

# Upsampling artifacts in neural audio synthesis

—  
**Jordi Pons** (@jordiponsdotme – [www.jordipons.me](http://www.jordipons.me))

work with Santiago Pasqual, Giulio Cengarle and Joan Serrà

[arxiv.org/abs/1703.09452](https://arxiv.org/abs/1703.09452)

# MUSIC SOURCE SEPARATION IN THE WAVEFORM DOMAIN

Alexandre Défossez  
Facebook AI Research  
INRIA  
École Normale Supérieure  
PSL Research University  
defossez@fb.com

Francis Bach  
INRIA  
École Normale Supérieure  
PSL Research University  
francis.bach@en

Nicolas Usunier  
Facebook AI Research  
usunier@fb.com

Léon Bottou  
Facebook AI Research  
leonb@fb.com

# MelGAN: Generative Adversarial Networks for Conditional Waveform Synthesis

Lyrebird AI  
Kundan Kumar

## SEGAN: Speech Enhancement Generative Adversarial Network

Santiago Pascual<sup>1</sup>, Antonio Bonafonte<sup>1</sup>, Joan Serra<sup>2</sup>

<sup>1</sup>Universitat Politècnica de Catalunya, Barcelona, Spain

<sup>2</sup>Telefónica Research, Barcelona, Spain

santi.pascual@upc.edu, antonio.bonafonte@upc.edu joan.serra@telefonica.com

### Abstract

Current speech enhancement techniques

is not im

ch enhancement [17]. However, further significant improvements of speech enhancement are needed when a clean phase spectrum is available. A deep network

Rithesh Kumar\*  
Lyrebird AI  
rithesh@descript.com

Jose Sotelo  
Lyrebird AI, Mila

Program Co-director

# Conv-TasNet: Surpassing Ideal Time-Frequency Magnitude Masking for Speech Separation

Yi Luo, Nima Mesgarani

# WAVE-U-NET: A MULTI-SCALE NEURAL NETWORK FOR END-TO-END AUDIO SOURCE SEPARATION

Daniel Stoller  
Queen Mary University of London  
d.stoller@qmul.ac.uk

Sebastian Ewert

Simon Dixon  
Queen Mary University of London  
s.e.dixon@qmul.ac.uk

What do these architectures have in common?

Abstract—Single-channel, source separation methods have recently achieved high accuracy, latency, and computational efficiency. The main limitation of these methods is the main insufficient. The major contribution of this paper is to formulate the separation problem in terms of the decoupling of the phase and magnitude of the mixed signal.

accuracy of the mask estimation method and the mask estimation method. The proposed method is evaluated using the inverse

investigate end-to-end source separation for the spectral phase estimation. The proposed method is evaluated using the inverse

ch has several limitations. It is dependent on many parameters, such as the number of audio frames, which can affect the frequency resolution.

# MUSIC SOURCE SEPARATION IN THE WAVEFORM DOMAIN

Alexandre Défossez  
Facebook AI Research  
INRIA  
École Normale Supérieure  
PSL Research University  
defossez@fb.com

Francis Bach  
INRIA  
École Normale Supérieure  
PSL Research University  
francis.bach@en

Nicolas Usunier  
Facebook AI Research  
usunier@fb.com

Léon Bottou  
Facebook AI Research  
leonb@fb.com

# MelGAN: Generative Adversarial Networks for Conditional Waveform Synthesis

Lyrebird AI  
Kundan Kumar

## SEGAN: Speech Enhancement Generative Adversarial Network

Santiago Pascual<sup>1</sup>, Antonio Bonafonte<sup>1</sup>, Joan Serra<sup>2</sup>

<sup>1</sup>Universitat Politècnica de Catalunya, Barcelona, Spain

<sup>2</sup>Telefónica Research, Barcelona, Spain

santi.pascual@upc.edu, antonio.bonafonte@upc.edu joan.serra@telefonica.com

### Abstract

Current speech enhancement techniques

is not im

ch enhancement [17]. However, further significant improvements of speech enhancement are needed when a clean phase spectrum is available. A deep network

Rithesh Kumar\*  
Lyrebird AI  
rithesh@descript.com

Jose Sotelo  
Lyrebird AI, Mila

Program Co-director

# Conv-TasNet: Surpassing Ideal Time-Frequency Magnitude Masking for Speech Separation

Yi Luo, Nima Mesgarani

# WAVE-U-NET: A MULTI-SCALE NEURAL NETWORK FOR END-TO-END AUDIO SOURCE SEPARATION

Daniel Stoller  
Queen Mary University of London  
d.stoller@qmul.ac.uk

Sebastian Ewert

Simon Dixon  
Queen Mary University of London  
s.e.dixon@qmul.ac.uk

They use neural upsamplers

Abstract—Single-channel, source separation methods have recently achieved high accuracy, latency, and computational efficiency. The main limitation of these methods is the main insufficient. The major contribution of this paper is to formulate the separation problem in terms of the decoupling of the phase and magnitude of the mixed signal. This is achieved by representing the mixed signal in the frequency domain and then applying a mask estimation method and the mask estimation method to the mixed signal. The proposed method is evaluated using the inverse Fourier transform.

accuracy of the mask estimation method and the mask estimation method. The proposed method is evaluated using the inverse Fourier transform.

investigate end-to-end source separation for the spectral front-end. There are several limitations. The proposed method is evaluated using the inverse Fourier transform.

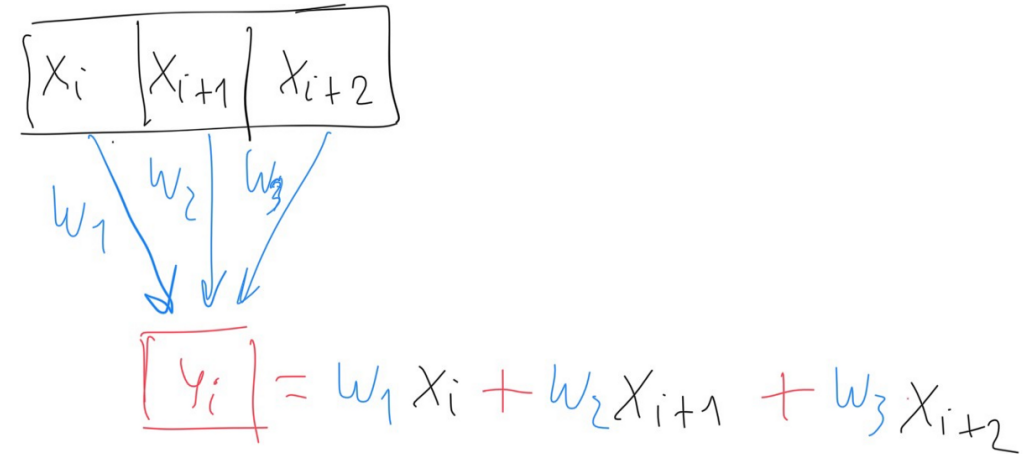
ch has several limitations. The proposed method is evaluated using the inverse Fourier transform.

