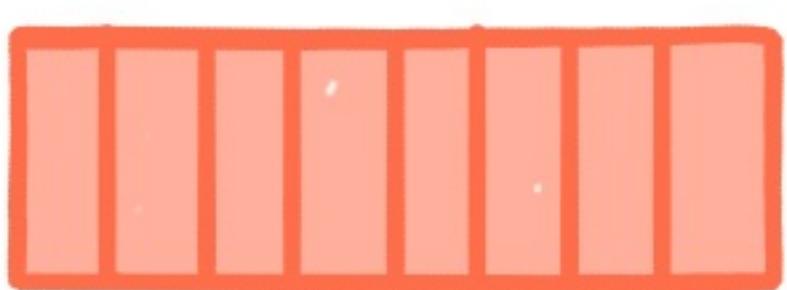
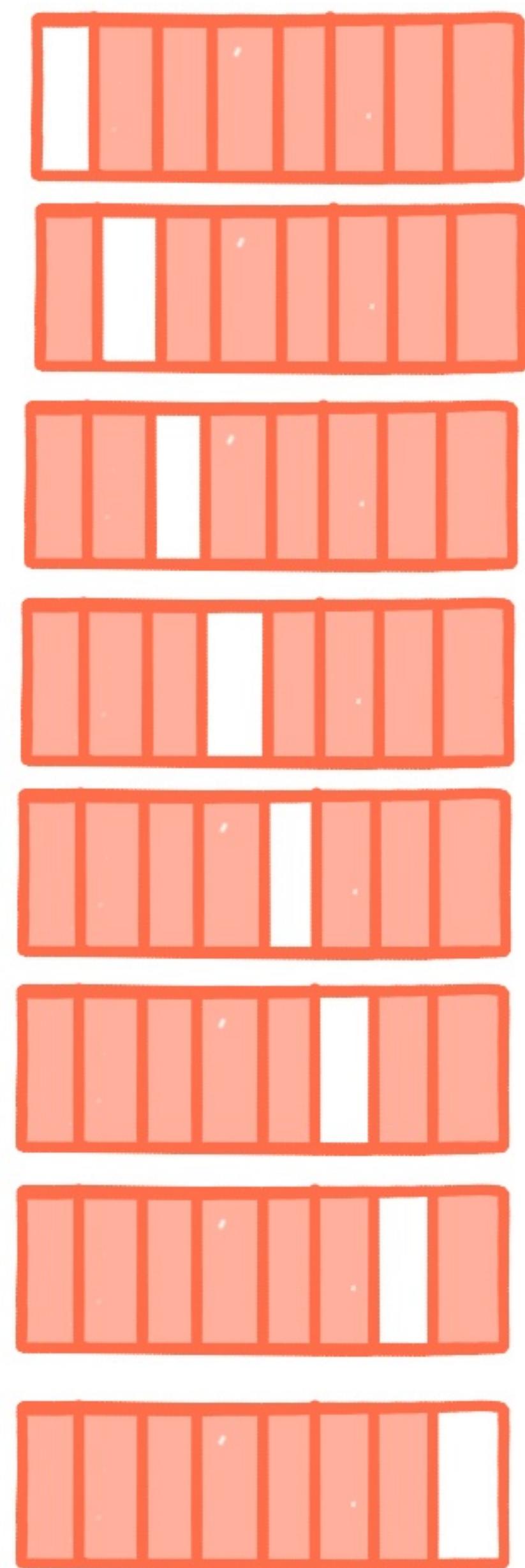


wave  
audio signal

read into  
buffer



process using  
sliding window



packing for fft result =  $\{ [DC, NY], C[1], C[2], \dots, C[N/2] \}$

vDSP\_fft\_zrip:

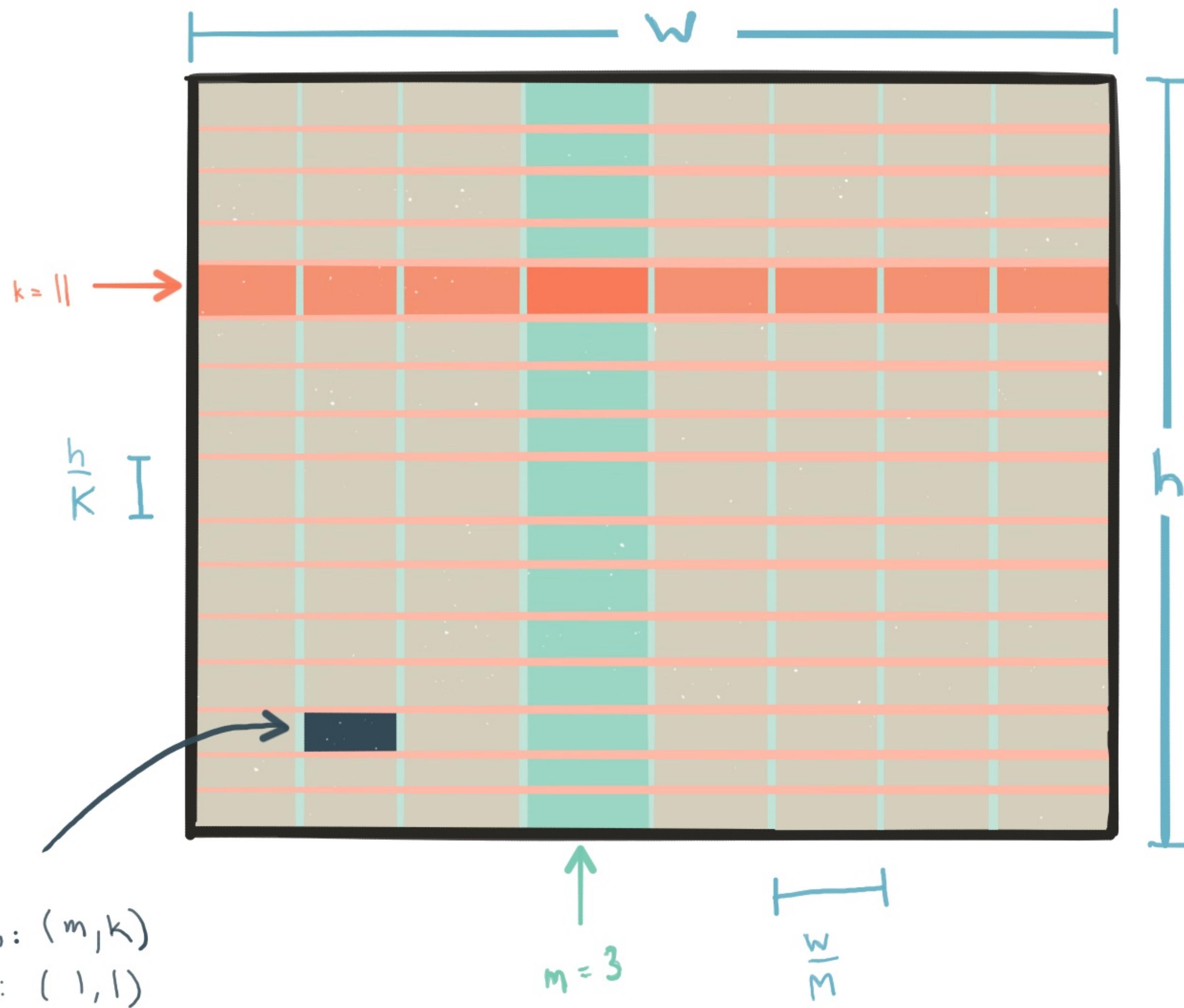
$$C_m = 2 \text{FDFT}(x_m)$$

$$\text{FDFT}(x_m) = \sum_{n=0}^{N-1} x_n e^{\frac{(-j 2\pi n m)}{N}}$$

vDSP\_zvmajs:

$$|C_n|_C = (\text{Re}[A_{nIA}])^2 + (\text{Im}[A_{nIA}])^2 \text{ for } n = \{0, N-1\}$$

$\text{Spec}(m, k) = v$   
 $m \in [0: M-1]$   
 $k \in [0: K-1]$   
 $N = \text{window size}$   
 $M = \text{total frames}$   
 $K = \frac{N}{2}$



origin:  $(m, k)$   
 size:  $(1, 1)$

Linear  $\gamma$  (frequency) Axis

$F_s$   
 $w(l) \quad l \in [0:N-1]$   
 $k \in [0:K]$   
 $K = \frac{N}{2}$   
 $n \in [0:M-1]$   
 $p \in [0:127]$

sample rate  
 window function's non-zero values  
 coefficient index  
 Nyquist frequency  
 STFT frame  
 MIDI pitch

$$F_{\text{coef}}(k) = \frac{k F_s}{N}$$

$$F_{\text{pitch}}(p) = 2^{\frac{(p-69)}{12}} \cdot 440$$

$$P_{\text{freq}}(f) = 12 \log_2 \left( \frac{f}{440} \right) + 69$$

$$\text{let } f = F_{\text{pitch}}(p)$$

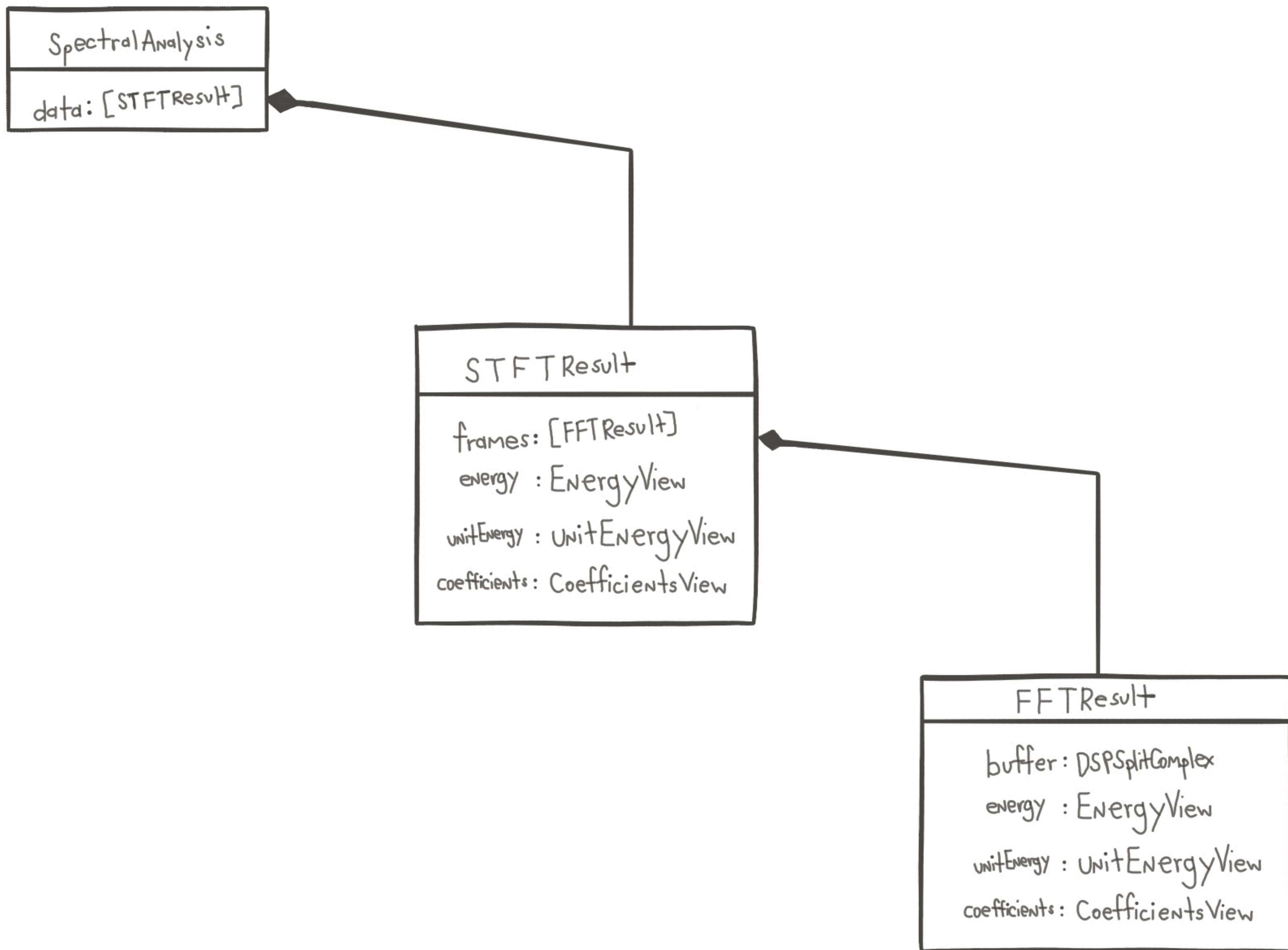
$$\text{then } f = 2^{\frac{(p-69)}{12}} \cdot 440$$

