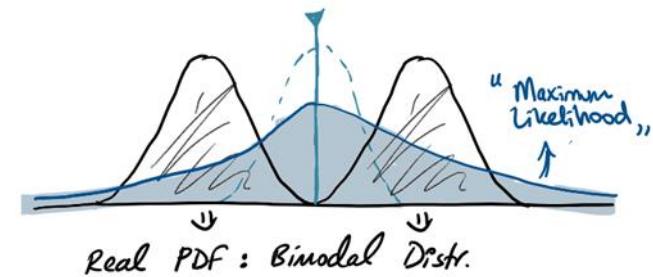
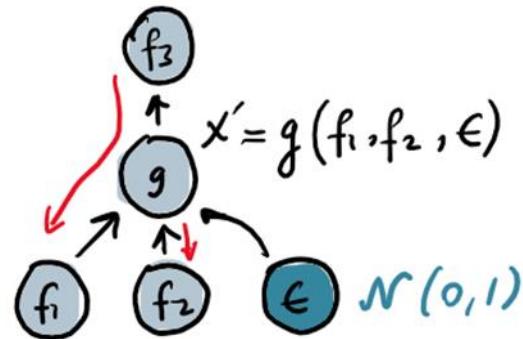
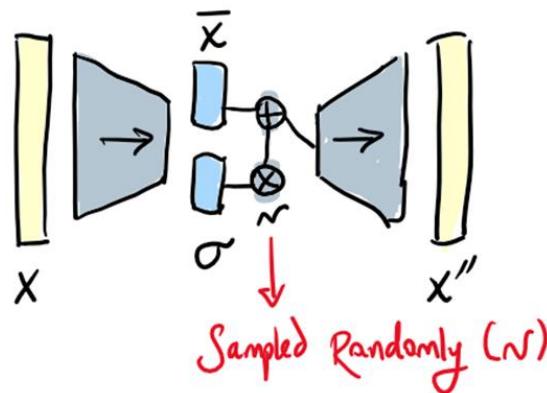


# Data Driven Engineering I: Machine Learning for Dynamical Systems

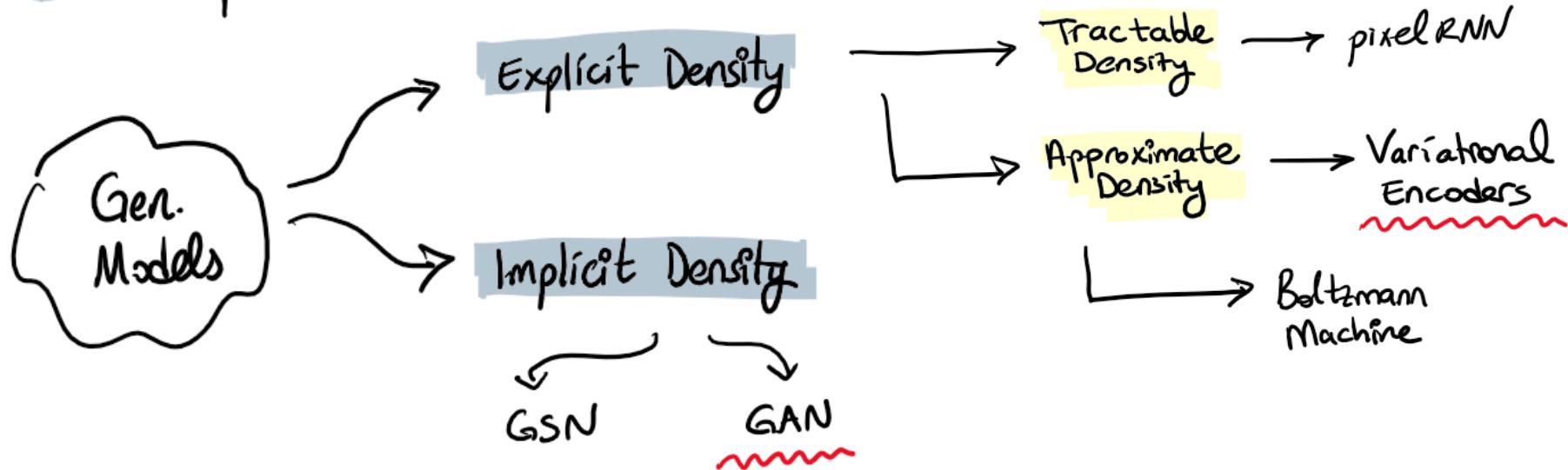
## Introduction to Generative Learning: VAEs and GANs

Institute of Thermal Turbomachinery  
Prof. Dr.-Ing. Hans-Jörg Bauer



# UL → Generative Models :

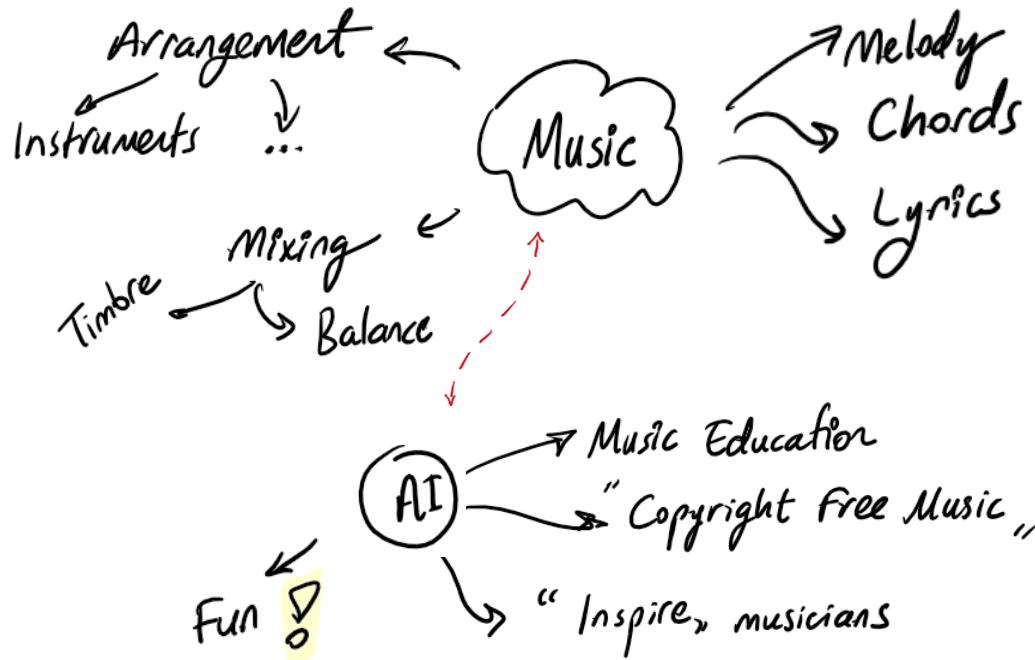
- \* It is probabilistic in nature



# ♪ Case Study: Composing Music with Gen. Methods



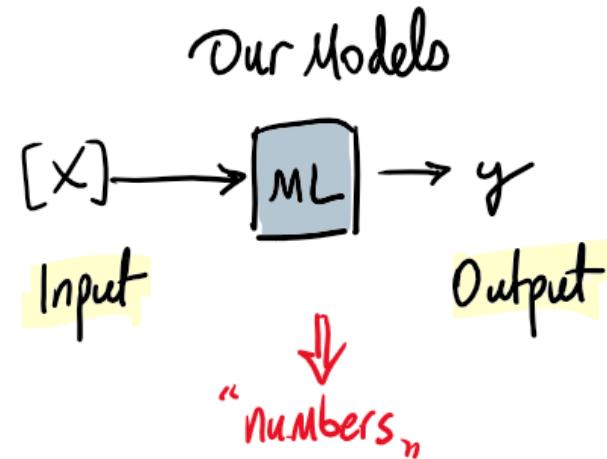
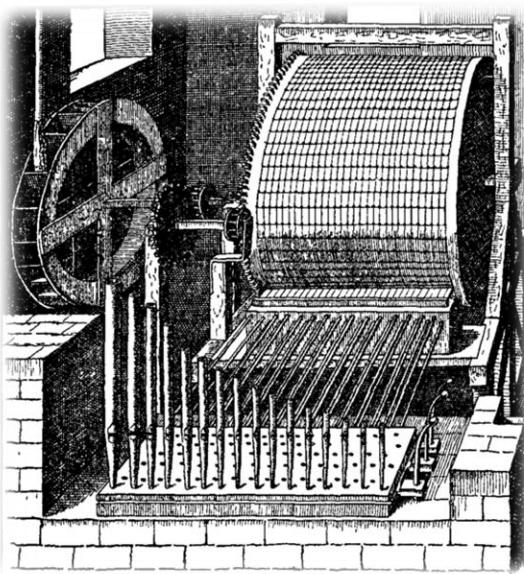
Bach  
↓  
Solo Violin



Industry :

- \* Google Magenta
- \* Spotify
- \* amper
- \* IBM
- \* Sony
- \* ...