# Main upsamplers:

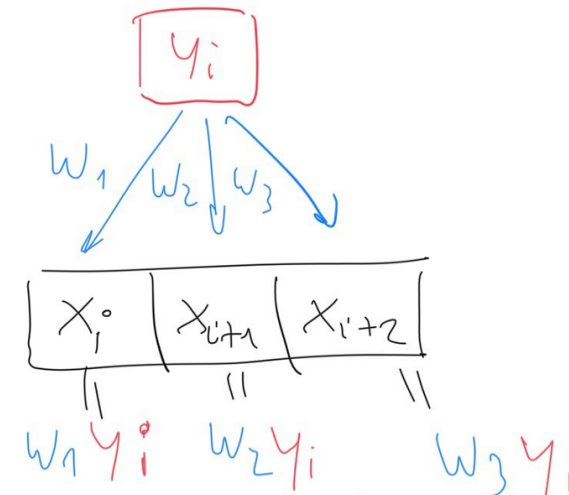
## Transposed convolutions

- Widely used

## Convolutions ("collapse")



## Transposed convolutions ("expand")





# Main upsamplers:

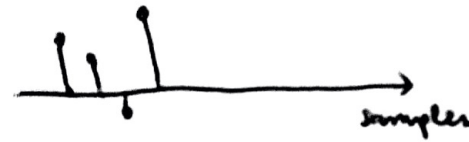
## Transposed convolutions

- Widely used

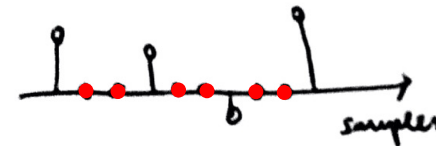
## Interpolation + convolution

- Often-times used

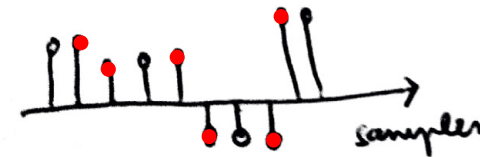
ORIGINAL SIGNAL



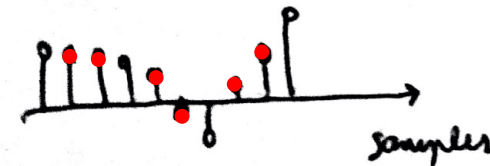
UPSAMPLE  $\times 3$



STRETCH  
INTERPOLATION



NEAREST NEIGHBOR  
INTERPOLATION



LINEAR  
INTERPOLATION

# Main upsamplers:

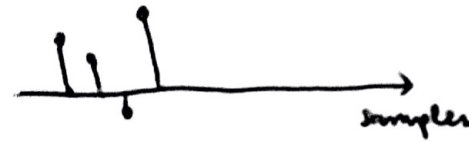
## Transposed convolutions

- Widely used

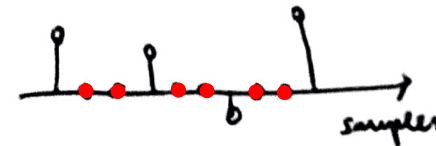
## Interpolation + convolution

- Often-times used

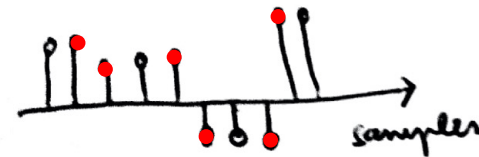
ORIGINAL SIGNAL



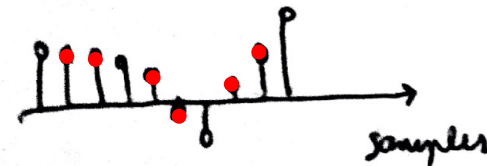
UPSAMPLE  $\times 3$



STRETCH INTERPOLATION + CONVOLUTION



NEAREST NEIGHBOR INTERPOLATION + CONVOLUTION



LINEAR INTERPOLATION + CONVOLUTION

# Main upsamplers:

## Transposed convolutions

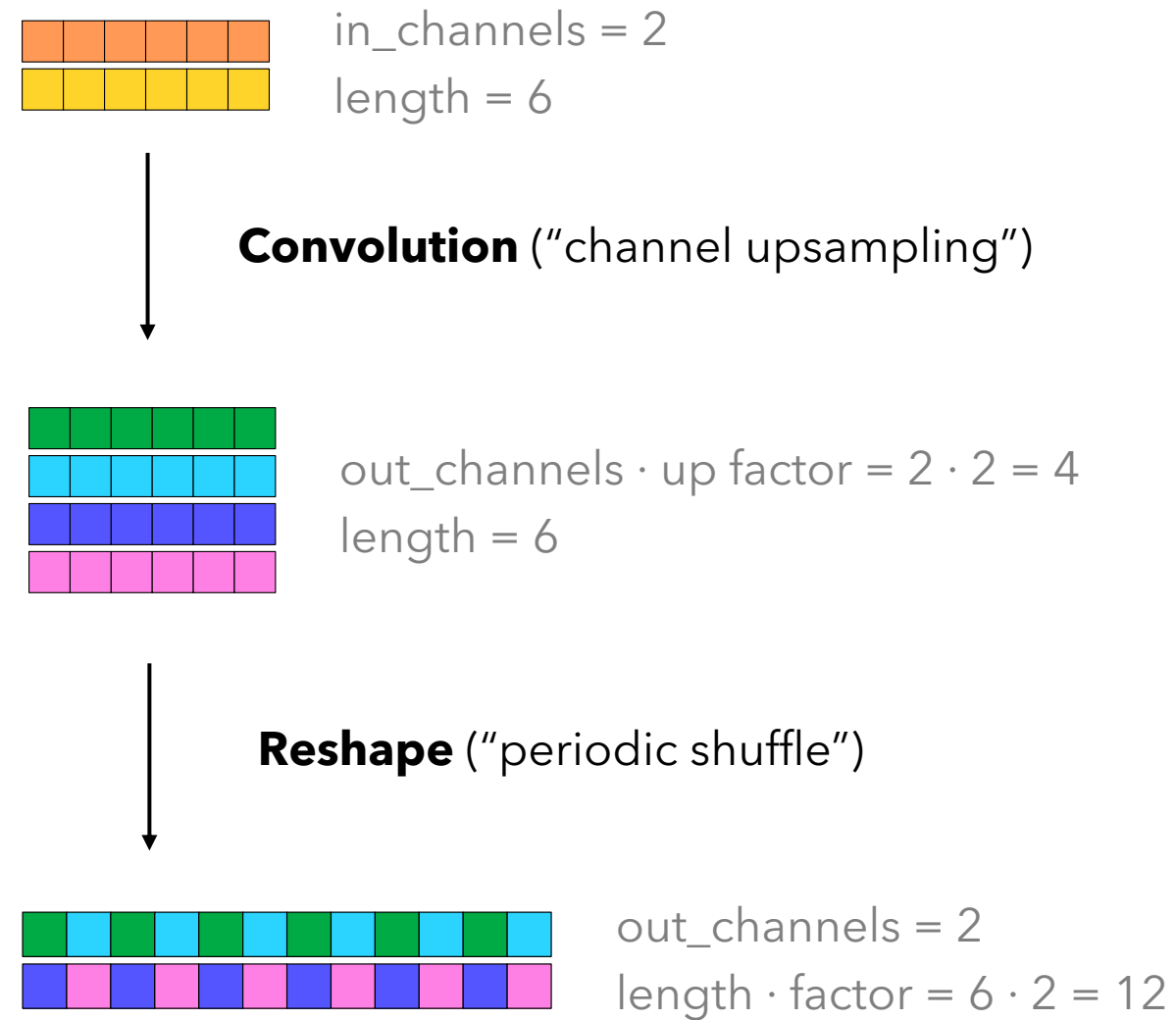
- Widely used

## Interpolation + convolution

- Often-times used

## Subpixel convolutions

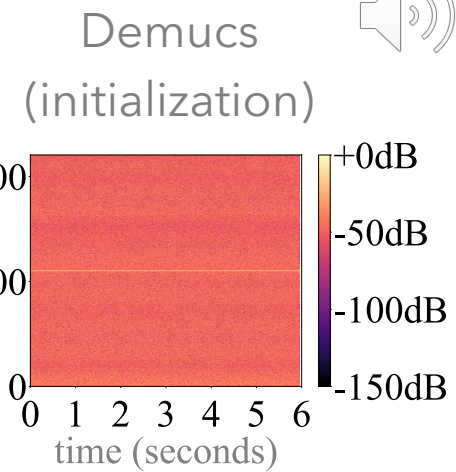
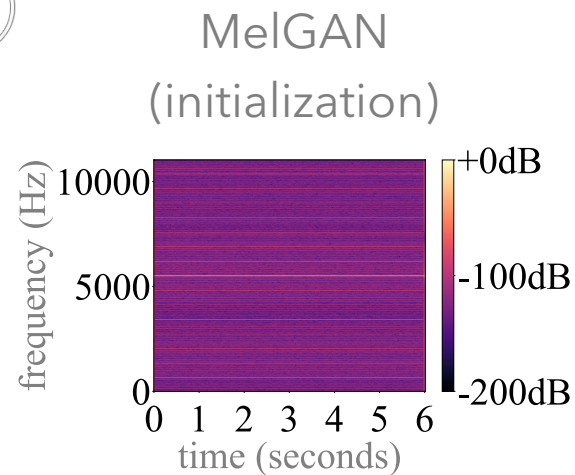
- Rarely used



# Upsampling artifacts:

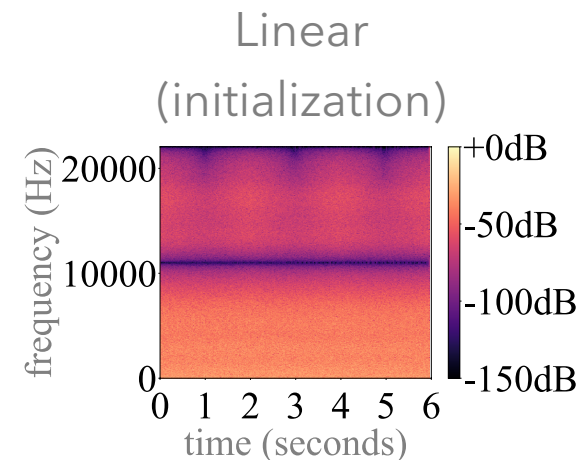
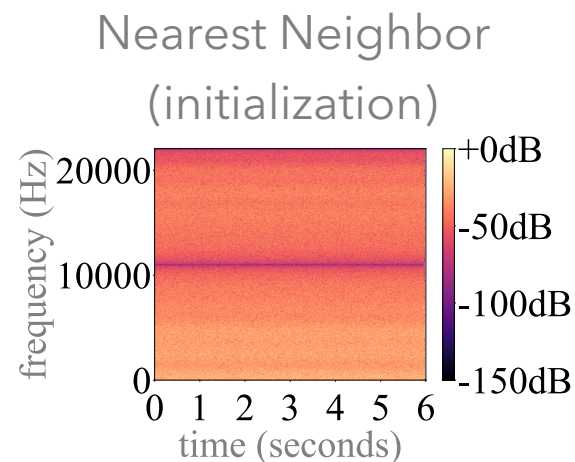
**Transposed convolutions**

Tonal artifacts



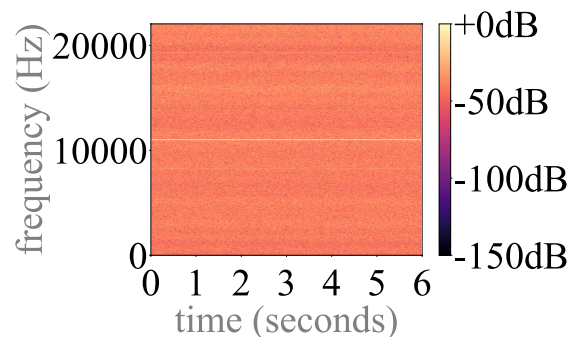
**Interpolation + convolution**

Filtering artifacts



**Subpixel convolutions**

Tonal artifacts



Subpixel Convolution  
(initialization)

# Agenda:

## **Transposed convolutions**

- Why do they introduce tonal artifacts?

## **Interpolation + convolution**

- Why do they introduce filtering artifacts?

## **Subpixel convolutions**

- Why do they introduce tonal artifacts?

