

19/11/24

OPERA

- Begun reading in the github documents for the OPERA dataset which appears to be useful for the machine learning side of things and for datasets related to respiratory data.
 - ↳ can arrange to speak to Evelyn once I have a better idea of what I want to do with this.
- Looked over the first parts of the OPERA paper and stopped at the spectrogram mention. Producing spectrograms will be something I need to do in order to conduct analysis on audio datasets of my choice.

- Followed referenced papers (stored there in notes) until I reached one that spoke of the Lobroa library in Python that is used for audio analysis and producing the spectrograms

Lobroa

- Followed instructions to install Lobroa library along with other modules it requires to perform all of its functionality.
- Have encountered research using machine learning to identify various environmental sounds (~~fauna~~ Animals, Natural landscapes, Human, Interior/Outside, Behavior/Utter) perhaps I could test machine learning methods on spectrograms to identify instruments / genres?
 - ↳ final paper doing just that, will read and see where I could advance. Submitted 1st November 2024

Ideas

- Should focus on conclusions of paper for inspiration for extension ideas on what has already been done
 - ↳ Spark to Pietro about these ideas.

XGBoost

- Mentioned by Jack as more powerful than the Convolutional Neural Networks he studied for his dissertation. Extreme Gradient Boosting has been applied to multiclass classification. 2022 study

Gravitational Angular Field

- Scale the step values $X = \{x_i\}$ to be within the range $[-1, 1]$ then convert these values into $\cos(\theta)$ of polar angle θ .

- plot the values of θ with τ increasing with the time to get a polar plot from your time series

summation

$$\text{GASF} = \cos(\theta_i + \theta_j) \quad (\text{matrix } [j, j] \text{ storage})$$

$$\text{GADF} = \sin(\theta_i - \theta_j) \quad (\text{difference})$$

Time average of GAF (spectrogram) as function of delay time

- Forward our what picture suggests is to compute these spectrograms and then somehow take the difference between them and then average it. Then do this method at multiple delay times (not the spectrogram and take averages!) and plot the time average as a function of delay time.
- Can be classified with image classification
- Maybe also try the GAF with this as would be a time series? then do image classification with the GAFs.

Progress

- implemented Evelyn's framework on the WI and managed to run tasks 10 and 11, both returned sensible results similar to those she was meaning:
 - ↳ importantly different though as these are my own iterations of the process.

OPERA - in detail

Audio recordings can be used to estimate respiratory rate and lung function, detect snoring and open events during sleep, assess the effect of smoking on health and diagnose diseases like flu and asthma.

- Supervised deep acoustic models have been proposed, but their performance relies heavily on the volume and quality of available labels → not easy to have for respiratory data like flu and asthma.
- ↳ maybe just generally:
 - ↳ training models with large amounts of labelled data has high potential to improve performance through transfer learning and supervised fine-tuning

Supervised vs. Unsupervised.

- Supervised learning uses labelled data, meaning each input (feature) comes with a corresponding output (label).
- Unsupervised data learning uses unlabelled data, meaning only input features are provided without corresponding output labels.

Feature: measurable property or characteristic of data that is used as input to a model.

OPERA creates unlabeled ~~or~~ separating dataset and its problem. Then it generates your dataset models, and evaluates them against existing pre-trained acoustic models across various applications.

- ~ 136K samples, 400+ hours orders of magnitude greater than number of respiratory audio samples in datasets used for training existing open acoustic models.
- 3 generalizable acoustic models pre-trained with the unlabeled data (contrastive, generative).
using hard Augmentor / CNN
- 10 labelled datasets to formulate 19 respiratory health task

Results Show the OPERA models are generalizable across multiple downstream tasks, including new datasets and unseen respiratory acoustic modalities

- e.g. clinical imaging, electronic health records and medical time series .

- Contrastive model is better for classification-based downstream tasks.
While the generative model performs better in regression tasks.

Classification - discrete data categorizing
Regression - predicting trends in continuous data like stock markets

Before pre-training all recordings are resampled to 16 kHz and resaved into a mono channel. They are then transformed into spectrograms using 64 Mel filter banks with a 64 ms Hamming window that shifts every 32 ms. → Spectrogram of $1 \times 126 \times 64$ dimensions.

epoch: complete pass through the entire training dataset.

- if a dataset contains 10,000 samples and you train for 5 epochs, the model will see each sample 5 times but not necessarily in order that are other pre-trained acoustic models : OpenSmile, VGGish and Andri's MAE and CLAP.

- How do the tasks use: Split the data up of that could be many distinguishing factors in the audio file?

Plan Sketches (Plan is due 06/12/24!)

(Prior speak about motivations)

- ① You have identified a state of the art game work (**OPERA**) developed by collaborators.

OPERA creates a large dataset of unlabeled respiratory data

- next currently available black black blab.
- also trains

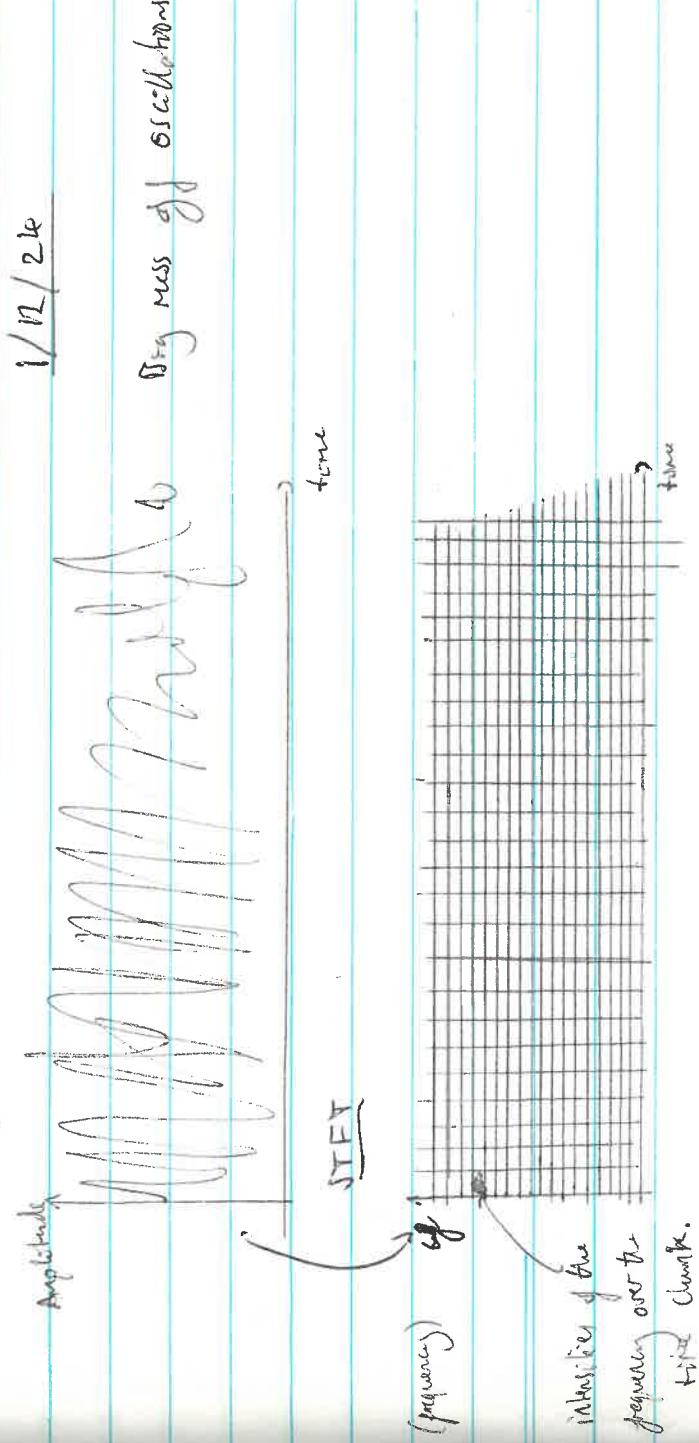
- ② **OPERA** uses features in an optimal way, within a ml framework

- ③ We have discussed some new features, based on average differences of spectrograms. Which we'd like to add to the standard features.

- ④ We'd like to know on samples of real audios and basket of classification / discrimination, what role our new features can have, and see if they can have an impact in the field.

- ⑤ There is a dataset connected to the development of OPERA, want to do a like for like comparison. I have found other data and if we see value on these respiratory applications, could be good to test in wider applications.

1/12/24

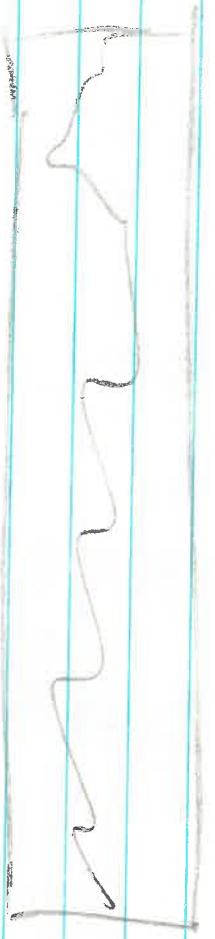


intensity of the frequency over the time chunk.

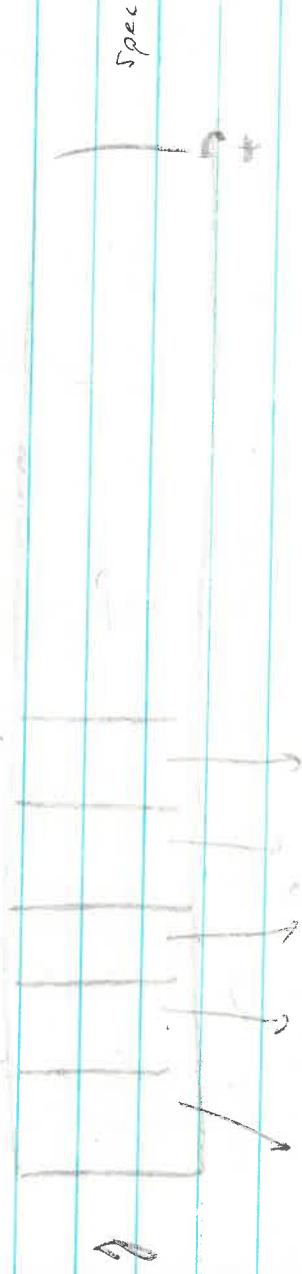
02/12/24

Gantt chart.

1. Creation of new images + replacing mel spectrograms in Kwoklyn process. 21/01 → 04/02 (2 weeks)
2. Pretraining models on the OPERA inventory dataset
28/01 → ~~04/02~~ 01/02 + 1 week
3. Completing tasks with the pretrained models and new images.
01/02 → 28/02 + 2 weeks →
4. Collation of results and comparison with mel spectrogram results.
18/02 → 11/03 + 2 weeks
5. Application of new images to other datasets with different objectives. 11/03 → 29/04
6. Report writing. (Deadline ~~2nd June~~)
29/04 → 02/05
7. Present above preparation (~~7th week of last~~)
01/05 → 01/05
8. Poster creation. (Deadline 12th June)
02/05 → 02/05.



• war



Spec

① ② ③ ④ ⑤ → n

Potential functions to modify:

→ from src / util.h , pre-process - audio - ml - +

→ Conversion with Click of :

might need to do some tuning to get the best results with this new form of spectrogram

- Reading code, clearly pre - process - audio - ml - + gets used in other important functions
 - get - split - Signal - librosa ①
 - get - entire - Signal - librosa ②
 - get - individual - segment - librosa ③

- ① Code an audio file, from silence, and split it into smaller segments.

Adds • filter to pick out frequencies within a certain range (Butterpass) → between 900 and 1800 Hz.

→ Can return either the split audio chunks or the spectrograms.

- ② Processes entire signal and uses license plate detection to remove selected parts. Puts 'Pass' onto the sample to get to output file (8s)
 - c) removing license plates sample too slow.
- ③ This function takes out the license after produces a package of a specified length. This must be very careful for the job at hand, can split & clip into very small segments and then do delay time measurements.
 - ↳ difference between ② start this makes those audio file continuous with whereas ③ splits into the sections where there is noise and makes them into separate playing cards (Dictionaries).

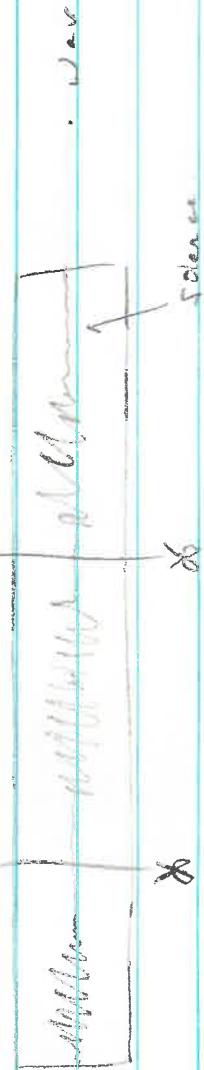
- Questions

- Can it just create a load of these delay images and feed them onto 'extract-opera-feature', and see what happens?