# Modeling Music



? Representation  
of music

? Conversion of  
music

# I/O in Music

## ① Symbolic Representation

\* Piano Rolls

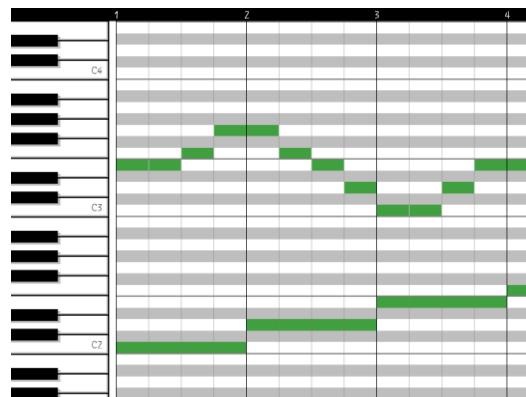


"image like,"

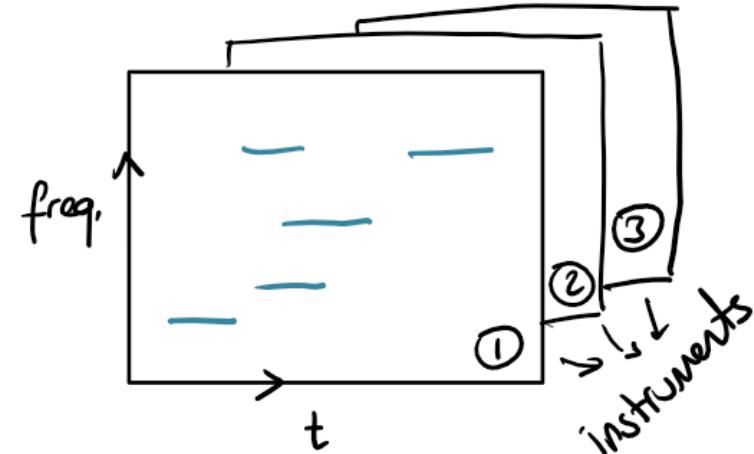
~~Eg:~~

\* MuseGAN

\* MidiNET



## ② Audio



# I/O in Music

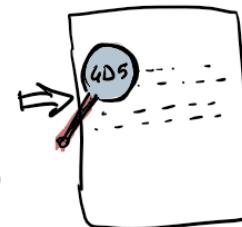
## ① Symbolic Representation

\* MIDI events

↓  
"text",

Musical  
Instrument  
Digital  
Interface

} "standard  
format",



## ② Audio

music21



```
(0.0) <music21.instrument.Violin Violin>
{0.0} <music21.tempo.MetronomeMark Quarter=250.0>
{0.0} <music21.key.Key of D major>
{0.0} <music21.meter.TimeSignature 4/4>
{0.0} <music21.note.Rest rest>
{3.75} <music21.tempo.MetronomeMark largamente Qu
{3.75} <music21.chord.Chord F#5>
{4.0} <music21.chord.Chord B3 F#4 D5 F#5>
{4.75} <music21.chord.Chord B5>
{5.0} <music21.chord.Chord E5 G5>
```

# I/O in Music

## ① Symbolic Representation

\* Piano Rolls



“image like”

\* Score



\* MIDI events



“text”

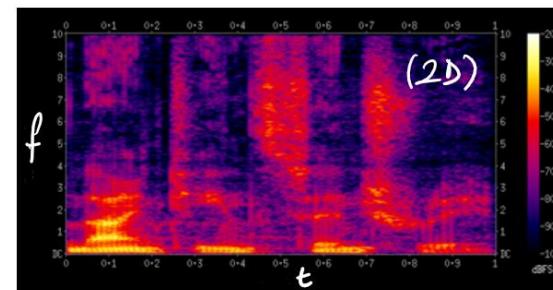
## ② Audio

\* Spectrogram



“image-like”

\* ~Wave form



# Understanding the Data

① ~metadata [X]

```

{0.0} <music21.instrument.Violin Violin>
{0.0} <music21.tempo.MetronomeMark Quarter=250.0>
{0.0} <music21.key.Key of D major>
{0.0} <music21.meter.TimeSignature 4/4>
{0.0} <music21.note.Rest rest>
{3.75} <music21.tempo.MetronomeMark largamente Quarter=32.0>
{3.75} <music21.chord.Chord F#5>
{4.0} <music21.chord.Chord B3 F#4 D5 F#5>
{4.75} <music21.chord.Chord B5>
{5.0} <music21.chord.Chord E5 G5>
  
```

- Starts @ beat 4
- Duration  $\Rightarrow 0.75$
- Chords  $\Rightarrow \begin{pmatrix} B \\ D F \end{pmatrix}$

# Vast Music Space

②  $[x] := \text{MIDI} \Rightarrow \text{Text} \Rightarrow \text{Numbers}$

? Multiple Notes

$\Rightarrow$  MIDI: 128 Single notes

$\Leftrightarrow$  Combinations  $\Rightarrow [E5 G5]$

~~Eg~~ Classes; A, B, C  $\Rightarrow \{\{A\}, \{B\}, \{C\}\}$

↳ examples  $f_A \ f_B \ f_C$

By  
TF

	$\{A\}$	$\{B\}$	$\{C\}$
$\{A\}$	1	0	0
$\{C\}$	0	0	1
$\{B\}$	0	1	0

# Vast Music Space

⇒ MIDI: 128 Single notes

↳ Combinations ⇒ [E5 G5]

~~Eg~~ {1 note, 2 notes}

```

{0.0} <music21.instrument.Violin Violin>
{0.0} <music21.tempo.MetronomeMark Quarter=250.0>
{0.0} <music21.key.Key of D major>
{0.0} <music21.meter.TimeSignature 4/4>
{0.0} <music21.note.Rest rest>
{3.75} <music21.tempo.MetronomeMark largamente Quarter=32.0>
{3.75} <music21.chord.Chord F#5>
{4.0} <music21.chord.Chord B3 F#4 D5 F#5>
{4.75} <music21.chord.Chord B5>
{5.0} <music21.chord.Chord E5 G5>
  
```

$$\hookrightarrow \frac{128!}{(128-2)! \cdot 2!} = \frac{128 \times 127}{2} = 8128 \text{ dimensions ?}$$

Solution  
 Notes = 128 ⇒ Reduce the number; (filter to C Major key)

③

## One-hot encoding:

- \* text  $\Leftrightarrow$  number ; reversible
- k Use dictionaries ~ look-up tables

(i) (Chord)  $\rightarrow$  get unique chords  $\Rightarrow$  Give an integer to each unique chord

(ii) (Duration)  $\rightarrow$  get unique durations  $\Rightarrow$  Give an integer to each unique duration

(iii) Create reverse dict. as well.

```
{0.0} <music21.instrument.Violin Violin>
{0.0} <music21.tempo.MetronomeMark Quarter=250.0>
{0.0} <music21.key.Key of D major>
{0.0} <music21.meter.TimeSignature 4/4>
{0.0} <music21.note.Rest rest>
{3.75} <music21.tempo.MetronomeMark largamente Quarter=32.0>
{3.75} <music21.chord.Chord F#5>
{4.0} <music21.chord.Chord B3 F#4 D5 F#5>
{4.75} <music21.chord.Chord B5>
{5.0} <music21.chord.Chord E5 G5>
```

{'A3': 0,	{0: 'A3',
'A3.A4': 1,	1: 'A3.A4',
'A3.A4.C#5': 2,	2: 'A3.A4.C#5',
'A3.A4.C5': 3,	3: 'A3.A4.C5',
'A3.A4.C5.A5': 4,	4: 'A3.A4.C5.A5',
'A3.A4.C5.B5': 5,	5: 'A3.A4.C5.B5',
'A3.A4.C5.C6': 6,	6: 'A3.A4.C5.C6',
'A3.A4.C5.E5': 7,	
'A3.A4.C5.F#5': 8,	
'A3.A4.C5.F5': 9,	

