



FRANCO CASPE
& ANDREA MARTELLONI

DEEP LEARNING FOR DSP ENGINEERS:

*CHALLENGES AND TRICKS TO BE PRODUCTIVE
WITH AI IN REAL-TIME AUDIO*

Who we are

centre for digital music



Andrea Martelloni

PhD Candidate at Artificial Intelligence and Music CDT, Queen Mary University of London.
DSP engineer, guitar player/songwriter.



Franco Caspe

PhD Candidate at Artificial Intelligence and Music CDT, Queen Mary University of London. Electronic & Computer Vision Engineer, guitar player.



What this talk is not about

Not about Deep Learning frameworks.

Not a hands-on DL tutorial.

Not about model optimisation for efficiency.

Not about (offline) Music Information Retrieval (MIR).

Not about LLMs.

Not about symbolic music.

What this talk is about

We want to make DL for Audio **more approachable** from your point of view.

We will focus on the use of DL for **real-time musical audio**, e.g. a plugin in a DAW.

We will build upon **DSP engineering** knowledge.

We will present **our own experiences** in building DL-based audio tools.

Outline

Uses of DL in audio

Basic working principles of DL

Steps of the design of a DL model through our experiences:

- The HITar

- Bessel's Trick

Takeaways and lessons learned

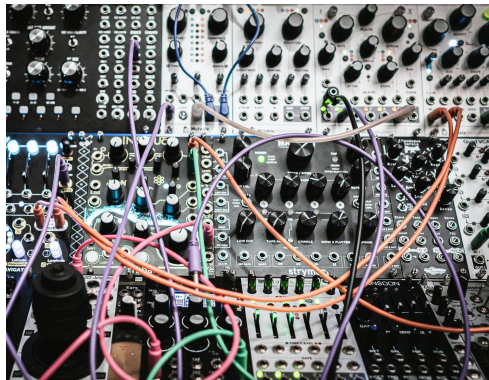
Uses of Deep Learning in Audio

What is Deep Learning good for?

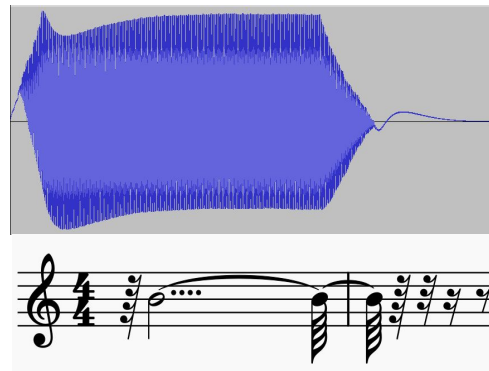
Complex
Problem



Complex
Solution



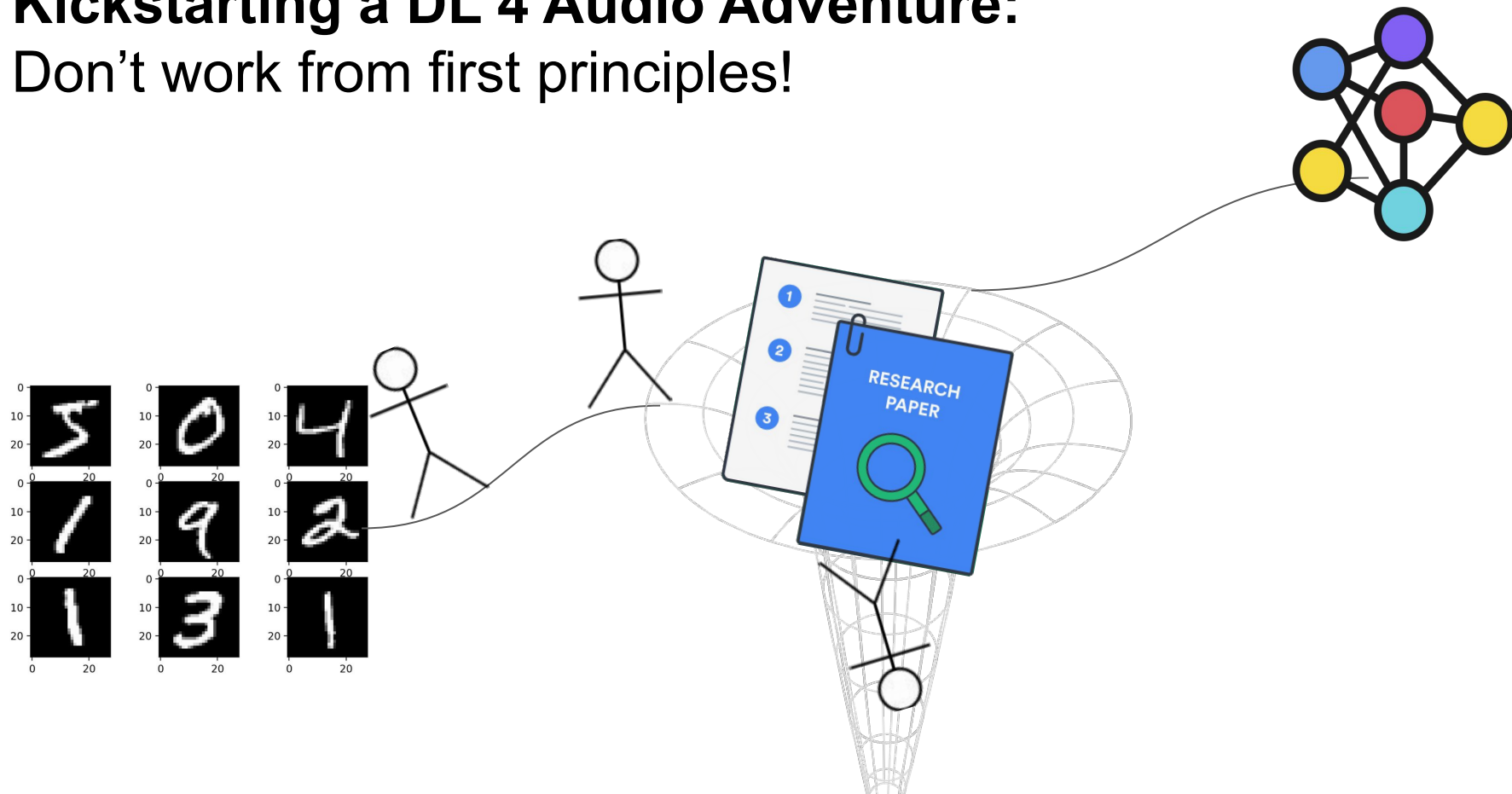
Brittle
Solution



We optimize a system with **data**, without explicitly defining **how** to fit that data.

Kickstarting a DL 4 Audio Adventure:

Don't work from first principles!



How to start making sense of Deep Learning?

- We are going to use **informed intuition** from DSP engineering knowledge.
- We are going to walk you through our **projects** and **experiences**,
- Recalling how we used our **DSP background** to **navigate** the complexities of DL.

Bessel's Trick

FM Tone Transfer



In Gain

Pitch



RMS



Envelope
Model

TWINCLE

+

Algorithm



3

Oscillators

1

2

3

4

5

6

Coarse



Fine



Boost



$f = 2$



Out Gain



Status: Ready to play!



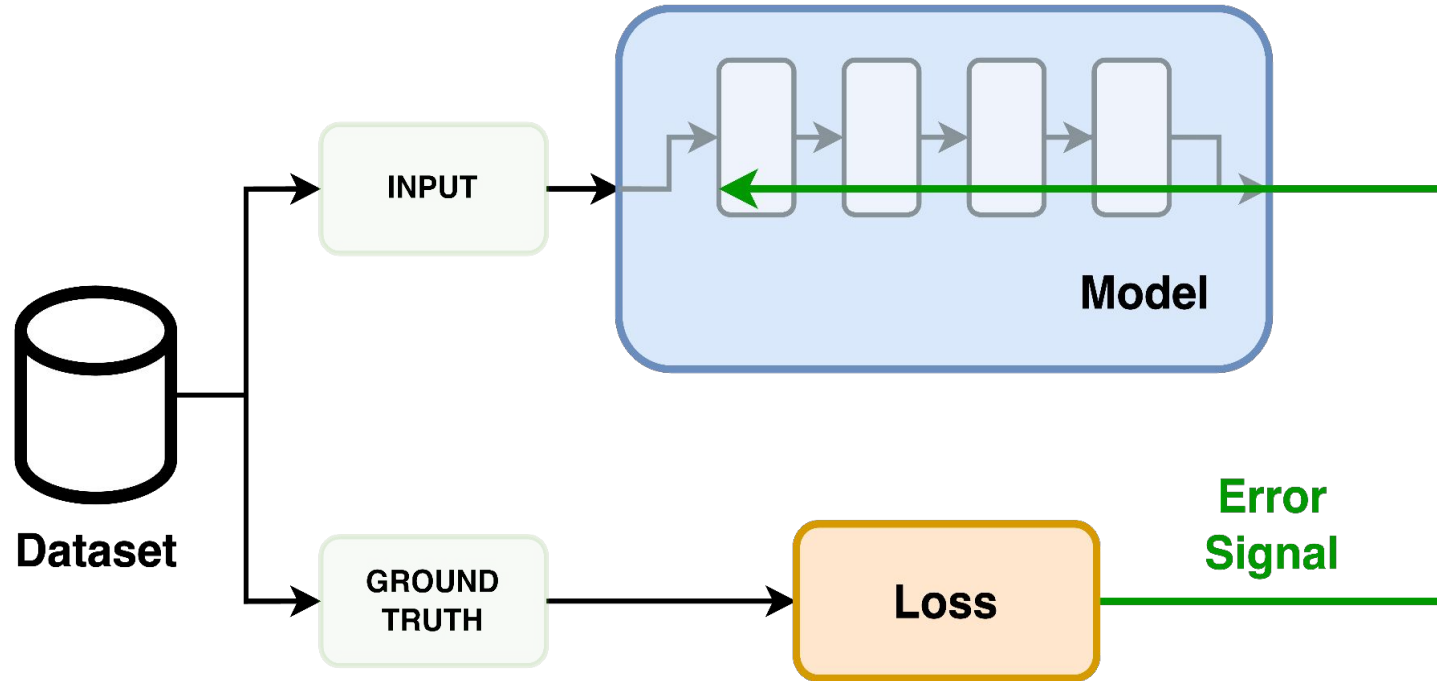
Augmented guitar capturing percussive techniques on the acoustic guitar.

Expressive gesture representation:
complex/brittle deterministic solutions.