$$\frac{f}{440} = 2^{\frac{(p-69)}{12}}$$

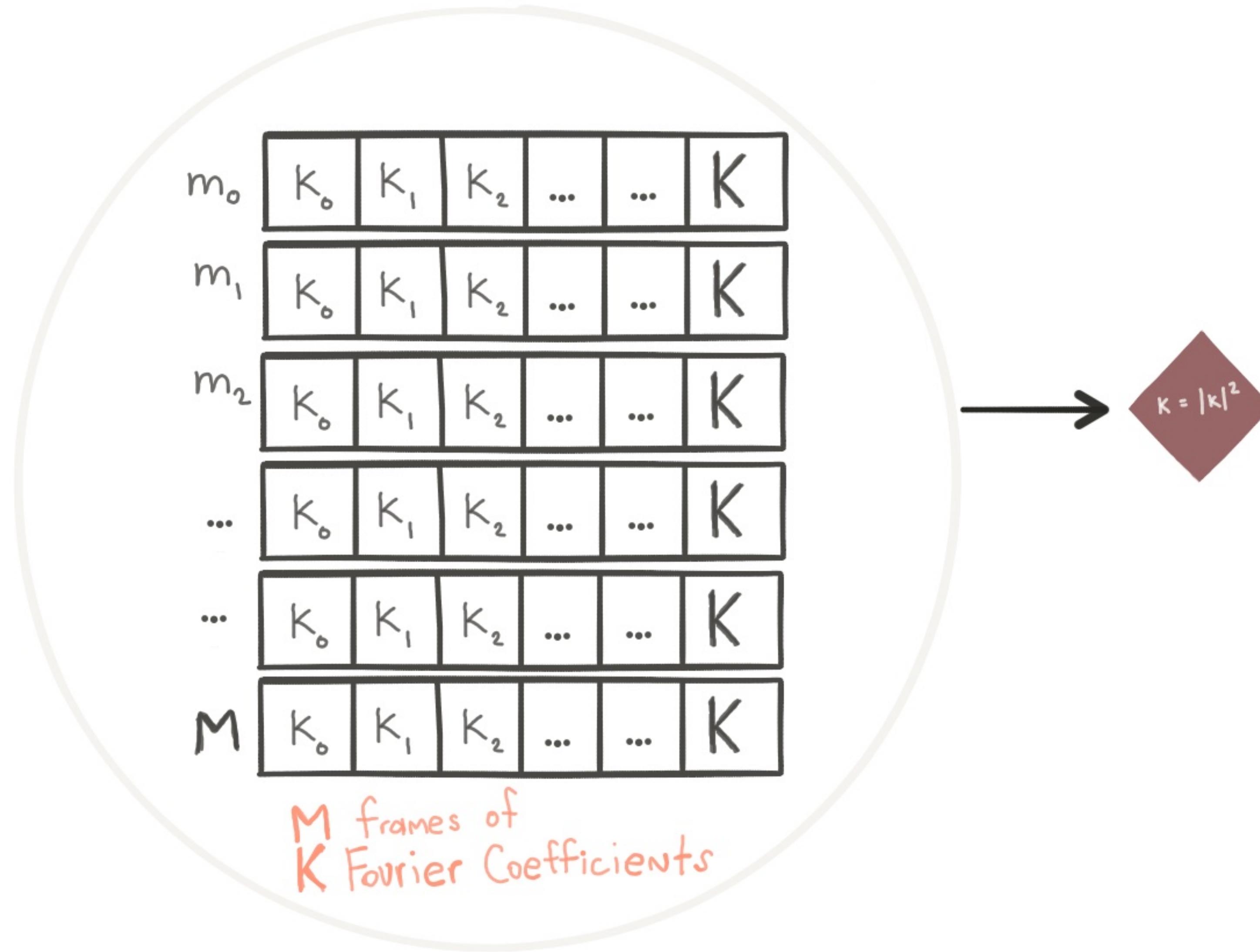
$$\log_2 \left( \frac{f}{440} \right) = \frac{p-69}{12}$$

$$12 \log_2 \left( \frac{f}{440} \right) = p - 69$$

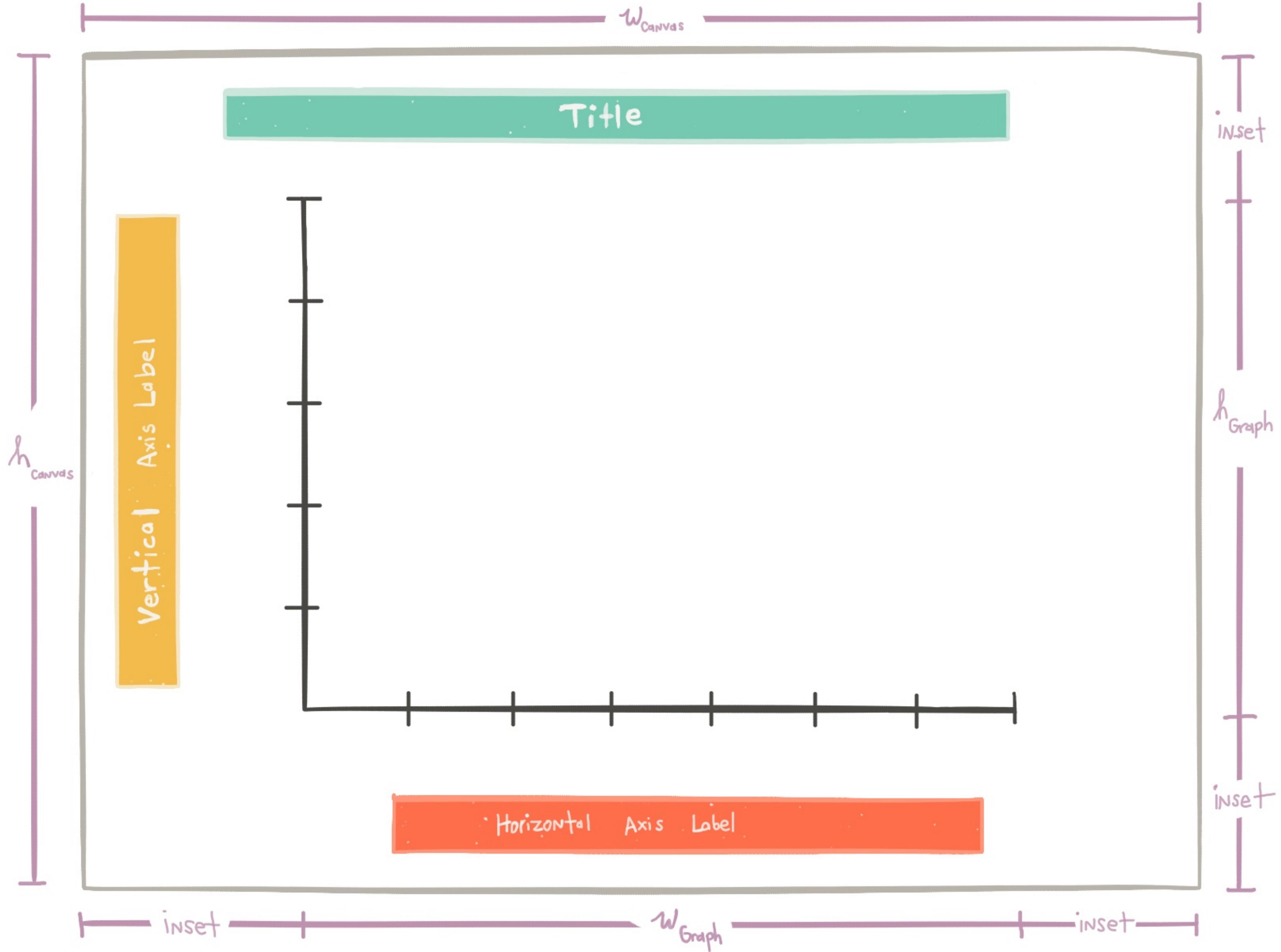
$$12 \log_2 \left( \frac{f}{440} \right) + 69 = p$$



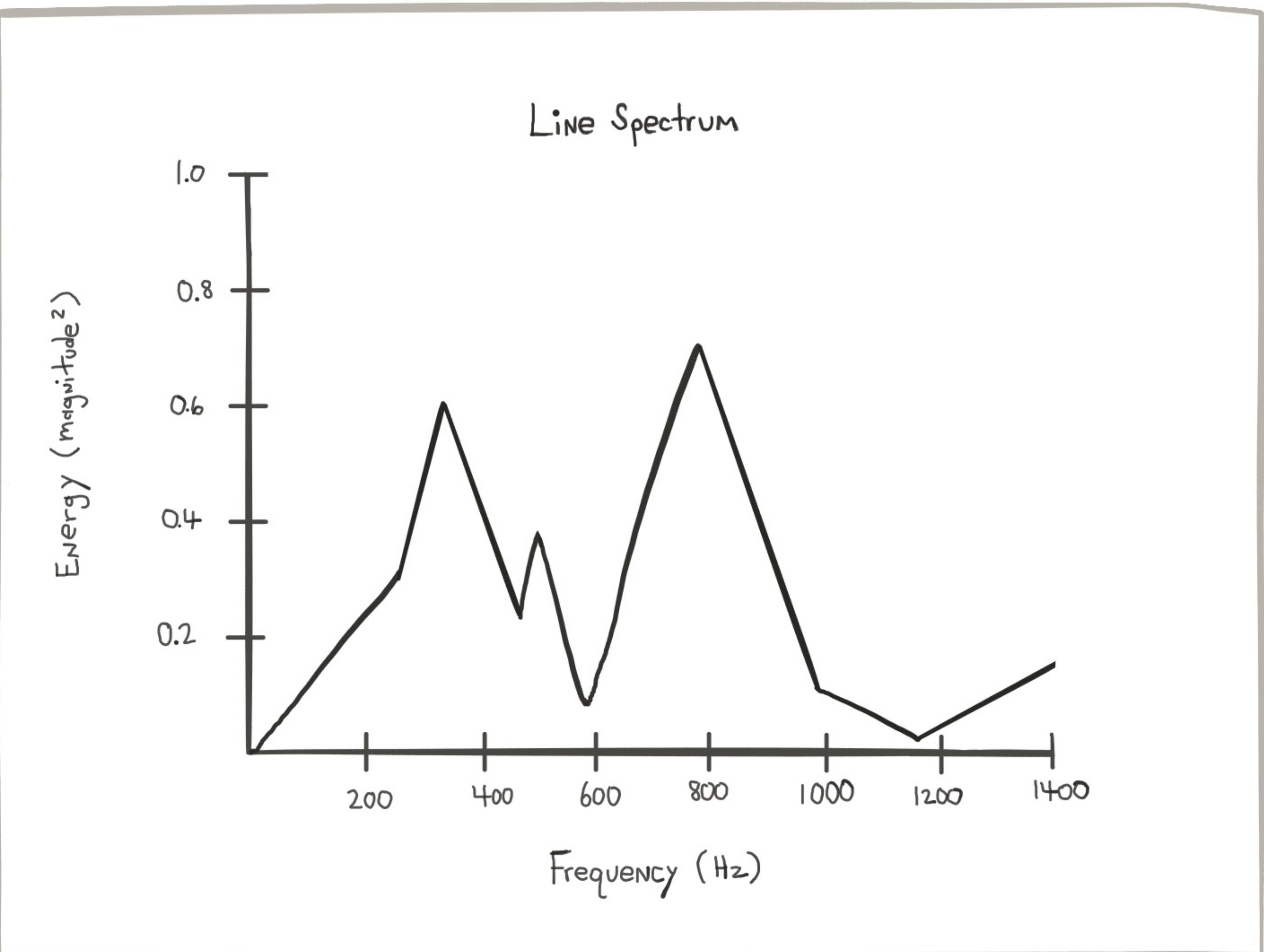
**Fill buffers until samples amount to a specific length of time.**