## ④ Creating [x] from Sequential Data

- \* Chord = [4, 11, 46, 122, 5, 7, ...]      Duration = [2, 2, 1, 1, 4, 3, ...]       $\Rightarrow$  "Sliding Windows"
- \* Sequence size = 32       $\Rightarrow$   $C_0 [0 \rightarrow 31]$        $D_0 [0 \rightarrow 31]$   
 $C_1 [1 \rightarrow 32]$        $D_1 [1 \rightarrow 32]$   
 $\dots$        $\dots$

```
{
  'A3': 0,
  'A3.A4': 1,
  'A3.A4.C#5': 2,
  'A3.A4.C5': 3,
  'A3.A4.C5.A5': 4,
  'A3.A4.C5.B5': 5,
  'A3.A4.C5.C6': 6,
  'A3.A4.C5.E5': 7,
  'A3.A4.C5.F#5': 8,
  'A3.A4.C5.F5': 9,
}
```

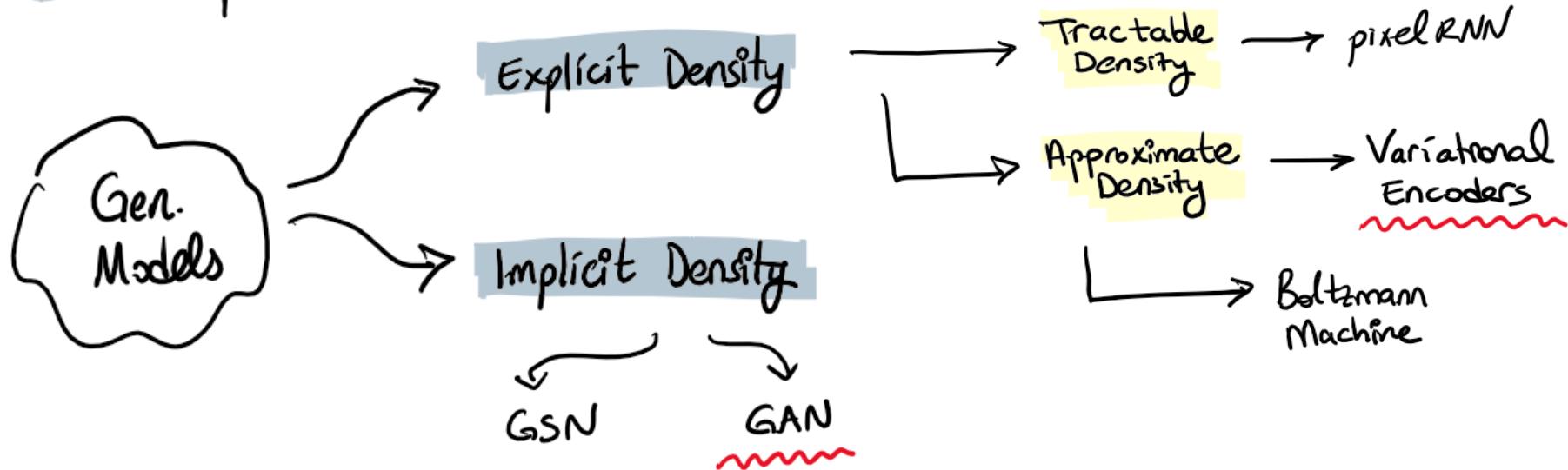
Ateliers & Saveurs in Montreal



# colab

# UL → Generative Models :

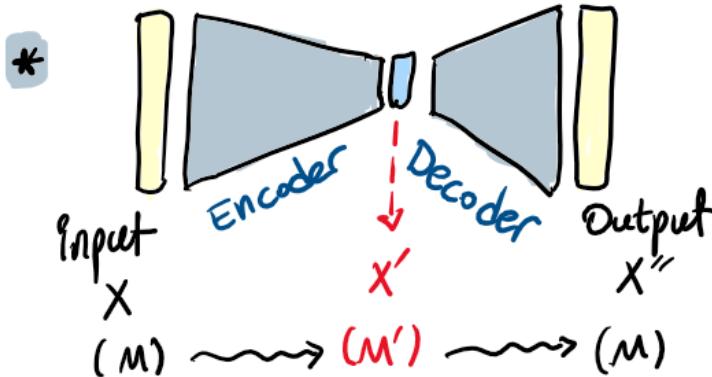
- \* It is probabilistic in nature



## Variational Encoders:

- \* "AE + Probability"
- \* introduced in 2013

AE:



Deterministic function

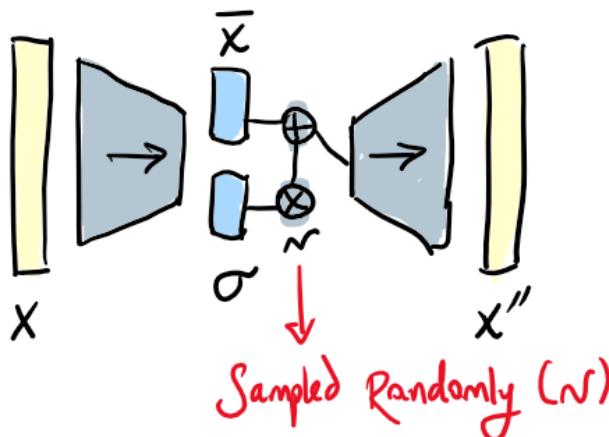
$$\alpha // x \rightarrow x' \rightarrow x''$$

will always be the same

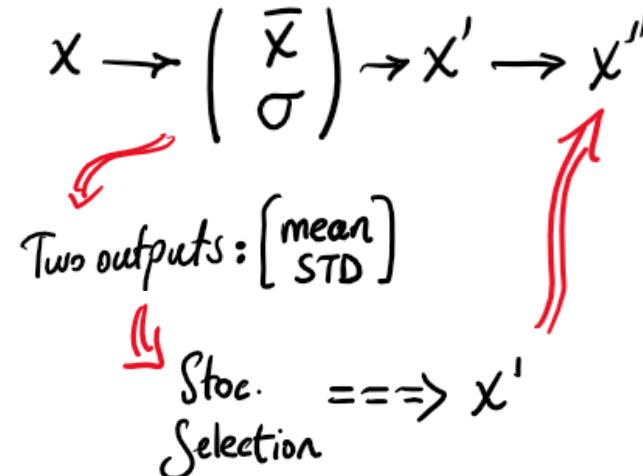
## Variational Encoders:

- \* "AE + Probability"
- \* introduced in 2013

VAE: Generative Models



Probabilistic function



Encoder

$$p_{\theta}(x'|x)$$

Decoder

$$p_{\theta'}(x''|x')$$

## Variational Encoders:

\* How do we enforce learning?

#1: creativity

#2: Backprop.

\* Loss := Reconstruction loss + "regularization,"

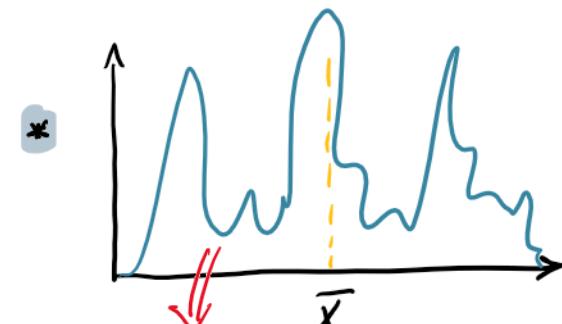
$$\|x - \hat{x}\|^2$$

~ as usual ~

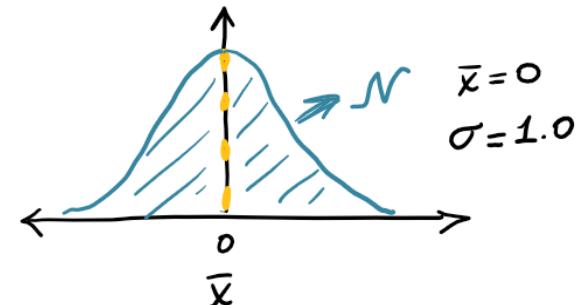
↓  
Confine creativity!

$$G = H - TS$$

\* Solution: Force it to follow a prior distribution.



if you let it free, it usually "cheats,"

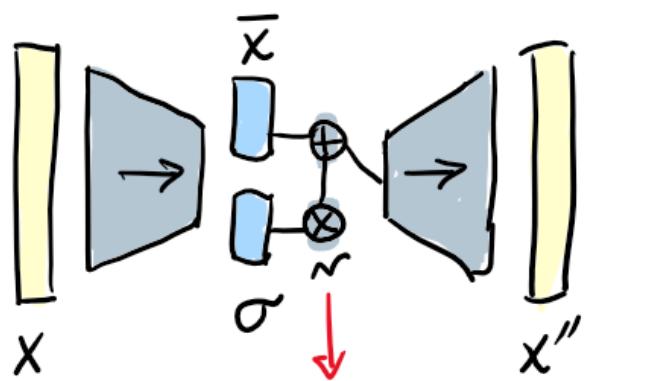


## Variational Encoders:

\* Gaussian Sampling



[enforced via cost function]