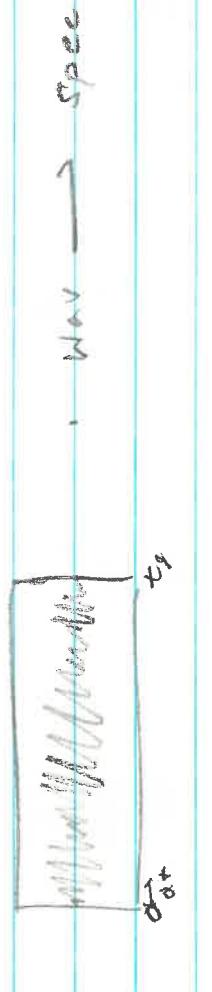
IC / 1 / 2025

Micro Works:

- Take average of the difference of spectrograms as a function of delay time between spectrograms.



Take 1 snippet



Split sample into N slices of the spectrogram as
 $dt = \frac{t_2}{N}$

Take the difference of each snippet with all other
snippets. Keeping total of which the delay time (how many dt's)
the two snippets are away from each other.

Group the differences between the spectrograms into
their respective delay times and average over
the time axis. Now plot the time averaged signal
as a function of delay time.

Variables → 1. length of snippet

→ other ways we discussed mentioned nearest
neighbours.

- have managed to successfully plot Specograms using these partitions.
- Need to see how the removal of silence will affect how I intend to do the delay time plots.
- Preprocess and spent seems to be done at the end of all the manipulations on partitions ① ② and ③. Merge of 5 only modify and open to do what I want them there seems to make sense.

31/1/25

- Have managed to create scatter plots of spectrograms, ready to create delay plots with.

First going to try Method ① outlined 30/1/25
outline of function:

- loop through snippets, take off support, subtract small supports and expand to row / column of matrix delay time $i \rightarrow (N)$ average along time axis then add.
- time i
 i
 $i+1$
 $i+2$
 $i+3$
 $i+4$
 $i+5$
 $i+6$
 $i+7$
 $i+8$
 $i+9$
 $i+10$
 $i+11$
 $i+12$
 $i+13$
 $i+14$
 $i+15$
 $i+16$
 $i+17$
 $i+18$
 $i+19$
 $i+20$
 $i+21$
 $i+22$
 $i+23$
 $i+24$
 $i+25$
 $i+26$
 $i+27$
 $i+28$
 $i+29$
 $i+30$
 $i+31$
 $i+32$
 $i+33$
 $i+34$
 $i+35$
 $i+36$
 $i+37$
 $i+38$
 $i+39$
 $i+40$
 $i+41$
 $i+42$
 $i+43$
 $i+44$
 $i+45$
 $i+46$
 $i+47$
 $i+48$
 $i+49$
 $i+50$
 $i+51$
 $i+52$
 $i+53$
 $i+54$
 $i+55$
 $i+56$
 $i+57$
 $i+58$
 $i+59$
 $i+60$
 $i+61$
 $i+62$
 $i+63$
 $i+64$
 $i+65$
 $i+66$
 $i+67$
 $i+68$
 $i+69$
 $i+70$
 $i+71$
 $i+72$
 $i+73$
 $i+74$
 $i+75$
 $i+76$
 $i+77$
 $i+78$
 $i+79$
 $i+80$
 $i+81$
 $i+82$
 $i+83$
 $i+84$
 $i+85$
 $i+86$
 $i+87$
 $i+88$
 $i+89$
 $i+90$
 $i+91$
 $i+92$
 $i+93$
 $i+94$
 $i+95$
 $i+96$
 $i+97$
 $i+98$
 $i+99$
 $i+100$
 $i+101$
 $i+102$
 $i+103$
 $i+104$
 $i+105$
 $i+106$
 $i+107$
 $i+108$
 $i+109$
 $i+110$
 $i+111$
 $i+112$
 $i+113$
 $i+114$
 $i+115$
 $i+116$
 $i+117$
 $i+118$
 $i+119$
 $i+120$
 $i+121$
 $i+122$
 $i+123$
 $i+124$
 $i+125$
 $i+126$
 $i+127$
 $i+128$
 $i+129$
 $i+130$
 $i+131$
 $i+132$
 $i+133$
 $i+134$
 $i+135$
 $i+136$
 $i+137$
 $i+138$
 $i+139$
 $i+140$
 $i+141$
 $i+142$
 $i+143$
 $i+144$
 $i+145$
 $i+146$
 $i+147$
 $i+148$
 $i+149$
 $i+150$
 $i+151$
 $i+152$
 $i+153$
 $i+154$
 $i+155$
 $i+156$
 $i+157$
 $i+158$
 $i+159$
 $i+160$
 $i+161$
 $i+162$
 $i+163$
 $i+164$
 $i+165$
 $i+166$
 $i+167$
 $i+168$
 $i+169$
 $i+170$
 $i+171$
 $i+172$
 $i+173$
 $i+174$
 $i+175$
 $i+176$
 $i+177$
 $i+178$
 $i+179$
 $i+180$
 $i+181$
 $i+182$
 $i+183$
 $i+184$
 $i+185$
 $i+186$
 $i+187$
 $i+188$
 $i+189$
 $i+190$
 $i+191$
 $i+192$
 $i+193$
 $i+194$
 $i+195$
 $i+196$
 $i+197$
 $i+198$
 $i+199$
 $i+200$
 $i+201$
 $i+202$
 $i+203$
 $i+204$
 $i+205$
 $i+206$
 $i+207$
 $i+208$
 $i+209$
 $i+210$
 $i+211$
 $i+212$
 $i+213$
 $i+214$
 $i+215$
 $i+216$
 $i+217$
 $i+218$
 $i+219$
 $i+220$
 $i+221$
 $i+222$
 $i+223$
 $i+224$
 $i+225$
 $i+226$
 $i+227$
 $i+228$
 $i+229$
 $i+230$
 $i+231$
 $i+232$
 $i+233$
 $i+234$
 $i+235$
 $i+236$
 $i+237$
 $i+238$
 $i+239$
 $i+240$
 $i+241$
 $i+242$
 $i+243$
 $i+244$
 $i+245$
 $i+246$
 $i+247$
 $i+248$
 $i+249$
 $i+250$
 $i+251$
 $i+252$
 $i+253$
 $i+254$
 $i+255$
 $i+256$
 $i+257$
 $i+258$
 $i+259$
 $i+260$
 $i+261$
 $i+262$
 $i+263$
 $i+264$
 $i+265$
 $i+266$
 $i+267$
 $i+268$
 $i+269$
 $i+270$
 $i+271$
 $i+272$
 $i+273$
 $i+274$
 $i+275$
 $i+276$
 $i+277$
 $i+278$
 $i+279$
 $i+280$
 $i+281$
 $i+282$
 $i+283$
 $i+284$
 $i+285$
 $i+286$
 $i+287$
 $i+288$
 $i+289$
 $i+290$
 $i+291$
 $i+292$
 $i+293$
 $i+294$
 $i+295$
 $i+296$
 $i+297$
 $i+298$
 $i+299$
 $i+300$
 $i+301$
 $i+302$
 $i+303$
 $i+304$
 $i+305$
 $i+306$
 $i+307$
 $i+308$
 $i+309$
 $i+310$
 $i+311$
 $i+312$
 $i+313$
 $i+314$
 $i+315$
 $i+316$
 $i+317$
 $i+318$
 $i+319$
 $i+320$
 $i+321$
 $i+322$
 $i+323$
 $i+324$
 $i+325$
 $i+326$
 $i+327$
 $i+328$
 $i+329$
 $i+330$
 $i+331$
 $i+332$
 $i+333$
 $i+334$
 $i+335$
 $i+336$
 $i+337$
 $i+338$
 $i+339$
 $i+340$
 $i+341$
 $i+342$
 $i+343$
 $i+344$
 $i+345$
 $i+346$
 $i+347$
 $i+348$
 $i+349$
 $i+350$
 $i+351$
 $i+352$
 $i+353$
 $i+354$
 $i+355$
 $i+356$
 $i+357$
 $i+358$
 $i+359$
 $i+360$
 $i+361$
 $i+362$
 $i+363$
 $i+364$
 $i+365$
 $i+366$
 $i+367$
 $i+368$
 $i+369$
 $i+370$
 $i+371$
 $i+372$
 $i+373$
 $i+374$
 $i+375$
 $i+376$
 $i+377$
 $i+378$
 $i+379$
 $i+380$
 $i+381$
 $i+382$
 $i+383$
 $i+384$
 $i+385$
 $i+386$
 $i+387$
 $i+388$
 $i+389$
 $i+390$
 $i+391$
 $i+392$
 $i+393$
 $i+394$
 $i+395$
 $i+396$
 $i+397$
 $i+398$
 $i+399$
 $i+400$
 $i+401$
 $i+402$
 $i+403$
 $i+404$
 $i+405$
 $i+406$
 $i+407$
 $i+408$
 $i+409$
 $i+410$
 $i+411$
 $i+412$
 $i+413$
 $i+414$
 $i+415$
 $i+416$
 $i+417$
 $i+418$
 $i+419$
 $i+420$
 $i+421$
 $i+422$
 $i+423$
 $i+424$
 $i+425$
 $i+426$
 $i+427$
 $i+428$
 $i+429$
 $i+430$
 $i+431$
 $i+432$
 $i+433$
 $i+434$
 $i+435$
 $i+436$
 $i+437$
 $i+438$
 $i+439$
 $i+440$
 $i+441$
 $i+442$
 $i+443$
 $i+444$
 $i+445$
 $i+446$
 $i+447$
 $i+448$
 $i+449$
 $i+450$
 $i+451$
 $i+452$
 $i+453$
 $i+454$
 $i+455$
 $i+456$
 $i+457$
 $i+458$
 $i+459$
 $i+460$
 $i+461$
 $i+462$
 $i+463$
 $i+464$
 $i+465$
 $i+466$
 $i+467$
 $i+468$
 $i+469$
 $i+470$
 $i+471$
 $i+472$
 $i+473$
 $i+474$
 $i+475$
 $i+476$
 $i+477$
 $i+478$
 $i+479$
 $i+480$
 $i+481$
 $i+482$
 $i+483$
 $i+484$
 $i+485$
 $i+486$
 $i+487$
 $i+488$
 $i+489$
 $i+490$
 $i+491$
 $i+492$
 $i+493$
 $i+494$
 $i+495$
 $i+496$
 $i+497$
 $i+498$
 $i+499$
 $i+500$
 $i+501$
 $i+502$
 $i+503$
 $i+504$
 $i+505$
 $i+506$
 $i+507$
 $i+508$
 $i+509$
 $i+510$
 $i+511$
 $i+512$
 $i+513$
 $i+514$
 $i+515$
 $i+516$
 $i+517$
 $i+518$
 $i+519$
 $i+520$
 $i+521$
 $i+522$
 $i+523$
 $i+524$
 $i+525$
 $i+526$
 $i+527$
 $i+528$
 $i+529$
 $i+530$
 $i+531$
 $i+532$
 $i+533$
 $i+534$
 $i+535$
 $i+536$
 $i+537$
 $i+538$
 $i+539$
 $i+540$
 $i+541$
 $i+542$
 $i+543$
 $i+544$
 $i+545$
 $i+546$
 $i+547$
 $i+548$
 $i+549$
 $i+550$
 $i+551$
 $i+552$
 $i+553$
 $i+554$
 $i+555$
 $i+556$
 $i+557$
 $i+558$
 $i+559$
 $i+560$
 $i+561$
 $i+562$
 $i+563$
 $i+564$
 $i+565$
 $i+566$
 $i+567$
 $i+568$
 $i+569$
 $i+570$
 $i+571$
 $i+572$
 $i+573$
 $i+574$
 $i+575$
 $i+576$
 $i+577$
 $i+578$
 $i+579$
 $i+580$
 $i+581$
 $i+582$
 $i+583$
 $i+584$
 $i+585$
 $i+586$
 $i+587$
 $i+588$
 $i+589$
 $i+590$
 $i+591$
 $i+592$
 $i+593$
 $i+594$
 $i+595$
 $i+596$
 $i+597$
 $i+598$
 $i+599$
 $i+600$
 $i+601$
 $i+602$
 $i+603$
 $i+604$
 $i+605$
 $i+606$
 $i+607$
 $i+608$
 $i+609$
 $i+610$
 $i+611$
 $i+612$
 $i+613$
 $i+614$
 $i+615$
 $i+616$
 $i+617$
 $i+618$
 $i+619$
 $i+620$
 $i+621$
 $i+622$
 $i+623$
 $i+624$
 $i+625$
 $i+626$
 $i+627$
 $i+628$
 $i+629$
 $i+630$
 $i+631$
 $i+632$
 $i+633$
 $i+634$
 $i+635$
 $i+636$
 $i+637$
 $i+638$
 $i+639$
 $i+640$
 $i+641$
 $i+642$
 $i+643$
 $i+644$
 $i+645$
 $i+646$
 $i+647$
 $i+648$
 $i+649$
 $i+650$
 $i+651$
 $i+652$
 $i+653$
 $i+654$
 $i+655$
 $i+656$
 $i+657$
 $i+658$
 $i+659$
 $i+660$
 $i+661$
 $i+662$
 $i+663$
 $i+664$
 $i+665$
 $i+666$
 $i+667$
 $i+668$
 $i+669$
 $i+670$
 $i+671$
 $i+672$
 $i+673$
 $i+674$
 $i+675$
 $i+676$
 $i+677$
 $i+678$
 $i+679$
 $i+680$
 $i+681$
 $i+682$
 $i+683$
 $i+684$
 $i+685$
 $i+686$
 $i+687$
 $i+688$
 $i+689$
 $i+690$
 $i+691$
 $i+692$
 $i+693$
 $i+694$
 $i+695$
 $i+696$
 $i+697$
 $i+698$
 $i+699$
 $i+700$
 $i+701$
 $i+702$
 $i+703$
 $i+704$
 $i+705$
 $i+706$
 $i+707$
 $i+708$
 $i+709$
 $i+710$
 $i+711$
 $i+712$
 $i+713$
 $i+714$
 $i+715$
 $i+716$
 $i+717$
 $i+718$
 $i+719$
 $i+720$
 $i+721$
 $i+722$
 $i+723$
 $i+724$
 $i+725$
 $i+726$
 $i+727$
 $i+728$
 $i+729$
 $i+730$
 $i+731$
 $i+732$
 $i+733$
 $i+734$
 $i+735$
 $i+736$
 $i+737$
 $i+738$
 $i+739$
 $i+740$
 $i+741$
 $i+742$
 $i+743$
 $i+744$
 $i+745$
 $i+746$
 $i+747$
 $i+748$
 $i+749$
 $i+750$
 $i+751$
 $i+752$
 $i+753$
 $i+754$
 $i+755$
 $i+756$
 $i+757$
 $i+758$
 $i+759$
 $i+760$
 $i+761$
 $i+762$
 $i+763$
 $i+764$
 $i+765$
 $i+766$
 $i+767$
 $i+768$
 $i+769$
 $i+770$
 $i+771$
 $i+772$
 $i+773$
 $i+774$
 $i+775$
 $i+776$
 $i+777$
 $i+778$
 $i+779$
 $i+780$
 $i+781$
 $i+782$
 $i+783$
 $i+784$
 $i+785$
 $i+786$
 $i+787$
 $i+788$
 $i+789$
 $i+790$
 $i+791$
 $i+792$
 $i+793$
 $i+794$
 $i+795$
 $i+796$
 $i+797$
 $i+798$
 $i+799$
 $i+800$
 $i+801$
 $i+802$
 $i+803$
 $i+804$
 $i+805$
 $i+806$
 $i+807$
 $i+808$
 $i+809$
 $i+810$
 $i+811$
 $i+812$
 $i+813$
 $i+814$
 $i+815$
 $i+816$
 $i+817$
 $i+818$
 $i+819$
 $i+820$
 $i+821$
 $i+822$
 $i+823$
 $i+824$
 $i+825$
 $i+826$
 $i+827$
 $i+828$
 $i+829$
 $i+830$
 $i+831$
 $i+832$
 $i+833$
 $i+834$
 $i+835$
 $i+836$
 $i+837$
 $i+838$
 $i+839$
 $i+840$
 $i+841$
 $i+842$
 $i+843$
 $i+844$
 $i+845$
 $i+846$
 $i+847$
 $i+848$
 $i+849$
 $i+850$
 $i+851$
 $i+852$
 $i+853$
 $i+854$
 $i+855$
 $i+856$
 $i+857$
 $i+858$
 $i+859$
 $i+860$
 $i+861$
 $i+862$
 $i+863$
 $i+864$
 $i+865$
 $i+866$
 $i+867$
 $i+868$
 $i+869$
 $i+870$
 $i+871$
 $i+872$
 $i+873$
 $i+874$
 $i+875$
 $i+876$
 $i+877$
 $i+878$
 $i+879$
 $i+880$
 $i+881$
 $i+882$
 $i+883$
 $i+884$
 $i+885$
 $i+886$
 $i+887$
 $i+888$
 $i+889$
 $i+890$
 $i+891$
 $i+892$
 $i+893$
 $i+894$
 $i+895$
 $i+896$
 $i+897$
 $i+898$
 $i+899$
 $i+900$
 $i+901$
 $i+902$
 $i+903$
 $i+904$
 $i+905$
 $i+906$
 $i+907$
 $i+908$
 $i+909$
 $i+910$
 $i+911$
 $i+912$
 $i+913$
 $i+914$
 $i+915$
 $i+916$
 $i+917$
 $i+918$
 $i+919$
 $i+920$
 $i+921$
 $i+922$
 $i+923$
 $i+924$
 $i+925$
 $i+926$
 $i+927$
 $i+928$
 $i+929$
 $i+930$
 $i+931$
 $i+932$
 $i+933$
 $i+934$
 $i+935$
 $i+936$
 $i+937$
 $i+938$
 $i+939$
 $i+940$
 $i+941$
 $i+942$
 $i+943$
 $i+944$
 $i+945$
 $i+946$
 $i+947$
 $i+948$
 $i+949$
 $i+950$
 $i+951$
 $i+952$
 $i+953$
 $i+954$
 $i+955$
 $i+956$
 $i+957$
 $i+958$
 $i+959$
 $i+960$
 $i+961$
 $i+962$
 $i+963$
 $i+964$
 $i+965$
 $i+966$
 $i+967$
 $i+968$
 $i+969$
 $i+970$
 $i+971$
 $i+972$
 $i+973$
 i