$$\rightarrow \diamondsuit \quad k = |k|^2$$

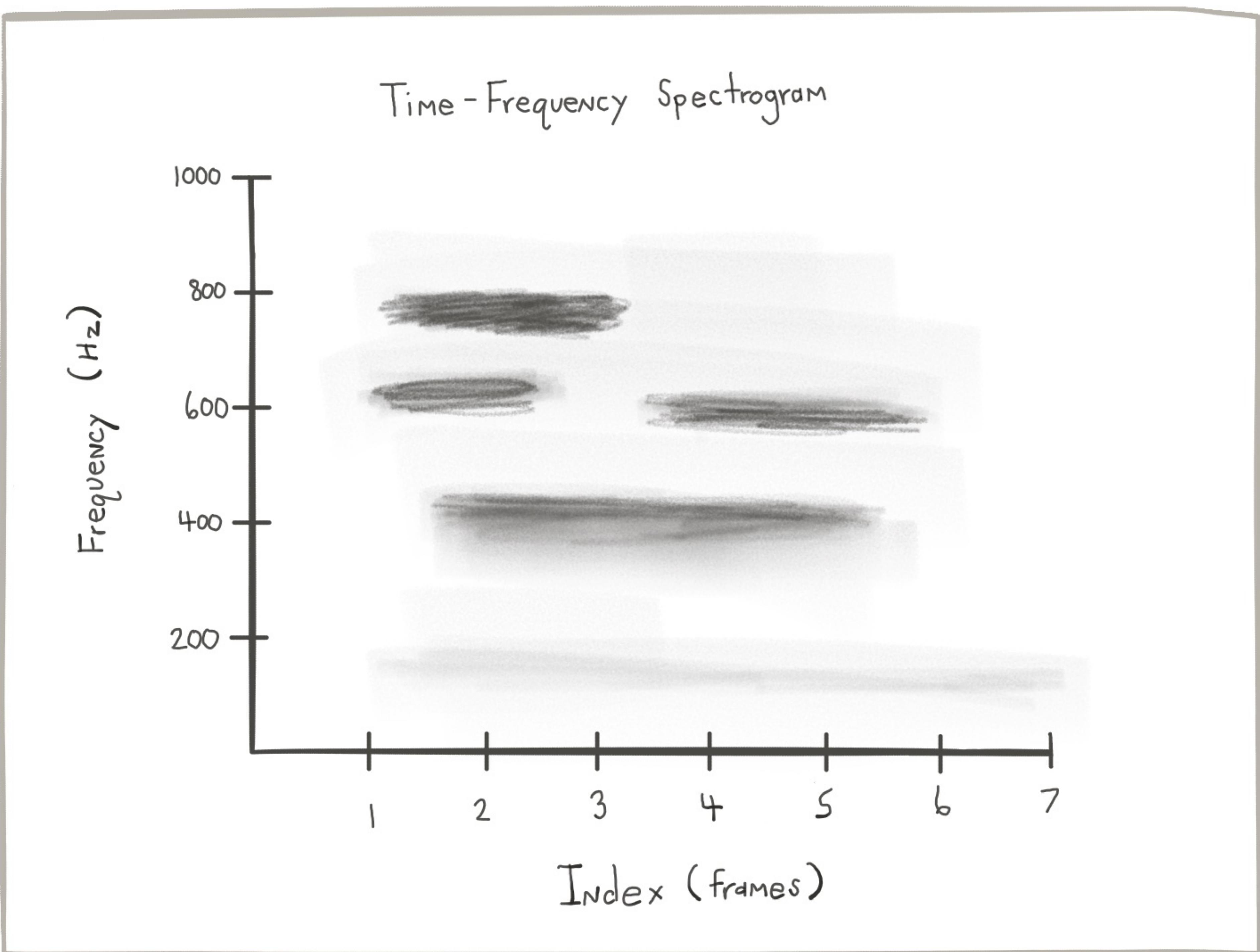


$$E_{\text{coef}}(k) = |x(k)|^2$$

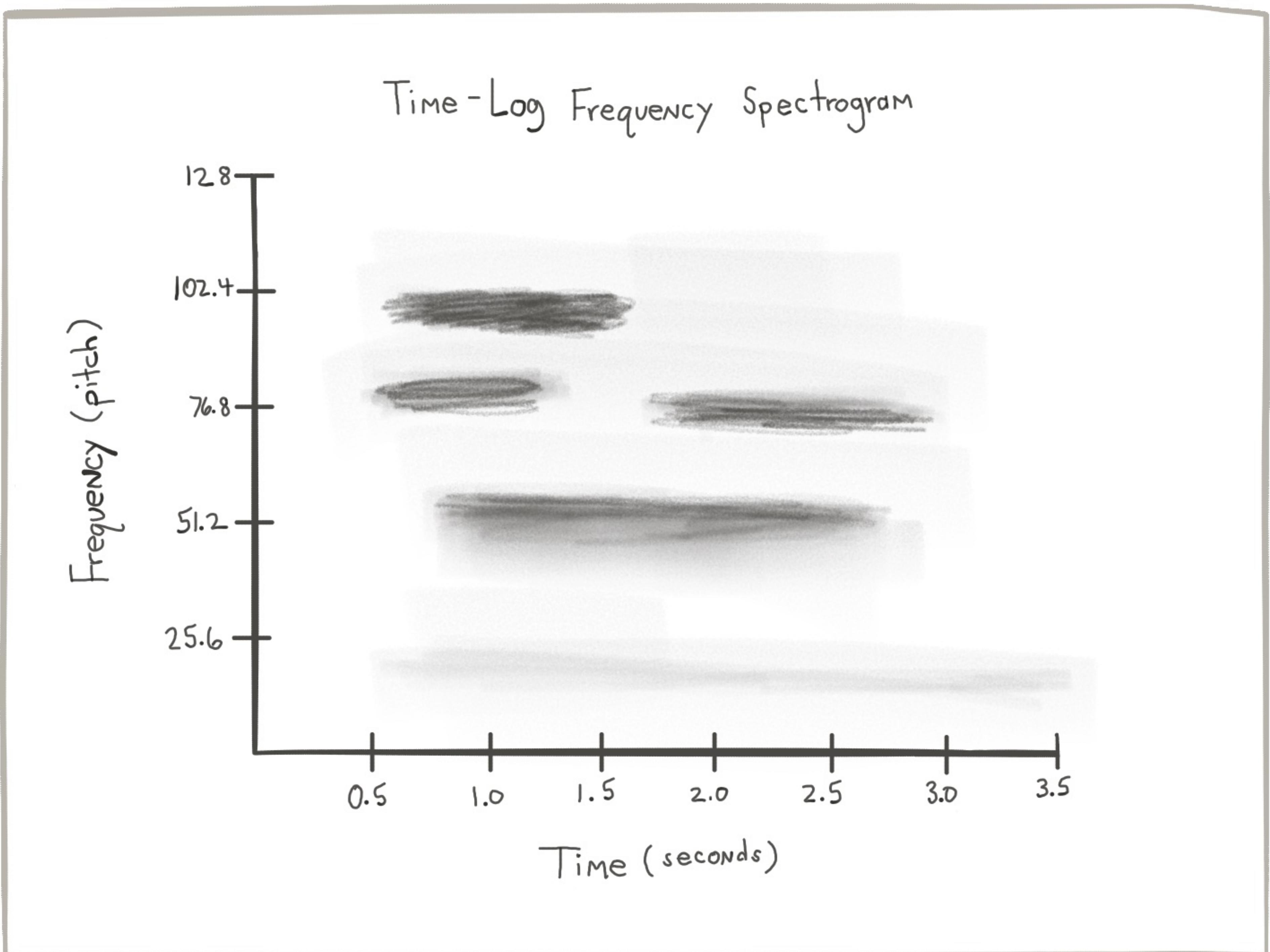


$$F_{\text{coef}}(k) = \frac{k \cdot F_s}{N}$$

$$F_{\text{coef}}(k) = \frac{k \cdot F_s}{N}$$



$$P_{\text{coef}}(k) = 12 \log_2 \left( \frac{k \cdot F_s}{440 \cdot N} \right) + 69$$



$$T_{\text{coef}}(m) = \frac{m \cdot H}{F_s}$$

$$F_{\text{pitch}}(p) = 2^{\frac{p-69}{12}} \cdot 440$$

$L_{\text{axis}}$  = the length of the axis in points

$$\frac{F_{\text{pitch}}(p)}{440} = 2^{\frac{p-69}{12}}$$

$L_{\text{values}}$  = the diameter of the range of values in unspecified units

$$\log_2\left(\frac{F_{\text{pitch}}(p)}{440}\right) = \frac{p-69}{12}$$

$$\frac{1}{L_{\text{axis}}} = \frac{x}{L_{\text{values}}}$$

$$x = \frac{L_{\text{values}}}{L_{\text{axis}}} = \text{ratio of unspecified value units per point}$$

$$\frac{Y}{L_{\text{axis}}} = \frac{1}{L_{\text{values}}}$$

$$y = \frac{L_{\text{axis}}}{L_{\text{values}}} = \text{ratio of points per unspecified value unit}$$

$$\begin{aligned} y \cdot n &= 200 \\ n &= \lceil \frac{200}{y} \rceil \end{aligned}$$

$$P_{\text{coef}}(k) = 12 \log_2\left(\frac{k \cdot F_s}{440 \cdot N}\right) + 69$$

given unit values of 1 Hz,  $L_{\text{axis}} = 1600 \text{ pt}$ ,  $L_{\text{values}} = 22,039.234375 \text{ Hz}$

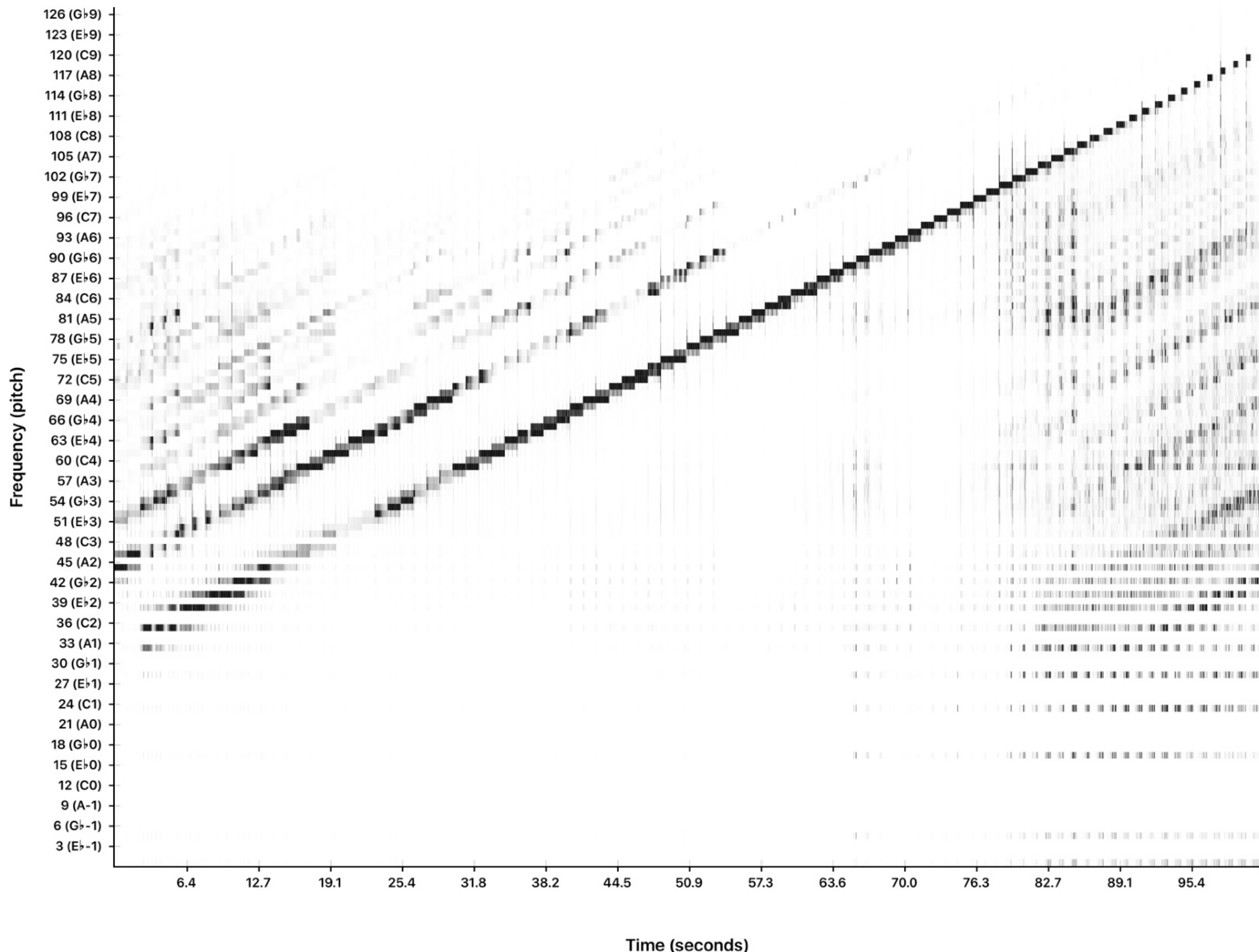
$$x = \frac{22,039.234375 \text{ Hz}}{1600 \text{ pt}} = 13.77452148 \text{ Hz/pt}$$