## **Artifacts due to spectral replicas**

- Signal processing perspective

## **The role of training**

- Learning from data reduces artifacts



- 

# Transposed convolutions

# Transposed convolutions: tonal artifacts

## Main sources of tonal artifacts:

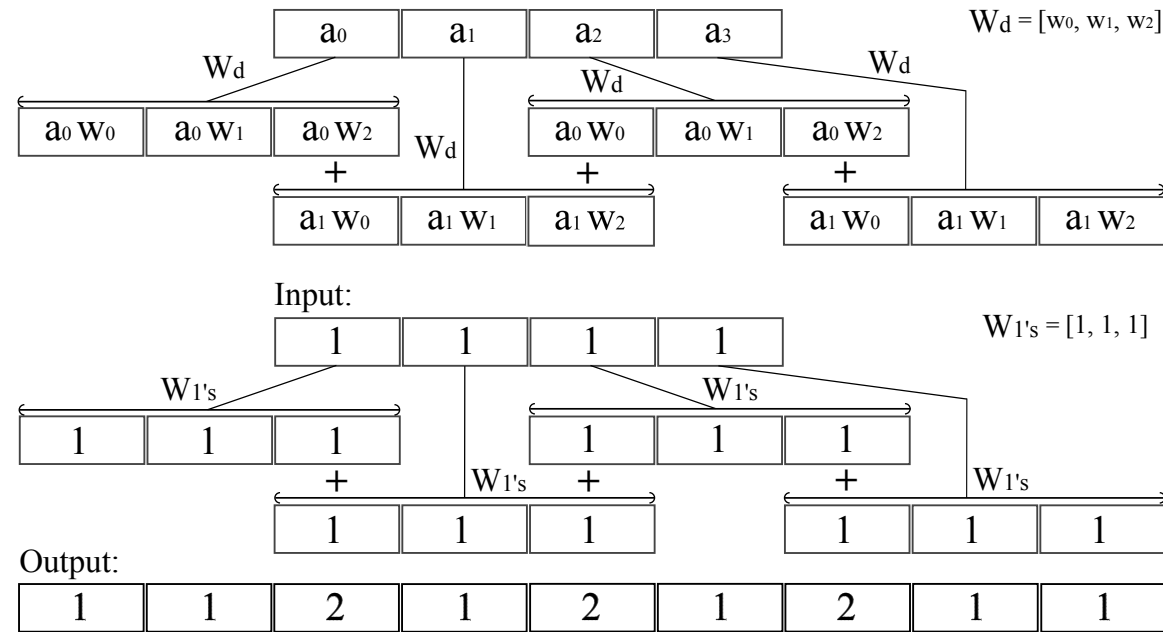
- Weights initialization
- Overlap issues

## Examples for discussion:

- No-overlap: stride = length.
- Partial-overlap: length is *not* a multiple of stride.
  - *Example: filter length = 5, and stride = 4.*
- Full-overlap: length is a multiple of stride.
  - *Example: filter length = 8, and stride = 4.*

Odena et al., 2016: "Deconvolution and Checkerboard Artifacts" in Distill.

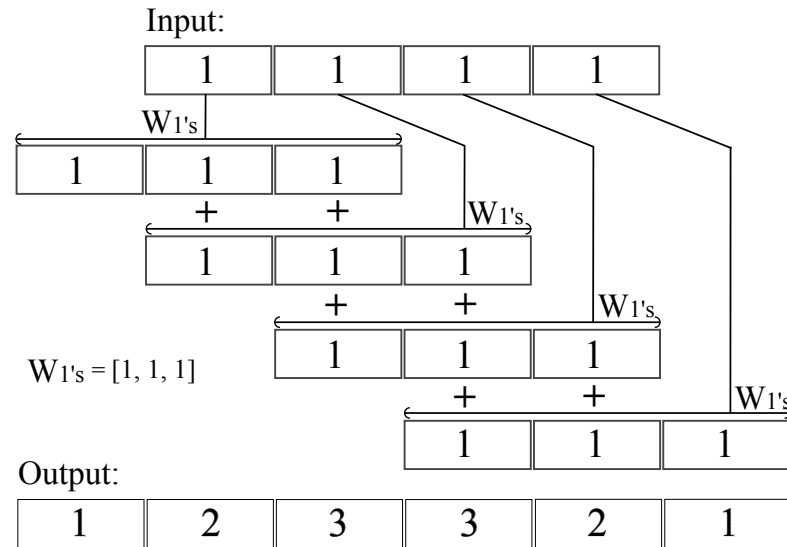
# Transposed convolutions: partial-overlap case



Example: length=3, stride=2

Note the periodicities due to **overlap issues**

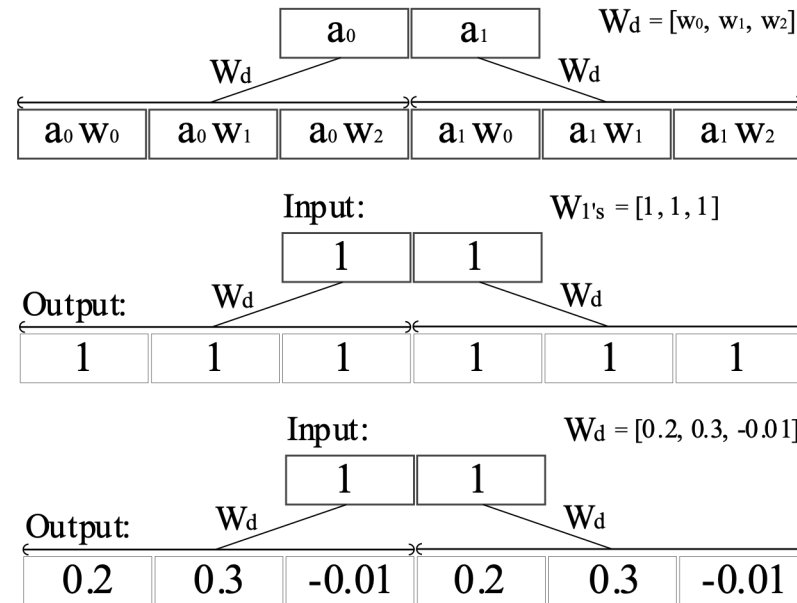
# Transposed convolutions: full-overlap case



Example: length=3, stride=1

**NO periodicities** due to constant overlap

# Transposed convolutions: no-overlap case



Example: length=3, stride=3.

**NO periodicities** due to overlap  
Note the **weights initialization** issue



# Transposed convolutions

## Main sources of tonal artifacts:

- Weights initialization
- Overlap issues

## Transposed convolution categories:

- No-overlap: stride = length.
- Partial-overlap: length is *not* a multiple of stride.
- Full-overlap: length is a multiple of stride.

## Important remark:

Even though you solve the overlap issue, the weights initialization issue remains due to random initialization!

# Agenda:

## **Transposed convolutions**

- ~~Why do they introduce tonal artifacts?~~

## **Interpolation + convolution**

- Why do they introduce filtering artifacts?

## **Subpixel convolutions**

- Why do they introduce tonal artifacts?

## **Artifacts due to spectral replicas**

- Signal processing perspective

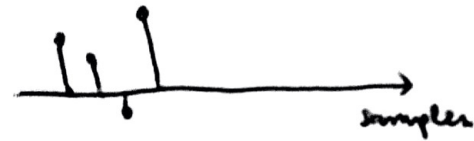
## **The role of training**

- Learning from data reduces artifacts

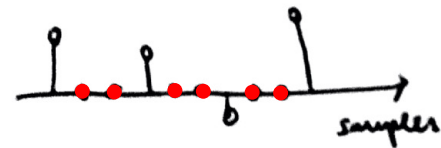
- # Interpolation + convolution

# Interpolation + convolution: filtering artifacts

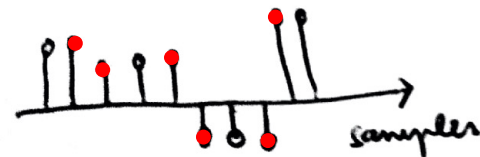
ORIGINAL SIGNAL



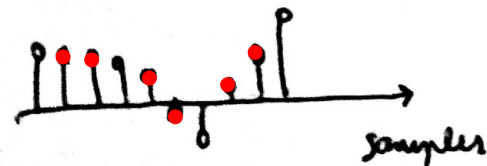
UPSAMPLE  $\times 3$



STRETCH INTERPOLATION + CONVOLUTION



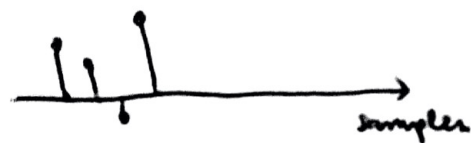
NEAREST NEIGHBOR INTERPOLATION + CONVOLUTION