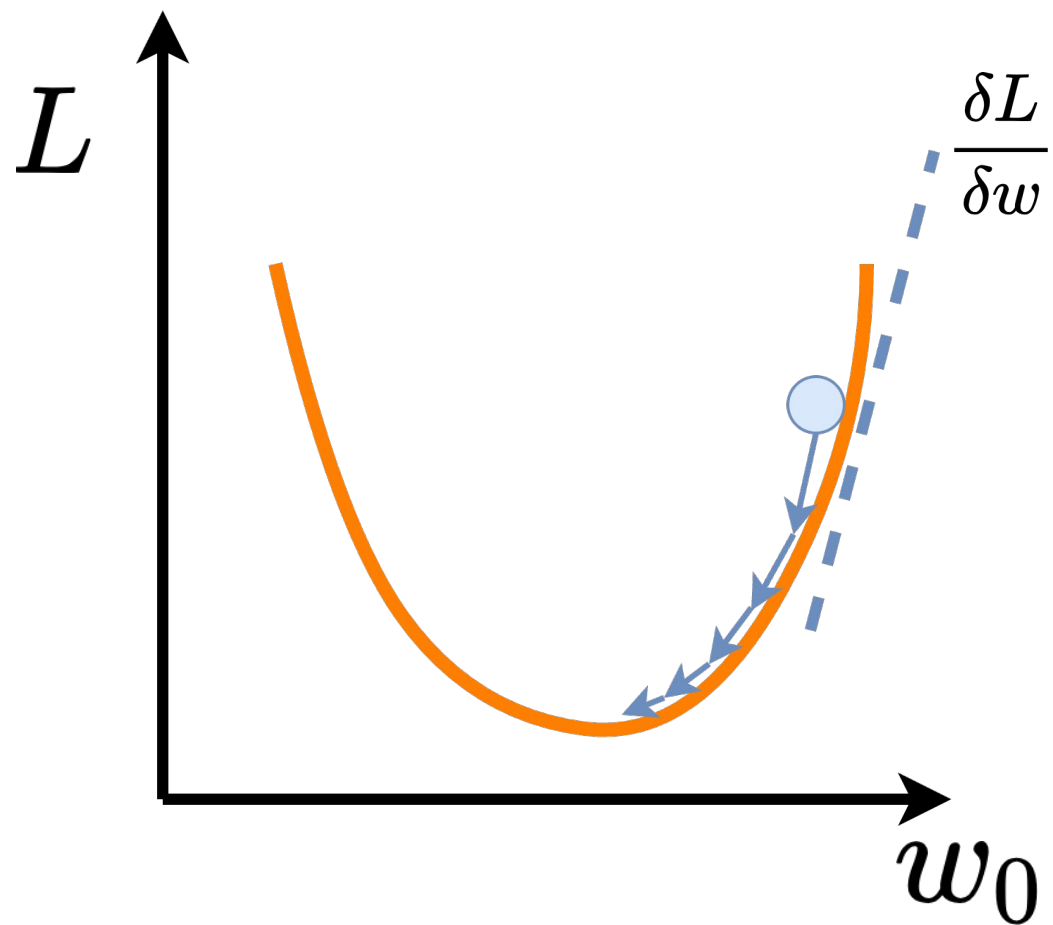


How Deep Learning works

Overview

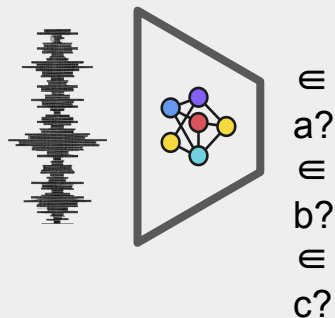


Training process



Relevant ML Tasks for Real-Time Audio

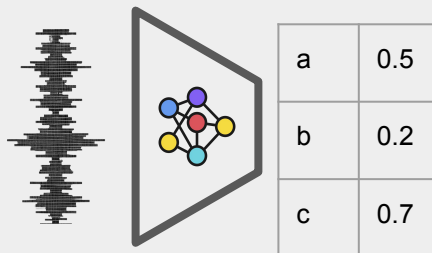
Transcription, discrete
gesture mapping.



Classification

Alignment to a particular
taxonomy

Parameter estimation,
waveform prediction



Regression

Estimation of quantities

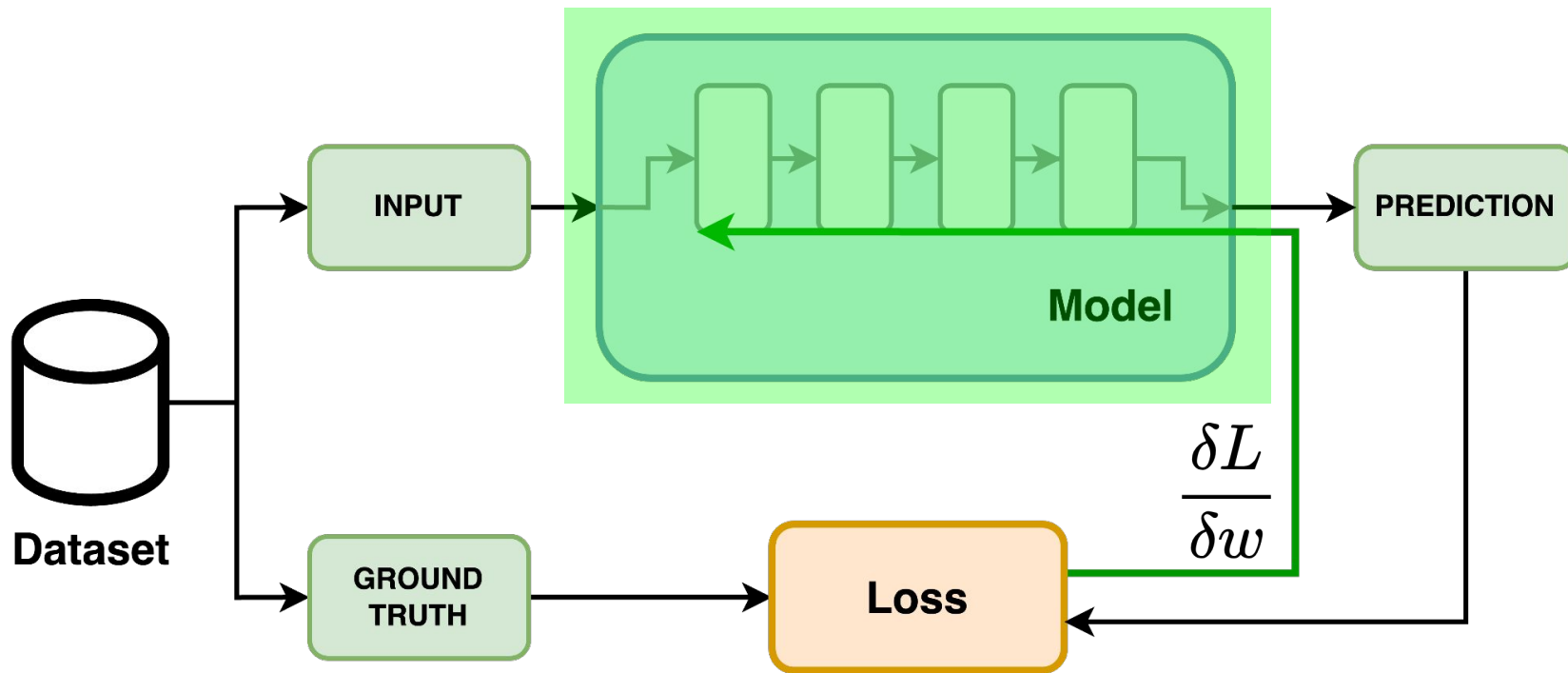
Compression, Generation,
Sampling



Representation Learning

Fit a distribution on a dataset

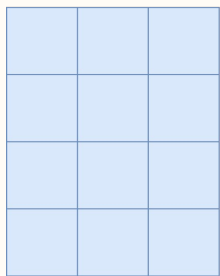
Overview



The Building Blocks

Matrix Multiplication

Weights



•

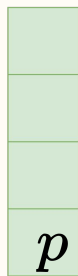


+

Biases



=

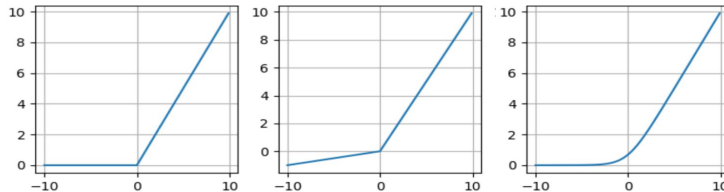


 Learnable
Parameters

 Layer's I/O

$$p = Ax + b$$

Activation Functions

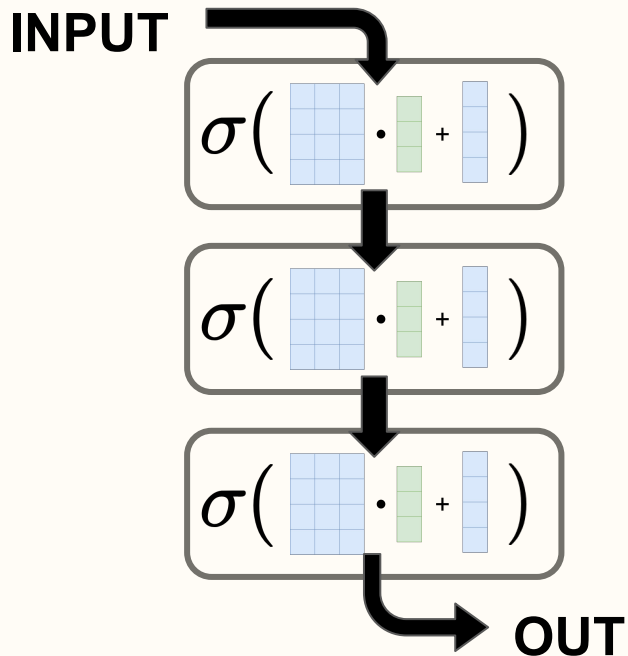


$$\begin{matrix} \sigma \\ \sigma \\ \sigma \end{matrix} \left(\begin{matrix} \text{green box} \\ \text{green box} \\ p \end{matrix} \right) = \begin{matrix} \text{green box} \\ \text{green box} \\ y \end{matrix}$$

$$y = \sigma(Ax + b)$$

Architectures

Stacking Operations



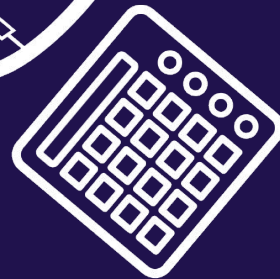
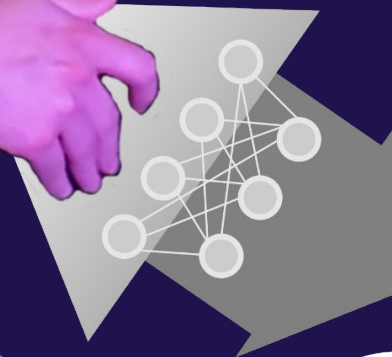
Layer Design

Clever ways
of doing
matrix multiplication
to exploit data structures.

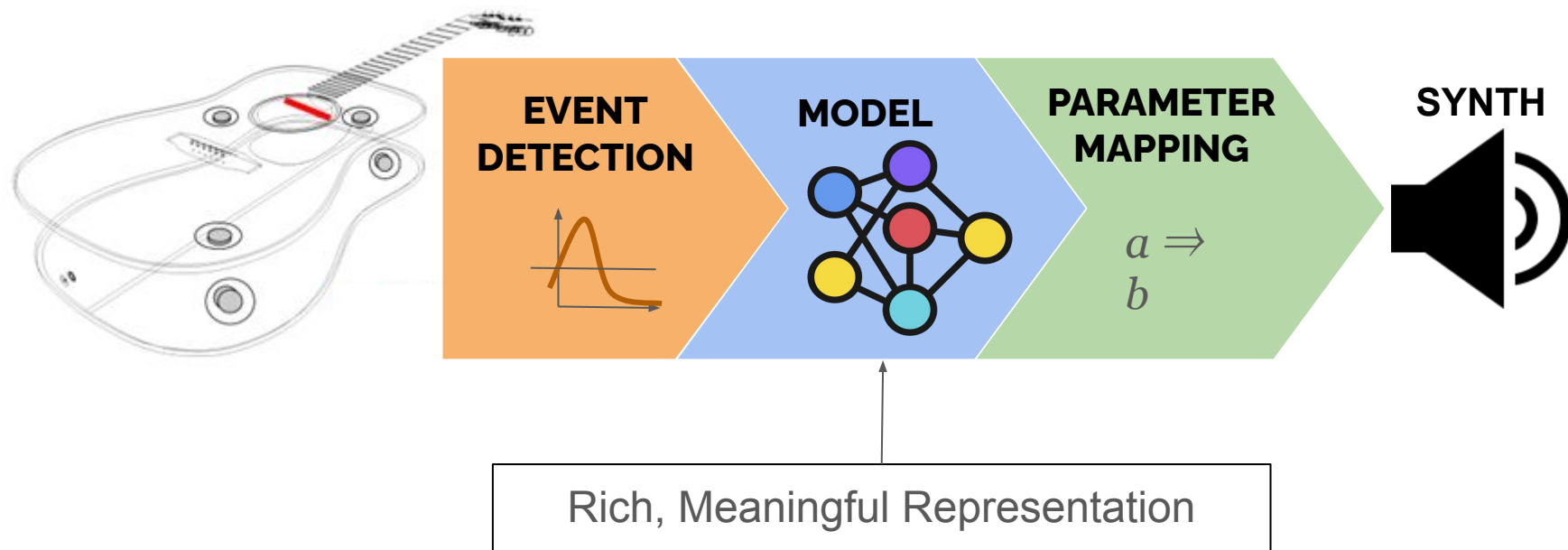


- Dense Layer
- Convolution
- Recurrency

The HITar



The HlTar's pipeline



Where I started from

Representation of musical/artistic gestural language with a DL model.