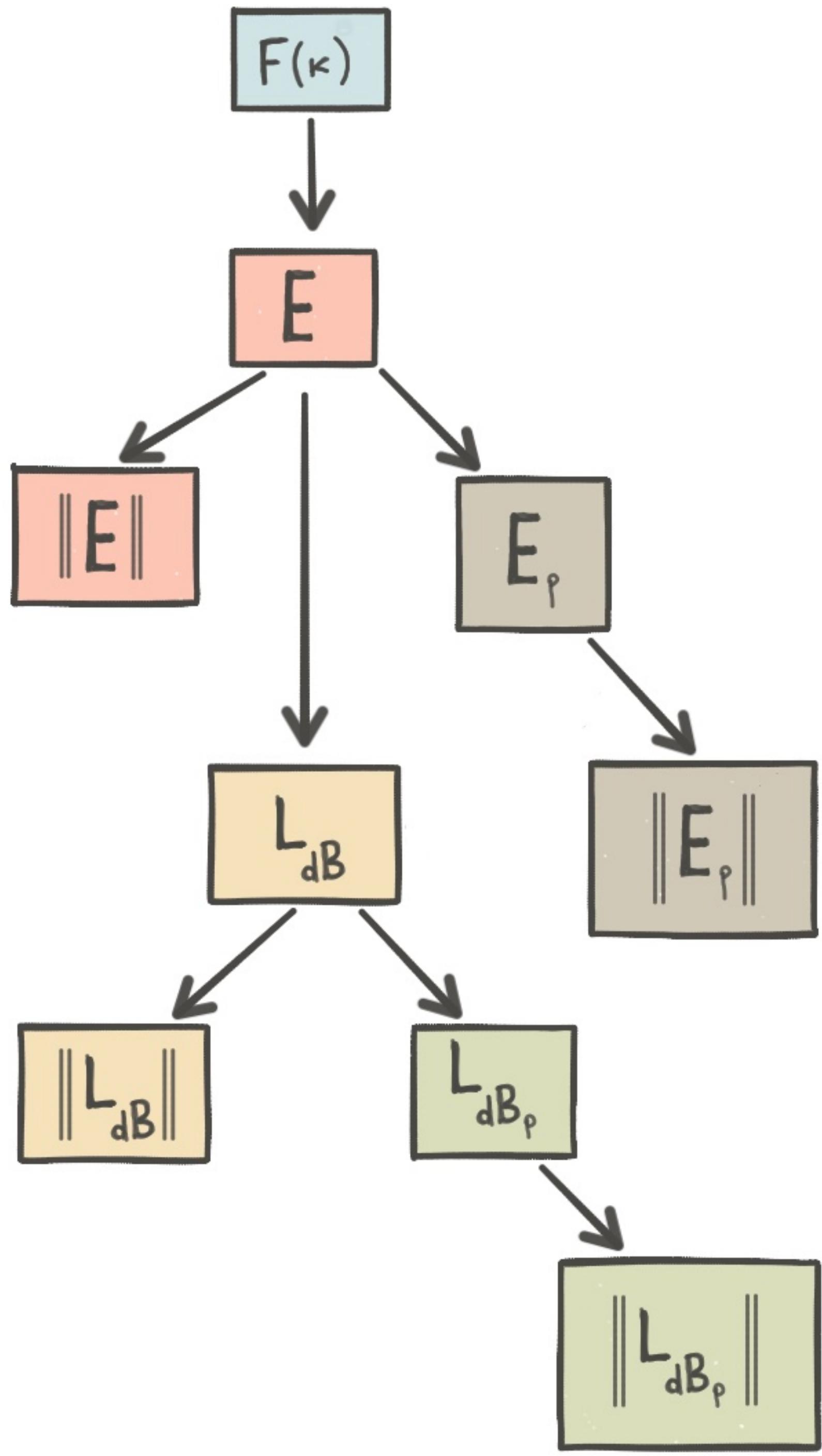
$$y = \frac{1600 \text{ pt}}{22,039.234375 \text{ Hz}} = 0.07259780321 \text{ pt/Hz}$$

$$\frac{200 \text{ pt}}{0.07259780321 \text{ pt/Hz}} = \frac{200 \text{ pt} \cdot 1 \text{ Hz}}{0.07259780321 \text{ pt}} = 2,754.904297 \text{ Hz}$$

$$2,754.904297 \text{ Hz} \cdot 0.07259780321 \text{ pt/Hz} = 200.0069478 \text{ pt}$$