$$\sum_{n=1}^N n = \frac{N(N+1)}{2}$$

$$= \frac{N(N+1)}{2}$$

- Average each entry along a row along the time axis and then append the frequency vectors together → Spectrogram as a function of delay time.
- doesn't make sense.

— Call with Pichro:

- ① only take delay time averages where the number of snapshots averaged is $\geq \frac{1}{3} N$
 - ↳ include the diagonal of zeros.
 - ↳ characteristic on the video analogue of this technique.
- ② After taking this, there is the idea of windows
 - two different time scales, episode length and snapshot length.
 - For long audio clips makes sense.
 - do normal delay averages but separately for each episode.
 - length of episode can vary - something to explore.
- ③ Another angle is the windowing function length, parameter involved in creating the spectrogram aspects how pronounced certain frequencies are by taking the window length as a third axis, has the effect of similar to holding at wavelength (left?) transform of audio instead of FT spectrogram.

6/2/25

Delay plots

- Had query over whether more sensible to leave conversion of digitization to device until before or after the delay line averaging process.
- Pietro believes that device conversion before is fine, the averages look better this way, waits at some point putting some solid evidence behind this.

- Now have capability to create delay plots of any window length for any audio file \rightarrow all functions written on spec - delay file.

DFT and Multi-DFT process

DFT algorithm takes the algebraic difference of several copies of frames separated by a lag time τ .

The differences are then Fourier transferred in space and the results averaged

- Fourier transform after taking algebraic difference? It can not correctly do this.

- The idea of episodes in the video content is wished to be described as multi - multi - DFT varying box size and also provide length (subset of time).
- Have downloaded Lorca's Matlab code for DFTM with Mobile Cilia, will need significant modification if I wish to use.

Main things to check coming up

- ① Fourier transform before or after time delay
manipulation (or just to lists).
- ② Need to see how to get feature from raw new data.

7/2/25

New problem: if an audio clip going in with a sample rate of 48000 gets put through pre process method with sample rate 16000, that appears to be reducing the time of the sample (may be just my visualisation).

→ Now big misunderstanding over sample rate and how spectrogram manipulation is reducing the length of the audio clips - need to investigate.

Have created a process to unpack a directory of wav files and transform them into delay images in a new specified directory alongside another folder holding images of the spectrograms.

10/2/25

- Making Petro Wednesday, emailed Evelyn want to get training
- Have implemented a delay msg output to get entire signal, whose not entirely confident it will help much.

11/2/25

→ Delay episode functions:

- 1st divide audio clip into time cubs.
- need to figure how to extract time length.

- 2nd perform samples functions on each episode

3rd just the same as delay image functions.

(Probably best to make two separate functions.)

12/2/25

Meeting

- convert the delay plots to take the modulus of the snippet difference rather than the actual values.
- Get stuck onto writing an orange classification process using google thing (PyTorch or TensorFlow). Should be relatively simple basically step by step on line.
- Definitely worth investigating what is happening to the ignore dimension, generate random .WAV file and see what happens.

14/2/28

Have located a step by step walkthrough on GitHub for an example of audio classification using Tensorflow CNN model.

Have been able to directly substitute my delay time function to generate the spectrograms and am now following through to replace the CNN machine learning.

16/2/25

Progress

Have replicated the use of a CNN to distinguish the sounds of chainsaws, background noise, engines and stars.

To solve this using the delay representation of the image with a size of 1 pixel.
(arie: next vary the size of window and see how this impacts the accuracy of the CNN). The results on the whole were that using the delay image had very little impact on the performance of the CNN, but importantly it didn't seem to perform worse. These conclusions, I have made by comparing the accuracy curves as a function of epoch for Spectrogram representation and its delay representation.

Something to note is that the spectrograms for background noise vs those in which there are chainsaws / engines / storm surges) are a uniform dist. (ish) with an overlaying patch of intensity in same to one face where the intensity starts to higher frequency -

changes



As the sample size increases
and the impact of big out
liers being more apparent in
a long & flat distribution in
this case.

Usually will be
average out
→ similarly uniform across
the time clip at those frequencies

Perhaps delay smages are less useful in this context as there is so much noise so some presence of a disturbing noise such as a clapping out just appears as a disturbance to the before low frequencies and here a band or three high frequencies like I've drawn rather than exposing 'hidden' others in an isolated audio recording → episodes might be removed by end value to this analysis.

There was one instance on which there was a big failure in identifying stories from stories which was a bit worse than working with spectrograms. This was just one time though and plots of spectrograms would be needed for a conclusion on whether my method struggles with this distinction.

To do

- try analogous again with larger windows longer and see how this affects performance.
- Check some concrete examples displaying what it actually does to the images and find out about what is happening to the interpreted time when to run and/or does change Evelyn's Specalog rare generating functions.
- Follow up email to Evelyn showing her the accuracy curves and ask her to go to do this with new models (good short and long) allow me to figure out how the code works a bit better.) ✓

17/12/28

Evelyn's response:

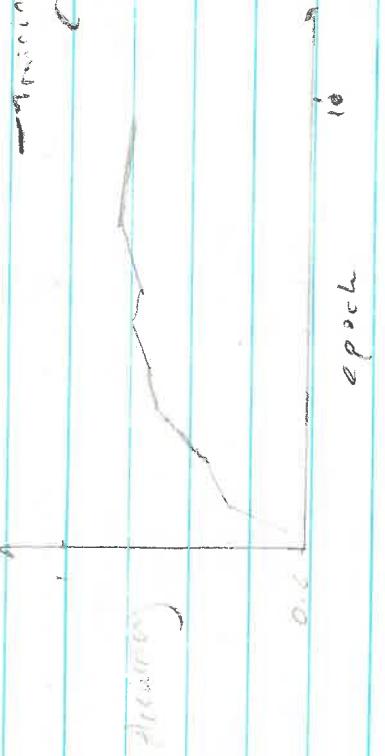
OPERA is a pretraining framework. It's for obtaining an end-to-end training done feasible for testing the effectiveness of the new images.

- Since advice as before to abolish the pretraining scripts.
also scripts / eval-all.ssh for downstream evaluation.

Testing larger window sizes

- the effects of increasing the window size are surprising by small.
- going from delay line steps of 1 (\times angular size) to 100 (\times size) still results in fairly good success rate

training Accuracy (1)
roughly the same - 100)



- tested window sizes from 1, 5, 10, 20, 50, 100
- Big window size ~~sample~~¹⁰⁰, the window size was ~~of~~ resulted in only a delay time every being suitable including the initial 0 steps so the images were $\frac{1}{5}$ 0 intensity.
- interesting that despite this the drawing curves didn't become significantly worse.
- The one more notable difference was that the difficulty or distinguishing characters / anyone gets with window size and for the largest window size, the sharpness was badly affected for the background.

Generating wav files

- ① Since sample rate 22050 Hz, 105 samples
440 Hz
- ② Top hat, sample rate 44100 Hz, 105 on between 3 and 7 seconds (44100 Hz is CP quality standard good for sharp edges (top hat))
 - Used $\frac{5}{2}$ for sine to save memory
- ③ White noise, SR: 44100 Hz, 105 covering half an half of (could try $\frac{1}{4}$ / period).

Plus the time and resulting spectrogram shape is handled:

Time - axis size of the spectrogram is given by
 \downarrow
(time \times sample rate)

$$\text{frames} = \frac{\text{samples}}{\text{hop length}} + 1.$$

hop length could be something to vary.

- ⇒ Solved time issue
- ↳ Problem was on the sample rate in the spectrogram generation, this must be the sample rate of the original audio in order for the image to be labelled with a correct time axis.