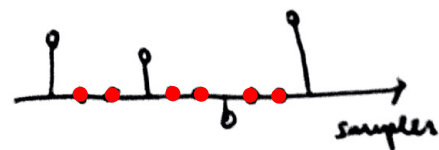
LINEAR INTERPOLATION + CONVOLUTION

# Interpolation: stretch + (non-learnable) convolution

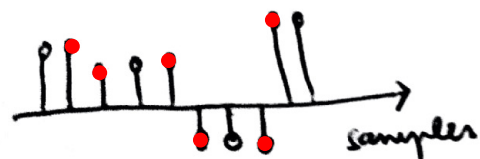
ORIGINAL SIGNAL



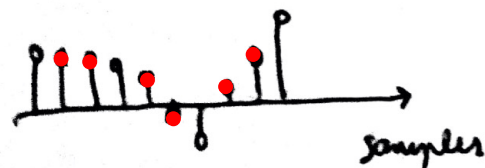
UPSAMPLE  $\times 3$



STRETCH INTERPOLATION + CONVOLUTION



NEAREST NEIGHBOR INTERPOLATION + CONVOLUTION



LINEAR INTERPOLATION + CONVOLUTION

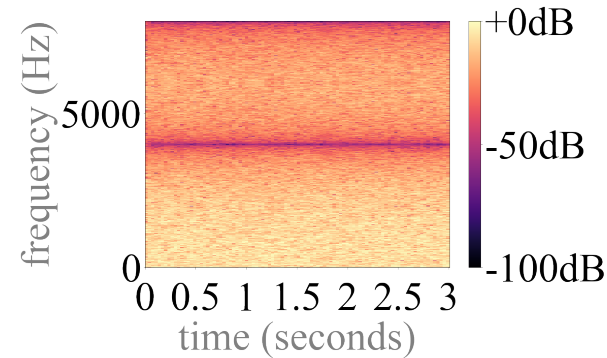


**TABLE 3.1**  
Short Table of Fourier Transforms

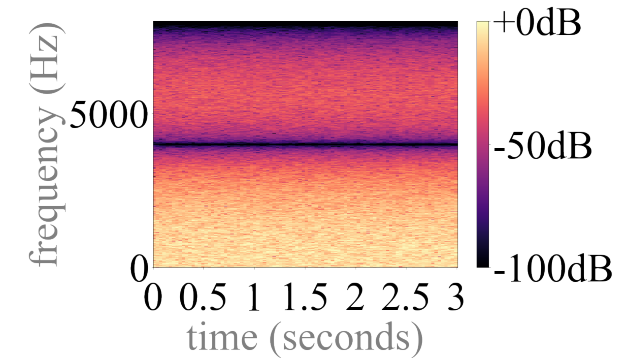
$g(t)$	$G(f)$	
1 $e^{-at}u(t)$	$\frac{1}{a + j2\pi f}$	$a > 0$
2 $e^{at}u(-t)$	$\frac{1}{a - j2\pi f}$	$a > 0$
3 $e^{-a t }$	$\frac{2a}{a^2 + (2\pi f)^2}$	
4 $te^{-at}u(t)$	$\frac{1}{(a + j2\pi f)^2}$	$a > 0$
5 $t^n e^{-at}u(t)$	$\frac{n!}{(a + j2\pi f)^{n+1}}$	
6 $\delta(t)$	$\delta(f)$	
7 $1$	$\delta(f)$	
8 $e^{j2\pi f_0 t}$	$\delta(f - f_0)$	
9 $\cos 2\pi f_0 t$	$0.5 [\delta(f + f_0) + \delta(f - f_0)]$	
10 $\sin 2\pi f_0 t$	$j0.5 [\delta(f + f_0) - \delta(f - f_0)]$	
11 $u(t)$	$\frac{1}{j2\pi f} + \frac{1}{2}\delta(f)$	
12 $\text{sgn } t$	$\frac{2}{j2\pi f}$	
$\cos 2\pi f_0 t u(t)$	$\frac{1}{4} [\delta(f - f_0) + \delta(f + f_0)] + \frac{j2\pi f}{(2\pi f_0)^2 - (2\pi f)^2}$	
14 $\sin 2\pi f_0 t u(t)$	$\frac{1}{4j} [\delta(f - f_0) - \delta(f + f_0)] + \frac{2\pi f_0}{(2\pi f_0)^2 - (2\pi f)^2}$	
15 $e^{-at} \sin 2\pi f_0 t u(t)$	$\frac{2\pi f_0}{(a + j2\pi f)^2 + 4\pi^2 f_0^2}$	
16 $e^{-at} \cos 2\pi f_0 t u(t)$	$\frac{a + j2\pi f}{(a + j2\pi f)^2 + 4\pi^2 f_0^2}$	
$\Pi\left(\frac{t}{\tau}\right)$	$\tau \text{sinc}(\pi f \tau)$	
18 $2B \text{sinc}(2\pi Bt)$	$\Pi\left(\frac{f}{2B}\right)$	
19 $\Delta\left(\frac{t}{\tau}\right)$	$\frac{\tau}{2} \text{sinc}^2\left(\frac{\pi f \tau}{2}\right)$	
20 $B \text{sinc}^2(\pi Bt)$	$\Delta\left(\frac{f}{2B}\right)$	
21 $\sum_{n=-\infty}^{\infty} \delta(t - nT)$	$f_0 \sum_{n=-\infty}^{\infty} \delta(f - nf_0)$	
22 $e^{-t^2/2\sigma^2}$	$\sigma\sqrt{2\pi}e^{-2(\sigma\pi f)^2}$	

$NW \rightarrow \text{rect} \xrightarrow{F} \text{sinc}(\cdot) \rightarrow \text{plot of sinc function}$   
 $LINEAR \rightarrow \text{tri} \xrightarrow{F} \text{sinc}^2(\cdot) \rightarrow \text{plot of sinc squared function}$

Upsample white noise at 4kHz by 4

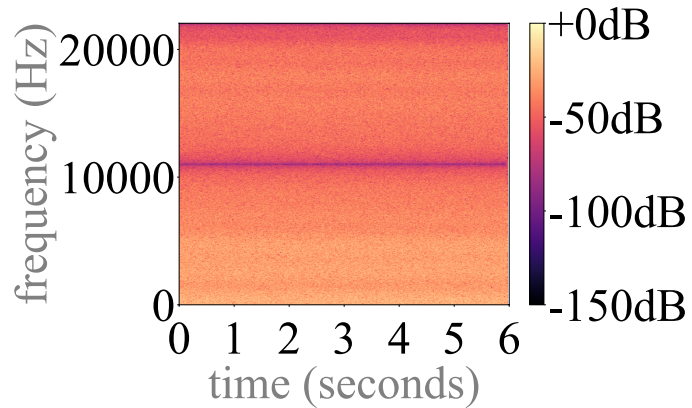


Nearest Neighbor

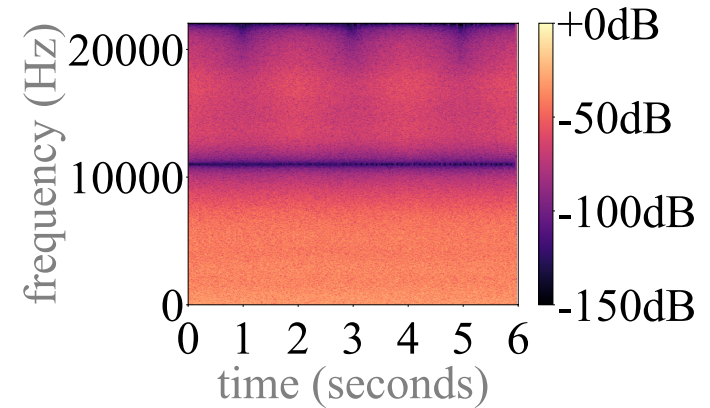


Linear interpolation

# Interpolation + convolution: filtering artifacts



Demucs-like  
Nearest Neighbor  
(initialization)



Demucs-like  
Linear interpolation  
(initialization)

## Important remark:

Filtering artifacts emerge because the frequency response of each interpolation colors the signal.

# Agenda:

## **Transposed convolutions**

- ~~Why do they introduce tonal artifacts?~~

## **Interpolation + convolution**

- ~~Why do they introduce filtering artifacts?~~

## **Subpixel convolutions**

- Why do they introduce tonal artifacts?

## **Artifacts due to spectral replicas**

- Signal processing perspective

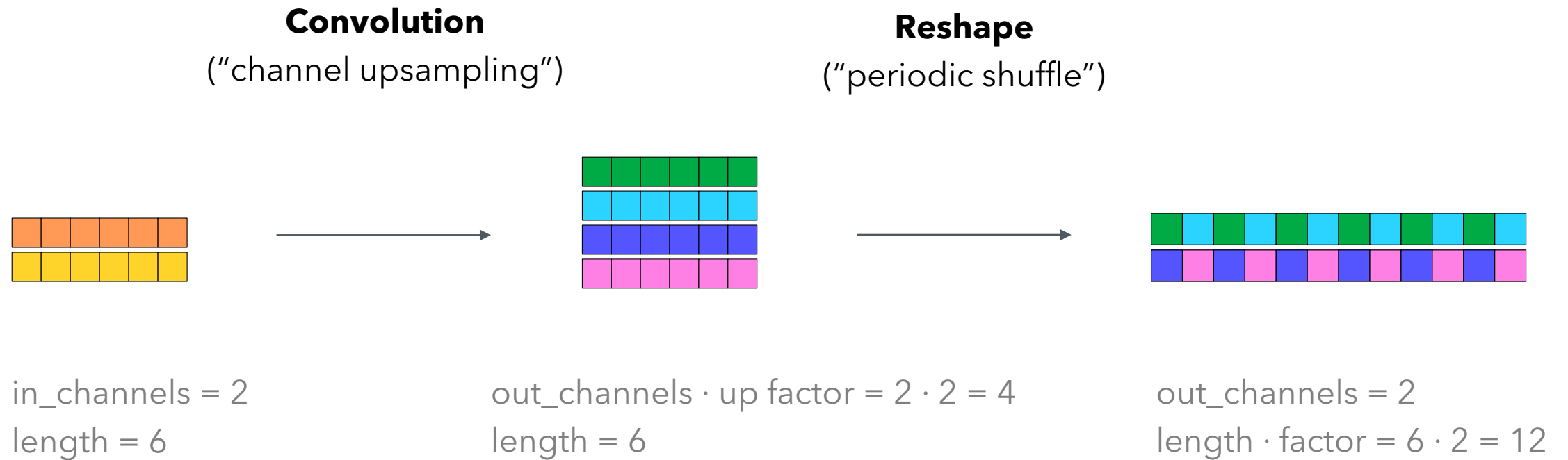
## **The role of training**

- Learning from data reduces artifacts

- 

# Subpixel convolution

# Subpixel convolution: tonal artifacts





# Agenda:

## **Transposed convolutions**

~~- Why do they introduce tonal artifacts?~~

## **Interpolation + convolution**

~~- Why do they introduce filtering artifacts?~~

## **Subpixel convolutions**

~~- Why do they introduce tonal artifacts?~~

## **Artifacts due to spectral replicas**

- Signal processing perspective

## **The role of training**

- Learning from data reduces artifacts

# — Artifacts due to spectral replicas

# Signal processing review

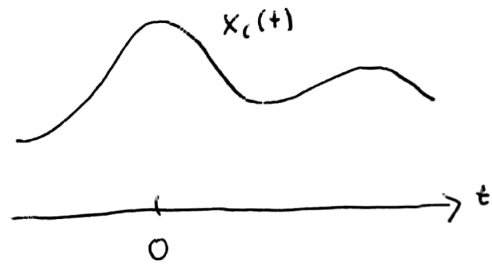
## **IDEA 1:**

Spectral replicas emerge when sampling/discretizing a signal!

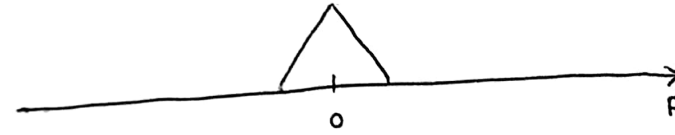
## **IDEA 2:**

When upsampling, one performs bandwidth extension – be aware of spectral replicas!