Sampled Randomly ( $\mathcal{N}$ )

Kullback - Leibler  
 (KL) Divergence  
 $(\mathcal{N} \& p')$

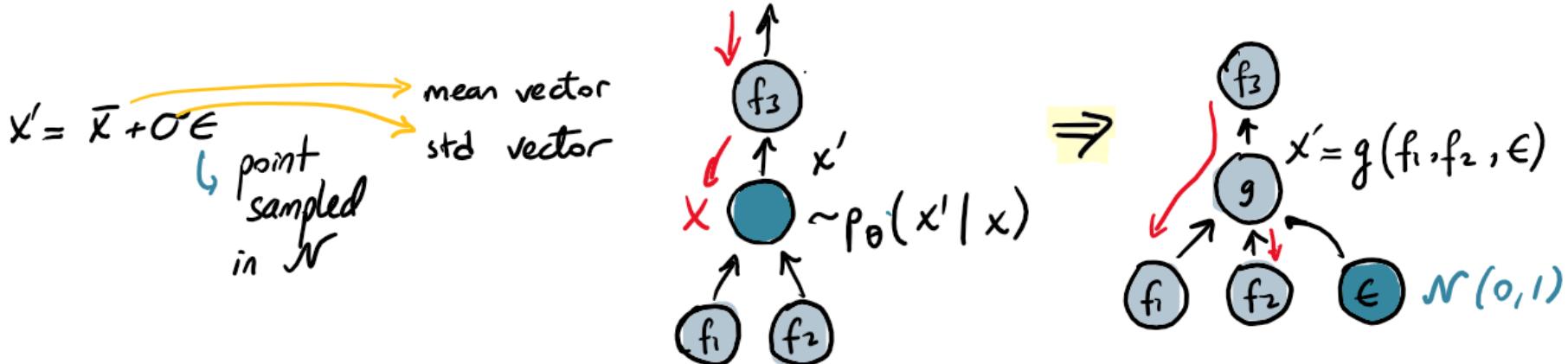
$$D_{KL} [p'(\bar{x}, \sigma) || N(0, 1)]$$

- $\mathcal{L}_{KL} = -\frac{1}{2} \sum_i 1 + \gamma_i - \exp(\gamma_i) - \bar{x}_i^2$
- $\gamma_i = \ln(\sigma_i^2)$



## Variational Encoders:

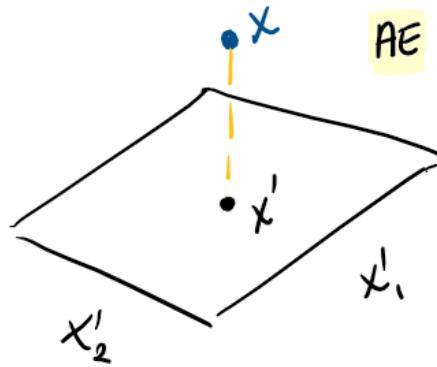
- \* How do we enforce learning?
  - #1: creativity
  - #2: Backprop.
- \* We cannot backprop. a stoc. layer  $\delta$



## Variational Encoders:

\* Gaussian Sampling

! It creates a cont'd. latent space.



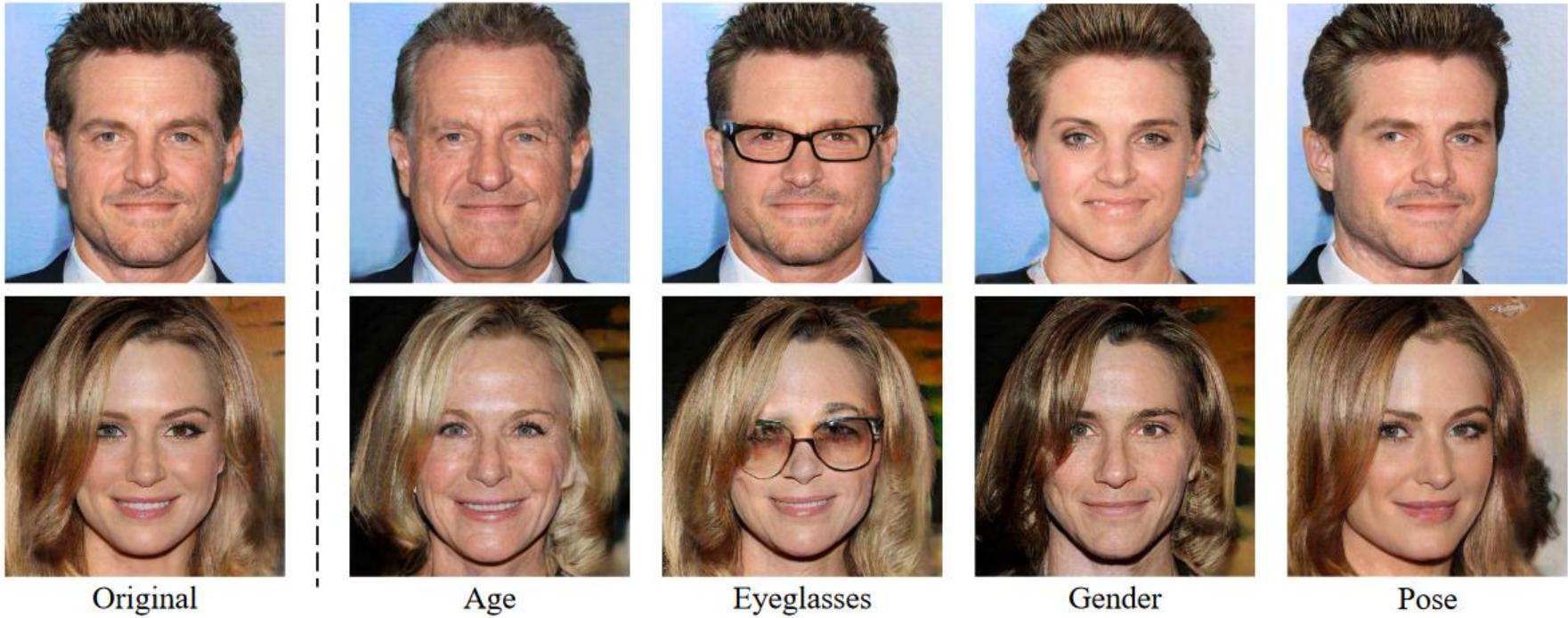
$[x'_1, x'_2]$   
 $\sim$   
 covariance  
 matrix  
 is  
 diagonal.

! You can perform semantic interpolation.

\*  $\left. \begin{array}{l} \text{Case 1: } [\bar{x}, \sigma], \\ \text{Case 2: } [\bar{x}, \sigma]_2 \\ \text{Case 3: } [\bar{x}, \sigma]_3 \end{array} \right\}$

" Case 4 := Case 1 - Case 2 + Case 3,  
 Case 5 :=  $\alpha$  (Case 1) +  $(1-\alpha)$  Case 2