In order to modify delay - need to take time or frames as input:

$$\text{frames} = \frac{\text{sample rate} \times \text{time}}{\text{hop-length}} + 1$$

↳ necessary time = $\frac{\text{hop-length}(\text{frames} - 1)}{\text{sample rate}}$

so far, you take time and then convert to the index using this equation.

From the +1 when specifying a number of indices to the place

$$\Delta t = \text{int}(\text{np.round}\left(\frac{55 \times \text{time_slice}}{\text{hop}}\right))$$

(this could cause an issue if the rounds to zero)
→ some delay examples

some were looks record (on personant) likely this
was sum build up of rounding errors?

Or issue with the throwing at the end
of dt or not an exact multiple of the
length of the array.
→ think must be this. if method shows prints
will need to address this.

→ note sample rate 16 000 does not seem to be consistent
in anything? → maybe will only increase for not decrease?
top hat is not a sensible audio signal
for this!

It is impacting the resolution of the data / there is small
scale variation that appears when recorded. So intuition says
keep at 16 000 but this is not trivial.
The other sample rates should be tracked as they help
the time axis sensible.

18/2/25

Progress :

1. Example plots

- have implemented example mainly of modulated white noise such that the frequency bins are all populated. Those are modulated with periodic step functions or sinusoidal envelopes.
 - ↳ could try to think of more examples
 - The results of the delay image transformation do preserves the frequency of the patterns in the spectrogram representation but appears to smoothen step junctions and general fluctuation.
Some appear to go to one
 - One major thought was that the act of due to keeping the averaging sample size a lot enough diminishes $\sim 35\%$ of the original bins over. When the patterns are so obvious such as constant frequency or off, it is clear to predict / you could extrapolate the trend further out to delay time = time.
 - nice to see if there was an option to match axes.
- #### 2. Code/Ok download
- in progress, each file \rightarrow 2GB and there are ~ 2.0 as it will take some time.

19/02/

Meeting

- ① Conference with current sponsor until 12/03/25
(and/or event day)
- ② No point in increasing window size (not in budget)
- ③ try to see if can activate success to the cutoff sample size (right now $\frac{1}{3}$)
- ④ After 12/03/25 pivot point of whether to look into health and/or delay time stuff
or move to the separable approach to:
see if can find some golden time size
of episode to make the algorithm perform
best!

PIPERA

I believe it have just one task to work with the data
images neither other normal and far than real
model: Open source VGGish Andro.ML CLAP et cetera

AUC $0.494 \pm 0.590 \pm 0.510 \pm 0.550 \pm 0.725 = 0.60$
 $0.054 \quad 0.053 \quad 0.021 \quad 0.049 \quad 0.787 \pm 0.00$

Also/
against a lot better

20/02/25

- tried Task 10 with the new images and have confirmed that they are giving different results to the OPERA paper.
- Task 10 performed worse with the new images it seems data is logged on the OPERA file
- still only able to use OPERATOR and C&T for the delay time images as CT uses different functions that would require some talk with Silvyn to implement my current method into.
- Figured out how to write and execute .sh scripts which is useful to execute the other tasks.
- Struggling with installing the Code UK dataset as initially got the wrong one and the size of the dataset is about 120 GB so trying to find a place to put it.
- Should look at the other datasets and tasks.

Progress

- lots of difficulty down loading and sorting datasets.
- we have cracked it now.
- still only able to evaluate OPERA CT/CE
- we now got logstore account and have moved OPERA directory here

Task	OPERA-CT	OPERA-CE
# 1	0.551 ± 0.002	0.513 ± 0.010
# 2	0.644 ± 0.005	0.636 ± 0.001
# 3		
# 4		
# 5	0.571 ± 0.005	0.599 ± 0.006
# 6	0.765 ± 0.001	0.694 ± 0.001
# 7	0.789 ± 0.021	0.758 ± 0.006
# 8	0.715 ± 0.021	
# 9	similar: 0.877 ± 0.002	
# 10	0.725 ± 0.005	0.707 ± 0.004
# 10	0.685 ± 0.009	0.688 ± 0.032

26/02/25

Definition of AUROC :

ROC curve : A plot of True Positive Rate vs False Positive Rate
at different classification thresholds

AUC (Area Under Curve) : The total area under the ROC curve,
which represents the model's ability to
rank positive instances higher than
negative ones.

AUROC = 0.7 → Fair to good classification
AUROC = 0.5 → No discrimination, random chance.

Take only a few next weekend. ✓

Focus on buying delay time down and see if the
performance (time) improves. → ↓

Make sure to do the thing where its gather that
gets thrown away increases and compare with
Spectrogram → seems decently promising to be fair.

The plan

Make it so that I can do classification with a
dataset → look in an specific task may be
count UK.

→ find one and understand it of course some
and subsequently run time? →

→ have this by next wednesday.

27/02/25

→ Many difficulties OPERA CE Tasks import see
→ 2 but the framework is not creating files
their give the input sec = 2 part.

→ Need to get an little creation this way of
testing the said / accuracy with image size .

28/02/25

Coming up with the pretraining stuff for the tree
being focusing on testing for sizeup .

My cleavage code should be easily adapted as
in the OPERA analysis the audio files are all
precessed and saved so I have arrays with
corresponding label arrays and train test
sets .
↓

lot of this in linear - eval . py and {labelfit} - processing

process

- ① the pre processing of the spectrograms
- ② the assigning of labels
- ③ full straight into cleavage code .

Now successfully preprocessed set of covid file
data (Task #1) and fed through training
with labels .

The labels were taken from covid - covid / and
potentially not changing as expected as the

results are giving just a fit curve,
① Many need to read CSV file myself and create
a label list.

② Also might have to re-jig all the checkboxes
etc because that was designed for few
categories and potentially letter options for
longer boxes

③ Still going array → pg → array which
could likely kill the speed up benefits of
delay image and so will need to do
further in depth with this.

01/03/23

- checked number boxes and - was on the
dataset currently using there are
884 positive cases and 1616 negative
cases → could be issue?

→ about how write spectrograms at delay!

02

Tried using spectrograms and running onto same
file, can't seem to break past the ~0.64
mark which indicates the model is just
picking one category every time to maximize
accuracy

→ try with training effects?

Meeting

- Joined by Evelyn → hadn't trained a model so up to now what had been doing was fairly hopeless.
- The google net training approach hadn't really made any ground, still unable to see growth on the validation data accuracy indicating overfitting
- Need to go deep onto multiple-protein and prepare-data - st. ↴
- Have been going through each dataset to pre-generate delivery times → quite likely to will have to modify this and do again.
- ④ Problem with Covid UK, only have downsample file names and supposed to have normal.
- ④ Problem (potentially) with big long → big blocked out (small value) region at top of each check delivery string, maybe those frequencies aren't present in the big long data? ↴
- might also need to reduce the target delay time. (relative max delay time to free-thrown!). ↴
- ④ → Meantained not to have gap between windows be larger than 0.4s

Have done initial configuration and logging so far.

Need to check & delete files get individual cycles
you are producing so should see why that
is necessary.

① Problem Covid-19 Sounds is a config file so many
files to download → just plug away by 1.

② Echoes also does entire spectrograms → need to add
make files line:

09/03/25

- Nothing too major since last update except
it knew how all the data sets and have
converted all data delay spectrograms.
- Now focusing on multiple patterns.sh and
having issues, currently running
"EOF ISRTR". No data left on file.

↳ located the problem, covid-19: UHD - 26811 - t4RWR ..
↳ have removed this from necessary arrays.

- 10/3/25
- Number of epochs was not meant to be 512,
actually 100 / 5 exceeded this so it
will not take 4 days!

- There are copper elements to the
pre processing that do not shift delay
image!
i.e. if the audio clip is too long
it will take a random clip of the right

length so as to bypass this

- all images go in with shape (256, 64)
with delay or it could potentially reduce this potentially / at least remove the elements of random cropping that could improve the process.

(c) Note: already mention crop - first that could easily sub for random crop.

- Issues with OPERA CSE and OPERA CT doing random cropping as part of the method and encountering shapes lower than the max time length probably due to delay effect

OR if its a random crop and starting at the end, it might not have enough remaining time to get the correct length.

- Yes, whenever the length of the delay sample is not long enough, the random crop fails. Currently there are low postpone of these features for now so can just discard batch L could also reduce the max length by the greater getting thrown away as an idea

Probing 'bad files' (ones of less than optimal size):

Vocabs cycle : 1 in 5024) incorrect.
(Covid UK education : 1 in 1800.)

Covid 19 sounds breath : 1470 out of 33766 ~ 4.3%
→ fixed
~~1000~~ → 4652 out of 33766 ~ 1.3%.

Covid 19 sounds cough : 32 out of 49627 ~ 0.07%
→ fixed.

ICU bii - zero bad file

ICU cycle - 882 out of 8024 ~ 0.0218.7%

Covid UK estimation : 295 out of 1800 ~ 19.6%
→ fixed

Covid UK song bii : 129 out of 1499 ~ 8.61%
→ fixed.

ly lyrics : 0 out of 10554 ~ 0%

11/09/25

Progress

- Dog revelations : Tricky has sent the correct training files for covid UK.
- Need to retain the input seq and padding generalizing on the get_entire_signal lib since because otherwise drops get through that are too short and random crops on OPERA CT and OPERA CE jail.

Audio Day

- Spoke to someone saying that using mel spectrograms was detrimental in machine learning and that is a pred position from a human perspective that biases certain frequencies and the idea is for the cough and breath data we are trying to use ML to capture the hard to spot characteristics of respiratory data that can distinguish it.
 - ↳ try and see how linear spectrograms could fit onto everything.
 - Pietro saying we don't want any copper we are just going to give the delay time that we go to
 - I think I should direct my attention on the generative OPERA GT model first and then move on to see where it gets.
- Notes and links
- by long stethoscope ↓ data doesn't contain frequencies passed ~ 1000 Hz.
 - ictus too, seems all stethoscope data loss the high frequency data in the averaging process

14/03/25

Progress

- Finally OPERA CT has trained using delay images.
- In OPERA extract-audio feature announced the use of get split signal to get entire signal then cropped the part length to the desired length.
 - for OPERA CT.
- Article:
 - for the paper after corrections at audio day gave motivations / introduction:
 - the use of mul-spectograms in common place on the field of respiratory audio data and by encouraging AT we question whether their use is biological with AT and respiratory audio where AT does not hold the same basis as only detecting among human species and in order to scroll for hidden time frames we should it be using frequency time representation that are suited to human ears and instead be more general
- Stories with now extracting feature as requires 2800 audio clips of the correct length (8 seconds)
 - Now padding (with zeros & believe) to 1.5x the input second length, then the delay image will cut short back to the correct length (or greater)
 - will result with all images having out at correct length of delay seconds.

- the next big step will be to modify the size of the image used training set and of OPERA - GT down from 256 to small and smaller and see how performance holds up.

16/03/2018

The OPERA GT feature extraction and benchmarking ran but failed for most of the tests.

For now I believe debugging will be too time consuming for little reward as all the processes are designed to be orthogonal to spectrogram and are not general to images. It also takes a long time to run pre training and other tests which makes mistakes costly so I need to find another way.

- Have found an open baseline method for analysis of the oblique dataset (long sounds) that seems easy to play around with.
- Currently the accuracy isn't great nor does it seem to be growing so much for real - spectrograms but if downs sample sampler and gather less run and also looks at acoustic features so I think it has potential.

↳ Would be helpful to implement transfer learning by the current CNN and FNN methods

→ Is debugging worth to

17/03/25

Papers

One very similar to my work:

A Novel Mel-spectrogram Support Representation Learning Framework For Severity Detection of Chronic Obstructive Pulmonary Diseases

- It has been noticed that as spectrograms contain diverse information of audio signal it is difficult to extract promising feature from image based CNN models [45].
 - Yannick has been trained on the mel-spectograms extracted from the Andro signals of Andro set, which is the largest dataset for android deep learning.
- Idea
- In the acoustic feature vector creation aspect of the code, there is the option to add deltas.
 - deltas are how all of the feature values have changed from the window before to the next.
 - In the same style could add another 'deltas' set where instead of the difference between the features with the window before, do the deltas of fixed delay bins and do my averaging technique!

18/03/2

Feature extraction

Currently, there are 34 features evaluated at 155 windows

The windows are two the top length & believe me now there is $\frac{1}{2}$ overlap between one window to the next.

The implementation of deltas : Tree vs essentially to shortest delay time implementation of my idea.
→ All is have to do is add delta but one, delta but two, delta but three etc

18/03/20

→ avoided delta & summary

→ Return to OPERA

over GT

AUC

Task 1	0.592 \pm 0.006
Task 2	0.594 \pm 0.003
Task 3	0.673 \pm 0.001
Task 4	0.623 \pm 0.001

AUC

Task 1	0.549 \pm 0.006	- 0.056
Task 2	0.670 \pm 0.001	- 0.007
Task 3	0.594 \pm 0.003	- 0.019
Task 4	0.623 \pm 0.001	\pm 0
Task 5	0.521 \pm 0.002	+ 0.019
Task 6	0.761 \pm 0.000	+ 0.026
Task 7	0.646 \pm 0.007	- 0.095
Task 8	0.684 \pm 0.020	+ 0.004
Task 9	0.867 \pm 0.000	+ 0.042

OWN GT - OPERA GT

AUC

Task 10 $0.693 \pm 0.012 = 0.010$

Task 11 $0.716 \pm 0.024 \rightarrow$ Way better again.
+ 0.110

→ Note lots of signatures on Corwn delay images.
→ Could be cut down and retain only.

