# DrumGAN: Synthesis of Drum Sounds With Timbral Feature Conditioning Using GANs
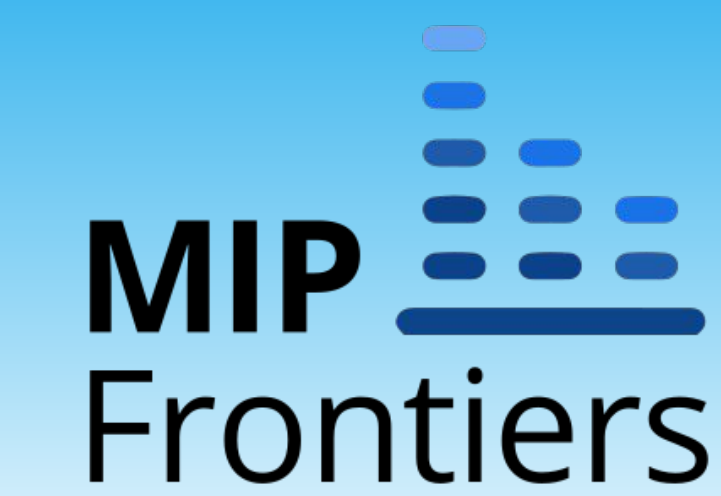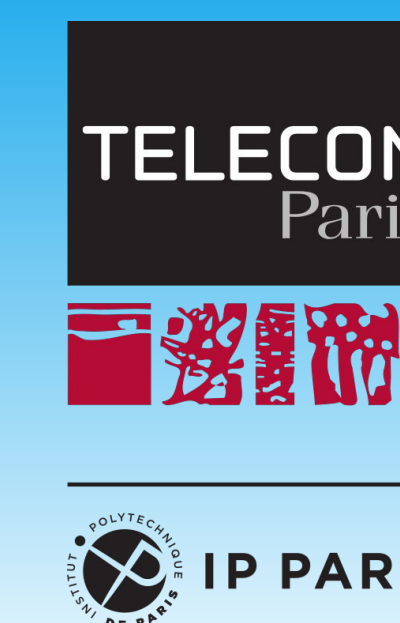
Javier Nistal, Stefan Lattner, and Gaël Richard

Sony CSL · TELECOM Paris · IP PARIS · MIP Frontiers

## Introduction

Audio synthesizers have **complicated parameters** with little **perceptual correspondence nor musical meaning**. Also, the type of **sounds** they can produce **are limited by the synthesis method** (e.g. additive, subtractive).