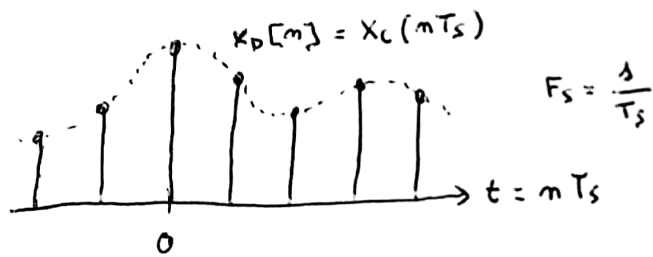
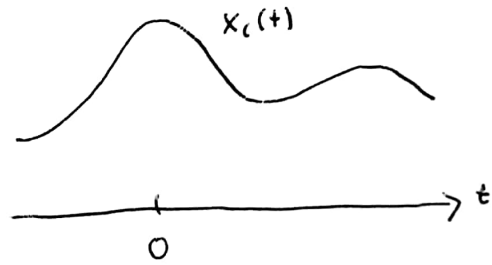
**Time domain**



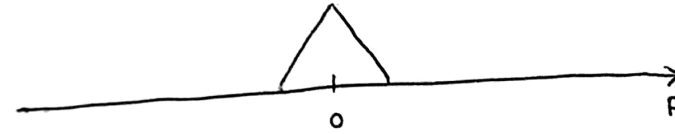
**Frequency domain**



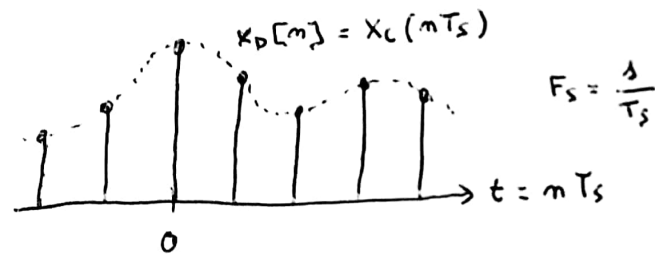
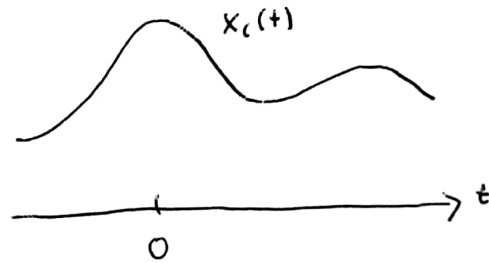
### Time domain



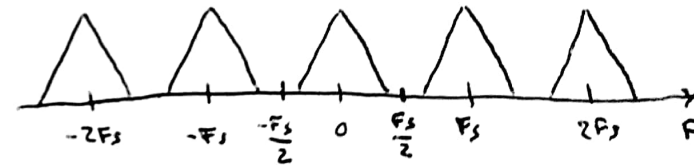
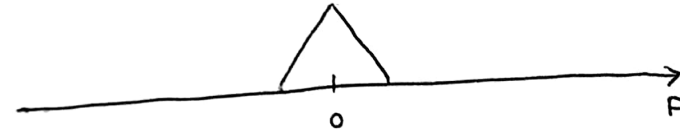
### Frequency domain



Time domain



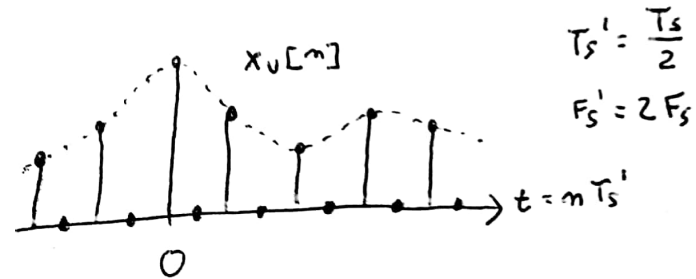
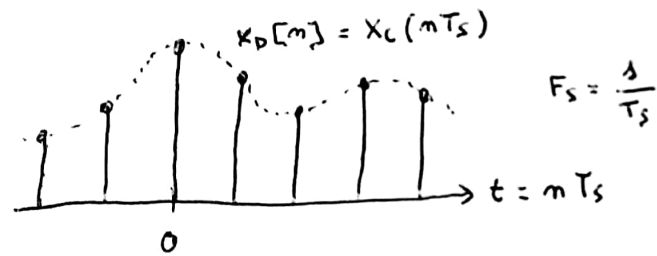
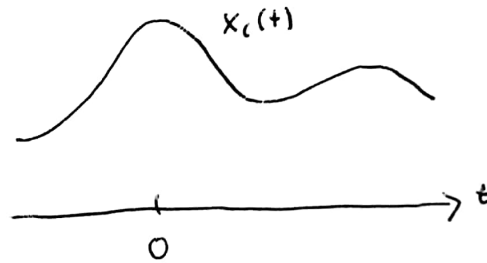
Frequency domain



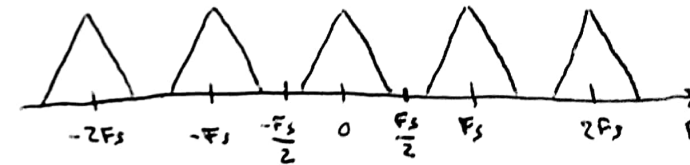
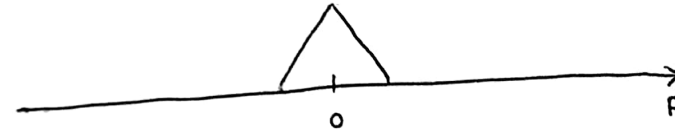
**IDEA 1:**

Spectral replicas emerge when sampling/discretizing a signal!

Time domain



Frequency domain

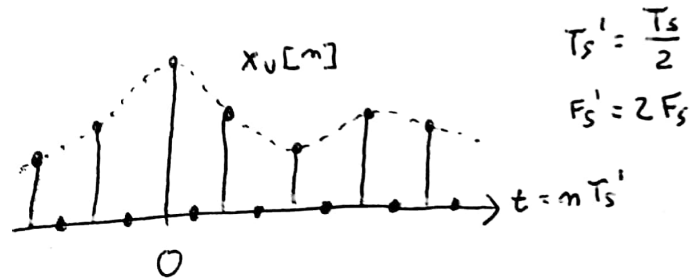
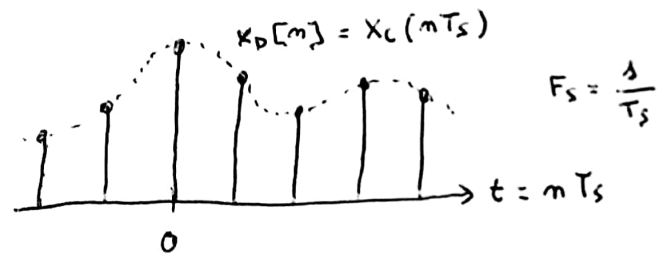
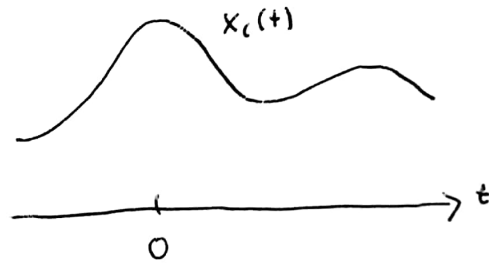


**Stretch interpolation x2**  
(upsampling with zeros)

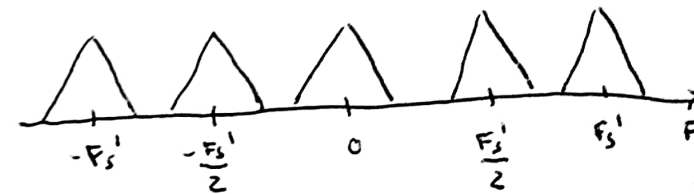
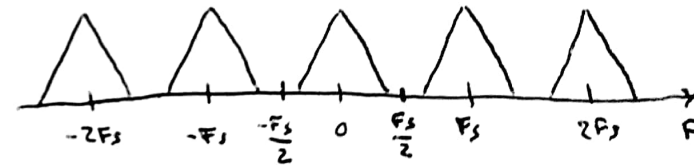
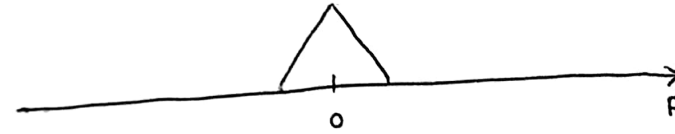
### IDEA 1:

Spectral replicas emerge when sampling/discretizing a signal!

Time domain



Frequency domain



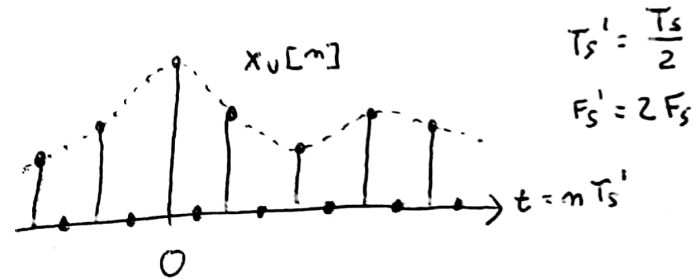
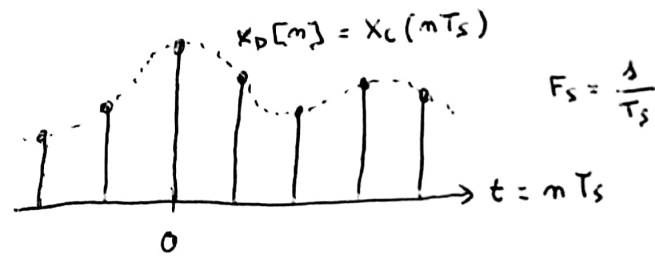
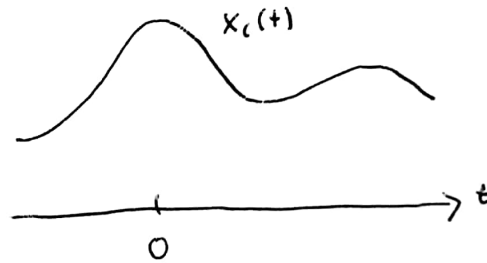
**Stretch interpolation x2**  
(upsampling with zeros)

**IDEA 1:**

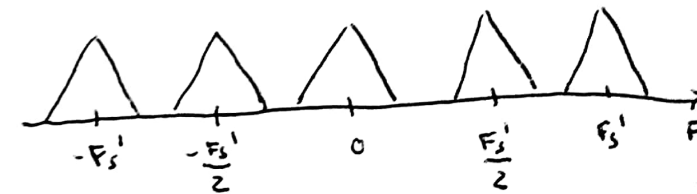
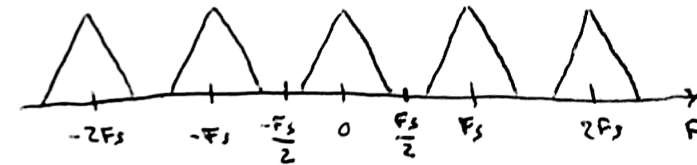
Spectral replicas emerge when sampling/discretizing a signal!



### Time domain



### Frequency domain



#### IDEA 1:

Spectral replicas emerge when sampling/discretizing a signal!

#### IDEA 2:

When upsampling, one performs bandwidth extension - be aware of spectral replicas!

**Stretch interpolation x2**  
(upsampling with zeros)

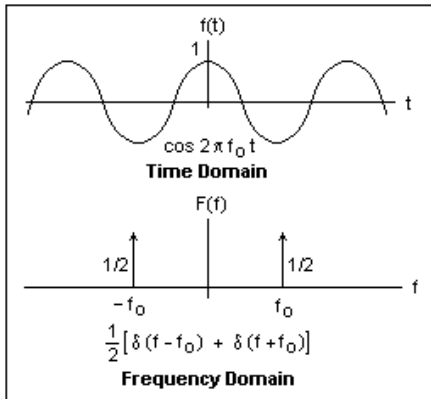
# Signal processing review

## IDEA 1:

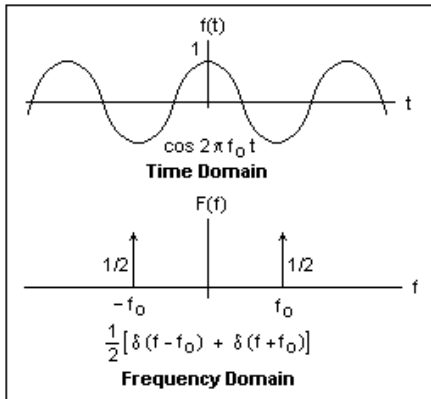
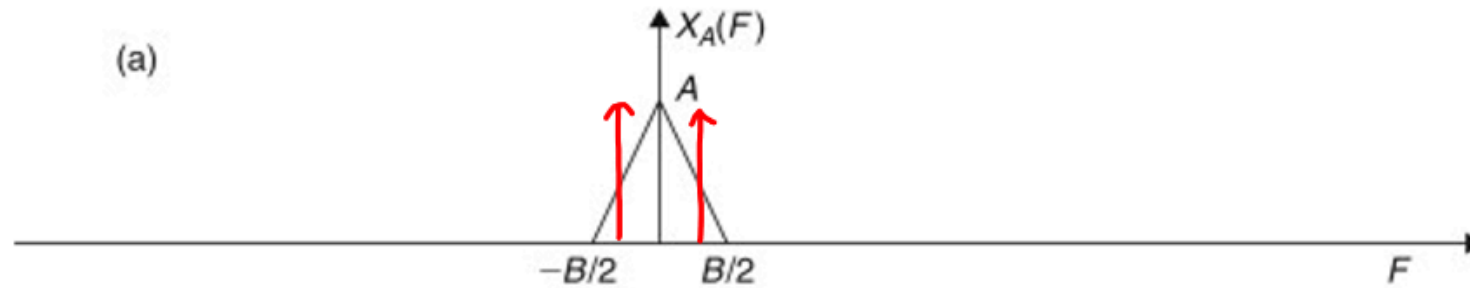
Spectral replicas emerge when sampling/discretizing a signal!

## IDEA 2:

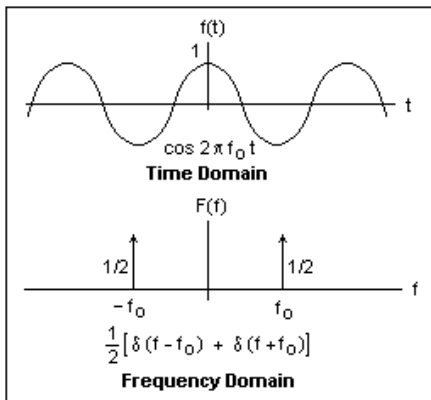
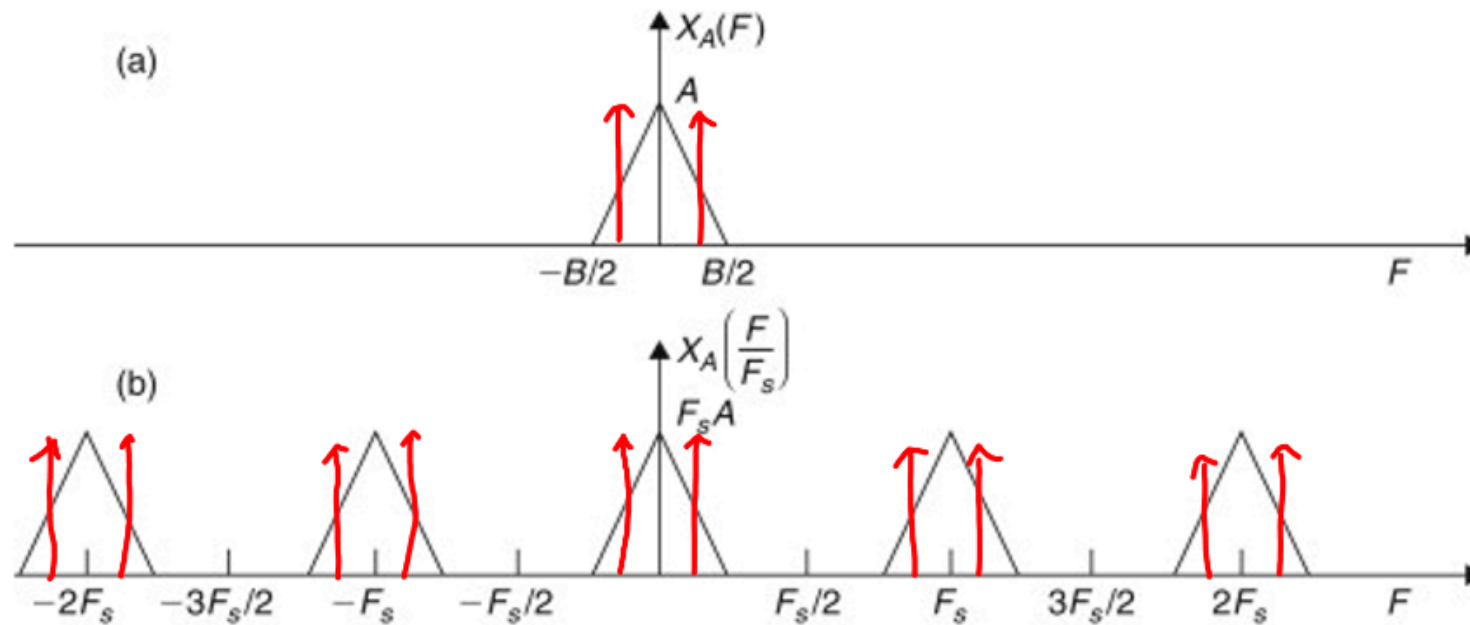
When upsampling, one performs bandwidth extension – be aware of spectral replicas!



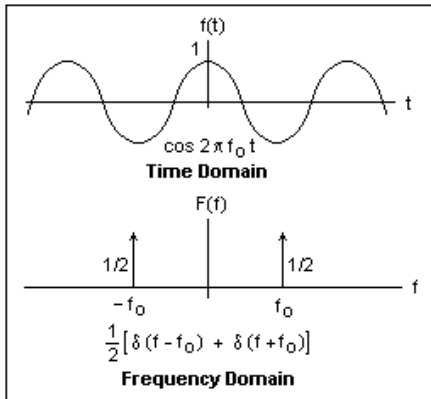
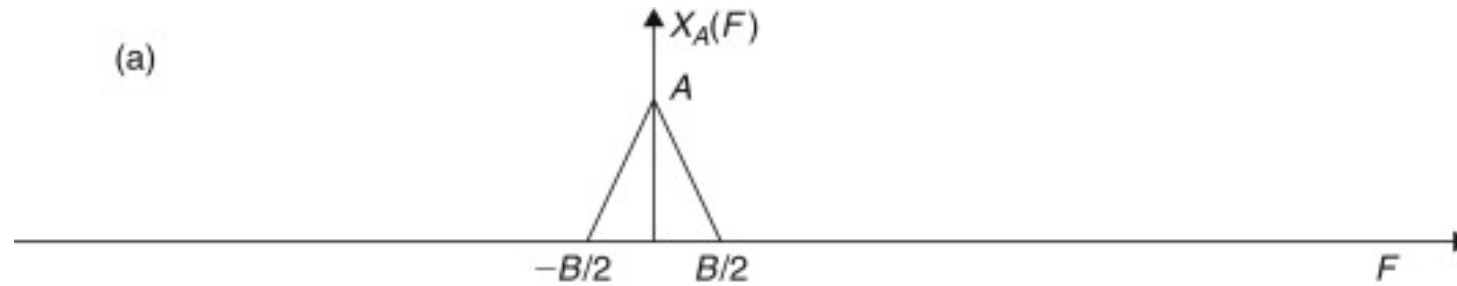
# Spectral replicas of tonal artifacts: multilayered case



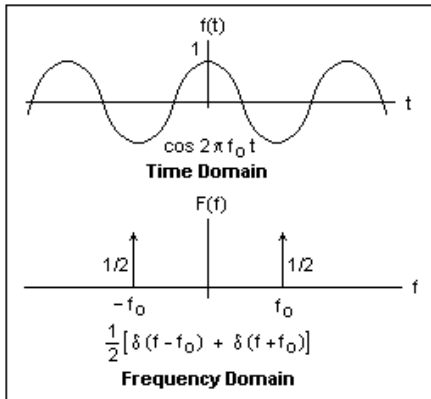
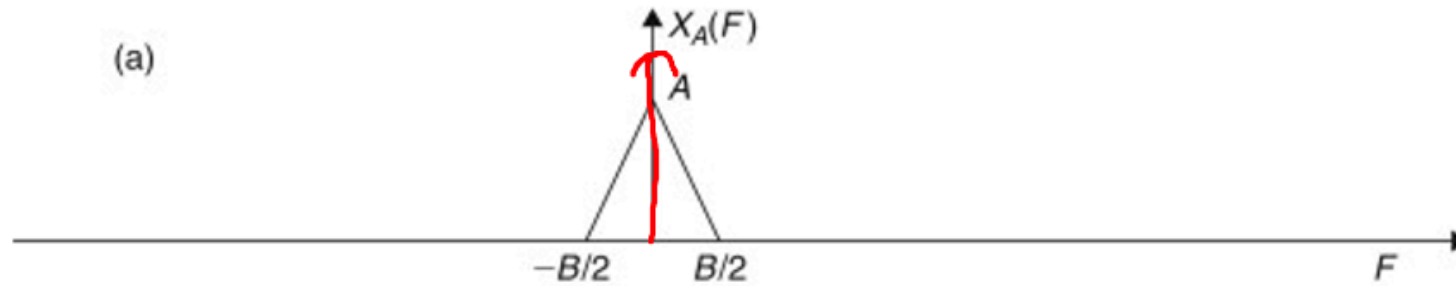
# Spectral replicas of tonal artifacts: multilayered case



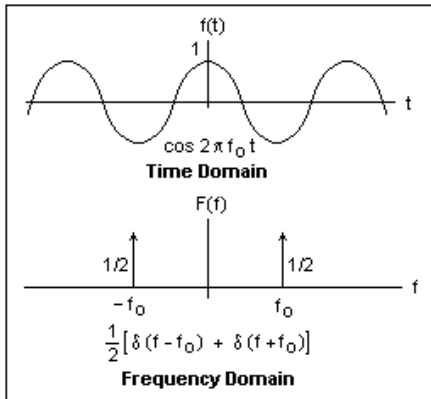
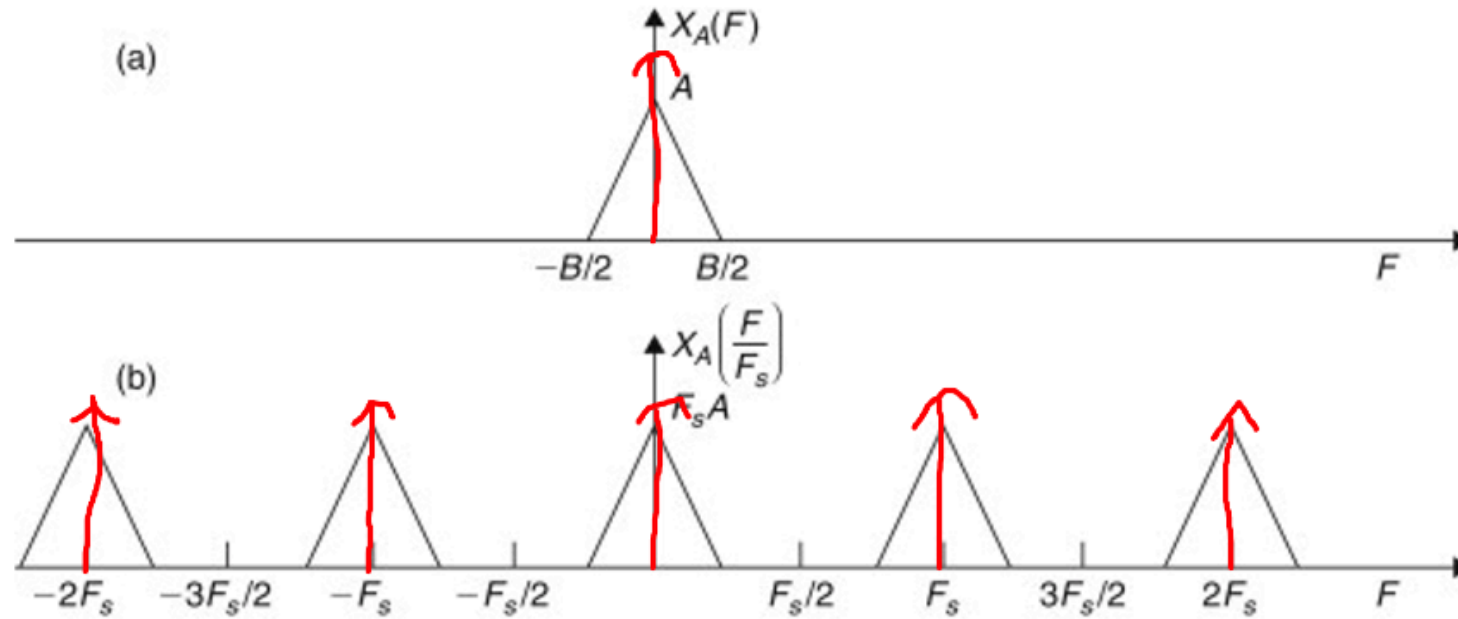
# Spectral replicas of signal offsets



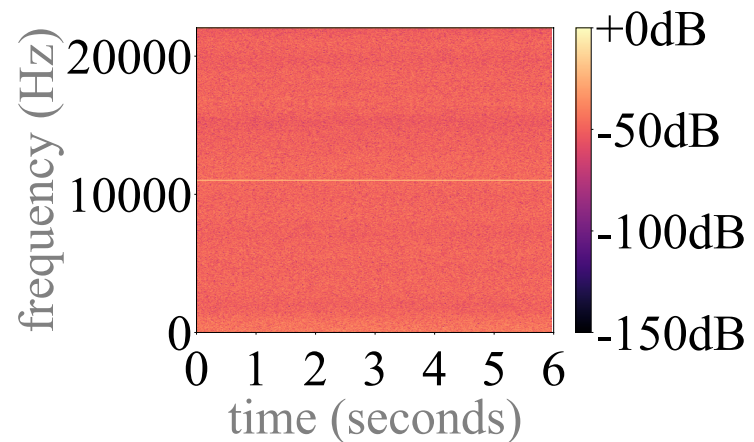
# Spectral replicas of signal offsets



# Spectral replicas of signal offsets

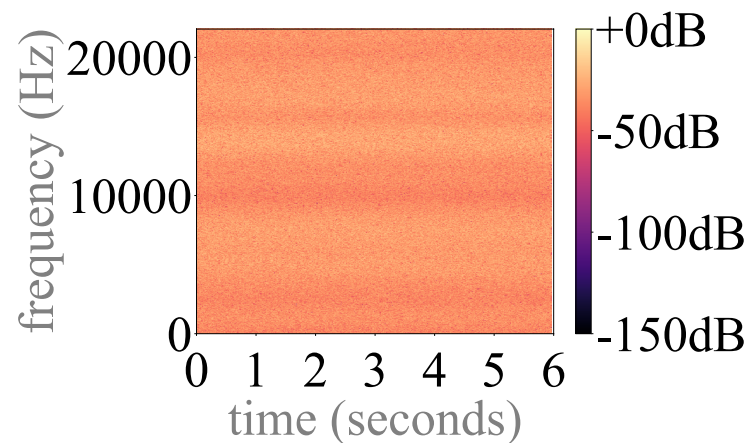
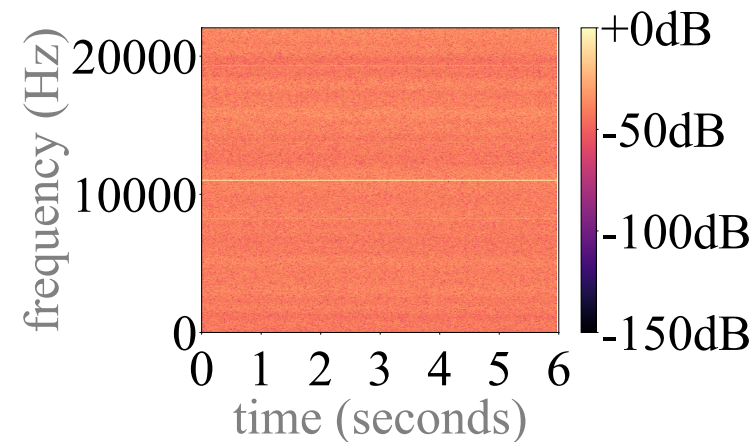


Demucs  
(initialization)

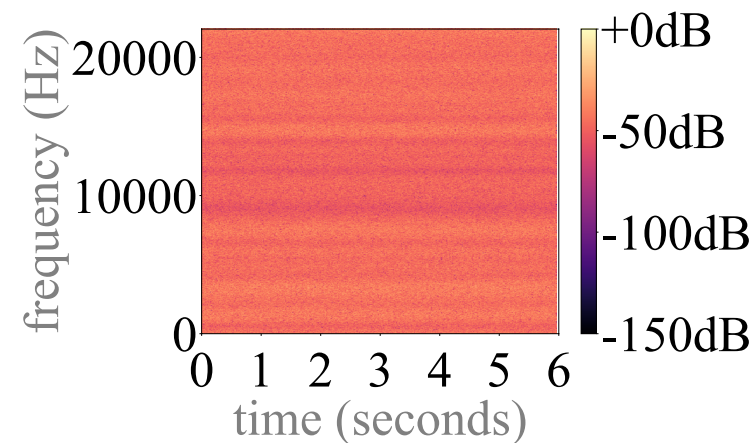


**with**  
biases & ReLUs  
(introduce signal offsets)

Subpixel Convolution  
(initialization)



**without**  
biases & ReLUs





# Artifacts due to spectral replicas

## **Additional sources of upsampling artifacts:**

- Spectral replicas of tonal artifacts
- Spectral replicas of filtering artifacts
- Spectral replicas of signal offsets