20/08/25

Task 12

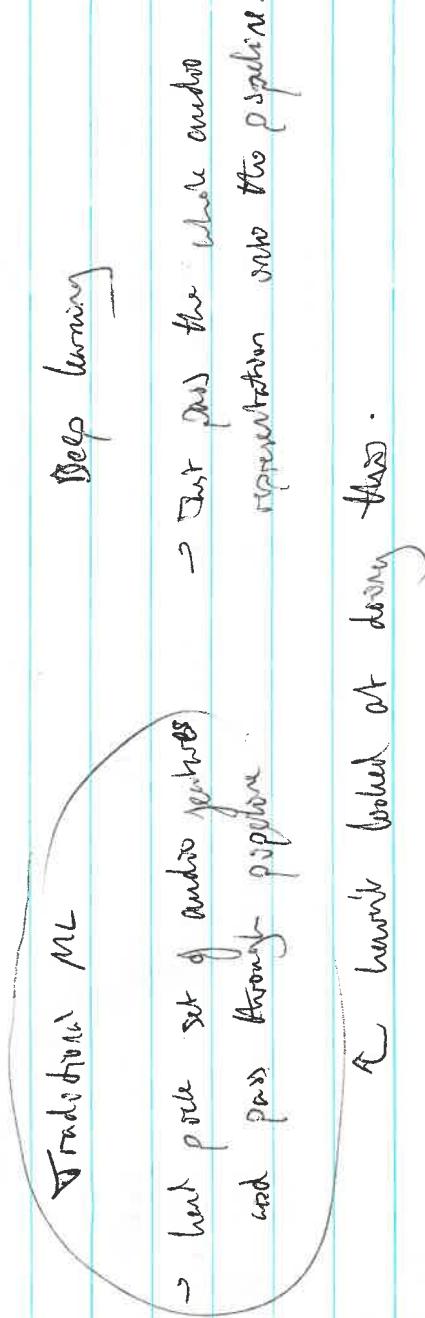
Task 13	$0.906 \pm 0.575 + 0.014$
Task 14	$0.764 \pm 0.595 - 0.061$
Task 15	$0.124 \pm 0.137 - 0.004$
Task 16	$0.859 \pm 0.570 - 0.019$
Task 17	$0.851 \pm 0.581 + 0.077$
Task 18	$0.133 \pm 0.145 + 0.003$
Task 19	

26/08/25

Set training running over weekend with max-len fixed for all datasets at 64.
→ made various changes to code accompanying all input - as to be 2 seconds of delay time

→ Query on the image size (256, 64) being independent of the max-len chosen.
→ will need to check this, have enabled Stylize.

Youtube video on audio features



→ Just pass the whole audio representation into the pipeline.

↑ hasn't looked at doing this.

Valentí Vélez dor signal channel

- something to consider, The delay images aren't normal.
Where 0 and 1.

29/5/25

- Pichó happy with 2-3 delay times + keeping
OPERA random for my images just crop!
- followed up with ~~to~~ today and Tong. To
make sure they are happy with their
choice of action.

- In the mean time going to create a
function to engineer some features.

Currently:

- function stores previous feature vector values to
compare with current feature vector and then
loops
- need to create a 2x feature - down
array that stores the previous 2
features and then compares the current
feature vector with current and then
re assign the later feature vector
back one step.

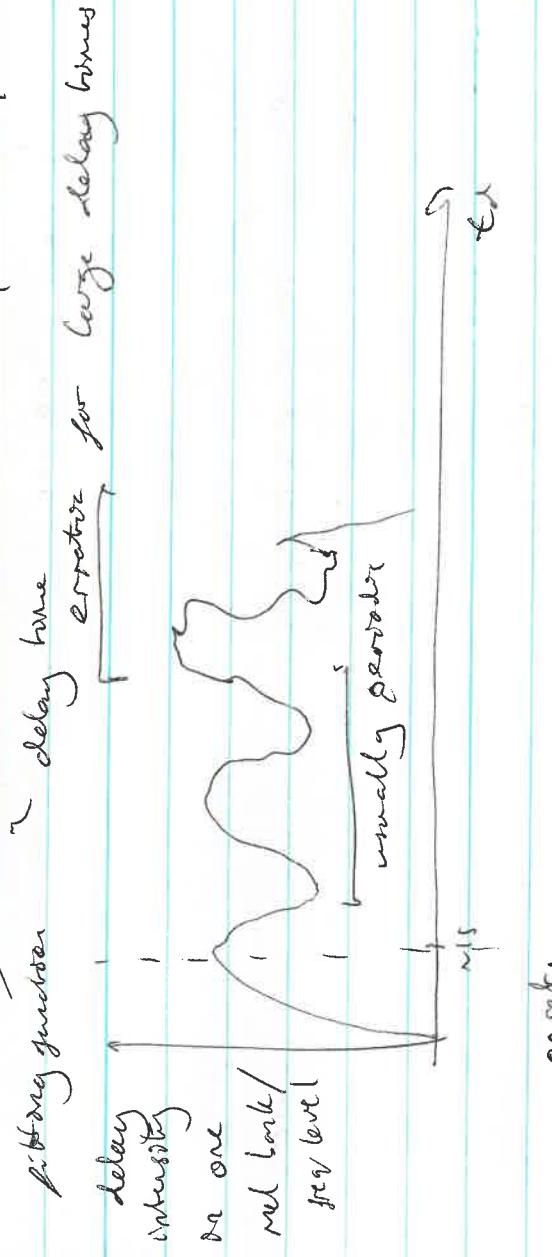
- ↳ doesn't make sense to change the generate
extraction function - make a new one
to manipulate the network output of feature
extraction.

8 / 4 / 25

Week 5 update

- I created the MFCe delayed A style feature and saw minimal increase in accuracy or trend with delay time.
 - Have been reading around papers and found that the Cividis Sounds original paper has an accessible git hub so is ideal to use for my feature engineering
 - Spoke to Pichars about feature engineering and we would like to see an exponential fitted to the onset of the delay signal up to the first peak. From this two parameters can be extracted

$$F(t_\alpha) = \alpha_0 \left(1 - \exp \left(-\frac{t_\alpha}{\alpha_1} \right) \right)$$



- The OPERA models are still on progress but training has sped up so hope to have all by end of the week.

9/4/25

Call with Picture

- To deal with difficult time scale values can just pass a filter / only keep a certain range of frequencies

$$\hookrightarrow SR = 22050 \downarrow \div 10 \\ nfft = 2205$$

$$\text{max frequency} = \frac{SR}{2} = 11025 \text{ Hz}$$

$$\text{Frequency resolution} = \frac{SR}{nfft} = 10 \text{ Hz}$$

- Some objects only have to ~ 3800 Hz and others more like 2600 Hz.

- Obviously going to just use 0 - 3800 Hz values. e.g. only take first 3800 values.

- Also just use the tables & n. producing would be better to create figures like a scatter graph.

- Could just increase time resolution and lower frequency in order to get less features and better values of samples for fitting. Consider their time scale is what we're after.

OSRA problems with cap lengths, have now resolved to say if exact sec is 8.18 or 8, $t = 2.5e$

2 sec \rightarrow 64

5 sec \rightarrow 160 target the discrepancy between 8 and 8.18 for now.

11 / 4 / 25

\rightarrow A couple images appear later to start with [0, 0, 0, ...] vector, interesting...

14 / 4 / 25

+ one result soon \sim $\frac{\text{hop length}}{\text{sample rate}}$ (new samples).

Have decided on your delay features to best:

$$\text{yut}: \quad y = \alpha_1(1 - \exp(-\frac{t}{\alpha_2}))$$

order: 10 (0.255 either side).

$$\text{hop-length} : \text{int}\left(\frac{5R}{20}\right)$$

$$\text{Frame length} : \text{int}\left(\frac{5R}{10}\right) = \text{window length}$$

sample-rate: 22050 (5R).

delay - seconds: 7.5.

will append (mean, median, std) + 2 to
the extracted features

↑
time scale + Multigrid
(car) (car)

- 6 extra features for now.
- never mind, i) to do Stargate features
on the multigrids and time scales, the
to get 11 features from each layer,
- 22 extra features.

Seems to be very clear that other features are
making the performance worse...

20/09/25

- KDD paper wasn't ideal for what I wanted to
do
- not a great baseline as lots of parameters
were varied and then the results were
dramy pooled.

However newigs covid19 sounds paper has a good
openly accessible code and easy to modify
opensource + SVM pipeline that I can
add my delay features to.

- the experimental gifts went working as
nicely as I'd like → need to speak
with Picto.

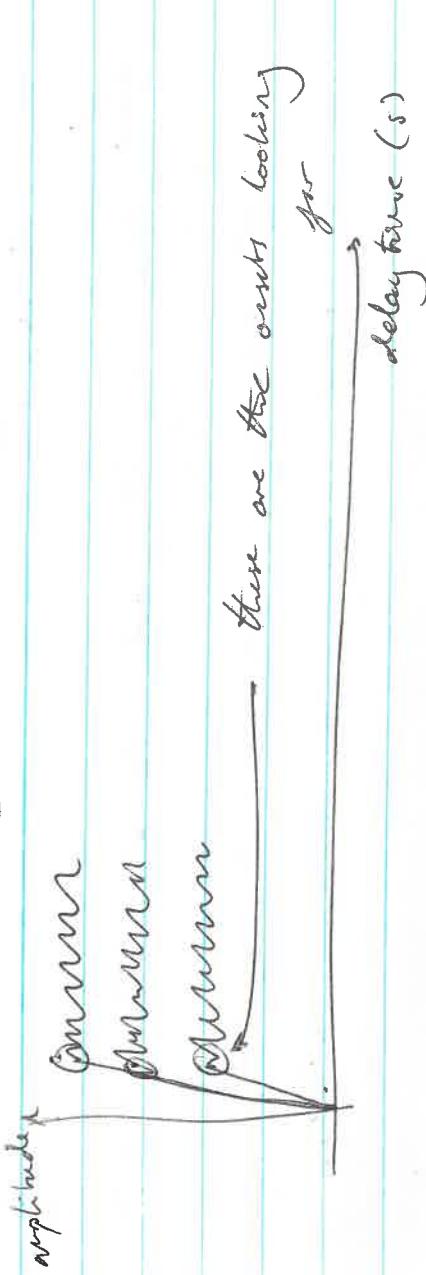
Lots of updates in the Powerpoint from this
week. ↗ Figures and comments.

29/4/25

Call with Pictro

- ① Can't say much more on OPERA joint.
→ Train 3 models for each larger
(2 extra per) to try and probe
if the nature of evolution on 2 and
② a. Current

- ② Feature engineering: Focus on voice data



- was this style that Pictro had been expecting
on the beginning
- will need to reduce order parameter
- evaluate the frequency dependence of parameters.

29/4/25

Protonic breath indices - 142, 2757, 7240

Post-minors update

- Implementing timescale multipliers and power features has considerably introduced improvements to the SVM classification albeit by only a 0.5% AUC in the ~~higher~~ ^{lower} case
- Not too convinced that time scale is doing a lot for the classification, the are very outliers in values to fit the curves
 - ↳ Need to try and combat this.
- Potentially over parameterizing by allowing power and time scale variables.
- Haven't explored smoothing the delay sing and to avoid picking out the wrong maxima.
- Proposed to take the delay and time (end-order) as the time scale and then allow the power fitting to give the curvature required.
- ↳ Have implemented these changes as well as modifying process to look for a threshold gradient ratio than early turning point
 - ↳ So far the values look much more reasonable in comparison to the raw data and the jobs are good.

→ Going to try and start writing up sections

Dashed layout

Table of contents

Introduction (How does it fit into the bigger picture) → may be try and take some physios like how the sounds generated by the throat and how they are modified by the prostate of men when eating like a filter.
↳ learn more about this.