# Interpolation $\Rightarrow$ Latent Space Arithmetics



Original

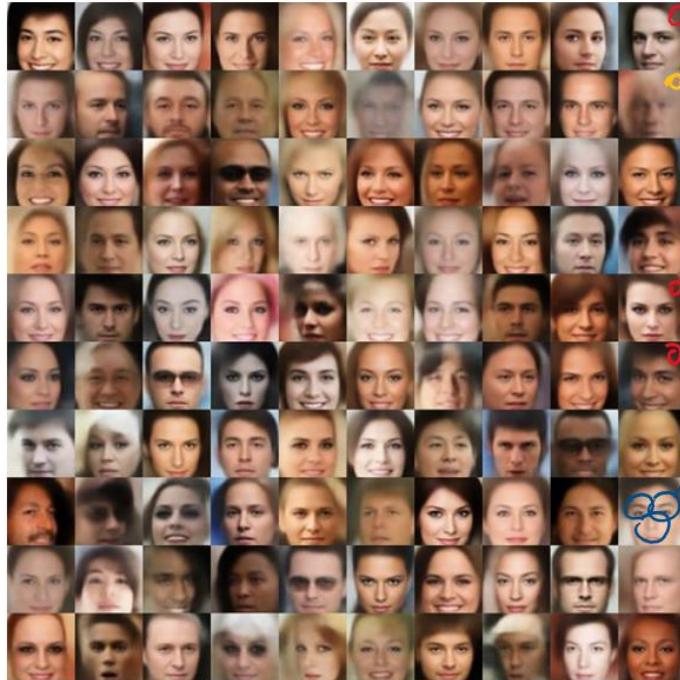
Age

Eyeglasses

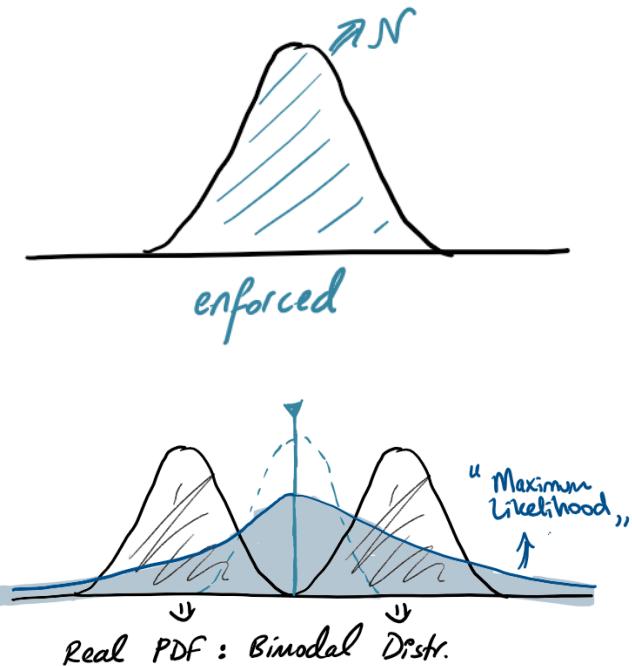
Gender

Pose

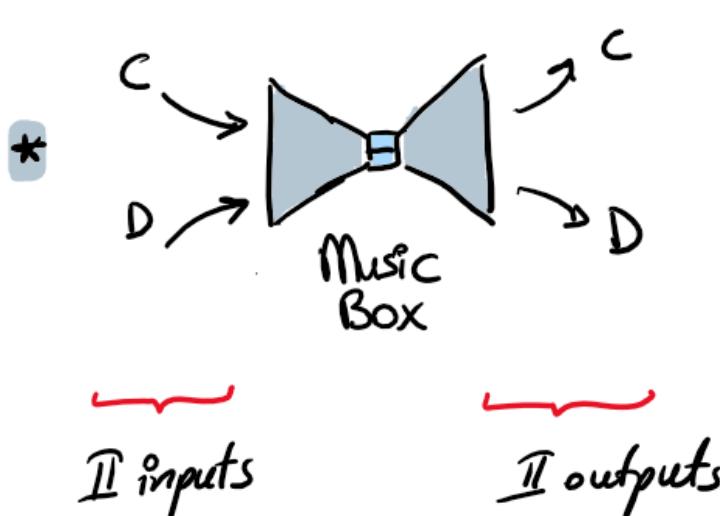
# Variational Encoders:



Blurry Edges  
 Blended in the background (safe pixel values)  
 Features are captured: eyes, nose, ...



## Our Strategy:



Case #1:  $C, D \xrightarrow{\text{VAE}}$ :

$$= [[\text{samples}], \text{features}]$$

Case #2:  $C, D \xrightarrow{\text{RNN}}$ :

$$= [\text{samples}, \text{seq.}, \text{features}]$$

\*

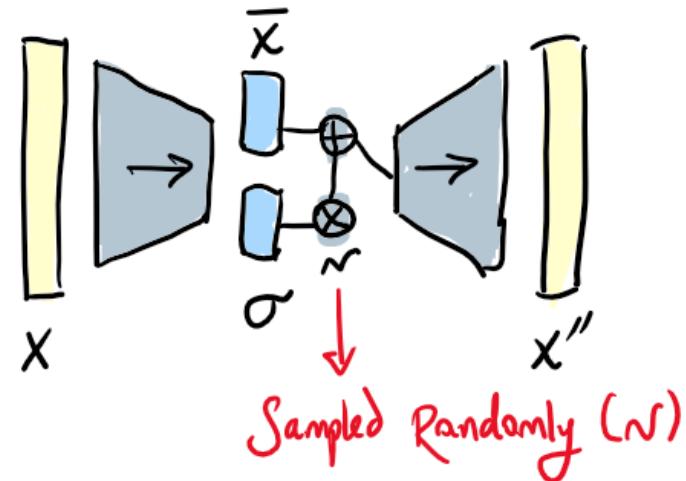
## Building our Graph. Network:

\* Prev. we used Sequential API in TF.

↳ Dense connections  
↳ Automated

\*

Custom. Structure  $\Rightarrow$  functional API

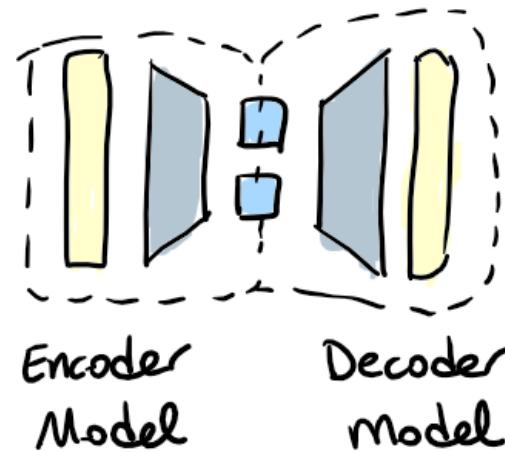
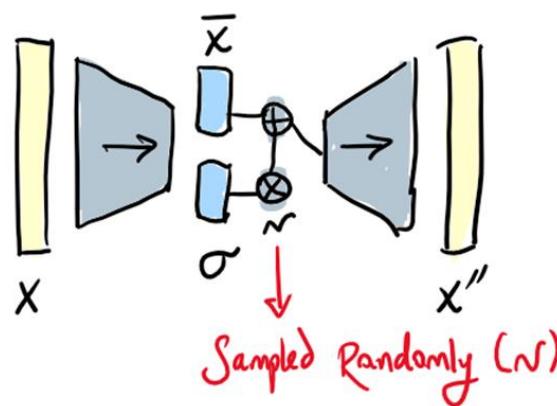


\*

## Building our Graph. Network:

\*

Custom. Structure  $\Rightarrow$  functional API

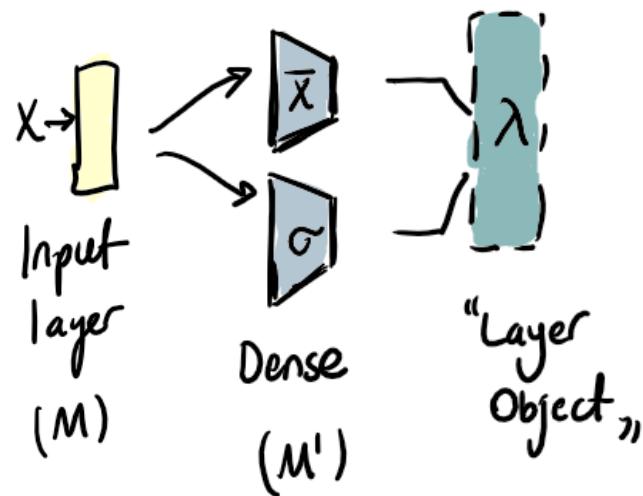


- Train  $\Rightarrow (E) + (D)$
- Gen  $\Rightarrow (D)$



## Building our Graph. Network:

Encoder :

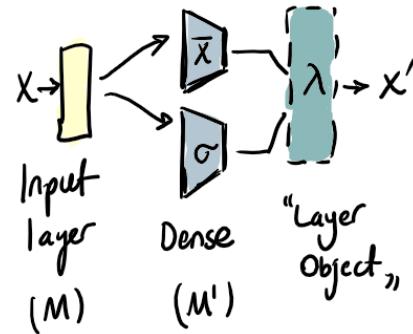
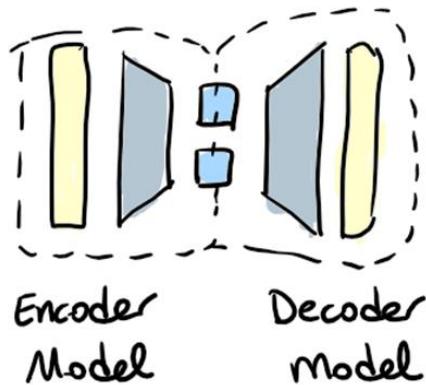


Sampling :

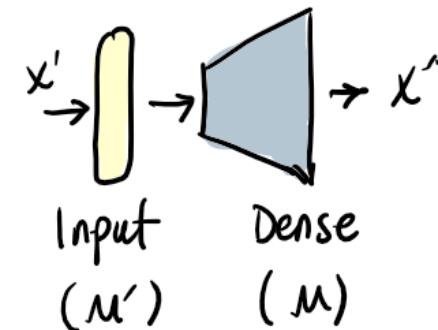
- $x' = \bar{x} + \sigma \epsilon \rightarrow x'$
- $\epsilon \leftarrow \mathcal{N}$

