = nn.Sequential(
            nn.Linear(28 * 28, 128), nn.ReLU(True),
            nn.Linear(128, 64), nn.ReLU(True),
            nn.Linear(64, 32*2)
        self.decoder = nn.Sequential(
            nn.Linear(32, 64), nn.ReLU(True),
            nn.Linear(64, 128), nn.ReLU(True),
            nn.Linear(128, 28 * 28))

    def reparameterize(self, mu, logvar):
        std = torch.exp(0.5*logvar)
        eps = torch.randn_like(std)
        return mu + eps*std

    def forward(self, x):
        enc_x = self.encoder(x)
        mu = enc_x[:, :32]
        logvar = enc_x[:, 32:]
        z = self.reparameterize(mu, logvar)
        return self.decoder(z), mu, logvar
```

} Reparameterization  
trick

# Autoencodeur variationnel jouet MNIST : d=32 dim

```
class VAE(nn.Module):
    def __init__(self):
        super(VAE, self).__init__()

        self.encoder = nn.Sequential(
            nn.Linear(28 * 28, 128), nn.ReLU(True),
            nn.Linear(128, 64), nn.ReLU(True),
            nn.Linear(64, 32*2) # Espace latent de 32 dimensions
        self.decoder = nn.Sequential(
            nn.Linear(32, 64), nn.ReLU(True),
            nn.Linear(64, 128), nn.ReLU(True),
            nn.Linear(128, 28 * 28))

    def reparameterize(self, mu, logvar):
        std = torch.exp(0.5*logvar)
        eps = torch.randn_like(std)
        return mu + eps*std

    def forward(self, x):
        enc_x = self.encoder(x)
        mu = enc_x[:, :32]
        logvar = enc_x[:, 32:]
        z = self.reparameterize(mu, logvar)
        return self.decoder(z), mu, logvar
```

} Reparameterization  
trick

# Autoencodeur variationnel jouet MNIST : d=32 dim

```
def loss_function(recon_x, x, mu, logvar):
    L2 = torch.sum((recon_x-x).pow(2))

    KLD = -0.5 * torch.sum(1 + logvar - mu.pow(2) - logvar.exp())

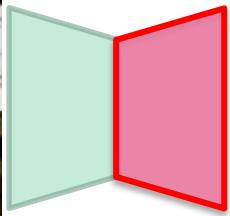
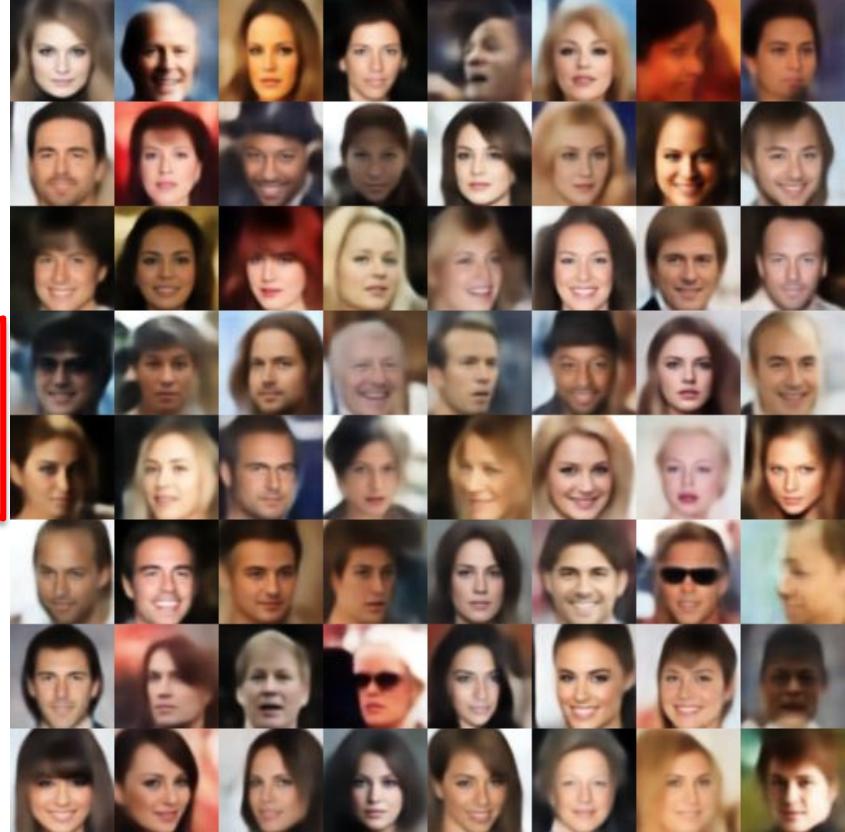
    return L2 + KLD
```

# Ex.: base de données *CelebA*

$x$

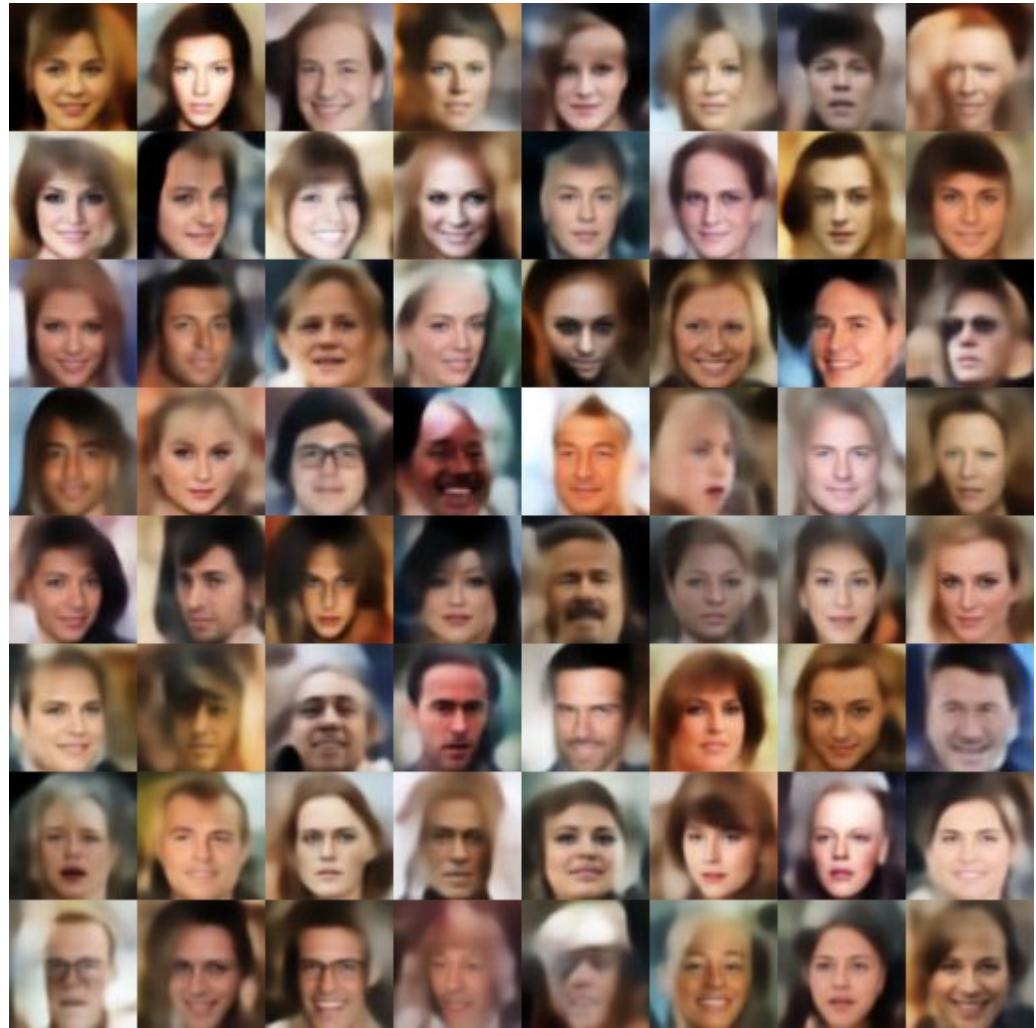
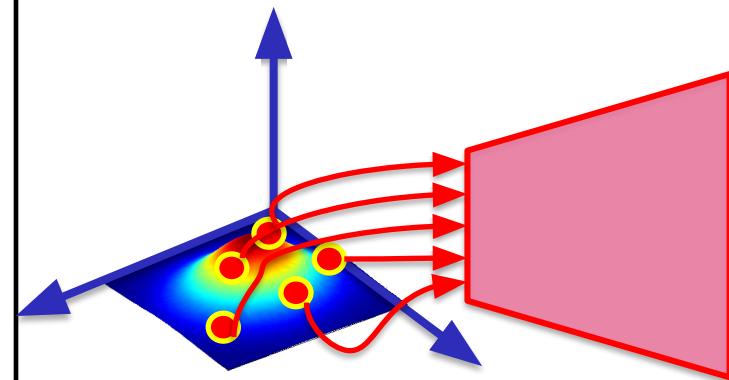


$x'$



# Ex.: base de données *CelebA*

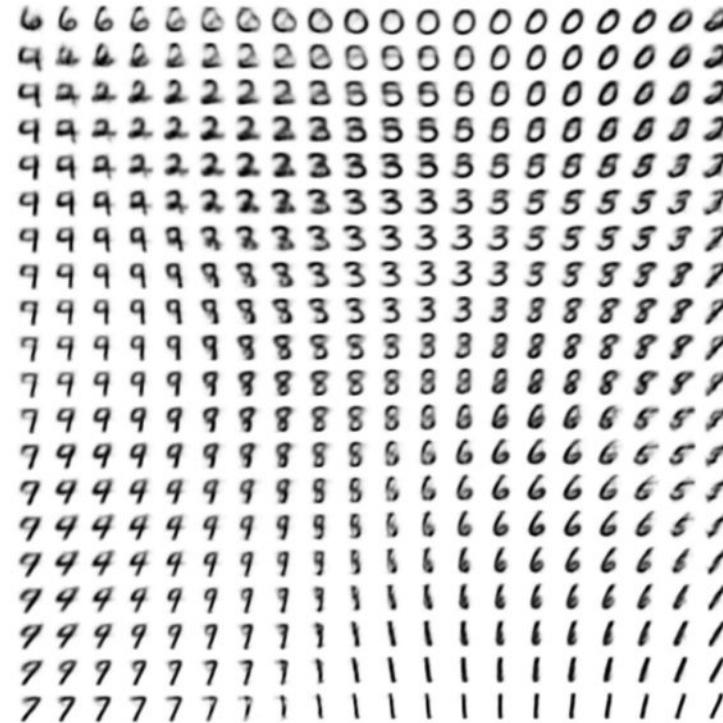
Décodage d'échantillons aléatoires z



# Ex.: VAE



(a) Learned Frey Face manifold

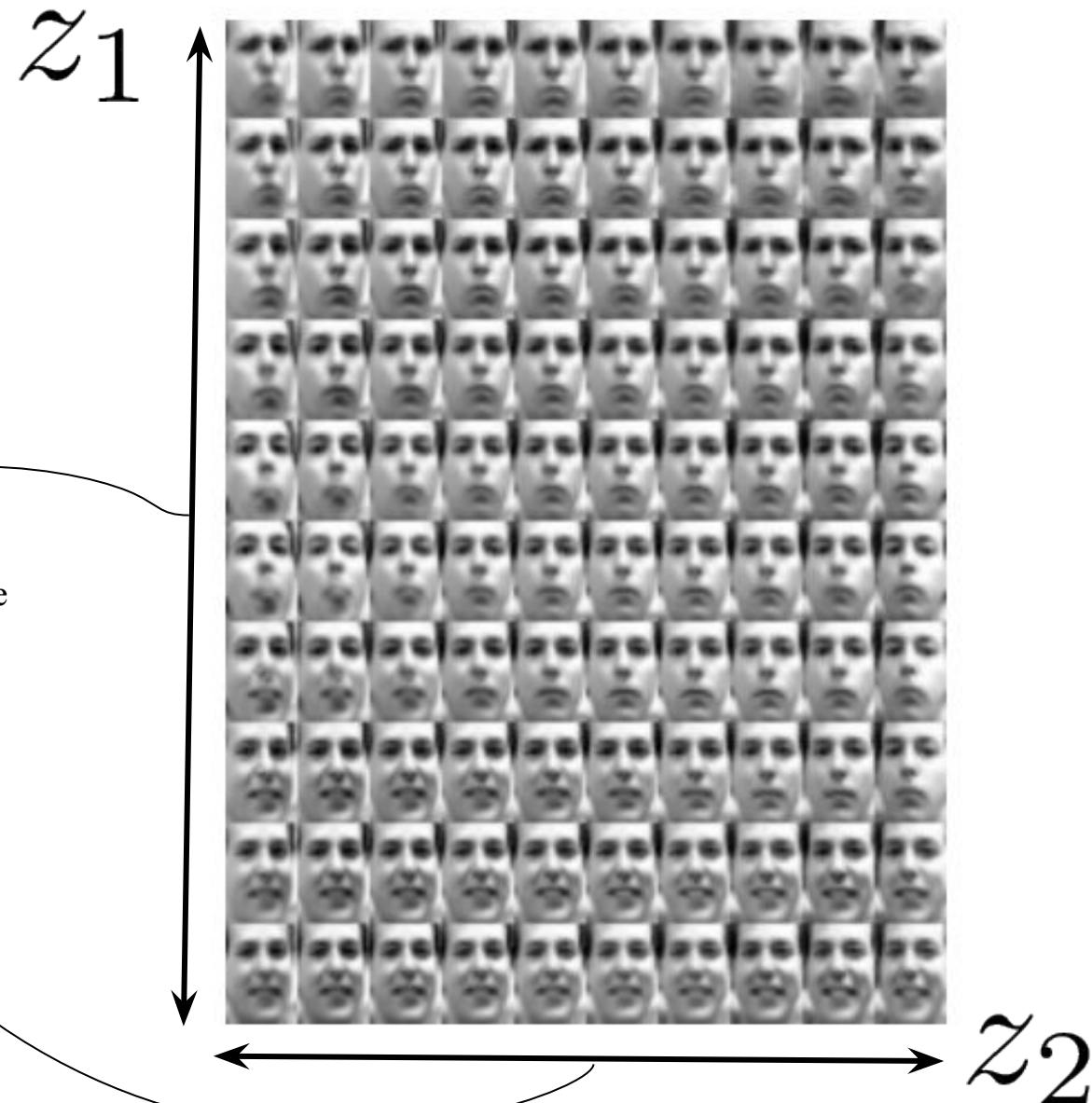


(b) Learned MNIST manifold

Figure 4: Visualisations of learned data manifold for generative models with two-dimensional latent space, learned with AEVB. Since the prior of the latent space is Gaussian, linearly spaced coordinates on the unit square were transformed through the inverse CDF of the Gaussian to produce values of the latent variables  $\mathbf{z}$ . For each of these values  $\mathbf{z}$ , we plotted the corresponding generative  $p_{\theta}(\mathbf{x}|\mathbf{z})$  with the learned parameters  $\theta$ .

# Ex.: VAE

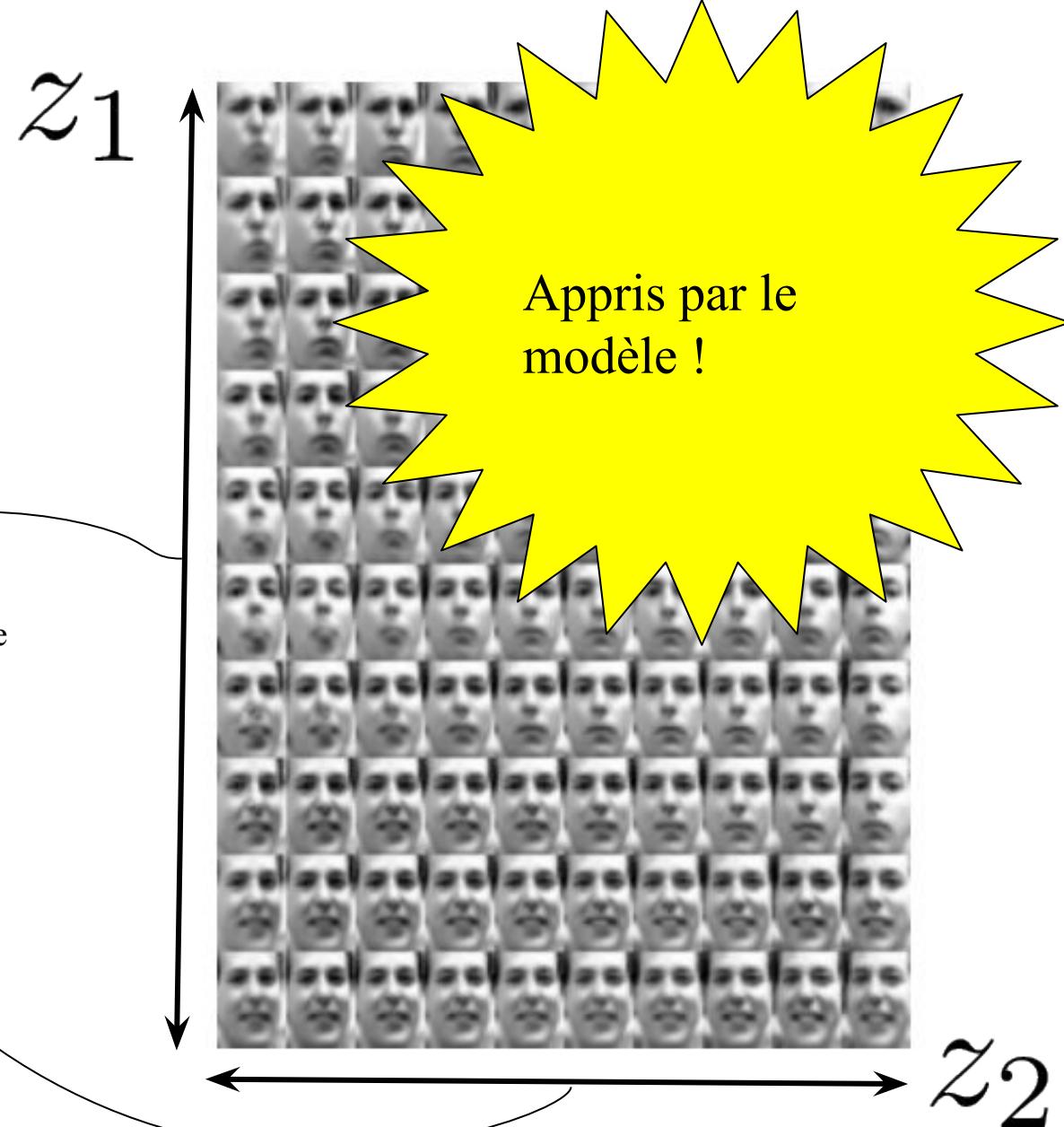
Dimension du “sourire”  
Dimension de l’orientation de la tête



Kingma, D. P., & Welling, M. (2013). Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.

# Ex.: VAE

Dimension du “sourire”  
Dimension de l’orientation de la tête

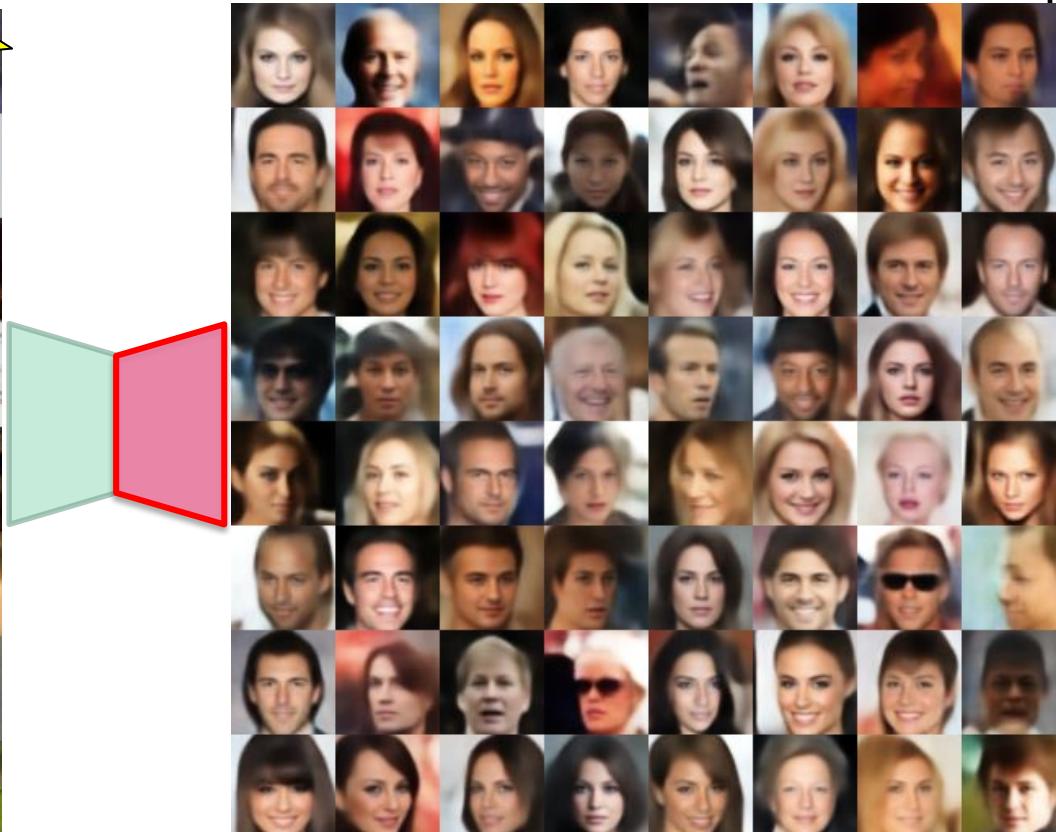


Images floues.  
Pourquoi ?



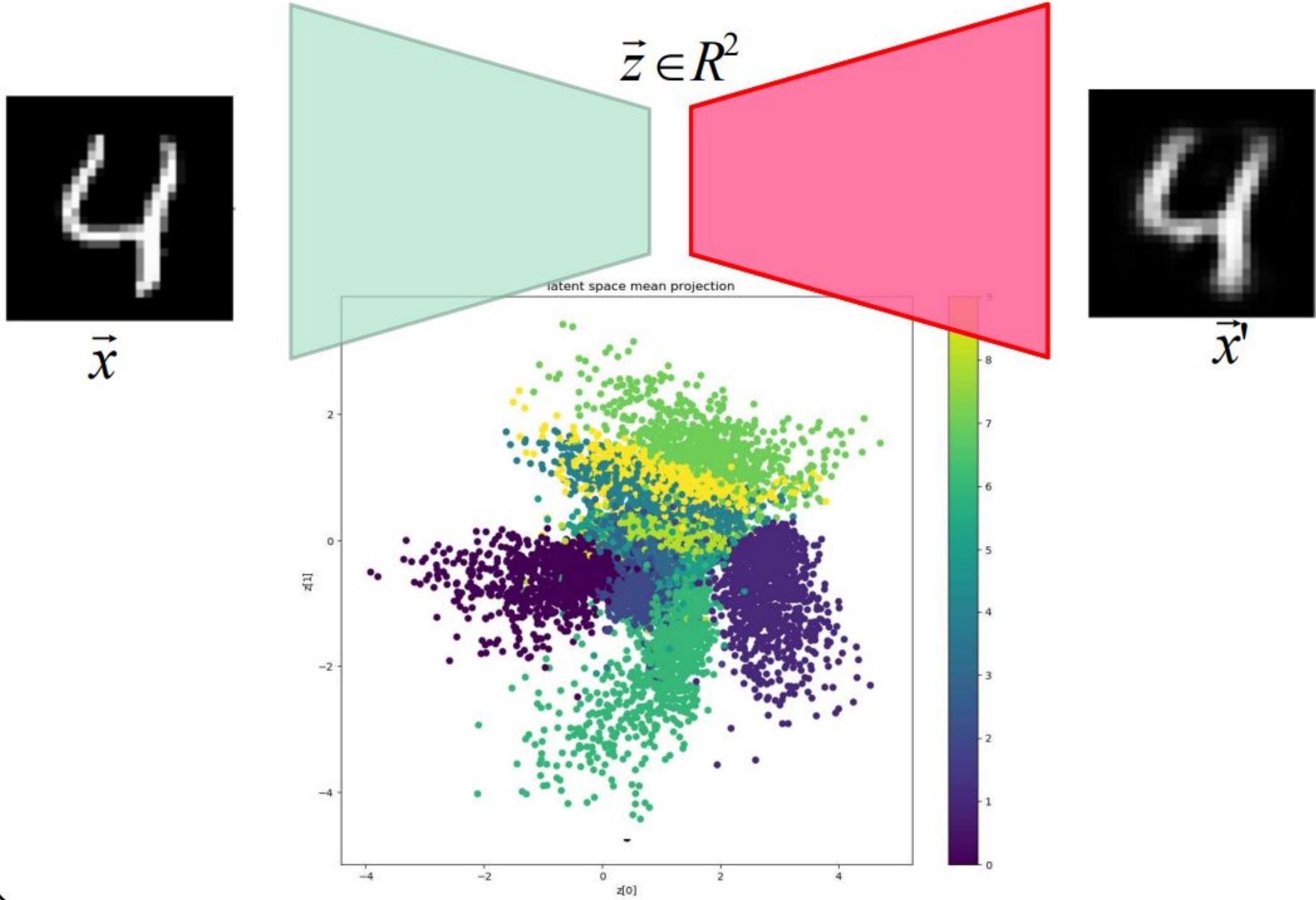
# Les données *CelebA*

$x^!$

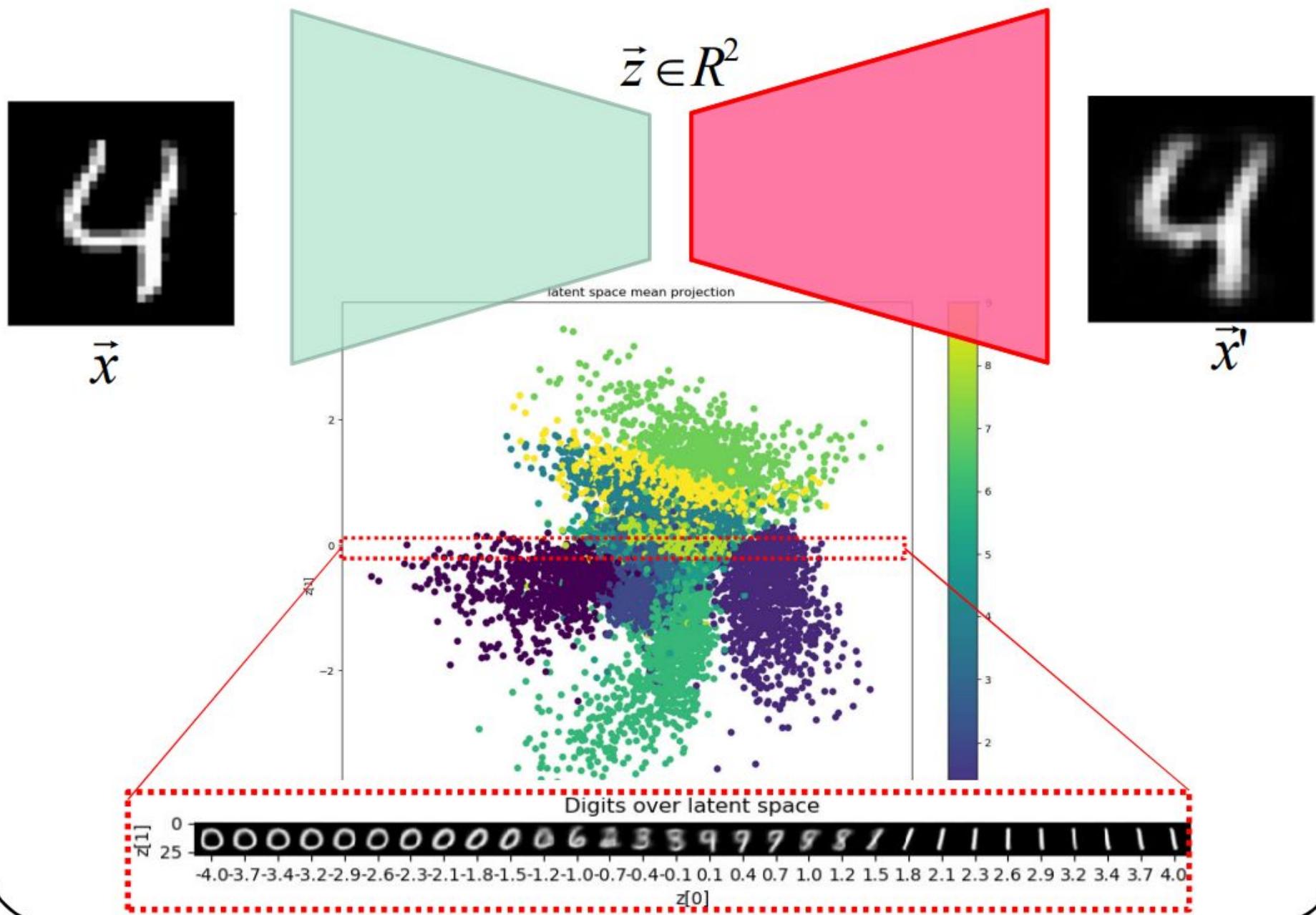


NOTE: on peut également conditionner un espace latent en y  
« **injectant les gradients d'un 2<sup>e</sup> réseau de neurones** »

# Autoencodeur de base



## Autoencodeur de base

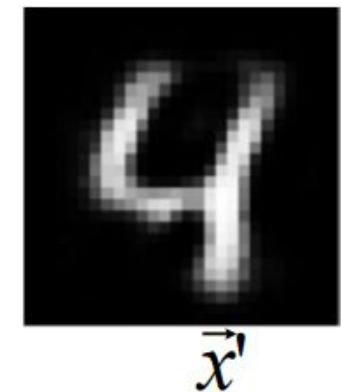
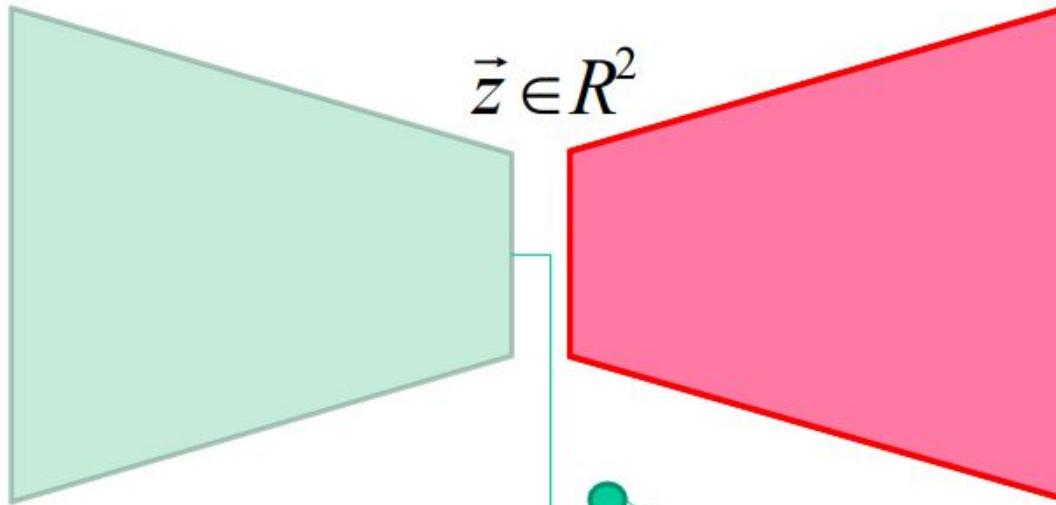


## Autoencodeur contraint



$\vec{x}$

$t=4$



$\vec{x}'$

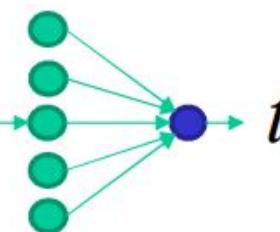
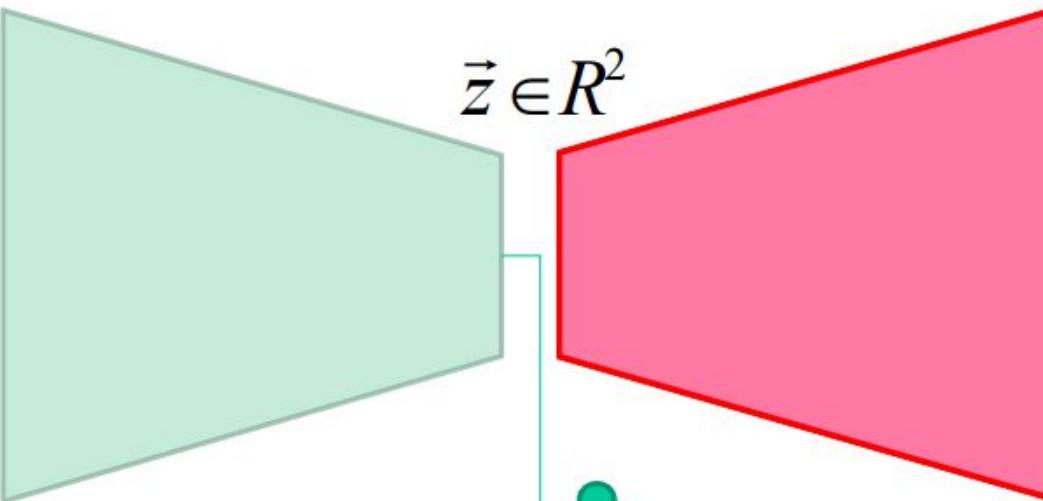
Régression linéaire  
(perceptron sans couche cachée)

## Autoencodeur contraint



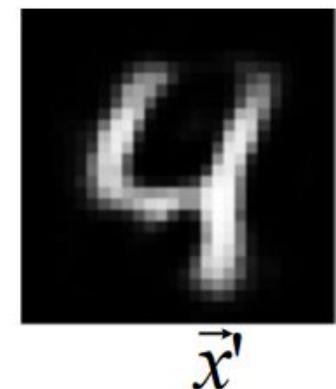
$\vec{x}$

$t=4$



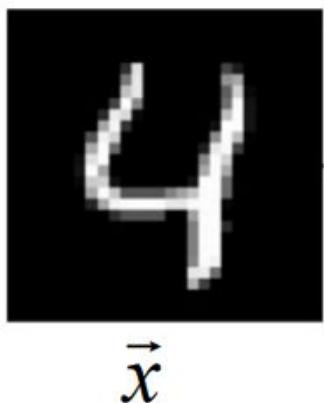
$$L_1 = \|\vec{x} - \vec{x}'\|^2$$

$$L_2 = \|\hat{t} - t\|^2$$

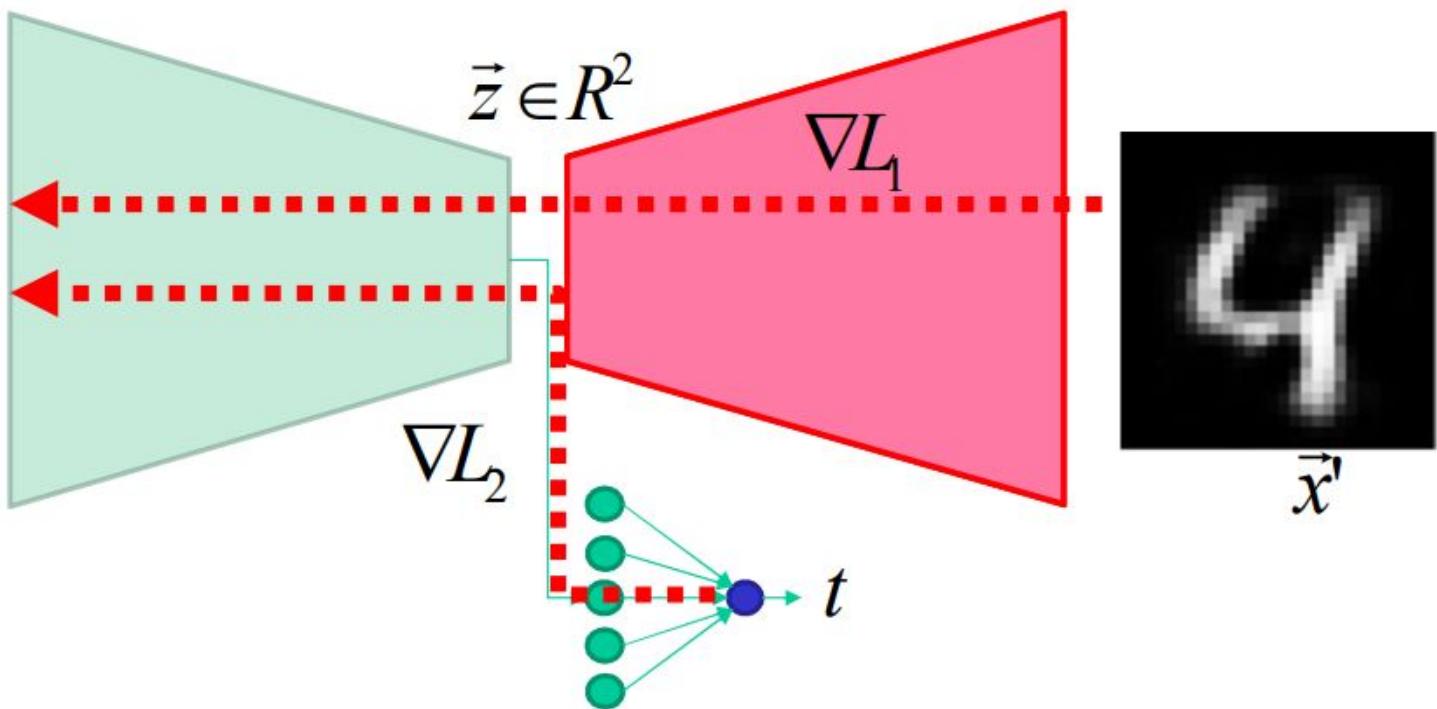


$\vec{x}'$

## Autoencodeur contraint



$\vec{x}$   
 $t=4$

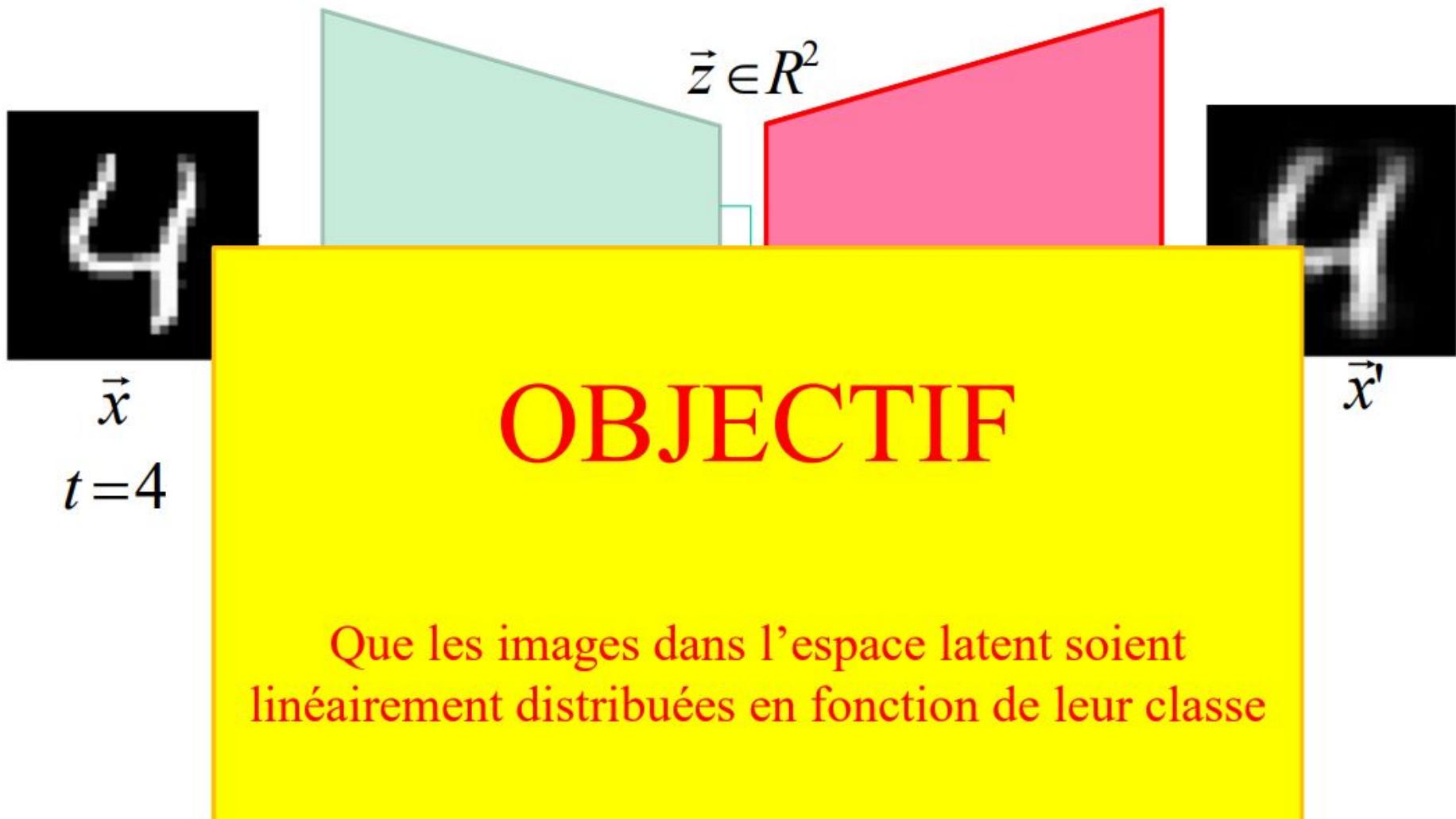


$$L_{\text{décodeur}} = L_1$$

$$L_{\text{reg}} = L_2$$

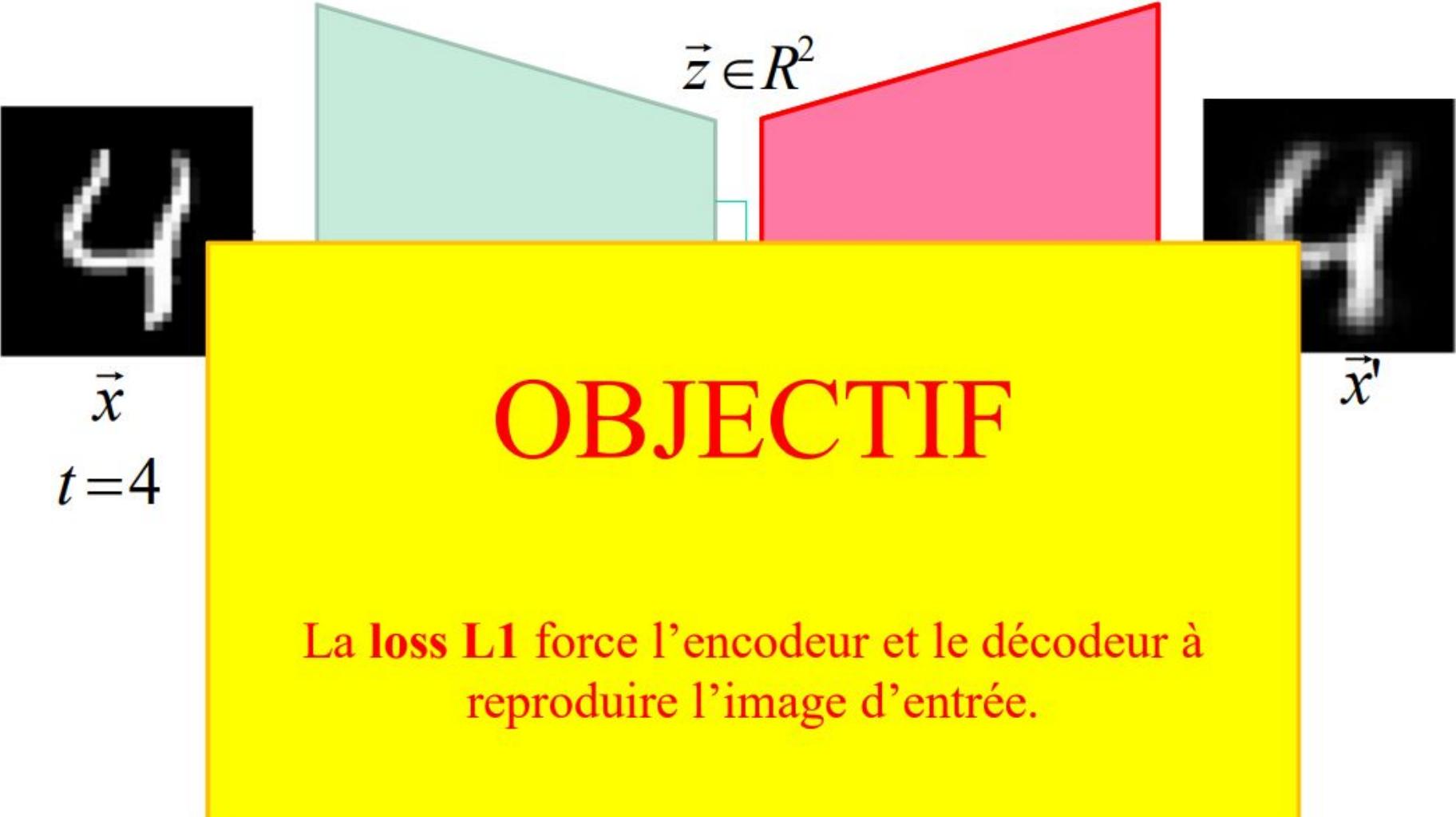
$$L_{\text{encodeur}} = L_1 + \lambda L_2$$

## Autoencodeur contraint



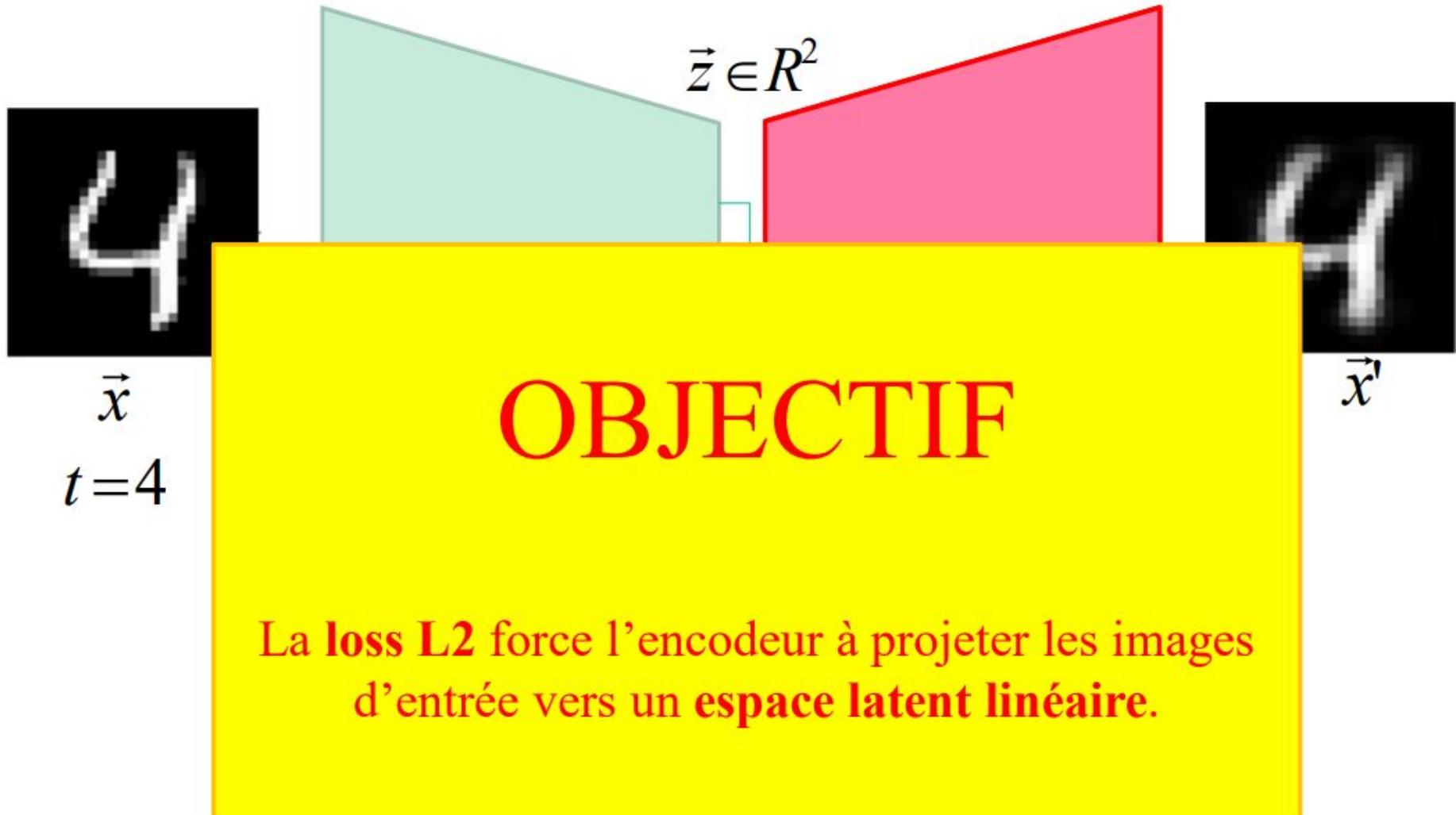
$$L_{encodeur} = L_1 + \lambda L_2$$

## Autoencodeur contraint



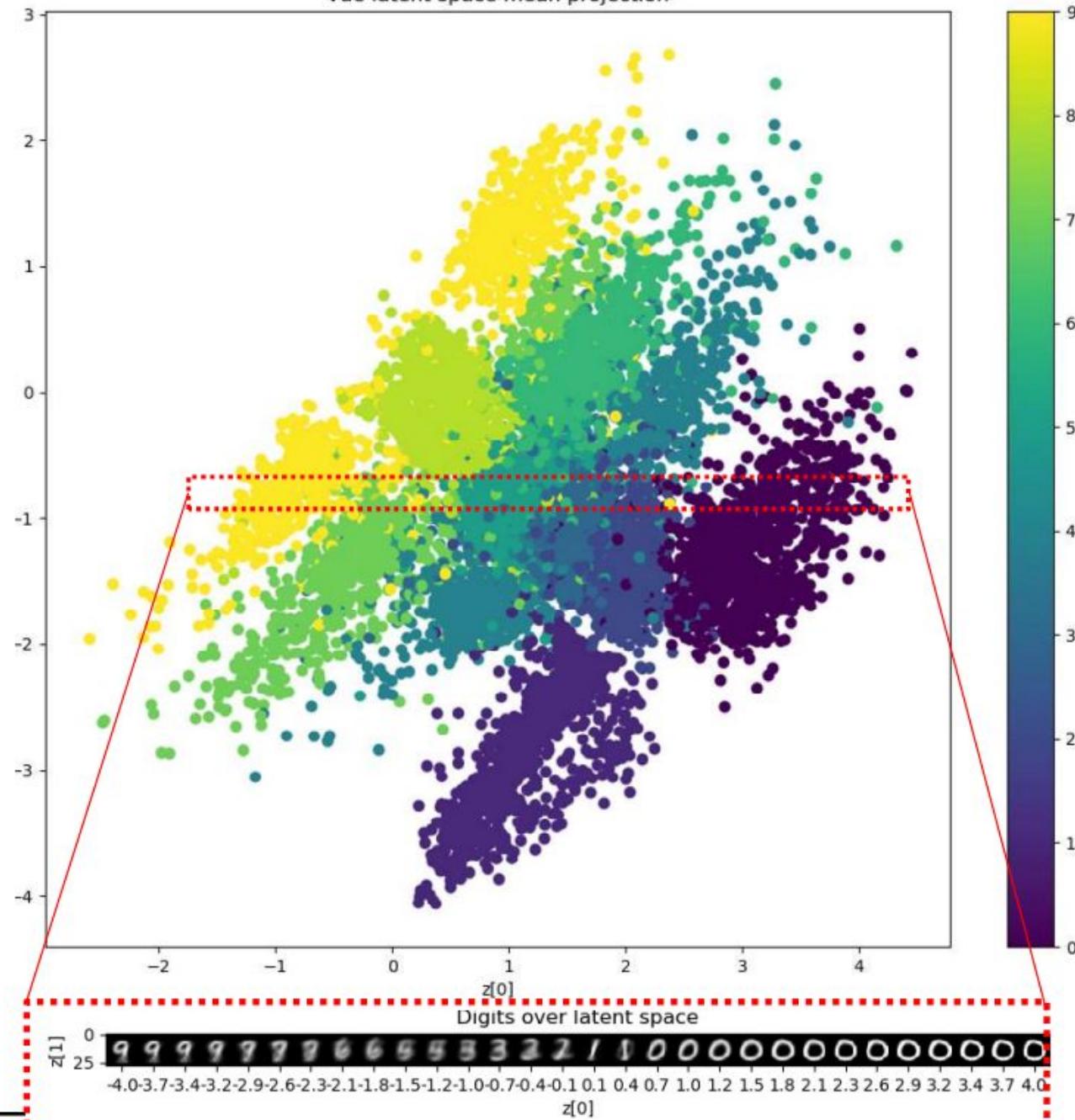
$$L_{encodeur} = L_1 + \lambda L_2$$

## Autoencodeur contraint



$$L_{encodeur} = L_1 + \lambda L_2$$

### Vae latent space mean projection

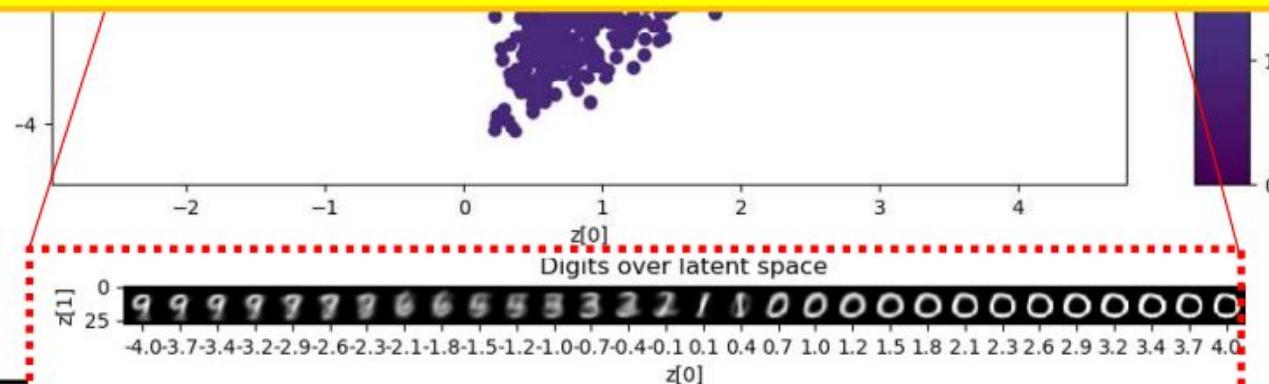


Vae latent space mean projection



## NOTE:

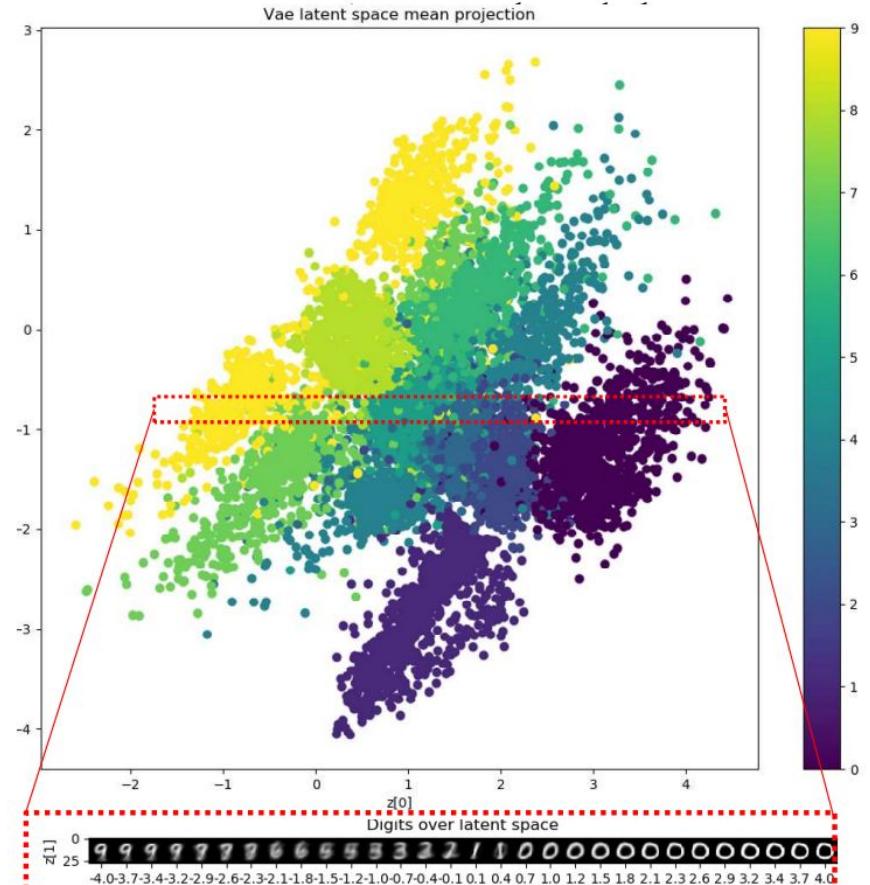
D'autres réseaux de neurones et d'autres loss peuvent être utilisés! Tout dépend de l'application ... et **de votre créativité.**



# VAE conditionnels

## Le problème :

Un autoencodeur (variationnel, ou pas) ne permet pas de générer une image d'une classe en particulier



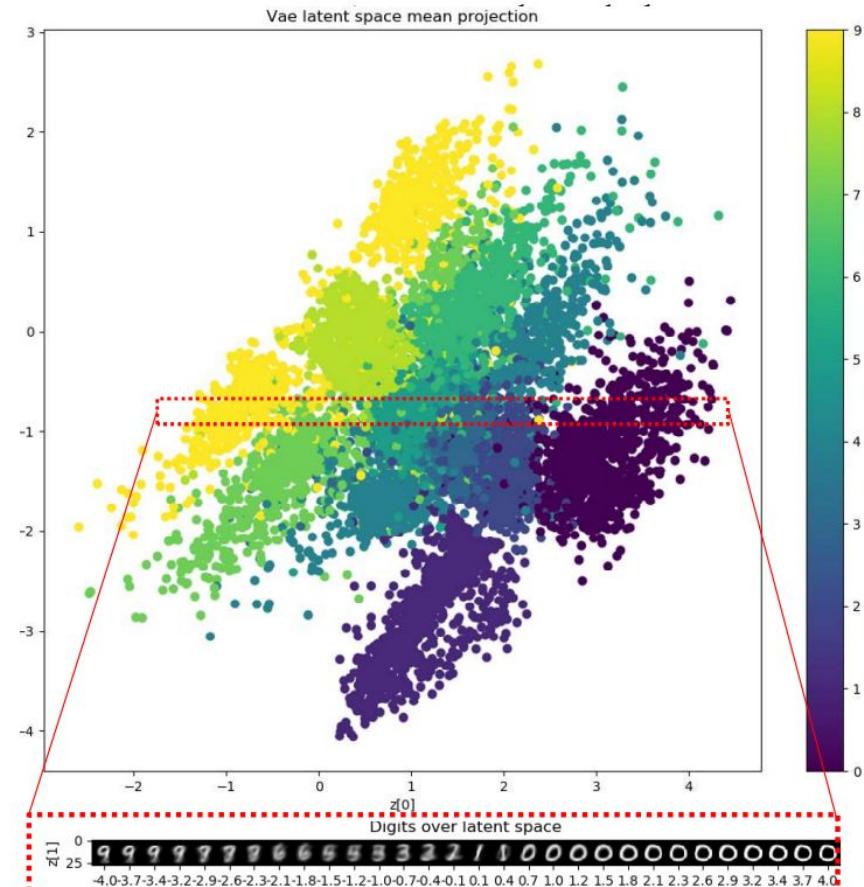
# VAE conditionnels

Le problème :

Un autoencodeur (variationnel, ou pas) ne permet pas de générer une image d'une classe en particulier

La solution:

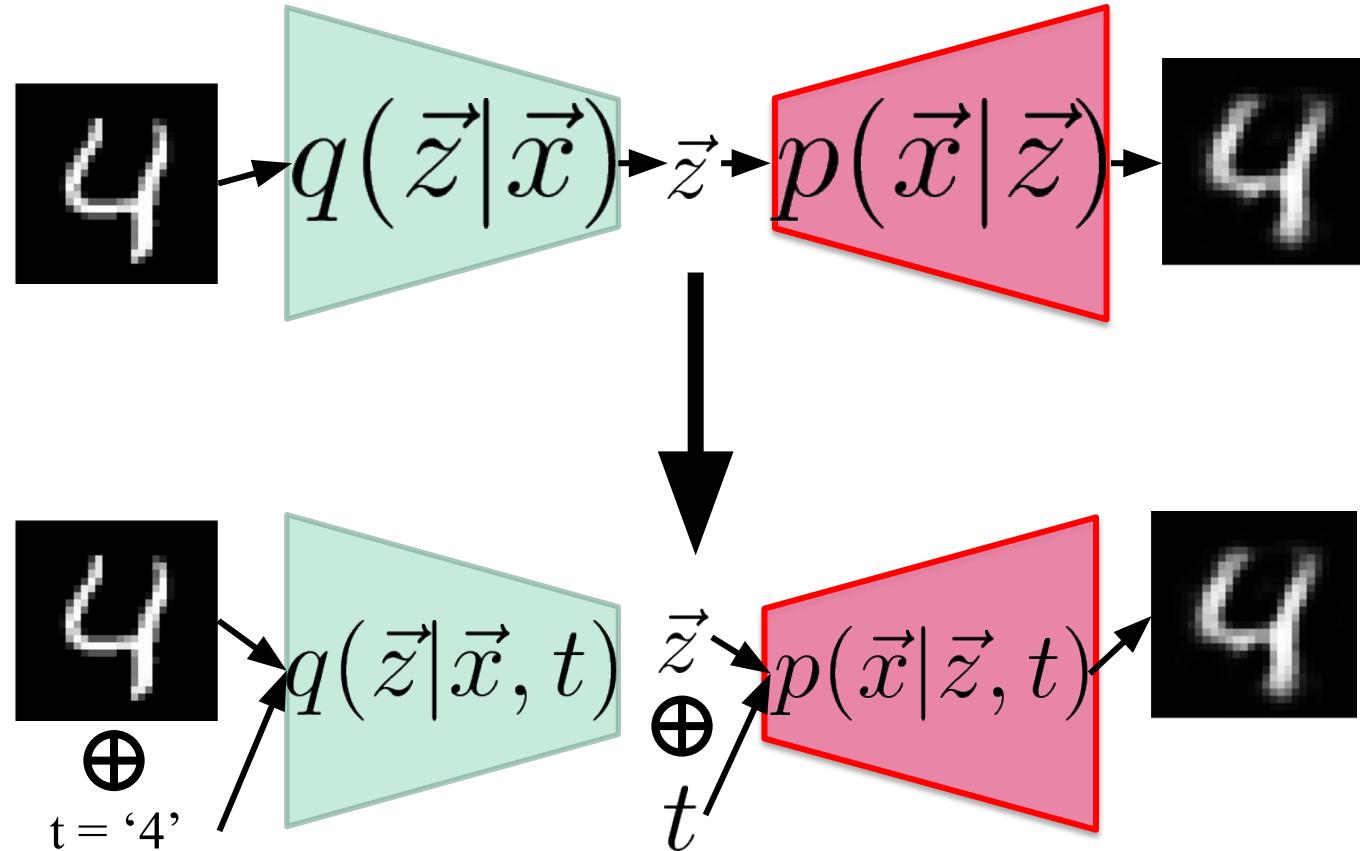
Injecter l'information de la classe dans notre encodeur et décodeur !



# VAE conditionnels

La solution:

Injecter l'information de la classe dans notre encodeur et décodeur !

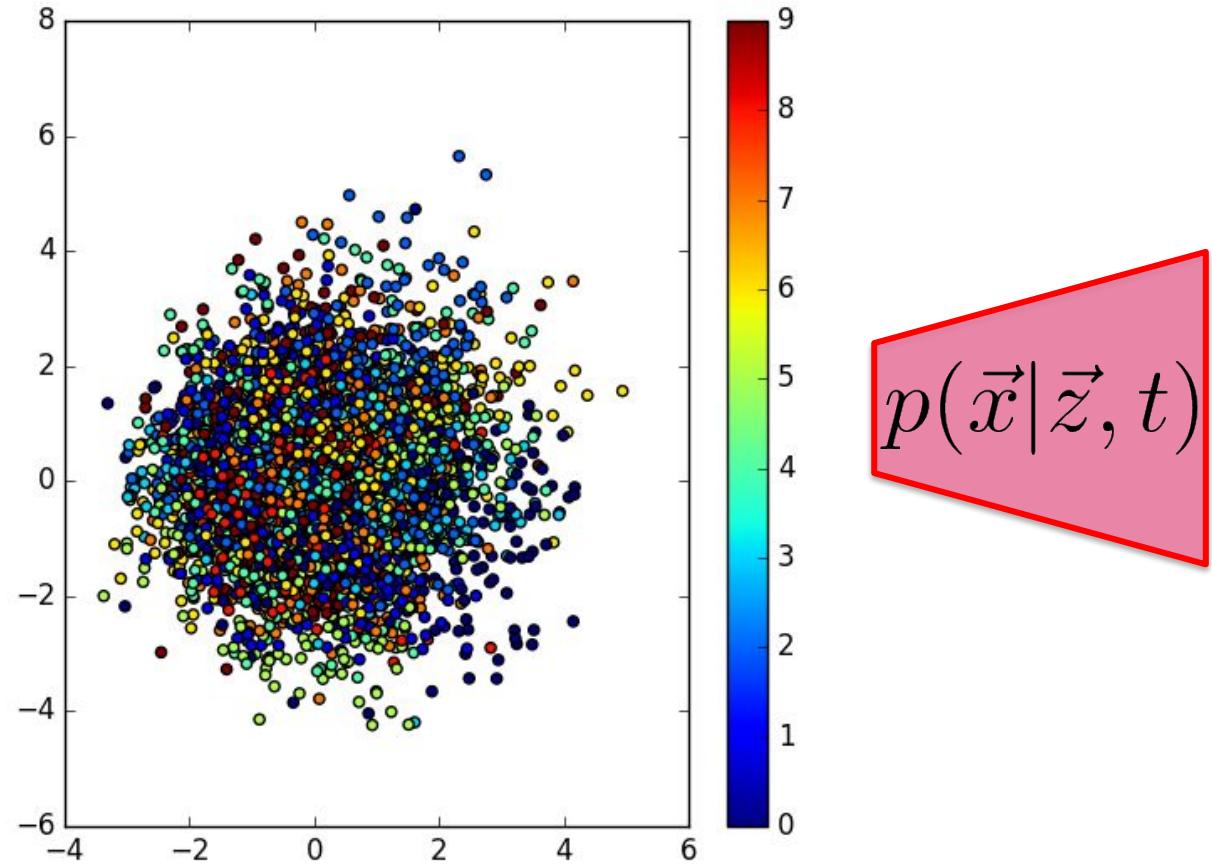


# VAE conditionnels

La solution:

Injecter l'information de la classe dans notre encodeur et décodeur !

$$q(\vec{z}|\vec{x}, t)$$

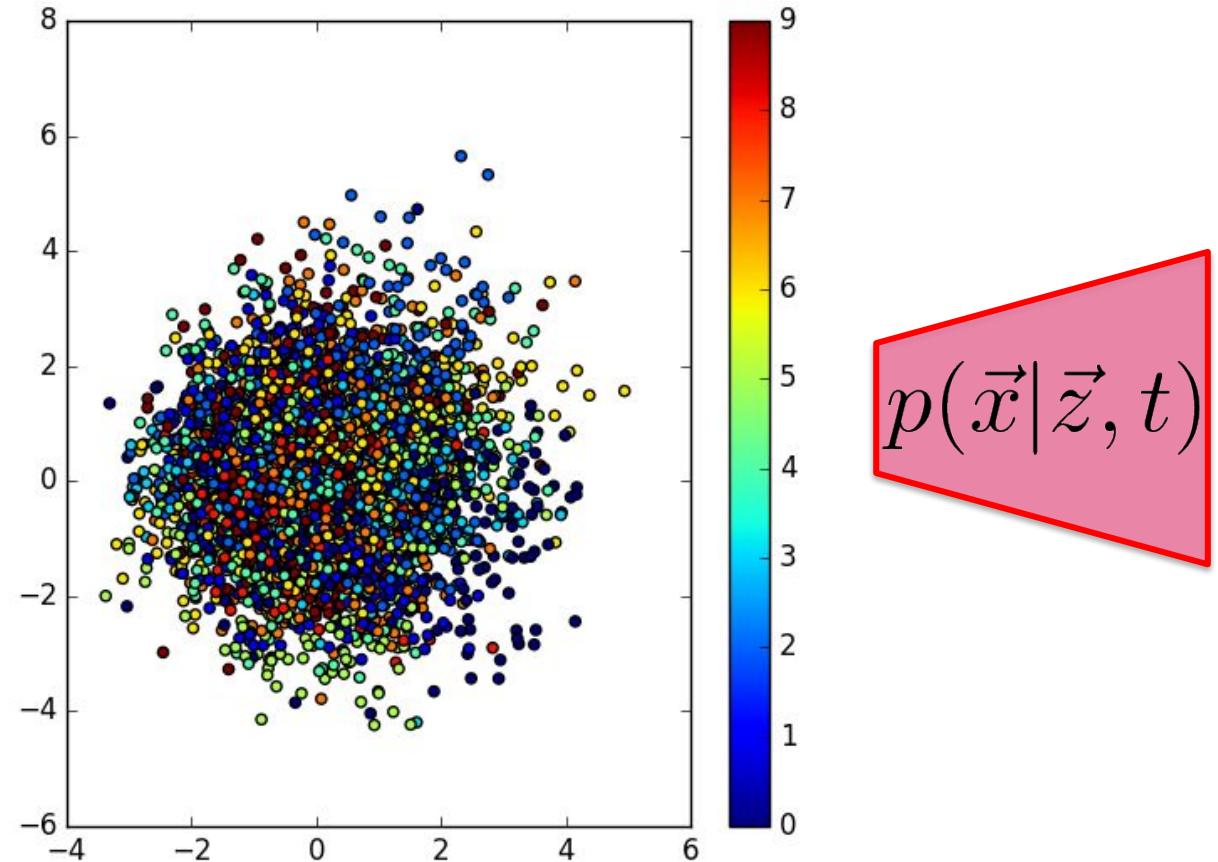


# VAE conditionnels

La solution: Injecter l'information de la classe dans notre encodeur et décodeur !

**Semble chaotique, mais ... strictement gaussien  $N(0, 1)$  pour toutes les classes !**

$$q(\vec{z}|\vec{x}, t)$$

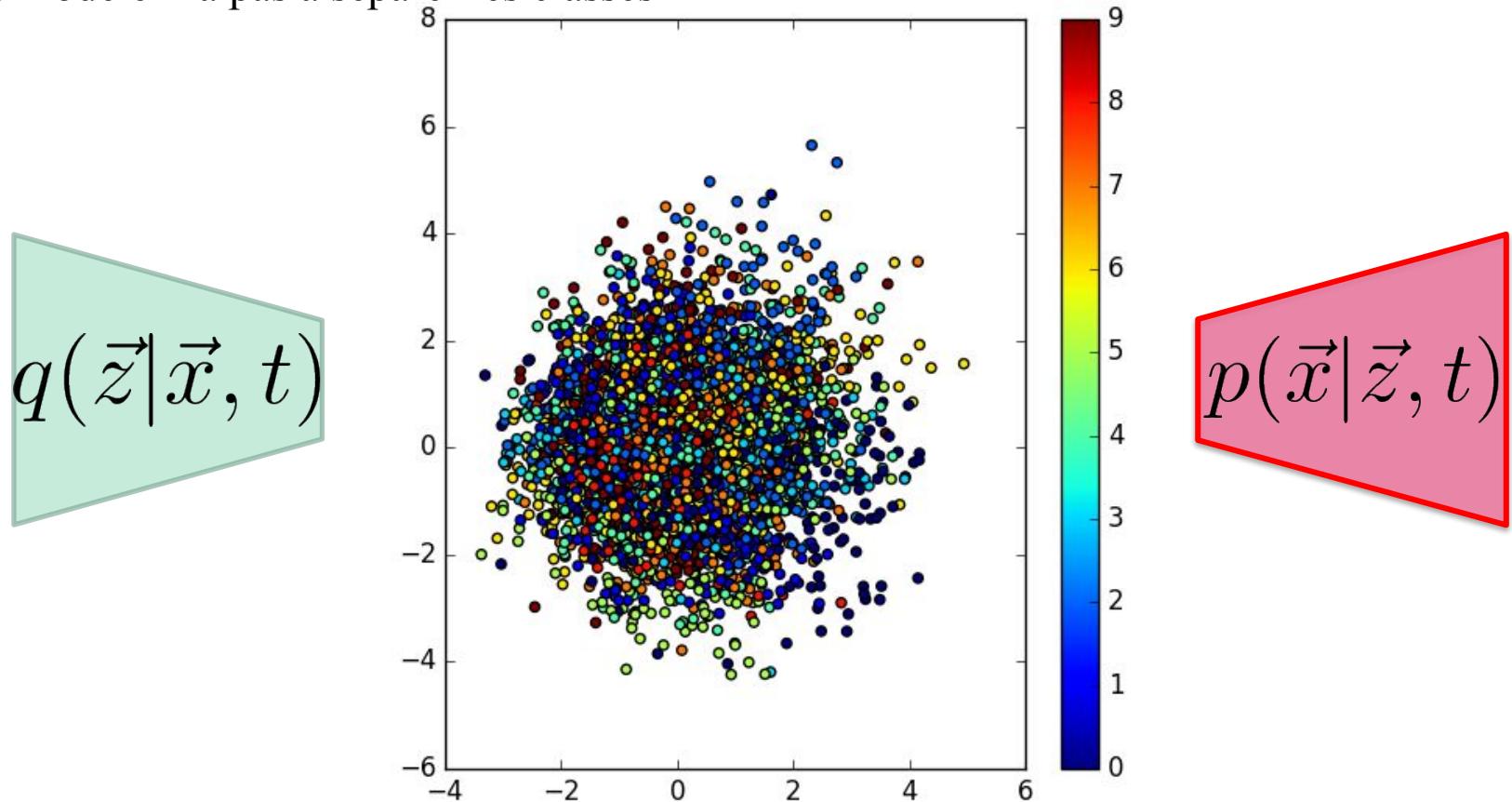


# VAE conditionnels

La solution: Injecter l'information de la classe dans notre encodeur et décodeur !

**Semble chaotique, mais ... strictement gaussien  $N(0, 1)$  pour toutes les classes !**

Le modèle n'a pas à séparer les classes



# VAE conditionnels

La solution:

Injecter l'information de la classe dans notre encodeur et décodeur !



# Plusieurs tutoriels, VAE

- <https://ijdykeman.github.io/ml/2016/12/21/cvae.html>
- <https://wiseodd.github.io/techblog/2016/12/10/variational-autoencoder/>
- <https://towardsdatascience.com/deep-latent-variable-models-unravel-hidden-structures-a5df0fd32ae2>
- C. Doersch, **Tutorial on Variational Autoencoders**, arXiv:1606.05908

## Autoencodeur variationnel

D.Kingma, M.Welling, Auto-Encoding Variational Bayes, arXiv:1312.6114v10 ([Annexe B](#))

## *ELBO loss : Evidence Lower Bound*

$$Loss = \underbrace{KL(N(z;0,1), N(z; \mu, \Sigma))}_{\text{Perte encodeur}} - \underbrace{\log(P(x | z))}_{\text{Perte décodeur}}$$

Si on suppose que  $P(x|z)$  est gaussien

$$Loss = \frac{1}{2} \sum_{i=1}^d (1 + \log(\sigma_i^2) + \mu_i^2 - \sigma_i^2) - \lambda ||\vec{x} - \vec{x}'||^2$$


Perte encodeur      Perte décodeur

# Autoencodeur variationnel

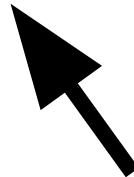
## *ELBO loss : Evidence Lower Bound*

On souhaite obtenir  $q(\vec{z}|\vec{x}) \approx p(\vec{z}|\vec{x})$ . Pour ce faire, on minimisera la KL-divergence entre les deux.

# Autoencodeur variationnel

## *ELBO loss : Evidence Lower Bound*

On souhaite obtenir  $q(\vec{z}|\vec{x}) \approx p(\vec{z}|\vec{x})$ . Pour ce faire, on minimisera la KL-divergence entre les deux.



Inconnu, impossible à estimer

On va devoir réorganiser quelques termes

# Autoencodeur variationnel

## *ELBO loss : Evidence Lower Bound*

On souhaite obtenir  $q(\vec{z}|\vec{x}) \approx p(\vec{z}|\vec{x})$ . Pour ce faire, on minimisera la KL-divergence entre les deux.

$$KL = [q(\vec{z}|\vec{x})||p(\vec{z}|\vec{x})] = - \int q(\vec{z}|\vec{x}) \log \frac{p(\vec{z}|\vec{x})}{q(\vec{z}|\vec{x})} d\vec{z}$$

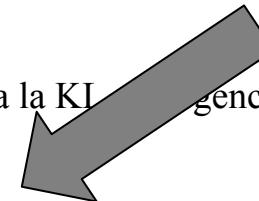
# Autoencodeur variationnel

Non symétrique !

## *ELBO loss : Evidence Lower Bound*

On souhaite obtenir  $q(\vec{z}|\vec{x}) \approx p(\vec{z}|\vec{x})$ . Pour ce faire, on minimisera la Kullback-Leibler divergence entre les deux.

$$KL = [q(\vec{z}|\vec{x})||p(\vec{z}|\vec{x})] = - \int q(\vec{z}|\vec{x}) \log \frac{p(\vec{z}|\vec{x})}{q(\vec{z}|\vec{x})} d\vec{z}$$



# Autoencodeur variationnel

Non symétrique !

## *ELBO loss : Evidence Lower Bound*

On souhaite obtenir  $q(\vec{z}|\vec{x}) \approx p(\vec{z}|\vec{x})$ . Pour ce faire, on minimisera la Kullback-Leibler divergence entre les deux.

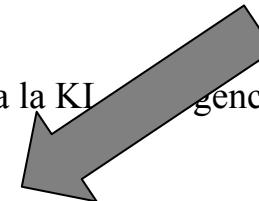
$$\begin{aligned} KL = [q(\vec{z}|\vec{x})||p(\vec{z}|\vec{x})] &= - \int q(\vec{z}|\vec{x}) \log \frac{p(\vec{z}|\vec{x})}{q(\vec{z}|\vec{x})} d\vec{z} \\ &= \int q(\vec{z}|\vec{x}) \log q(\vec{z}|\vec{x}) - \log p(\vec{z}|\vec{x}) d\vec{z} \end{aligned}$$

# Autoencodeur variationnel

## *ELBO loss : Evidence Lower Bound*

On souhaite obtenir  $q(\vec{z}|\vec{x}) \approx p(\vec{z}|\vec{x})$ . Pour ce faire, on minimisera la Kullback-Leibler divergence entre les deux.

Non symétrique !



Bayes:

$$p(\vec{z}|\vec{x}) = \frac{p(\vec{x}|\vec{z})p(\vec{z})}{p(\vec{x})}$$

$$\begin{aligned} KL &= [q(\vec{z}|\vec{x})||p(\vec{z}|\vec{x})] = - \int q(\vec{z}|\vec{x}) \log \frac{p(\vec{z}|\vec{x})}{q(\vec{z}|\vec{x})} d\vec{z} \\ &= \int q(\vec{z}|\vec{x}) \log q(\vec{z}|\vec{x}) - \log p(\vec{z}|\vec{x}) d\vec{z} \\ &= \int q(\vec{z}|\vec{x})(\log q(\vec{z}|\vec{x}) - \log p(\vec{x}|\vec{z}) - \log p(\vec{z}) + \log p(\vec{x})) d\vec{z} \end{aligned}$$

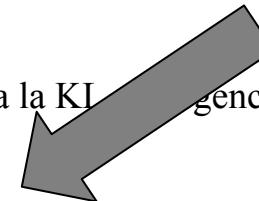
# Autoencodeur variationnel

## ELBO loss : Evidence Lower Bound

Non symétrique !

On souhaite obtenir  $q(\vec{z}|\vec{x}) \approx p(\vec{z}|\vec{x})$ . Pour ce faire, on minimisera la Kullback-Leibler divergence entre les deux.

$$KL = [q(\vec{z}|\vec{x})||p(\vec{z}|\vec{x})] = - \int q(\vec{z}|\vec{x}) \log \frac{p(\vec{z}|\vec{x})}{q(\vec{z}|\vec{x})} d\vec{z}$$



Bayes:

$$p(\vec{z}|\vec{x}) = \frac{p(\vec{x}|\vec{z})p(\vec{z})}{p(\vec{x})}$$

$$= \int q(\vec{z}|\vec{x}) \log q(\vec{z}|\vec{x}) - \log p(\vec{z}|\vec{x}) d\vec{z}$$

$$= \int q(\vec{z}|\vec{x})(\log q(\vec{z}|\vec{x}) - \log p(\vec{x}|\vec{z}) - \log p(\vec{z}) + \log p(\vec{x})) d\vec{z}$$

$$= \int q(\vec{z}|\vec{x}) \log \frac{q(\vec{z}|\vec{x})}{p(\vec{z})} d\vec{z} - \int q(\vec{z}|\vec{x}) \log p(\vec{x}|\vec{z}) d\vec{z} + \log p(\vec{z}) \int q(\vec{z}|\vec{x}) d\vec{z}$$

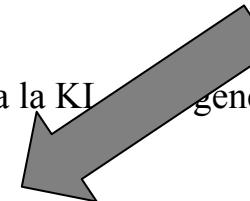
# Autoencodeur variationnel

## ELBO loss : Evidence Lower Bound

Non symétrique !

