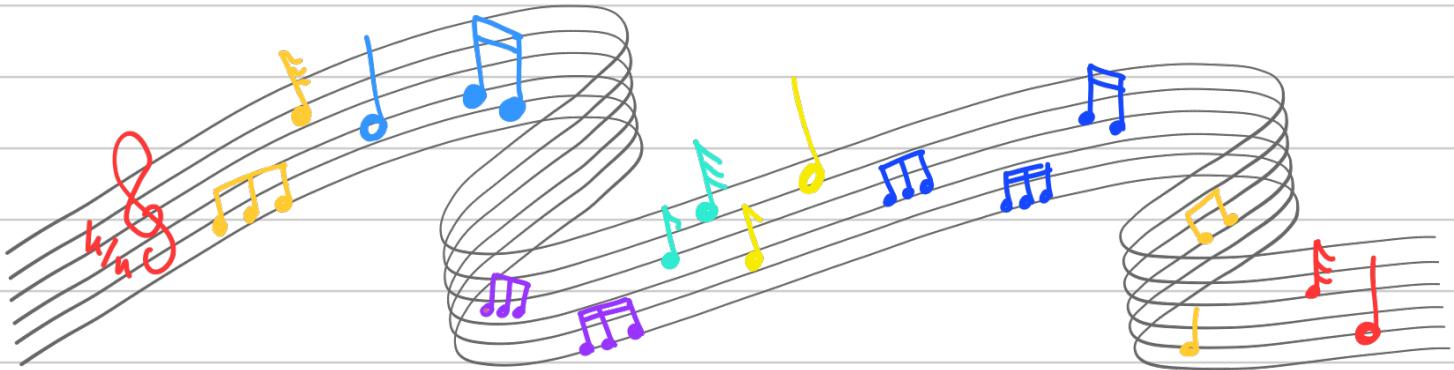


MUSIC GENERATION



WHAT

HOW

WHO *me*

AUTOENCODERS

RECURSIVE NETWORKS

TRANSFORMERS

THE MAESTRO DATASET



1250+



POLYPHONIC



THE MIDI FORMAT

1983

!!

PITCH ↑ ↓ 

midi → START 

DURATION ← → 

ALONG WITH OTHER DATA . . .

INSTRUMENTS

VOLUME

EFFECTS

META-DATA

CHANNELS

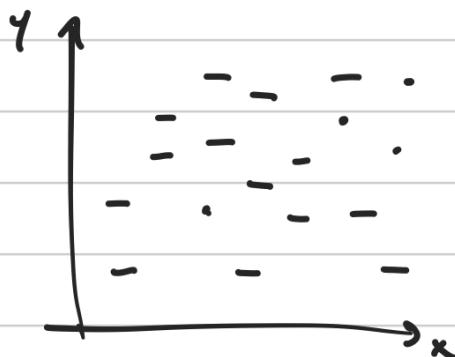
BUT WE KEEP IT SIMPLE !



AUTOENCODER

TRAINING DATASET

THE PIANO ROLL REPRESENTATION



where:

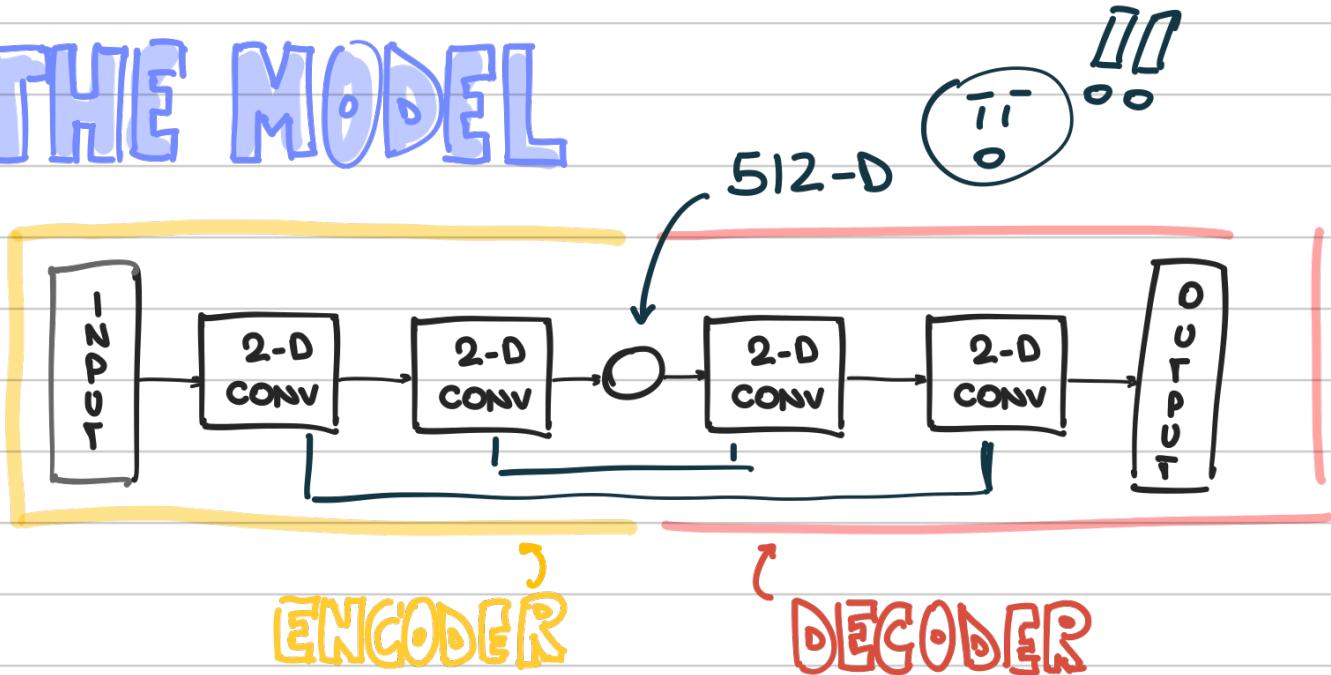
- x = time
- y = pitch



SIMPLE!

FRAGMENTS OF 50ms (time sampling)

THE MODEL



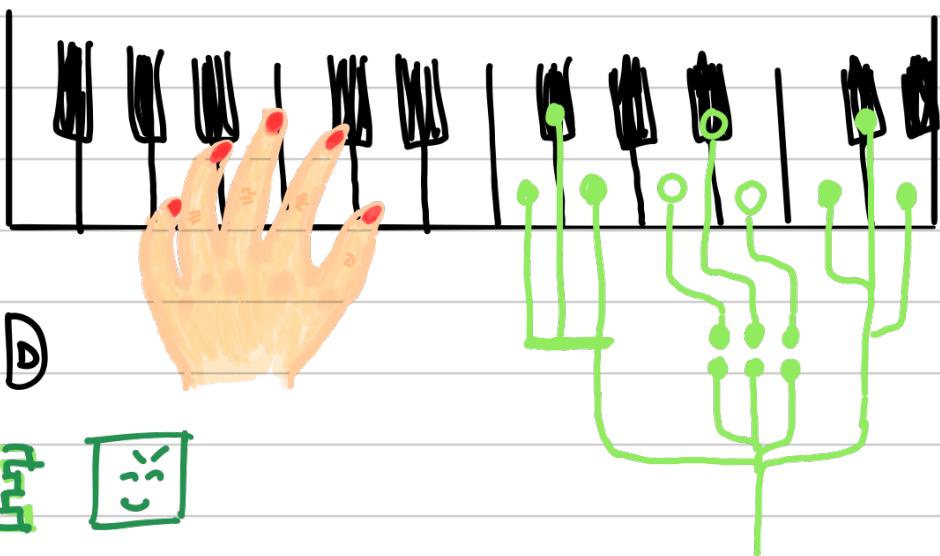
OUTPUT IS ACTIVATION OF PITCHES

\Rightarrow THRESHOLD IS SET TO GEN.