DrumGAN is a Progressive Growing GAN (PGAN) that can synthesize a wide variety of drum sounds and that enables steering the synthesis according to parameters that respond to human perception.
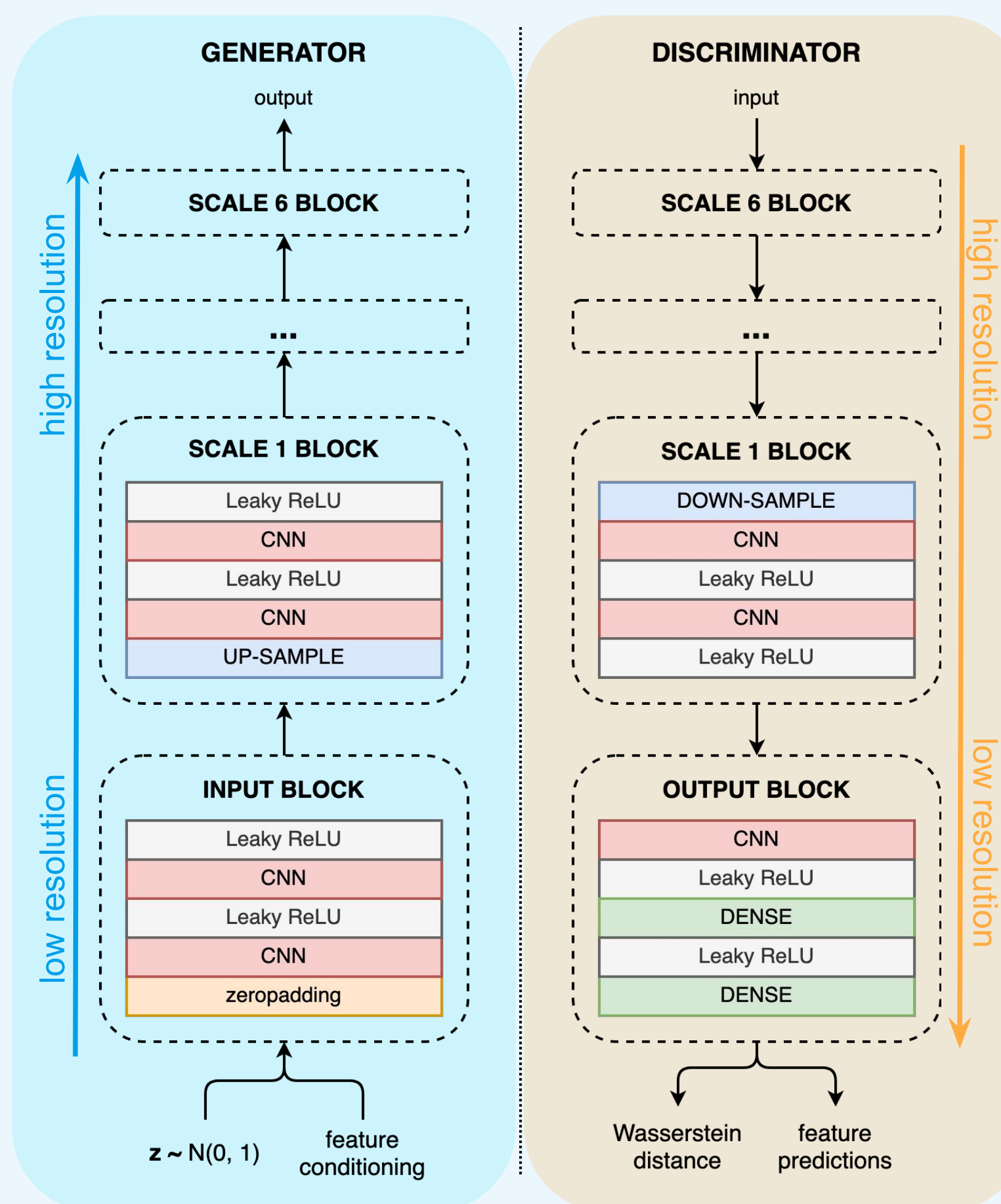
## Dataset

- ❖ **~300k** one-shot, 1s-long and aligned audio samples
- ❖ **Kicks (K), Snares (S) and Cymbals (C)** classes
- ❖ **16kHz sampling-rate**
- ❖ **90/10%** train-validation split
- ❖ **Complex STFT** representation
  - ➢ window size: 2048
  - ➢ hop size: 512

## Audio-Commons Features

- Audio Commons perceptual models → high-level timbral features of the sound
- Human ratings given to sounds from Freesound
- Linear regression models of spectral and temporal low-level features (e.g., spectral centroid, dynamic-range)
- All features are normalised to the range [0-1]