On souhaite obtenir  $q(\vec{z}|\vec{x}) \approx p(\vec{z}|\vec{x})$ . Pour ce faire, on minimisera la Kullback-Leibler divergence entre les deux.

$$KL = [q(\vec{z}|\vec{x})||p(\vec{z}|\vec{x})] = - \int q(\vec{z}|\vec{x}) \log \frac{p(\vec{z}|\vec{x})}{q(\vec{z}|\vec{x})} d\vec{z}$$



Bayes:

$$p(\vec{z}|\vec{x}) = \frac{p(\vec{x}|\vec{z})p(\vec{z})}{p(\vec{x})}$$

$$= \int q(\vec{z}|\vec{x}) \log q(\vec{z}|\vec{x}) - \log p(\vec{z}|\vec{x}) d\vec{z}$$

$$= \int q(\vec{z}|\vec{x})(\log q(\vec{z}|\vec{x}) - \log p(\vec{x}|\vec{z}) - \log p(\vec{z}) + \log p(\vec{x})) d\vec{z}$$

$$= \int q(\vec{z}|\vec{x}) \log \frac{q(\vec{z}|\vec{x})}{p(\vec{z})} d\vec{z} - \int q(\vec{z}|\vec{x}) \log p(\vec{x}|\vec{z}) d\vec{z} + \log p(\vec{z}) \int q(\vec{z}|\vec{x}) d\vec{z}$$

1



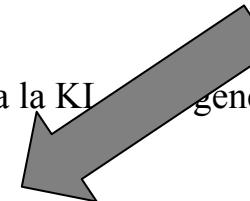
# Autoencodeur variationnel

## *ELBO loss : Evidence Lower Bound*

Non symétrique !

On souhaite obtenir  $q(\vec{z}|\vec{x}) \approx p(\vec{z}|\vec{x})$ . Pour ce faire, on minimisera la Kullback-Leibler divergence entre les deux.

$$KL = [q(\vec{z}|\vec{x})||p(\vec{z}|\vec{x})] = - \int q(\vec{z}|\vec{x}) \log \frac{p(\vec{z}|\vec{x})}{q(\vec{z}|\vec{x})} d\vec{z}$$



Bayes:

$$p(\vec{z}|\vec{x}) = \frac{p(\vec{x}|\vec{z})p(\vec{z})}{p(\vec{x})}$$

$$= \int q(\vec{z}|\vec{x}) \log q(\vec{z}|\vec{x}) - \log p(\vec{z}|\vec{x}) d\vec{z}$$

$$= \int q(\vec{z}|\vec{x})(\log q(\vec{z}|\vec{x}) - \log p(\vec{x}|\vec{z}) - \log p(\vec{z}) + \log p(\vec{x})) d\vec{z}$$

$$= \underbrace{\int q(\vec{z}|\vec{x}) \log \frac{q(\vec{z}|\vec{x})}{p(\vec{z})} d\vec{z}}_{KL[q(\vec{z}|\vec{x})||p(\vec{z})]} - \underbrace{\int q(\vec{z}|\vec{x}) \log p(\vec{x}|\vec{z}) d\vec{z}}_{\mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]} + \log p(\vec{x}) \int q(\vec{z}|\vec{x}) d\vec{z}$$

$$KL[q(\vec{z}|\vec{x})||p(\vec{z})]$$

$$\mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

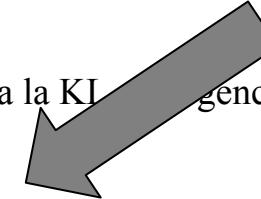
1

# Autoencodeur variationnel

## *ELBO loss : Evidence Lower Bound*

On souhaite obtenir  $q(\vec{z}|\vec{x}) \approx p(\vec{z}|\vec{x})$ . Pour ce faire, on minimisera la Kullback-Leibler divergence entre les deux.

$$KL = [q(\vec{z}|\vec{x})||p(\vec{z}|\vec{x})] = - \int q(\vec{z}|\vec{x}) \log \frac{p(\vec{z}|\vec{x})}{q(\vec{z}|\vec{x})} d\vec{z}$$



Non symétrique !

Bayes:

$$p(\vec{z}|\vec{x}) = \frac{p(\vec{x}|\vec{z})p(\vec{z})}{p(\vec{x})}$$

$$= \int q(\vec{z}|\vec{x}) \log q(\vec{z}|\vec{x}) - \log p(\vec{z}|\vec{x}) d\vec{z}$$

$$= \int q(\vec{z}|\vec{x})(\log q(\vec{z}|\vec{x}) - \log p(\vec{x}|\vec{z}) - \log p(\vec{z}) + \log p(\vec{x})) d\vec{z}$$

$$= \underbrace{\int q(\vec{z}|\vec{x}) \log \frac{q(\vec{z}|\vec{x})}{p(\vec{z})} d\vec{z}}_{KL[q(\vec{z}|\vec{x})||p(\vec{z})]} - \underbrace{\int q(\vec{z}|\vec{x}) \log p(\vec{x}|\vec{z}) d\vec{z}}_{\mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]} + \log p(\vec{x}) \int q(\vec{z}|\vec{x}) d\vec{z}$$

$$KL[q(\vec{z}|\vec{x})||p(\vec{z})]$$

$$\mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

1

En réorganisant les termes:

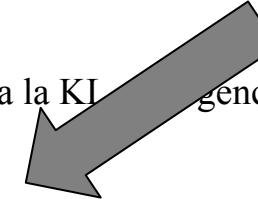
$$KL[q(\vec{z}|\vec{x})||p(\vec{z}|\vec{x})] - KL[q(\vec{z}|\vec{x})||p(\vec{z})] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})] = \log p(\vec{x})$$

# Autoencodeur variationnel

## *ELBO loss : Evidence Lower Bound*

On souhaite obtenir  $q(\vec{z}|\vec{x}) \approx p(\vec{z}|\vec{x})$ . Pour ce faire, on minimisera la Kullback-Leibler divergence entre les deux.

$$KL = [q(\vec{z}|\vec{x})||p(\vec{z}|\vec{x})] = - \int q(\vec{z}|\vec{x}) \log \frac{p(\vec{z}|\vec{x})}{q(\vec{z}|\vec{x})} d\vec{z}$$



Non symétrique !

Bayes:

$$p(\vec{z}|\vec{x}) = \frac{p(\vec{x}|\vec{z})p(\vec{z})}{p(\vec{x})}$$

$$= \int q(\vec{z}|\vec{x}) \log q(\vec{z}|\vec{x}) - \log p(\vec{z}|\vec{x}) d\vec{z}$$

$$= \int q(\vec{z}|\vec{x})(\log q(\vec{z}|\vec{x}) - \log p(\vec{x}|\vec{z}) - \log p(\vec{z}) + \log p(\vec{x})) d\vec{z}$$

$$= \underbrace{\int q(\vec{z}|\vec{x}) \log \frac{q(\vec{z}|\vec{x})}{p(\vec{z})} d\vec{z}}_{KL[q(\vec{z}|\vec{x})||p(\vec{z})]} - \underbrace{\int q(\vec{z}|\vec{x}) \log p(\vec{x}|\vec{z}) d\vec{z}}_{\mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]} + \log p(\vec{x}) \int q(\vec{z}|\vec{x}) d\vec{z}$$

$$KL[q(\vec{z}|\vec{x})||p(\vec{z})]$$

$$\mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

1

En réorganisant les termes:

$$KL[q(\vec{z}|\vec{x})||p(\vec{z}|\vec{x})] - \underbrace{KL[q(\vec{z}|\vec{x})||p(\vec{z})]}_{Evidence LOwer BOund (ELBO)} + \underbrace{\mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]}_{\text{Constante par rapport à } q!} = \log p(\vec{x})$$

*Evidence LOwer BOund (ELBO)*

Constante par rapport à  $q$  !

# Autoencodeur variationnel

## *ELBO loss : Evidence Lower Bound*

On souhaite obtenir  $q(\vec{z}|\vec{x}) \approx p(\vec{z}|\vec{x})$ . Pour ce faire, on minimisera la KL-divergence entre les deux.

$$KL[q(\vec{z}|\vec{x})||p(\vec{z}|\vec{x})] - \underbrace{KL[q(\vec{z}|\vec{x})||p(\vec{z})]}_{\text{Evidence Lower BOund (ELBO)}} + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})] = \underbrace{\log p(\vec{x})}_{\text{Constante par rapport à q !}}$$

KL divergence  $\geq 0$

8

ELBO

-12

$\log p(x) \leq 0$

-4

# Autoencodeur variationnel

## *ELBO loss : Evidence Lower Bound*

On souhaite obtenir  $q(\vec{z}|\vec{x}) \approx p(\vec{z}|\vec{x})$ . Pour ce faire, on minimisera la KL-divergence entre les deux.

$$KL[q(\vec{z}|\vec{x})||p(\vec{z}|\vec{x})] - \underbrace{KL[q(\vec{z}|\vec{x})||p(\vec{z})]}_{\text{Evidence Lower BOund (ELBO)}} + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})] = \underbrace{\log p(\vec{x})}_{\text{Constante par rapport à q !}}$$

KL divergence  $\geq 0$

8

3

ELBO

-12

-7

$\log p(\vec{x}) \leq 0$

-4

-4

# Autoencodeur variationnel

## *ELBO loss : Evidence Lower Bound*

On souhaite obtenir  $q(\vec{z}|\vec{x}) \approx p(\vec{z}|\vec{x})$ . Pour ce faire, on minimisera la KL-divergence entre les deux.

$$KL[q(\vec{z}|\vec{x})||p(\vec{z}|\vec{x})] - \underbrace{KL[q(\vec{z}|\vec{x})||p(\vec{z})]}_{Evidence Lower BOund (ELBO)} + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})] = \underbrace{\log p(\vec{x})}_{\text{Constante par rapport à } q !}$$

KL divergence  $\geq 0$

8

3

1

ELBO

-12

-7

-5

$\log p(\vec{x}) \leq 0$

-4

-4

-4

# Autoencodeur variationnel

## *ELBO loss : Evidence Lower Bound*

On souhaite obtenir  $q(\vec{z}|\vec{x}) \approx p(\vec{z}|\vec{x})$ . Pour ce faire, on minimisera la KL-divergence entre les deux.

$$KL[q(\vec{z}|\vec{x})||p(\vec{z}|\vec{x})] - \underbrace{KL[q(\vec{z}|\vec{x})||p(\vec{z})]}_{\text{Evidence Lower BOund (ELBO)}} + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})] = \underbrace{\log p(\vec{x})}_{\text{Constante par rapport à q !}}$$

KL divergence  $\geq 0$

8

3

1

0

ELBO

-12

-7

-5

-4

$\log p(\vec{x}) \leq 0$

-4

-4

-4

-4

On atteint notre objectif en minimisant la KL divergence ou en maximisant la ELBO

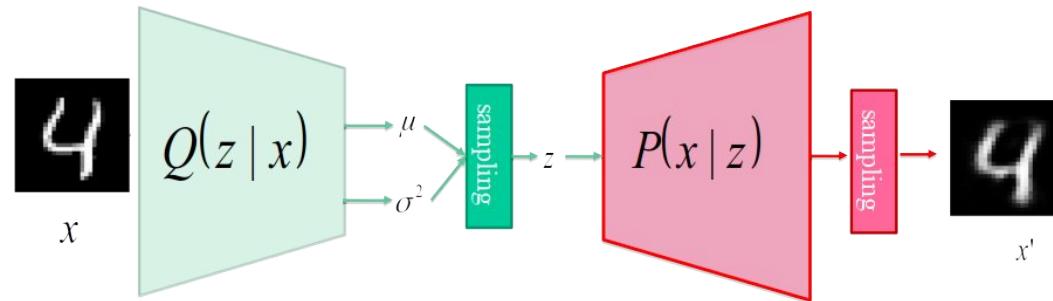
# Autoencodeur variationnel

*ELBO loss : Evidence Lower Bound*

On atteint notre objectif en minimisant la KL divergence ou en **maximisant la ELBO**

$$-KL[q(\vec{z}|\vec{x})||p(\vec{z})] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

*Evidence LOwer BOund (ELBO)*



On présume que  $p(\vec{z}) = \mathcal{N}(\vec{z}; 0, 1)$

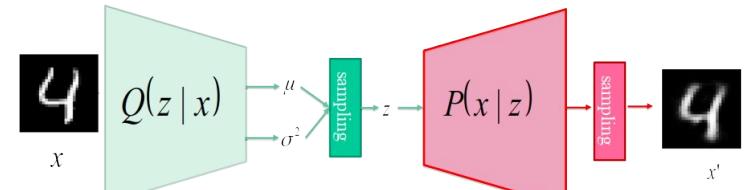
$$-KL[q(\vec{z}|\vec{x})||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

# Autoencodeur variationnel

*ELBO loss : Evidence Lower Bound*

On souhaite minimiser la loss

$$-KL[q(\vec{z}|\vec{x})||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$



# Autoencodeur variationnel

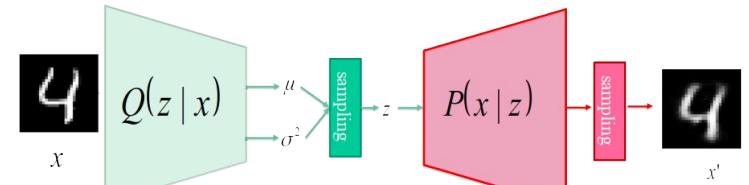
*ELBO loss : Evidence Lower Bound*

On souhaite minimiser la loss

$$-KL[q(\vec{z}|\vec{x})||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

On formule q par l'encodeur

$$-KL[\mathcal{N}(\vec{z}; \mu, \sigma^2)||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$



# Autoencodeur variationnel

*ELBO loss : Evidence Lower Bound*

On souhaite minimiser la loss

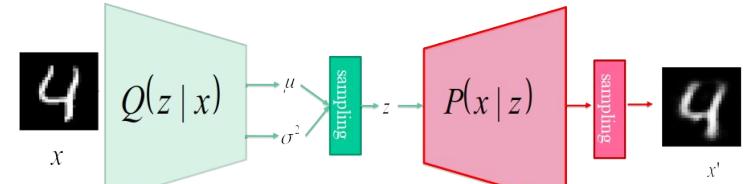
$$-KL[q(\vec{z}|\vec{x})||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

On formule q par l'encodeur

$$-KL[\mathcal{N}(\vec{z}; \mu, \sigma^2)||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

Définition de la KL-divergence pour deux gaussiennes:

$$KL[\mathcal{N}(\mu_1, \sigma_1^2)||\mathcal{N}(\mu_2, \sigma_2^2)] = \frac{(\mu_1 - \mu_2)^2}{2\sigma_2^2} + \frac{1}{2}\left(\frac{\sigma_1^2}{\sigma_2^2} - 1\right) \log \frac{\sigma_1^2}{\sigma_2^2}$$



# Autoencodeur variationnel

*ELBO loss : Evidence Lower Bound*

On souhaite minimiser la loss

$$-KL[q(\vec{z}|\vec{x})||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

On formule q par l'encodeur

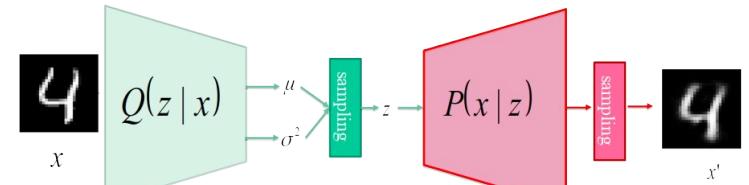
$$-KL[\mathcal{N}(\vec{z}; \mu, \sigma^2)||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

Définition de la KL-divergence pour deux gaussiennes:

$$KL[\mathcal{N}(\mu_1, \sigma_1^2)||\mathcal{N}(\mu_2, \sigma_2^2)] = \frac{(\mu_1 - \mu_2)^2}{2\sigma_2^2} + \frac{1}{2}\left(\frac{\sigma_1^2}{\sigma_2^2} - 1 - \log \frac{\sigma_1^2}{\sigma_2^2}\right)$$

En replaçant pour nos valeurs:

$$KL[\mathcal{N}(\vec{z}; \mu, \sigma^2)||\mathcal{N}(z; 0, 1)] = \frac{1}{2}(\mu - 0)^2 + \frac{1}{2}\left(\frac{\sigma^2}{1^2} - 1 - \log \frac{\sigma^2}{1^2}\right)$$



# Autoencodeur variationnel

*ELBO loss : Evidence Lower Bound*

On souhaite minimiser la loss

$$-KL[q(\vec{z}|\vec{x})||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

On formule q par l'encodeur

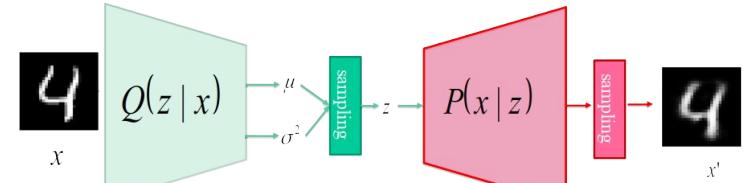
$$-KL[\mathcal{N}(\vec{z}; \mu, \sigma^2)||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

Définition de la KL-divergence pour deux gaussiennes:

$$KL[\mathcal{N}(\mu_1, \sigma_1^2)||\mathcal{N}(\mu_2, \sigma_2^2)] = \frac{(\mu_1 - \mu_2)^2}{2\sigma_2^2} + \frac{1}{2}\left(\frac{\sigma_1^2}{\sigma_2^2} - 1 - \log \frac{\sigma_1^2}{\sigma_2^2}\right)$$

En replaçant pour nos valeurs:

$$\begin{aligned} KL[\mathcal{N}(\vec{z}; \mu, \sigma^2)||\mathcal{N}(z; 0, 1)] &= \frac{1}{2}(\mu - 0)^2 + \frac{1}{2}\left(\frac{\sigma^2}{1^2} - 1 - \log \frac{\sigma^2}{1^2}\right) \\ &= \frac{1}{2}(\mu^2 + \sigma^2 - 1 - \log \sigma^2) \end{aligned}$$



# Autoencodeur variationnel

*ELBO loss : Evidence Lower Bound*

On souhaite minimiser la loss

$$-KL[q(\vec{z}|\vec{x})||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

On formule q par l'encodeur

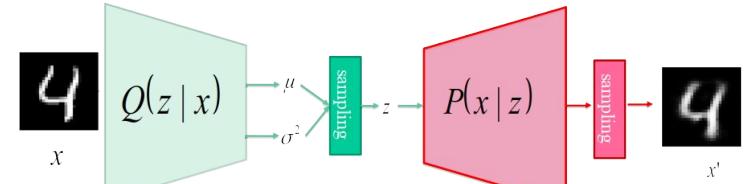
$$-KL[\mathcal{N}(\vec{z}; \mu, \sigma^2)||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

Définition de la KL-divergence pour deux gaussiennes:

$$KL[\mathcal{N}(\mu_1, \sigma_1^2)||\mathcal{N}(\mu_2, \sigma_2^2)] = \frac{(\mu_1 - \mu_2)^2}{2\sigma_2^2} + \frac{1}{2}\left(\frac{\sigma_1^2}{\sigma_2^2} - 1 - \log \frac{\sigma_1^2}{\sigma_2^2}\right)$$

En replaçant pour nos valeurs:

$$\begin{aligned} KL[\mathcal{N}(\vec{z}; \mu, \sigma^2)||\mathcal{N}(z; 0, 1)] &= \frac{1}{2}(\mu - 0)^2 + \frac{1}{2}\left(\frac{\sigma^2}{1^2} - 1 - \log \frac{\sigma^2}{1^2}\right) \\ &= \frac{1}{2}(\mu^2 + \sigma^2 - 1 - \log \sigma^2) \\ &= \frac{1}{2}(-1 - \log \sigma^2 - \mu^2 + \sigma^2) \end{aligned}$$



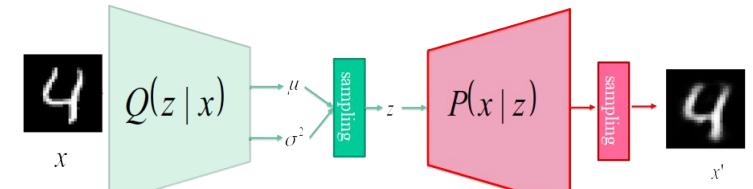
# Autoencodeur variationnel

*ELBO loss : Evidence Lower Bound*

On souhaite minimiser la loss

$$-KL[q(\vec{z}|\vec{x})||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

$$Loss = \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$



# Autoencodeur variationnel

*ELBO loss : Evidence Lower Bound*

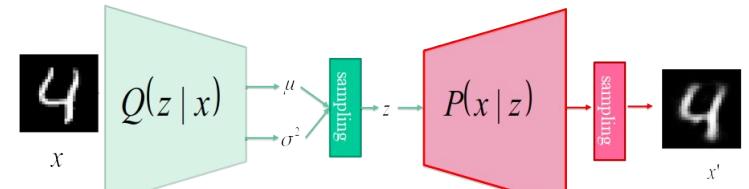
On souhaite minimiser la loss

$$-KL[q(\vec{z}|\vec{x})||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

$$Loss = \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

$$= \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \frac{1}{K} \sum_{i=1}^K [\log p(\vec{x}_i|\vec{z}_i)]$$

-



# Autoencodeur variationnel

*ELBO loss : Evidence Lower Bound*

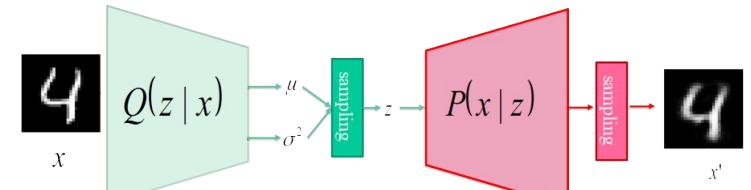
On souhaite minimiser la loss

$$-KL[q(\vec{z}|\vec{x})||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

$$Loss = \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

$$= \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \frac{1}{K} \sum_{i=1}^K [\log p(\vec{x}_i|\vec{z}_i)]$$

$$= \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \log p(\vec{x}|\vec{z})$$



# Autoencodeur variationnel

*ELBO loss : Evidence Lower Bound*

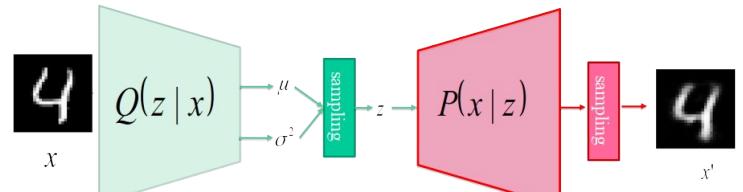
On souhaite minimiser la loss

$$-KL[q(\vec{z}|\vec{x})||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

$$Loss = \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

$$= \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \frac{1}{K} \sum_{i=1}^K [\log p(\vec{x}_i|\vec{z}_i)]$$

$$= \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \log p(\vec{x}|\vec{z})$$



$$p(\vec{x}|\vec{z}) = p(\vec{x}|\mathcal{N}(\vec{x}; \vec{x}', \sigma))$$

# Autoencodeur variationnel

*ELBO loss : Evidence Lower Bound*

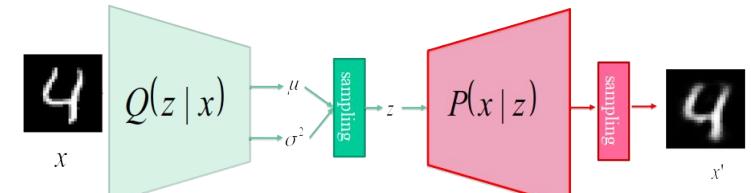
On souhaite minimiser la loss

$$-KL[q(\vec{z}|\vec{x})||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

$$Loss = \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

$$= \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \frac{1}{K} \sum_{i=1}^K [\log p(\vec{x}_i|\vec{z}_i)]$$

$$= \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \log p(\vec{x}|\vec{z})$$



$$p(\vec{x}|\vec{z}) = p(\vec{x}|\mathcal{N}(\vec{x}; \vec{x}', \sigma))$$

$$\log p(\vec{x}|\mathcal{N}(\vec{x}; \mu, \sigma)) = -\log \sigma - \frac{1}{2} \log(2\pi) - \frac{1}{2\pi\sigma^2} (\vec{x} - \mu)^2$$

# Autoencodeur variationnel

*ELBO loss : Evidence Lower Bound*

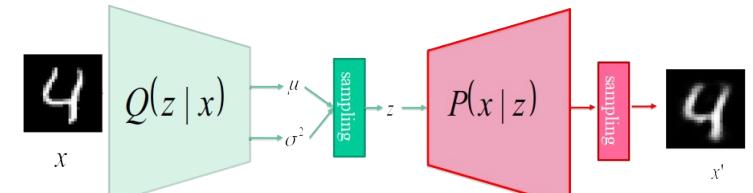
On souhaite minimiser la loss

$$-KL[q(\vec{z}|\vec{x})||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

$$Loss = \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

$$= \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \frac{1}{K} \sum_{i=1}^K [\log p(\vec{x}_i|\vec{z}_i)]$$

$$= \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \log p(\vec{x}|\vec{z})$$



$$p(\vec{x}|\vec{z}) = p(\vec{x}|\mathcal{N}(\vec{x}; \vec{x}', \sigma))$$

$$\log p(\vec{x}|\mathcal{N}(\vec{x}; \mu, \sigma)) = -\cancel{\log \sigma} - \frac{1}{2}\cancel{\log(2\pi)} - \frac{1}{2\pi\sigma^2}(\vec{x} - \mu)^2$$

# Autoencodeur variationnel

*ELBO loss : Evidence Lower Bound*

On souhaite minimiser la loss

$$-KL[q(\vec{z}|\vec{x})||\mathcal{N}(\vec{z}; 0, 1)] + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

$$Loss = \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \mathbb{E}_{z \sim q}[\log p(\vec{x}|\vec{z})]$$

$$= \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \frac{1}{K} \sum_{i=1}^K [\log p(\vec{x}_i|\vec{z}_i)]$$

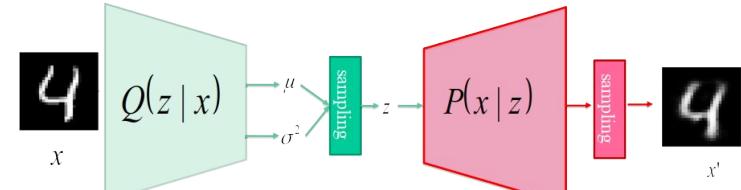
$$= \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) + \log p(\vec{x}|\vec{z})$$

$$= \frac{1}{2}(1 + \log \sigma^2 + \mu^2 - \sigma^2) - \lambda(\vec{x} - \vec{x}')^2$$

$$p(\vec{x}|\vec{z}) = p(\vec{x}|\mathcal{N}(\vec{x}; \vec{x}', \sigma))$$

$$\log p(\vec{x}|\mathcal{N}(\vec{x}; \mu, \sigma)) = -\cancel{\log \sigma} - \frac{1}{2}\cancel{\log(2\pi)} - \frac{1}{2\pi\sigma^2}(\vec{x} - \mu)^2$$

$$\log p(\vec{x}|\mathcal{N}(\vec{x}; \vec{x}', \sigma)) = -(\vec{x} - \vec{x}')^2$$



## Autoencodeur variationnel

D.Kingma, M.Welling, Auto-Encoding Variational Bayes, arXiv:1312.6114v10 ([Annexe B](#))

## *ELBO loss : Evidence Lower Bound*

$$Loss = \underbrace{KL(N(z;0,1), N(z; \mu, \Sigma))}_{\text{Perte encodeur}} - \underbrace{\log(P(x | z))}_{\text{Perte décodeur}}$$

Si on suppose que  $P(x|z)$  est gaussien

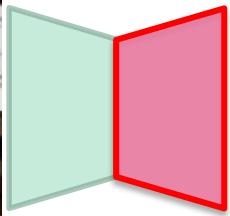
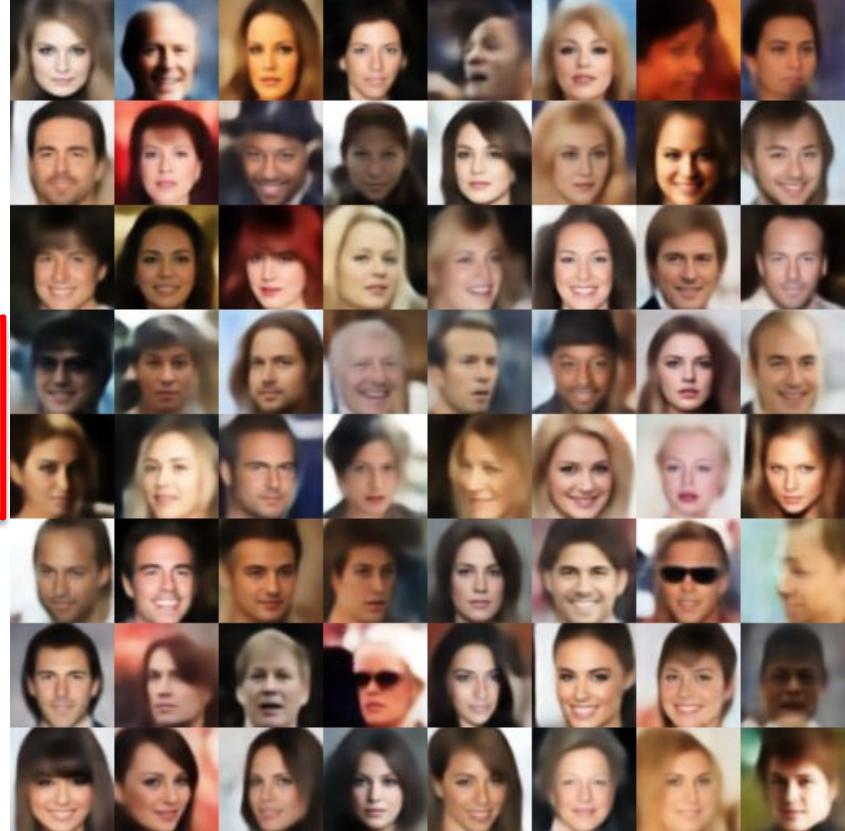
$$Loss = \frac{1}{2} \sum_{i=1}^d (1 + \log(\sigma_i^2) + \mu_i^2 - \sigma_i^2) - \lambda ||\vec{x} - \vec{x}'||^2$$


# Ex.: base de données *CelebA*

$x$



$x'$



# GAN

Generative Adversarial Nets

On voudrait générer des images  $\vec{x}$  en échantillonnant  $p(\vec{x})$

=> **TROP DIFFICILE** car  $p(\vec{x})$  trop complexe



Comme précédemment, pour simplifier le problème, on pourrait introduire une variable latente  $\vec{z}$  et ainsi modéliser

$$p(\vec{x}, \vec{z}) = p(\vec{x}|\vec{z})p(\vec{z})$$

Modèle génératif

Distribution *a priori*



Comme pour les VAE, on utilisera une **distribution *a priori*** facile à échantillonner : une **gaussienne!**

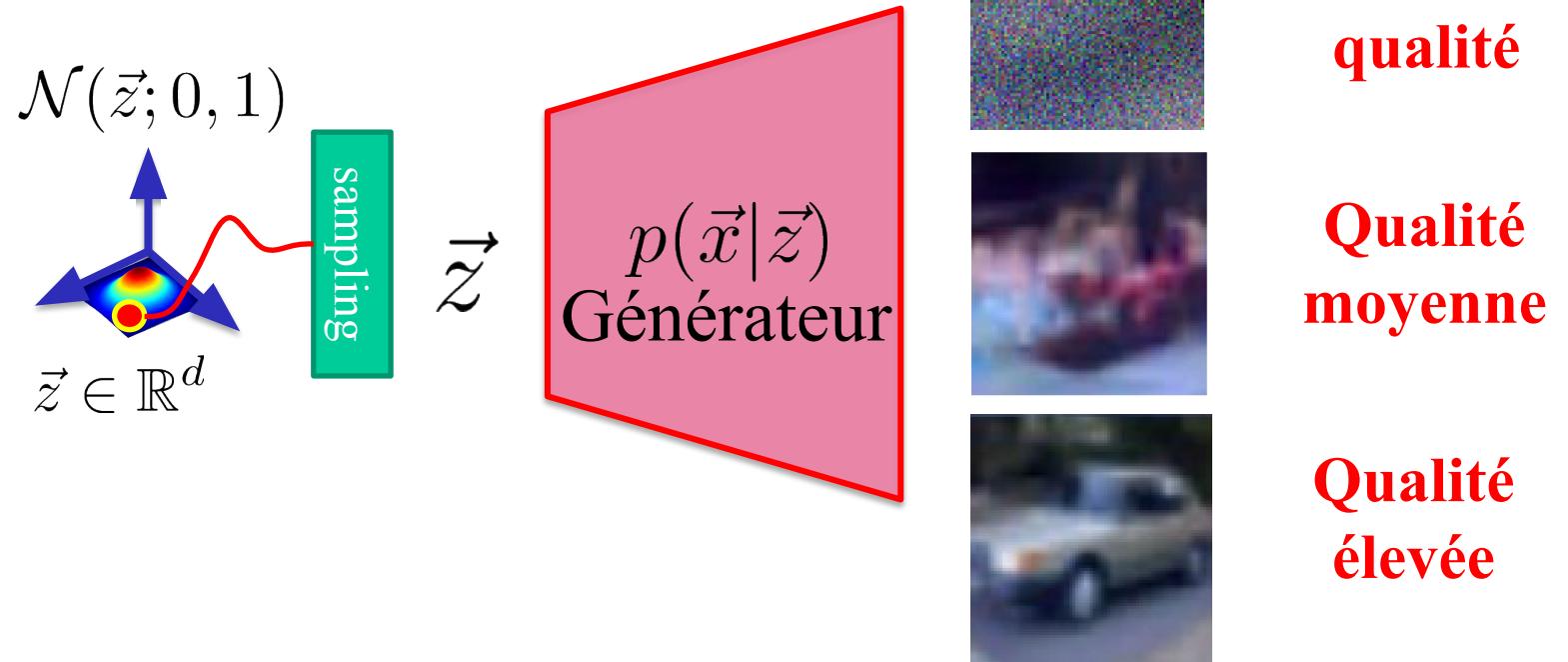
$$p(\vec{z}) = \mathcal{N}(\vec{z}; 0, 1)$$

Comment estimer  $p(\vec{x}|\vec{z})$  ?

À l'aide d'un réseau de neurones car ce sont **d'excellentes machines pour estimer des probabilités conditionnelles**



Dépendamment des performances du décodeur, les images générées  
**seront de qualité très variable.**



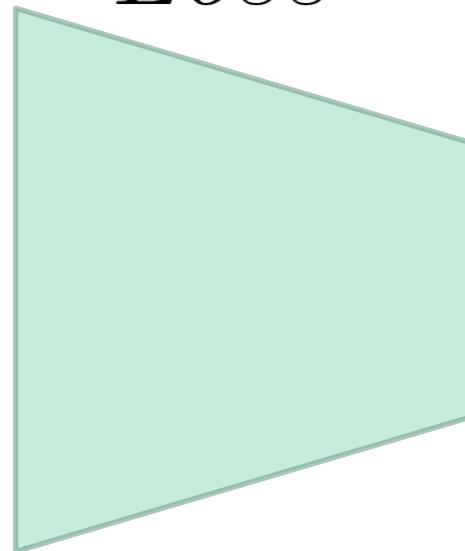
Pour entraîner un décodeur, il faut une **fonction de perte** (*loss*) qui **mesure la qualité (degré de réalisme) des images produites**

Pour un autoencodeur (variationnel ou non) c'est facile!  
car on a un encodeur et une image de référence

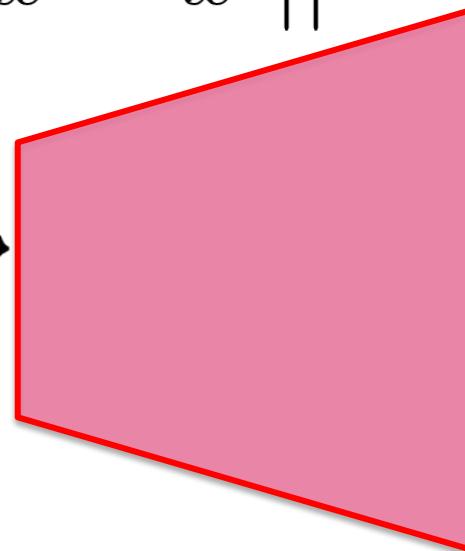
$$Loss = ||\vec{x} - \vec{x}'||^2$$



$\vec{x}$



$\vec{z}$

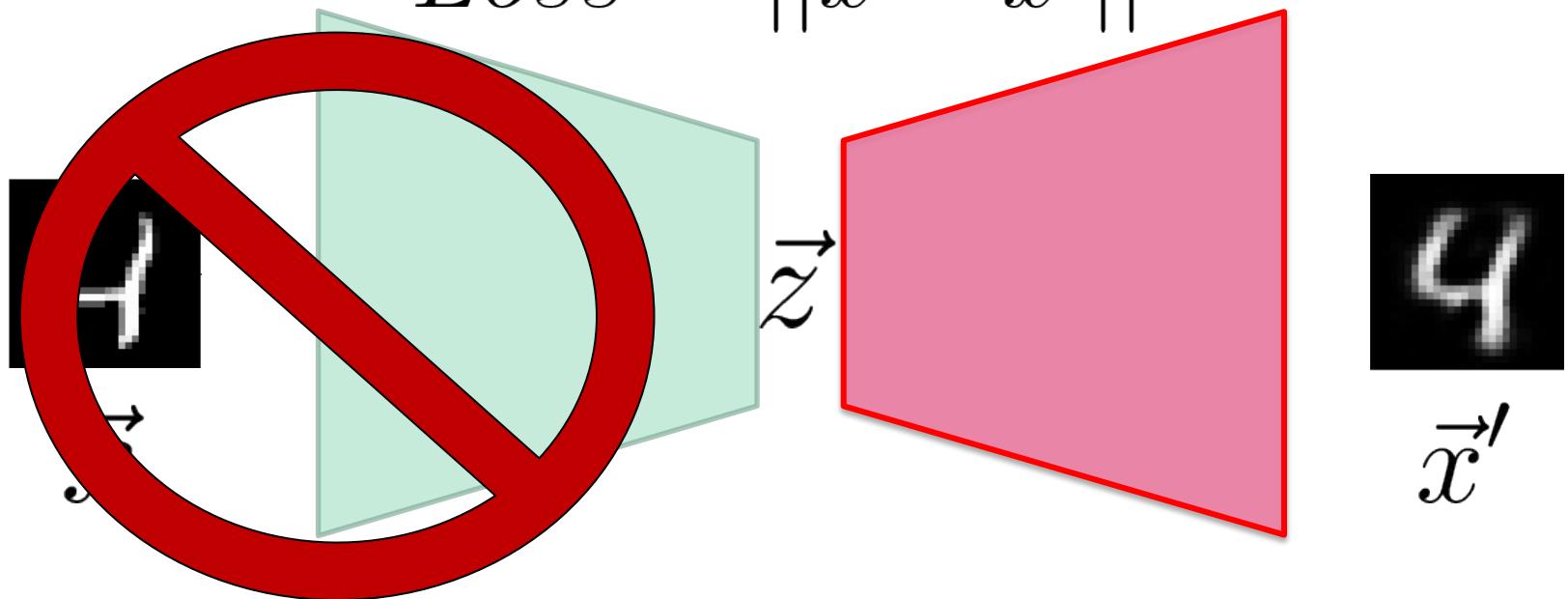


$\vec{x}'$

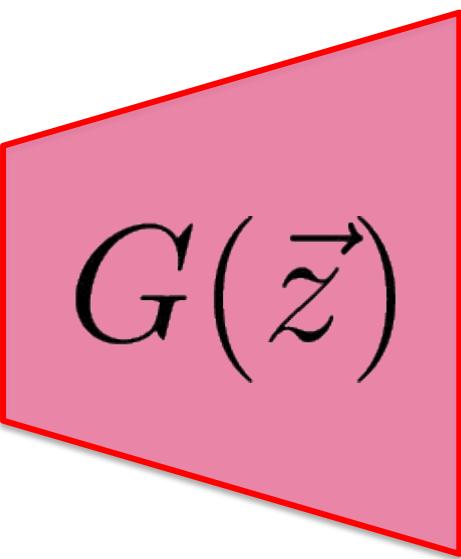
Pour entraîner un décodeur, il faut une **fonction de perte** (*loss*) qui **mesure la qualité (degré de réalisme) des images produites**

Comment faire pour un réseau **sans encodeur**?