- RELU ACTIVATIONS
 - SKIP CONNECTIONS
 - KERNELS 3×3
 - CONV-2D 32 AND 64 FILTERS
-

MUSICA MAESTRO!!!

LISTEN...



... AND FIND

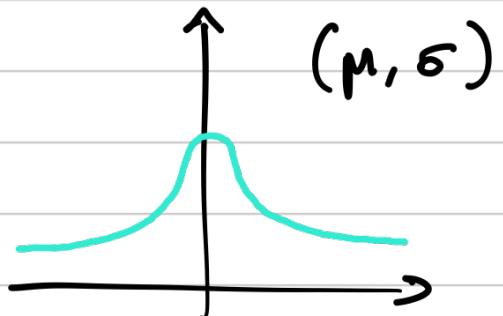
THE FAKE ☺

WELL... NOT WHAT WE WOULD CALL "NICE"

MAYBE A BETTER LATENT SPACE...

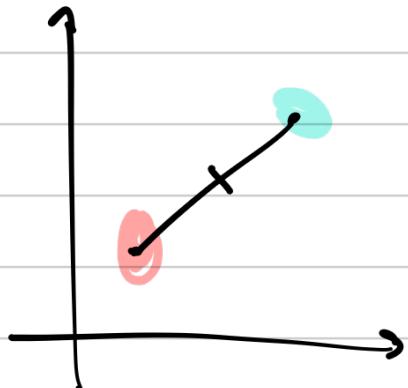
- Random value of latent space \Rightarrow meaningful output
- Give the decoder "generative" capabilities

• KL-DIVERGENCE



MAYBE INTERPOLATING CLASSES:

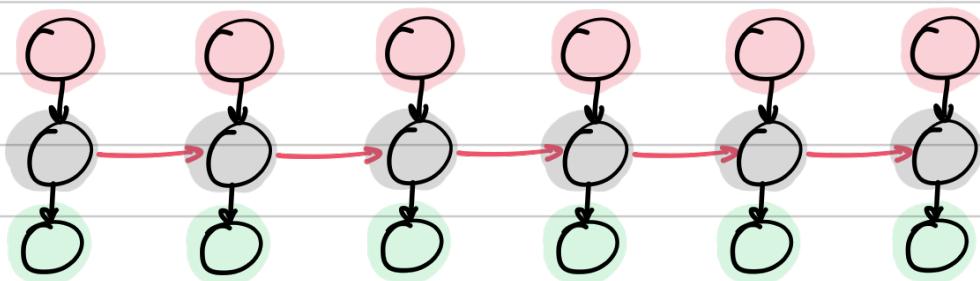
- Higher chance of a meaningful output!
- Both for AE and VAE



SOME CONSIDERATIONS

- PIANO-ROLL DOESN'T WORK
- THE AUTO-ENCODER IS NOT 'SPECIALIZED'

