# Types of generated sounds

- Speech (Text-to-Speech)

- Music

- Music notes (samples)
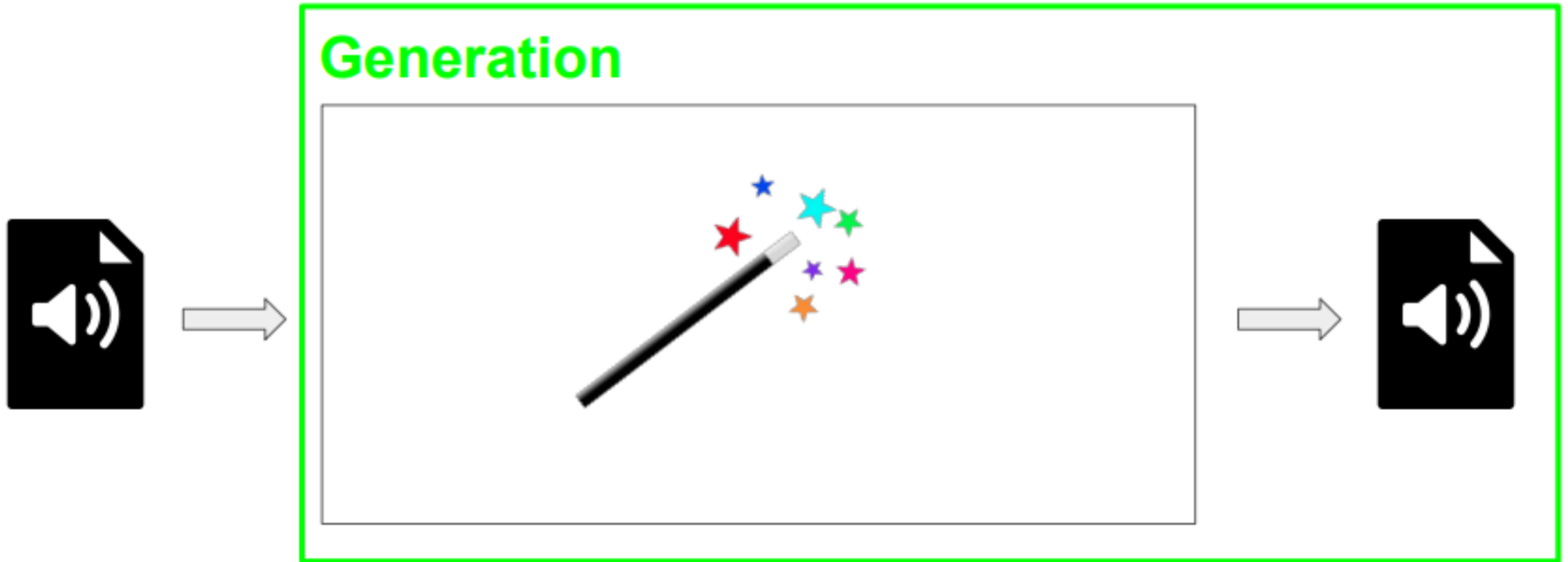
- Sound design

- ...