$$Loss = \|\vec{x} - \vec{x}'\|^2$$





$$\vec{z} \sim p(\vec{z})$$
A pink parallelogram with a red double-line border. Inside the parallelogram, the text  $G(\vec{z})$  is written in a black serif font.


$$\vec{x}$$

Loss sans cible de référence?



Approximer la perte à l'aide d'un

**2<sup>e</sup> réseau de neurones**

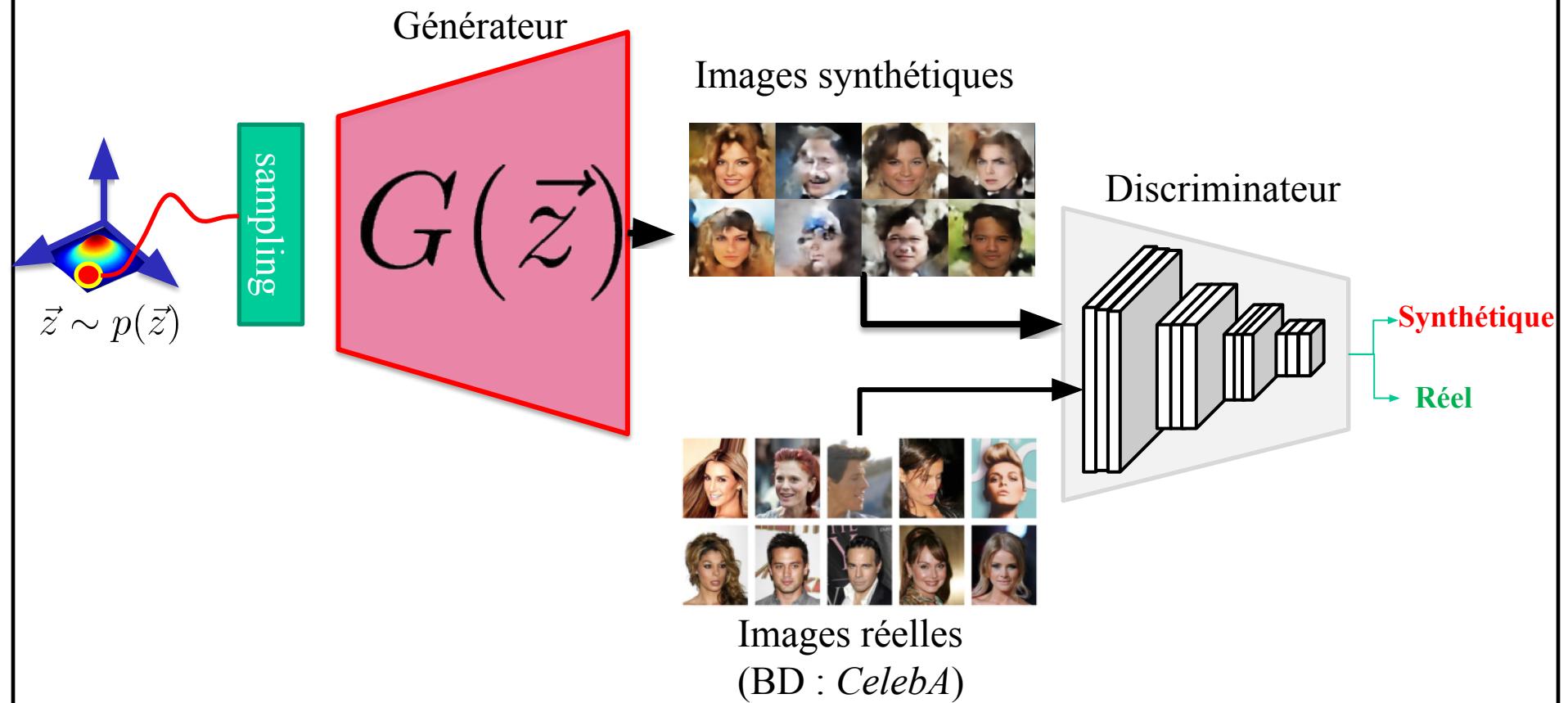
et d'une

**Base de données de référence**

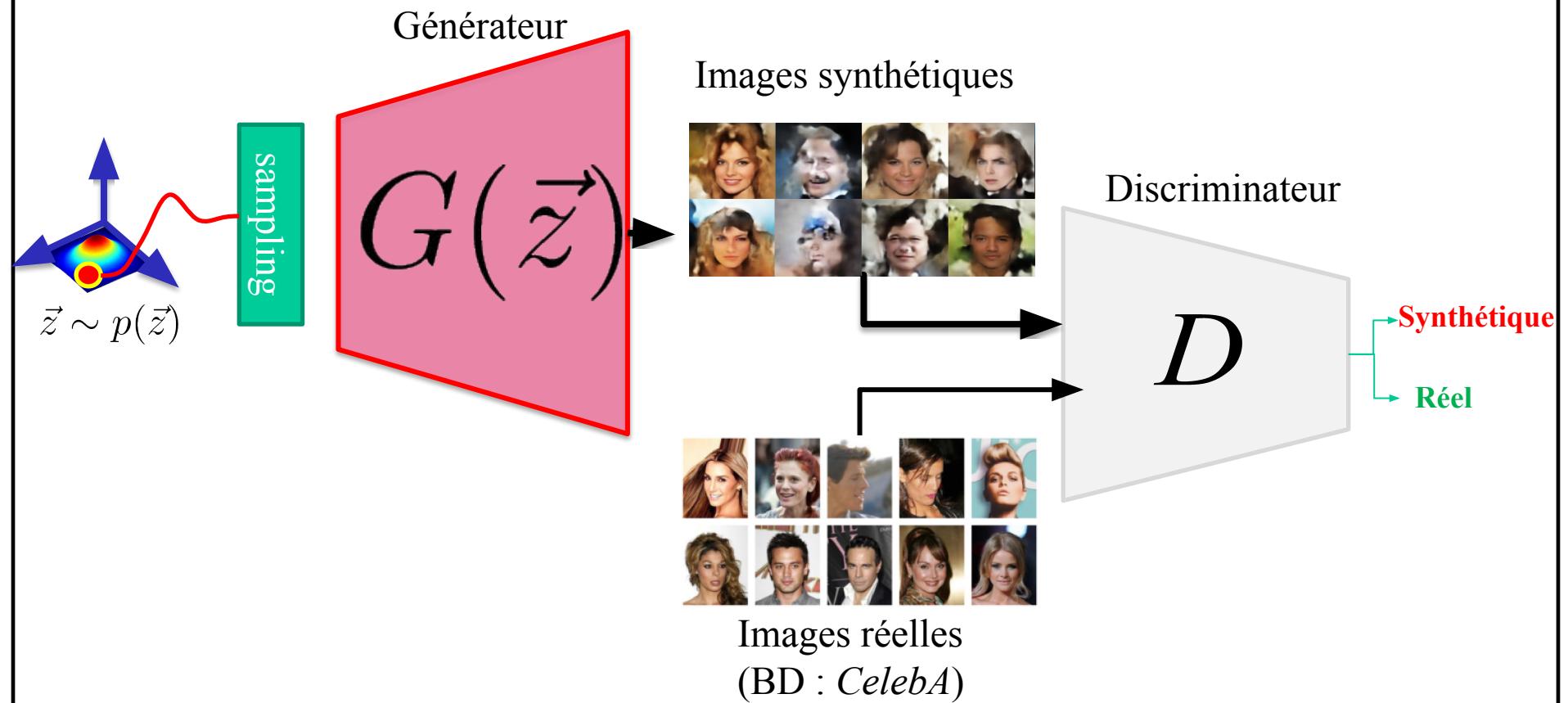
$$\vec{z} \sim p(\vec{z}$$

référence?

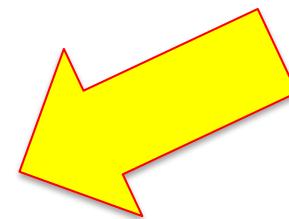
Pour entraîner un décodeur, il faut une **fonction de perte** (loss) qui mesure la qualité des images produites



Pour entraîner un décodeur, il faut une **fonction de perte** (loss) qui mesure la qualité des images produites



# Données étiquetées

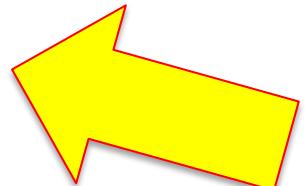


Produites par le  
générateur

$t = 0$  ('synthétique')



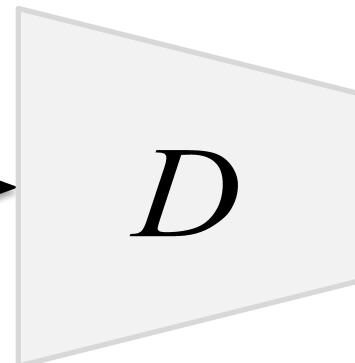
$t = 1$  ('réel')



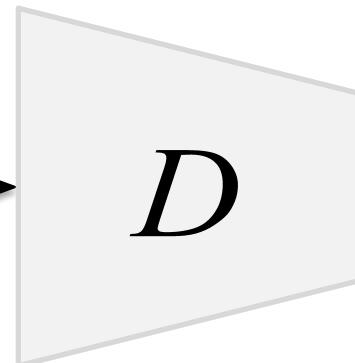
Issues d'une vraie BD

Deux réseaux aux **objectifs différents** :

**Discriminateur** : différencie les images synthétiques des images réelles



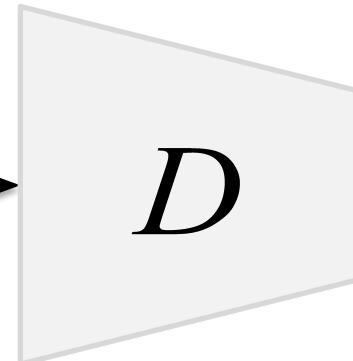
Synthétique



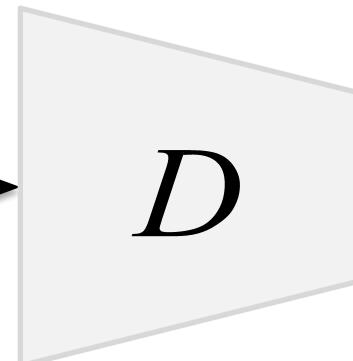
Réel

# Discriminateur : classifieur binaire (régression logistique)

=> Perte l'entropie croisée



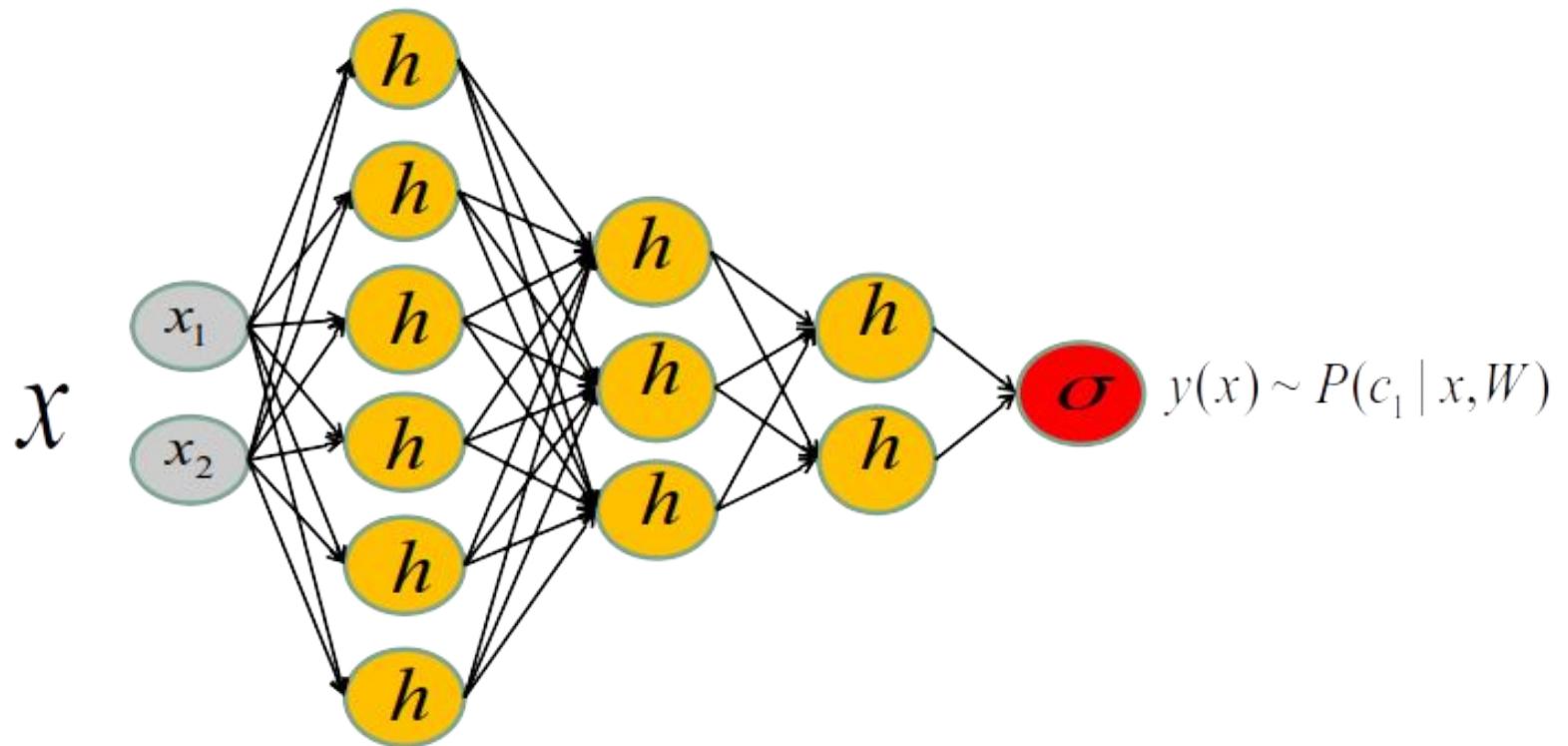
Synthétique



Réel

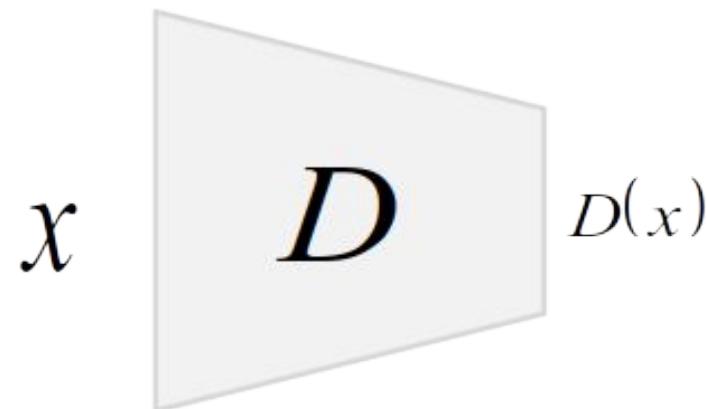
**Rappel, entropie croisée** pour une régression logistique binaire:

$$L_D = \frac{1}{N} \sum_i -t_i \ln(\vec{y}(x_i)) - (1 - t_i) \ln(1 - \vec{y}(x_i))$$



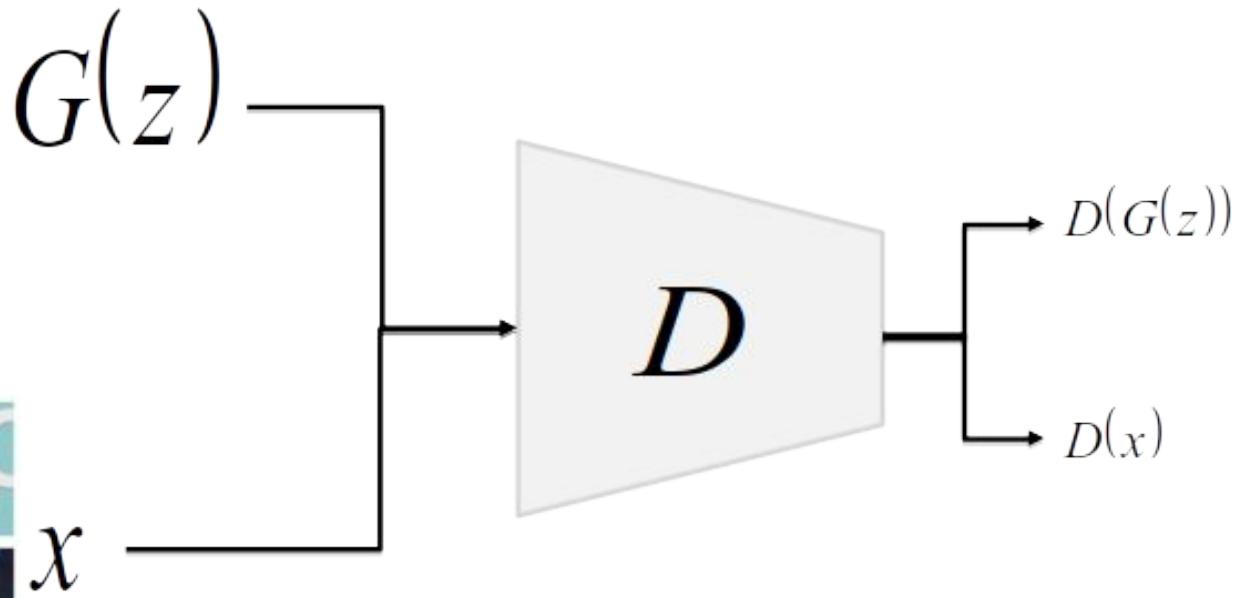
Le réseau discriminateur est représenté par la **lettre D**

$$L_D = \frac{1}{N} \sum_i -t_i \ln(\vec{D}(x_i)) - (1 - t_i) \ln(1 - \vec{D}(x_i))$$



Puisque les images **synthétiques** ont été générées par le **générateur**

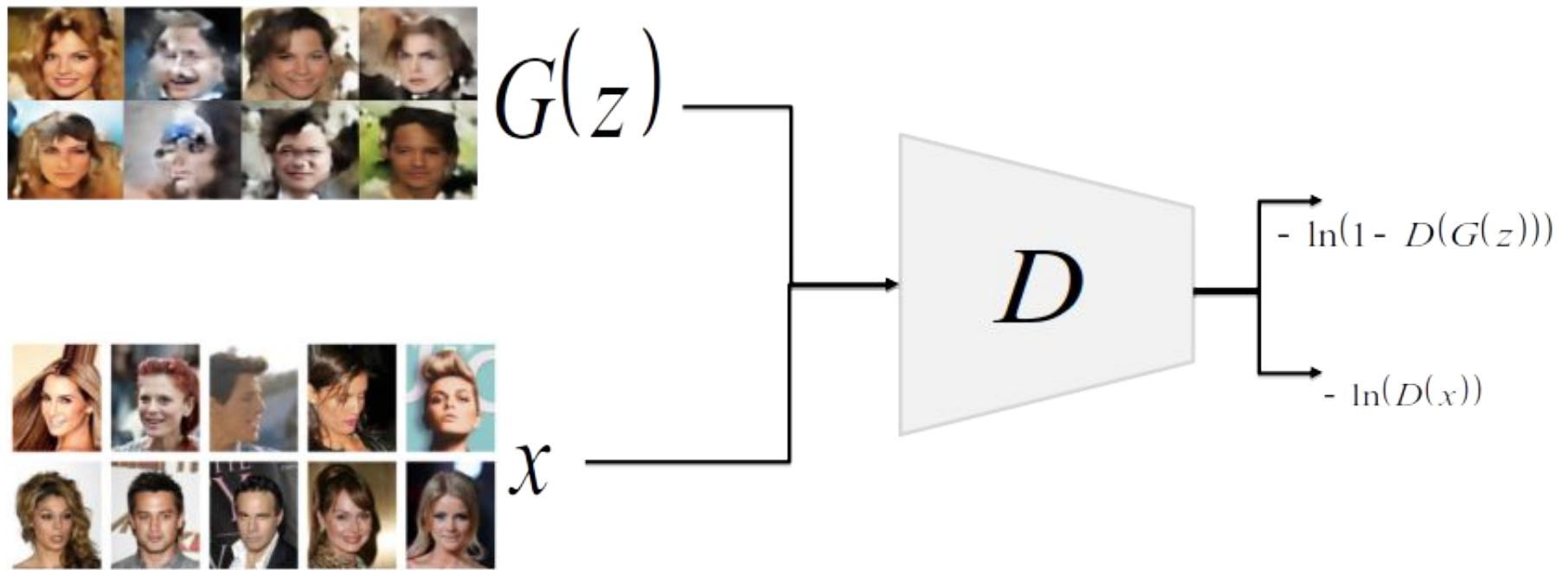
$$L_D = \frac{1}{N} \sum_i -t_i \ln(\vec{D}(\vec{x}_i)) - (1 - t_i) \ln(1 - \vec{D}(\vec{G}(\vec{z}_i)))$$



Sans perte de généralité, séparer la loss des images réelles et synthétiques

$$L_D = - \frac{1}{N_{reel}} \sum_i \ln(D(\vec{x}_i)) - \frac{1}{N_{syn}} \sum_j \ln(1 - D(\vec{G}(\vec{z}_j)))$$

Perte images réelles
Perte images synthétiques



## Rappel: Espérance mathématique et approximation Monte Carlo

$$IE[x] = \int xp(x)dx$$

$$IE[f(x)] = \int f(x)p(x)dx$$

## Rappel: Espérance mathématique et approximation Monte Carlo

$$IE[x] = \int xp(x)dx$$

$$\approx \frac{1}{N} \sum_{i=1}^N x_i \quad \text{où } x_i \sim p(x)$$

approximation  
Monte Carlo

$$IE[f(x)] = \int f(x)p(x)dx$$

$$\approx \frac{1}{N} \sum_{i=1}^N f(x_i) \quad \text{où } x_i \sim p(x)$$

## Rappel: Espérance mathématique et estimateur Monte Carlo

$$L_D = - \frac{1}{N_{real}} \sum_i \ln(D(\vec{x}_i)) - \frac{1}{N_{syn}} \sum_j \ln(1 - D(G(\vec{z}_j)))$$

Perte images réelles
Perte images synthétiques

$$L_D = -IE_{x \sim P_{real}}[\ln(D(x))] - IE_{z \sim P_z}[\ln(1 - D(G(z)))]$$

## Rappel: Espérance mathématique et estimateur Monte Carlo

$$L_D = -IE_{x \sim P_{real}}[\ln(D(x))] - IE_{z \sim P_z}[\ln(1 - D(G(z)))]$$

# Objectif du discriminateur

## Paramètres du discriminateur

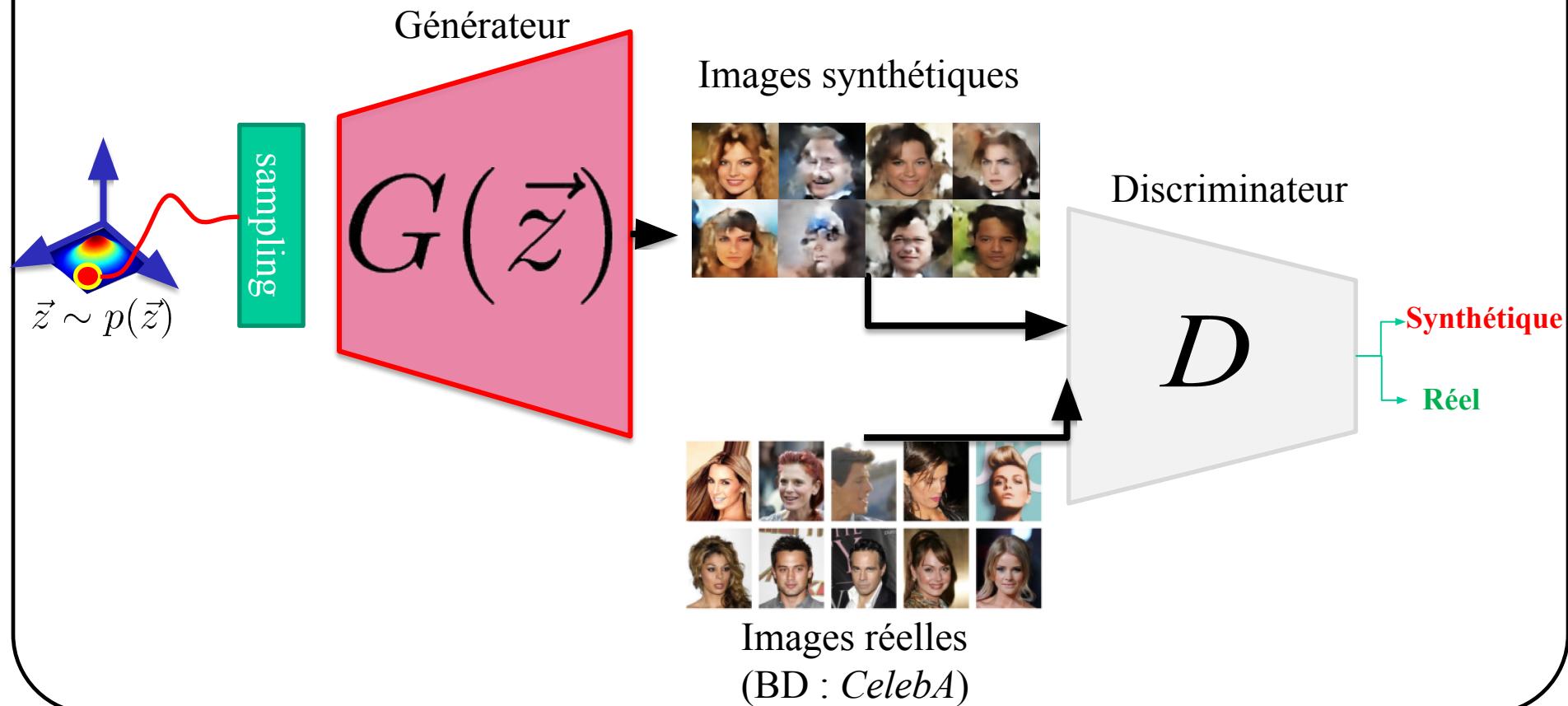
$$W_D = \arg \min_{W_D} -IE_{x \sim P_{real}} [\ln(D(x))] - IE_{z \sim P_z} [\ln(1 - D(G(z)))]$$

Ou encore, de façon équivalente

$$W_D = \arg \max_{W_D} IE_{x \sim P_{real}} [\ln(D(x))] + IE_{z \sim P_z} [\ln(1 - D(G(z)))]$$

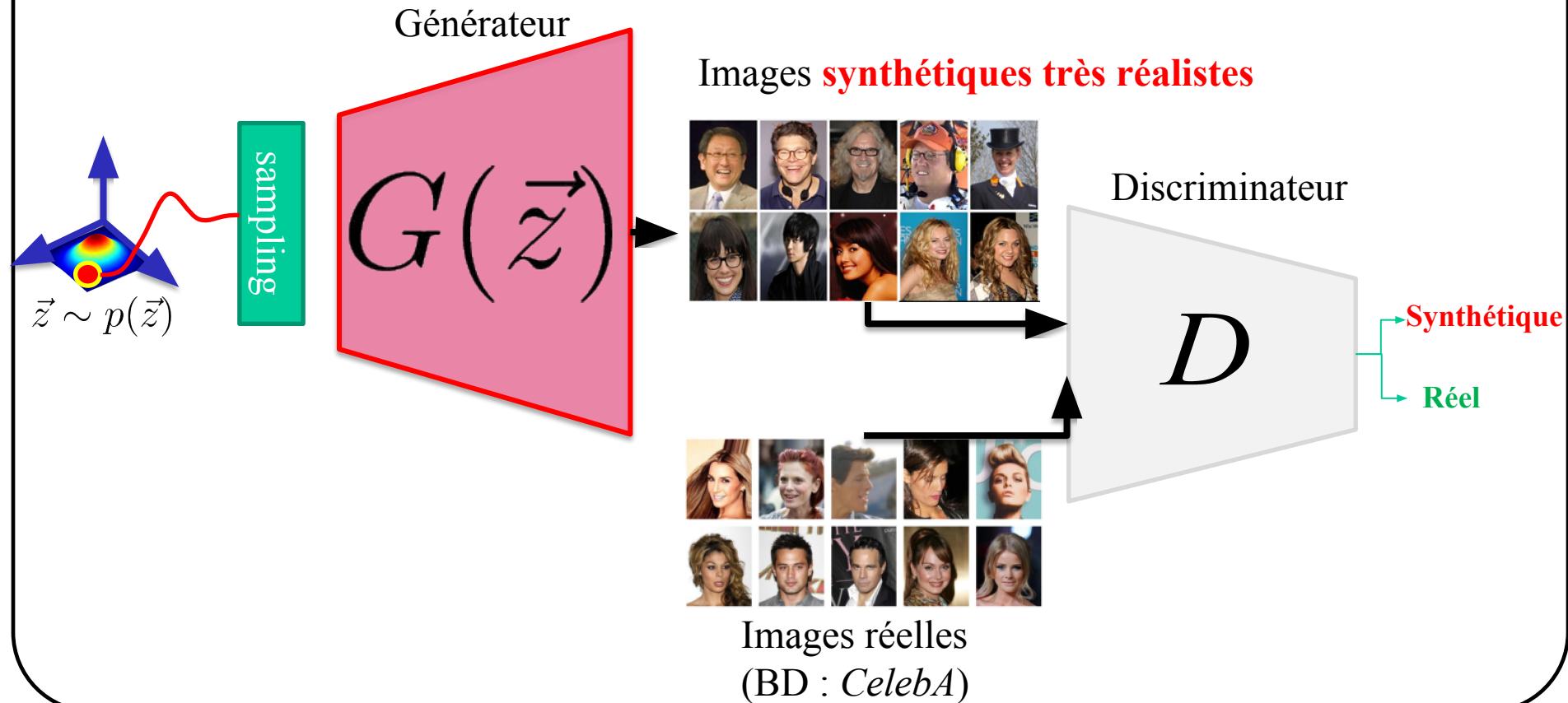
# Objectif du générateur

Produire des images aussi réalistes que celle de la BD de référence



## Objectif du générateur

S'il y parvient, le discriminateur ne pourra plus les distinguer des images réelles. La loss du discriminateur sera alors élevée.



## Objectif du discriminateur :

bien distinguer les images réelles des images synthétiques

$$W_D = \arg \max_{W_D} IE_{x \sim P_{real}} [\ln(D(x))] + IE_{z \sim P_z} [\ln(1 - D(G(z)))]$$

## Objectif du générateur :

produire des images synthétiques indistinguables des images réelles

$$W_G = \arg \min_{W_G} IE_{z \sim P_z} [\ln(1 - D(G(z)))]$$

# « Two player » mini-max game

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))].$$

Ian Goodfellow et al., “Generative Adversarial Nets”, NIPS 2014

# « Two player » mini-max game

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))].$$

Discriminateur veux  
 $D(x) = 1$  pour les vrais données

Discriminateur veux  
 $D(G(x)) = 0$  pour les données synthétiques

Générateur veux  
 $D(G(x)) = 1$  pour les données synthétiques

Ian Goodfellow et al., “Generative Adversarial Nets”, NIPS 2014

## NOTE

dans les faits, on ne minimise pas cette loss

$$W_G = \arg \min_{W_G} IE_{z \sim P_z} [\ln(1 - D(G(z)))]$$

on maximise plutôt celle-ci

$$W_G = \arg \max_{W_G} IE_{z \sim P_z} [\ln(D(G(z)))]$$

**for** number of training iterations **do**

**for**  $k$  steps **do**

- Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_g(\mathbf{z})$ .
- Sample minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from data generating distribution  $p_{\text{data}}(\mathbf{x})$ .
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[ \log D_{\theta_d}(\mathbf{x}^{(i)}) + \log(1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z}^{(i)}))) \right]$$

**end for**

- Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_g(\mathbf{z})$ .
- Update the generator by ascending its stochastic gradient (improved objective):

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(D_{\theta_d}(G_{\theta_g}(\mathbf{z}^{(i)})))$$

**end for**



a)



b)



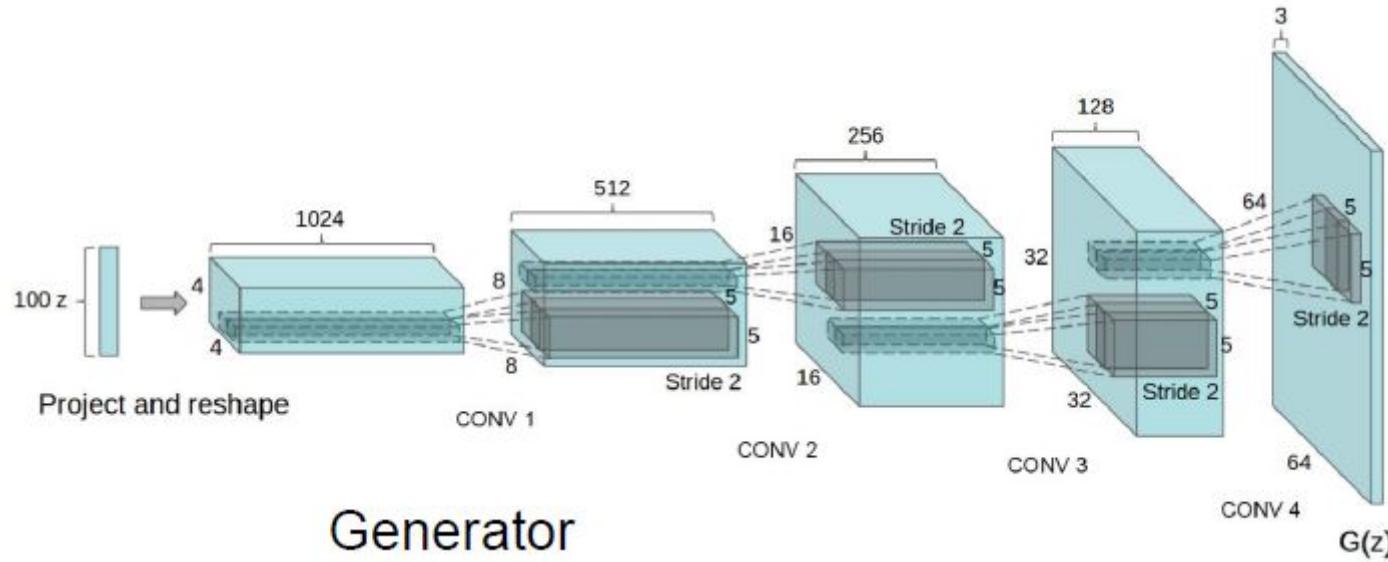
c)



d)

Figure 2: Visualization of samples from the model. Rightmost column shows the nearest training example of the neighboring sample, in order to demonstrate that the model has not memorized the training set. Samples are fair random draws, not cherry-picked. Unlike most other visualizations of deep generative models, these images show actual samples from the model distributions, not conditional means given samples of hidden units. Moreover, these samples are uncorrelated because the sampling process does not depend on Markov chain mixing. a) MNIST b) TFD c) CIFAR-10 (fully connected model) d) CIFAR-10 (convolutional discriminator and “deconvolutional” generator)

# Deep Convolution Generative Adversarial Net (DCGAN)



Radford et al, “Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks”, ICLR 2016

# Deep Convolution Generative Adversarial Net (DCGAN)

## Recommandations discriminateur

- Conv stride>1 au lieu des couches de pooling
- ReLU partout sauf en sortie : tanh

## Recommandations générateur

- Conv transpose au lieu de upsampling
- LeakyReLU partout

## Autre recommandations

- BatchNorm partout
- Pas de FC, juste des conv

Radford et al, “Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks”, ICLR 2016

# Deep Convolution Generative Adversarial Net (DCGAN)

## Recommandations discriminateur

- Conv
- Pool

<https://github.com/soumith/ganhacks>

## Recommendations

- Conv
- Leaky ReLU partout

## Autre recommandations

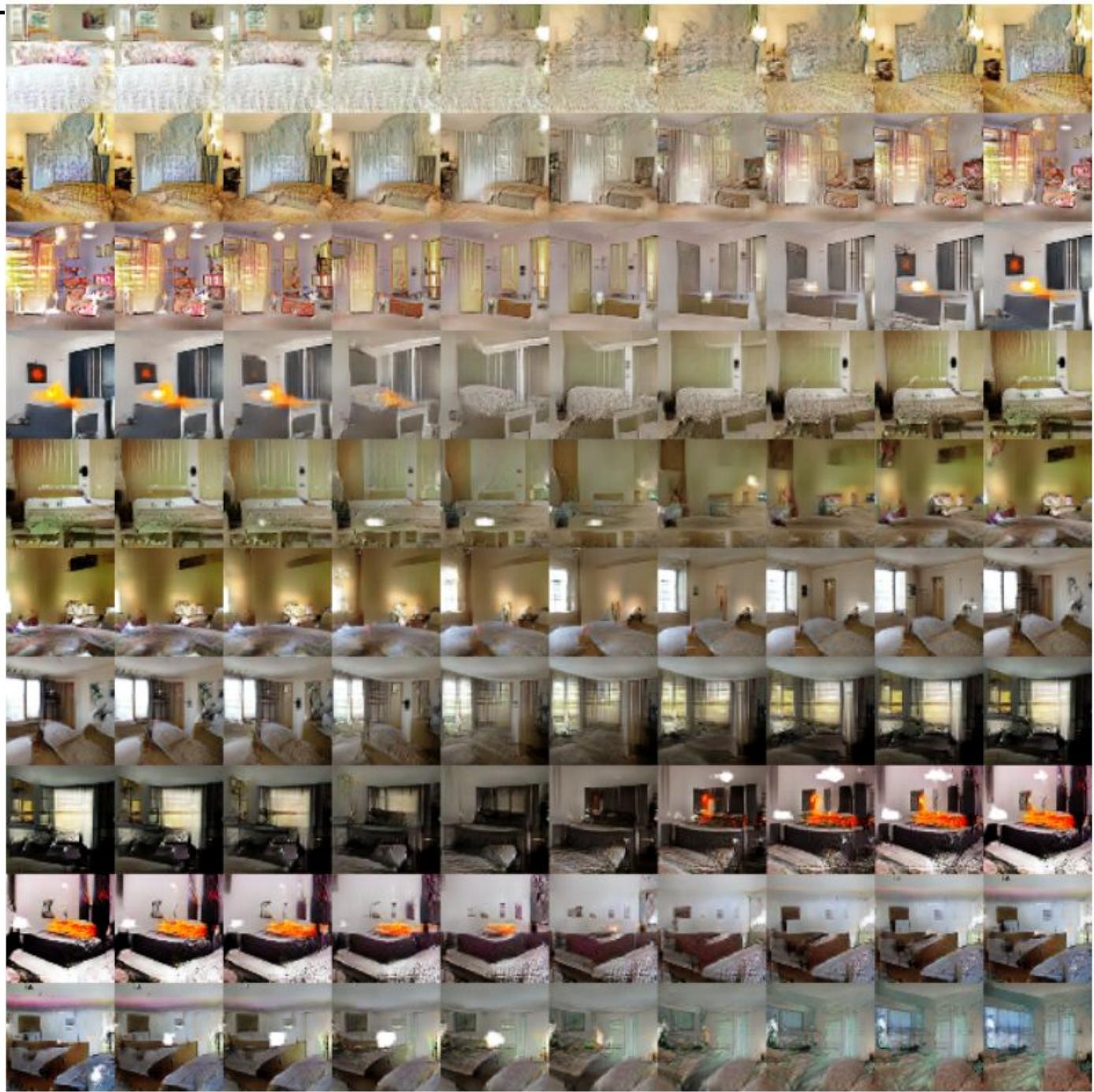
- BatchNorm partout
- Pas de FC, juste des conv

Radford et al, “Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks”, ICLR 2016

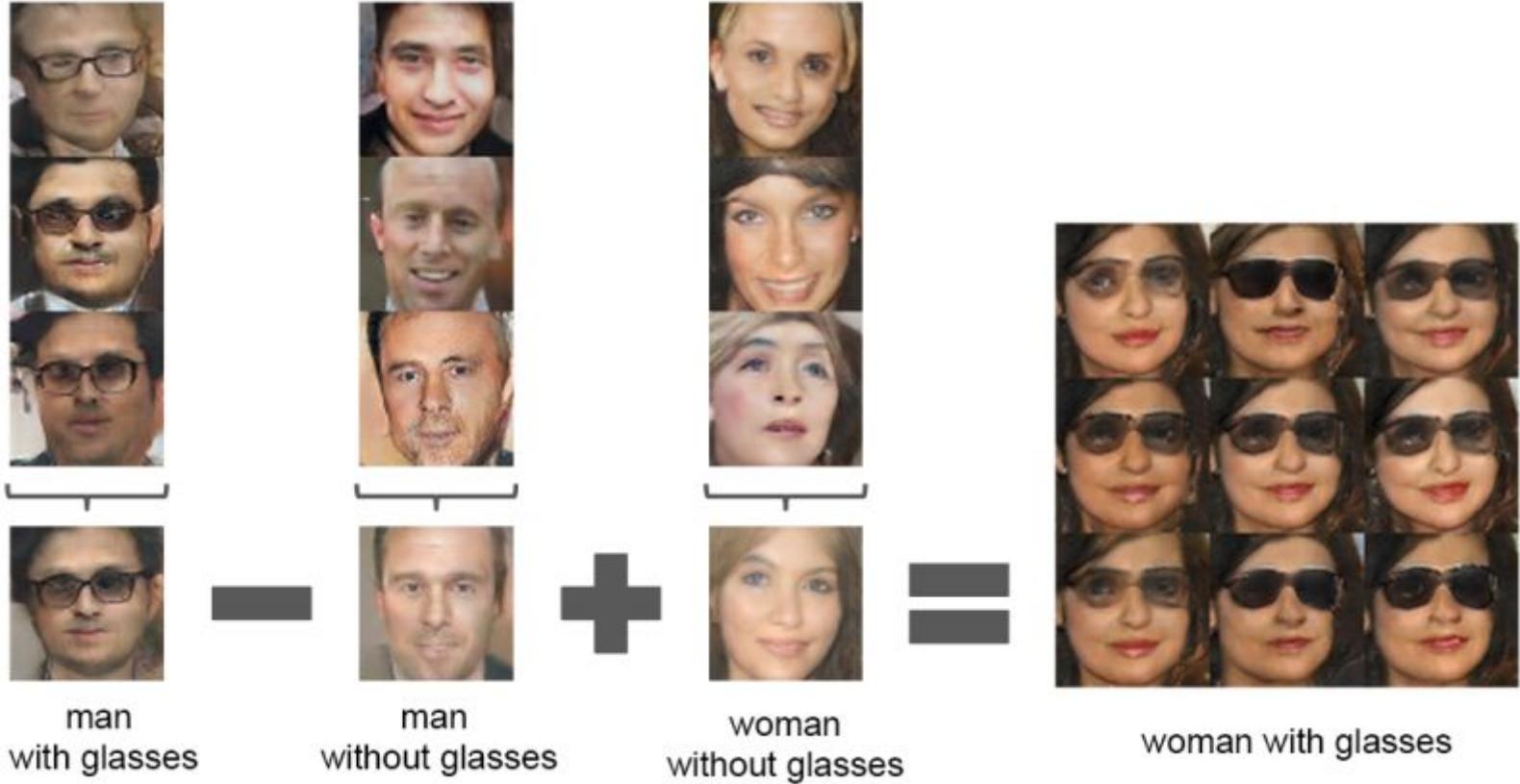
# Deep Convolution Generative Adversarial Net (DCGAN)

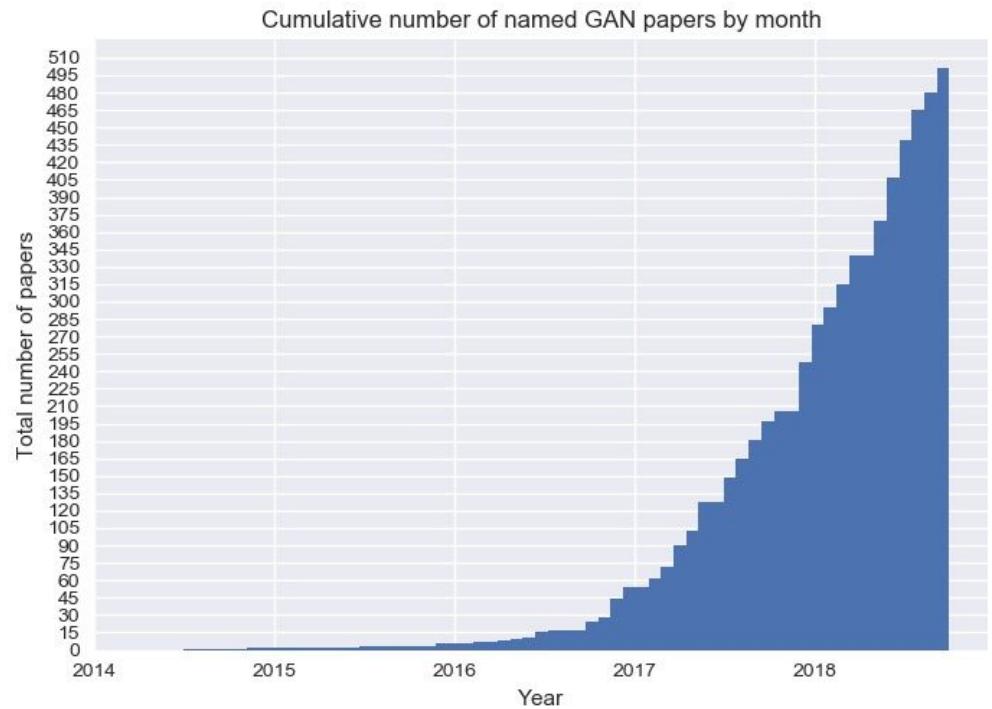


Interpolation  
entre 9 vecteurs  
latents aléatoires



# *“Vector arithmetic for visual concepts”*





<https://github.com/hindupuravinash/the-gan-zoo>

# Problèmes d'instabilité

- Si discriminateur et générateur et n'apprennent pas ensemble:
  - disparition des gradients
  - effondrement des modes
  - on ne peut générer d'images à haute résolution
- Plusieurs solutions proposées:
  - *Wasserstein GAN* (*utilise “earth mover distance”*)
  - *Least Squares GAN* (*utilise distance d'erreur quadratique*)
  - *Progressive GAN*
  - ....