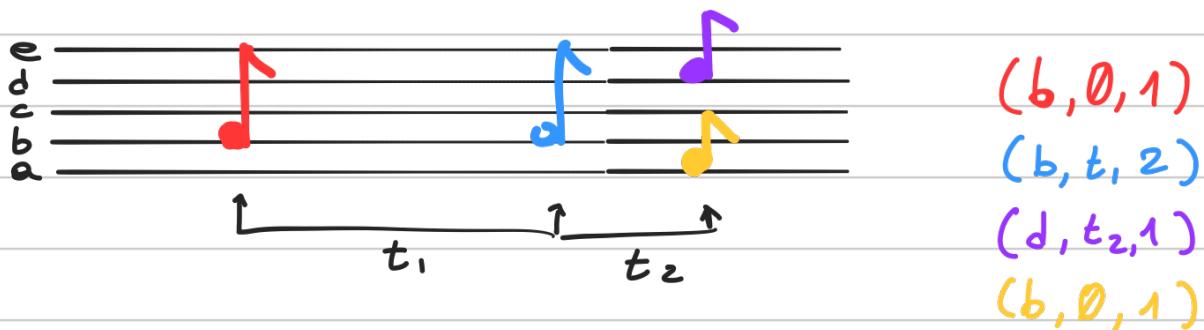
RECURSIVE NETWORKS



MUSIC IS RELATED TO TIME !!

TRAINING DATASET

- TUPLES: (PITCH, STEP, DURATION)



- SEQUENCE LENGTH: 25
 - INPUT SHAPE: (25, 3)
- ↓
- OUTPUT SHAPE: (1 , 3)

THE MODEL

STEP DURATION PITCH
↓ 1.0 * ↓ 1.0 * ↓ 0.05 *

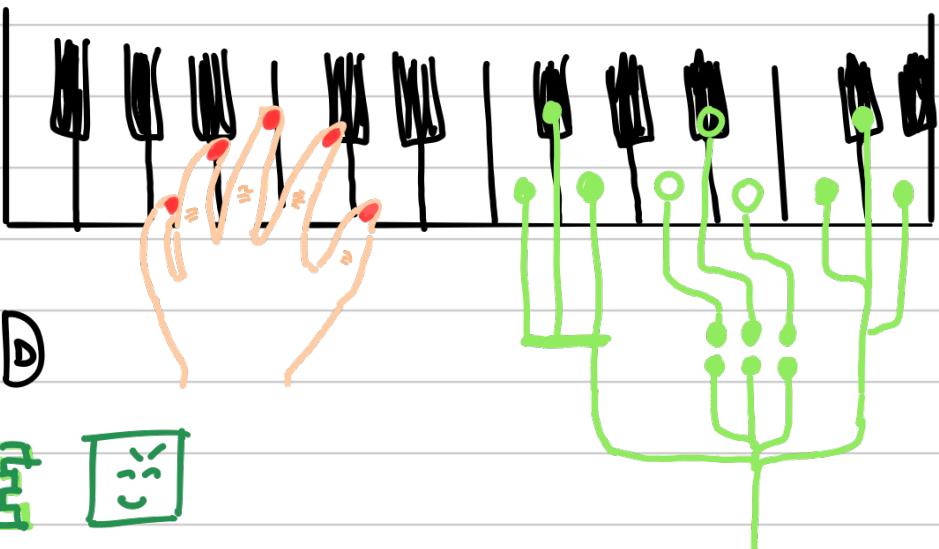
POSITIVE PRESSURE SPARSE CATEGORICAL
 $\text{mse} + 10 \cdot \max(-y, 0)$ CROSS ENTROPY
ONLY IF y IS NEGATIVE

