ual: **creative computing institute**

# LSTM Modelling of FFT Magnitudes for Neural Audio Generation with Minimal Data

Louis McCallum

CCI, UAL, London

MusicRAI 2024

Training LSTM models on sequences of **spectral data** can allow for **high quality audio** to be generated from **minimal datasets** of raw audio

This, in conjunction with **short training times**, provides significant creative and accessibility advantages that make it a more than **viable alternative** to other state of the art audio generation methods.

ual: creative computing
institute

# Heavier Models

Whilst audio quality from trained models is high, it requires significant **equipment**, **data** and **time** to acquire these trained models


Train for 1-2 days, min 3 hours of audio


Train for 1-2 hours, 12 mins audio


the tokenizers and the autoregressive models for the semantic and acoustic modeling stages are trained on a dataset containing five million audio clips, amounting to 280k hours of music at 24 kHz. Each of the stages is trained with multi-

IMPOSSIBLE?

# Heavier Models

Although using other's **pretrained models** may be enough for some musicians (essentially using an novel off the shelf synthesiser), **being involved in the training** process can be a significant part of the creative process
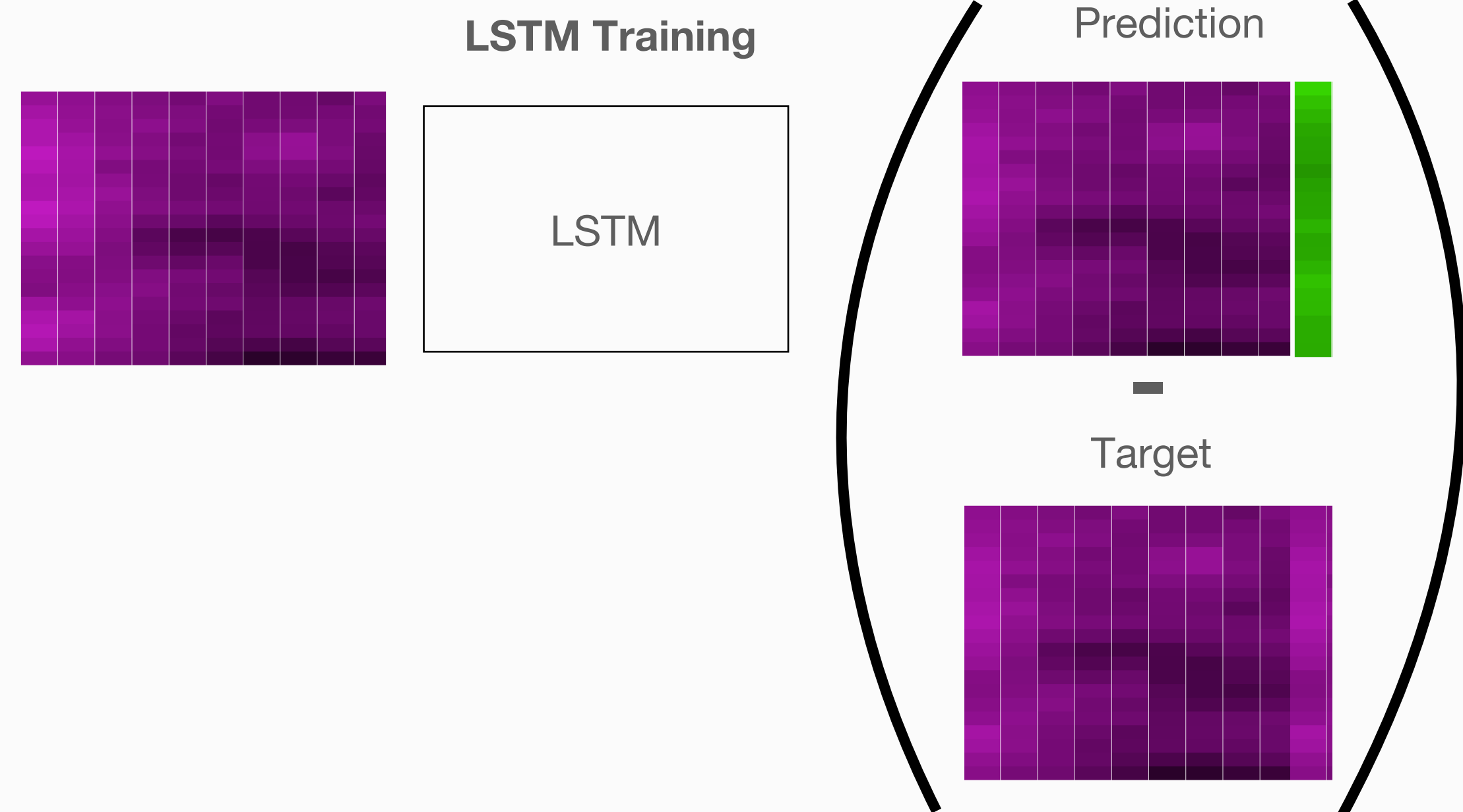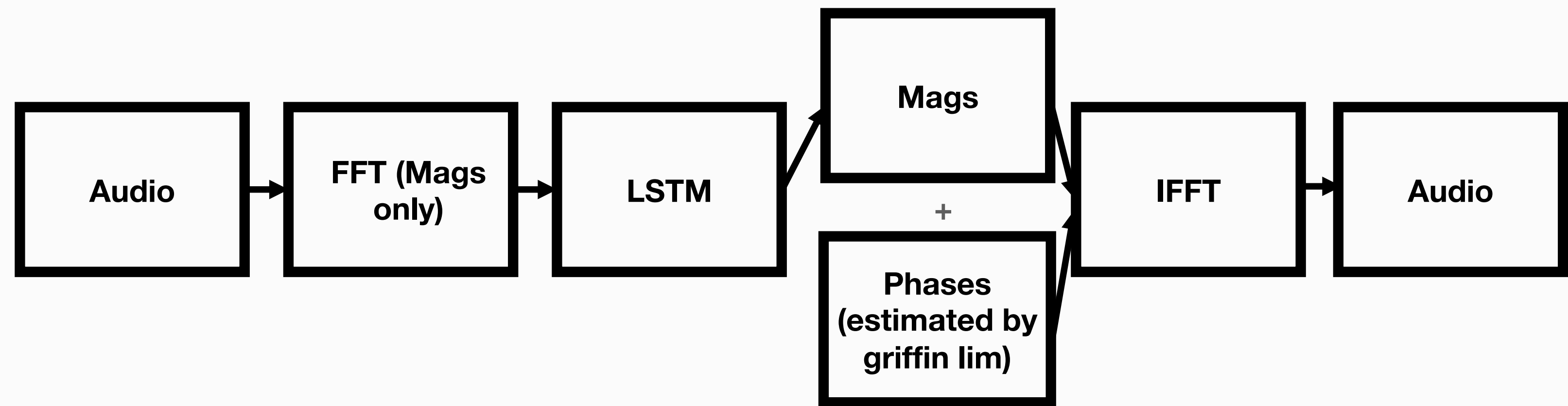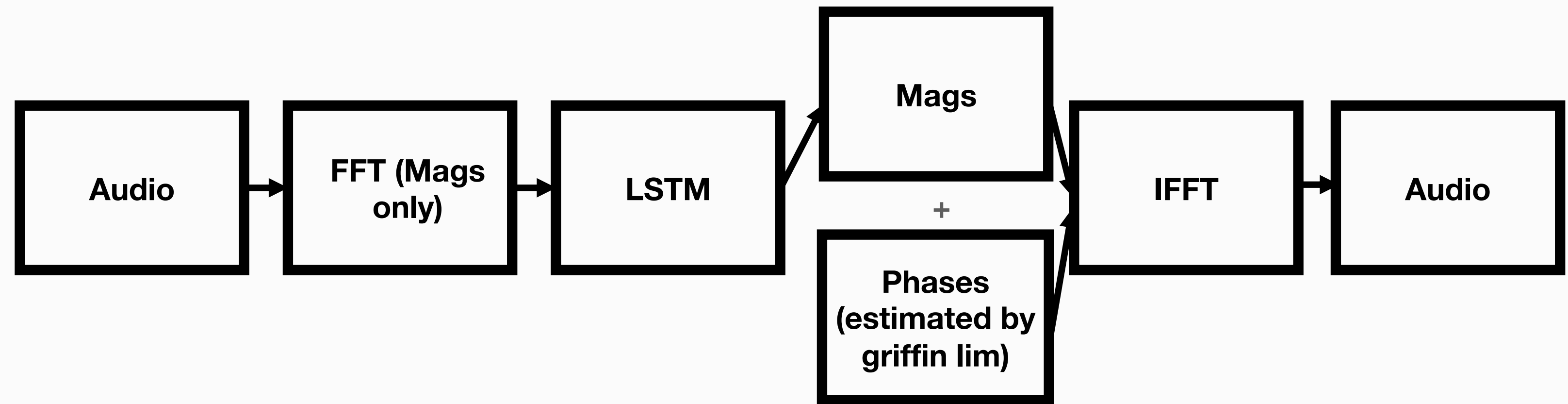
# MAGNet

A small network with as little as **3 LSTM layers** is used to model the sequence of FFT **magnitudes**. A **phase vocoding** method is then used to estimate the phases and allow for conversion back into **CD quality** audio