# Problèmes d'instabilité

- Si discriminateur et générateur et n'apprennent pas ensemble:
  - **disparition des gradients**
  - effondrement des modes
  - on ne peut générer d'images à haute résolution

Si le discriminateur apprend trop vite, le générateur sera systématiquement battu, et n'apprendra rien

# Problèmes d'instabilité

- Si discriminateur et générateur et n'apprennent pas ensemble:
  - disparition des gradients
  - **effondrement des modes**
  - on ne peut générer d'images à haute résolution

Le générateur peut apprendre à générer tout le temps la même image qui bat le discriminateur

# Problèmes d'instabilité

- Si discriminateur et générateur et n'apprennent pas ensemble:

# Le génie de la disci

## Epoch 7

6	1	1	1	2	1	1	2	6	2
1	3	2	1	1	3	3	1	3	1
1	1	4	7	4	3	1	0	1	1
7	6	2	0	2	1	4	9	1	1
3	7	7	1	1	7	0	3	1	1
6	1	2	9	3	1	1	1	4	9
9	1	1	3	1	7	1	1	1	7
9	1	1	7	1	1	1	1	7	2
1	9	3	1	4	1	9	1	7	1
1	1	0	1	7	3	1	4	1	1

## Epoch 21

bat

<https://datascience.stackexchange.com/questions/29485/gan-discriminator-converging-to-one-output>

# Problèmes d'instabilité

- Si discriminateur et générateur et n'apprennent pas ensemble:
  - disparition des gradients
  - effondrement des modes
  - **on ne peut générer d'images à haute résolution**

Beaucoup d'artéfacts

# LS GAN

## Problème des GANs de base

“**sigmoïde** de sortie” **oublie** les exemples correctement classifiés et loin du plan de séparation

$$\max_D V(D) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

$$\min_G V(G) = \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))] \text{ or } \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} - [\log D(G(\mathbf{z}))]$$

# LS GAN

## Problème des GANs de base

“sigmoid”

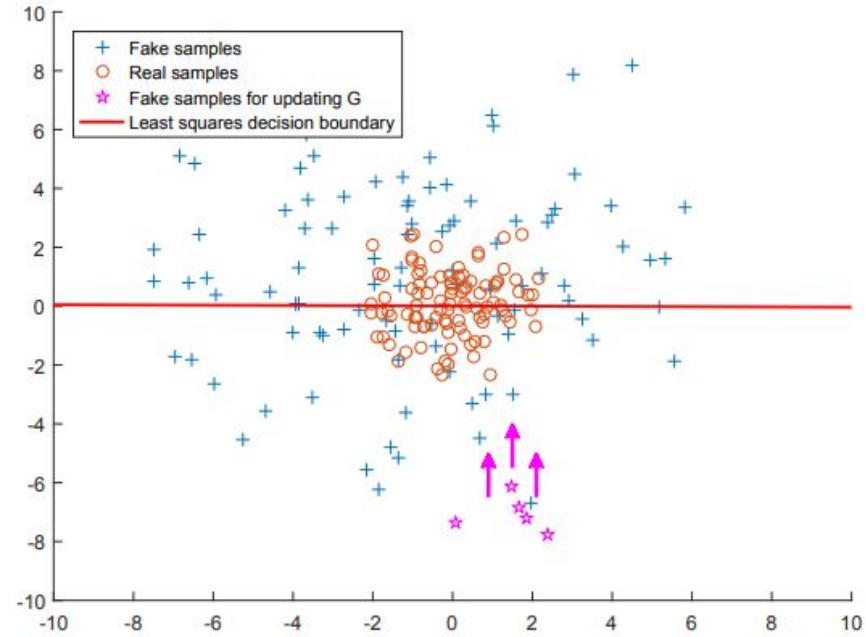
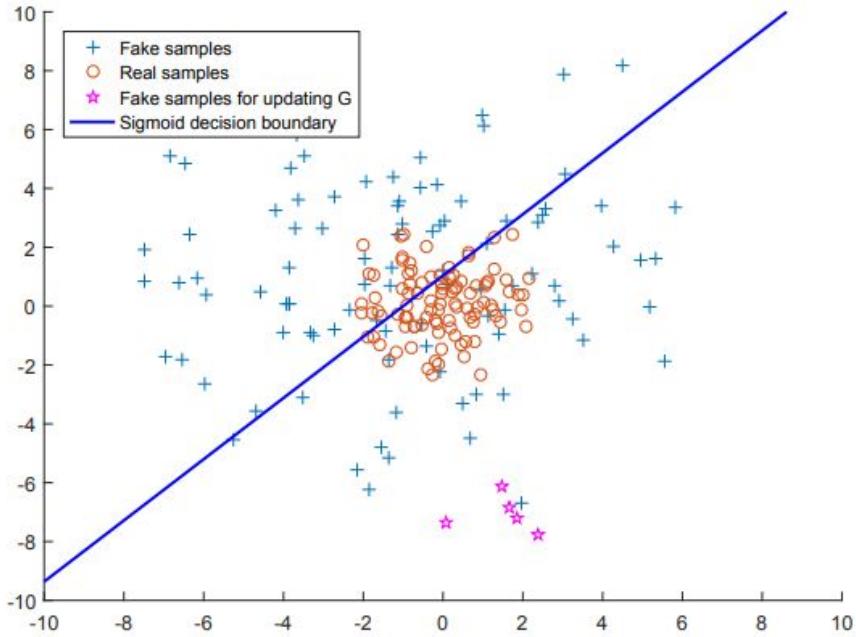
séparation

Le discriminateur ne s'entraîne plus lorsque les images synthétiques sont très différentes des images réelles

$$\min_G V(G) = \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))] \text{ or } \mathbb{E}_{z \sim p_z(z)} - [\log D(G(z))]$$

# LS GAN

## Problème des GANs de base



# LS GAN

Quand on utilise **l'erreur quadratique**, même les exemples « trop bien classifiés » contribuent aux gradients du générateur. Le but est de rapprocher les images synthétiques des images réelles

$$\max_D V(D) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

$$\min_G V(G) = \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))] \text{ or } \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} - [\log D(G(\mathbf{z}))]$$

**Pour LS GAN, la sortie du réseau n'est plus une sigmoïde**

$$\min_D V_{\text{LSGAN}}(D) = \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [(D(\mathbf{x}) - 1)^2] + \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [(D(G(\mathbf{z})))^2]$$

$$\min_G V_{\text{LSGAN}}(G) = \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [(D(G(\mathbf{z})))^2],$$

“Least Squares GAN” Mao et al. ICCV’17

# LS GAN



(a) Generated by LSGANs.

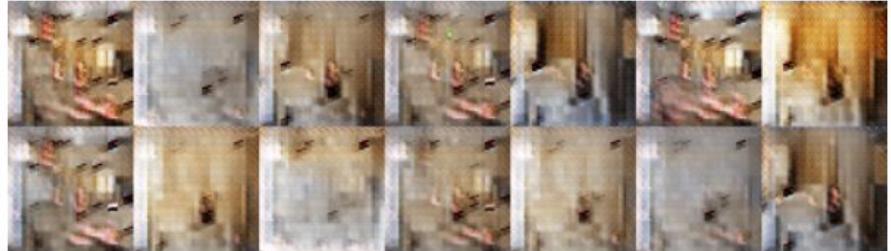


(b) Generated by DCGANs (Reported in [11]).

# LS GAN



(a) LSGANs.



(b) Regular GANs.



(c) LSGANs.



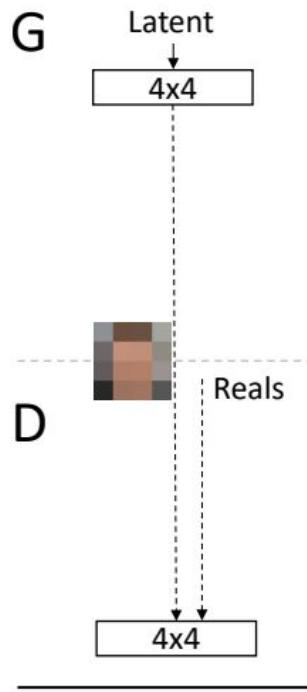
(d) Regular GANs.

“Least Squares GAN” Mao et al. ICCV’17

# progressive GAN

On veut générer des images à **haute résolution**

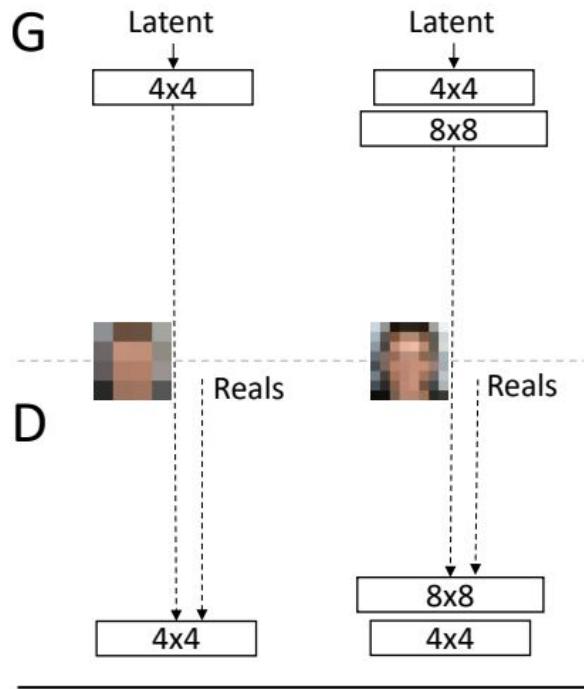
On commence avec des images de **faible résolution : 4x4 pixels**



*“Progressive GAN” Karras et al. ICLR’18*

# progressive GAN

Et progressivement, on augmente la résolution de l'image



*“Progressive GAN” Karras et al. ICLR’18*

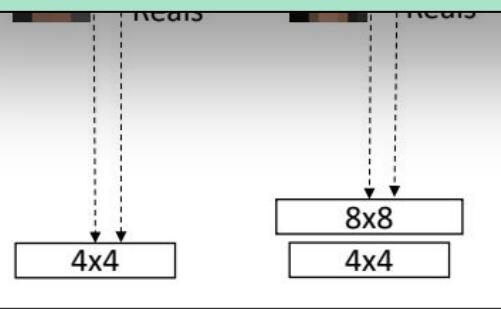
# progressive GAN

Et progressivement on augmente la résolution de l'image

G

*“Progressive Growing GAN requires that the capacity of both the generator and discriminator model be expanded by adding layers during the training process”*

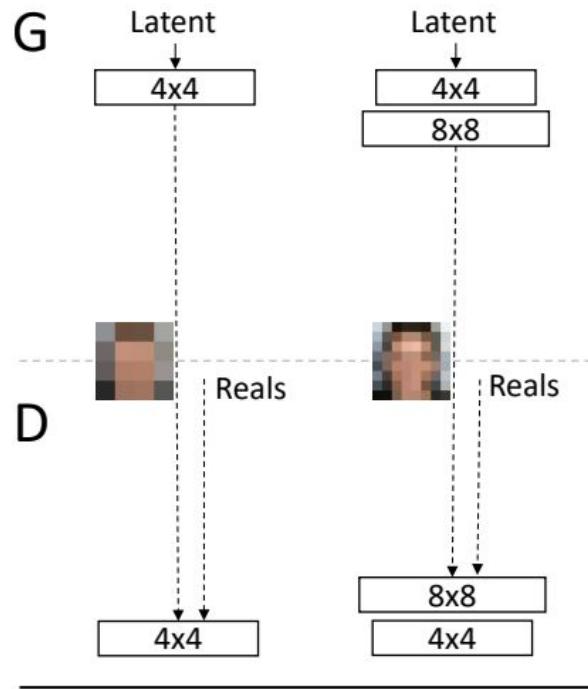
D



*“Progressive GAN” Karras et al. ICLR’18*

# progressive GAN

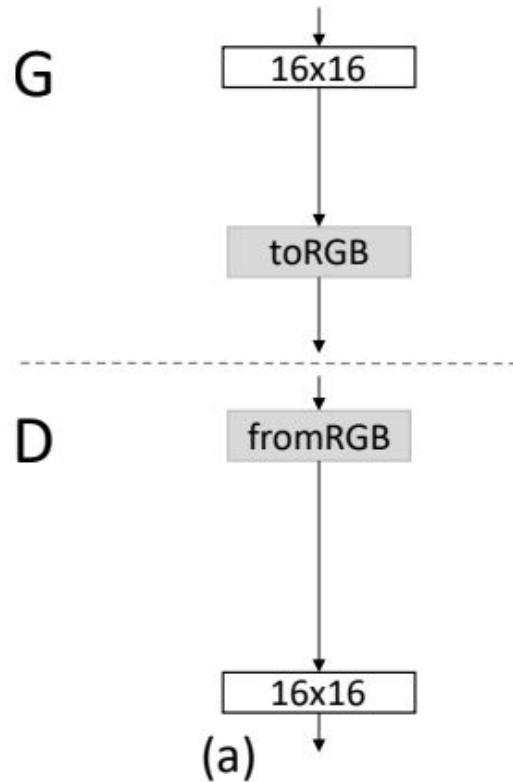
Et progressivement, chaque couche qu'on ajoute vient bonifier la couche précédente : cela se fait à l'aide d'une **opération « résiduelle »**.



“Progressive GAN” Karras et al. ICLR’18

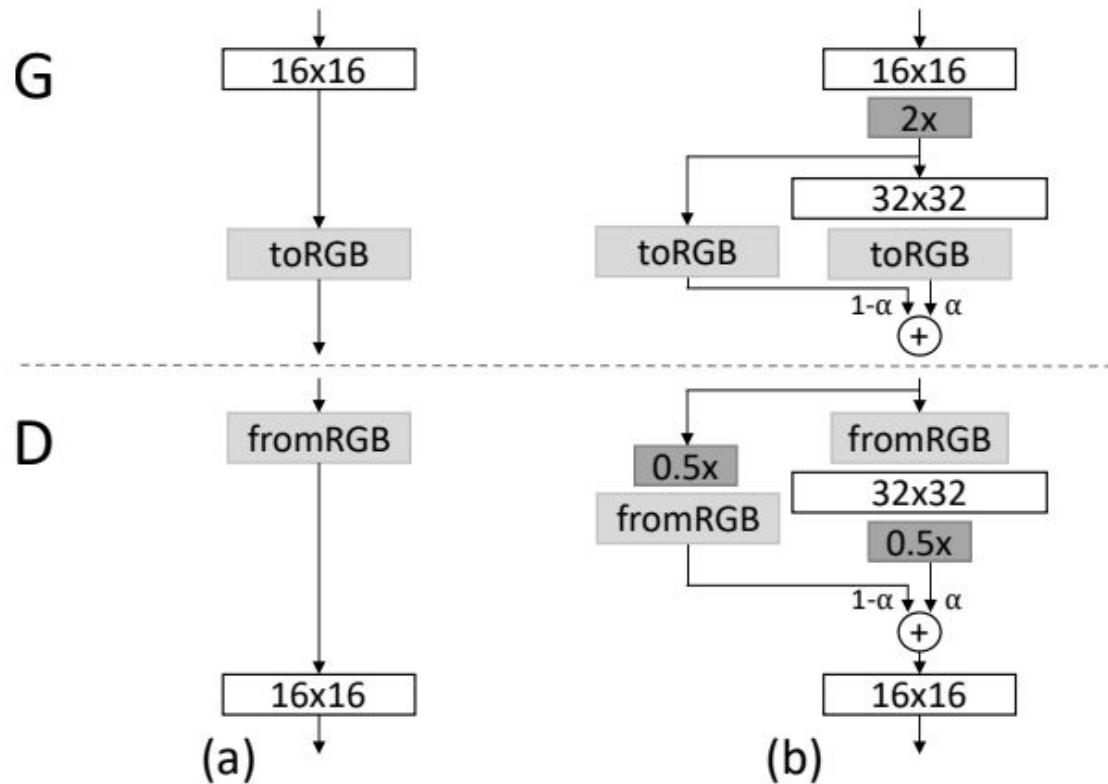
# Ajout de couches

Lorsque **l'entraînement** d'une couche de résolution RxR (ici 16x16) est **terminé**...



# Ajout de couches

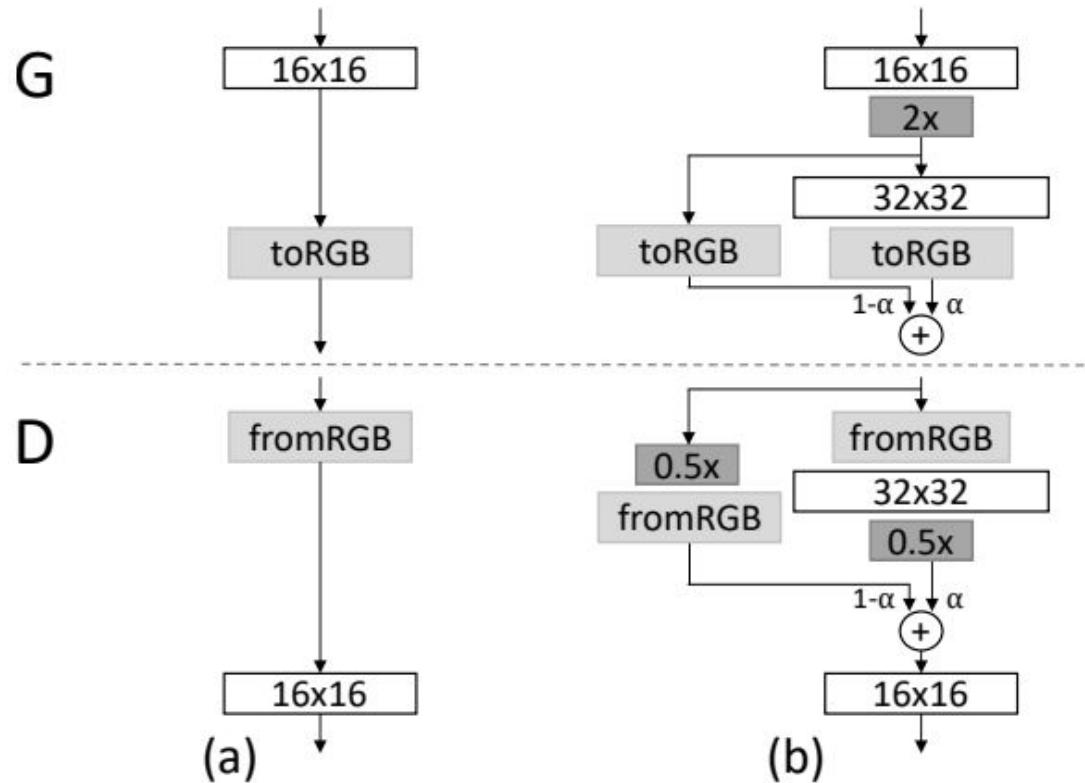
... On **ajoute une nouvelle couche** de résolution 2Rx2D (ici 32x32) au générateur ET au discriminateur.



# Ajout de couches

... mais pour éviter un choc, on ajoute une **composante résiduelle** comprenant un facteur  $\alpha$

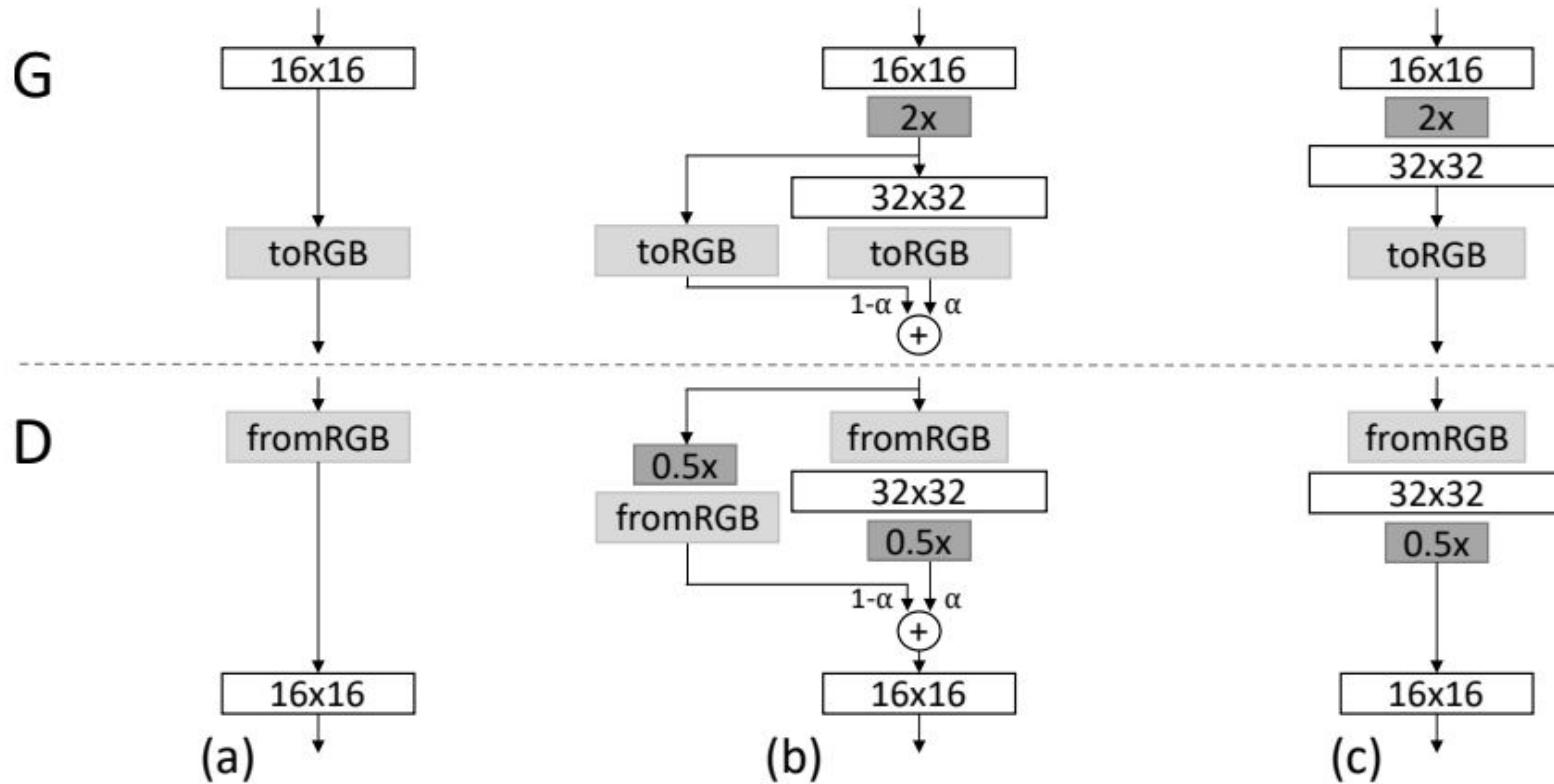
**Au début de l'entraînement,  $\alpha=0$  et progressivement  $\alpha$  augmente pour atteindre  $\alpha=1$  à la fin**



“Progressive GAN” Karras et al. ICLR’18

# Ajout de couches

Lorsque l'entraînement est terminé, on enlève la composante résiduelle.



# “Progressive GAN” Karras et al. ICLR’18

<b>Generator</b>	Act.	Output shape	Params		<b>Discriminator</b>	Act.	Output shape	Params
Latent vector	–	512 × 1 × 1	–		Input image	–	3 × 1024 × 1024	–
Conv 4 × 4	LReLU	512 × 4 × 4	4.2M		Conv 1 × 1	LReLU	16 × 1024 × 1024	64
Conv 3 × 3	LReLU	512 × 4 × 4	2.4M		Conv 3 × 3	LReLU	16 × 1024 × 1024	2.3k
Upsample	–	512 × 8 × 8	–		Conv 3 × 3	LReLU	32 × 1024 × 1024	4.6k
Conv 3 × 3	LReLU	512 × 8 × 8	2.4M		Downsample	–	32 × 512 × 512	–
Conv 3 × 3	LReLU	512 × 8 × 8	2.4M		Conv 3 × 3	LReLU	32 × 512 × 512	9.2k
Upsample	–	512 × 16 × 16	–		Conv 3 × 3	LReLU	64 × 512 × 512	18k
Conv 3 × 3	LReLU	512 × 16 × 16	2.4M		Downsample	–	64 × 256 × 256	–
Conv 3 × 3	LReLU	512 × 16 × 16	2.4M		Conv 3 × 3	LReLU	64 × 256 × 256	37k
Upsample	–	512 × 32 × 32	–		Conv 3 × 3	LReLU	128 × 256 × 256	74k
Conv 3 × 3	LReLU	512 × 32 × 32	2.4M		Downsample	–	128 × 128 × 128	–
Conv 3 × 3	LReLU	512 × 32 × 32	2.4M		Conv 3 × 3	LReLU	128 × 128 × 128	148k
Upsample	–	512 × 64 × 64	–		Conv 3 × 3	LReLU	256 × 128 × 128	295k
Conv 3 × 3	LReLU	256 × 64 × 64	1.2M		Downsample	–	256 × 64 × 64	–
Conv 3 × 3	LReLU	256 × 64 × 64	590k		Conv 3 × 3	LReLU	256 × 64 × 64	590k
Upsample	–	256 × 128 × 128	–		Conv 3 × 3	LReLU	512 × 64 × 64	1.2M
Conv 3 × 3	LReLU	128 × 128 × 128	295k		Downsample	–	512 × 32 × 32	–
Conv 3 × 3	LReLU	128 × 128 × 128	148k		Conv 3 × 3	LReLU	512 × 32 × 32	2.4M
Upsample	–	128 × 256 × 256	–		Conv 3 × 3	LReLU	512 × 32 × 32	2.4M
Conv 3 × 3	LReLU	64 × 256 × 256	74k		Downsample	–	512 × 16 × 16	–
Conv 3 × 3	LReLU	64 × 256 × 256	37k		Conv 3 × 3	LReLU	512 × 16 × 16	2.4M
Upsample	–	64 × 512 × 512	–		Conv 3 × 3	LReLU	512 × 16 × 16	2.4M
Conv 3 × 3	LReLU	32 × 512 × 512	18k		Downsample	–	512 × 8 × 8	–
Conv 3 × 3	LReLU	32 × 512 × 512	9.2k		Conv 3 × 3	LReLU	512 × 8 × 8	2.4M
Upsample	–	32 × 1024 × 1024	–		Conv 3 × 3	LReLU	512 × 8 × 8	2.4M
Conv 3 × 3	LReLU	16 × 1024 × 1024	4.6k		Downsample	–	512 × 4 × 4	–
Conv 3 × 3	LReLU	16 × 1024 × 1024	2.3k		Minibatch stddev	–	513 × 4 × 4	–
Conv 1 × 1	linear	3 × 1024 × 1024	51		Conv 3 × 3	LReLU	512 × 4 × 4	2.4M
Total trainable parameters			<b>23.1M</b>		Conv 4 × 4	LReLU	512 × 1 × 1	4.2M
					Fully-connected	linear	1 × 1 × 1	513
					Total trainable parameters			<b>23.1M</b>

Table 2: Generator and discriminator that we use with CELEBA-HQ to generate  $1024 \times 1024$  images.