• LOSS WEIGHTS *

• HOT OR COLD: TEMPERATURE
ACTS ON PITCH

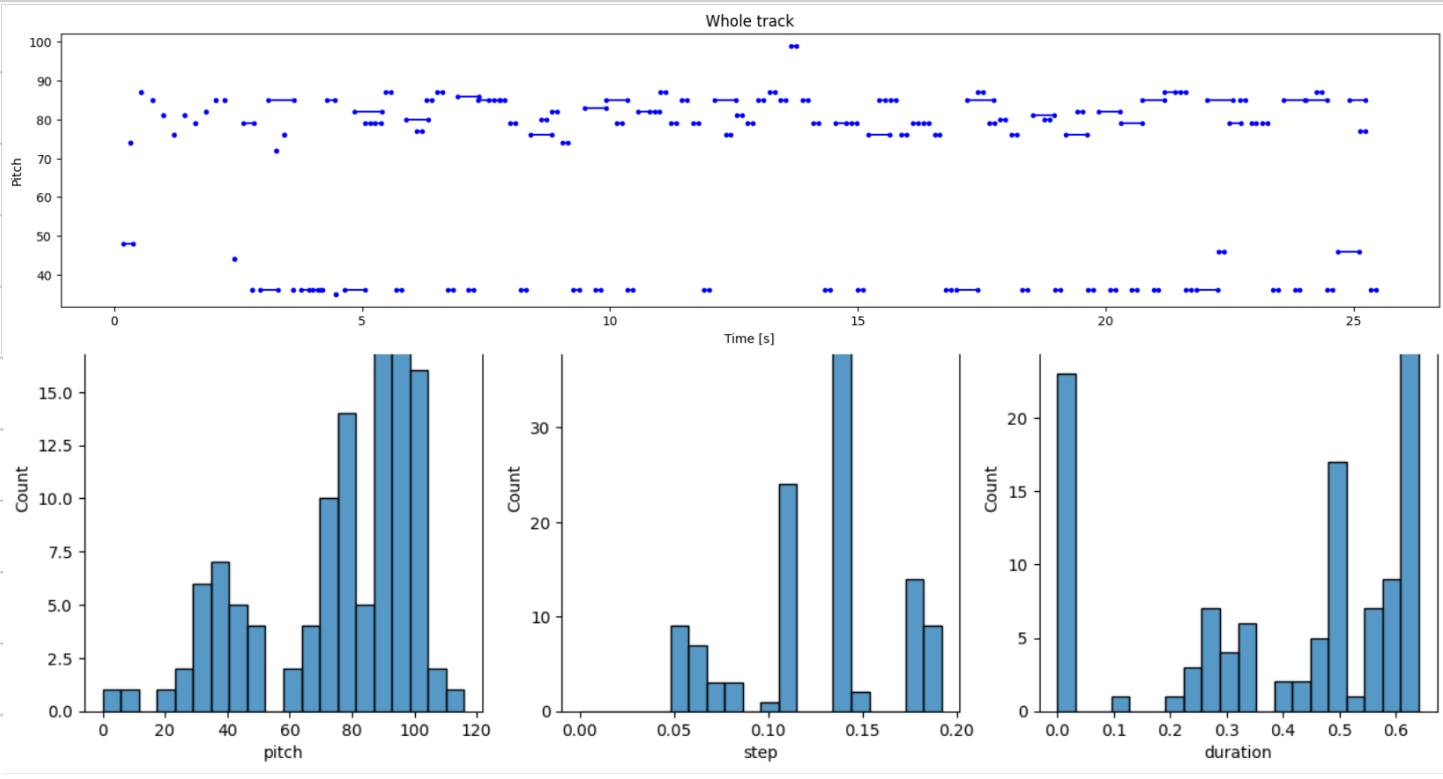
MUSICA MAESTRO!!!

LISTEN...