# MAGNet

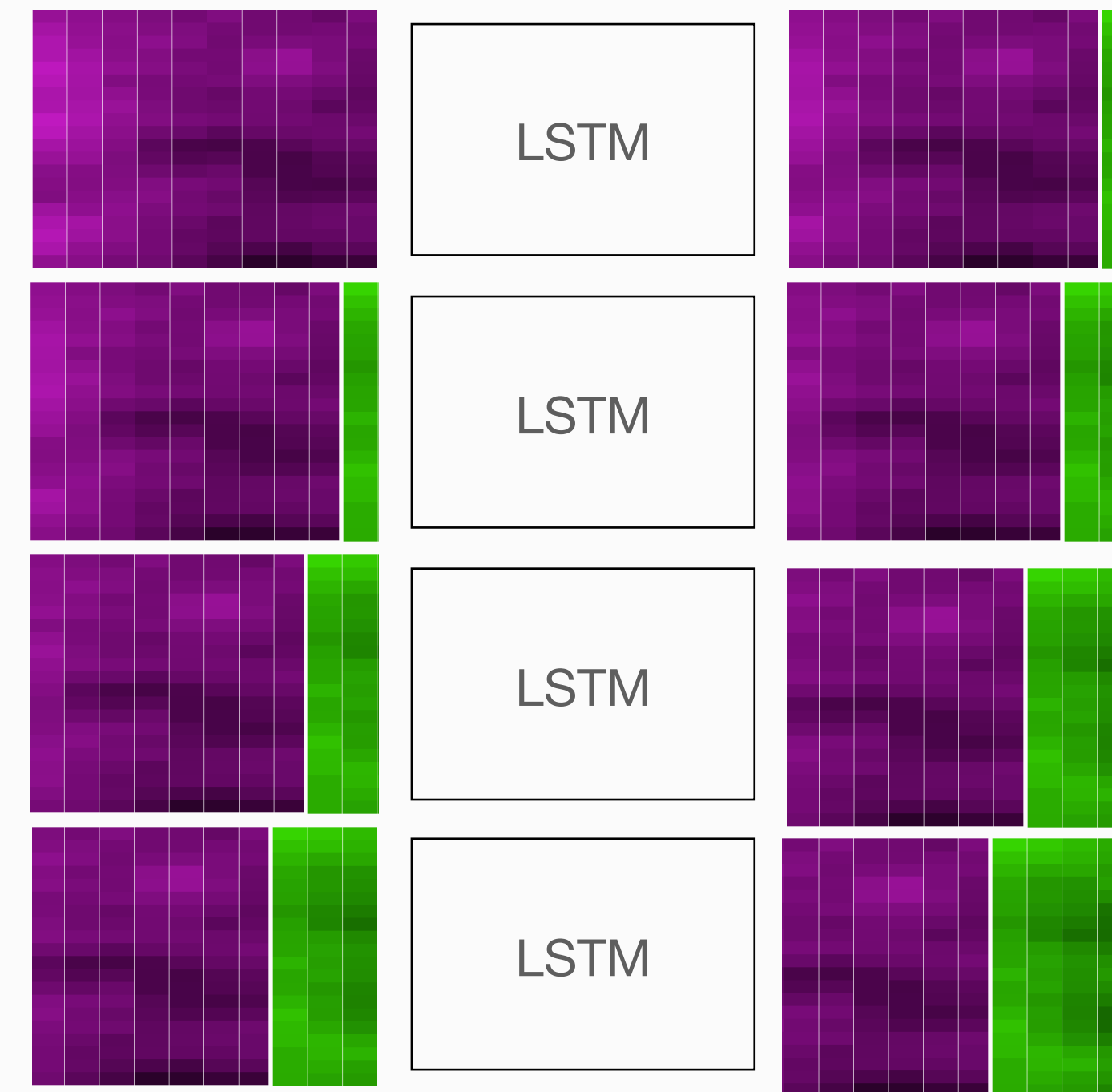Generation techniques involve **autoregressive generation** and run in **realtime,** allowing for creative intervention



Auto Regressive Generation

Datasets can be as small as **30 seconds**. Training a model for **300 iterations** can allow for good results and can take **3 minutes** on free cloud platforms such as Google Colab, or personal computers such as the M1 Macbook Pro

ual: creative computing institute

# Controllable Palette



This allows for musicians to not only be **very specific** in the audio palette they are looking to use, they also only need to find a **small amount** of it

# Interactive Training

They can also train a model, **listen to the output**, see how it fits into the music they are working on and **update their dataset** or training parameters and embark again at little cost.

Train models
https://github.com/Louismac/MAGNet

Use in models in realtime with audio reactive drawing with Dorothy

https://github.com/Louismac/dorothy

ual: creative computing
institute

# Thank you

l.mccallum@arts.ac.uk

**ual:** creative computing
institute