\*

# Building our Graph. Network:



+



$$\text{Encoder} = [$$

④ { input,  
mean, log var,  
lambda ]

+

$$\text{Decoder} = [$$

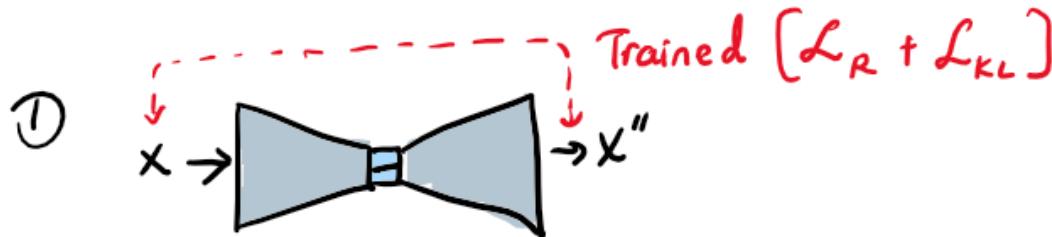
② { input,  
decoder ]

Ateliers & Saveurs in Montreal

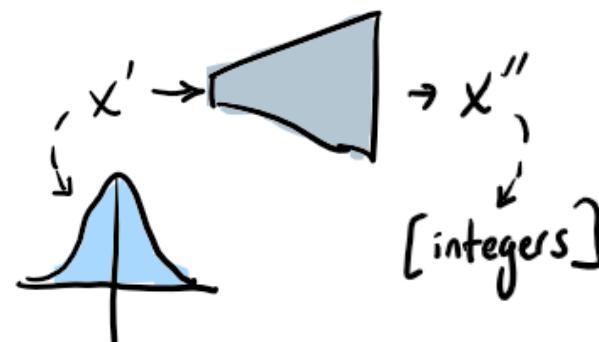


# colab

# Generating Music :



- i Generate noise via  $\mathcal{N}$
- ii Get  $x''$  [integers]
- iii (int)  $\rightarrow$  MIDI  $\rightarrow$  

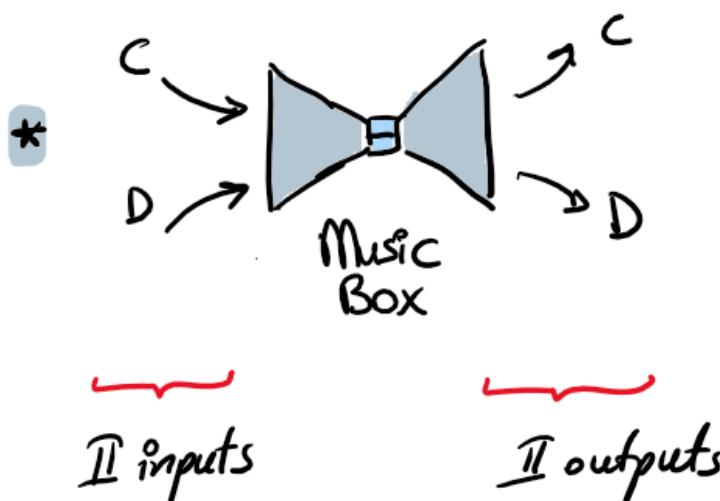


Ateliers & Saveurs in Montreal



# colab

## Our Strategy:

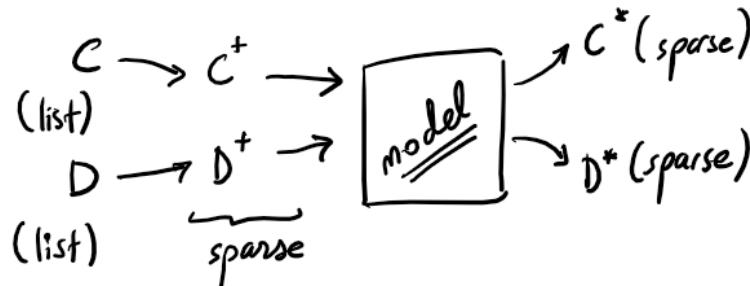


Case #2:  $C, D \rightarrow RNN$ :  
 (LSTM)  
 = [samples, seq., features]

\* LSTM  $\rightarrow$   $[x]$  data  
 $y$  labels  
 ↗ next chord  
 ↗ next duration

# Our Strategy:

\* Idea:



⚠ Memory-intensive

? Alternative  $\Rightarrow$  "Embeddings"

Eg Possibility =  $[0, \dots, 1000]$

i Encoding  $\Rightarrow$  "4"  $\Rightarrow$  Vector of size 1000  
 $[1 \rightarrow 1, 999 \rightarrow 0]$

ii Embedding  $\Rightarrow$  Integers  $\Rightarrow$  vector of fixed size (float)

Eg vector-size = 3

0  $\rightarrow$  [1.1, 2.6, 2.4]

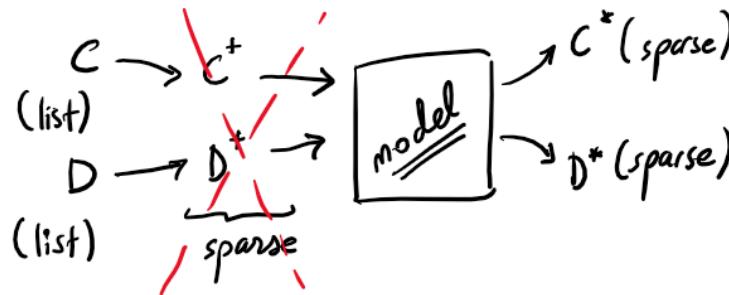
1  $\rightarrow$  [0.1, 4.2, 1.5]

:

} TF  
optimize it  
during the training

# Our Strategy:

## \* New Idea:



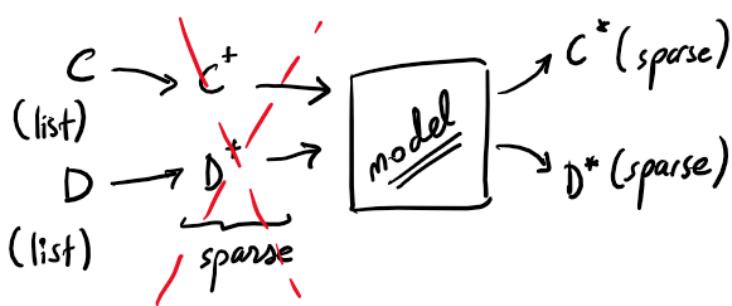
chord input

```
[164 164 149 164 182 164 157 149 114 92 108 108 90 114 108 95 92 164
 184 150 154 48 48 1 16 255 208 224 233 422 417 451]
```

duration input

```
[4 4 4 4 2 2 2 2 4 4 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 4]
```

## Our Strategy:



\* Still Sparse output?

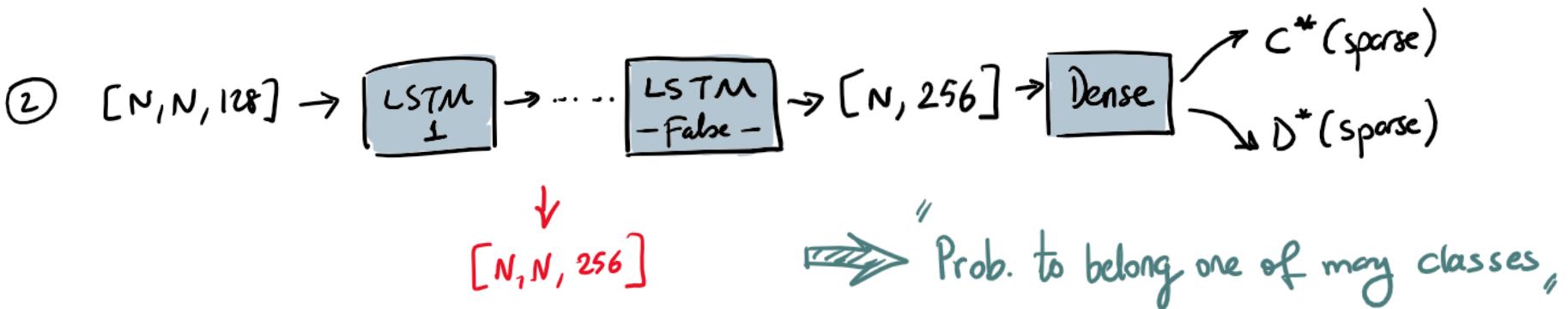
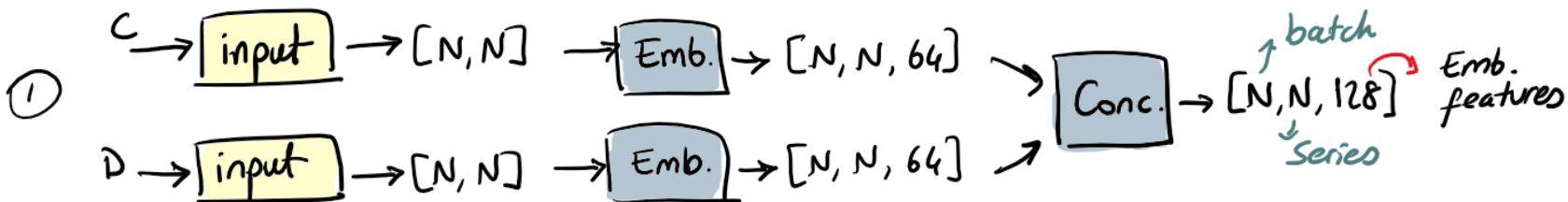
## chord output