# Sound representations

- Raw-audio

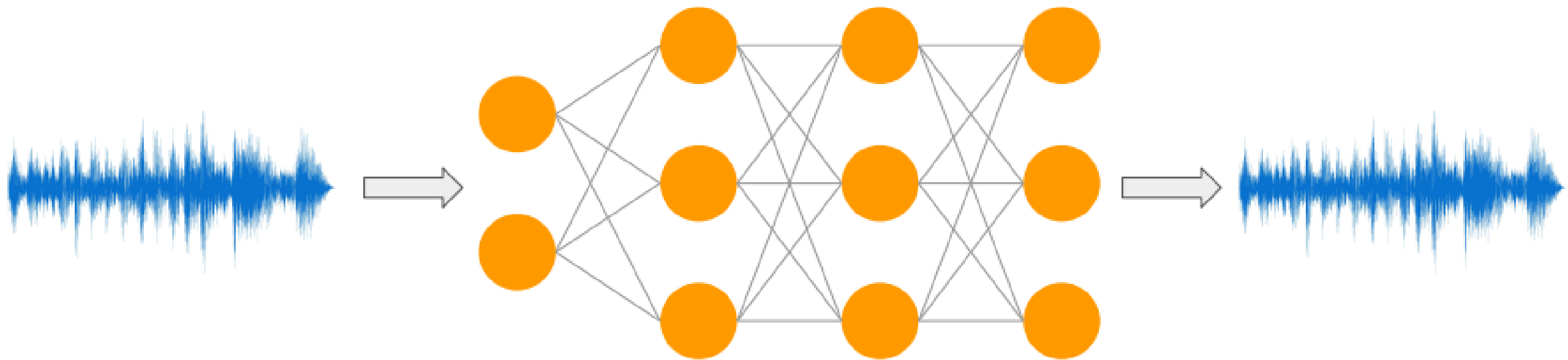- Spectrograms

# Generation from raw audio: Challenges

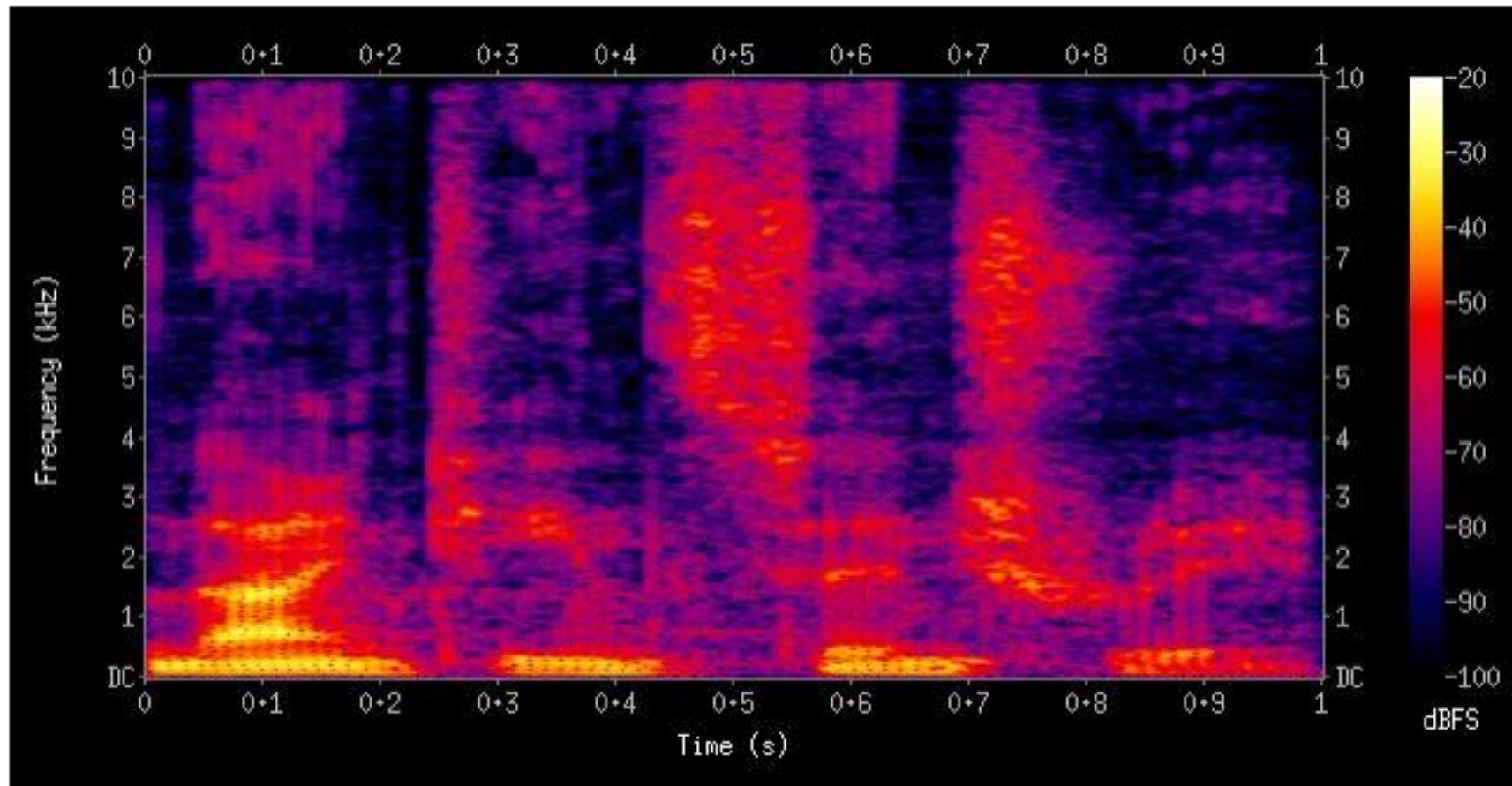- Difficult to capture long-range dependencies