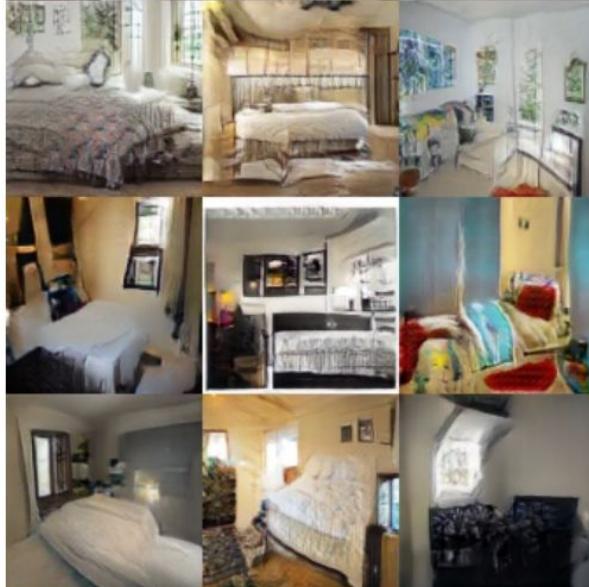
*“Progressive GAN” Karras et al. ICLR’18*



<https://youtu.be/XOxxPcy5Gr4>

201

# *“Progressive GAN” Karras et al. ICLR’18*



Mao et al. (2016b) ( $128 \times 128$ )

LSGAN



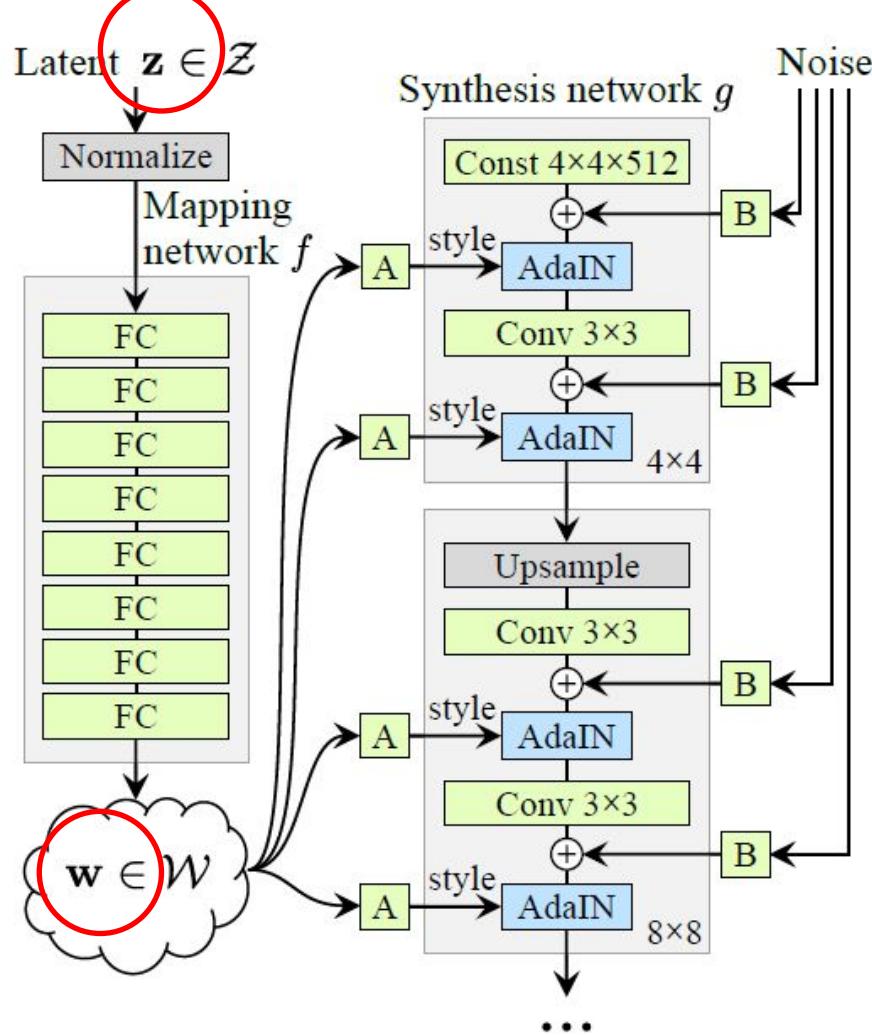
Gulrajani et al. (2017) ( $128 \times 128$ )

PGAN



Our ( $256 \times 256$ )

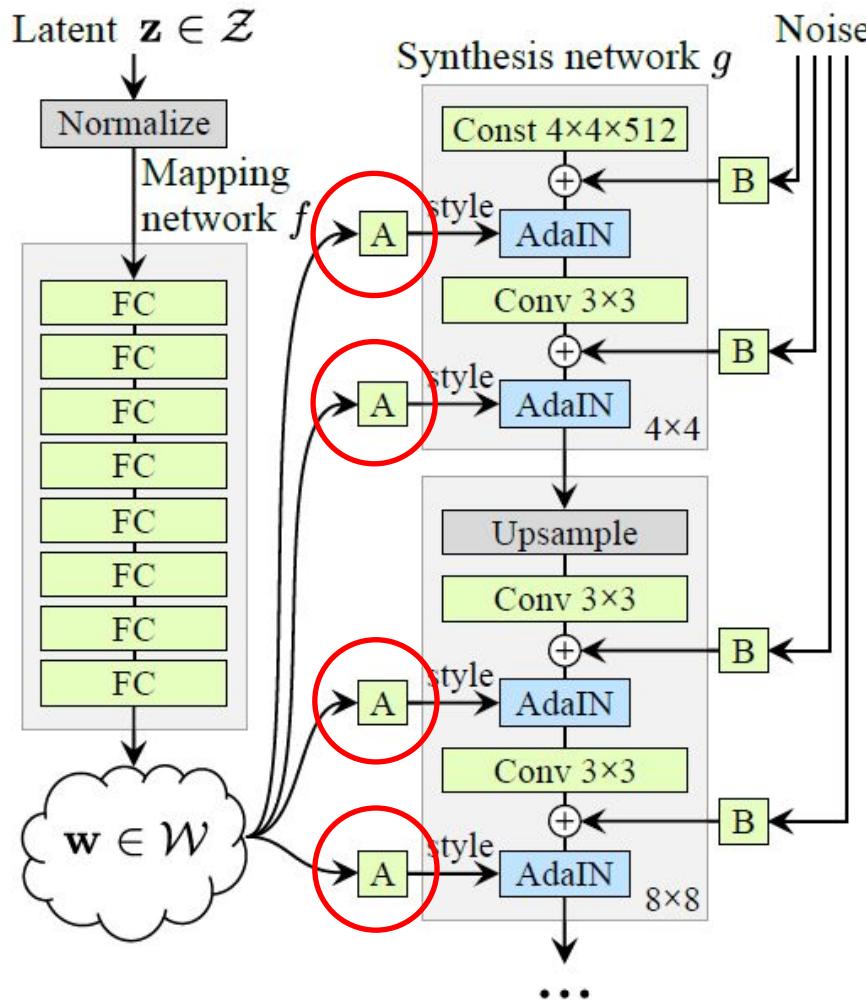
# Style GAN



(b) Style-based generator

Le Générateur reçoit une version modifiée “ $w$ ” par 8 couches FC du vecteur latent “ $z$ ”

# Style GAN

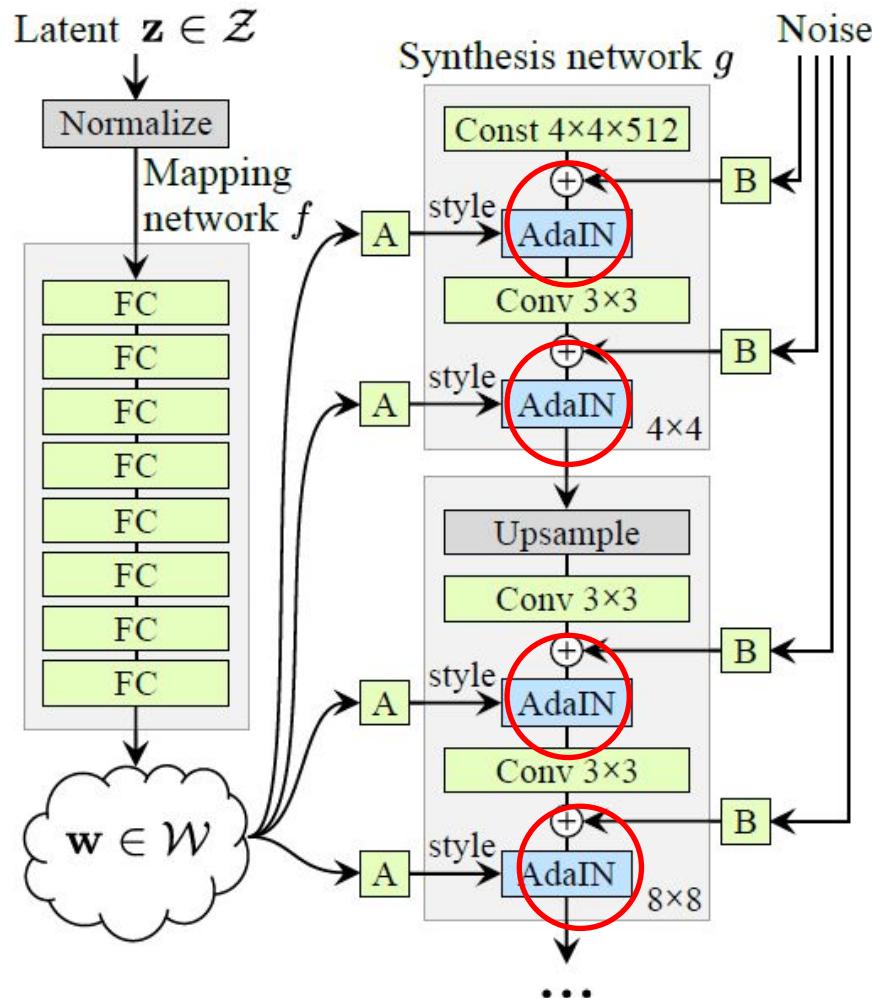


(b) Style-based generator

Le Générateur reçoit une version modifiée “ $w$ ” par 8 couches FC du vecteur latent “ $z$ ”

Et ce, à **chaque couche** du réseau

# Style GAN



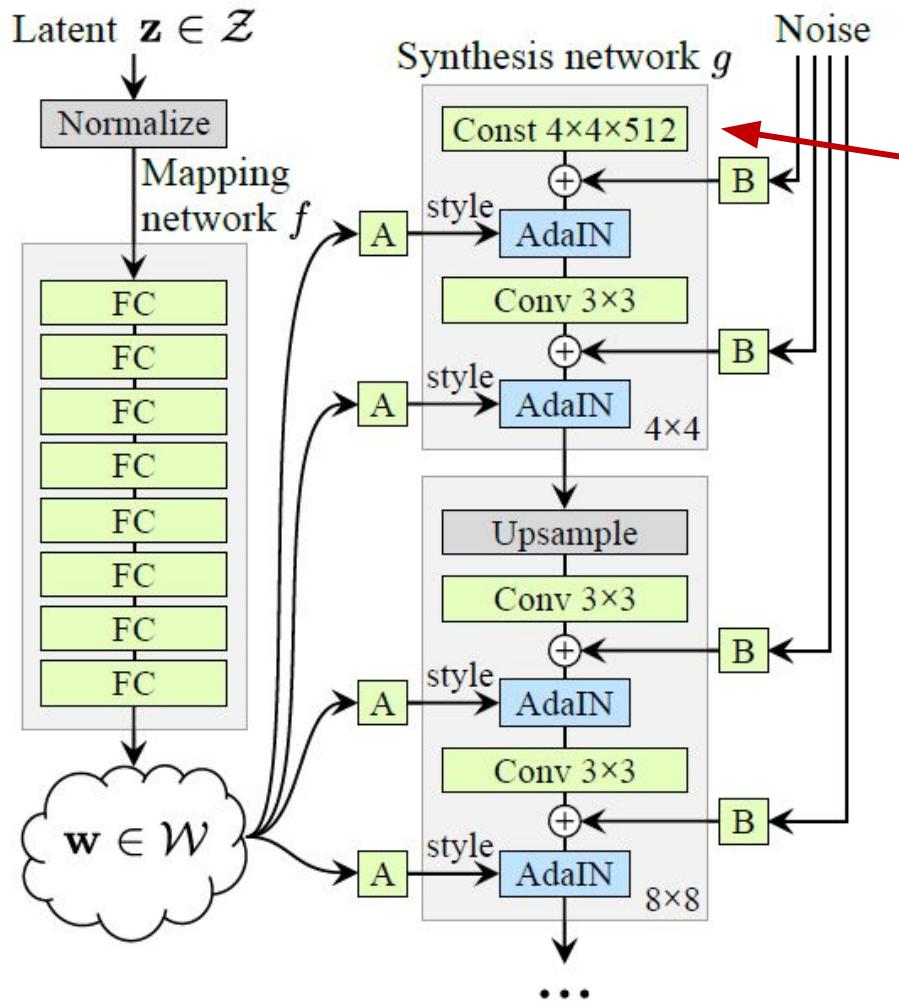
(b) Style-based generator

AdaIN: adaptive instance normalization

$$AdaIN(x, y) = \sigma(y) \frac{x - \mu(x)}{\sigma(x)} + \mu(y)$$

Comme du batchNorm, mais dont les 2 opérateurs affines sont fournis par  $\mathbf{w}$

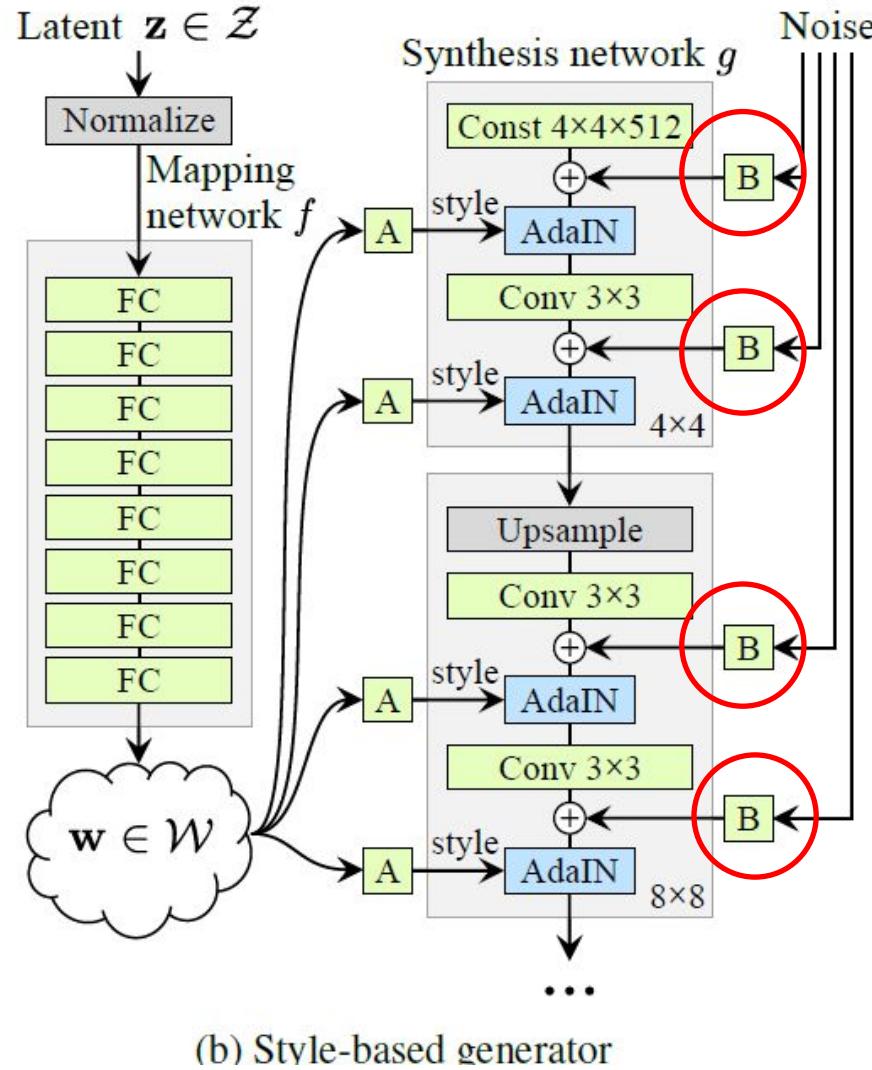
# Style GAN



(b) Style-based generator

L'entrée du réseau est un **tenseur constant** appris par entraînement

# Style GAN



Du **bruit** est multiplié à chaque carte d'activation de **chaque couche** du réseau

# Effet du bruit



(a) Generated image

(b) Stochastic variation

(c) Standard deviation

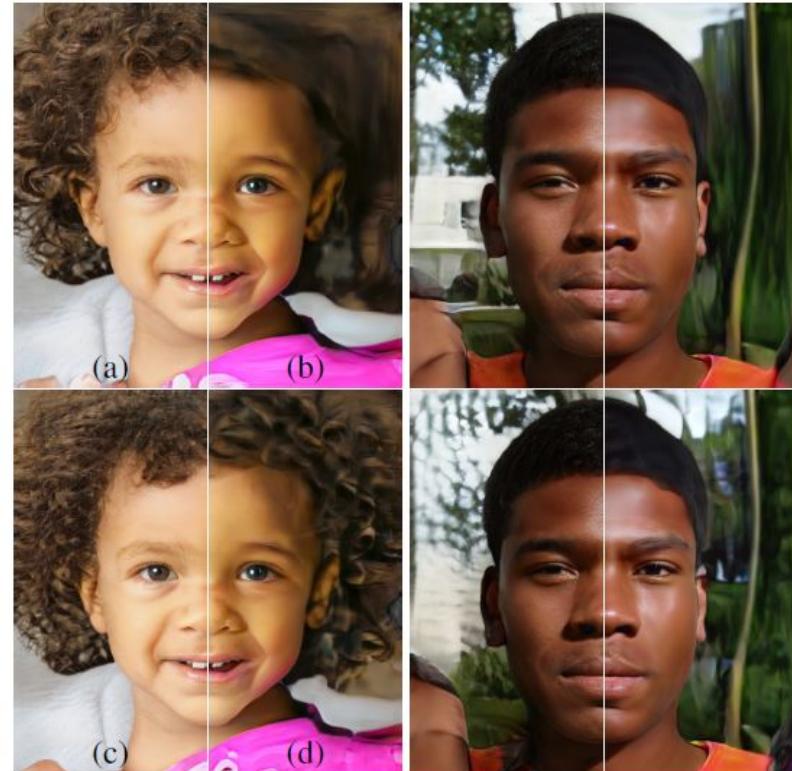
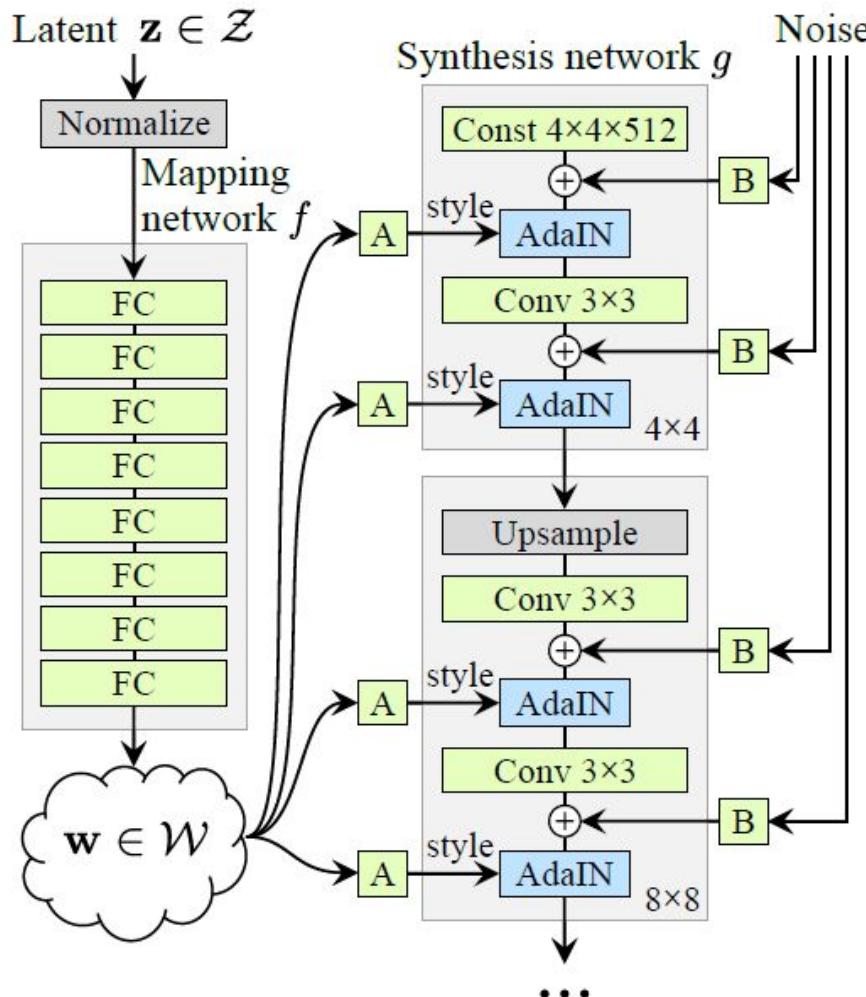


Figure 5. Effect of noise inputs at different layers of our generator. (a) Noise is applied to all layers. (b) No noise. (c) Noise in fine layers only ( $64^2 - 1024^2$ ). (d) Noise in coarse layers only ( $4^2 - 32^2$ ). We can see that the artificial omission of noise leads to featureless “painterly” look. Coarse noise causes large-scale curling of hair and appearance of larger background features, while the fine noise brings out the finer curls of hair, finer background detail, and skin pores.

# Style GAN



(b) Style-based generator

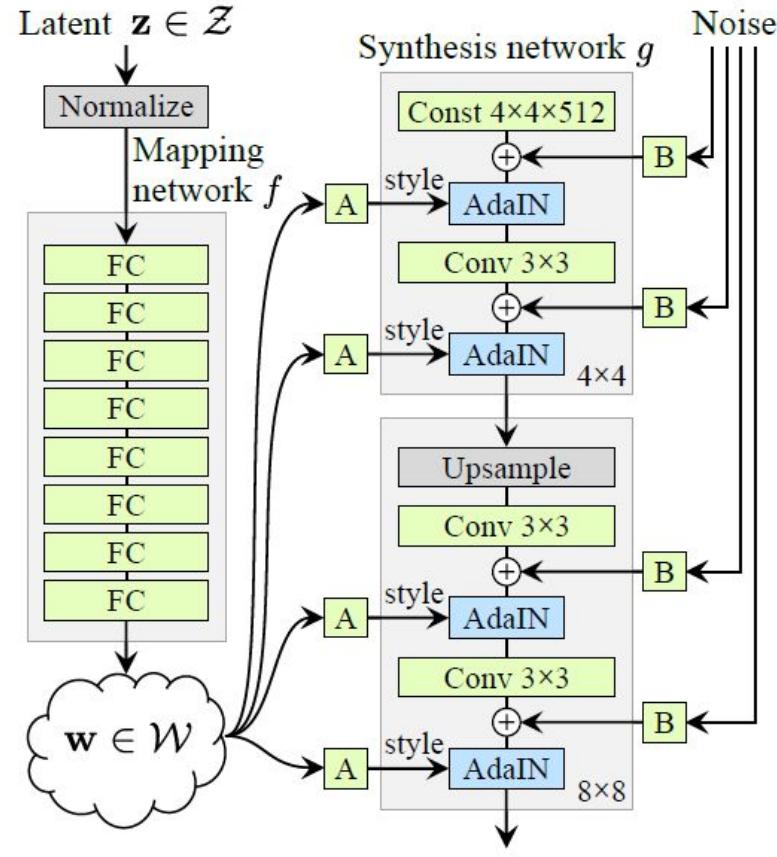
Du **bruit** est multiplié à chaque carte d'activation de **chaque couche** du réseau

Entraînement progressif comme pour **progressive GAN**

# Style GAN



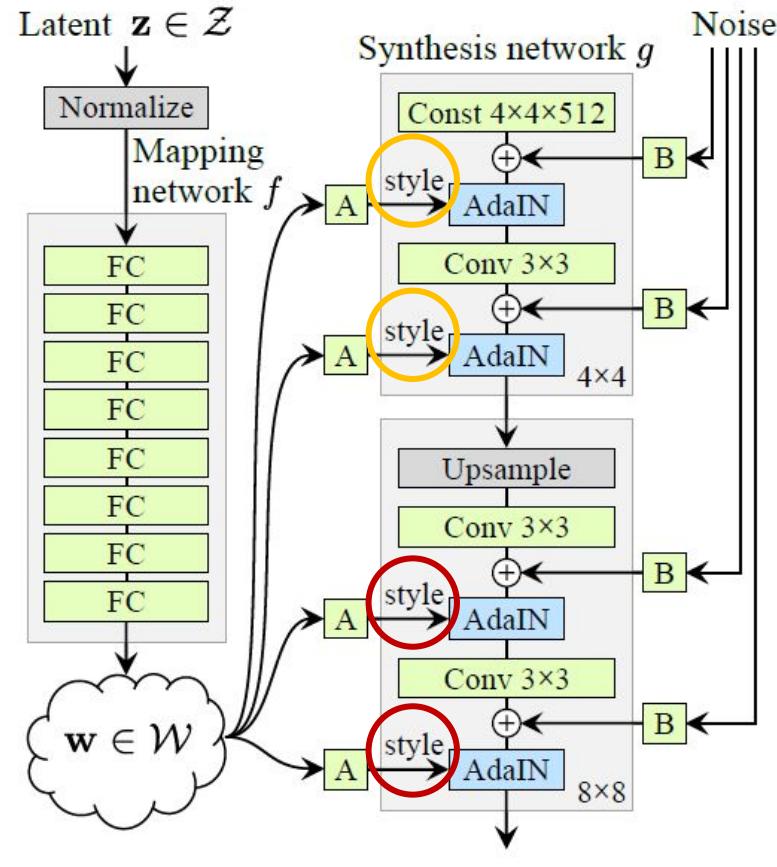
# Style GAN



(b) Style-based generator

Qu'est-ce qui arrive si  
on fait varier le  $\mathbf{z}$  entre  
les couches ?

# Style GAN



(b) Style-based generator

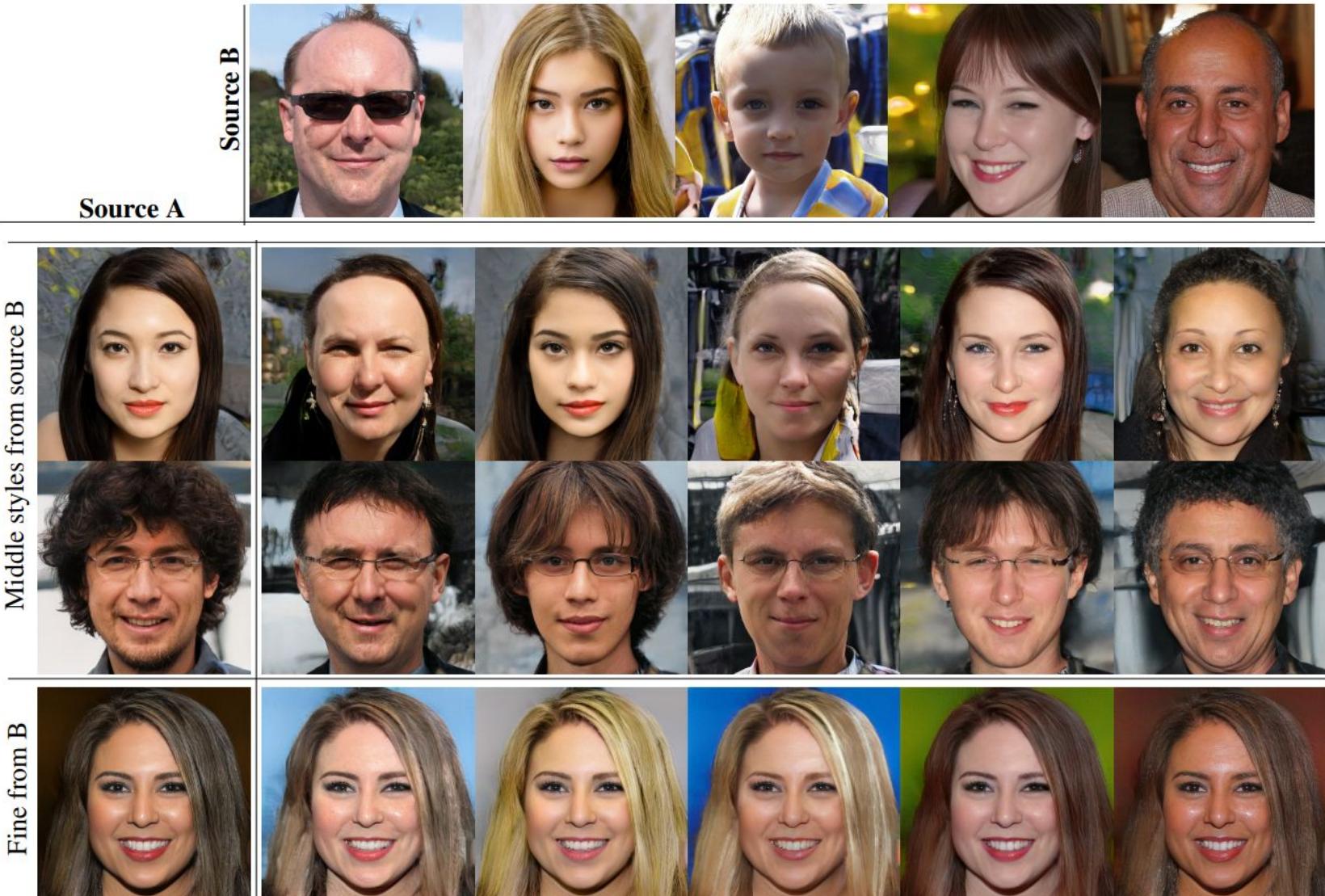
Qu'est-ce qui arrive si  
on fait varier le  $\mathbf{z}$  entre  
les couches ?

# Style GAN

Coarse styles from source B

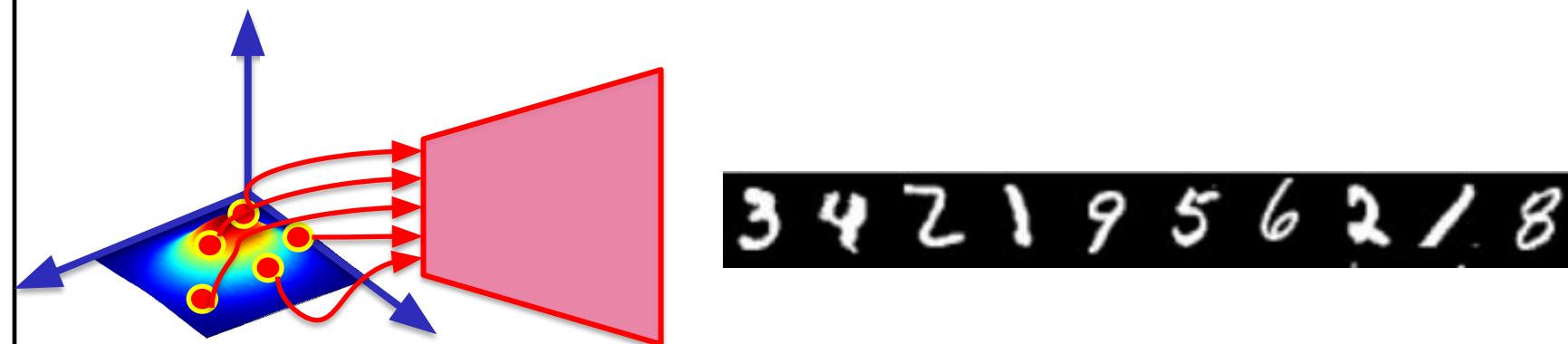


# Style GAN



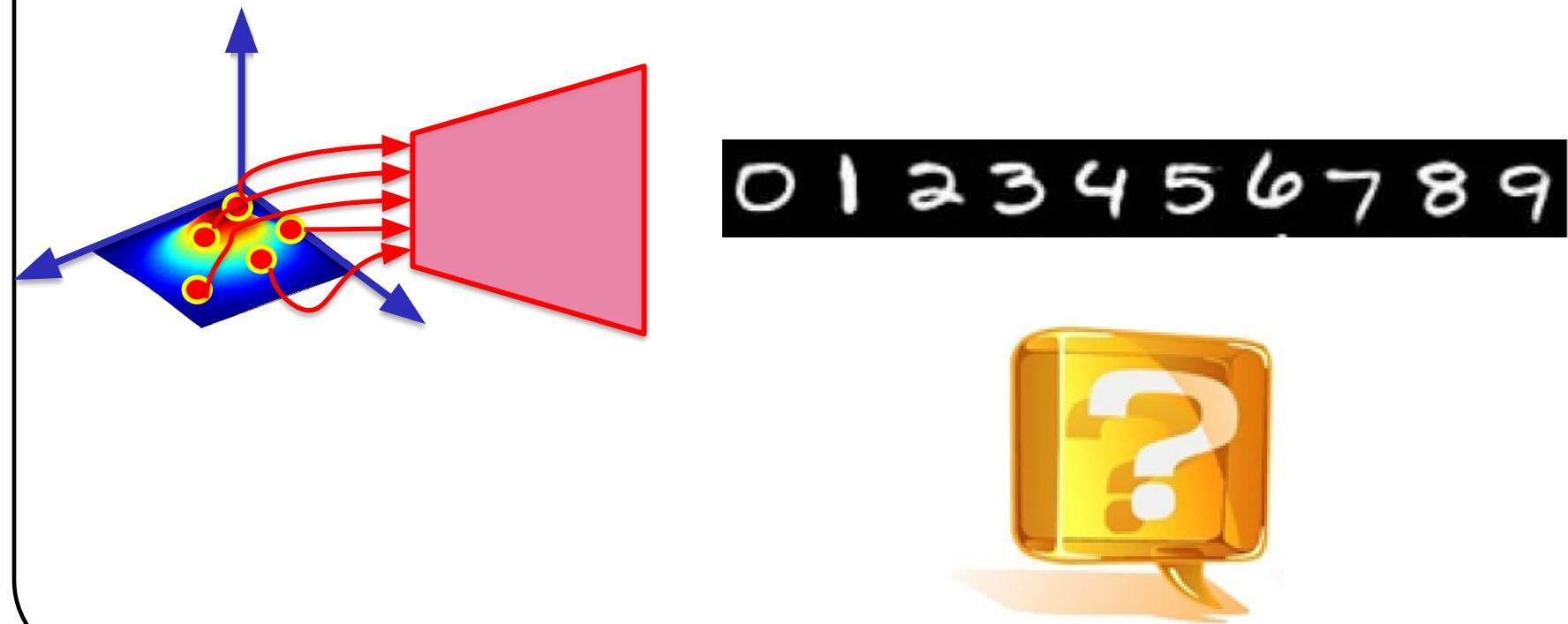
# Défi avec les GAN

Soit un GAN entraîné sur MNIST, si je décode **10 vecteurs latents pris au hasard**, j'aurai les images de **10 caractères aléatoires**.



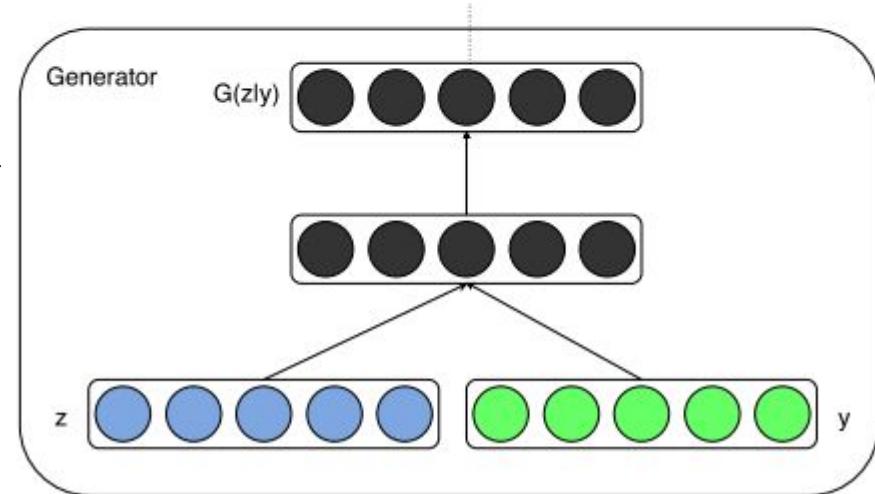
# Défi avec les GAN

**Question:** comment générer des images de catégories prédéterminées? Ex. comment sélectionner 10 vecteurs latent afin de produire la séquence de caractères : 0,1,2,3,4,5,6,7,8,9?

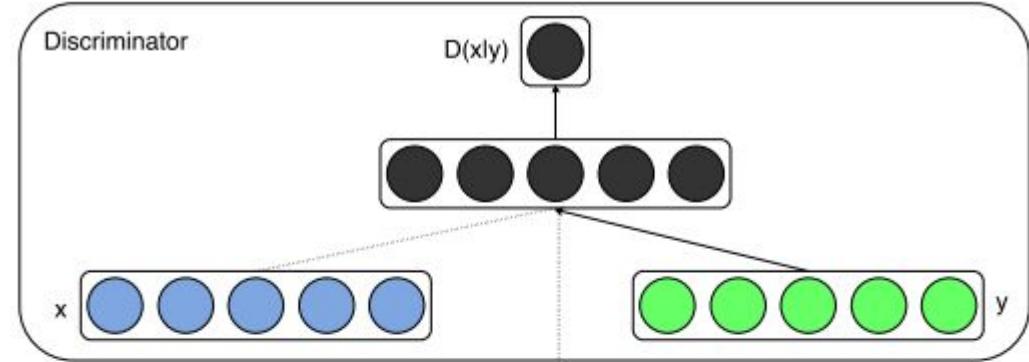


# Gan conditionnel

L'idée est d'encoder un vecteur latent  $\vec{z}$  ainsi qu'un **vecteur de classe** « one-hot »  $\vec{y}$

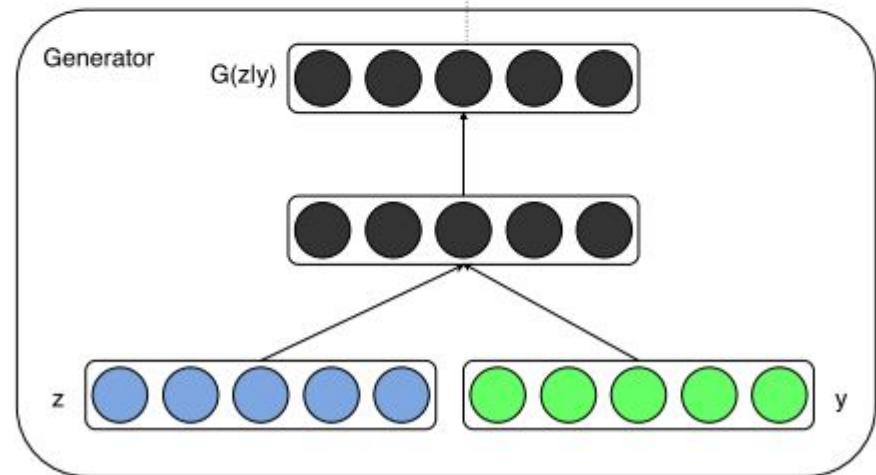


# Gan conditionnel

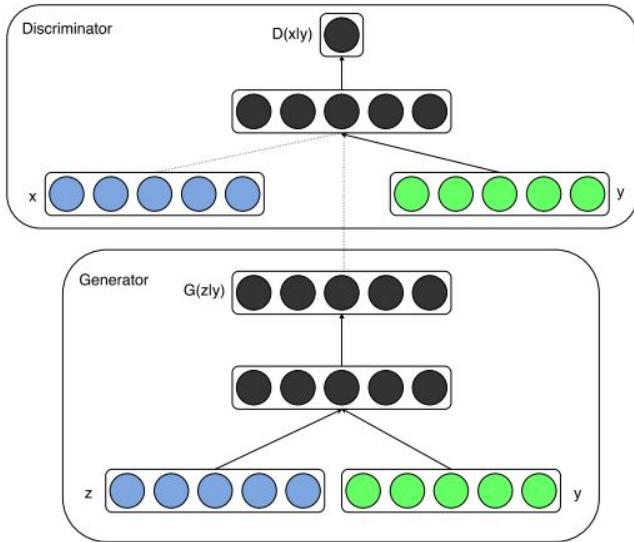


Et de discriminer une image  $\vec{x}$

avec **le même « one-hot »**  $\vec{y}$



# Gan conditionnel



## GAN de base

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))].$$

## GAN conditionnel

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x}|\mathbf{y})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z}|\mathbf{y})))].$$

Mirza, Mehdi & Osindero, Simon. (2014). Conditional Generative Adversarial Nets. arXiv:1411.1784v1

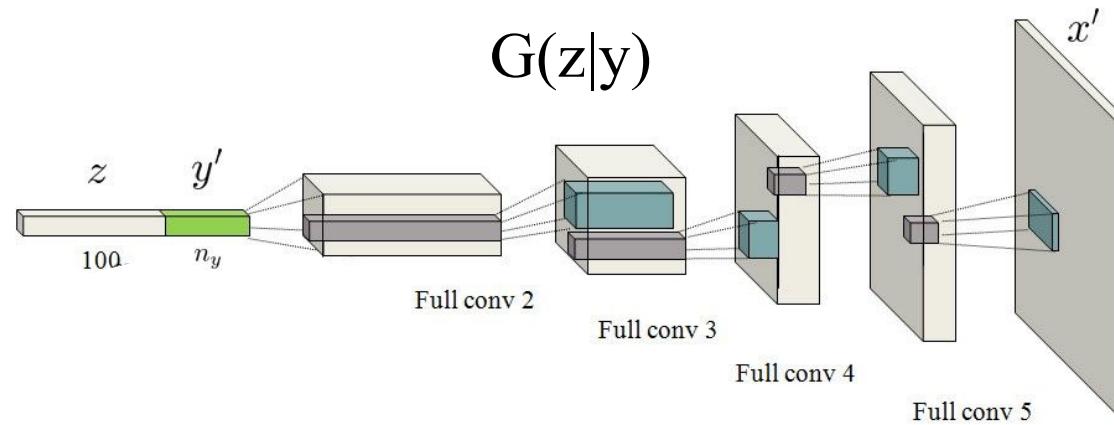
## “Fashion MNIST”



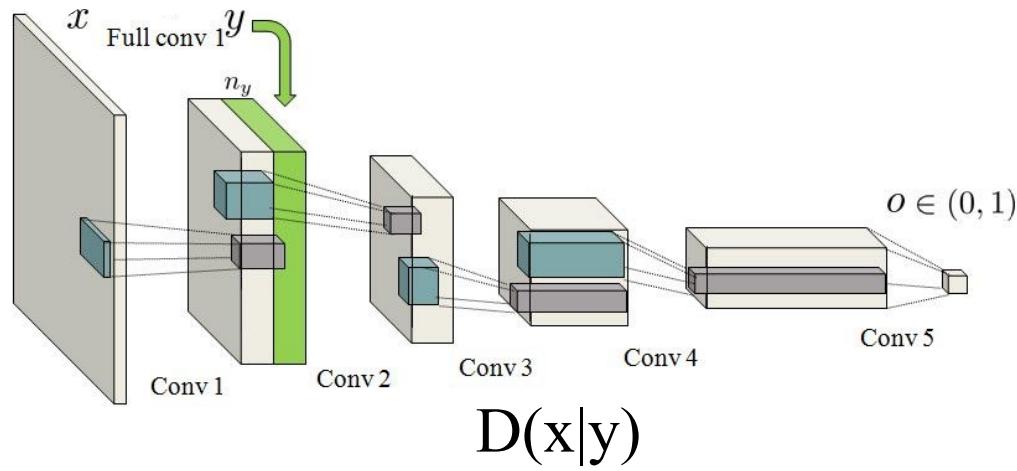
“t-shirt”, ‘pants’, ‘pullover’, ‘dress’, ‘coat’, ‘sandal’, ‘shirt’, ‘sneaker’, ‘bag’, ‘ankle boot’.

Mirza, Mehdi & Osindero, Simon. (2014). Conditional Generative Adversarial Nets. arXiv:1411.1784v1

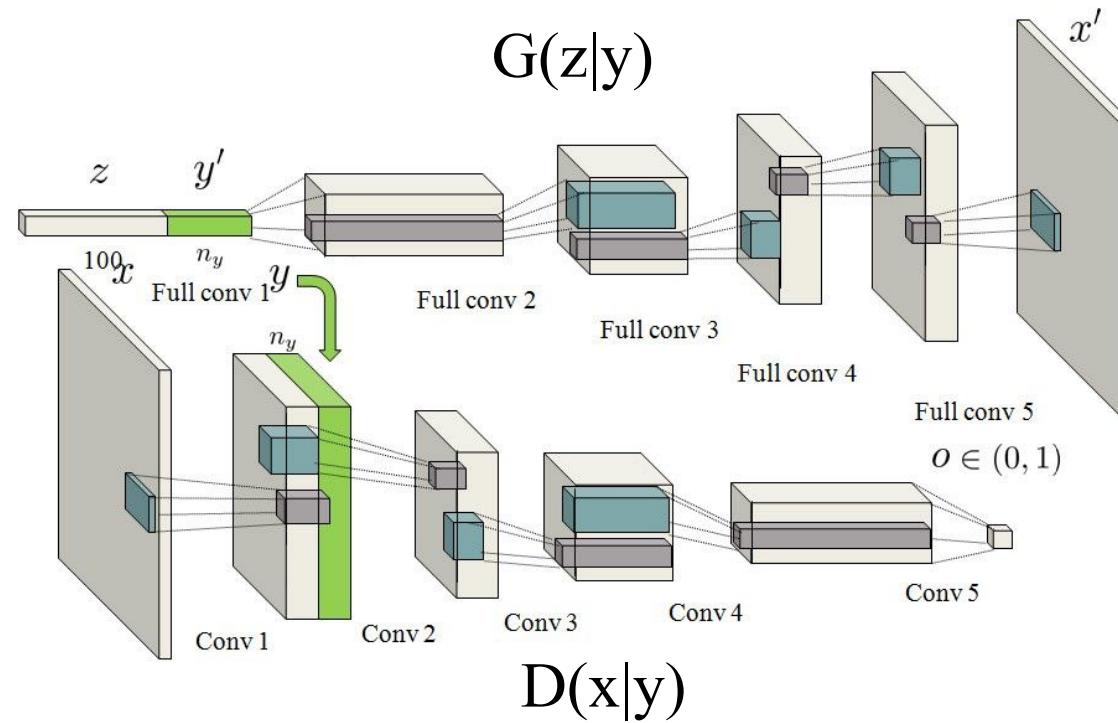
# Version convolutionnelle



# Version convolutionnelle



# Version convolutionnelle



### Gray whale



### Welsh springer spaniel



### Persian cat



### Tiger



### Chiffonier



### Fire truck



Miyato, T., Kataoka, T., Koyama, M., & Yoshida, Y. (2018). Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*.