- Real-time Automatic Drum Transcription. (Jacques 2018, MA 2021)
- Semantic description of musical gestures and dance. (Murray-Browne 2021)

Jacques and Roebel, 'Automatic Drum Transcription with Convolutional Neural Networks'. DAFx 2018.

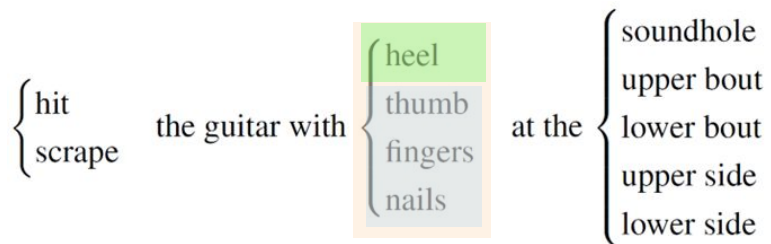
MA, Bhattacharjee, and Rao, 'Four-Way Classification of Tabla Strokes with Models Adapted from Automatic Drum Transcription'. ISMIR 2021.

Murray-Browne and Tigas, 'Latent Mappings'. NIME 2021.

The task behind the HITar

The “language”

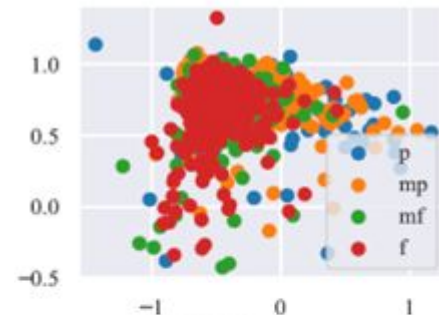
Taxonomy of percussive playing



CLASSIFICATION

The “nuance”

Represent sound events uniquely based on what they mean



REPRESENTATION LEARNING

Both things at the same time!

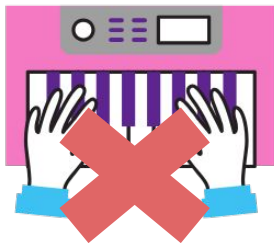
The Data for the HITar

Making good data was the **first**, **longest** and most **expensive** part of my project.

High Variance



Realistic Setting



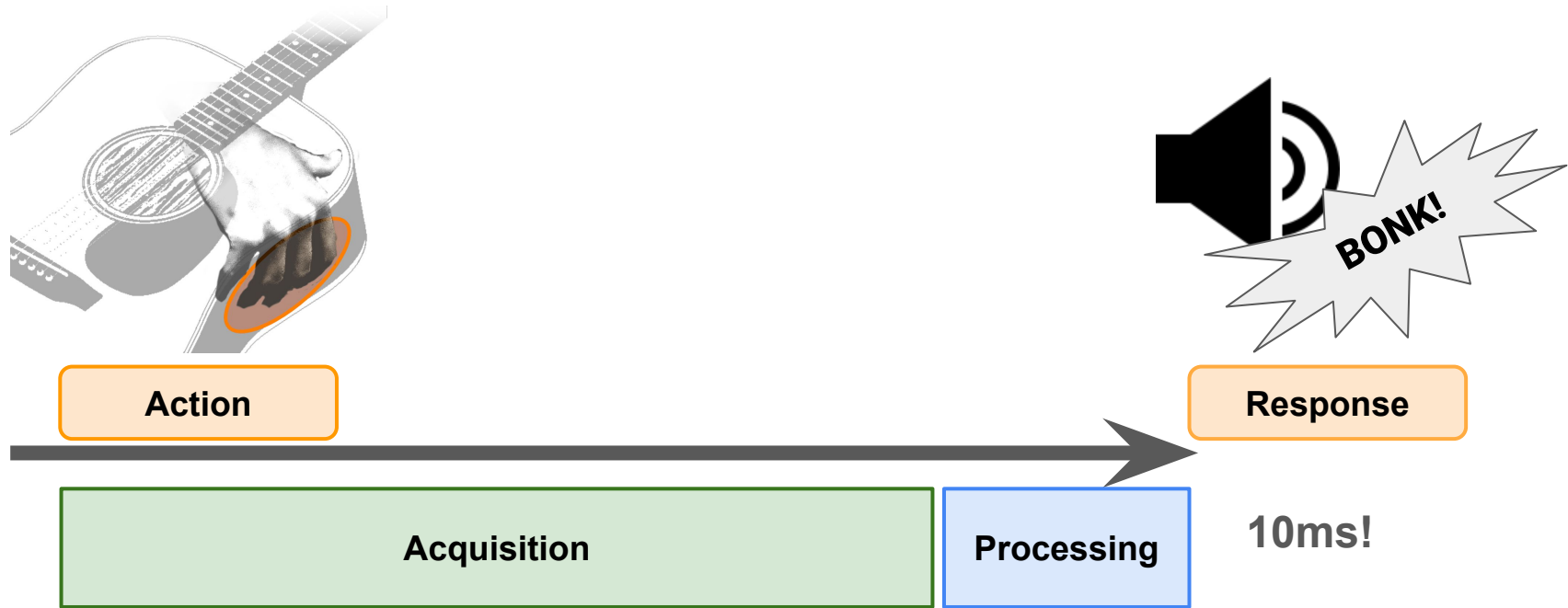
No Confounding Factors



Make **detailed labels**!

REAScript (and reapy) are your friends... ;)

The Real-Time Constraint of the HITar



Data Representation

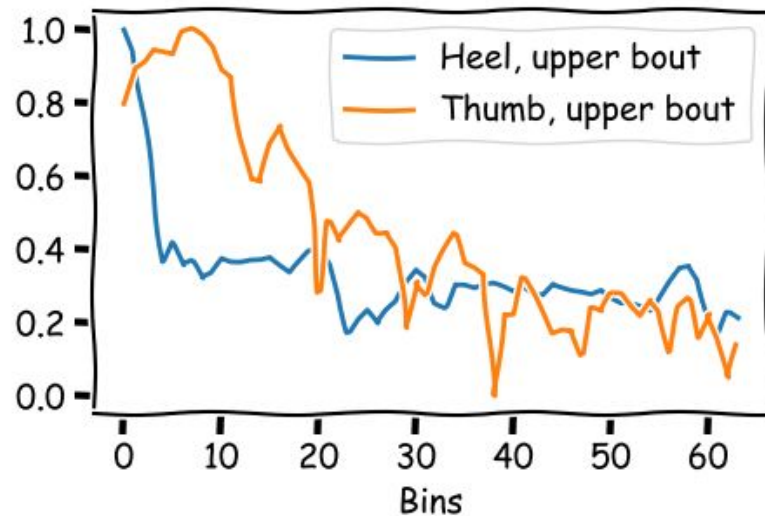
Representation

Input data:

- First 10 ms of the event
- Input representation:
Downsampled magnitude spectrum

Data augmentation

Phase inversion, filtering, clipping...



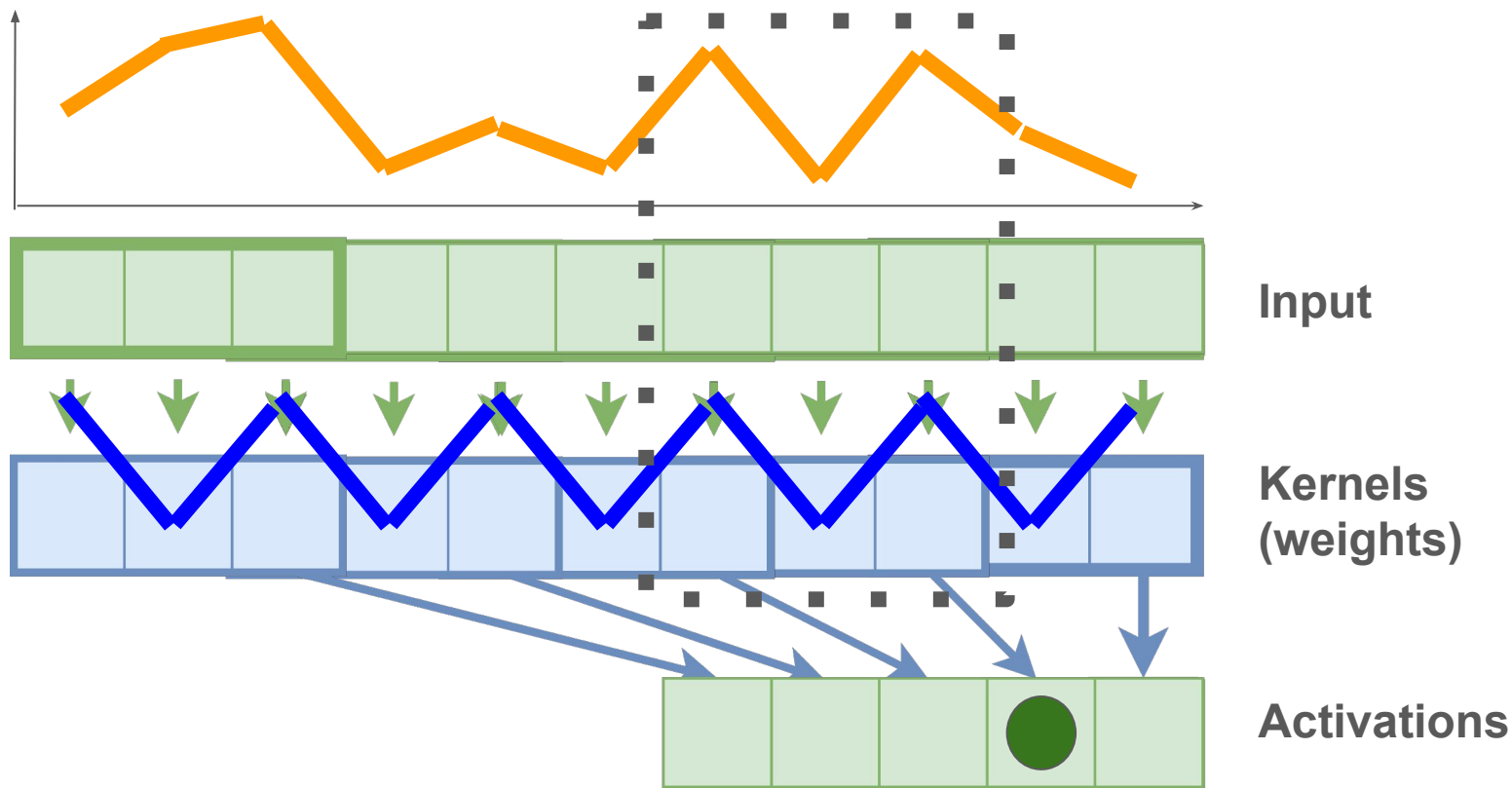
Intuitions behind the HITar's model

Task: learn a meaningful, readable representation of guitar body hits.