## duration output

[0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0]

## Model Graph

Which (chord dur.) := Classification



Ateliers & Saveurs in Montreal



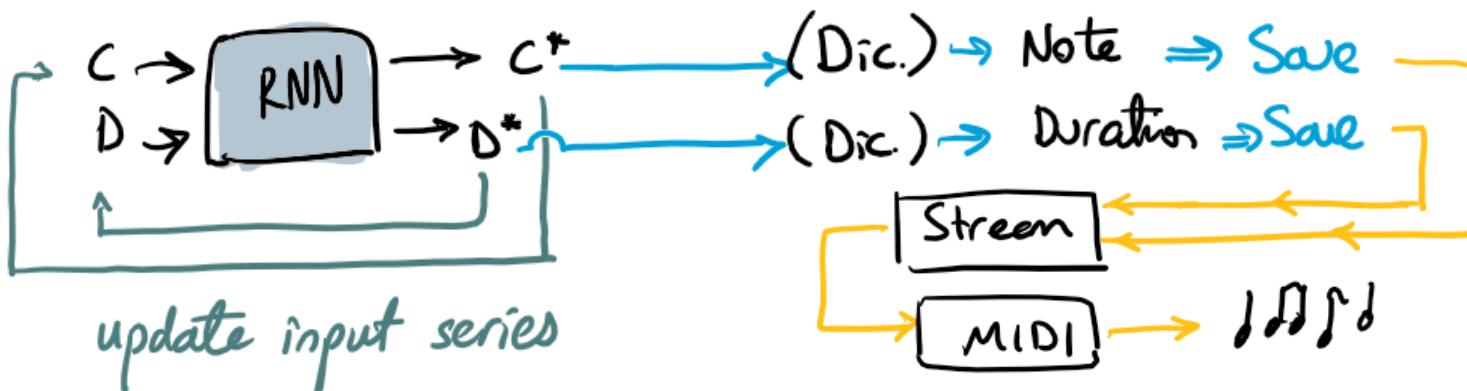
# colab

## Back-conversion :

- ①  $C \rightarrow$    $\rightarrow C^* \xrightarrow{\text{Pick most probable one}}$   
 $D \rightarrow$    $\rightarrow D^* \xrightarrow{\text{Pick most prob. one}}$
- ②  $C^* \rightarrow \{ \text{Int-to-Chord} \} \rightarrow \text{Note (str.)}$   
 $D^* \rightarrow \{ \text{Int-to-Duration} \} \rightarrow \text{Duration (str.)}$

## Back-conversion :

- ③ For generating long music;



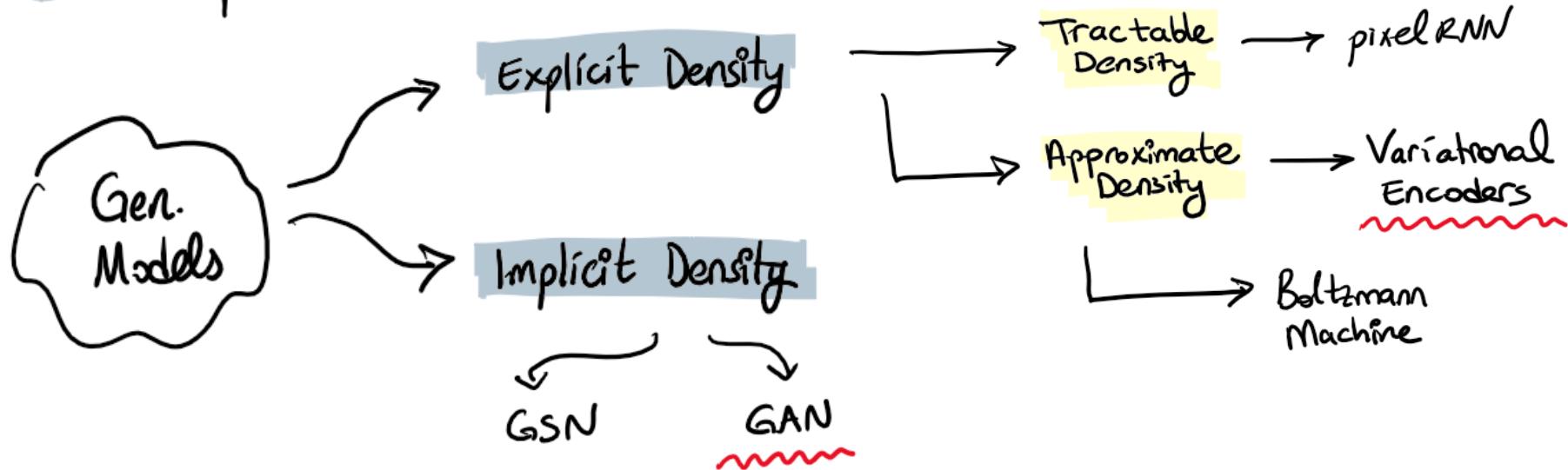
Ateliers & Saveurs in Montreal



# colab

# UL → Generative Models :

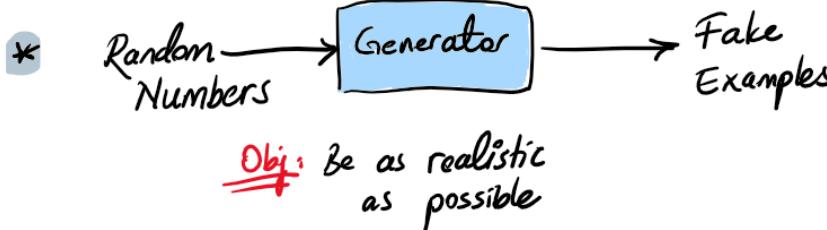
- \* It is probabilistic in nature



# Generative Adversarial Networks : GANs

- \* Generative  $\Rightarrow$  creating non-existing data
- \* Adversarial  $\Rightarrow$  Competitive dynamics (game-like)
- \* Network  $\Rightarrow$  Neural networks

\* GANS  
 (2014) 

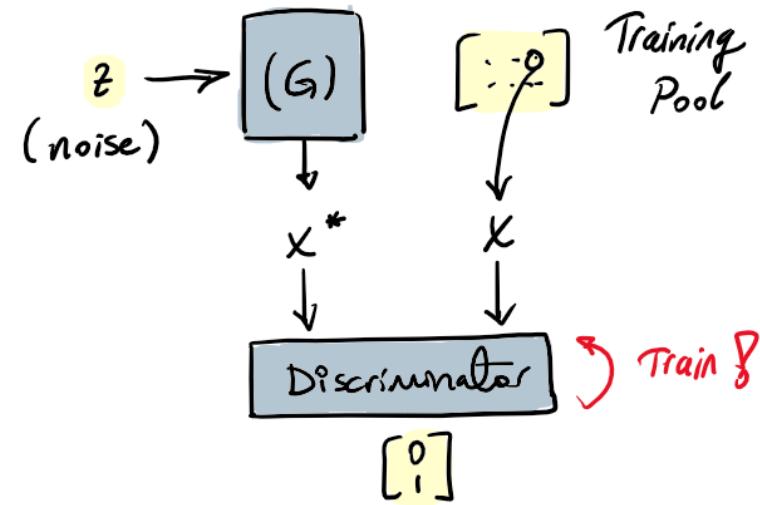


## Training Algorithm:

For each training do:

# Train ( $D$ ):

- (1) Take a random real example from Training data,  $x$
- (2) Get a fake example from Generator,  $x^*$
- (3) Use Discriminator to classify  $x$  &  $x^*$ .
- (4) Compute the class. error.
- (5) Backprop. error & update Discriminator trainable parameters.