# Agenda:

## **~~Transposed convolutions~~**

- ~~- Why do they introduce tonal artifacts?~~

## **~~Interpolation + convolution~~**

- ~~- Why do they introduce filtering artifacts?~~

## **~~Subpixel convolutions~~**

- ~~- Why do they introduce tonal artifacts?~~

## **~~Artifacts due to spectral replicas~~**

- ~~- Signal processing perspective~~

## **The role of training**

- Learning from data reduces artifacts

# — The role of training

# Is training dealing with the problematic initializations?

Music source separation (MUSDB [28] benchmark)	SDR $\uparrow$	epoch	#parm
Demucs-like: transposed CNN (full-overlap)	5.35	319 s	703M
Demucs-like: nearest neighbor interpolation	5.17	423 s	716M
Demucs-like: linear interpolation	4.62	430 s	716M
Demucs-like: subpixel CNN	5.38	311 s	729M

Mixture



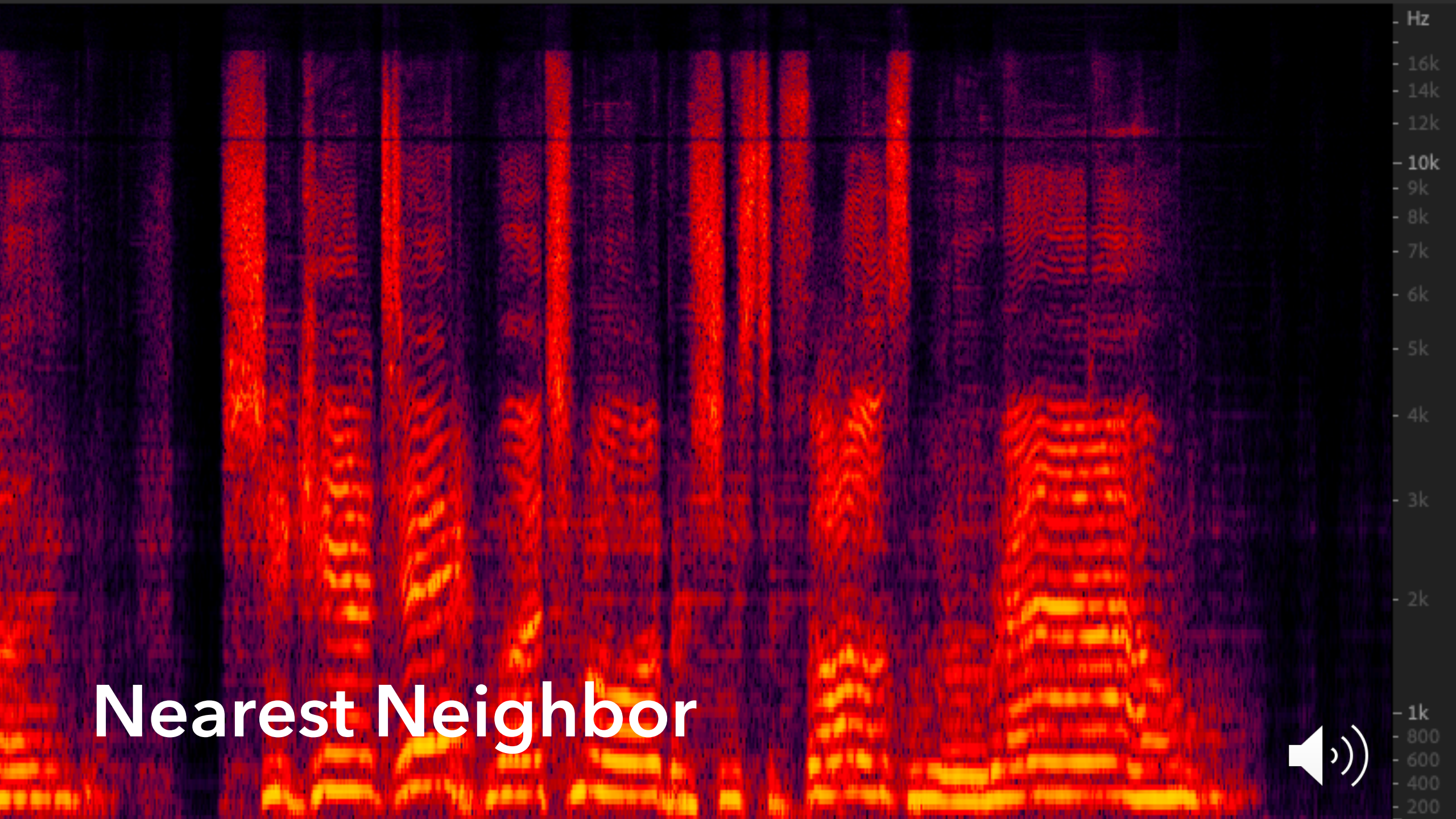
Hz  
16k  
14k  
12k  
10k  
9k  
8k  
7k  
6k  
5k  
4k  
3k  
2k  
1k  
800  
600  
400  
200



# Subpixel convolution







Nearest Neighbor



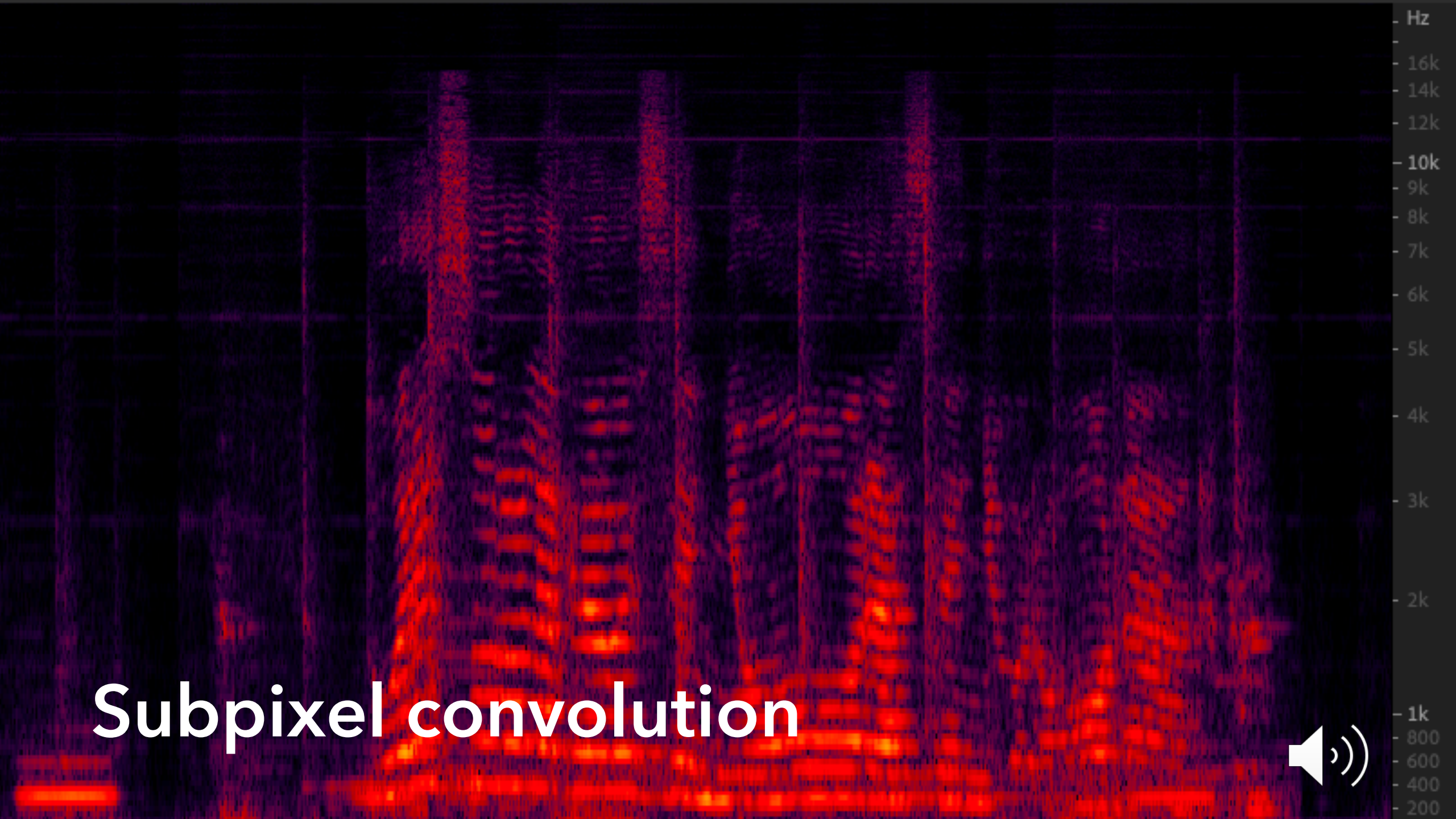


**Mixture: out of distribution**





# Subpixel convolution



Nearest Neighbor



# **The role of training: helps overcoming the noisy initializations to get state-of-the-art results**

## **Formal evaluation:**

- Transposed and subpixel CNNs achieve the best SDR scores.
  - Despite their poor initialization!
- Nearest neighbour upsampler follows closely!

## **Informal listening:**

- Upsampling artifacts can emerge even after training!
  - Tonal artifacts: silent parts and with out-of-distribution data.
  - Filtering artifacts: they are not that perceptually annoying.

# Agenda:

## **~~Transposed convolutions~~**

~~- Why do they introduce tonal artifacts?~~

## **~~Interpolation + convolution~~**

~~- Why do they introduce filtering artifacts?~~

## **~~Subpixel convolutions~~**

~~- Why do they introduce tonal artifacts?~~

## **~~Artifacts due to spectral replicas~~**

~~- Signal processing perspective~~

## **~~The role of training~~**

~~- Learning from data reduces artifacts~~





[arxiv.org/pdf/2010.14356.pdf](https://arxiv.org/pdf/2010.14356.pdf)

@jordiponsdotme - [www.jordipons.me](http://www.jordipons.me)