Pitch    Melody

Structure
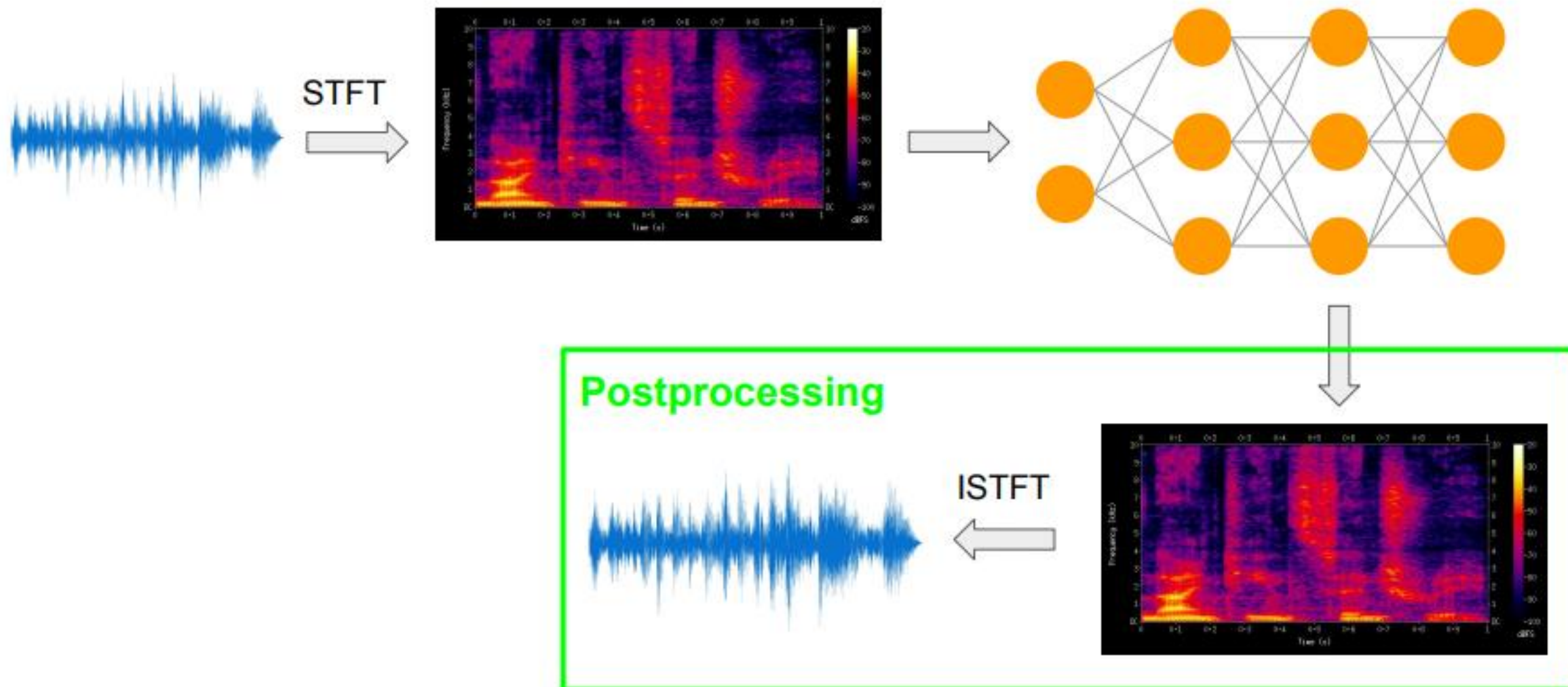
Rhythm    Timbre

Harmony

# Sound generation task

# Generation from raw audio

# Use a more compact representation of sound

# Generation from spectrograms



**Postprocessing**

ISTFT
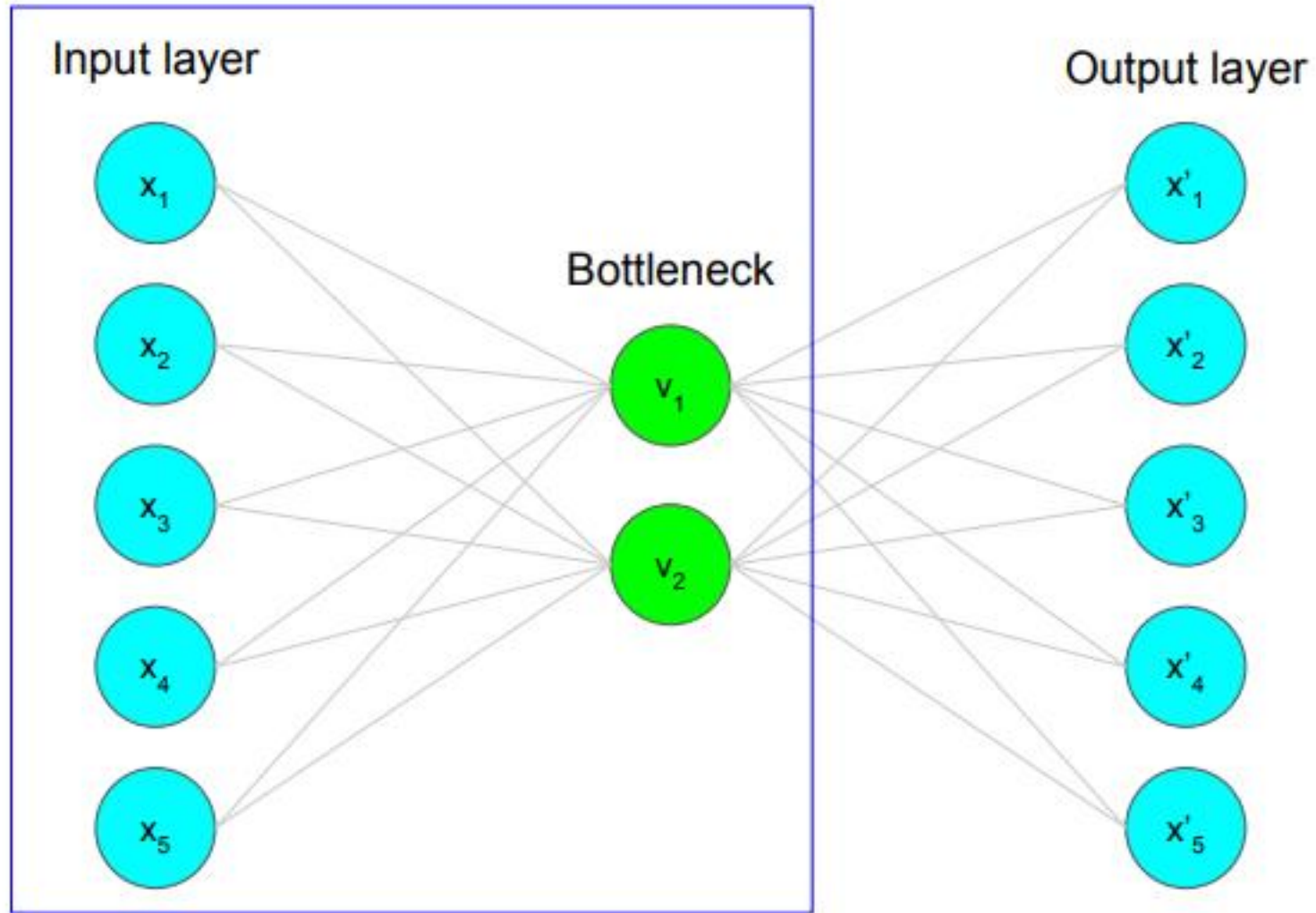
# Autoencoders: The sneaky idea

Create an architecture with a bottleneck, which ensures a lower-dimensional representation of the original data.

# Autoencoder = Encoder + Decoder

**Encoder** = compress data into lower-dimensional representation (*latent* space)

**Decoder** = Decompress representation back to original domain

Input layer

Bottleneck

Output layer

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$

$v_1$ $v_2$

$x'_1$ $x'_2$ $x'_3$ $x'_4$ $x'_5$

**Original data**

**Reconstruction**

https://github.com/musikalkemist/generating-sound-with-neural-networks

# Generation with AEs

# Generation with VAEs

# This content of this slides was obtained by the following repository:

https://github.com/musikalkemist/generating-sound-with-neural-networks