## Log-Frequency Spectrogram of Chromatic\_Scale-Steinway\_Grand.aif



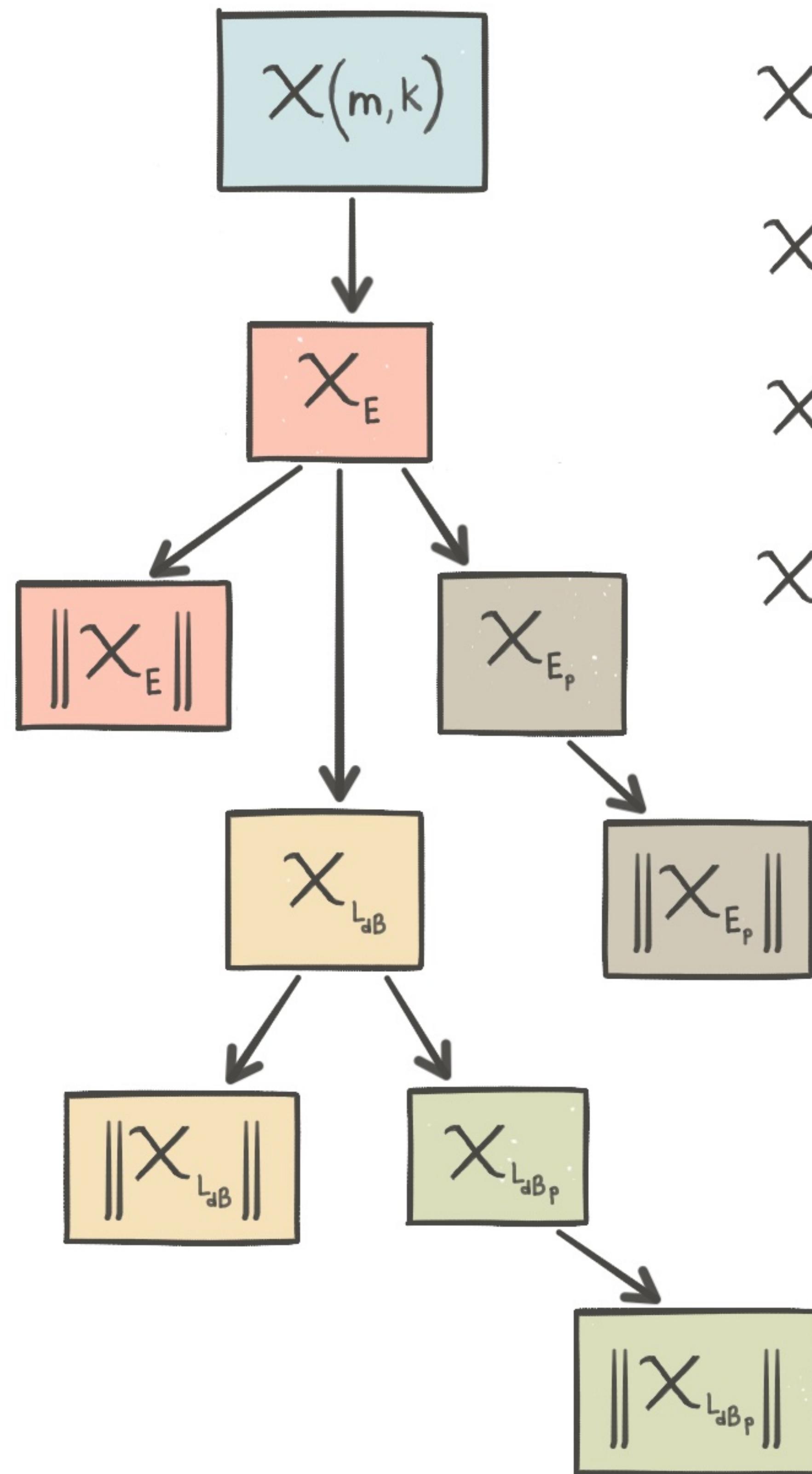


$$E = |F(k)|^2$$

$$E_p = \sum_{k \in P(p)} E \quad \text{where } p \in [0:127]$$

$$L_{dB} = 2 \log_{10} E$$

$$L_{dB_p} = \sum_{k \in P(p)} L_{dB} \quad \text{where } p \in [0:127]$$



$$X_E = X(m,k) = \sum_0^{m-1} E$$

$$X_{E_p} = X(m,k) = \sum_0^{m-1} E_p$$

$$X_{LdB} = X(m,k) = \sum_0^{m-1} L_{dB}$$

$$X_{LdB_p} = X(m,k) = \sum_0^{m-1} L_{dB_p}$$

$$\|f(x)\| = f(x) \rightarrow [0:1]$$

$$F_{\text{pitch}}(p) = 2^{\frac{p-69}{12}} \cdot 440$$

$$\frac{F_{\text{pitch}}(p)}{440} = 2^{\frac{p-69}{12}}$$

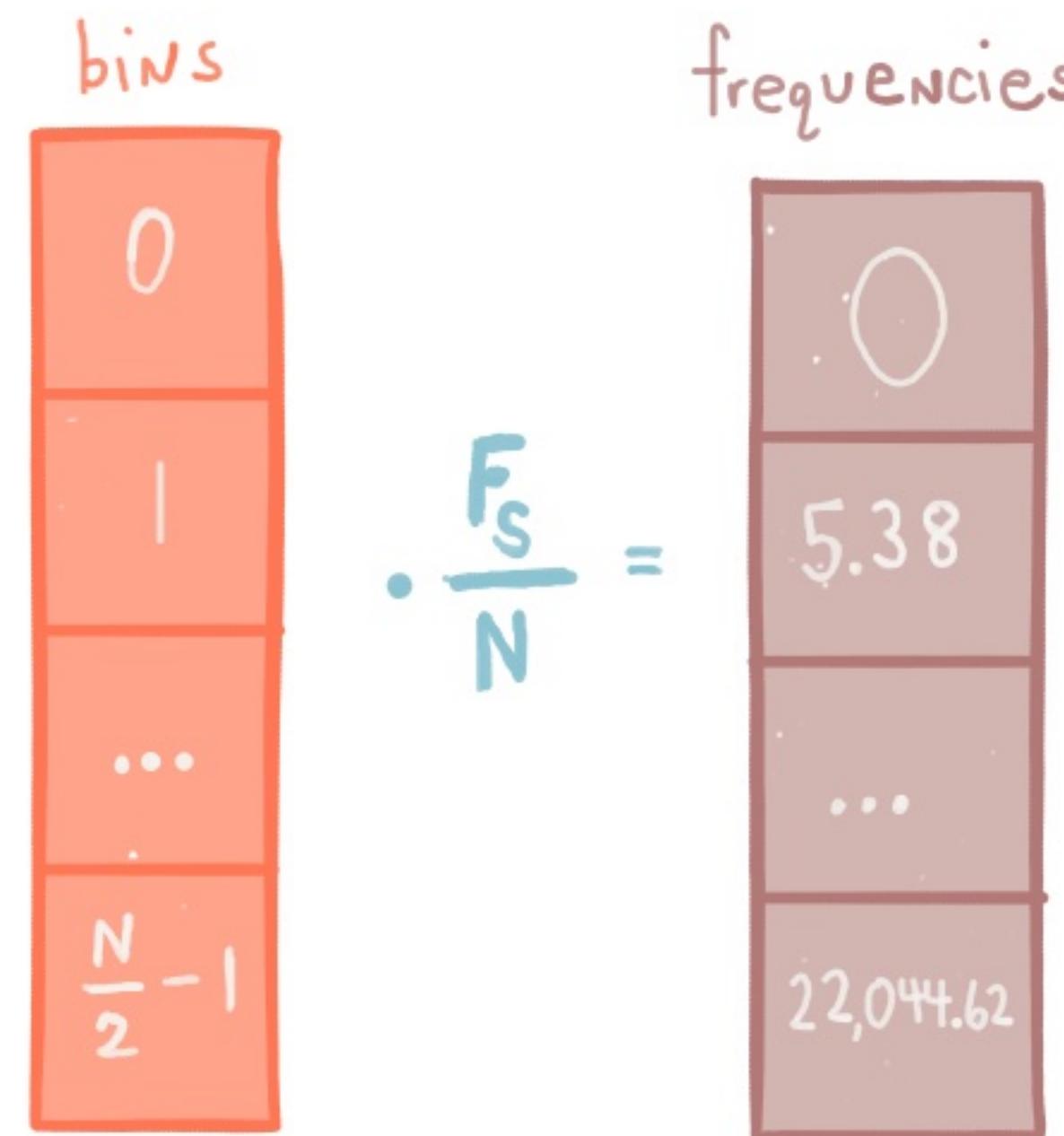
$$\log_2\left(\frac{F_{\text{pitch}}(p)}{440}\right) = \frac{p-69}{12}$$

$$p = 12 \log_2\left(\frac{F_{\text{pitch}}(p)}{440}\right) + 69$$

$$F_{\text{coef}}(k) = \frac{k \cdot F_s}{N}$$

$$F_{\text{pitch}}(p) = F_{\text{coef}}(k)$$

$$P_{\text{coef}}(k) = 12 \log_2\left(\frac{k \cdot F_s}{440 \cdot N}\right) + 69$$



$$\text{let } w = 8192, h = 4096, F_s = 44100$$

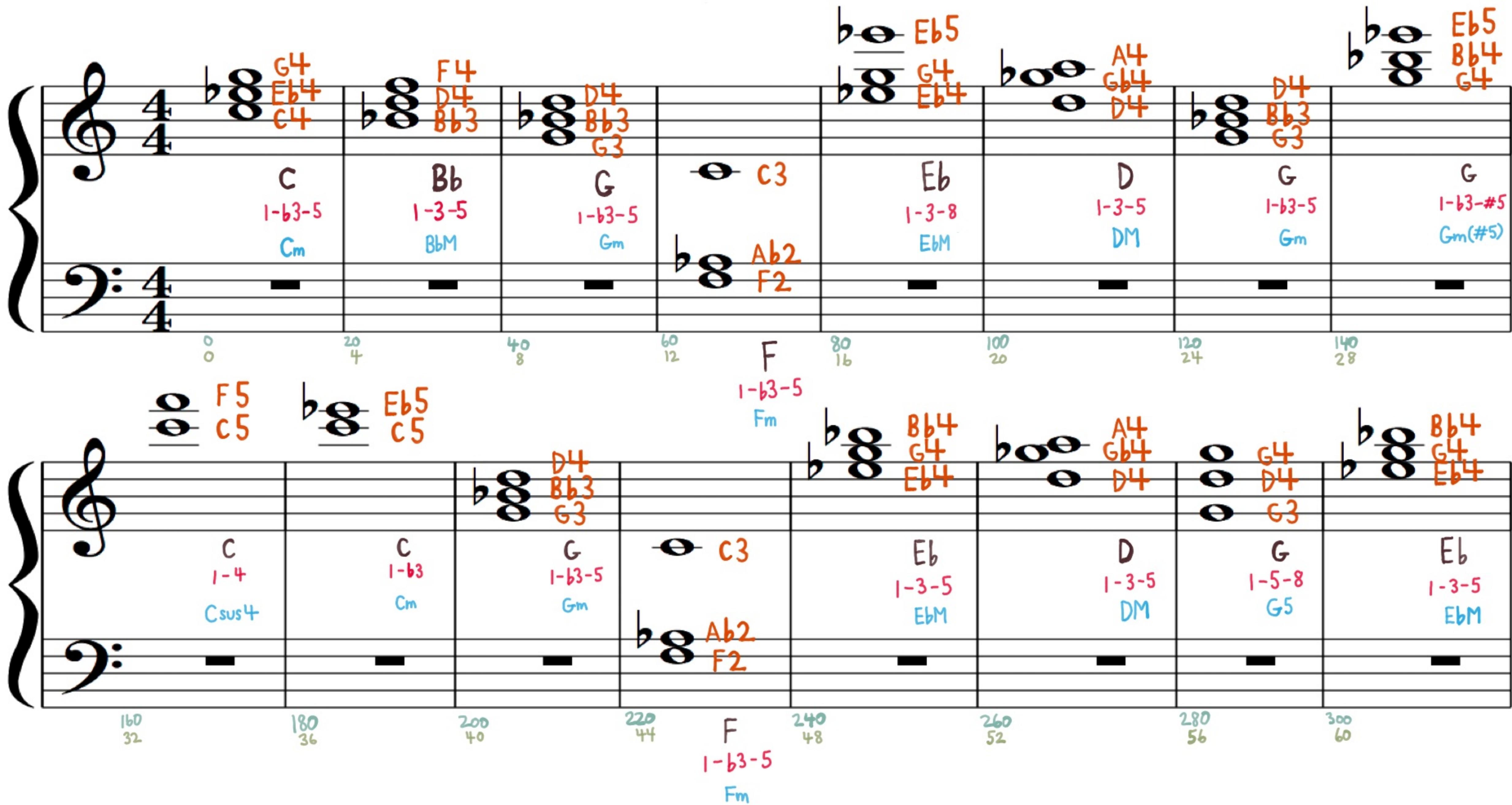
$$2^{-\frac{68}{12}}$$

Note

Root

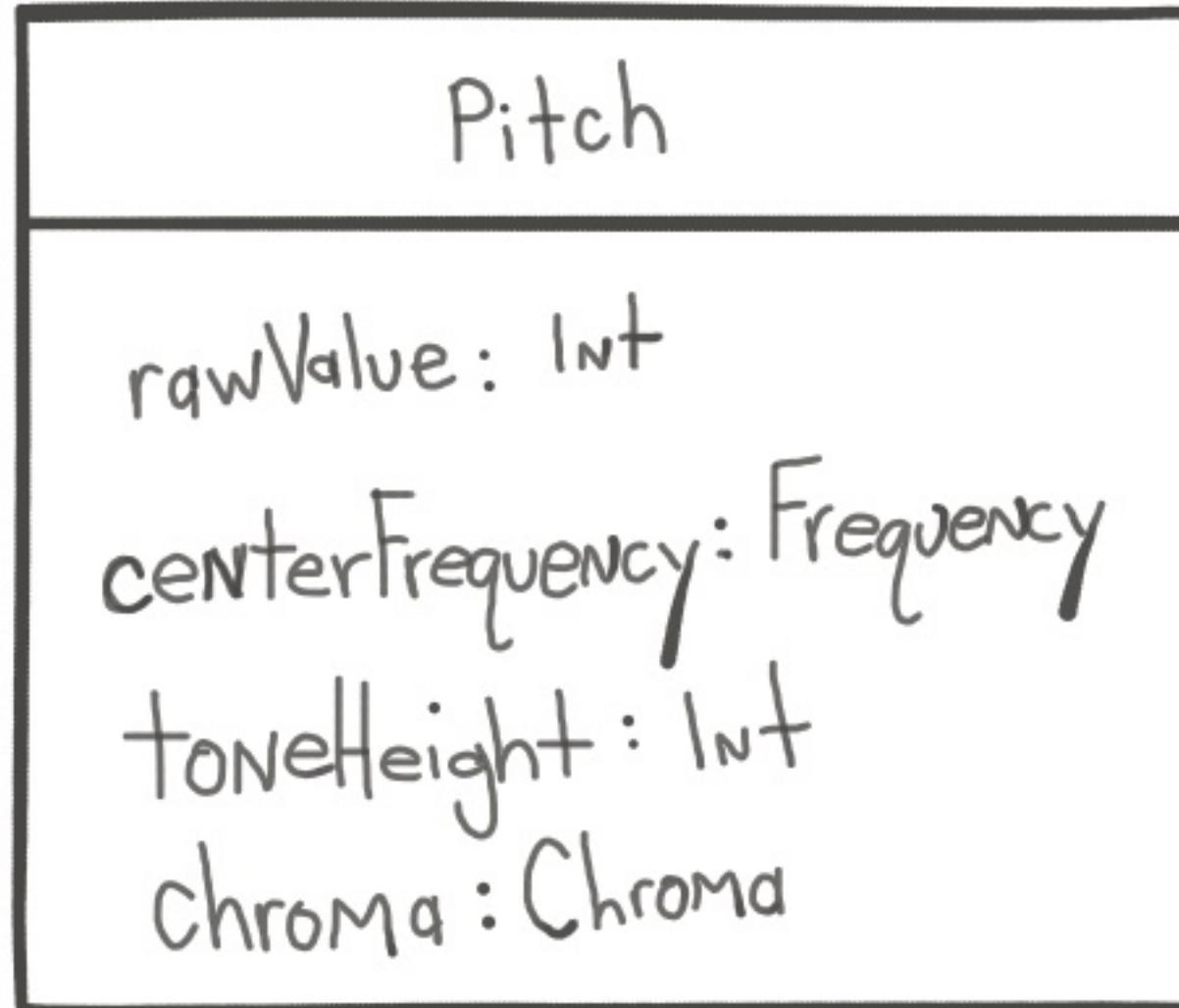
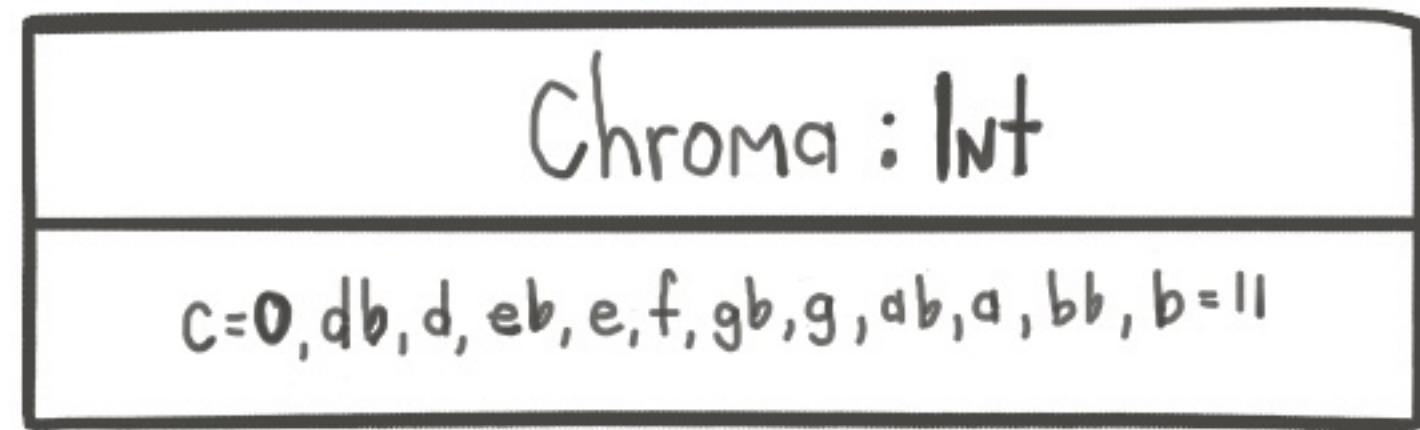
Intervals