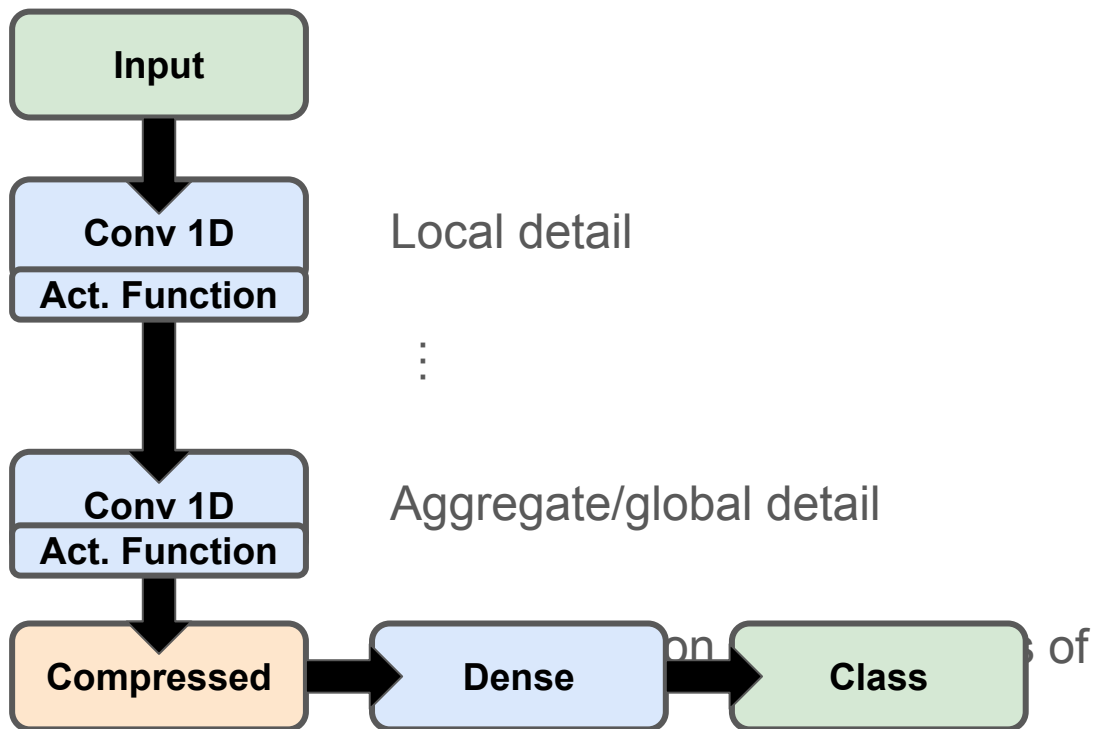
Input: magnitude spectrum.

Intuition: convolution to find patterns in spectra.

Extracting features with Convolutional Layers



Convolutional Encoding



HiTar - features of each unique hit?

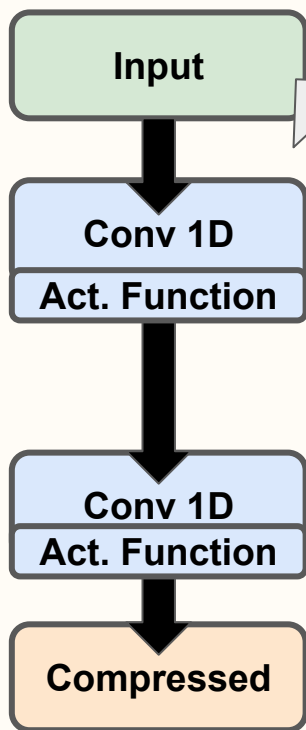
In a musical instrument, there are **many ways to play the same thing**.

Can we capture that in a NN?

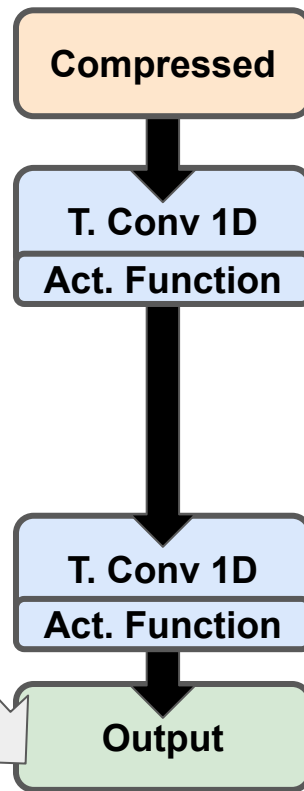
Intuition: expand the compressed representation back, to reproduce the input.

We can make a CNN de-compress data!

Convolutional Encoding

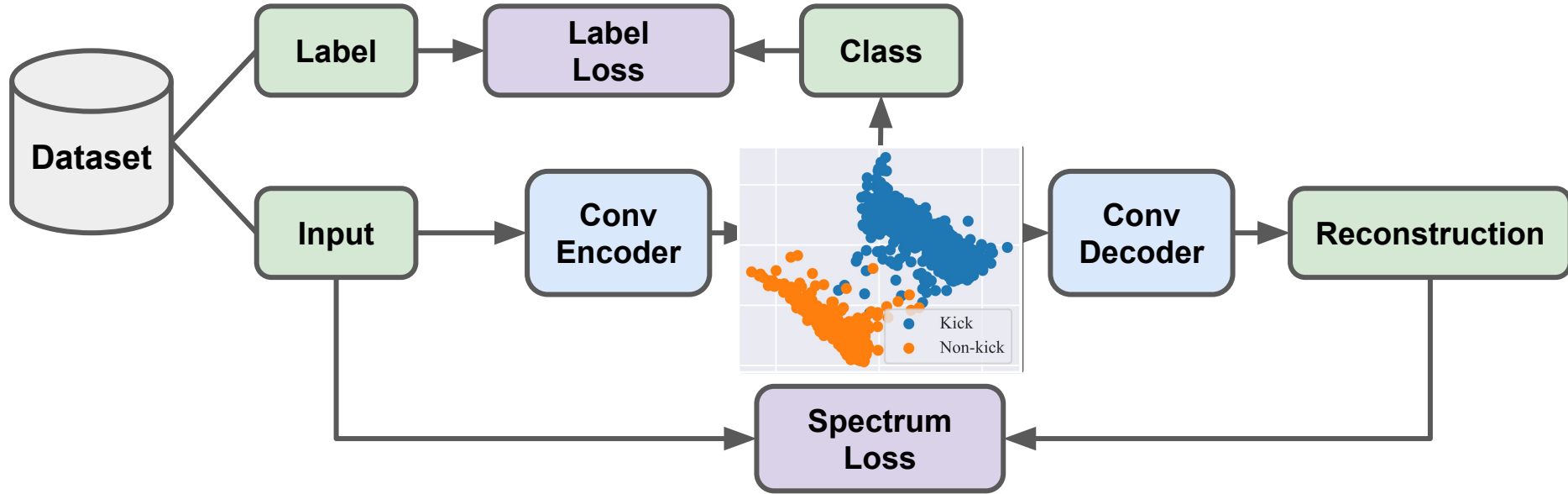


Convolutional Decoding

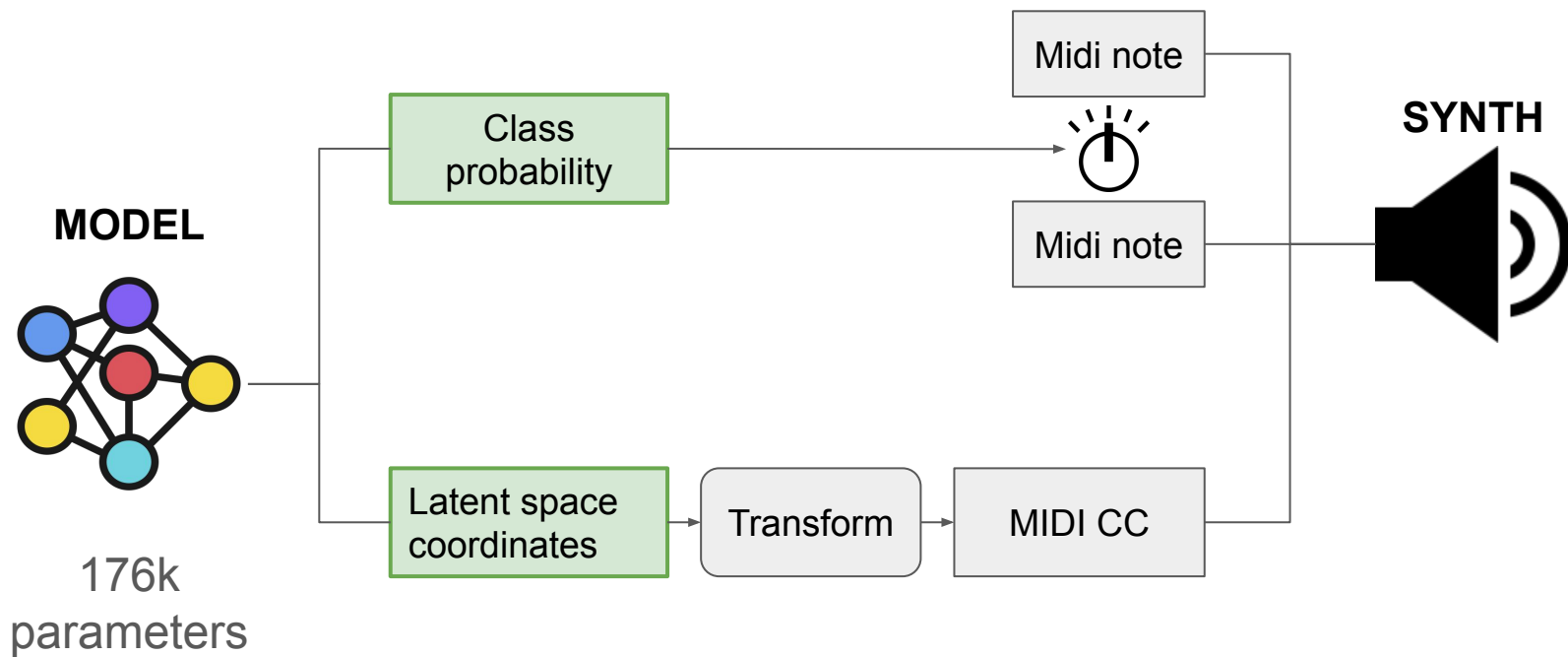


Same Size!

CNN for AutoEncoding

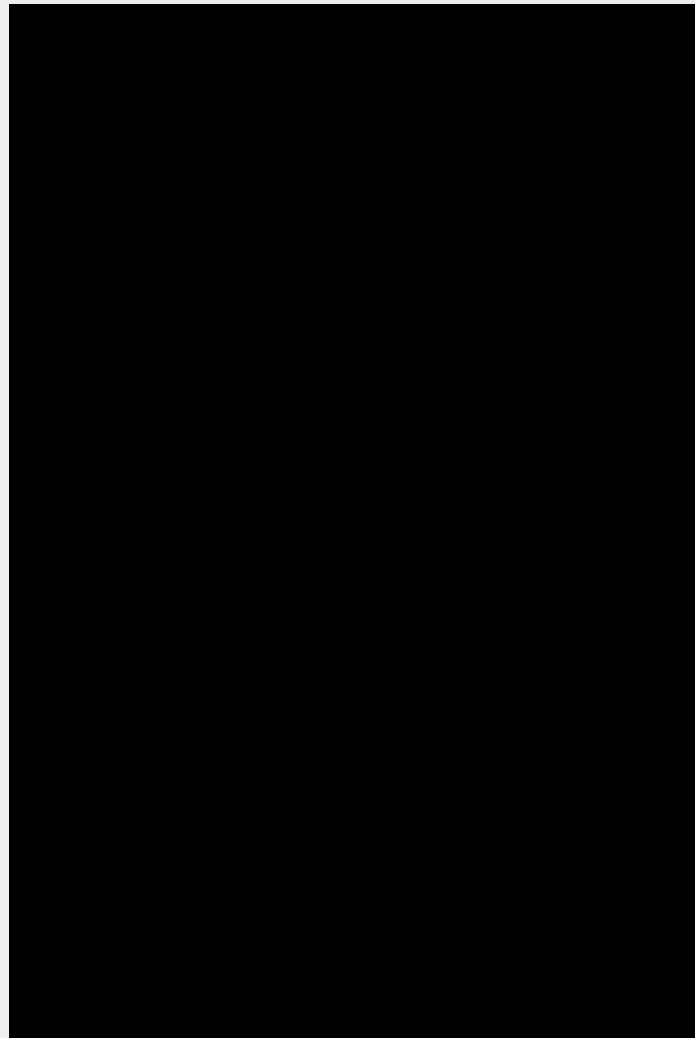
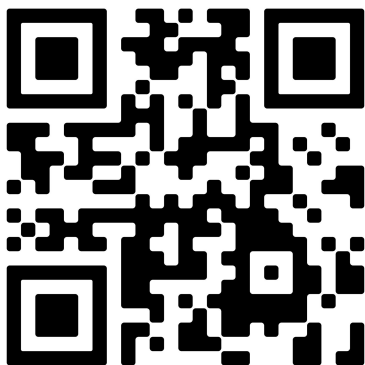


The Mapping





...in action!