Literature review → maybe introduce the papers from following and how they fit into everything.

now → Methods

- Now I make the classes (general)
get a generic explanation how to go from sound chip to delay spec.
- Feature engineering
→ ISO9 (maybe 6773)
→ time scales, averages, max/min values
(parameter values) →
 - OPERA
 - MAT, voice evaluation → (will take some work to break this down)
 - (parameter values)

Results & Discussion

- Feature engineering

- OPERA

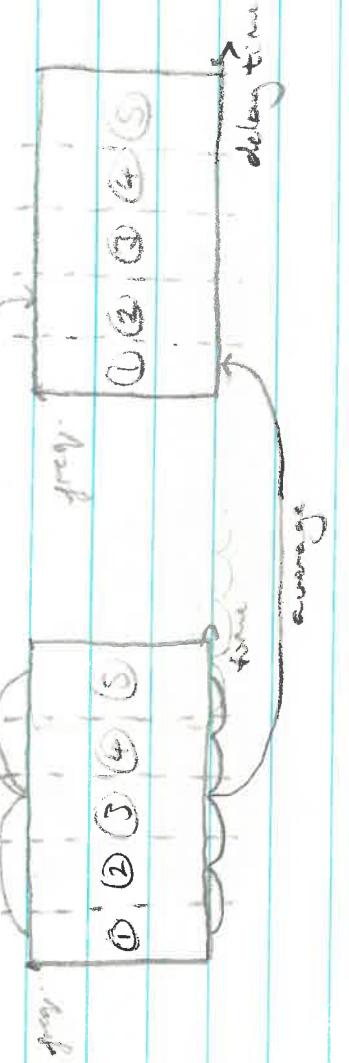
Conclusion

References

Appendices

- Might be an idea to separate two chapters one for future engineering (methods + results + discussion) and one for OPS same style.

Method: jogre. average.



8/5/2025

- by making changes to fitting yesterday performance has improved.
- just realised that by using different loss functions (original vs. open smile) are slightly worse - should probably change these back for better comparability.
- currently extracting VGGish feature to concatenate with my delay feature and see if they add there.

10/5/25

Sometimes D overlooked is the extraction of
open mouth and other feature sets on OpenRAT
itself that D can use as baselines.

- ↳ D have implemented very delayed features but
waste this frame and keep length.
the frame resolution is not good
- D'll have to play with these parameters
- gradient threshold been modified to 0.5

- ↳ The figure is nearly complete - make
the and labels bigger and make the
list list reader.

Inhibition on delay images

- on your delay the image, the decrease in
going to grow and grow.
- D+ tree is limited to the max amplitude - min
amplitude on this can (2 degrees of freedom)
- This is why the decorrelation saturates.
- The timescale of this saturation (the time
at which the and/or or completely uncorrelated)
it could hold information like some time
scale related to the random nature of oscillation
in vocal tract / lungs or something
- This is related to anxiety or things that
affect how the sounds we make, hence we
propose that finding this time scale may
add to classifying emotion.

11/5/25

Feature engineering + OPERA datasets.

- tried engineering feature for covid UK covid 19 sounds and covid world.
- fit works better for breathing raw with the smoothing approach potentially
- have been decline when implemented to covid uk and covid 19 sounds it's fair which is weird.
- seems to have slightly improved AUC for gender classification with covid world sounds again. multiclass vs strongest geek

$$\begin{array}{ll} \text{AUC before :} & 0.672 \pm 0.007 \\ \text{AUC after :} & 0.679 \pm 0.004 \end{array}$$

- costs factor. slower better I think its definitely worth changing to real - spectra 1. for this experiment → more focus on lower frequencies

Notes on the generation of sound.

'air coming from the lungs' generates 'sustaining total pressure' which controls the resistance given by the vocal folds (muscles and ligaments attached to vocal chords).

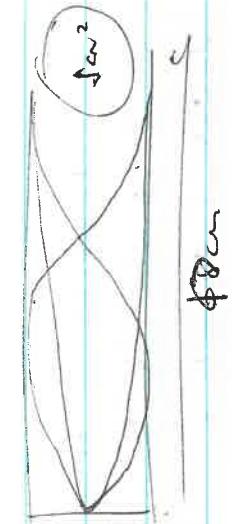
The alternating opening and closing of the vocal folds due to the force balance forms three bands of sound generated from (longer).

The factors that affect the frequency of vibration are the mass, length and tension of the vibrating string.

Laryngeal muscles modify the amplitude, fundamental frequency and waveform of the glottal pulses. The use lower frequency (more mass) \rightarrow slower oscillation, uses higher larynx high head.

Considering a more complicated model of coupled oscillators due to the interconnected muscles involved in generating sound introduces non-linearity and 'chaos'.

Mouth:



One end open tube

$$V_{sound} = 340 \text{ m/s}$$

$$\text{Fundamental frequency} = \frac{340}{(17 + 4)\times 10^{-2}} = 5000 \text{ Hz}$$

You get resonances due to the nature of the vocal tract \rightarrow it know this.

All the sounds are the fundamental frequency + harmonics modulated by the laryngeal muscles to change the amplitude of those present.

We introduce a periodic source commonly by introducing turbulence to the air escaping the tract for example the is sound by constricting the tongue to have a narrow gap (turbulent flow).

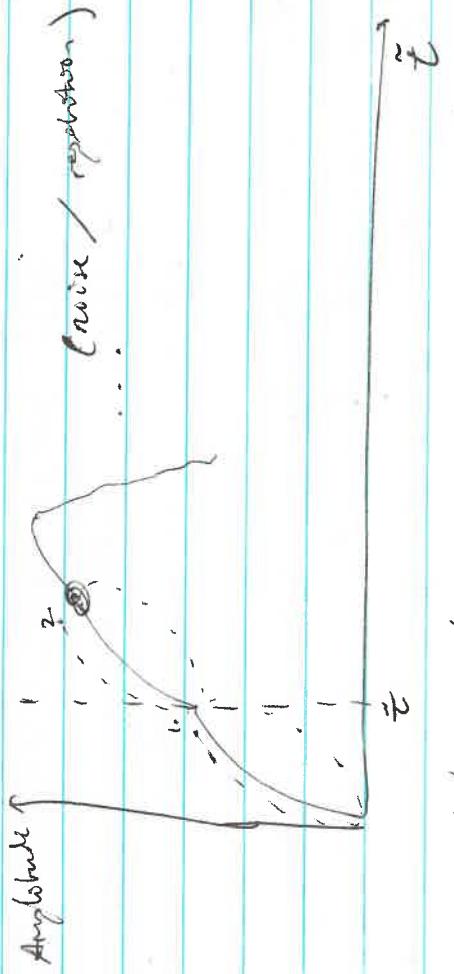
AVC, SPEC, SEN

- Have included sensitivity and specificity to the OPERAT metrics.
 - ScORI is much better in all metrics.

Med Delay Images

Made the transition to real delay images. It seems logical to target the lower frequencies / weight the tonescales there because there are where the sounds come from

- Main changes:
 - o created a cap on successive time scales to prevent the peak under for locks an to successive assets.
 - o rough assets look like so

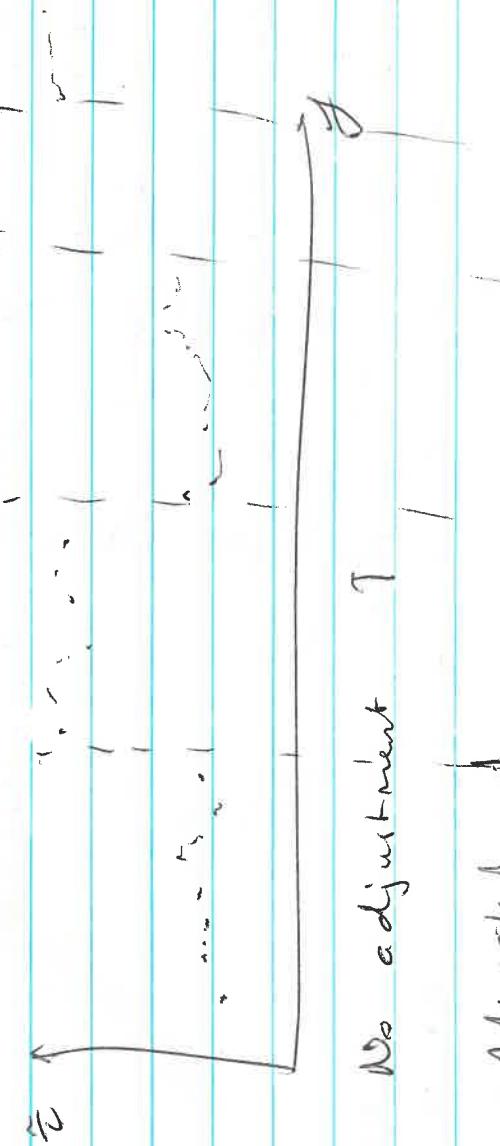


I believe one answer ~~should~~ due to some
time pressure on the date and the most important
to celebrate is the first / need to be
consistent.

I think this because the homescale with frequency scattered plots have large jumps also mixed with calipers other measures

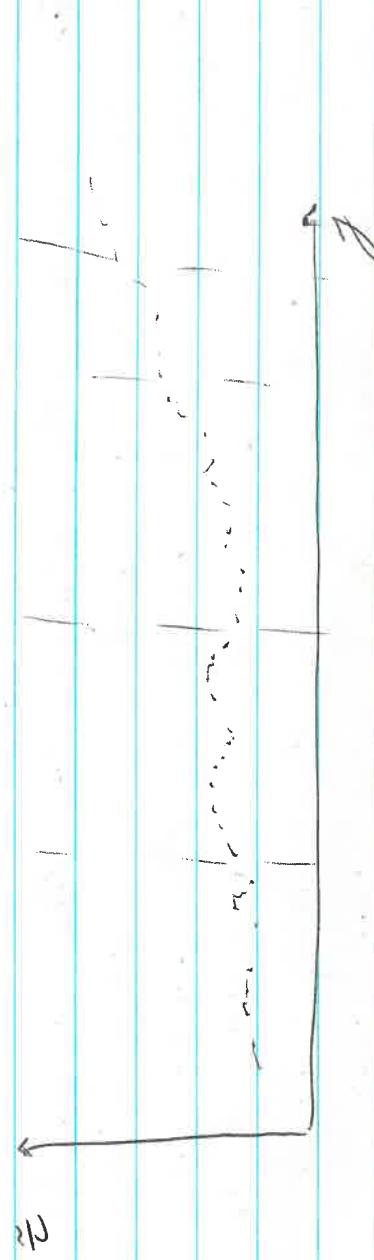
and if these were cut out there would be a more natural trend, hence I constrain the location of consecutive tonescales

of consecutive tonescales



No adjustment ↑

Adj. used ↓



14/5/25

Meeting with Pietro

→ Add random labels to `SCDF` to make sure AUC improvement isn't random.

→ Place the `skelag` profile shapes in `results`.
(example of each.)

→ Find where the tonescale values correspond to the frequency and any liberative or what they correspond to.
↳ Think the tonescales are on the order of how the layered muscles will fire to adjust the formants' position.

fundament of frequency (as an option to produce these frequencies themselves.)

- take average values of timescales for both classes
- box plot!

Some more situation on delay

Phenomenon at delay tones from our audience drop correspond to a timescale within the data.

timescale in which

For example, the changes to the long and muscles occur in order to clean artifacts different sounds by modifying the position of the formants or producing voice less sounds would manifest itself in the delay profile at some timescale at which these changes usually occur etc.

↳ Paper suggesting average syllable length 226 ms
≈ average timescale it's been fading.

After At this timescale, we propose that the delay - amplitude difference should saturate as when displayed by this answer. the delayed spectrogram is likely to be no more decorrelated from the original.

→ Breathy timescales are very short, very periodic behaviour not changing the breathing rate at frequencies between 0.5 - 80 Hz
→ unlikely to gain much info

However, you can imagine that by combining

Sounds or speaking sounds, the number of flex and nodding strokes are multiple times per second i.e. between 1 - 80 Hz.

- ! Make emphasis on formants and the role of flapping mouth muscle to convey social changes.
 - This should be a powerful emotion recognition feature.

One change to delay feature I'm extracting:

- 90% overlap quoted in a paper makes little difference but increase pixel density (going for that which minimizes the tradeoff in frequency). & 4kHz low pass.
- Try dry and with NewTTS.

15/5/25

Example condition delay could pick up:

Acousticology - abnormal increase in the clarity of whispered sounds during articulation due to the presence of long vocalization, which increases the sound speed and thus lowers the damping of the vocal signal due to the efficiency of transmission through fluid.

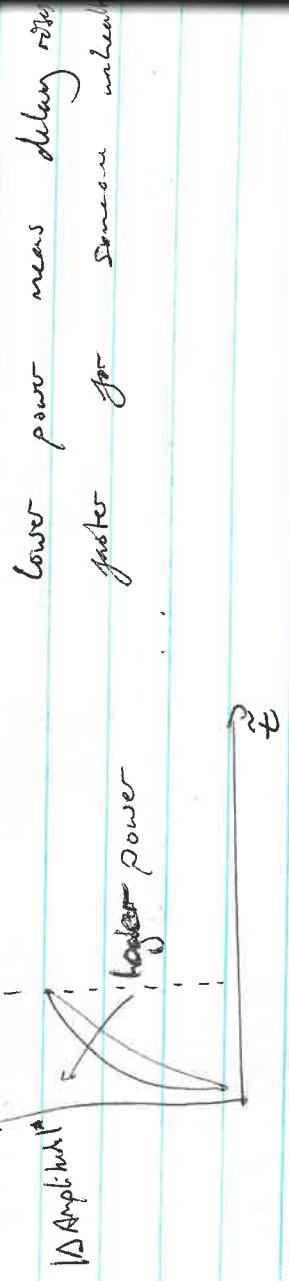
↳ lowered damping resulting in higher amplitudes of higher frequency voice components which would alter the character of the delay concert at those levels.