Chord

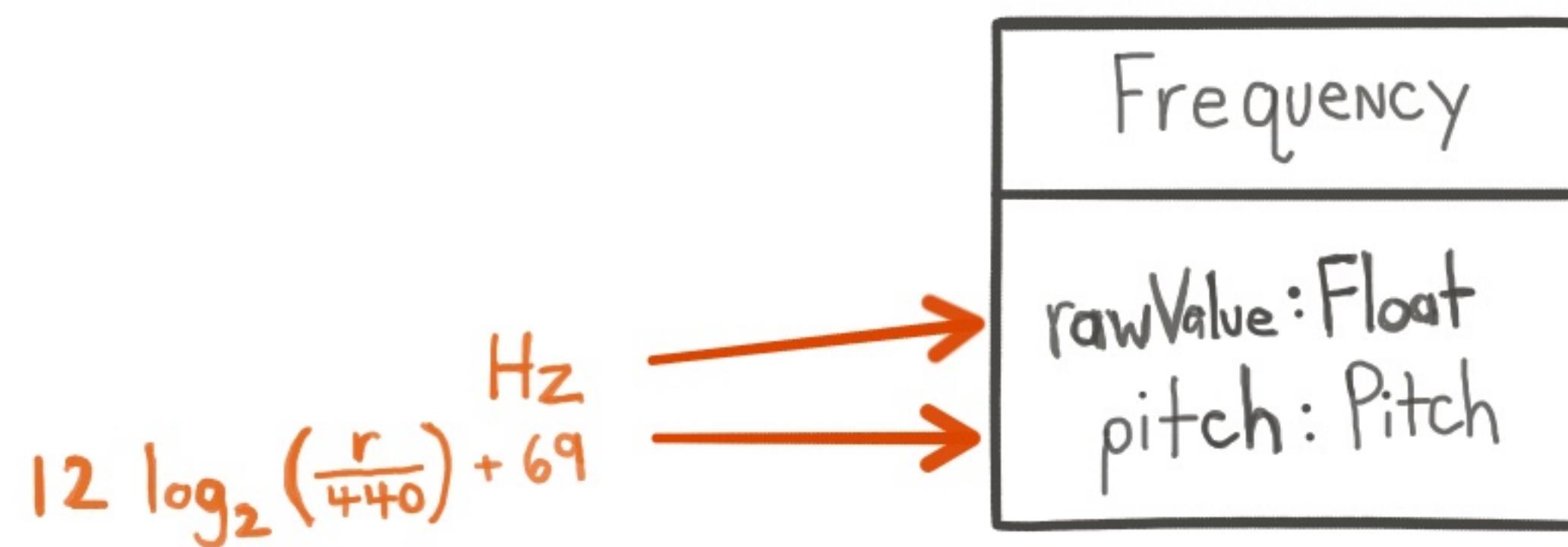
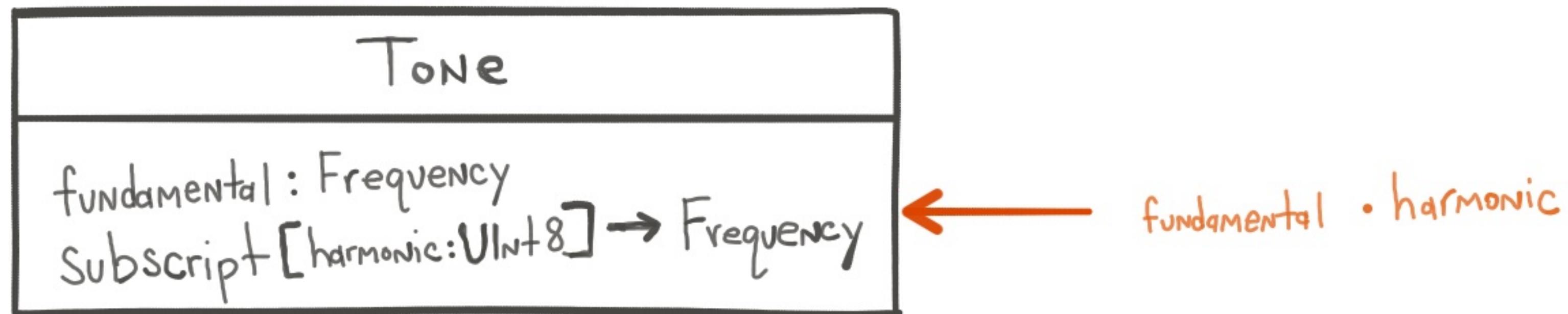


Frame (feature rate = 10 Hz)

Frame (feature rate = 2 Hz)



$$\begin{aligned}
0 - 127 &= \text{MIDI notes} \\
2^{\left(\frac{r-69}{12}\right)} \cdot 440 & \\
\frac{r}{12} - 1 & \\
r \bmod 12 &
\end{aligned}$$



$$12 \log_2 \left( \frac{r}{440} \right) + 69 \text{ Hz}$$

### Rational Transfer Function

$$Y(z) = \frac{b(1) + b(2)z^{-1} + \dots + b(n_b+1)z^{-n_b}}{1 + a(2)z^{-1} + \dots + a(n_a+1)z^{-n_a}} X(z) \quad \text{where}$$

$n_a$  = feedback filter order  
 $n_b$  = feedforward filter order

or

$$a(1)y(n) = b(1)x(n) + b(2)x(n-1) + \dots + b(n_b+1)x(n-n_b) - a(2)y(n-1) - \dots - a(n_a+1)y(n-n_a)$$

$$a_0 y_n = b_0 x_n + b_1 x_{n-1} + \dots + b_8 x_{n-8} - a_1 y_{n-1} - \dots - a_8 y_{n-8}$$

filterCoefficients [60]

	a
0	1.0
1	-7.3704400113886903
2	24.356331306218301
3	-47.053025057578601
4	58.0734443623397
5	-46.877802066188003
6	24.175265864191601
7	-7.2884051256710896
8	0.98518731451036901

	b
0	0.0031436437997640599
1	-0.023195340053373301
2	0.076751499507004203
3	-0.14850000730407001
4	0.18360216818219399
5	-0.14850000730407001
6	0.076751499507004398
7	-0.023195340053373401
8	0.0031436437997640699

	x
0	0.0
1	-3.05176e-05
2	-6.10352e-05
3	-6.10352e-05
4	-0.00012207
5	-9.15527e-05
6	-0.000152588
7	-0.000152588
8	-0.000152588
9	-0.000213623
10	-0.000213623
11	-0.000183105
12	-0.000244141
13	-0.000152588
14	-0.000244141
15	-0.000152588
16	-0.000213623
17	-0.000152588
18	-0.00012207
19	-0.00012207
20	3.05176e-05
21	-3.05176e-05
22	9.15527e-05
23	9.15527e-05
24	0.000152588
25	0.000244141
26	0.000305176
27	0.000366211
28	0.000396729
29	0.000396729
30	0.000457764
31	0.000457764
32	0.000427246
33	0.000488281

$$y_0 = 0$$

$$y_1 = -9.593646402e-08 + 0 - 0$$

## MatLab ChromaToolbox Routine

Demo  
convert audio file  
estimate tuning ← ignore for now  
STMSP window length = 4410  
convert audio to pitch using filter bank

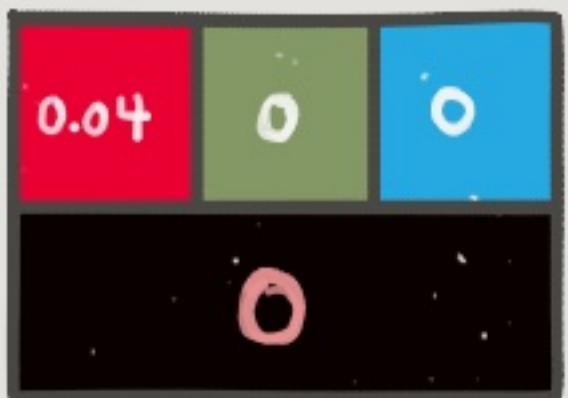
wav\_to\_audio  
downmix stereo to MONO  
resample to 22050 Hz with an antialiasing filter of order 8,820,000

### audio\_to\_pitch\_via\_FB

assume sample rate of 22050 Hz  
assume pitch range of [21:108]  
initialize 'fs\_pitch' array with 128 zeros  
initialize 'fs\_index' array with 128 zeros  
 $fs\_pitch[21:59] = 882$   
 $fs\_pitch[60:95] = 4410$   
 $fs\_pitch[96:120] = 22050$   
 $fs\_index[21:59] = 3$   
 $fs\_index[60:95] = 2$   
 $fs\_index[96:120] = 1$   
create copies of signal at 4410 Hz and 882 Hz  
calculate the total number of frames and their start/stop locations  
iterate through the range of pitch values (p) performing:  
1) get signal 'x' according to 'fs\_index'  
2) forward-back filter with corresponding a and b coefficients in filter bank  
3) square all values in the filtered signal storing as 'f\_square'  
4) iterate through range ' $w=1:winLenSTMSP$ ' performing:  
a) calculate 'factor' as input sample rate divide by the 'fs\_pitch' entry corresponding to the current pitch  
b) iterate through range ' $k=1:seg\_pcm\_num\{w\}$ ' performing:  
i) calculate start as  $\lceil \frac{seg\_pcm\_start\{w\}(k)}{22050} \cdot fs\_pitch(p) \rceil$   
ii) calculate stop as  $\lfloor \frac{seg\_pcm\_stop\{w\}(k)}{22050} \cdot fs\_pitch(p) \rfloor$   
iii) set 'f\_pitch\_energy\{w\}(p,k)' by calculating the sum of f\_square from start to stop and multiplying by 'factor'

index

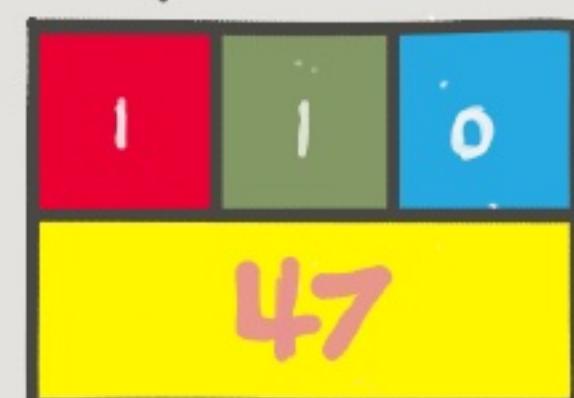
## Color Map



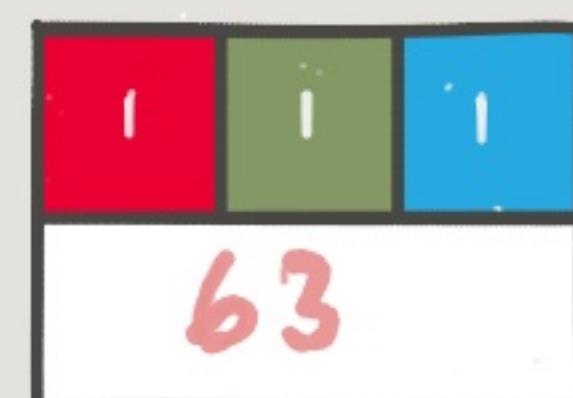
data = 0



data = 0.3594



data = 0.7344



data = 0.9844

Data	R	G	B
0.5911	1	0.5804	0
0.137	0.3725	0	0
0.3523	0.9569	0	0
0.175	0.498	0	0
0.2519	0.7059	0	0
0.01684	0.08235	0	0
0.05741	0.1647	0	0
0.6173	1	0.6667	0
0.05508	0.1647	0	0
0.03619	0.1216	0	0
0.07619	0.2078	0	0
0.1372	0.3725	0	0

```
func color(for value: Float64) -> UIColor {  
    let index = (value * 64).rounded(.down)  
    return colorMap[index]  
}
```