My Takeaways

Good data is fundamental.

...but DSP engineering knowledge helps with **data augmentation**.

Real-time is a **constraint**, not an optimisation.

Bessel's Trick

FM Tone Transfer



In Gain

Pitch



RMS



Envelope
Model

TWINCLE

+

Algorithm



3

Oscillators

1

2

3

4

5

6

Coarse



Fine



Boost



$f = 2$



Out Gain



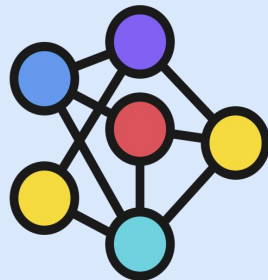
Status: Ready to play!

Tone Transfer

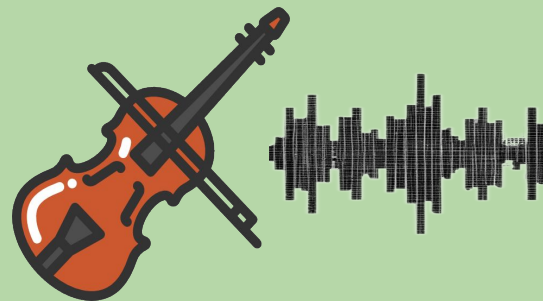
AUDIO INPUT



MODEL



OUTPUT

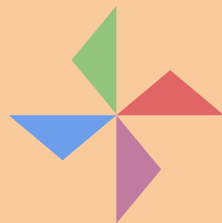


Design Premises



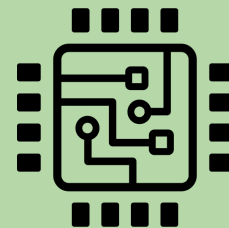
PERFORMANCE ORIENTED

- *Focus on music phrasing and articulation*



DIVERSE & INTERVENABLE

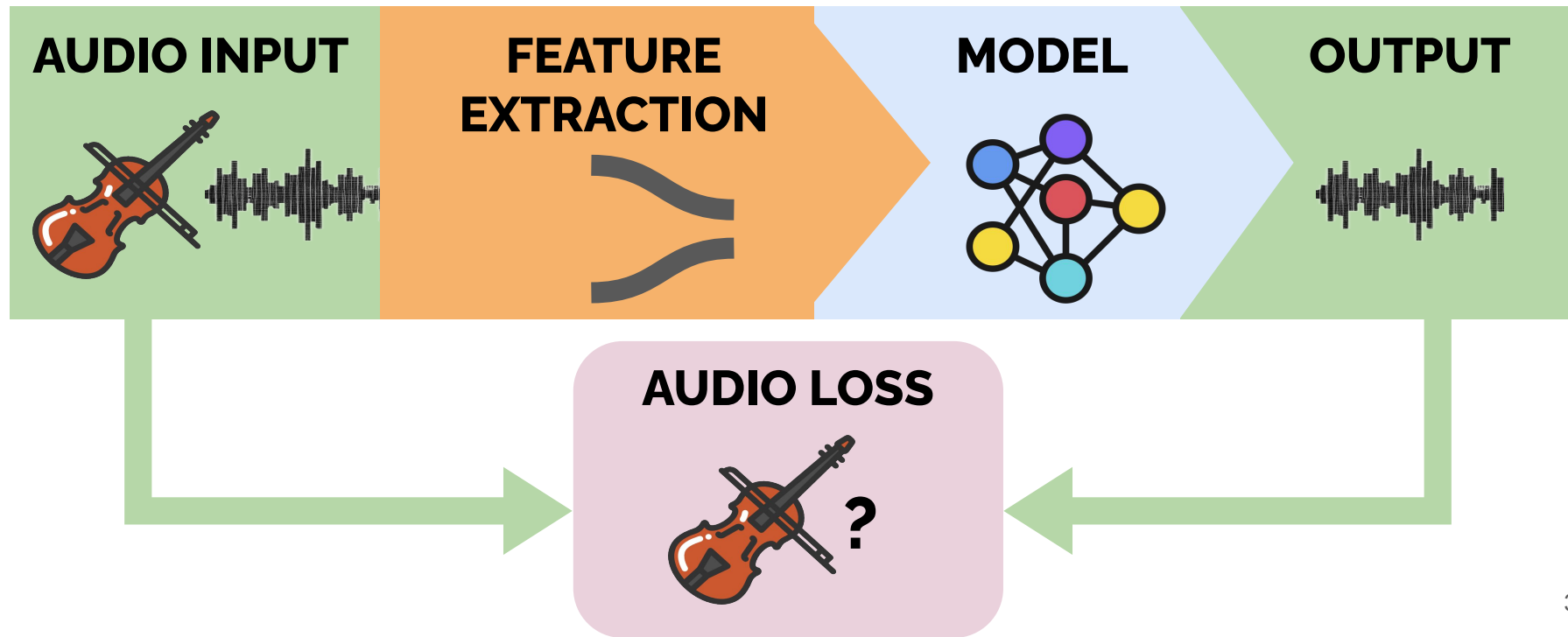
- *Variety of timbres*
- *Ability to edit sounds*



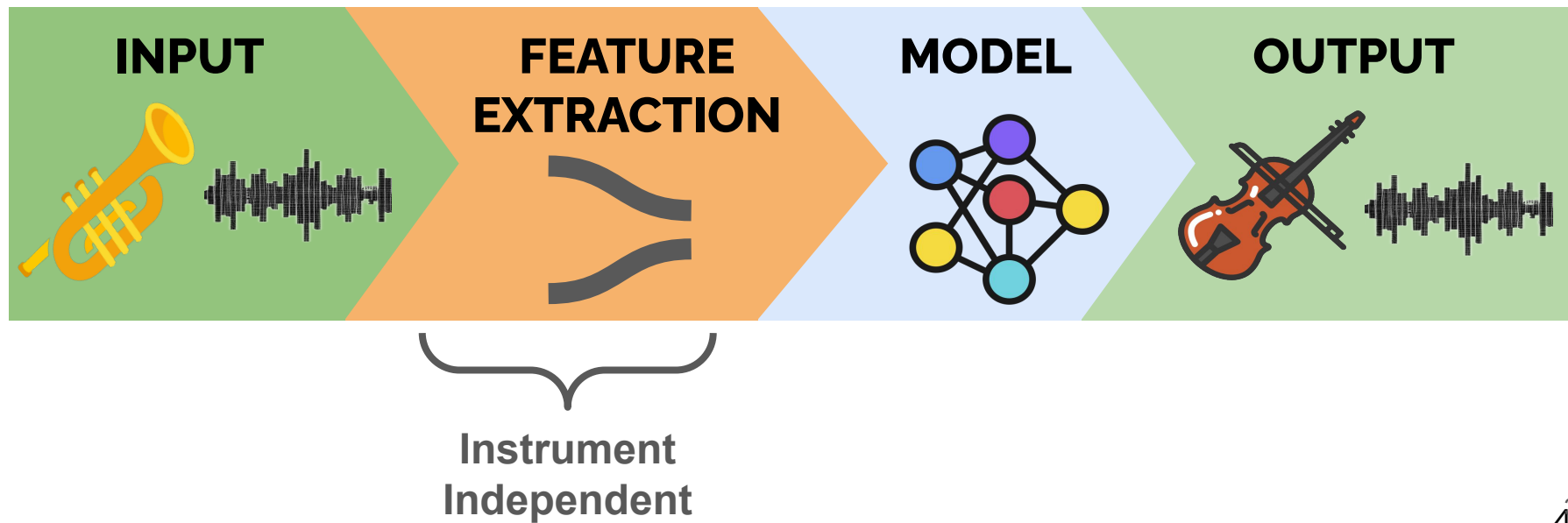
REAL TIME !

- *Low-latency inference*

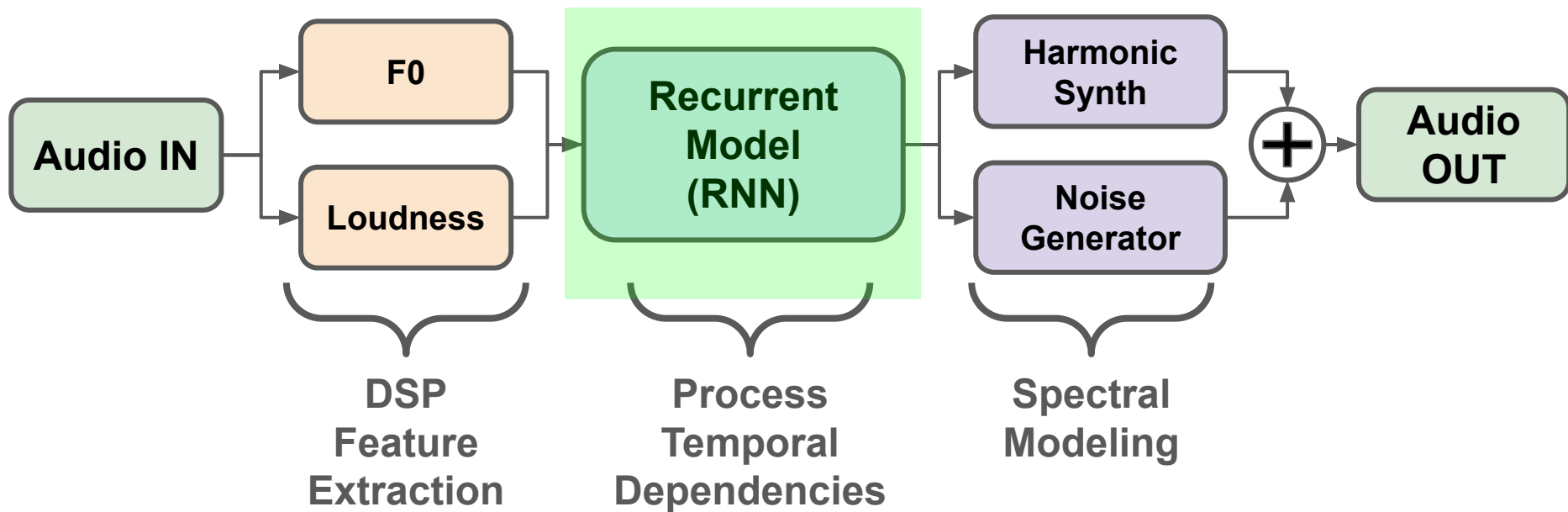
Tone Transfer Pipeline



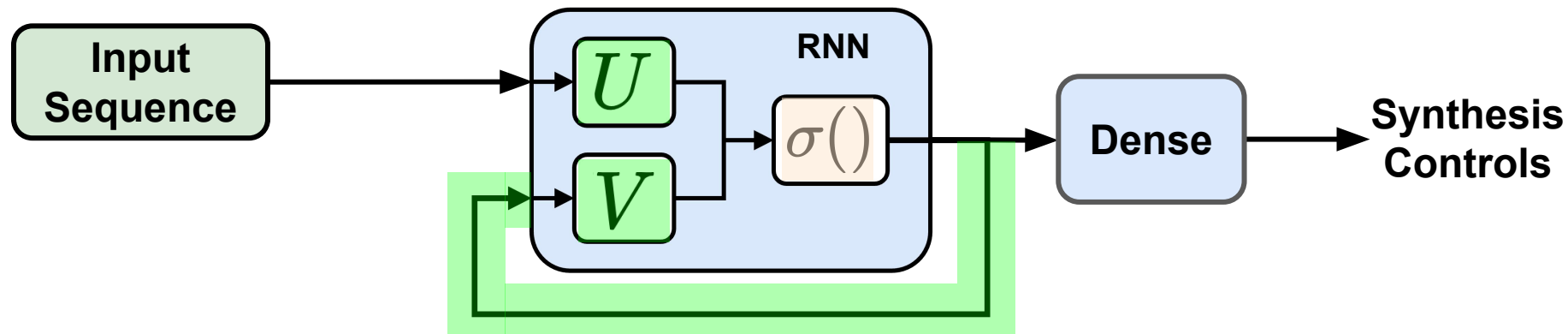
The claim



Our Reference Model: The DDSP Decoder (simplified)