Boominess · Brightness · Warmth · Depth · Hardness · Sharpness · Roughness

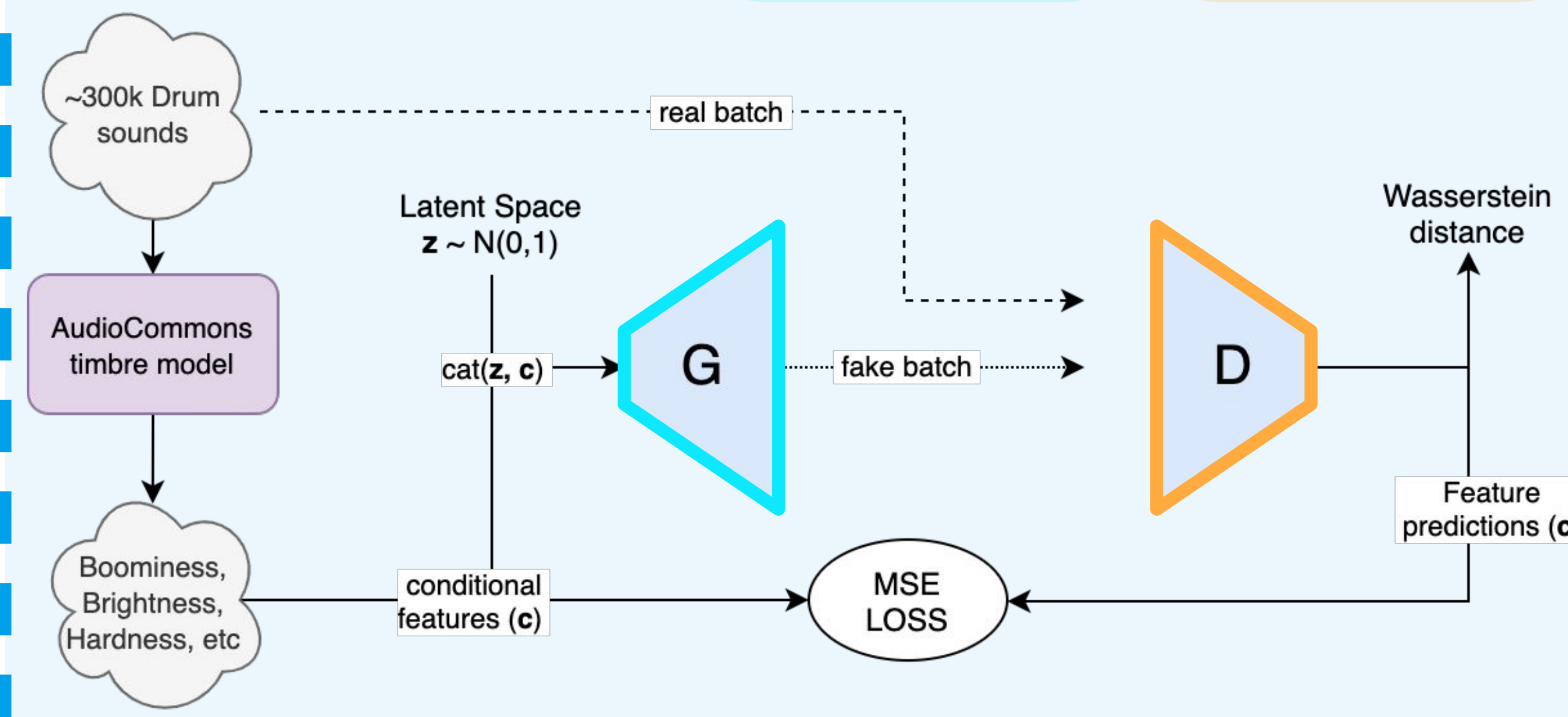## Architecture & Training Procedure

### Progressive Growing GAN

**Input Block**:
$$[z, c] \rightarrow [batch, ch, f_{s0}, t_{s0}]$$

**Scale blocks**: progressively added to the network while training. CNN out channels for each scale:
$$[256, 128, 128, 128, 128, 64]$$

**Training**:
- ❖ 1.1M iterations (~200k i/scale)
- ❖ batch-size: [30, 30, 20, 20, 12, 12]
- ❖ Adam optimizer
- ❖ learning rate: 1e-3.

**GENERATOR**

output

SCALE 6 BLOCK

...