Also. Bronchophony check or the same but for spoken words.

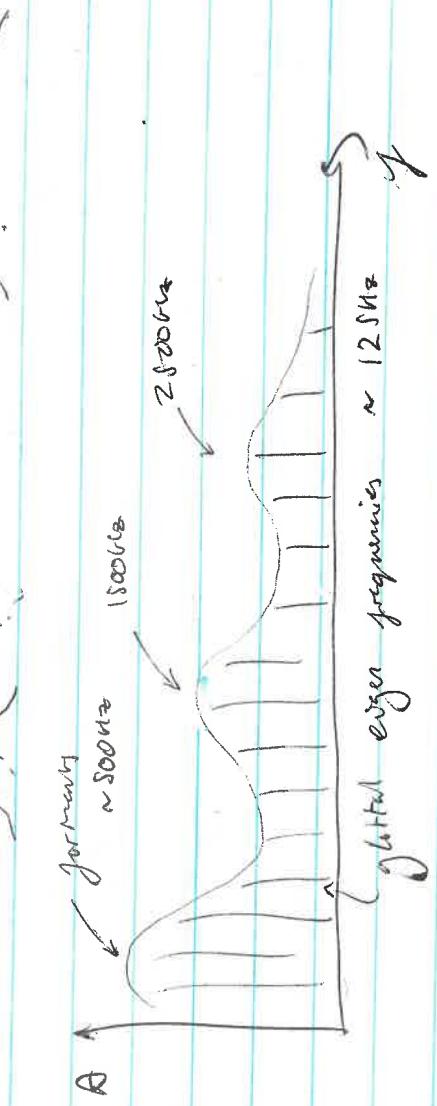
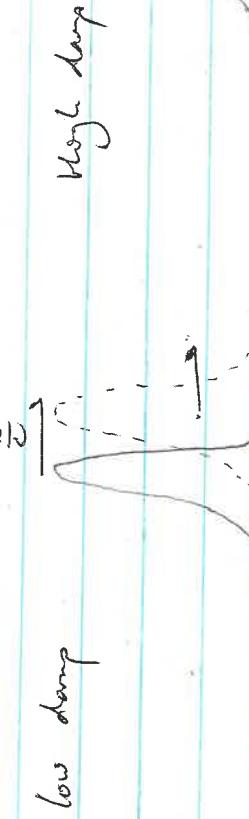
In these cases delay could pick out from a sentence whether someone was healthy or affected by larynx or the lung / vocal tract by the nature of the delay added at certain frequencies. (This checks that the low alters the efficiency of transmission in these frequencies.)

General reference from Poetru ↑ very good.

Why might Power parameter be higher for someone healthy vs. someone not healthy for cough?



If the higher frequencies are more damped, variation in their presence / not presence will be lower



glottal edges frequencies ~ 1250Hz

If you sharper the higher frequency harmonics and then move them around on the order of 0.25, the correlation between the frequency present before and after will be a sharper difference hence the decorrelation would now factor for more pronounced garment peaks

↳ the presence of liquid or general changes to the structure of the throat during illness could then influence the parameters extracted from the delay onset.

Generally associate lower power with more rapid rise in decorrelation

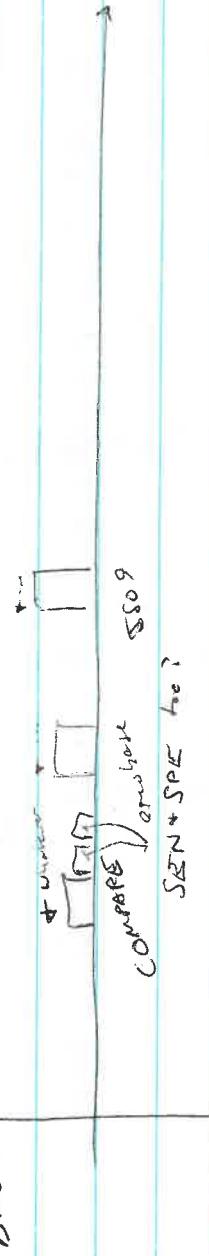
↳ could plot this / note on appendix

→ gather info on decorrelation → more stop start / what was described earlier.

ECGMI → stethoscopic data,

How to display the AUC values with and without delay:

Task → Overall AUC



→ Tables of values in appendix

Features to extract

Covid Urk : 6373 ✓ + 988 ✓ + delay [cough
break]

Covid 19 sounds: 6373 + 988 ✓ + delay [cough
break]

Coughvid : 6373 + 988 ✓ + delay

HF lung : 6373 + 988 ✓ + delay

SC EMI : 6373 + 988 ✓ + delay

Coughware : 6373 + 988 ✓ + delay

KAUH : 6373 + 988 ✓ + delay

(For same 988's need to remember old delay
features are appended still.

Got all AUC S&N STS values for New DPs,

ST809 C = 0.001 (Only power parameter)

probase C = 0.001 (only power param)

ComPAK C = 0.0001 (all powers).

16/5/25

Covid like cough $c = 0.0001$ Task 1

Covid like breath $c = 0.0001$ Task 2

Covid like cough $c = 0.001$ (0.0001 for 6 tasks)
Task 3

Covid like breath $c = 0.001$ Task 4

Cough and sound board $c = 0.0001$ Task 5

cough and gender

.. Task 6

vellos → only power of value. $c = 0.0001$ Task 7

Cougher sex $c = 0.0001$ for smo Task 8
 $c = 0.0001$ for Gout.

(smoker omitted)

ICUH $c = 0.001$ for end. Task 9

$c =$

Modelling based ~~at~~ by hand

Cough : Amplitude - 0 - 0.8
Breath : Amplitude - 0 - 100
Power - 0 - 1.2

Breath : frequency - 0 - 0.8
Amplitude - 0 - 80
Power - 0 - 1.2
COPD : Amplitude - 0 - 15
Power - 0 - 0.8.

17/5/20

Line Plots

For coughs:

(multiple of 1000)

take 1 at 120 Hz (~glottal pulse J)

take 1 at 520 Hz (~F0)

take 1 at 1600 Hz (~F1)

5 per plot.

Take 1 at

400ms resolution so order {J, IJ, IJ, 100

The reason that there is fine variability in the order onset is that controls were asked to perform three coughs or 5 or whatever. The spacing was not very even or all and in some cases it all blends onto one.

Poer's References

- 1) phases of oscillation and the amplitude of vocal oscillation have proven to be a good indicator of covid-19 / more likely just respiratory disease.
Bottom line: covid-19 causes asymmetries in vocal fold oscillation (during environment) resulting in discriminative features for covid-19 in symptomatic people.

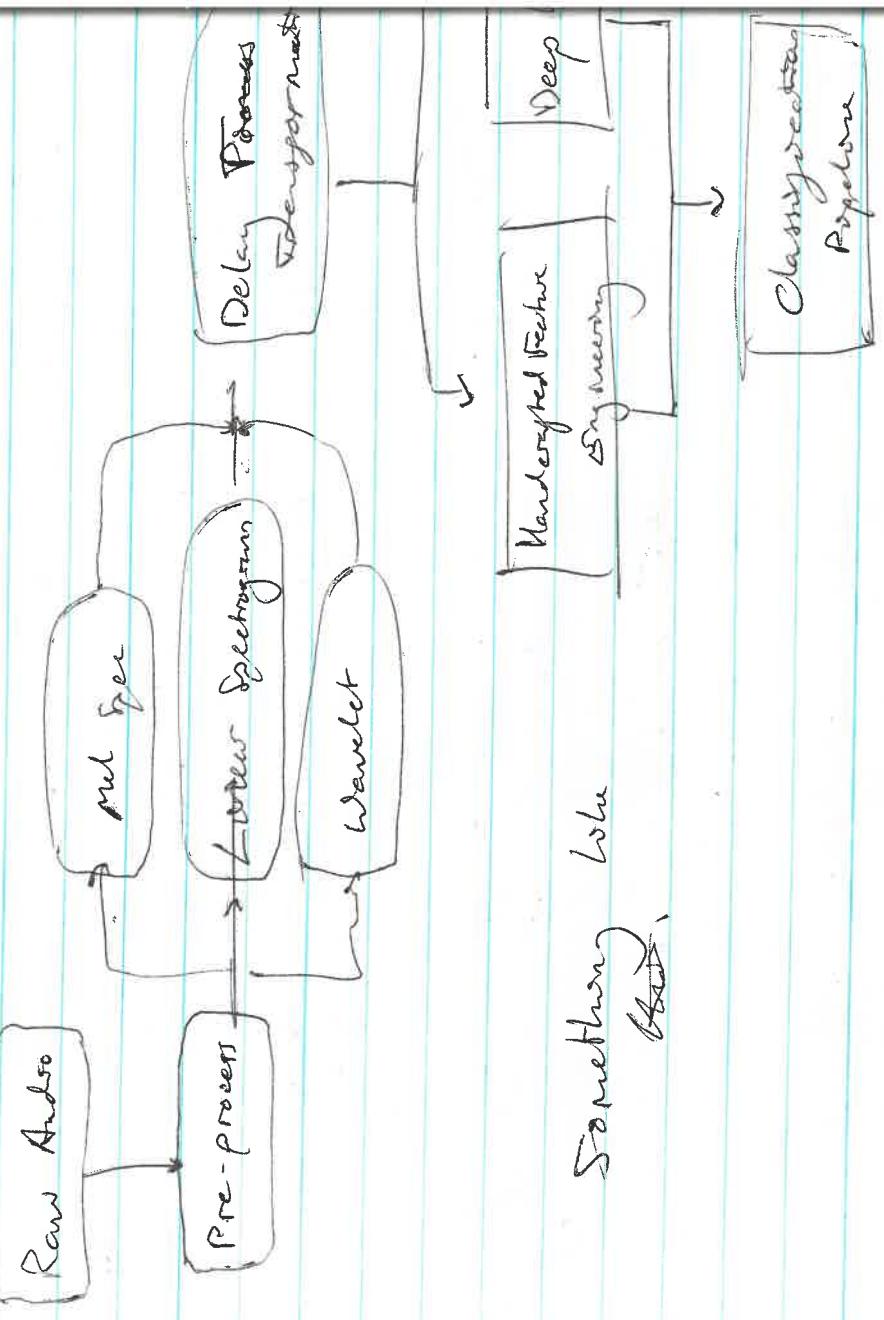
↳ Could the asynchronisation manifest as a periodicity on a timescale \rightarrow if found that could increase the rate of onset in delay feature

Vocalic sound plots in choosing:

- ii) tension asymmetry between the two vocal cords results in modulated signals due to the coupling between the vocal cords causing non linear behaviour
 - ↳ evidenced that vocal cord coupling can create global signal modulation of order 5 global pulse periods.
 - ↳ tension asymmetry as a result of tension fluctuation (back to (i)).
- 120 Hz → modulation frequency of 240 Hz within our range.
- ↳ How would this non linearly modify in a power spectrum?
 - One massive flume or three approach of the non uniformity across databases
 - Vocal differences between users make it clouded what delay gestures are from someone saying something different from other people and the sounds they make.
 - After The uniformity of the audio task can vary a lot for example if someone didn't use a quiet environment someone songs 3 or 4 or 5 times or the same saying between their song to all result in unique characteristics to the delay feature making it

hard to notice friends.

Flow Chart



Something like
(that).

18 / 5 / 25

How could I attempt to probe the
~200 Hz region of delay.
time resolution has to be $\frac{1}{200} = 0.005$

4 s clip \rightarrow 800 samples = doable.

Capped at ~ 2000 samples. (45 seconds)

$$\frac{16000}{200} = 80 \text{ hop lengths.}$$

19/05/25

~ 20% increase in AUC observed by
increasing from resolution and having in on

0 → 0.1 time scale for all model fitting

↳ this could be tracking model bias
in the order of the global pulse!

20/05/25

→ call with Archiv

→ try range of hop and frame length to evaluate the AUC increase.
[CNN (orange classif) /
LSTM (green classif)]

- What is the best of operable?
- ↳ might save (1024) and hop is (160).
- It's okay to put OPERA for each job.

→ try 512, 286, 128, 64, 32

Choosing hop length of non-sensored means that the fitting moves to theough length type peaks where I anticipate there isn't so much information.

will try 1024 instead

↳ 1024, 812, 286, 64)

Delay image equation?

$$F(\omega, \tilde{\epsilon}) = \int_0^t |F(\omega, t') - F(\omega, \tilde{\epsilon})| dt$$

$$\tilde{F}(\omega, \tilde{\epsilon}) = \sum_t |F(\omega, t) - F(\omega, \tilde{\epsilon})|$$

↳ delay feature → spectrogram

$$\tilde{F}(\omega, \tilde{\epsilon}) = \sum_t |F(\omega, t) - F(\omega, \tilde{\epsilon})|$$

X $\tilde{\epsilon}$ T - $\tilde{\epsilon}$
↳ length of clip

21/5/25

OpenSmile vs features on some instances
↳ drops

↳ not sure when only seems to be
for GSVS features, and ggs4

↳ would conclude that gets the hop-length
change but other why does it work on
~ 3000 clips before it?.