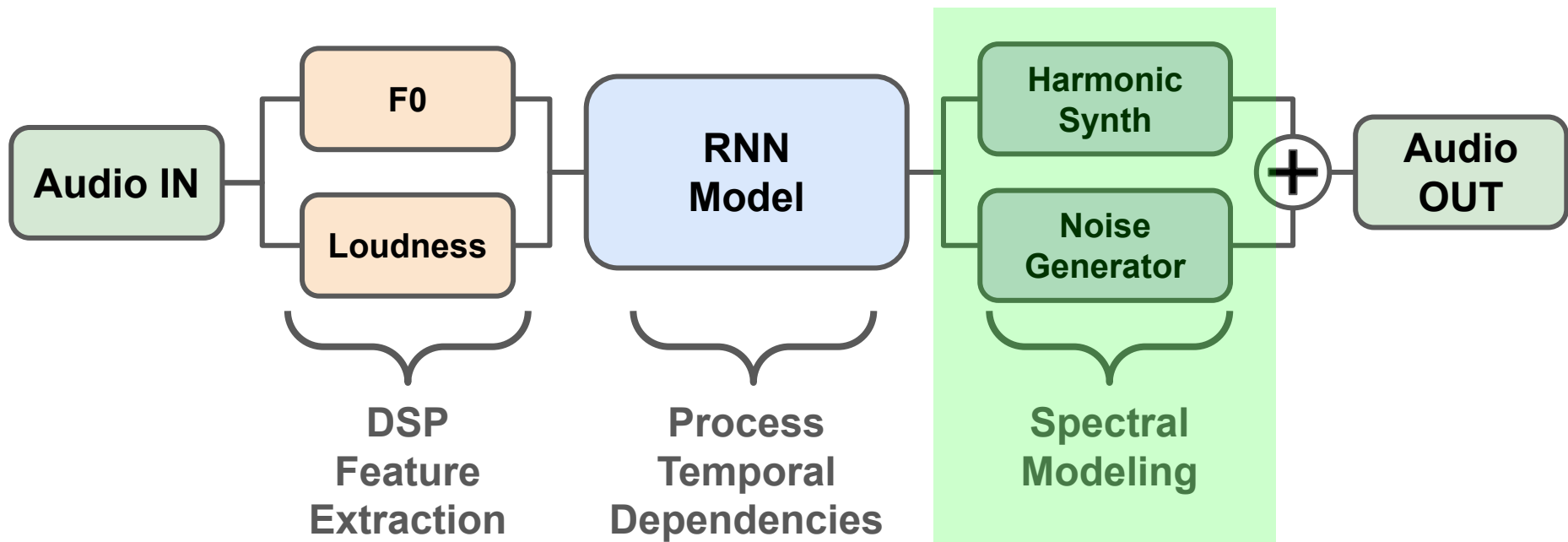


Temporal modelling with RNNs

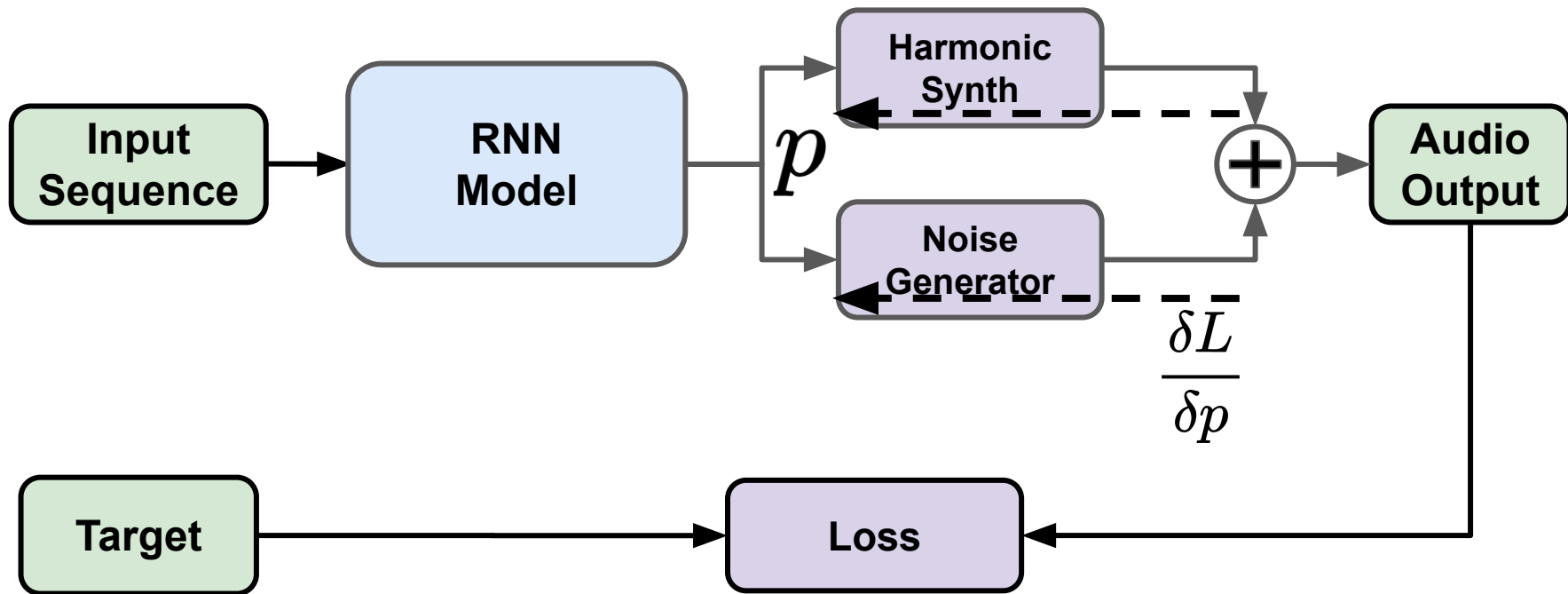


- Temporal aggregation through feedback, like an IIR filter

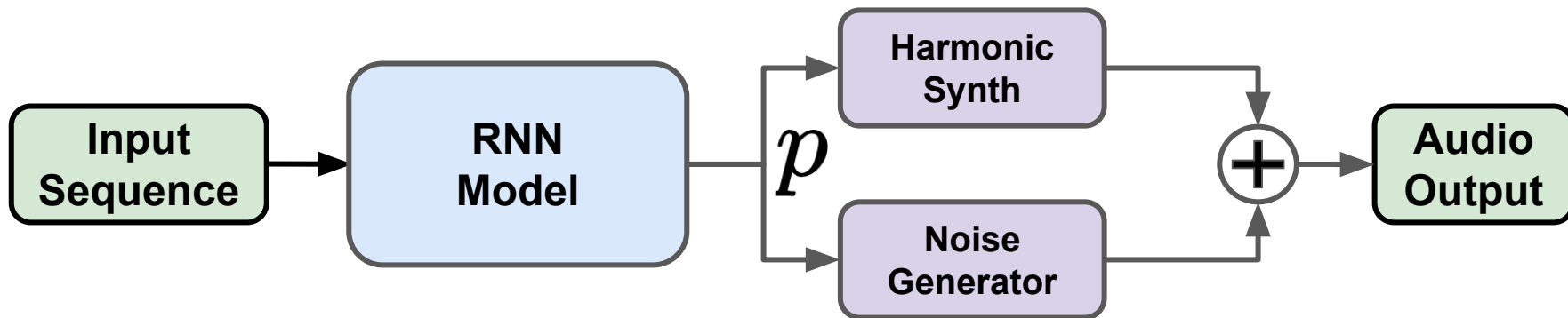
Reference Tone Transfer Model



Differentiable Signal Processing

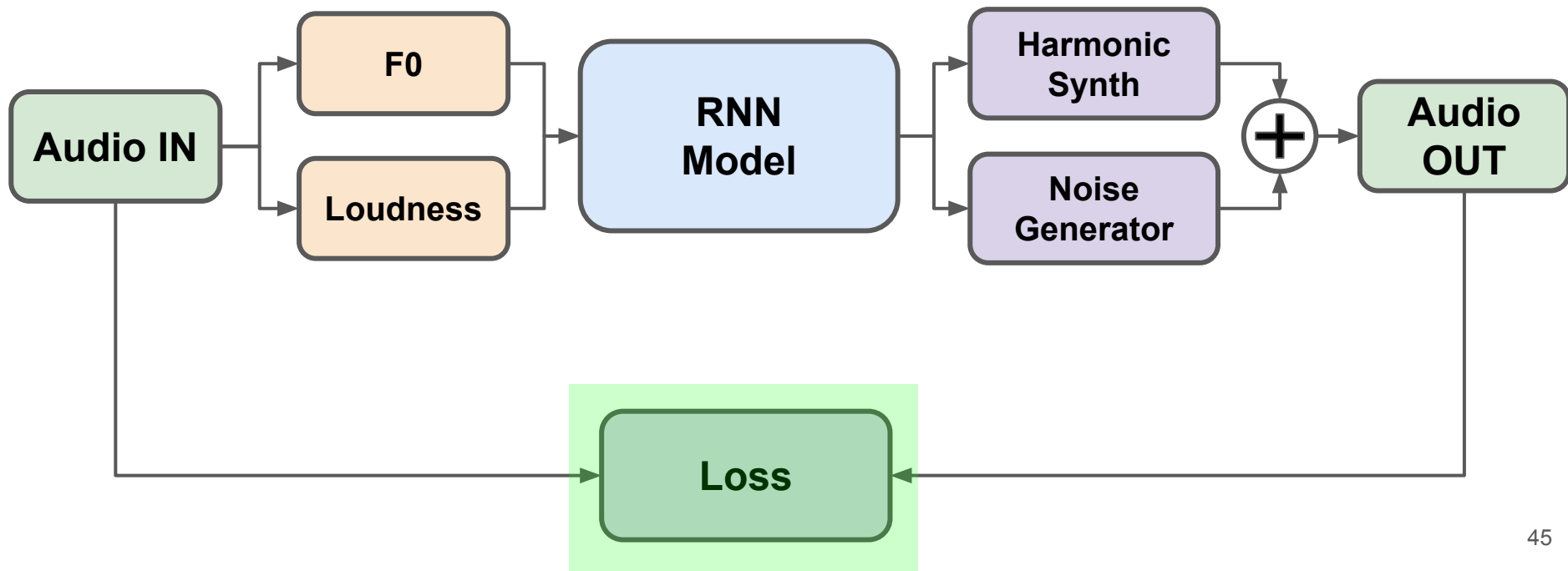


Differentiable Signal Processing

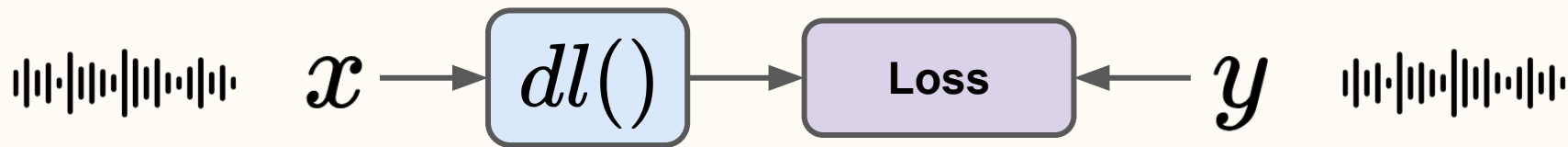



- Inductive bias towards Signal Generation / Processing.
- Oscillators, Filter Banks, Compressors, Waveshapers, FM, Reverb, FFT, Denoisers . . .