... AND FIND

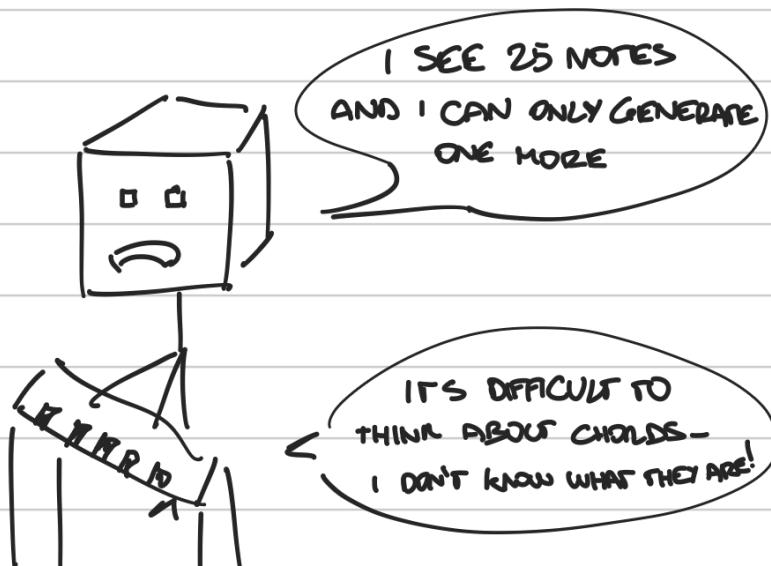
THE FAKE ☹



CONSIDERATIONS

- MODEL CREATES PAUSES
- SIMPLE ARMONIES
- BUT IT FEELS A BIT MONOPHONIC...

CHORDS ARE NOT ENFORCED ↗



SOLUTIONS:

- EMBED CHORDS AS INPUTS
- GENERATE MORE THAN 1 NOTE ON EACH ITERATION

TRANSFORMERS

TRAINING DATASET

TOKENIZER: 340+ EVENTS 😮

TRAIN 90% TEST 10%

NOTES BATCHING → SEQUENCING

THE MODEL

"MUSIC TRANSFORMER" 12 SEP 2018

↑ GENERATING MUSIC USING
LONG-TERM STRUCTURE

• RELATIVE ATTENTION!

- STD. RELATIVE POSITION ENCODING
- ENABLE LONGER SEQUENCES - MINUTES

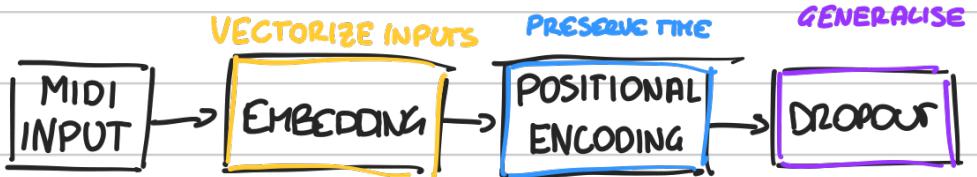
• TEMPORAL SENSITIVITY

- SINE POSITION ENCODING

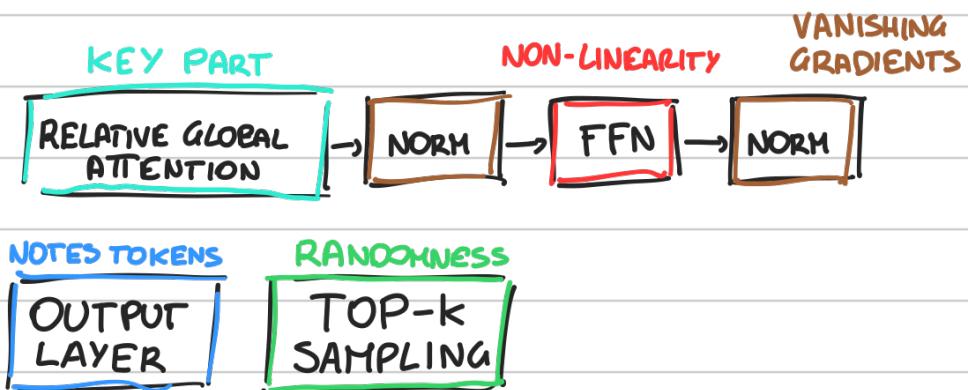
• TOP-K SAMPLING

- CONSTRAINED RANDOMNESS

PREPROCESS:



DECODER:



GENERATION:

¹ ↕ DURING TRAINING ↘²

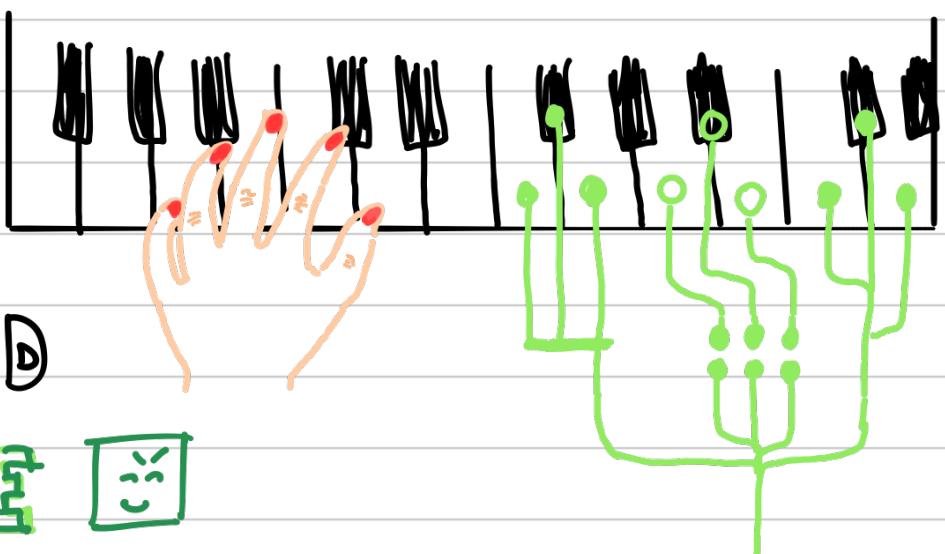
RELATIVE ATTENTION

MASKED CATEGORICAL CROSS ENTROPY

LISTEN...

... AND FIND

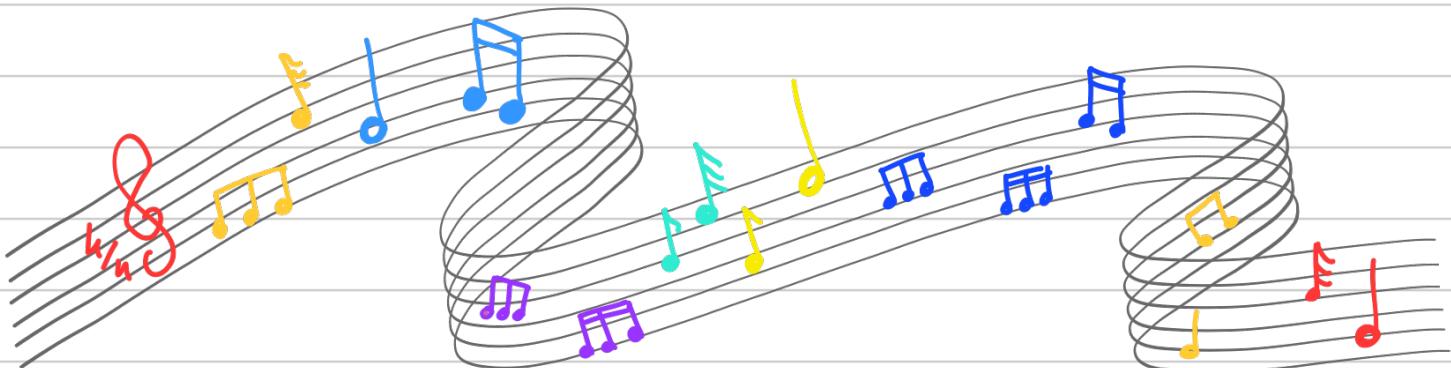
THE FAKE ☺



CONSIDERATIONS

- CHORDS ✓
 - ARMONIES ✓
 - PAUSES ✓
 - TOKENIZATION WORKS! 😊
-

CONCLUSION



A JOURNEY THROUGH:

NOT THE RIGHT
APPROACH

AEs

A JUMP FORWARD
ON TEMPORAL DATA

RNNs

BETTER RESULTS
ON DIFFERENT
SCENARIOS

VAEs

HIGH POTENTIAL
AND BEST RESULTS

TRANSFORMERS

AND I HAD A LOT OF FUN!

AI according to the news:



AI in real life:

