

Data Augmentation: Improving your Datasets

A decorative graphic on the left side of the slide. It features a large orange circle with a white border at the top left. Below it is a large blue circle containing the word 'MERIDIAN' in white capital letters. The blue circle is filled with a pattern of small, scattered numbers in white and orange. To the right of the blue circle is a series of vertical orange bars of varying heights, resembling a sound wave or a bar chart. Below the blue circle is a light blue circle, and to its left is a small orange circle. To the right of the blue circle is a medium-sized blue circle, and below it is a small orange circle. At the bottom left is a small orange circle.

MERIDIAN

Bruno Padovese

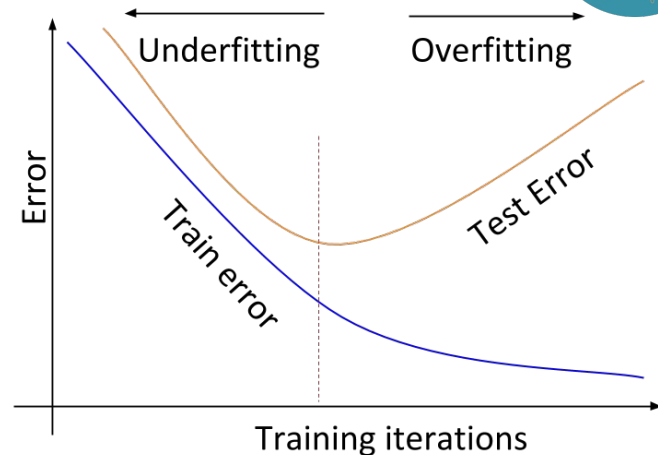
MERIDIAN, Institute for Big Data Analytics,
Dalhousie University, Halifax, Canada

Introduction



DNNs are smart, but not always...

- Tend to overfit the training set
- Can easily inherit and perpetuate biases



filename	sel_id	label	start	end
NOPP6_EST_20090329_084500.wav	0	1	890.100436	893.100436
NOPP6_EST_20090329_090000.wav	0	1	51.413506	54.413506
	1	1	41.592974	44.592974
	2	1	97.386199	100.386199
	3	1	115.234384	118.234384
	4	1	288.680821	291.680821

Benefit from large amounts of labeled data

- Costly to build a dataset
- Domain specific data may require input from experts

DNNs classifiers can often be sensitive to very minor changes

- Ex: Adversarial attacks



“panda”

57.7% confidence

+ .007 ×



noise

=



“gibbon”

99.3% confidence

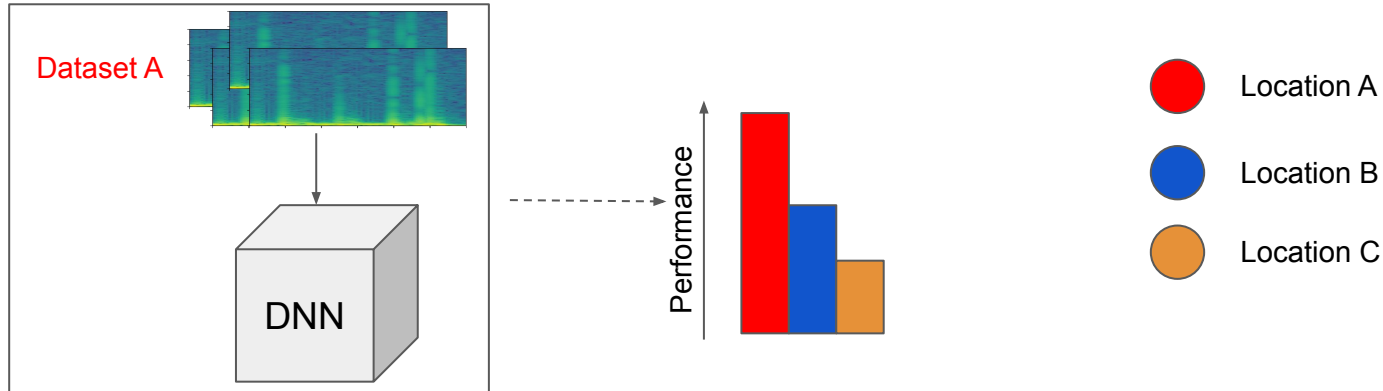
DNNs problems in Underwater Acoustics



It may happen naturally as well:

- Changes in amplitude
- Different types of hydrophones
- Distinct geographic locations

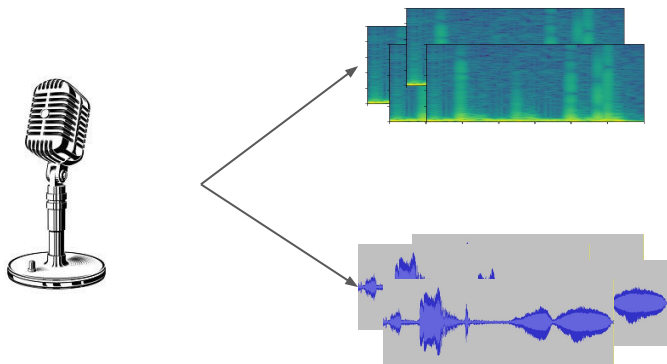
Thus, a model trained in one Dataset may perform poorly when tested in another location



Solutions...?



- More data...
 - It would always be helpful to simply have access to more data from all sources
 - Would require a lot of effort towards collecting and annotating more data



		label		start	end
filename	sel_id				
NOPP6_EST_20090329_084500.wav	0	1	890.100436	893.100436	
NOPP6_EST_20090329_090000.wav	0	1	51.413506	54.413506	
	1	1	41.592974	44.592974	
	2	1	97.386199	100.386199	
	3	1	115.234384	118.234384	
	4	1	288.680821	291.680821	

- What if we could artificially inflate the size of our dataset?

What is data augmentation?

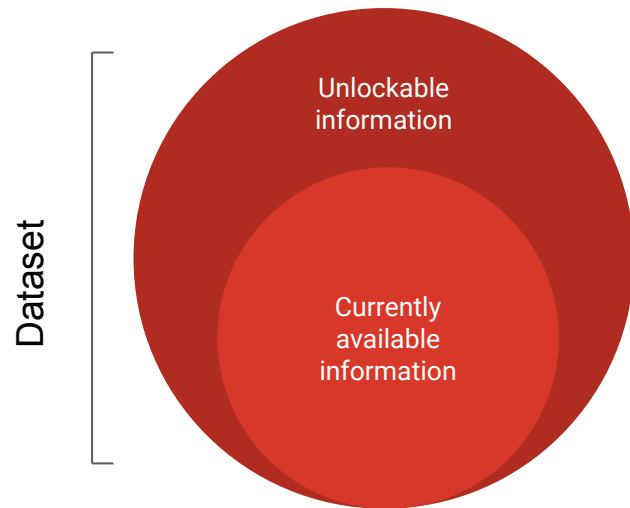


Data Augmentation is a data-space solution to the problem of data limitation

- Suit of techniques that enhance the size and quality of training datasets
- Inexpensive way to acquire more labeled data

Possible techniques include:

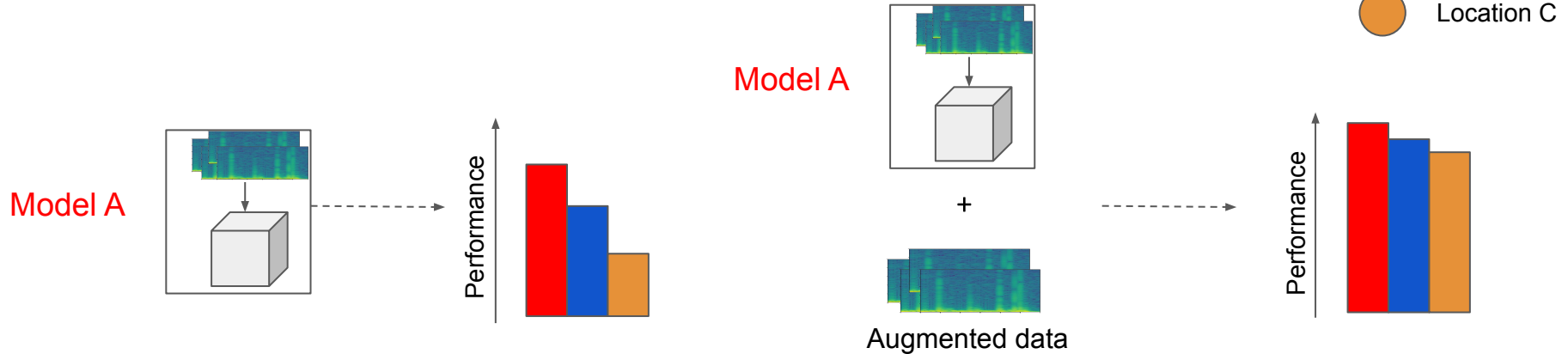
- Geometric transformations
- Color space augmentations
- Mixing
- Random erasing
- Deep learning based methods



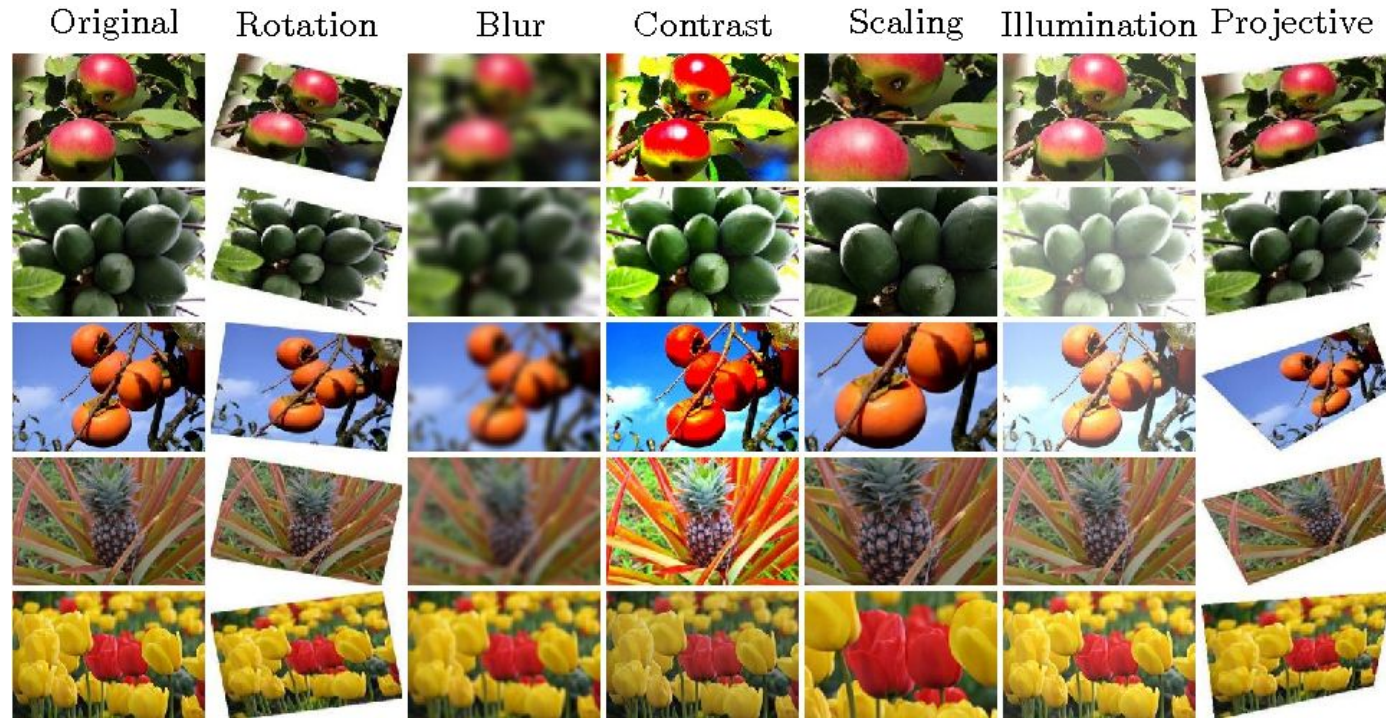
What is data augmentation?



- Data Augmentation can:
 - Enhance the quality of training sets
 - Leading to better classifiers
 - Expands model adaptability to other contexts (e.g. Other geographic locations)



Simple Data Augmentation



Source: Data Augmentation for Plant Classification

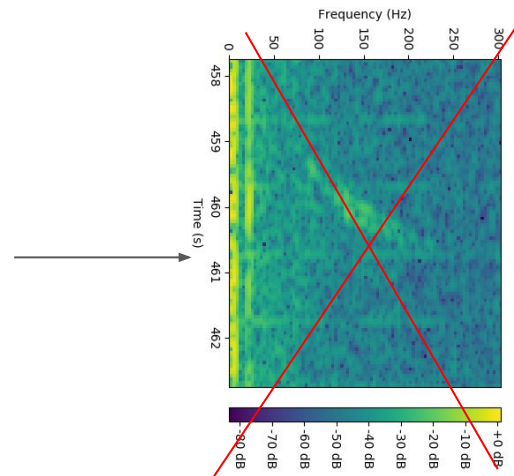
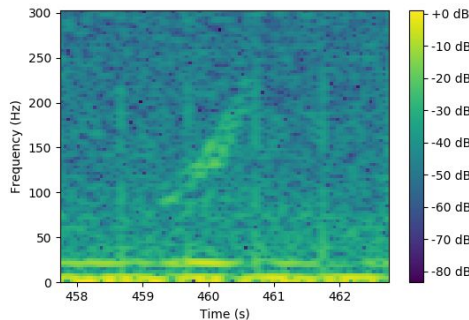
'Safety' of Data Augmentation



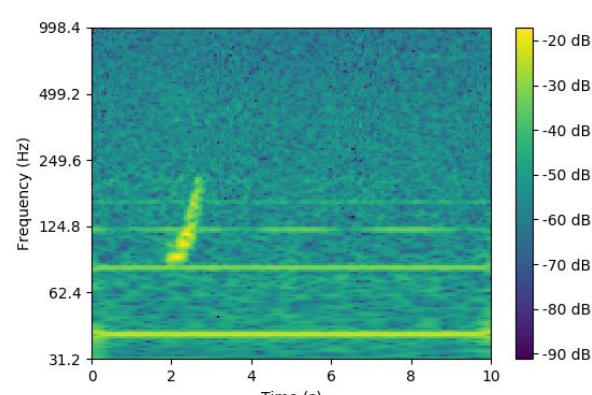
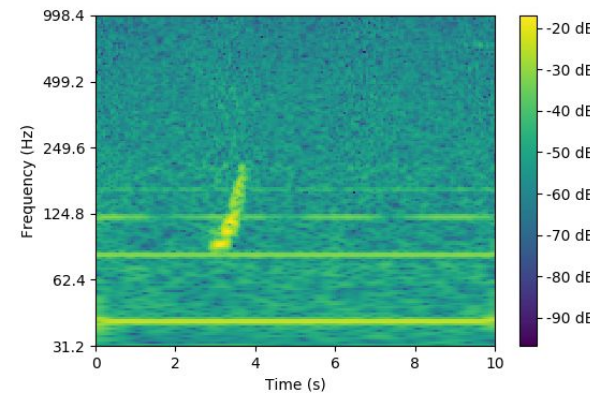
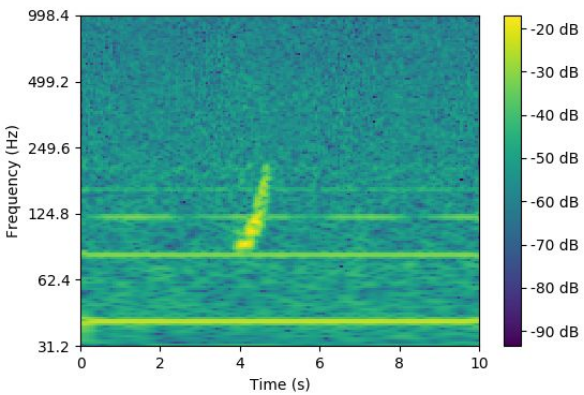
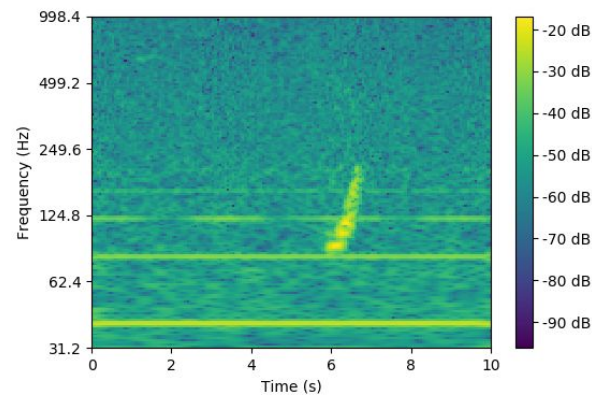
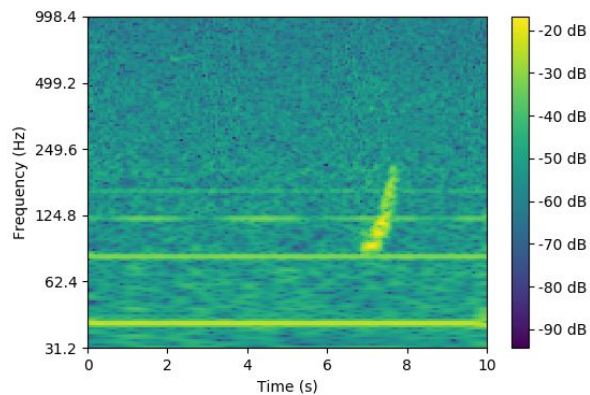
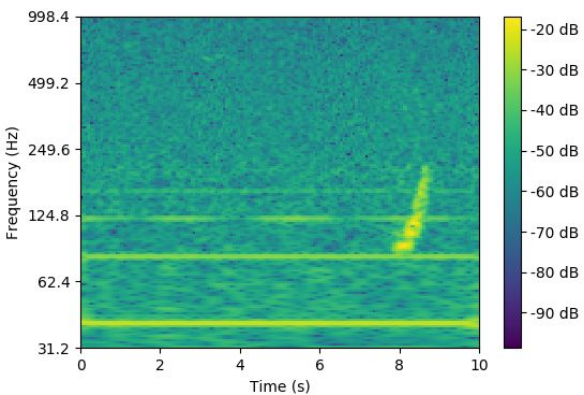
- The type of Augmentation will depend on your data
 - Is the label preserved post-transformation?
 - 'Unsafe' transformations are those that do not preserve the label

Rotations and flips

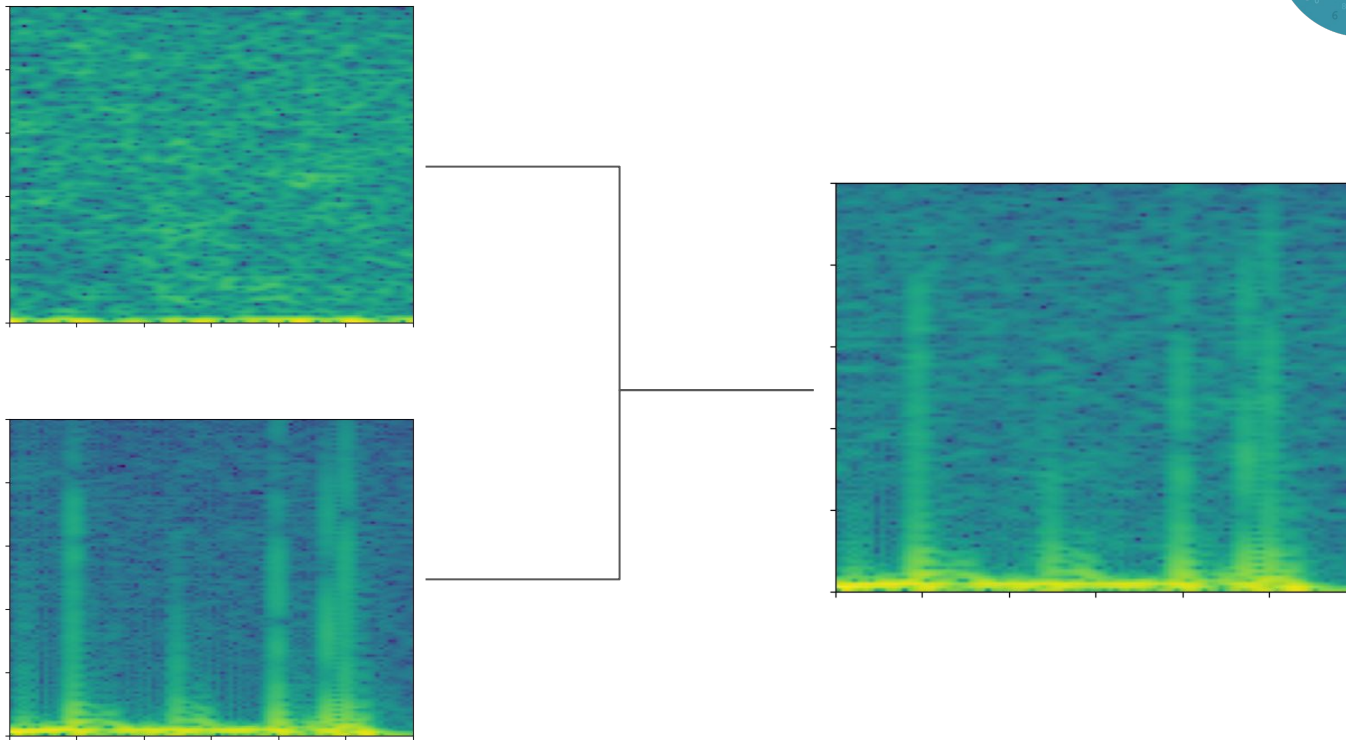
- Generally safe on ImageNet
- Problematic on spectrograms



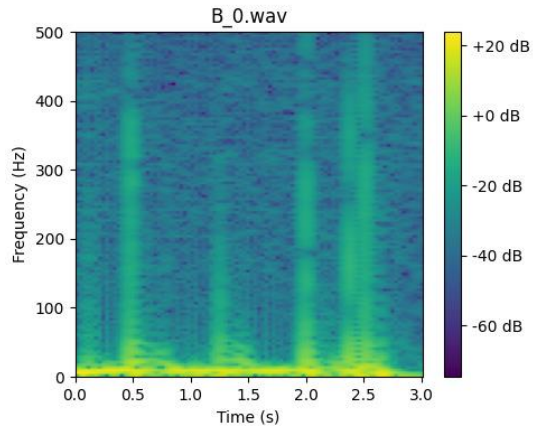
Temporal Shifting



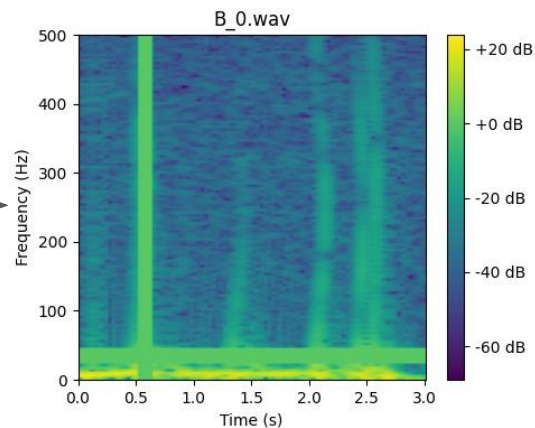
Mixup



SpecAugment



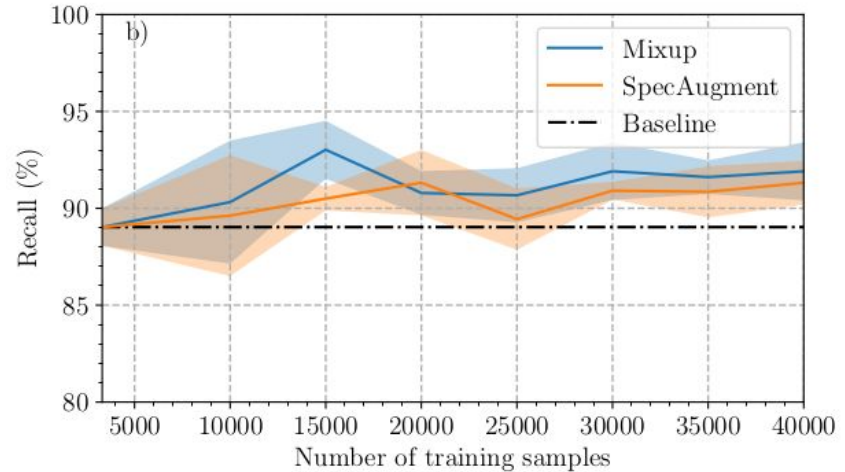
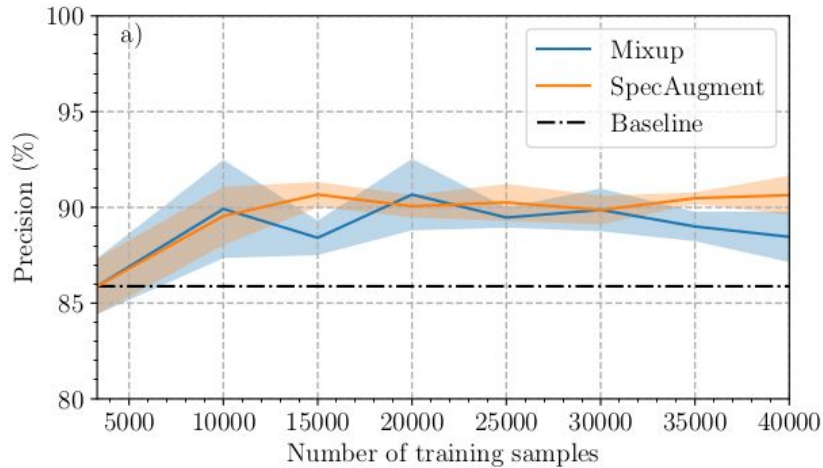
Masking + warping



Benefits of data Augmentation



- Original training dataset: 3309 samples



Kayra: A python package that provides implementation of several data augmentation techniques



- Simpler augmentation methods generate new samples in a very specific manner
 - Limited by the type of transformation
 - Won't generate completely new data
 - Generative methods based on deep learning are capable of modeling raw audio
- Limited by the data already available
 - Data augmentation does not replace real data

Deep Generative Models



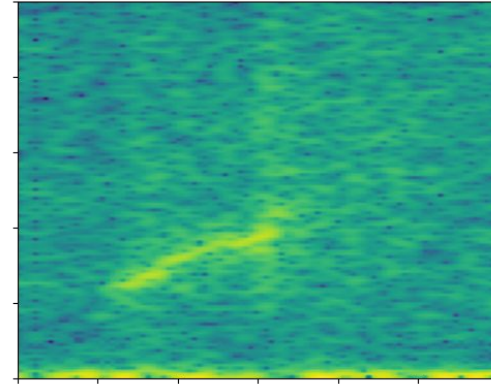
Generative Models - Audio



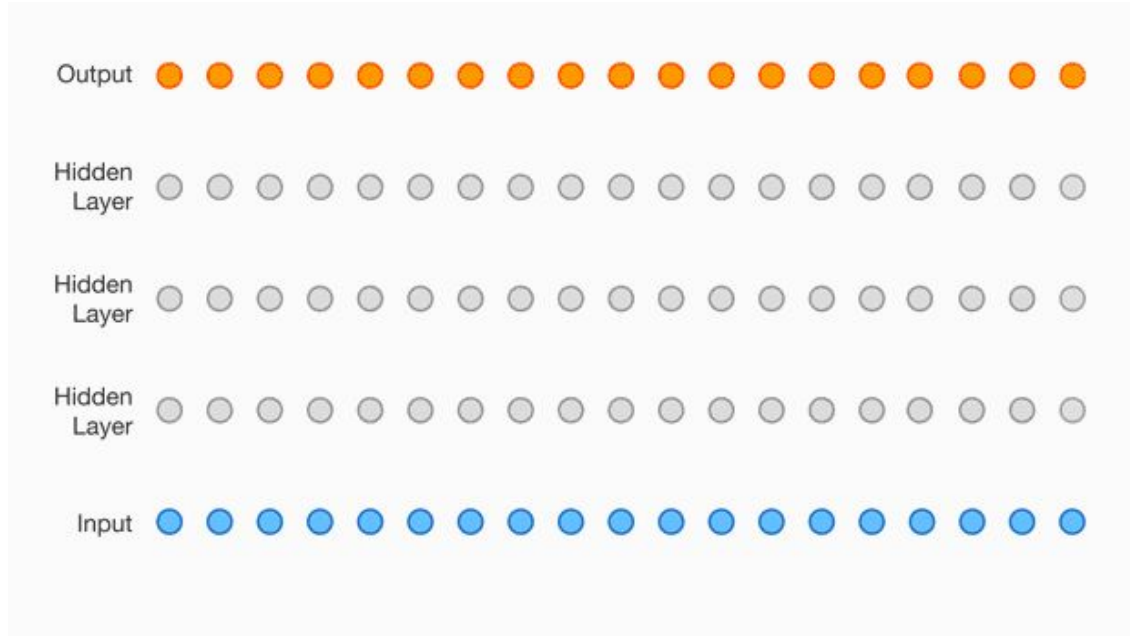
Human Speech



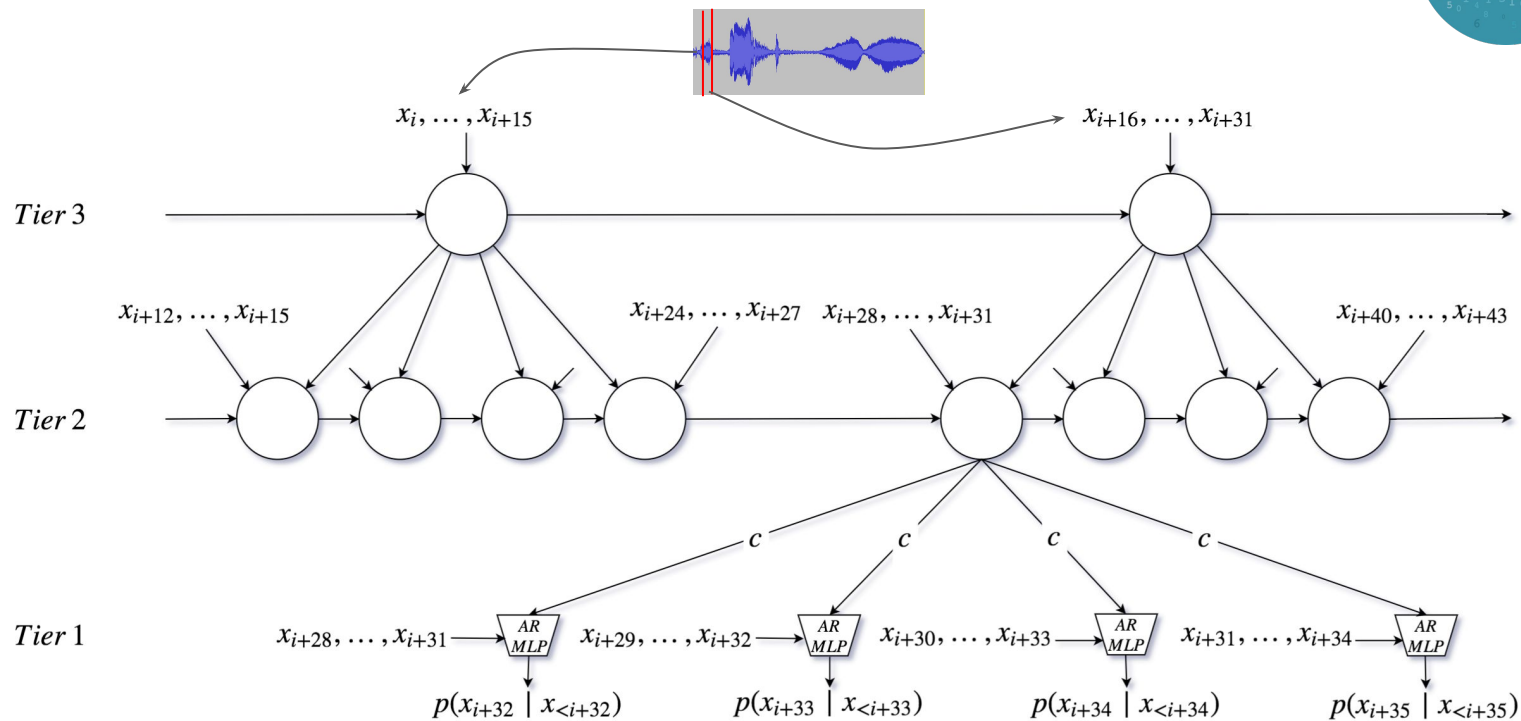
NARW upcall



Generative Models - WaveNets



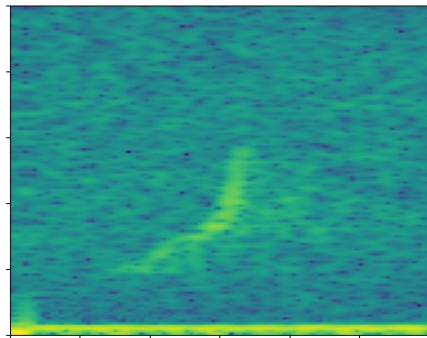
SampleRNN



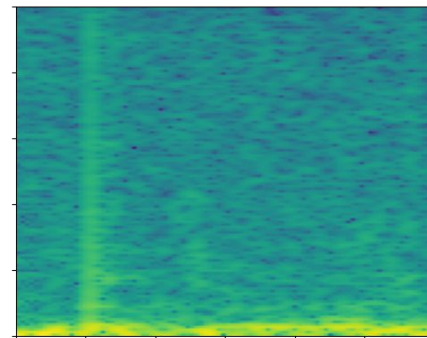
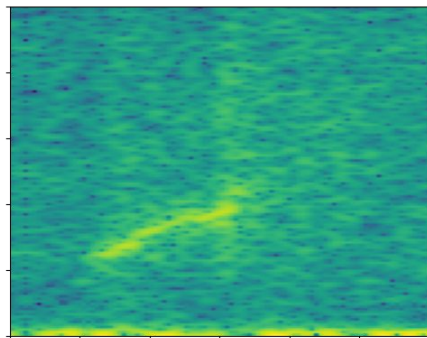
SampleRNN - Quality of generations



Original



Generated



So... what is good enough?



Are our generated samples good enough to be included in our dataset?

- Problem: labeling a new sample as positive when there is no vocalization
- How do we evaluate our generated samples?

There are several possibilities with trade-offs

- Conduct a manual labelling process of the generated samples
 - Would produce an accurate augmented set
 - Expensive
- Use a pre-trained model to classify each generation
 - Inexpensive but less accurate
 - Could inherit model bias



- A more balanced approach can be considered
 - Manually label some generated samples
 - Use these samples in conjunction with a pre-trained model to label the remaining generations
- Clustering methods can be used to group generated samples into classes
 - With a visualization tool, we can then ask for the user to label only the samples that the method is not confident about