Reference Tone Transfer Model

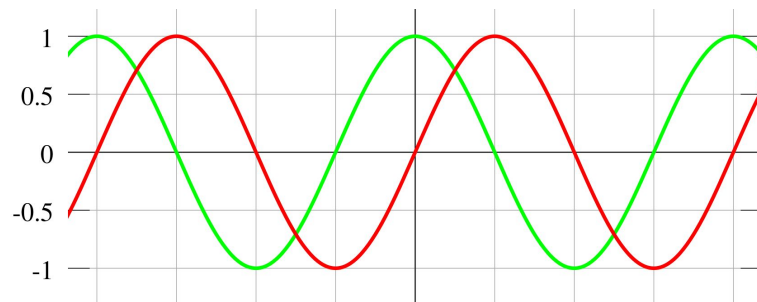


Loss Functions for Audio



$$L = ||dl(x) - y||$$


- Strong phase enforcement

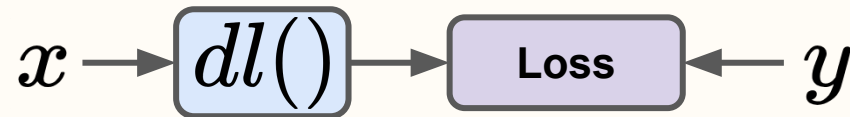


- Suitable for sample-by-sample generation

Loss Functions for Audio

Sample-by-sample Generation

- MAE / MSE



Block-by-Block Generation

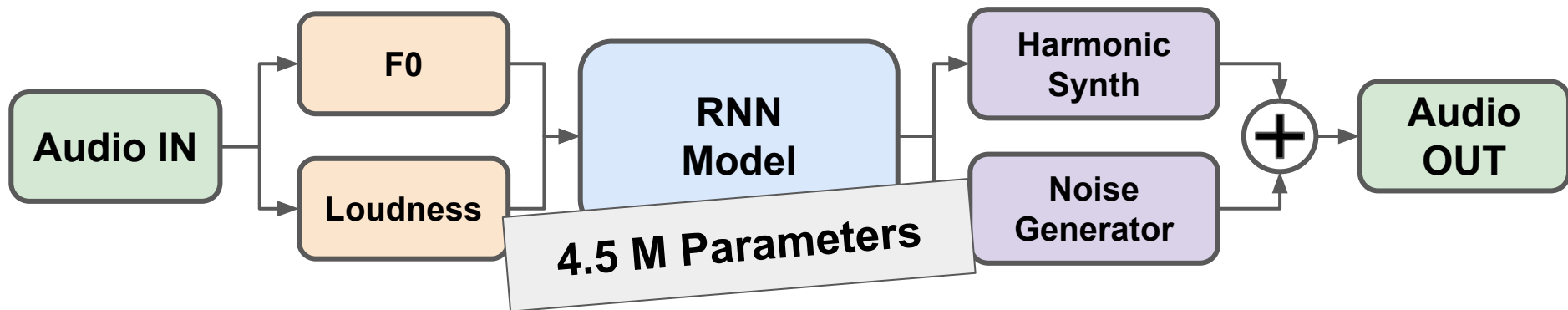
- Multi Scale
Spectral
Loss 

*Higher Frequency
Resolution*

*Higher Temporal
Resolution*

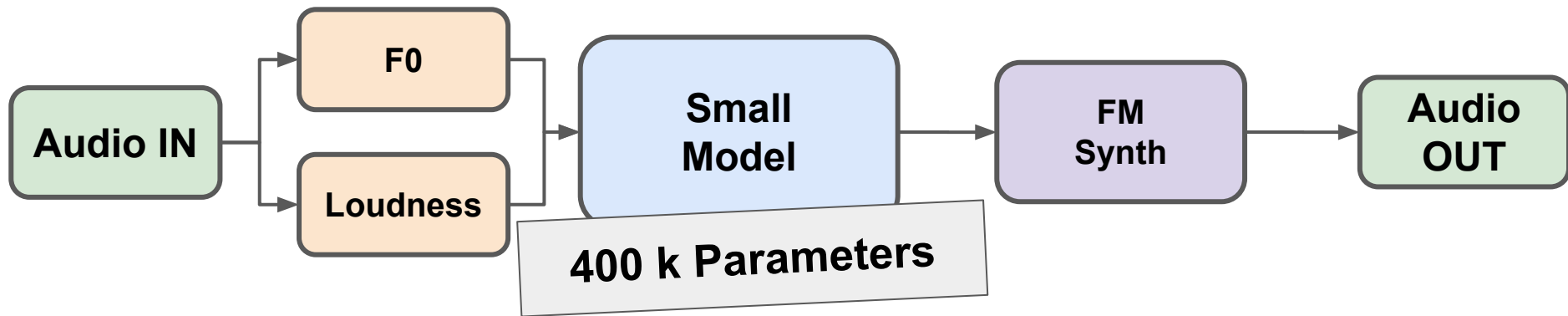
$$L = \begin{aligned} & \left| \left| S_{g(x)} - S_y \right| \right|_{1024} + \\ & \left| \left| S_{g(x)} - S_y \right| \right|_{256} + \\ & \left| \left| S_{g(x)} - S_y \right| \right|_{64} \end{aligned}$$

Reference Tone Transfer Model



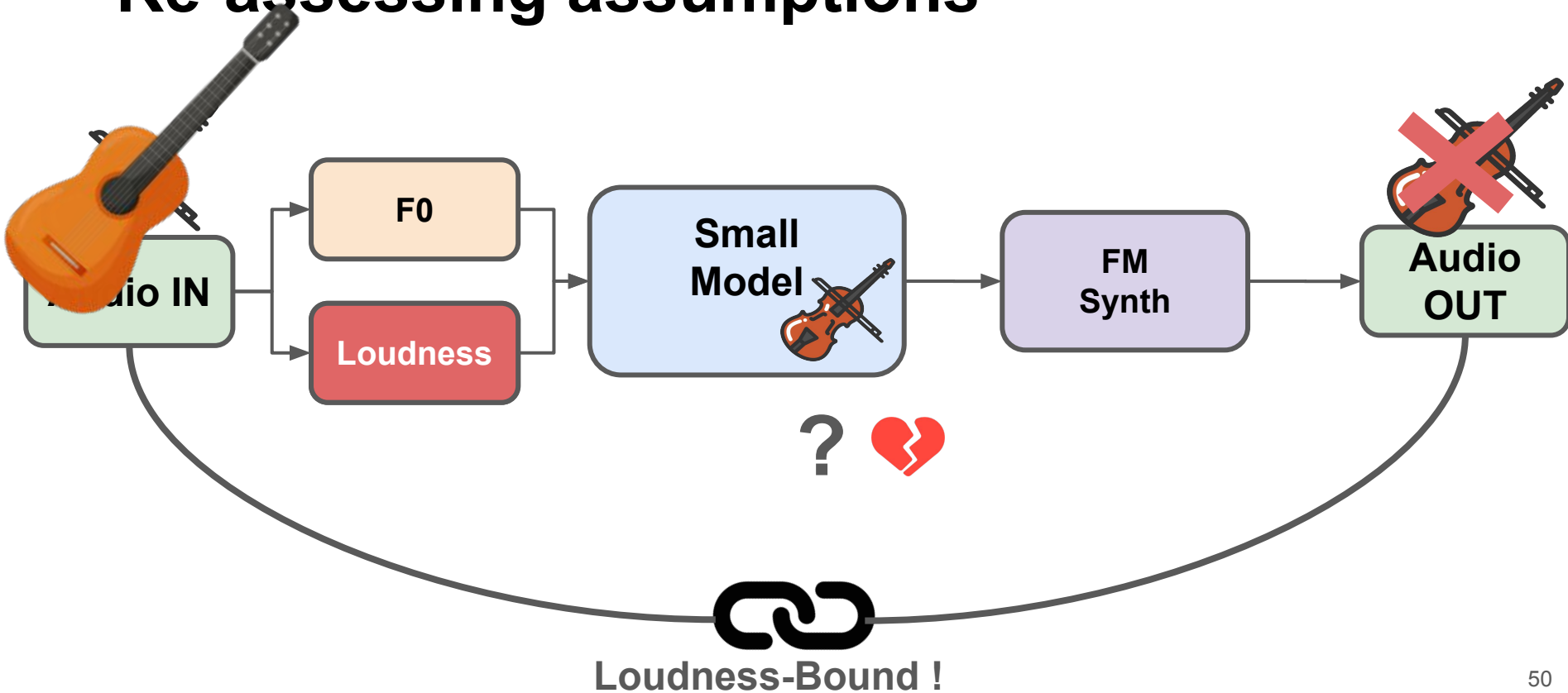
- **Block by Block Operation**
 - Suitable for Real-Time.

Idea: use a compact yet versatile synth (FM)

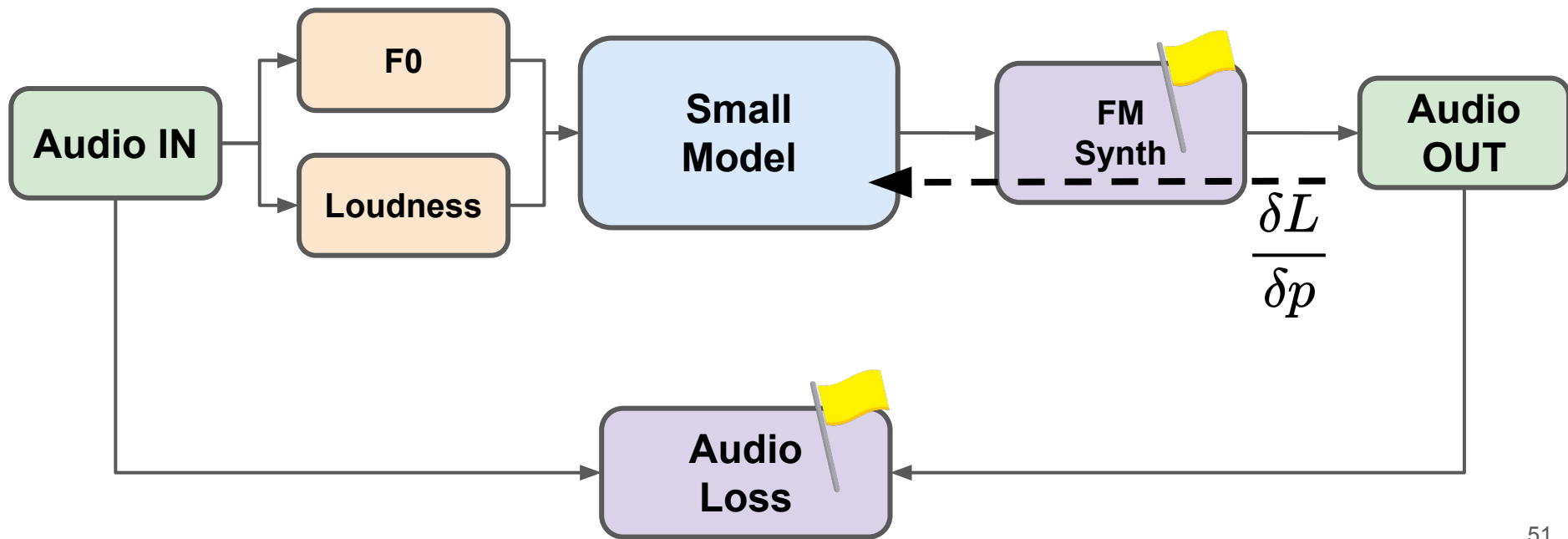


- Works great on the test set!
- 10x Smaller, works in real time!
- NOT Great during live use.

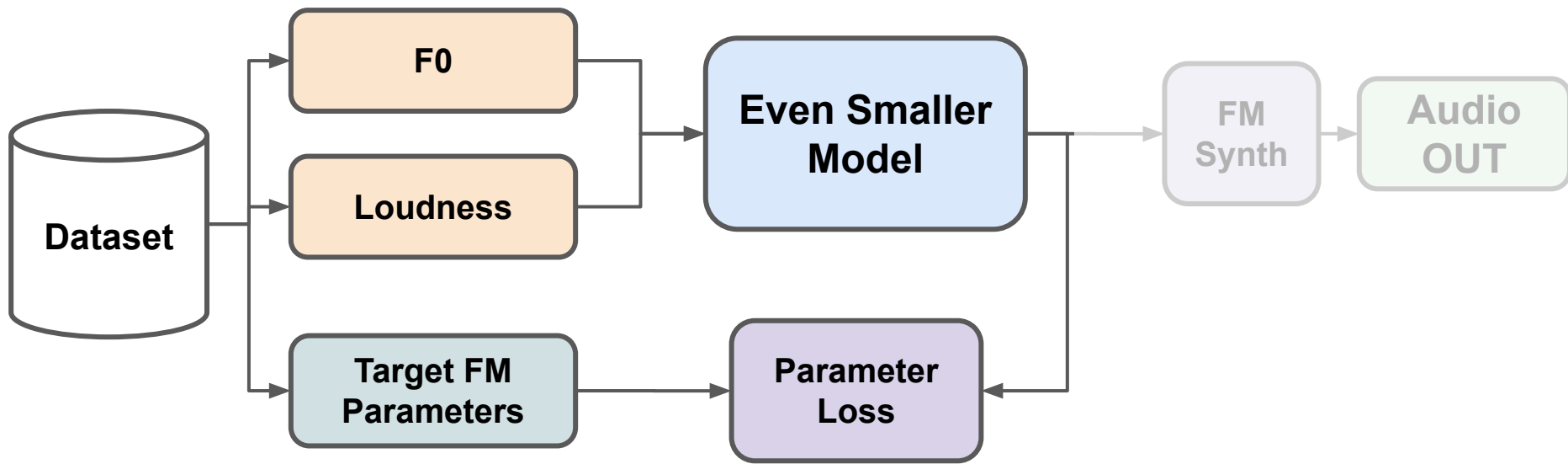
Re-assessing assumptions



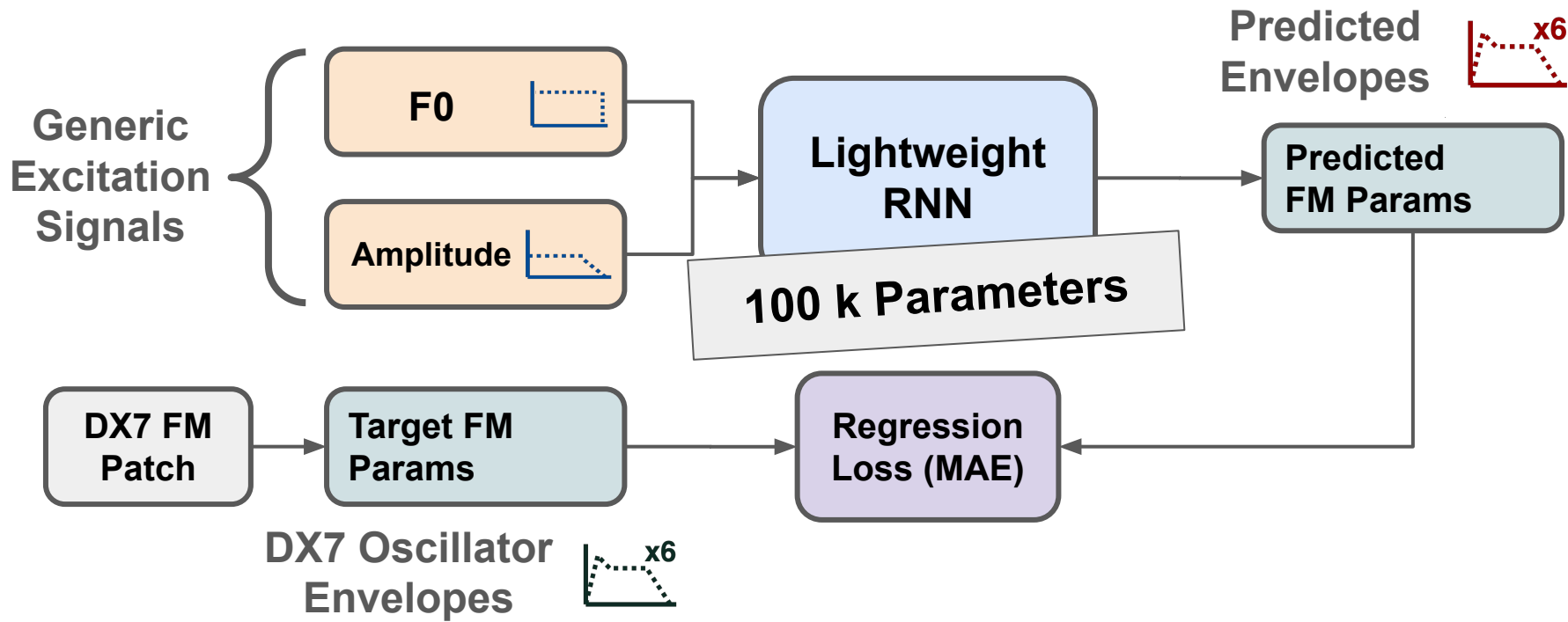
Using FM Synthesis with Audio Losses



Skipping the Synth during Training

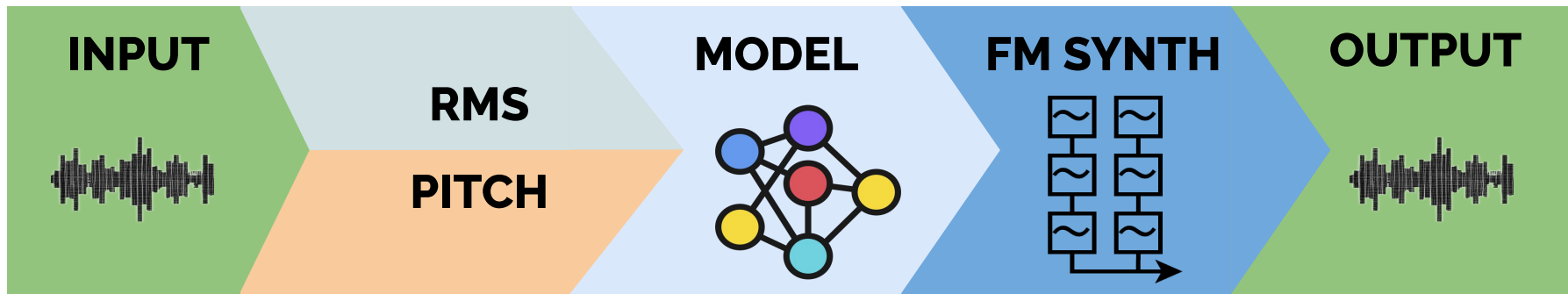


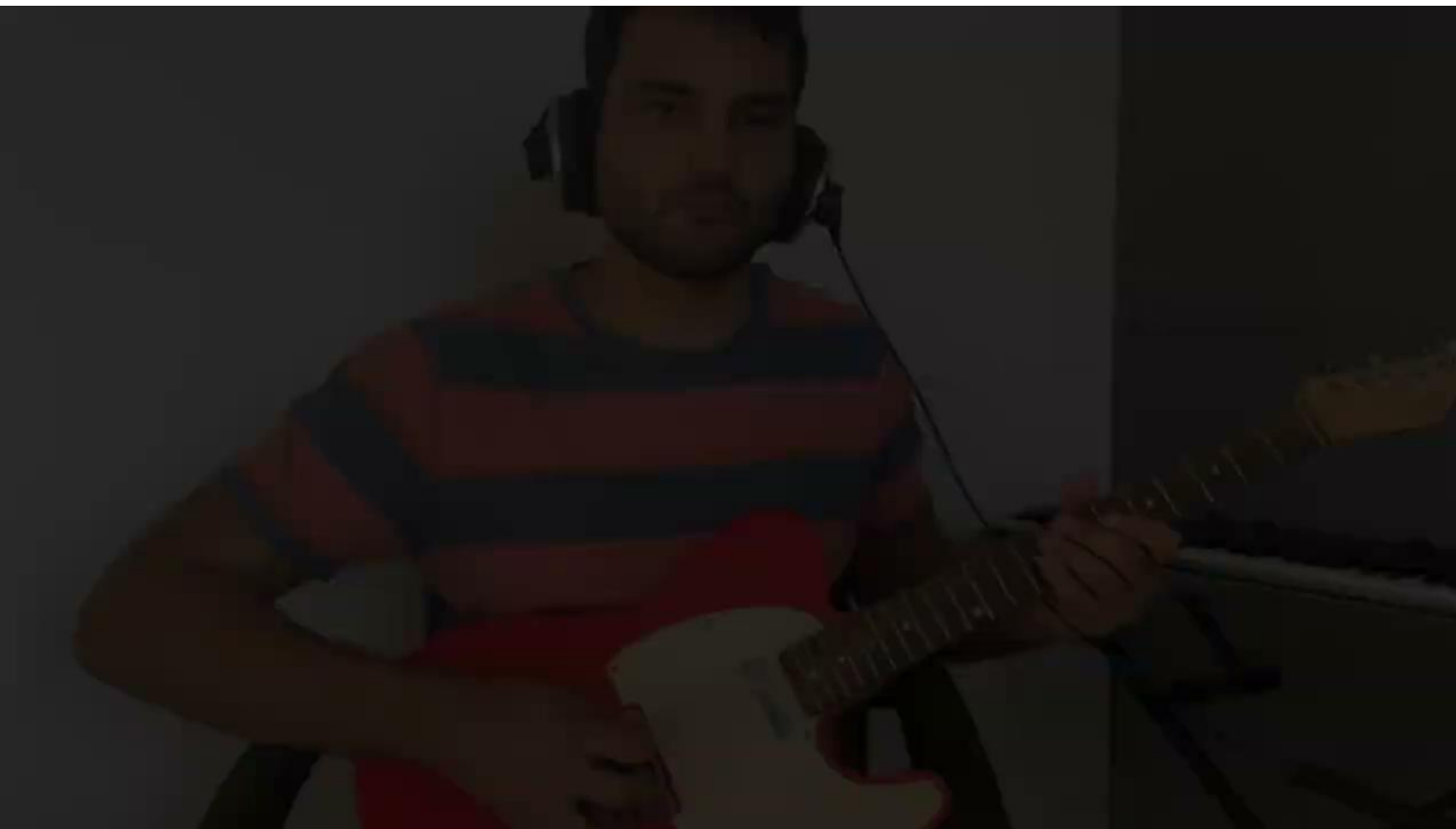
Skipping the Synth during Training



Bessel's Trick

FM Tone Transfer Pipeline





See it on
Github!



Takeaways

Have multiple loss functions at hand.

DSP Primitives can help simplify the network architecture.

Think of what is **explicitly** and **implicitly** learned in the model.

Final Remarks

Takeaways

- There's a strong DSP intuition behind DL for Audio.
- Clarify your real-time constraints early.
- Don't work from first principles.
- Be data-driven
 - Think about WHAT data.
 - Think about HOW the data is treated.

The Augmented Instruments Lab is recruiting!

1-2 PhD positions

Email: andrew.mcpherson@imperial.ac.uk

1 postdoc position

Areas of interest:

- Musical interaction
- Embedded hardware for audio and sensors
- Critical perspectives on technology and culture

Applications due
December/January

See instrumentslab.org



Thank you



@the_hitar



@caspefranco