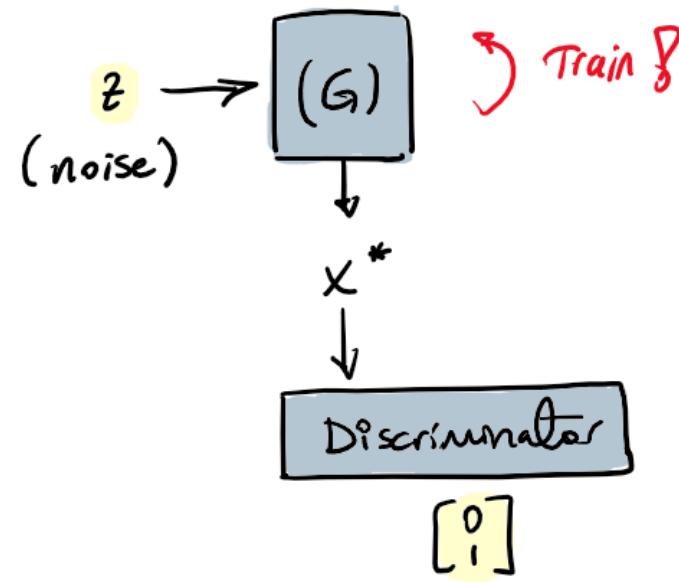


## Training Algorithm:

#Train ( $G$ ):

- (6) Generate a new fake  $x^*$ .
- (7) Use **Discriminator** to classify  $x^*$ .
- (8). Compute the error.
- (9). Update **Generator**'s trainable parameters via backprop.

end for



# Training GANs

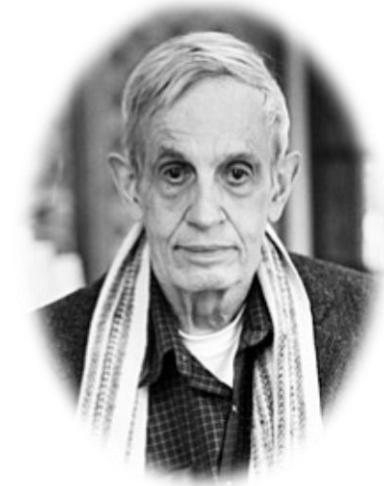
- \* In MLP, we have a clear goal & measure  
~~Ex~~ Minimize Cross-entropy loss.
- \* In GANs, two networks have competing obj. !  
 $(G) \uparrow ; (D) \downarrow \parallel (D) \uparrow ; (G) \downarrow$

# Training GANs

! Nash Equilibrium := Point where neither "player",  
can improve their situation

- $(G)$  := fakes are indistinguishable from Real data
- $(D)$  := at best randomy guess ( $F/R \Rightarrow 1$ )

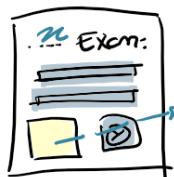
! In practice; ~impossible to achieve Nash Eq.



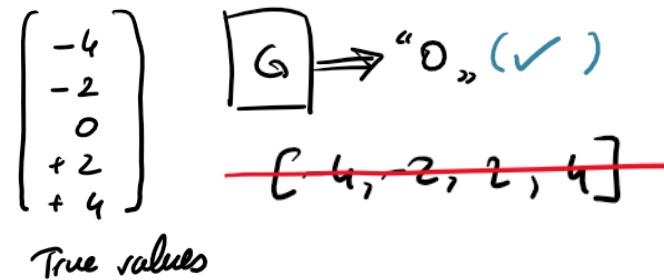
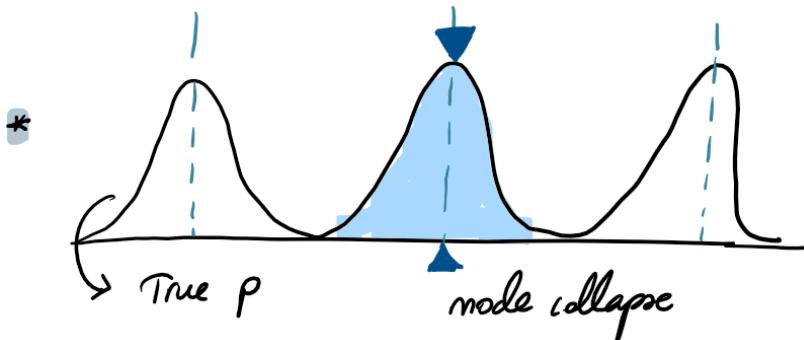
it still works ...

# Training GANs

## Training problems #1: Mode Collapse



“ Abracadabra 8 „ ~ I create as I speak 8 ”



# Training GANs

## Training Problem #2: Over generalization

- \* Models that should not exist, do exist.

\*  $\begin{bmatrix} -4 \\ -2 \\ 0 \\ +2 \\ +4 \end{bmatrix}$

$\boxed{G}$

$\rightarrow -2/3$

$\rightarrow 1/2$

...

[real numbers]

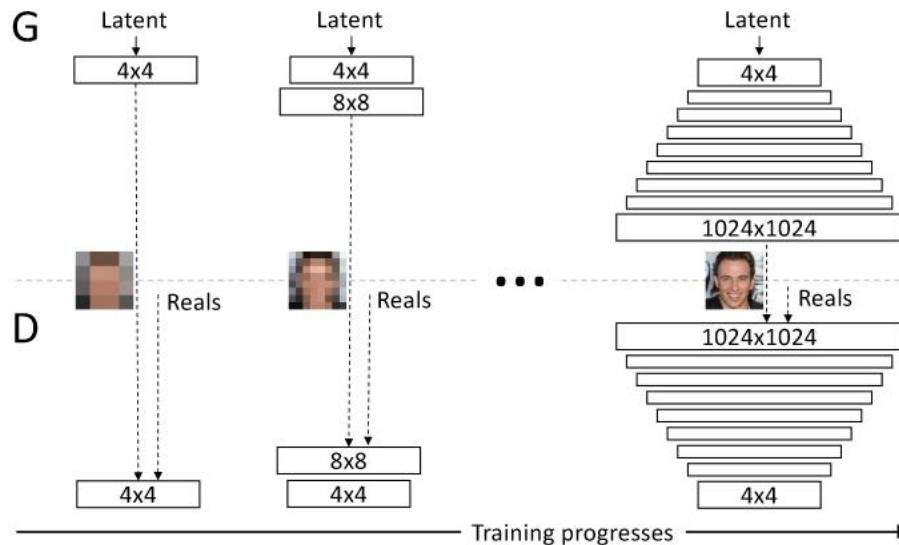
True values  
[integers]

\* Image generation



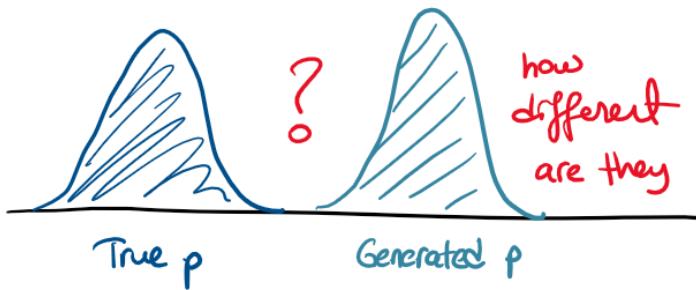
# Possible Remedies

## ① Growing the network gradually



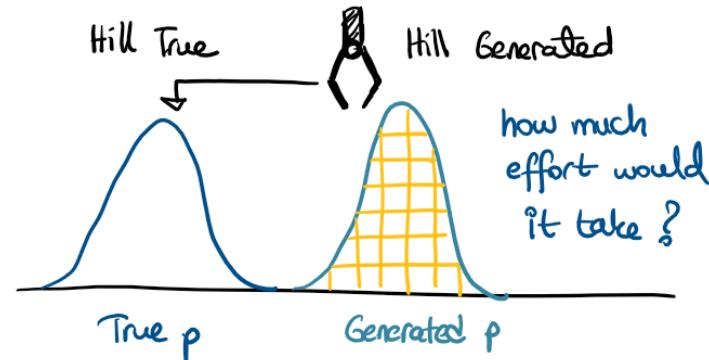
## Possible Remedies

### ② Alternative loss definitions $\Rightarrow$ Wasserstein GAN



\* Distance :=  $(\text{Similarity}) \begin{pmatrix} \text{TV distance} \\ \text{KL divergence} \\ \text{JS divergence} \\ \dots \end{pmatrix}$

$\Rightarrow$  Earth-mover Distance



## Possible Remedies

② Alternative loss definitions  $\Rightarrow$  Wasserstein GAN

\* Class.  $\Rightarrow [0, 1]$

$\hookrightarrow$  Train...  $\Rightarrow$  0.9971 || 0.0004  
 0.9969 || 0.00013  
 0.9981 || 0.00021  
 ...

  
*Vanishing gradient problem*

\* Stabilized learning