# Self-attention GAN

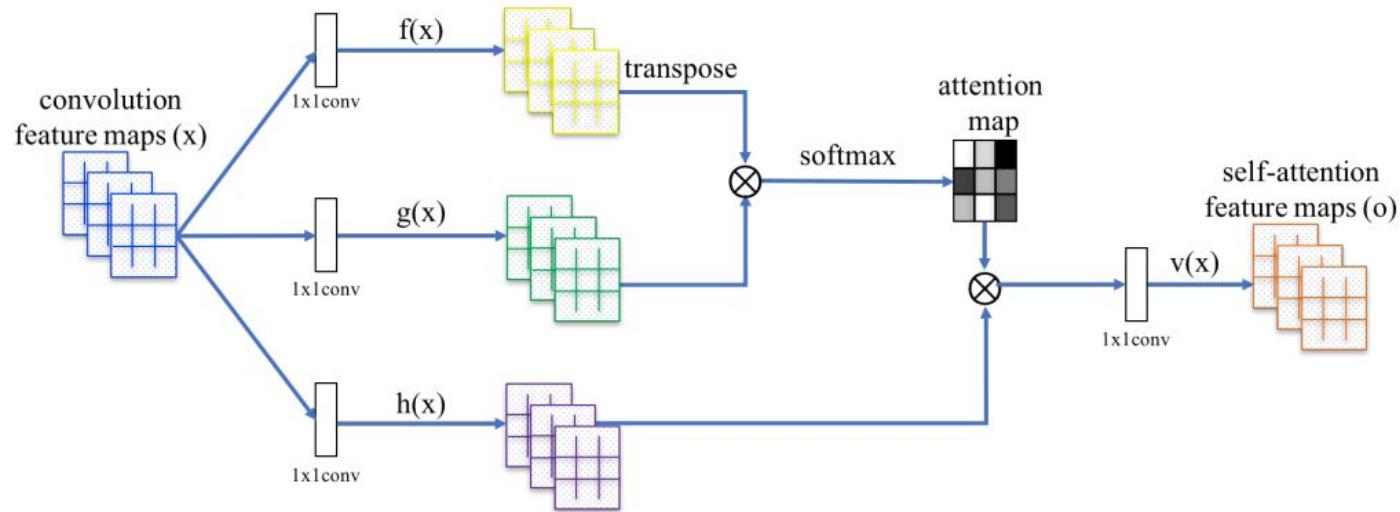
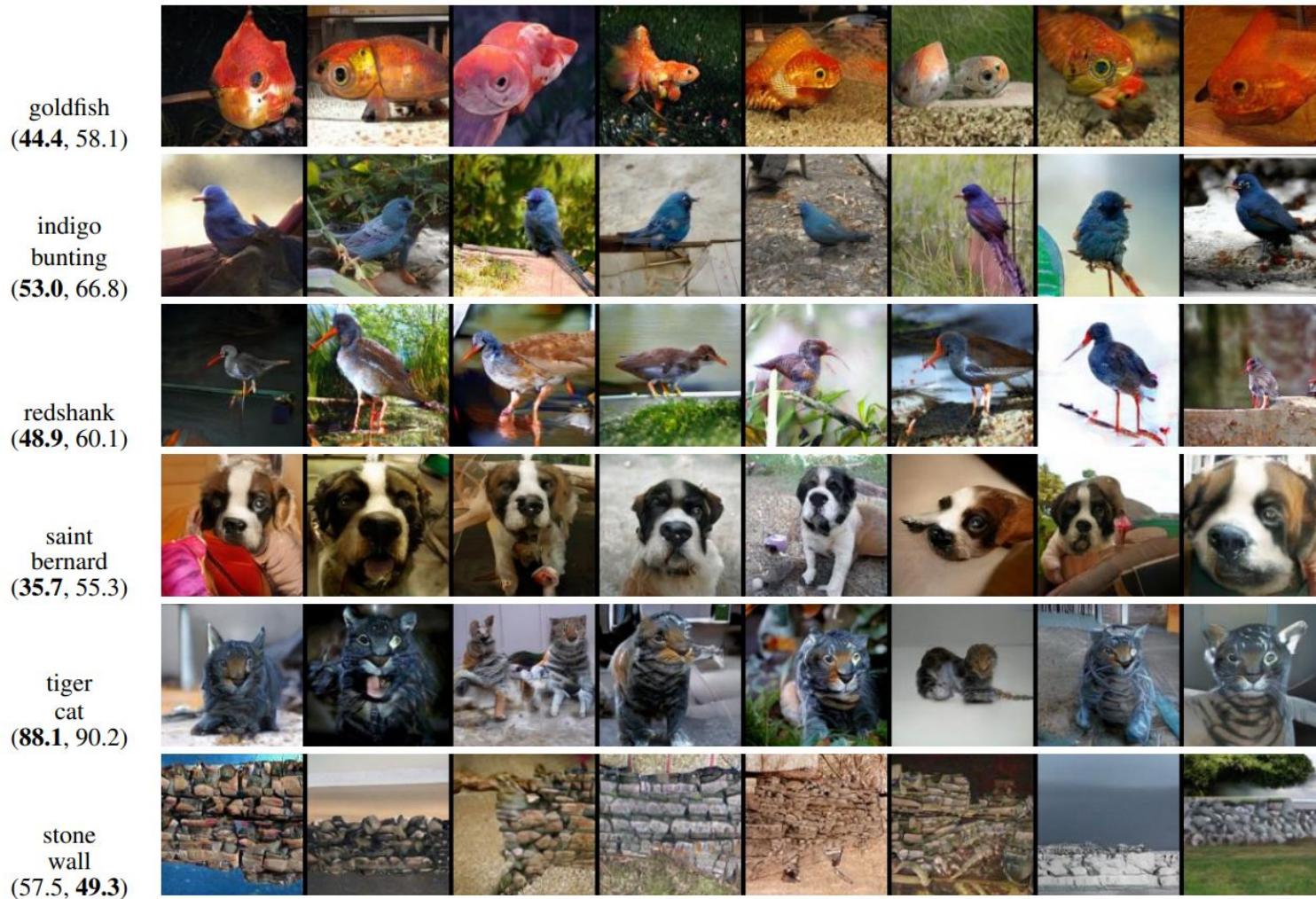


Figure 2. The proposed self-attention module for the SAGAN. The  $\otimes$  denotes matrix multiplication. The softmax operation is performed on each row.

GAN standard + *intra*-attention + d'autres détails moins intéressants

Zhang, H., Goodfellow, I., Metaxas, D., & Odena, A. (2019, May). Self-attention generative adversarial networks. In *International conference on machine learning* (pp. 7354-7363). PMLR.

# Self-attention GAN



Zhang, H., Goodfellow, I., Metaxas, D., & Odena, A. (2019, May). Self-attention generative adversarial networks. In *International conference on machine learning* (pp. 7354-7363). PMLR.

# BigGAN

Un plus gros self-attention GAN



Figure 1: Class-conditional samples generated by our model.

Brock, A., Donahue, J., & Simonyan, K. (2018). Large scale GAN training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*.

# BigGAN

Un plus gros self-attention GAN



Figure 6: Samples generated by our BigGAN model at  $512 \times 512$  resolution.

Brock, A., Donahue, J., & Simonyan, K. (2018). Large scale GAN training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*.

# BigGAN

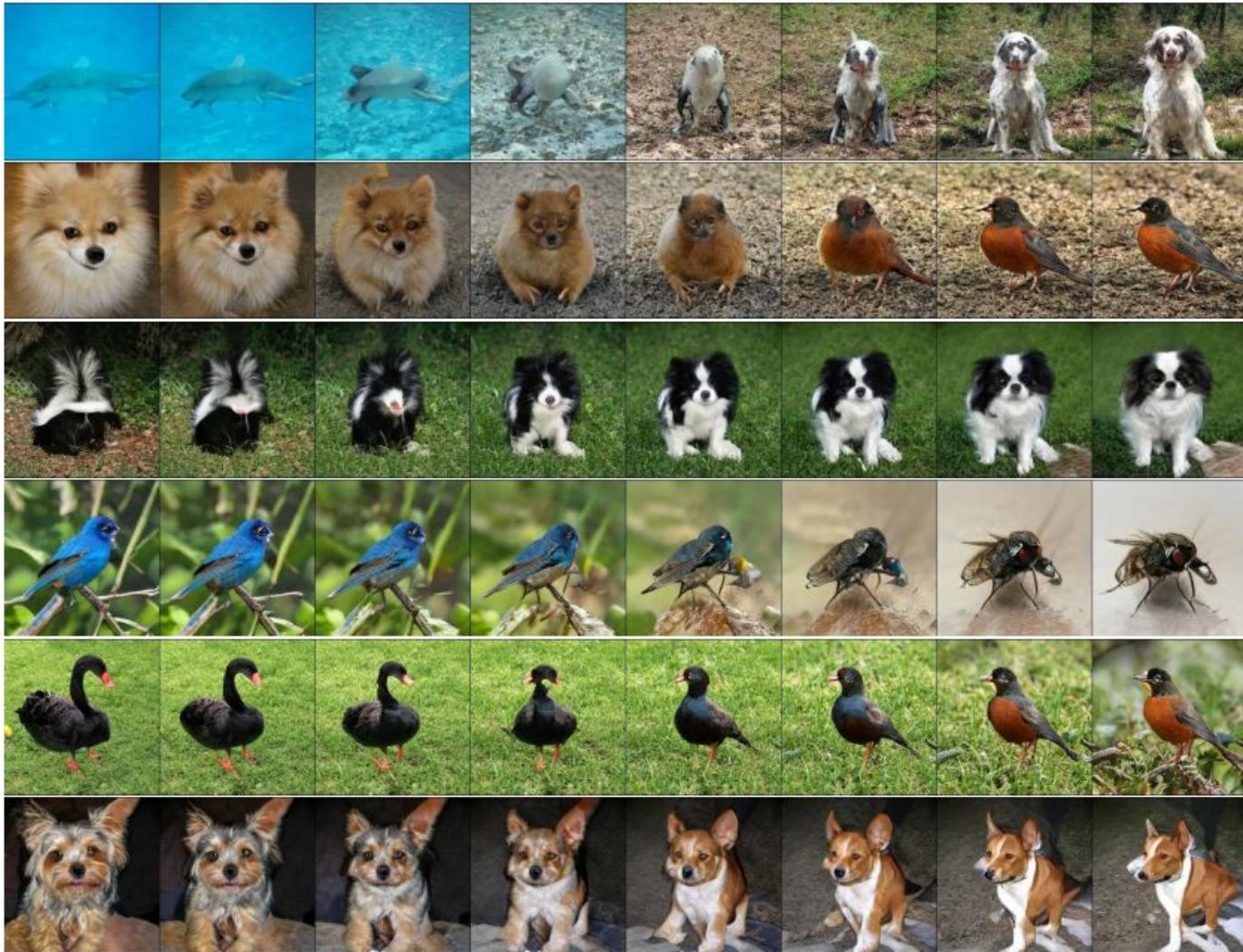


Figure 8: Interpolations between  $z, c$  pairs.

Est-ce qu'on doit conditionner  
seulement sur des classes ?

# Texte -> Images

Text description	This bird is red and brown in color, with a stubby beak	The bird is short and stubby with yellow on its body	A bird with a medium orange bill white body gray wings and webbed feet	This small black bird has a short, slightly curved bill and long legs	A small bird with varying shades of brown with white under the eyes	A small yellow bird with a black crown and a short black pointed beak	This small bird has a white breast, light grey head, and black wings and tail
256x256 StackGAN-v2							
Text description	This flower has a lot of small purple petals in a dome-like configuration	This flower is pink, white, and yellow in color, and has petals that are striped	This flower has petals that are dark pink with white edges and pink stamen	This flower is white and yellow in color, with petals that are wavy and smooth	A picture of a very clean living room	A group of people on skis stand in the snow	Eggs fruit candy nuts and meat served on white dish
256x256 StackGAN-v2							

Zhang, H., Xu, T., Li, H., Zhang, S., Wang, X., Huang, X., & Metaxas, D. N. (2018). Stackgan++: Realistic image synthesis with stacked generative adversarial networks. *IEEE transactions on pattern analysis and machine intelligence*, 41(8), 1947-1962.

# Pix2Pix

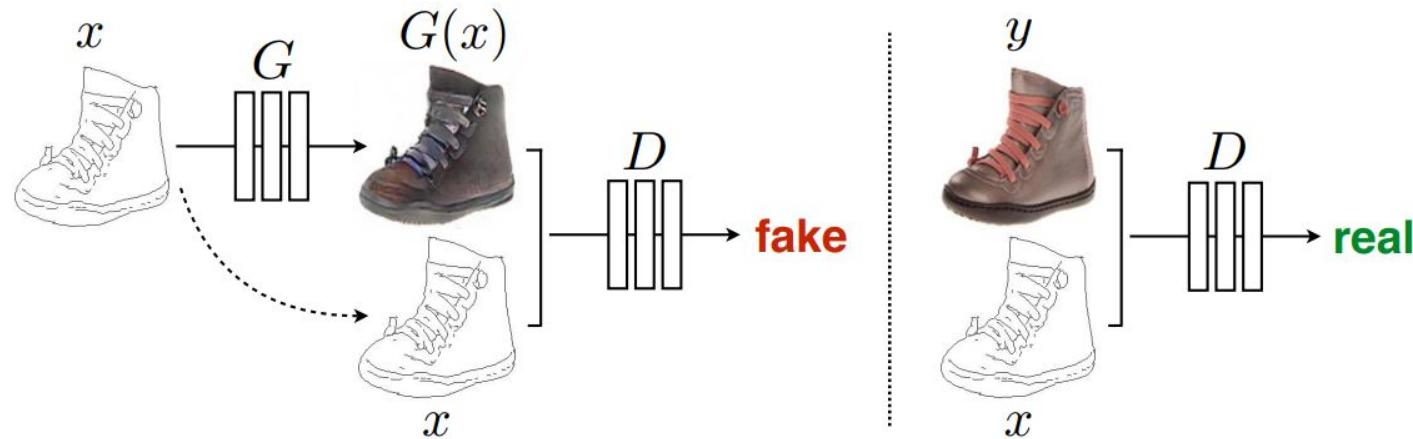


Figure 2: Training a conditional GAN to map edges→photo. The discriminator,  $D$ , learns to classify between fake (synthesized by the generator) and real {edge, photo} tuples. The generator,  $G$ , learns to fool the discriminator. Unlike an unconditional GAN, both the generator and discriminator observe the input edge map.

# Pix2Pix



Figure 16: Example results of our method on automatically detected edges→handbags, compared to ground truth.

Isola, P., Zhu, J. Y., Zhou, T., & Efros, A. A. (2017). Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1125-1134).

# Pix2Pix

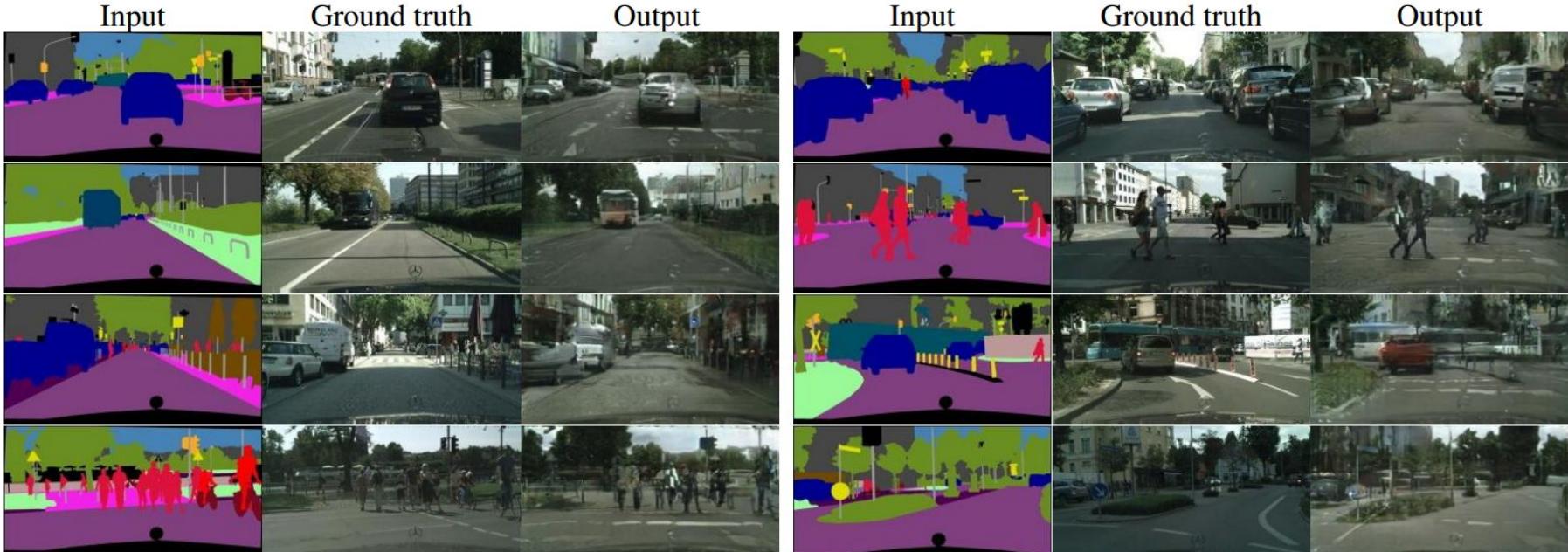


Figure 13: Example results of our method on Cityscapes labels→photo, compared to ground truth.

Isola, P., Zhu, J. Y., Zhou, T., & Efros, A. A. (2017). Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1125-1134).

# Pix2Pix

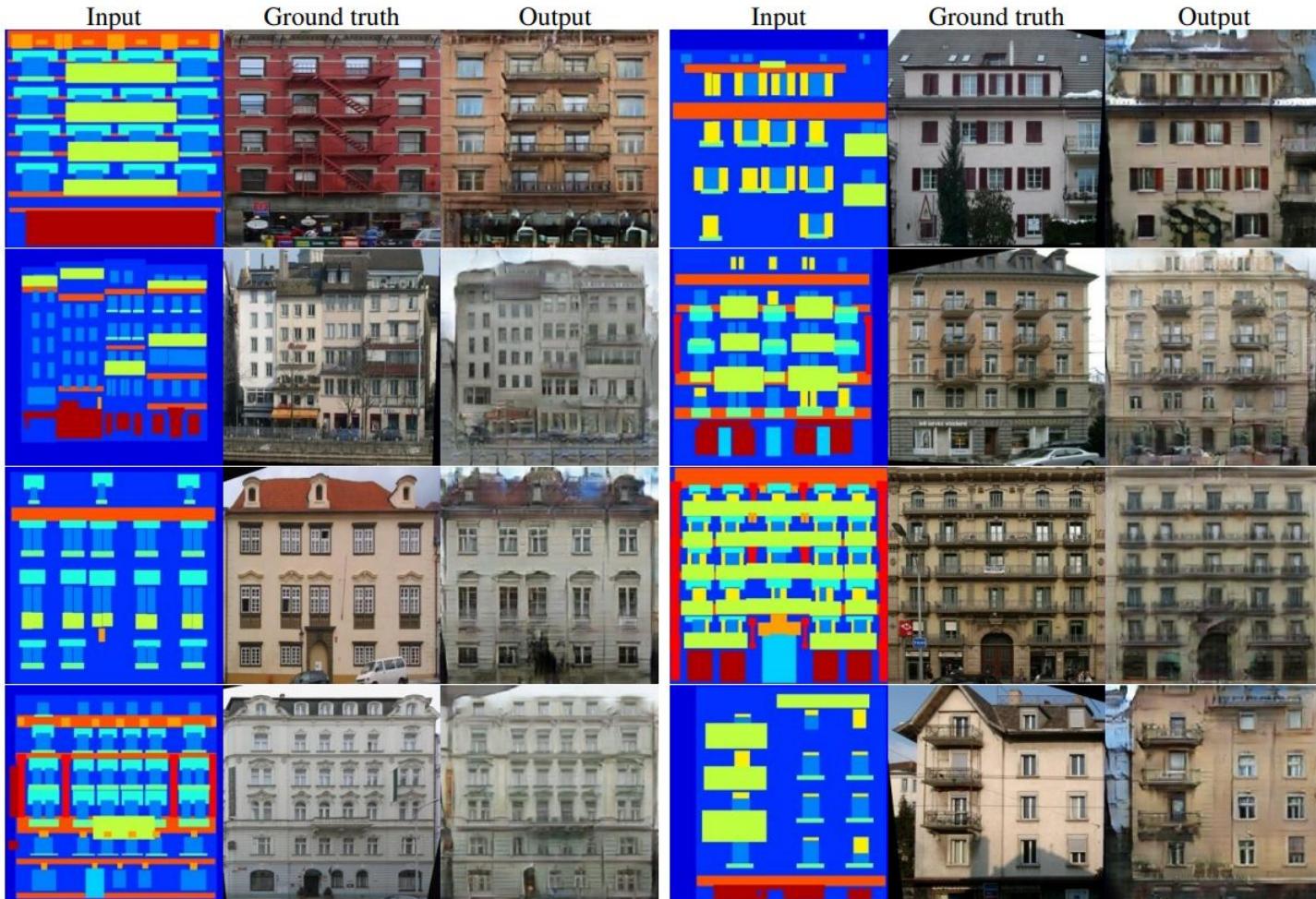


Figure 14: Example results of our method on facades labels→photo, compared to ground truth.

Isola, P., Zhu, J. Y., Zhou, T., & Efros, A. A. (2017). Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1125-1134).

# Pix2Pix

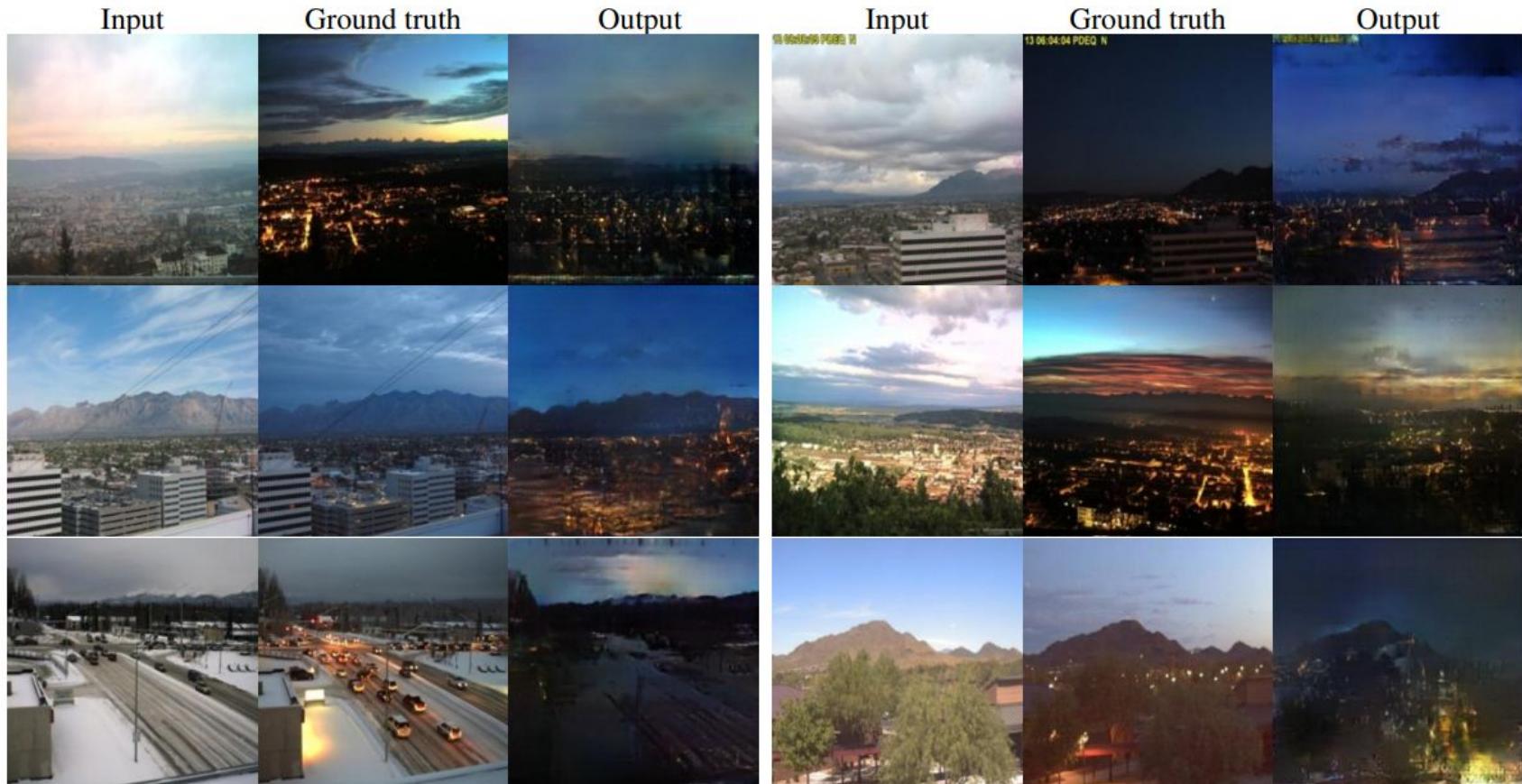


Figure 15: Example results of our method on day→night, compared to ground truth.

Isola, P., Zhu, J. Y., Zhou, T., & Efros, A. A. (2017). Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1125-1134).

# Pix2Pix



<https://affinelayer.com/pixsrv/>

Isola, P., Zhu, J. Y., Zhou, T., & Efros, A. A. (2017). Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1125-1134).

# Pix2Pix



Brannon Dorsey

@brannondorsey

...

More person-to-person video translation with  
**#machinelearning**, from me to Kurzweil. With puppet-  
like control. **#pix2pix**



6:58 PM · Dec 12, 2016 · Twitter Web Client

<https://twitter.com;brannondorsey/status/808461108881268736>

Isola, P., Zhu, J. Y., Zhou, T., & Efros, A. A. (2017). Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1125-1134).

# Semantic image synthesis

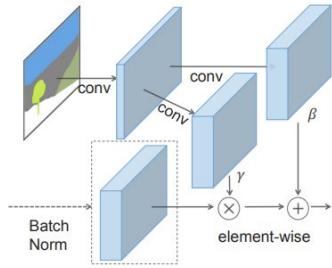


Figure 2: In the SPADE, the mask is first projected onto an embedding space and then convolved to produce the modulation parameters  $\gamma$  and  $\beta$ . Unlike prior conditional normalization methods,  $\gamma$  and  $\beta$  are not vectors, but tensors with spatial dimensions. The produced  $\gamma$  and  $\beta$  are multiplied and added to the normalized activation element-wise.

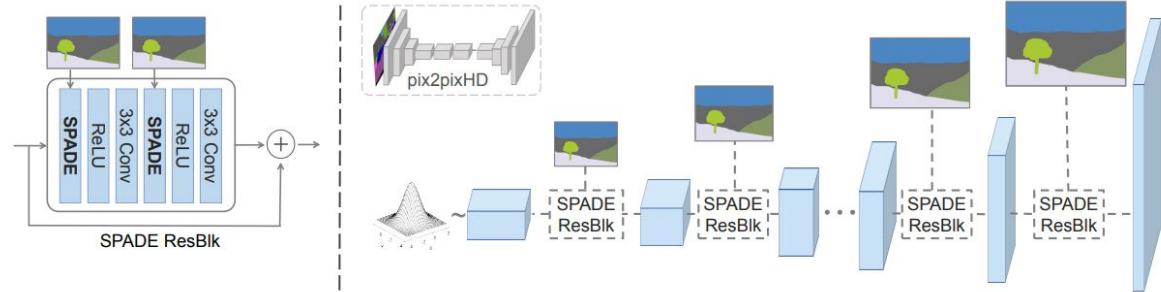


Figure 4: In the SPADE generator, each normalization layer uses the segmentation mask to modulate the layer activations. (left) Structure of one residual block with the SPADE. (right) The generator contains a series of the SPADE residual blocks with upsampling layers. Our architecture achieves better performance with a smaller number of parameters by removing the downsampling layers of leading image-to-image translation networks such as the pix2pixHD model [48].

# *Semantic image synthesis*

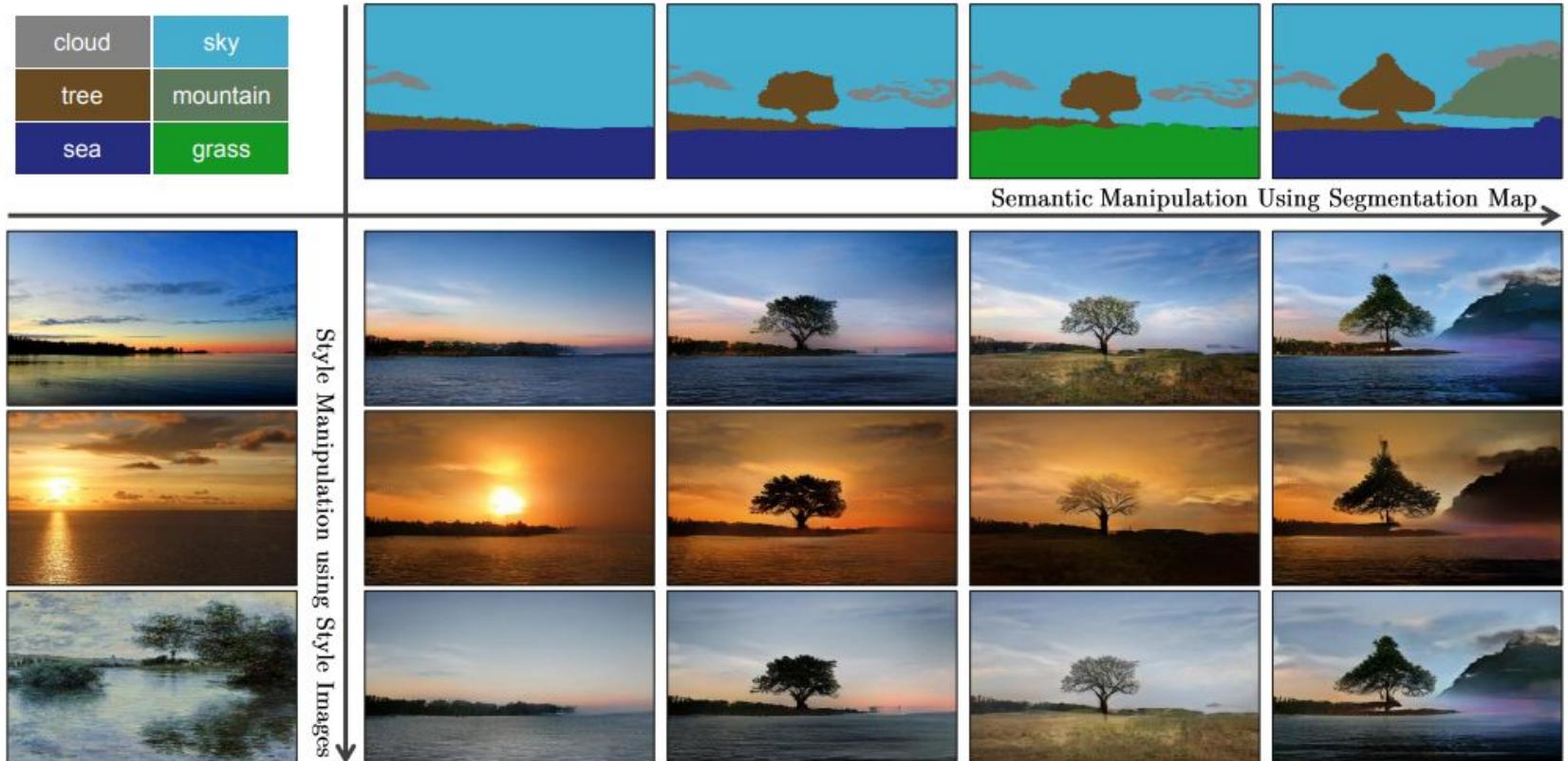


Figure 1: Our model allows user control over both semantic and style as synthesizing an image. The semantic (e.g., the existence of a tree) is controlled via a label map (the top row), while the style is controlled via the reference style image (the leftmost column). Please visit our [website](#) for interactive image synthesis demos.

# Semantic image synthesis

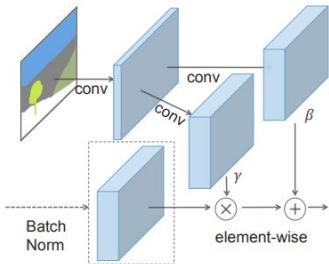


Figure 2: In the SPADE, the mask is first projected onto an embedding space and then convolved to produce the modulation parameters  $\gamma$  and  $\beta$ . Unlike prior conditional normalization methods,  $\gamma$  and  $\beta$  are not vectors, but tensors with spatial dimensions. The produced  $\gamma$  and  $\beta$  are multiplied and added to the normalized activation element-wise.

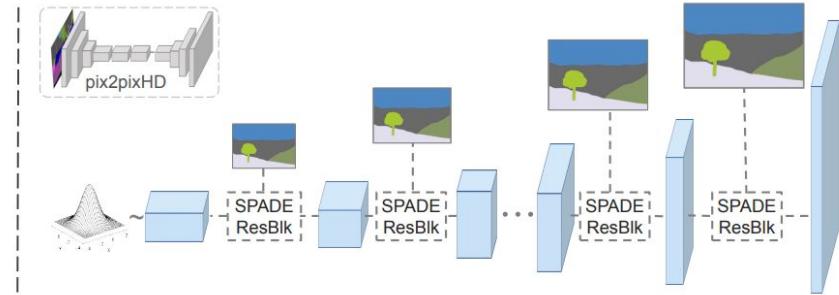
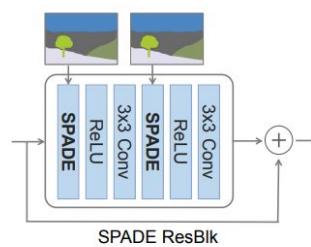
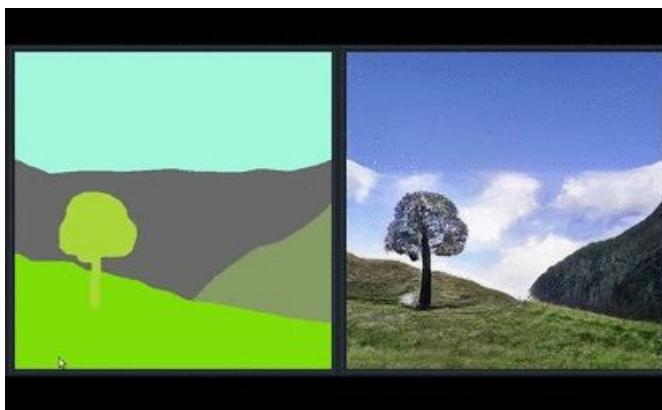


Figure 4: In the SPADE generator, each normalization layer uses the segmentation mask to modulate the layer activations. (left) Structure of one residual block with the SPADE. (right) The generator contains a series of the SPADE residual blocks with upsampling layers. Our architecture achieves better performance with a smaller number of parameters by removing the downsampling layers of leading image-to-image translation networks such as the pix2pixHD model [48].



<https://www.youtube.com/watch?v=MXWm6w4E5q0>

Park, T., Liu, M. Y., Wang, T. C., & Zhu, J. Y. (2019). Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 2337-2346).

Et si on a pas de paires d'images ?

# CycleGAN

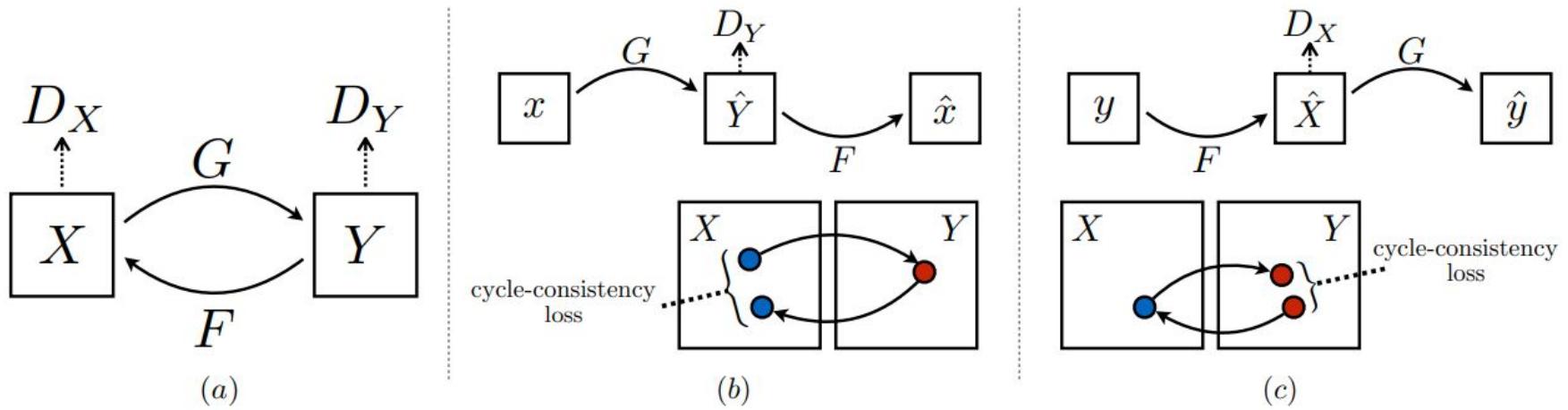
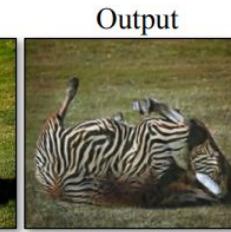
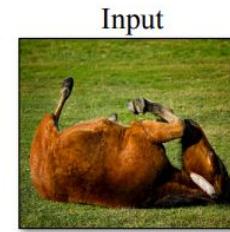
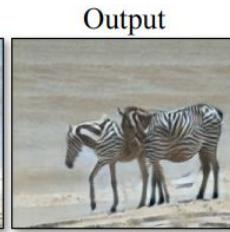
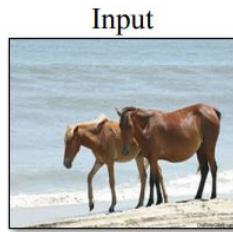
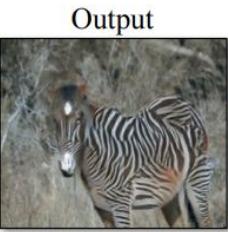
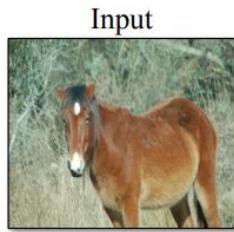


Figure 3: (a) Our model contains two mapping functions  $G : X \rightarrow Y$  and  $F : Y \rightarrow X$ , and associated adversarial discriminators  $D_Y$  and  $D_X$ .  $D_Y$  encourages  $G$  to translate  $X$  into outputs indistinguishable from domain  $Y$ , and vice versa for  $D_X$  and  $F$ . To further regularize the mappings, we introduce two *cycle consistency losses* that capture the intuition that if we translate from one domain to the other and back again we should arrive at where we started: (b) forward cycle-consistency loss:  $x \rightarrow G(x) \rightarrow F(G(x)) \approx x$ , and (c) backward cycle-consistency loss:  $y \rightarrow F(y) \rightarrow G(F(y)) \approx y$

# CycleGAN



horse → zebra



zebra → horse



winter Yosemite → summer Yosemite



summer Yosemite → winter Yosemite

# CycleGAN

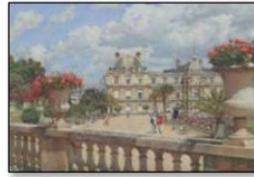


apple → orange



orange → apple

Input



# Code pytorch pour plus de 30 modèles de GANs

<https://github.com/eriklindernoren/PyTorch-GAN>

Mais les VAE ne sont pas obsolètes !

# Nouveau VAE (NVAE)

Vahdat, A., & Kautz, J. (2020). Nvae: A deep hierarchical variational autoencoder. *Advances in Neural Information Processing Systems*, 33, 19667-19679.

- VAE hiérarchique

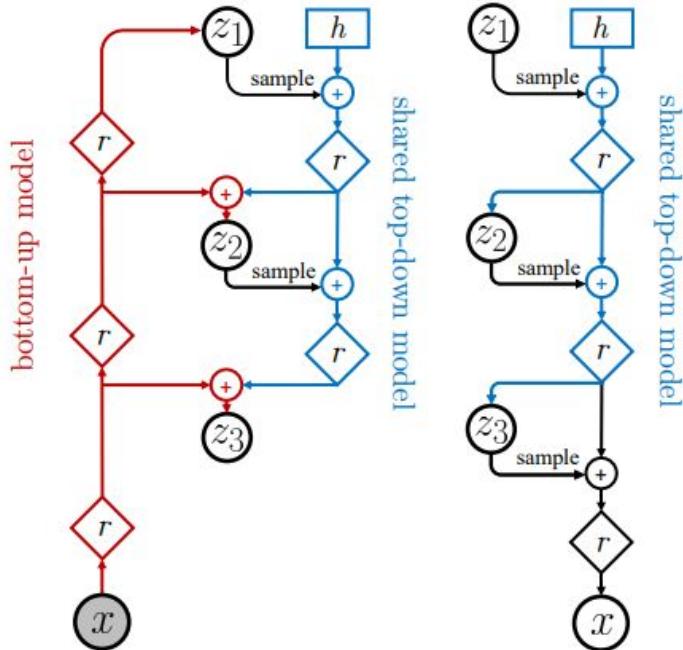
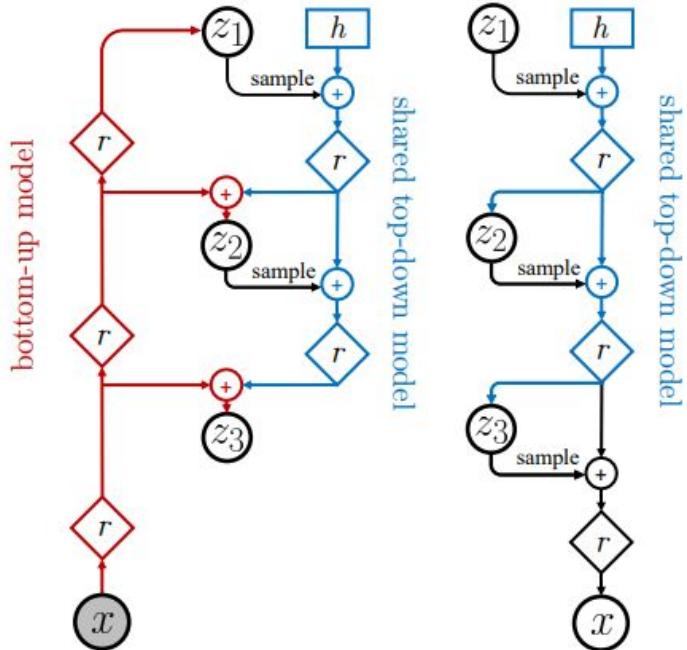


Figure 2: The neural networks implementing an encoder  $q(\mathbf{z}|\mathbf{x})$  and generative model  $p(\mathbf{x}, \mathbf{z})$  for a 3-group hierarchical VAE.  $\diamond^r$  denotes residual neural networks,  $\odot^+$  denotes feature combination (e.g., concatenation), and  $\square^h$  is a trainable parameter.

# Nouveau VAE (NVAE)

Vahdat, A., & Kautz, J. (2020). Nvae: A deep hierarchical variational autoencoder. *Advances in Neural Information Processing Systems*, 33, 19667-19679.



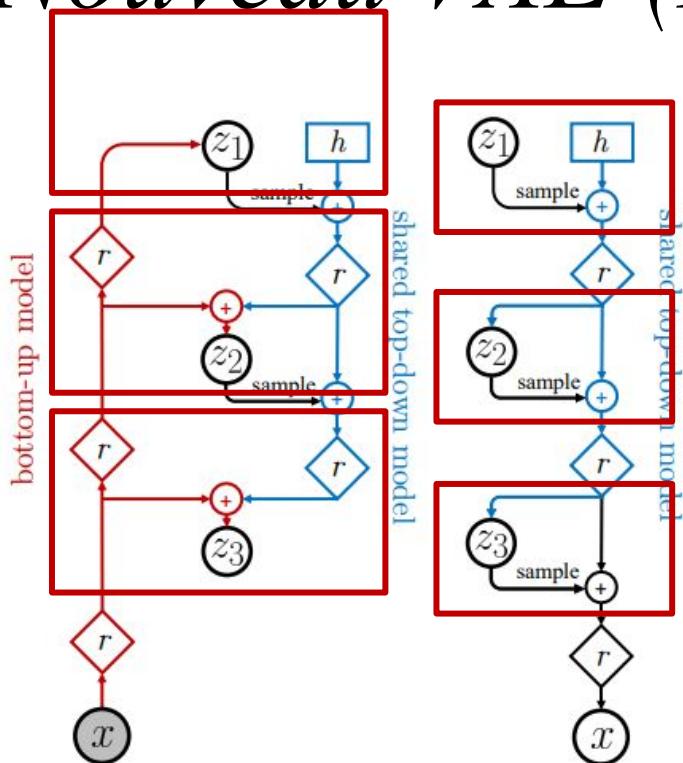
- VAE hiérarchique
- Les différents niveaux correspondent à des résolutions différentes

(a) Bidirectional Encoder (b) Generative Model

Figure 2: The neural networks implementing an encoder  $q(\mathbf{z}|\mathbf{x})$  and generative model  $p(\mathbf{x}, \mathbf{z})$  for a 3-group hierarchical VAE.  $\diamond^r$  denotes residual neural networks,  $\odot$  denotes feature combination (e.g., concatenation), and  $\square^h$  is a trainable parameter.

# Nouveau VAE (NVAE)

Vahdat, A., & Kautz, J. (2020). Nvae: A deep hierarchical variational autoencoder. *Advances in Neural Information Processing Systems*, 33, 19667-19679.



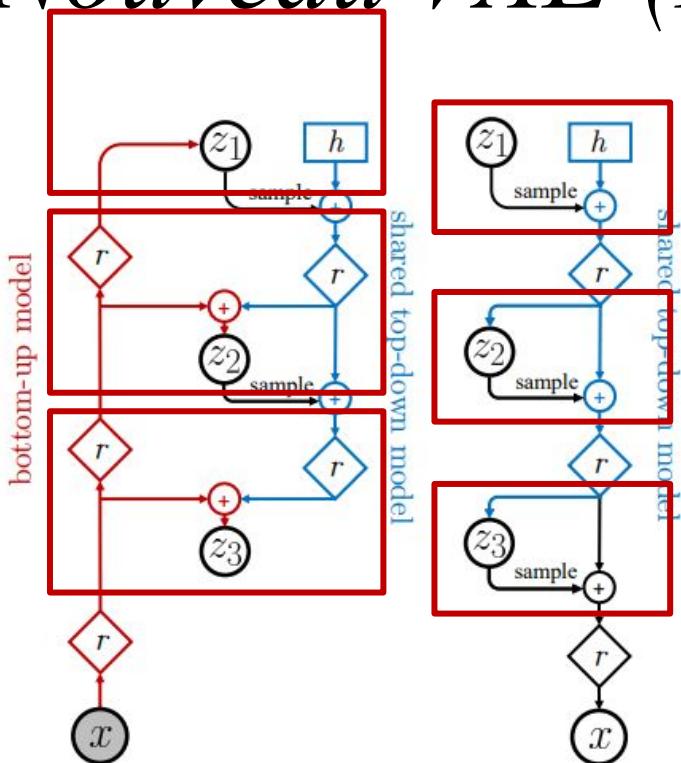
- VAE hiérarchique
- Les différents niveaux correspondent à des résolutions différentes
- Des  $z$  sont produits à chaque niveau

(a) Bidirectional Encoder (b) Generative Model

Figure 2: The neural networks implementing an encoder  $q(\mathbf{z}|\mathbf{x})$  and generative model  $p(\mathbf{x}, \mathbf{z})$  for a 3-group hierarchical VAE.  $\diamond^r$  denotes residual neural networks,  $\odot$  denotes feature combination (e.g., concatenation), and  $\square^h$  is a trainable parameter.

# Nouveau VAE (NVAE)

Vahdat, A., & Kautz, J. (2020). Nvae: A deep hierarchical variational autoencoder. *Advances in Neural Information Processing Systems*, 33, 19667-19679.



- VAE hiérarchique
- Les différents niveaux correspondent à des résolutions différentes
- Des  $z$  sont produits à chaque niveau
- Les différents niveaux de bruits correspondent à différents niveaux de détails

(a) Bidirectional Encoder (b) Generative Model

Figure 2: The neural networks implementing an encoder  $q(\mathbf{z}|\mathbf{x})$  and generative model  $p(\mathbf{x}, \mathbf{z})$  for a 3-group hierarchical VAE.  $\diamond^r$  denotes residual neural networks,  $\bigcirc +$  denotes feature combination (e.g., concatenation), and  $\square_h$  is a trainable parameter.

# *Nouveau VAE (NVAE)*

Vahdat, A., & Kautz, J. (2020). Nvae: A deep hierarchical variational autoencoder. *Advances in Neural Information Processing Systems*, 33, 19667-19679.

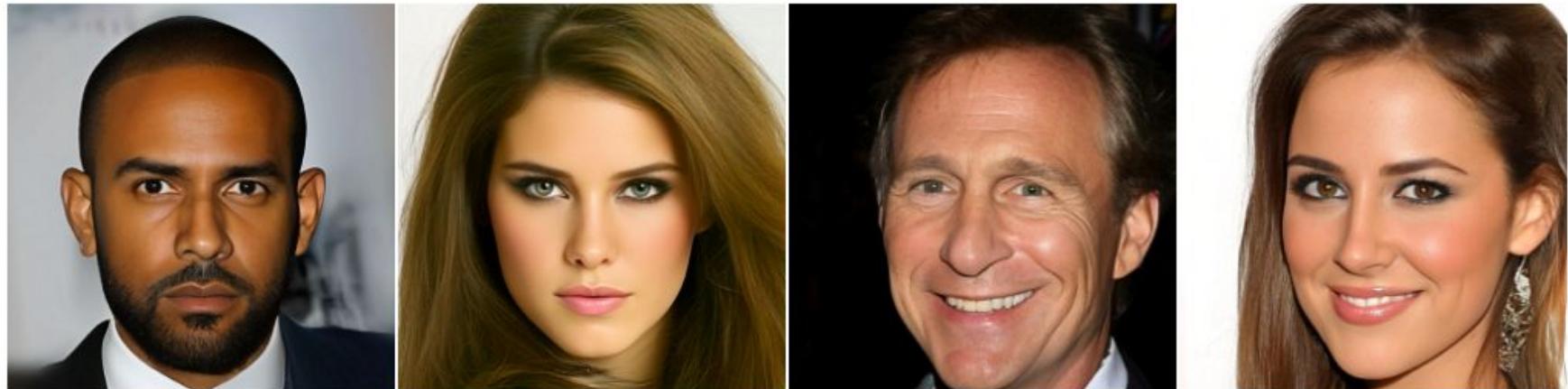


Figure 1:  $256 \times 256$ -pixel samples generated by NVAE, trained on CelebA HQ [28].

# *Nouveau VAE (NVAE)*

Vahdat, A., & Kautz, J. (2020). Nvae: A deep hierarchical variational autoencoder. *Advances in Neural Information Processing Systems*, 33, 19667-19679.



(a) MNIST ( $t = 1.0$ )



(b) CIFAR-10 ( $t = 0.7$ )



(c) CelebA 64 ( $t = 0.6$ )



(d) CelebA HQ ( $t = 0.6$ )

(e) FFHQ ( $t = 0.5$ )

# VQ-VAE

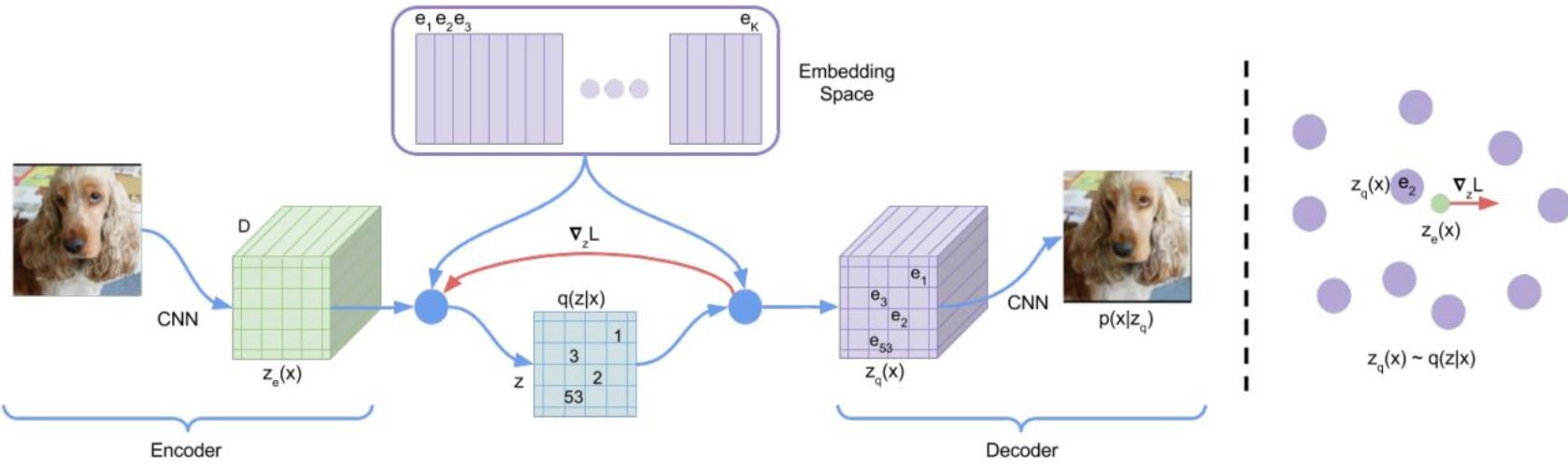


Figure 1: Left: A figure describing the VQ-VAE. Right: Visualisation of the embedding space. The output of the encoder  $z(x)$  is mapped to the nearest point  $e_2$ . The gradient  $\nabla_z L$  (in red) will push the encoder to change its output, which could alter the configuration in the next forward pass.

- Van Den Oord, A., & Vinyals, O. (2017). Neural discrete representation learning. *Advances in neural information processing systems*, 30.
- Razavi, A., Van den Oord, A., & Vinyals, O. (2019). Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems*, 32.

# VQ-VAE

Van Den Oord, A., & Vinyals, O. (2017). Neural discrete representation learning. *Advances in neural information processing systems*, 30.

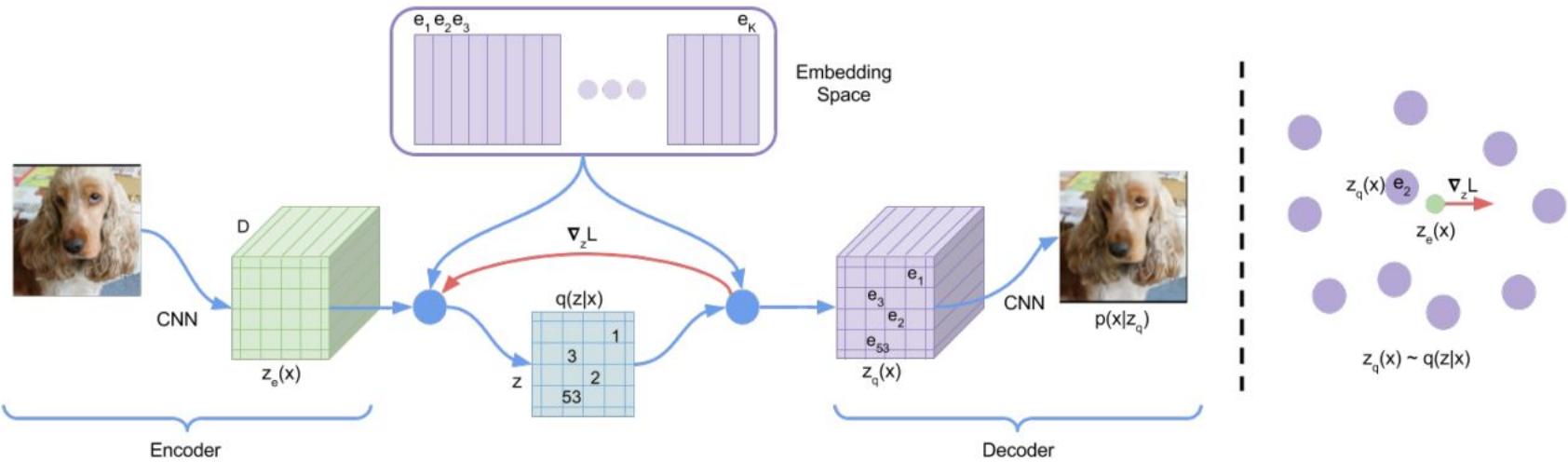


Figure 1: Left: A figure describing the VQ-VAE. Right: Visualisation of the embedding space. The output of the encoder  $z(x)$  is mapped to the nearest point  $e_2$ . The gradient  $\nabla_z L$  (in red) will push the encoder to change its output, which could alter the configuration in the next forward pass.

- L'espace latent est *discret* plutôt que continu
- L'encodeur apprend à projeter les images vers des centroïdes
- Chaque composante de  $z$  est un centroïde
  
- Van Den Oord, A., & Vinyals, O. (2017). Neural discrete representation learning. *Advances in neural information processing systems*, 30.
- Razavi, A., Van den Oord, A., & Vinyals, O. (2019). Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems*, 32.

# VQ-VAE

Van Den Oord, A., & Vinyals, O. (2017). Neural discrete representation learning. *Advances in neural information processing systems*, 30.

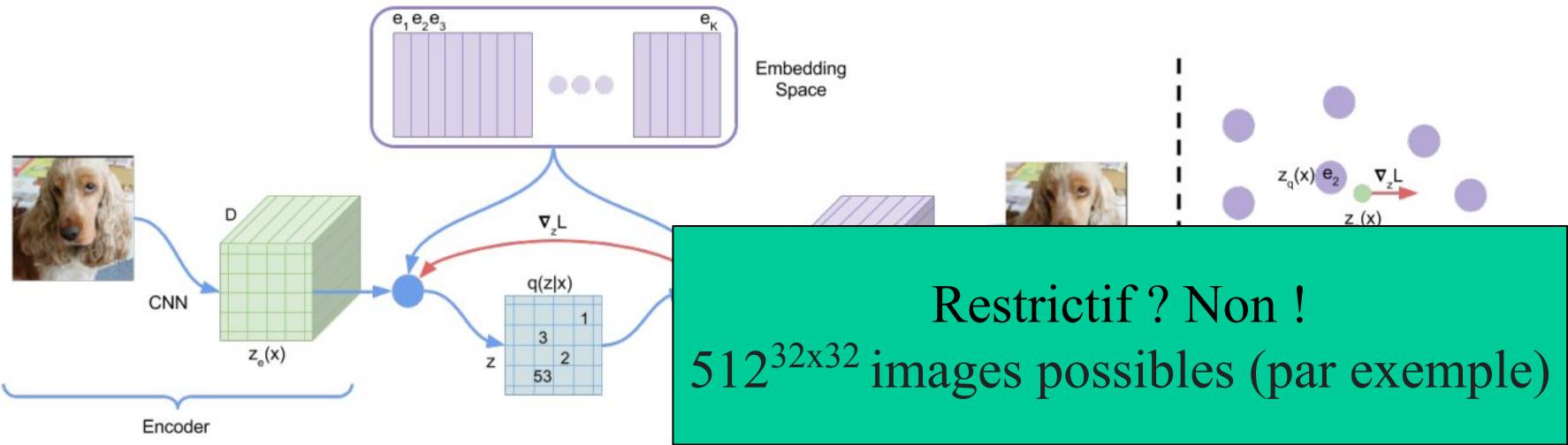
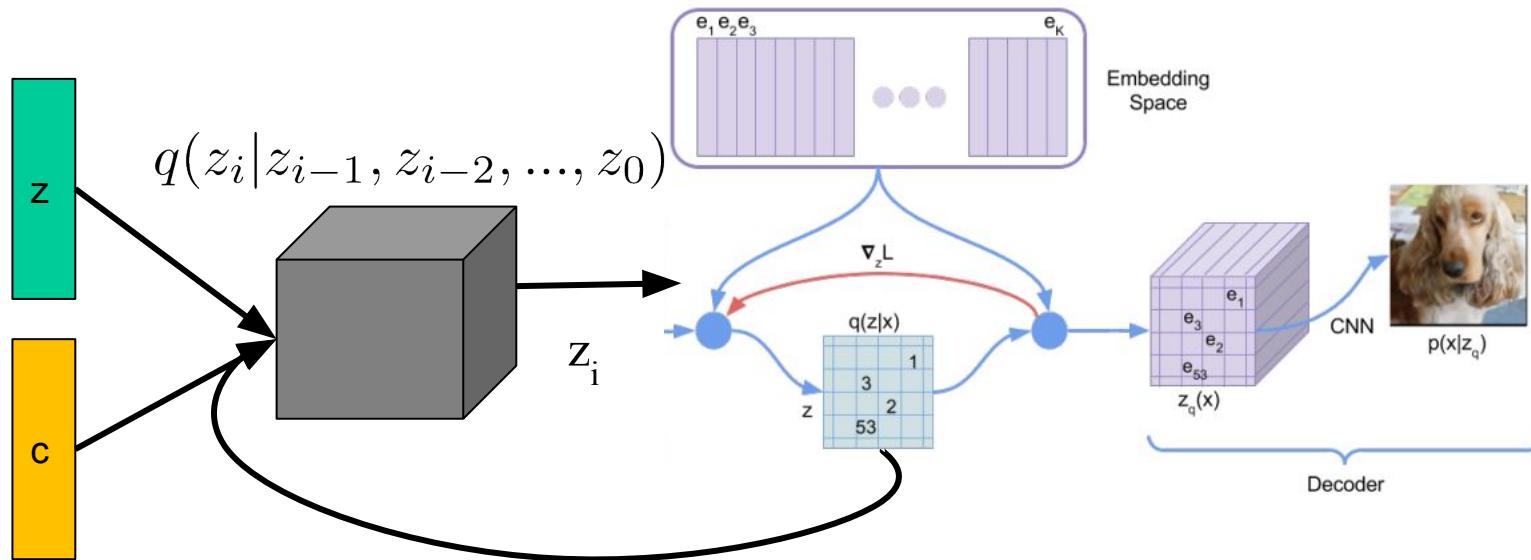


Figure 1: Left: A figure describing the VQ-VAE. Right: Visualisation of the embedding space. The output of the encoder  $z(x)$  is mapped to the nearest point  $e_2$ . The gradient  $\nabla_z L$  (in red) will push the encoder to change its output, which could alter the configuration in the next forward pass.

- L'espace latent est *discret* plutôt que continu
  - L'encodeur apprend à projeter les images vers des centroïdes
  - Chaque composante de  $z$  est un centroïde
  - Les centroïdes sont appris en étant rapprochés de la sortie de l'encodeur
- Van Den Oord, A., & Vinyals, O. (2017). Neural discrete representation learning. *Advances in neural information processing systems*, 30.
- Razavi, A., Van den Oord, A., & Vinyals, O. (2019). Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems*, 32.

# VQ-VAE

Van Den Oord, A., & Vinyals, O. (2017). Neural discrete representation learning. *Advances in neural information processing systems*, 30.



- En génération, un autre modèle apprend à choisir les  $e$

- Van Den Oord, A., & Vinyals, O. (2017). Neural discrete representation learning. *Advances in neural information processing systems*, 30.
- Razavi, A., Van den Oord, A., & Vinyals, O. (2019). Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems*, 32.

# VQ-VAE



Figure 1: Class-conditional 256x256 image samples from a two-level model trained on ImageNet.

- Van Den Oord, A., & Vinyals, O. (2017). Neural discrete representation learning. *Advances in neural information processing systems*, 30.
- Razavi, A., Van den Oord, A., & Vinyals, O. (2019). Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems*, 32.

# VQ-VAE



**VQ-VAE (Proposed)**

**BigGAN deep**

- Van Den Oord, A., & Vinyals, O. (2017). Neural discrete representation learning. *Advances in neural information processing systems*, 30.
- Razavi, A., Van den Oord, A., & Vinyals, O. (2019). Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems*, 32.

# DALL-E

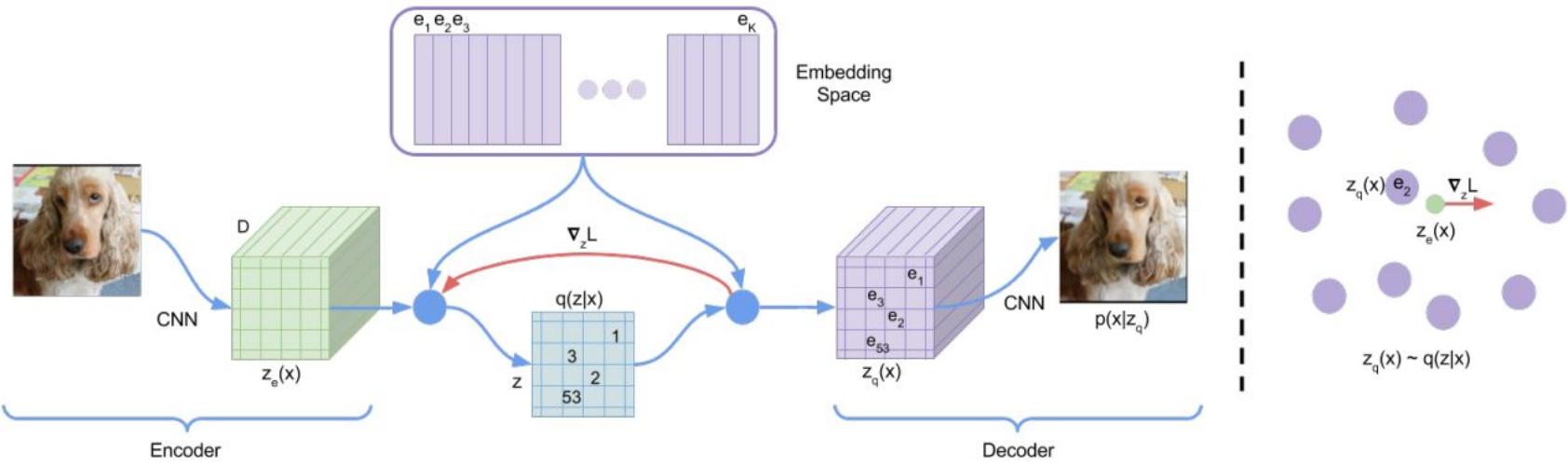
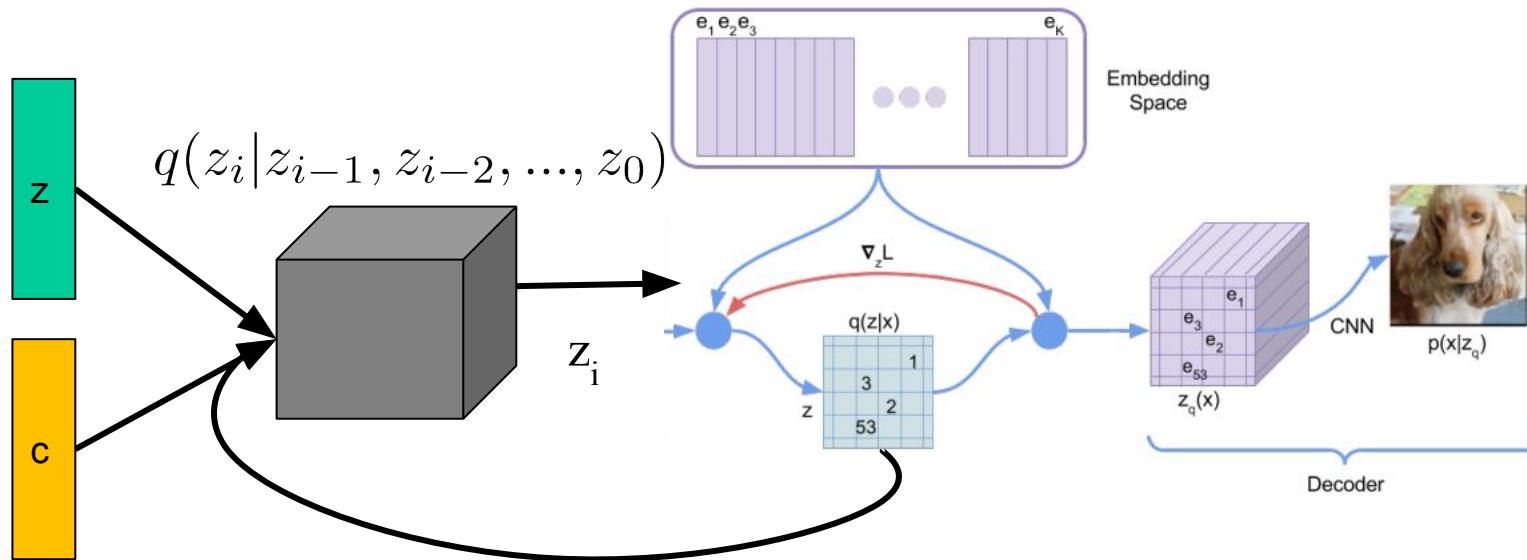


Figure 1: Left: A figure describing the VQ-VAE. Right: Visualisation of the embedding space. The output of the encoder  $z(x)$  is mapped to the nearest point  $e_2$ . The gradient  $\nabla_z L$  (in red) will push the encoder to change its output, which could alter the configuration in the next forward pass.

- Comme VQ-VAE, mais plus gros et avec quelques twists
- Voir le code: <https://github.com/openai/DALL-E>

# DALL-E

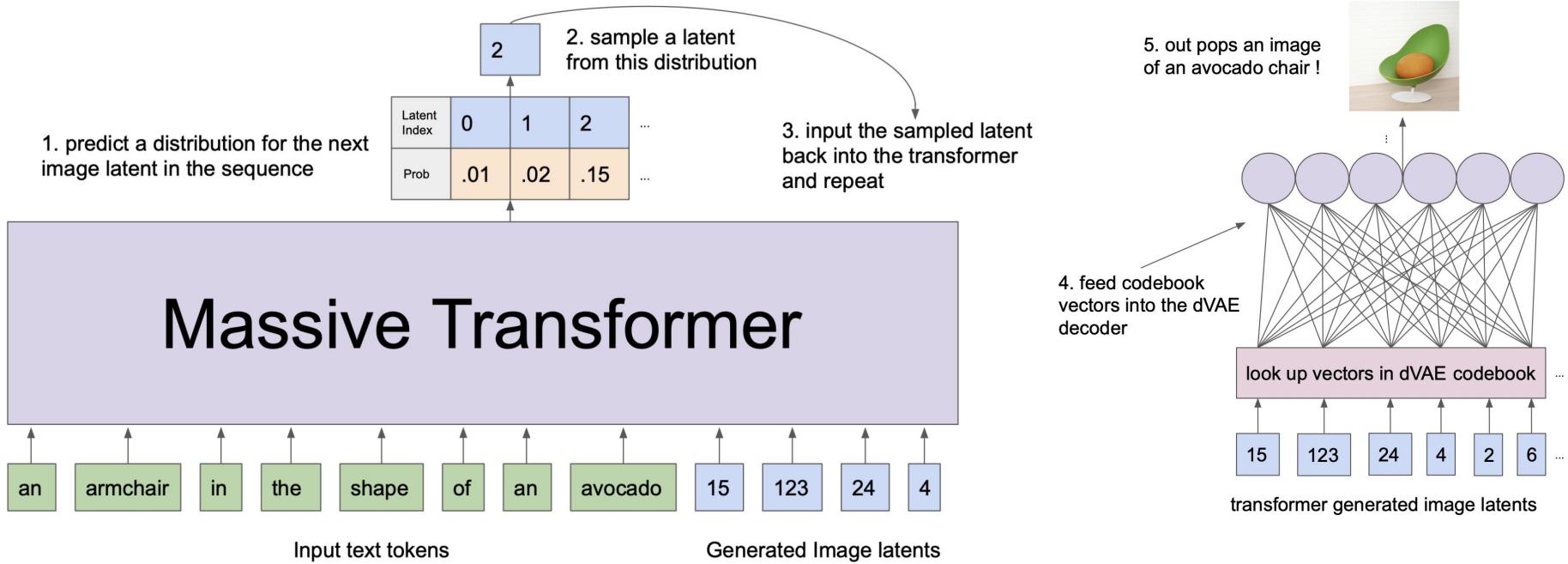
Van Den Oord, A., & Vinyals, O. (2017). Neural discrete representation learning. *Advances in neural information processing systems*, 30.



- Le “générateur” est un transformer (GPT-3 avec 12 millions de paramètres)

# DALL-E

Van Den Oord, A., & Vinyals, O. (2017). Neural discrete representation learning. *Advances in neural information processing systems*, 30.



source: <https://ml.berkeley.edu/blog/posts/dalle2/>

- Le “générateur” est un transformer (GPT-3 avec 12 millions de paramètres)

# DALL-E

Van Den Oord, A., & Vinyals, O. (2017). Neural discrete representation learning. *Advances in neural information processing systems*, 30.

TEXT PROMPT

an illustration of a baby daikon radish in a tutu walking a dog

AI-GENERATED  
IMAGES



TEXT PROMPT

an armchair in the shape of an avocado....

AI-GENERATED  
IMAGES



<https://openai.com/blog/dall-e/>

Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Radford, A., ... & Sutskever, I. (2021, July). Zero-shot text-to-image generation. In *International Conference on Machine Learning* (pp. 8821-8831). PMLR.

# DALL-E

Van Den Oord, A., & Vinyals, O. (2017). Neural discrete representation learning. *Advances in neural information processing systems*, 30.

TEXT PROMPT

a store front that has the word 'openai' written on it....

AI-GENERATED IMAGES



TEXT & IMAGE PROMPT

the exact same cat on the top as a sketch on the bottom

AI-GENERATED IMAGES



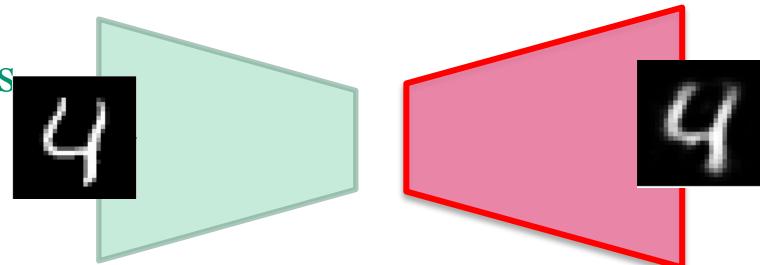
<https://openai.com/blog/dall-e/>

Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Radford, A., ... & Sutskever, I. (2021, July). Zero-shot text-to-image generation. In *International Conference on Machine Learning* (pp. 8821-8831). PMLR.

# Conclusion

Autoencodeurs:

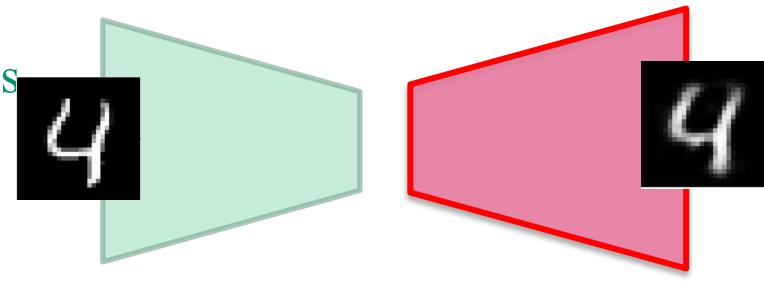
- Permettent d'apprendre des caractéristiques des données sans étiquettes de classe
- Peuvent être utilisés par des modèles supervisés
- Difficiles à échantillonner



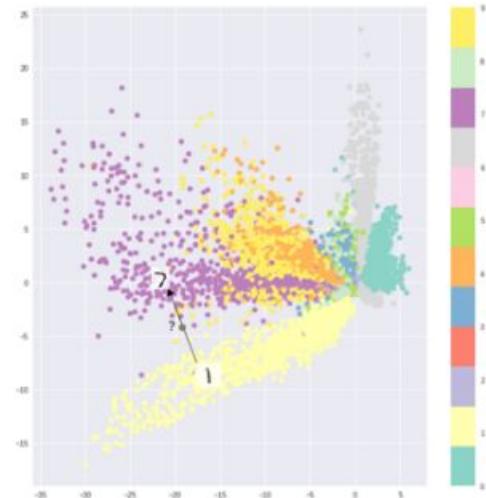
# Conclusion

Autoencodeurs:

- Permettent d'apprendre des caractéristiques des données sans étiquettes de classe
- Peuvent être utilisés par des modèles supervisés
- Difficiles à échantillonner



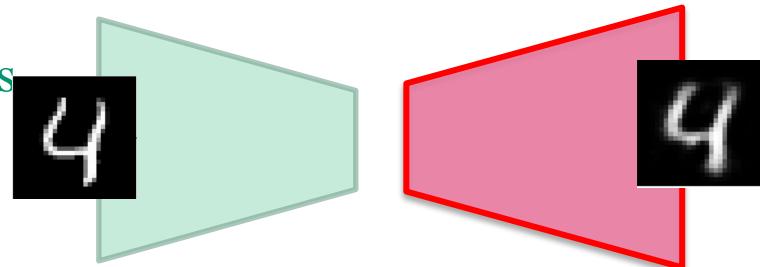
Only reconstruction loss



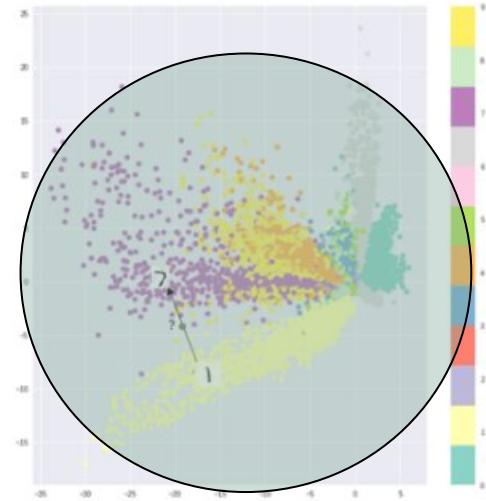
# Conclusion

Autoencodeurs:

- Permettent d'apprendre des caractéristiques des données sans étiquettes de classe
- Peuvent être utilisés par des modèles supervisés
- Difficiles à échantillonner



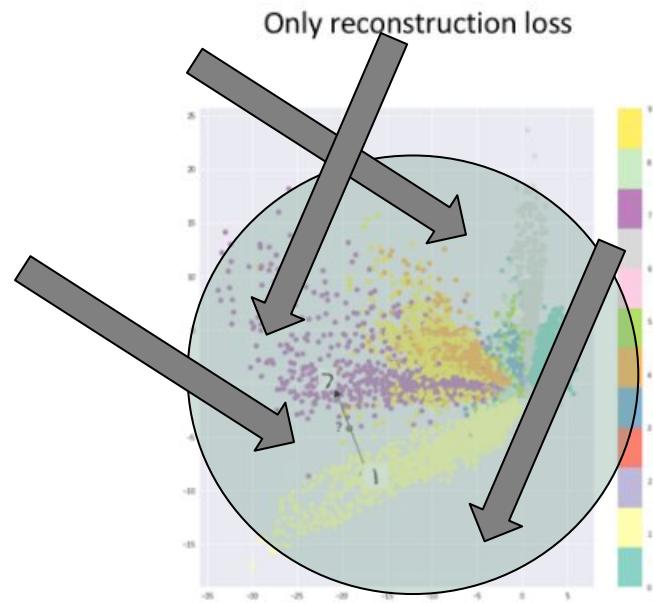
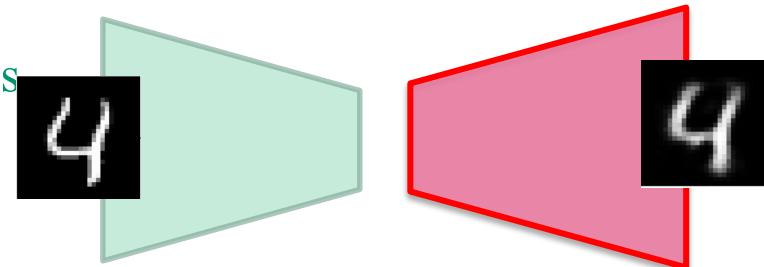
Only reconstruction loss



# Conclusion

Autoencodeurs:

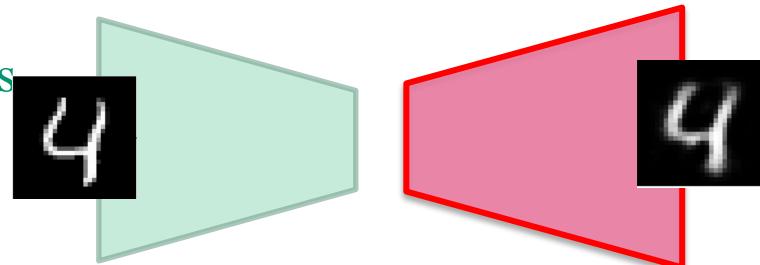
- Permettent d'apprendre des caractéristiques des données sans étiquettes de classe
- Peuvent être utilisés par des modèles supervisés
- Difficiles à échantillonner



# Conclusion

Autoencodeurs:

- Permettent d'apprendre des caractéristiques des données sans étiquettes de classe
- Peuvent être utilisés par des modèles supervisés
- Difficiles à échantillonner



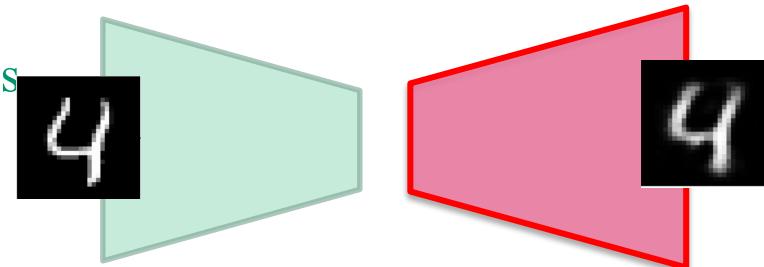
Autoencodeurs variationnels:

- Approche théoriquement intéressante aux AE
- $q(z|x)$  appris est échantillonnable et utile
- Génère des images floues

# Conclusion

Autoencodeurs:

- Permettent d'apprendre des caractéristiques des données sans étiquettes de classe
- Peuvent être utilisés par des modèles supervisés
- Difficiles à échantillonner



Autoencodeurs variationnels:

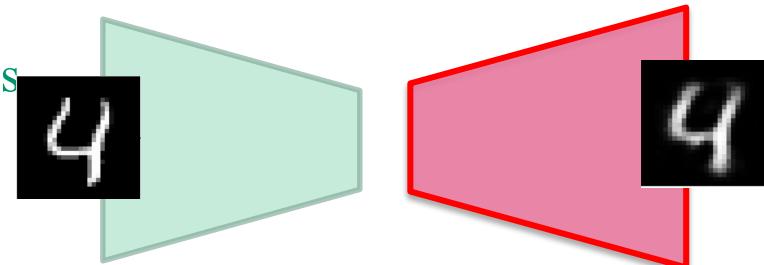
- Approche théoriquement intéressante aux AE
- $q(z|x)$  appris est échantillonnable et utile
- Génère des images floues

$$Loss = \frac{1}{2} \sum_{i=1}^d (1 + \log(\sigma_i^2) + \mu_i^2 - \sigma_i^2) - \lambda \|\vec{x} - \vec{x}'\|^2$$

# Conclusion

Autoencodeurs:

- Permettent d'apprendre des caractéristiques des données sans étiquettes de classe
- Peuvent être utilisés par des modèles supervisés
- Difficiles à échantillonner



Autoencodeurs variationnels:

- Approche théoriquement intéressante aux AE
- $q(z|x)$  appris est échantillonnable et utile
- Génère des images floues

*Generative adversarial networks (GANs):*

- Abandonne l'idée d'apprendre  $p(x)$
- Meilleurs résultats en génération
- Difficiles à entraîner