Potential RQ's to form my discussion around

Can these delay features become standard features?

↳ haven't seen generalization. To different
datasets / contexts - music, motion detection, ...

↳ having such fine scale images could
become problematic for computation
although it is not uncommon to use
different scales for different features)

- ↳ more rigorous method of testing required / are these the informed parameters → could potentially have log of divided asset be fit linear and extract power law?
 - include possibility of holding delay home region constraint i.e. 0-0.5 and going Wikang.
- Household Deep Learning contribute to their value.
- ↳ talk about what failed with OVRAT
 - ↳ currently the value is in the smallest part of the start of the range
 - maybe for extremely large passages of audio, this could be an effective way to reduce delta size and still find value.

Roberto's Applied sans

- Audio spoof detection - natural voice is very unique to people, thus feature may capture unique characteristics of voice not shared with other replaced or spoofs.

Whores for horses : 360 Hz 2000 Hz 4000 Hz
↓
120 280

22/5/25

Can delay features become a standard?

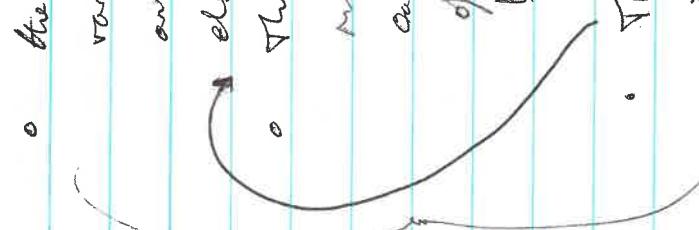
- They have proven to be strong distinguishing features.
- you chose a respiratory illness.
- discuss o achieving AUC's of up to 92% on classifying what you o previous maximum being 75%.
- think in o Clearly tells us value of being good weighted about that current Open source features + You the feature results from standard paper cannot make a take → All from Evaluation results.
 - o degradation ↑ → more likely respiratory illness.
- HOWEVER,

Generalization to other tasks and data

⇒ a

- Process of fitting needs to be refined
 - o the threshold of relevant approach or too variable due to problems of either low onto the length timescale or else.
- o The compensation between timescale and multiplier, clearly some reduction can be made. Maybe the ratio of two of the parameters or better still both.

Dependence of power per unit period



- The fine resolution spectrogram process is not optimal. Don't think that it is necessary to go so small hop-length,

~~Note~~ that because of the peak saturation
criticisms were raised and more constraints,
the same jobs may result from having
a smaller work force resolution.

Frame length of 1024 is same as CoarseST
so not needing to change that. Attempted
to modify hop for full comparability but it
failed on a few occasions.

With to this list on that I don't think that towards
is necessary if the gradient identification process
can be improved.

- Once these points are screened out, then
would need to generalize to multiple datasets
 - ↳ could reference Cleveland software
 - ↳ other data (COST - spectroscopic...)
- * Also, the frequency dependence might be
of importance for speech / song for
example. Just taking the just-noticeably-well
out of convergence in ODEs. There
might be organization in the positions
of the larger tone scales / features
could be derived up by them.
 - ↳ this could be echoing FO ...
junctants but clearly not & seeing
if there are better ways than OJ
junctants.

Once all of this addressed, then will know
but initial results suggest that they could
provide great benefit to respiratory and/or

→ Potentially deep learning discussion

- explain why I think my results aren't meaningful
- suggest ways in which it could be used to delay's advantage
- Applications outside of respiratory audio.
- Spoofer detection

- ~~should~~ looking away from the user at different feature timescales that could benefit me on other audio datasets.

24/5/25

- Big development - ~90% AUC was due to a change in process halfway through extra features? Thus invalidates claim and were unable to ~1% improvement.

25/5/25

will go down...

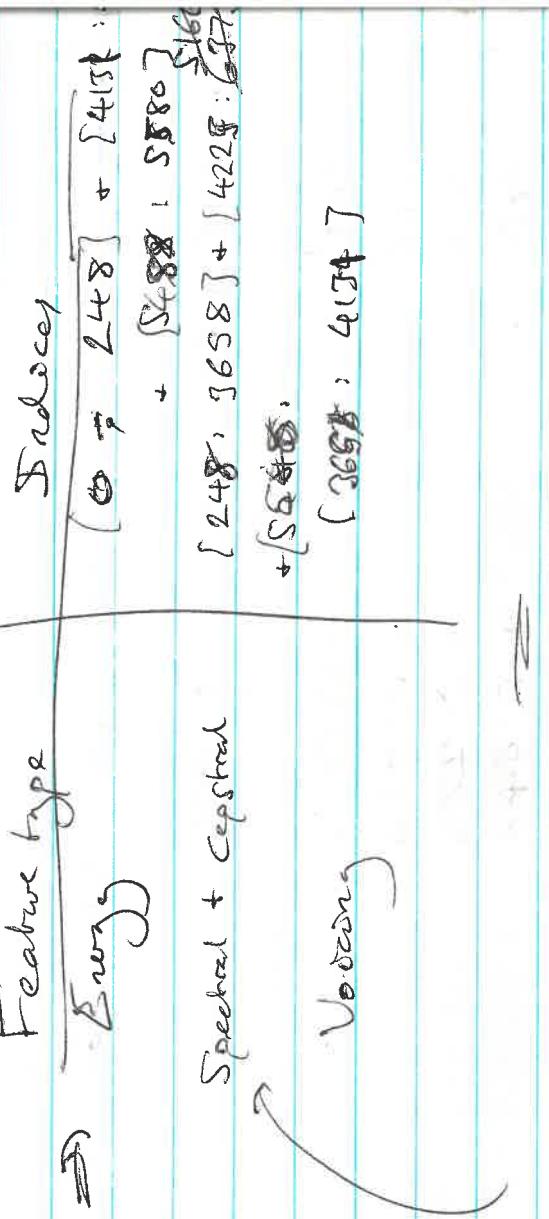
→ tomorrow hopefully back:

- try feature extraction fixing timescale at 20Hz → remove the feature entirely
- could refine approach achieve anything.
- try smooth gradient as a path

④ Back to Energy (diff functions?)

$$4135 - 4226 + [5488 : 4]$$

⑤ Back to Spectral + Cepstral
4227 - 6378

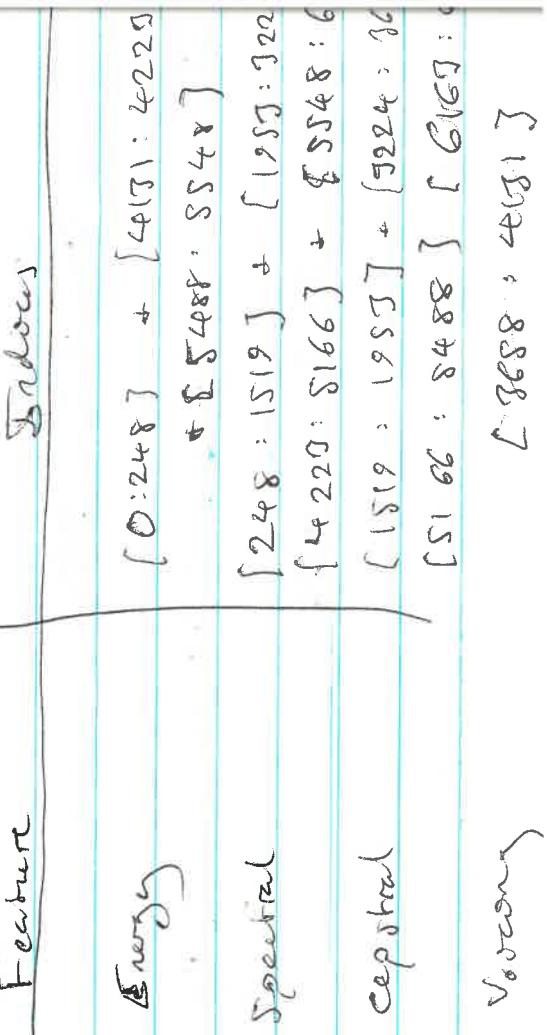


think this is most likely where information is contained.

separate Spectral and cepstral.

Spectral : [248 : 1519] + [1983 : 3224]
+ [4228 : 3168] + [5488 :
Cepstral : [1519 : 1959] + [8224 : 3658]
+ [3168 : 5488] +

Re + do :



28/5/25

Results: Spectral + Delay is best feature

What plots to include:

- Feature distributions on their own.
- Maybe baselines + feature types in same figure?
- Try vertical stack frequency distributions.

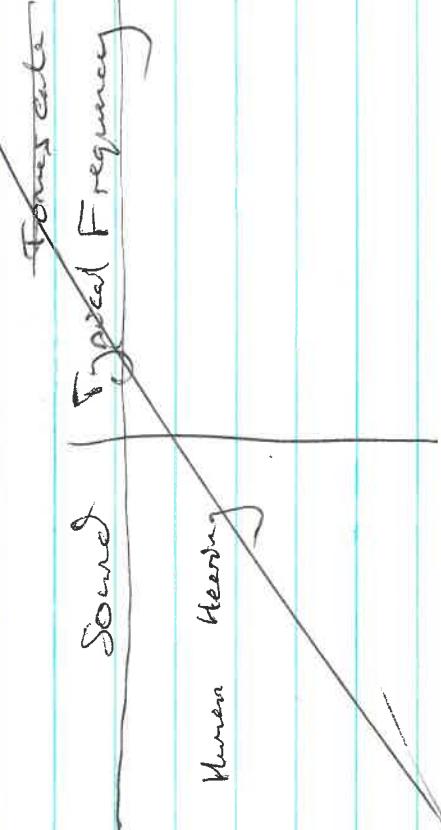
Report: Further recommendation, established features may not be the best for respiratory and asthmatic sets targeting in this case spectral features can out perform.

30/5/25

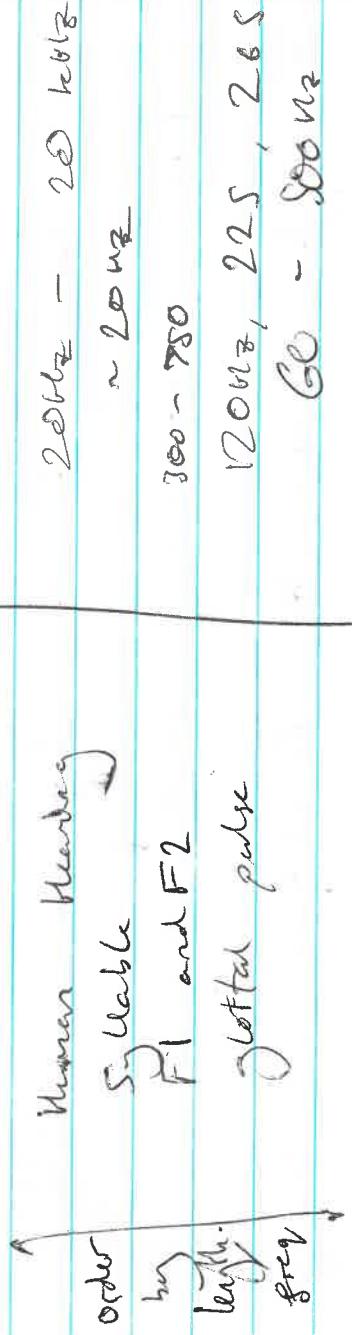
Sample: 23, 24, 25 - compare log (y) vs

Can find spectro data showing influence of log length and smooth or not necessary.

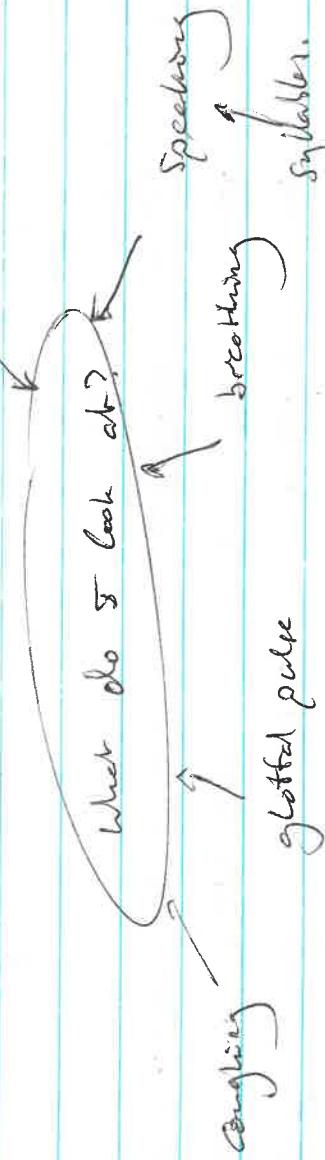
Table of diverse frequencies



Acoustic Phenomena | Typical Frequency



Fonants



16 / 25

Abstract:

- ① Respiratory and vocal features of coughing, breathing and speaking of organisms that can be harvested in medicine toward health-care applications. Currently, the acoustic features (harmonics) used to extract this information are outlined. My project.