SCALE 1 BLOCK
- Leaky ReLU
- CNN
- Leaky ReLU
- CNN
- UP-SAMPLE

INPUT BLOCK
- Leaky ReLU
- CNN
- Leaky ReLU
- CNN
- zeropadding

$z \sim N(0, 1)$ · feature conditioning

**DISCRIMINATOR**

input

SCALE 6 BLOCK

...

SCALE 1 BLOCK
- DOWN-SAMPLE
- CNN
- Leaky ReLU
- CNN
- Leaky ReLU

OUTPUT BLOCK
- CNN
- Leaky ReLU
- DENSE
- Leaky ReLU
- DENSE

Wasserstein distance · feature predictions

high resolution · low resolution

~300k Drum sounds

Latent Space $z \sim N(0,1)$

real batch

AudioCommons timbre model

$cat(z, c)$ → G → fake batch → D

Wasserstein distance

Boominess, Brightness, Hardness, etc

conditional features (c)

Feature predictions (c')

MSE LOSS

## Results

### Attribute coherence

A specific feature $f_i$ is set to 0.2, 0.5, and 0.8, keeping the other features and $z$ fixed. The outputs of **G** ($fx_{i0.2/0.5/0.8}$) are evaluated with the Audio Commons Models. Three conditions are examined:

**E1: $fx_{i0.2} < fx_{i0.5}$    E2: $fx_{i0.5} < fx_{i0.8}$    E3: $fx_{i0.2} < fx_{i0.5}$**

| | E1 (drumGAN) | E2 (drumGAN) | E3 (drumGAN) | E1 (baseline) | E2 (baseline) | E3 (baseline) |
|---|---|---|---|---|---|---|
| brightness | 0,74 | 0,71 | 0,7 | 0,99 | 0,99 | 1 |
| hardness | 0,64 | 0,64 | 0,62 | 0,64 | 0,65 | 0,59 |
| depth | 0,79 | 0,72 | 0,74 | 0,94 | 0,65 | 0,94 |
| roughness | 0,76 | 0,68 | 0,67 | 0,63 | 0,59 | 0,57 |
| boomines | 0,8 | 0,74 | 0,77 | 0,98 | 0,82 | 0,98 |
| warmth | 0,76 | 0,71 | 0,71 | 0,92 | 0,79 | 0,91 |
| sharpness | 0,84 | 0,82 | 0,82 | 0,63 | 0,77 | 0,45 |

**DISCUSSION**
*baseline* yields high accuracies for features describing the global frequency distribution (e.g., brightness, depth). *drumGAN* performs better on features describing complex frequential relationships (e.g., roughness, sharpness).

### Scores & Distances

**Inception Score** (IS), **Kernel Inception Distance** (KID), **Fréchet Audio Distance** (FAD)

→ *real data*: scores and distances on real data
→ *uncond*: unconditional drumGAN
→ *drumGAN_train*: conditioned on training labels
→ *drumGAN_val*: conditioned on validation labels
→ *drumGAN_rand*: conditioned on random labels
→ *baseline*: UNet conditioned on real labels

Legend: drumGAN_rand · drumGAN_train · uncond · baseline · drumGAN_val · real data

| | IS | KID | FAD |
|---|---|---|---|
| | 2,09 | 1,36 | 0,7 |
| | 2,18 | | 0,76 |
| | | 0,39 | |
| | 2,19 | 1,07 | 1 |
| | 2,18 | 1,45 | 3,09 |
| | 2,26 | 0,55 | 0,70 |
| | | 0,05 | 0 |

**DISCUSSION**
The **IS** of *drumGAN* is **close to that of real data** → outputs are assignable into {K, S, C}. *uncond* yields worse KID and FAD → the **features help generating more realistic samples**. *drumGAN* outperforms the *baseline* metrics for all conditional settings.

**SOUND EXAMPLES** ➡ [QR] · **PAPER** ➡ [QR] · [GitHub] ➡ [QR]