## Chroma Vector Quantization

Given a chroma vector  $v_n \in [0:1]^{12}$ , where  $n$  is the vector's frame, the quantization function  $\tau$  maps  $v_n$  such that  $v_n \in [0:4]^{12}$ .

let  $c_i = v_n[i]$  for  $i \in [0:11]$

$$\tau(c_i) = \begin{cases} 0 & \text{for } c_i < 0.05 \\ 1 & \text{for } c_i < 0.1 \\ 2 & \text{for } c_i < 0.2 \\ 3 & \text{for } c_i < 0.4 \\ 4 & \text{for } c_i \leq 1 \end{cases}$$

# pitch\_to\_CENS.m

seg\_num = 321

f\_chroma\_energy = chroma band energies for each frame

$$f_{\text{chroma\_energy\_distr}} = \sum_{c \in [1:12], k \in [1:\text{seg\_num}]} \left\{ \begin{array}{ll} \frac{f_{\text{chroma\_energy}_{ck}}}{\sum_{c=1}^{12} f_{\text{chroma\_energy}_{ck}}} & \text{for } \sum_{c=1}^{12} f_{\text{chroma\_energy}_{ck}} > \text{NormThresh} \\ 0 & \text{for } \sum_{c=1}^{12} f_{\text{chroma\_energy}_{ck}} < \text{NormThresh} \end{array} \right.$$

$$f_{\text{stat\_help}} = \sum_{c \in [1:12], k \in [1:\text{seg\_num}]} \left( \sum_{i=1}^4 \left\{ \begin{array}{ll} \text{quantWeights}_i & \text{for } f_{\text{chroma\_energy\_dist}_{ck}} > \text{quantSteps}_i \\ 0 & \text{for } f_{\text{chroma\_energy\_dist}_{ck}} < \text{quantSteps}_i \end{array} \right\} \right)$$

## Variables

midiMin = 21

midiMax = 108

winLengthSmooth = 21

downsampSmooth = 5

inputFeatureRate = 10

quantSteps = [0.4, 0.2, 0.1, 0.05]

quantWeights = [0.25, 0.25, 0.25, 0.25]

NormThresh = 0.001

f\_feature

smoothDownsampleFeature.m

stat\_window = normalized Hanning window of length winLengthSmooth

f\_feature\_stat = f\_feature applying stat\_window as a FIR filter with decimation factor downsampSmooth

$$\text{stat\_num} = \left\lceil \frac{\text{seg\_num}}{\text{downsampSmooth}} \right\rceil$$

$$\text{cut} = \left\lfloor \frac{\text{winLengthSmooth}^{-1}}{2 \cdot \text{downsampSmooth}} \right\rfloor$$

$$f_{\text{feature\_stat}} = f_{\text{feature}}[1:12][\text{cut}:\text{stat\_num}+\text{cut}]$$

$$\text{new\_feature\_rate} = \frac{\text{inputFeatureRate}}{\text{downsampleSmooth}}$$

f\_feature

normalizeFeature.m

$$f_{\text{featureNorm}} = \sum_{c,k} \left\{ \begin{array}{ll} \frac{f_{\text{feature}_{ck}}}{\|f_{\text{feature}}\|_k} & \text{for } \|f_{\text{feature}}\|_k > \text{NormThresh} \\ 0.2886751346 & \text{for } \|f_{\text{feature}}\|_k > \text{NormThresh} \\ \text{for } c \in [1:12], k \in [1:\text{stat\_num}] \end{array} \right.$$

f\_CENS

# Discrete Cosine Transform

Type 2 
$$\sum_{k=0}^{\ell} \left( \sum_{j=0}^{\ell} A_j \cos\left(\frac{k(j+0.5)\pi}{\ell}\right) \right)$$

Type 3 
$$\sum_{k=0}^{\ell} \left( \frac{A_0}{2} + \sum_{j=1}^{\ell} A_j \cos\left(\frac{(k+0.5)j\pi}{\ell}\right) \right)$$

Type 4 
$$\sum_{k=0}^{\ell} \left( \sum_{j=0}^{\ell} A_j \cos\left(\frac{k(j+0.5)\pi}{\ell}\right) \right)$$

## Elliptical Order

complete elliptic integral of the first kind  $K$ :

$$K(k) = \int_0^{\frac{\pi}{2}} \frac{d\theta}{\sqrt{1-k^2 \sin^2 \theta}} = \int_0^1 \frac{dt}{\sqrt{(1-t^2)(1-k^2 t^2)}}$$

as a power series:

$$K(k) = \frac{\pi}{2} \sum_{n=0}^{\infty} \left( \frac{(2n)!}{2^{2n} (n!)^2} \right)^2 k^{2n} = \frac{\pi}{2} \sum_{n=0}^{\infty} (P_{2n}(0))^2 k^{2n}$$

where  $P_n$  is the Legendre polynomials, which is equivalent to:

$$K(k) = \frac{\pi}{2} \left( 1 + \left(\frac{1}{2}\right)^2 k^2 + \left(\frac{1 \cdot 3}{2 \cdot 4}\right)^2 k^4 + \dots + \left(\frac{(2n-1)!!}{(2n)!!}\right)^2 k^{2n} + \dots \right)$$

computed efficiently in terms of the arithmetic-geometric mean:

$$K(k) = \frac{\pi}{2} \text{agm}(1, \sqrt{1-k^2})$$

arithmetic-geometric mean of two positive real numbers  $x$  and  $y$ :

compute the arithmetic ( $a_1$ ) and geometric ( $g_1$ ) means of  $x$  and  $y$ ,

$$a_1 = \frac{1}{2}(x+y)$$

$$g_1 = \sqrt{xy}$$

then compute with  $a_1$  substituted for  $x$  and  $g_1$  substituted for  $y$ . By repeating this process, two sequences  $a_n$  and  $g_n$  are defined:

$$a_{n+1} = \frac{1}{2}(a_n + g_n)$$

$$g_{n+1} = \sqrt{a_n g_n}$$

these two sequences converge to the same number, the arithmetic-geometric mean or  $\text{agm}(x, y)$

calculating elliptic filter order for given values  $\epsilon, A, \Omega_p, \Omega_s$ :

$$n = \frac{K(k) K(\sqrt{1-k^2})}{K(k_p) K(\sqrt{1-k_p^2})}$$

where  $K(\cdot)$  is the complete elliptic integral of the first kind

From Theory and Application of Digital Signal Processing by Rabiner and Gold

passband and stopband edge frequencies:  $\Omega_p, \Omega_s$

passband and stopband gains:  $G_p, G_s$   $G = \frac{1}{\sqrt{1+\epsilon^2}} = 10^{-\frac{A}{20}}$

ripple parameters:  $\epsilon_p, \epsilon_s$   $\epsilon = \sqrt{G^{-2}-1} = \sqrt{10^{-\frac{2A}{20}} - 1}$

attenuation (dB):  $A_p, A_s$   $A = -20 \log_{10} G = 10 \log_{10}(1 + \epsilon^2)$

selectivity parameter:  $k$   $k = \frac{\Omega_p}{\Omega_s}$

discrimination parameter:  $k_1$   $k_1 = \frac{\epsilon_p}{\epsilon_s}$

### Landen Transformations

$$k_n = \left( \frac{k_{n-1}}{1 + k'_{n-1}} \right)^2, \quad n = 1, 2, \dots, M$$

complete elliptic integral of the first kind:  $K = \int_0^{\frac{\pi}{2}} \frac{d\theta}{\sqrt{1 - k^2 \sin^2 \theta}}$

degree equation:  $F_N(k') = k_1^{-1}$  or  $N \frac{K'}{K} = \frac{K_1'}{K_1}$

$F_N(w) = cd(NwK_1, k_1)$ ,  $w = cd(vK, k)$

$cd(x, k)$  denotes the Jacobian elliptic function with modulus  $k$  and real quarter-period  $K$

magnitude response  $|H(\Omega)|^2 = \frac{1}{1 + \epsilon_p^2 F_N^2(\frac{\Omega}{\Omega_p})} = \frac{1}{1 + \frac{\epsilon_s^2}{F_N^2(\frac{\Omega_s}{\Omega})}}$

$$k' = \sqrt{1 - k^2}$$

$$K' = K \text{ with } k = k'$$

$$\Omega = \frac{2\pi f}{f_s}$$

# Chord Recognition

## Template - Based

Map the 12 chromas to a bit field

Measuring similarity using the inner product of normalized vectors:

$$s(x, y) = \frac{\langle x | y \rangle}{\|x\| \cdot \|y\|}$$

## Evaluating Results

TP = True Positive, Correctly recognized chord

FP = False Positive, Incorrectly recognized chord

FN = False Negative, Failure to recognize a chord

P = Precision      R = Recall      F = F-Measure

$$P = \frac{\# TP}{\# TP + \# FP} \quad R = \frac{\# TP}{\# TP + \# FN} \quad F = \frac{2 \cdot P \cdot R}{P + R}$$

## Incorporating Harmonics

Along with each note, consider the first 8 partials.

For example the first 8 partials for C would be C, C, G, C, E, G, B<sub>b</sub>, C.

Define constant  $\alpha$  and assume energy of the  $k^{\text{th}}$  partial is  $\alpha^{k-1}$ .  
A chroma vector for just a C note would then be the following:

$$t_C^h = (1 + \alpha + \alpha^3 + \alpha^7, 0, 0, 0, \alpha^4, 0, 0, \alpha^2 + \alpha^5, 0, 0, \alpha^6, 0)$$

C            D<sub>b</sub>    D    E<sub>b</sub>    E    F    G<sub>b</sub>    G    A<sub>b</sub>    A    B<sub>b</sub>    B

Add the vectors for each note to get the template for the composed chord.

$$t_{\text{cmaj}}^h = t_C^h + t_E^h + t_G^h$$

let  $\alpha = 0.9$

$$t_C^h = (3.1072969, 0, 0, 0, 0.6561, 0, 0, 1.40049, 0, 0, 0.531441, 0)$$

partials for E would be E E B E A B D E  
and for G: G G D G B D F G

## Bit Representation of Chord Patterns

12 intervals (interval 1 is implied)  
 Each interval has 4 possible states requiring 2 bits  
 possible states:

00 interval is not included

10 interval is included

01 flattened interval is included

11 sharpened interval is included

uses the 24 least significant bits in a UInt32 value

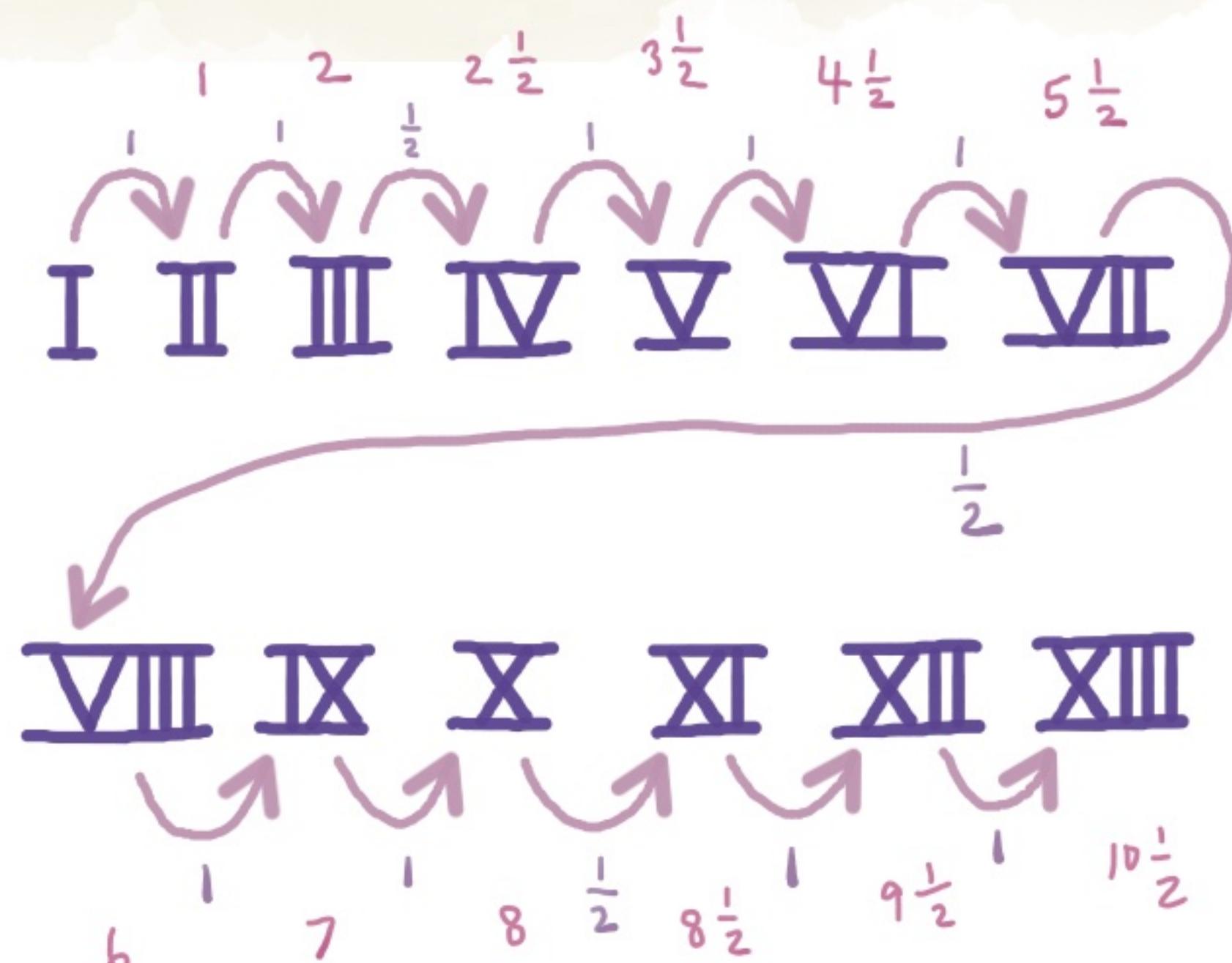
example 1: 7sus4 = (4, 5, b7) = 0000000000000000000010010100000

convert to chroma for a root of C: [C, F, G, Bb]

C	Db	D	Eb	E	F	Gb	G	Ab	A	Bb	B
I	II	III	IV	V	VI	VII					
C		F		G		Bb					

and for a root of Db: [Db, Gb, Ab, B]

C	Db	D	Eb	E	F	Gb	G	Ab	A	Bb	B
VII	I	II	III	IV	V	VI	VII				
Db			Gb		Ab		B				



example 2: Maj13#11 = (3, 5, 7, 9, #11, 13) = 00000000100011001000100010001000

convert to chroma for a root of C: [C, E, G, B, D, Gb, A]

C	Db	D	Eb	E	F	Gb	G	Ab	A	Bb	B
I	II	III	IV	V	VI	VII	VIII	IX	X	XI	XII
VIII	IX	X	XI	XII	XIII	XIV	XV	XVI	XVII	XVIII	XIX
C	E	G	Ab	A	Bb	B	D	Eb	F	Gb	

and for a root of Db: [Db, F, Ab, C, Eb, G, Bb]

C	Db	D	Eb	E	F	Gb	G	Ab	A	Bb	B
VII	I	II	III	IV	V	VI	VII	VIII	IX	X	XI
VII	Db	I	VIII	IX	X	XI	XII	XIII	XIV	XV	XVI

## Chroma Template Vectors

C vector = ( 3.107297 0 0 0 0.6560999 0 0 1.40049 0 0 0.5314409 0 )

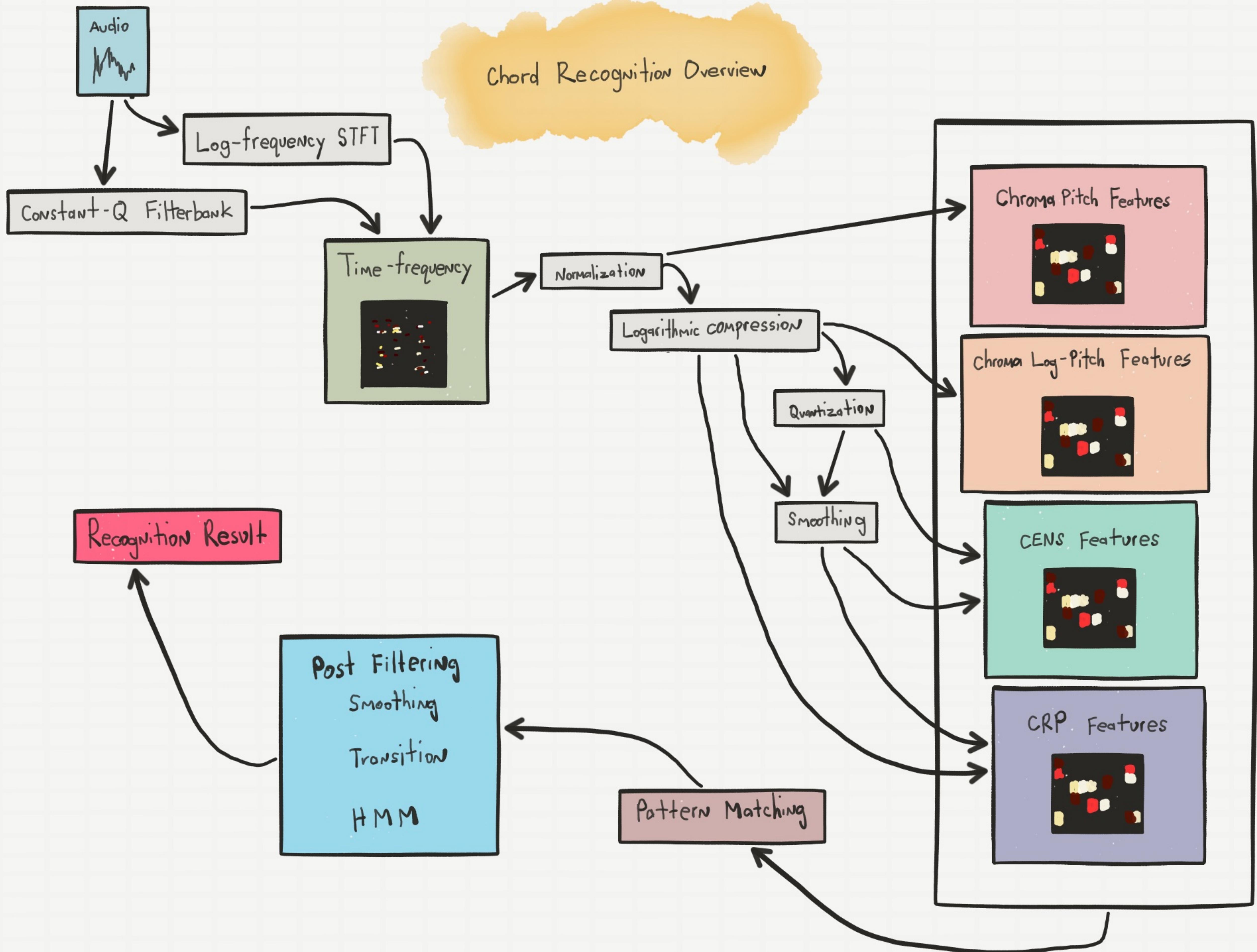
E vector = ( 0 0 0.5314409 0 3.107297 0 0 0 0.6560999 0 0 0 1.40049 )

G vector = ( 0 0 1.40049 0 0 0.5314409 0 3.107297 0 0 0 0 0.6560999 )

---

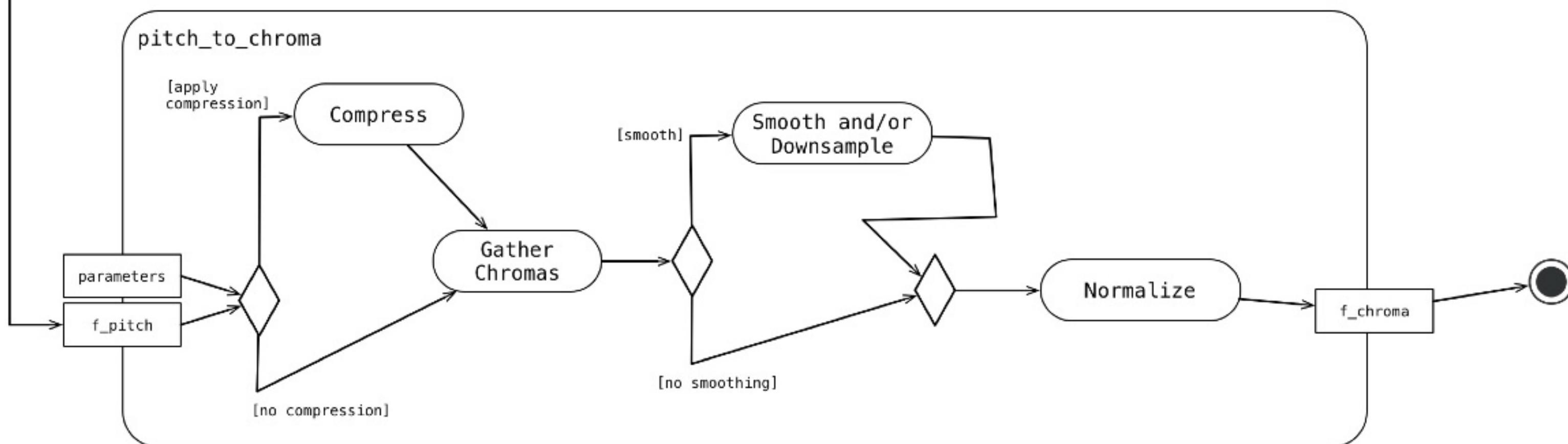
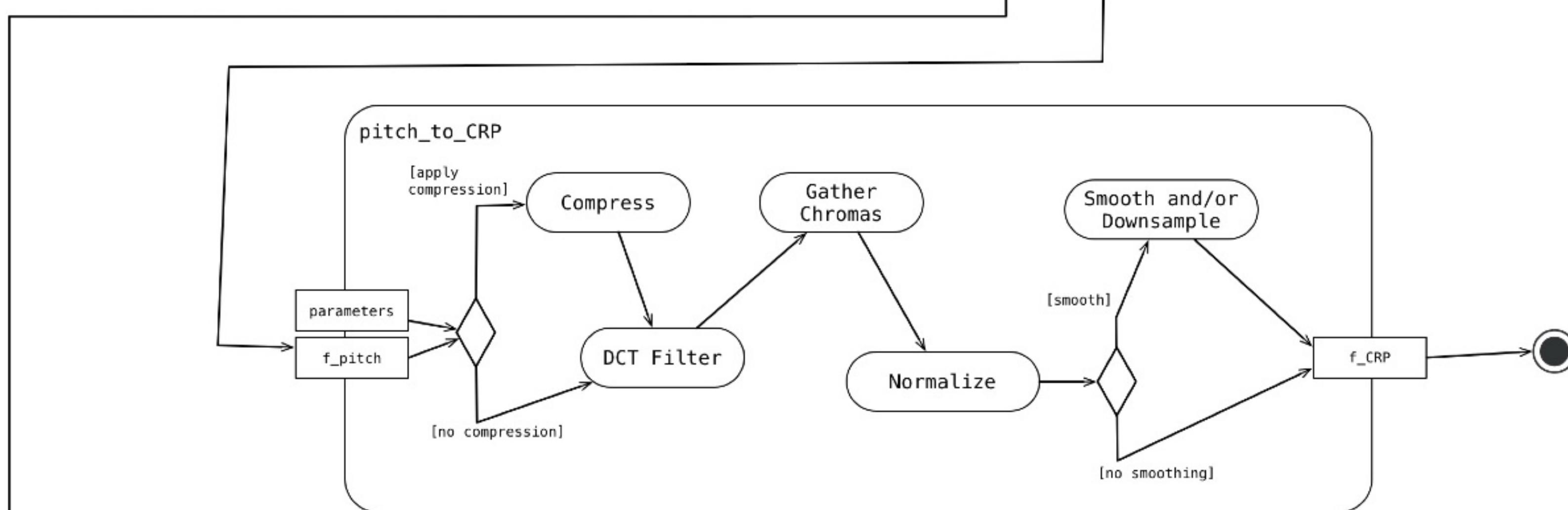
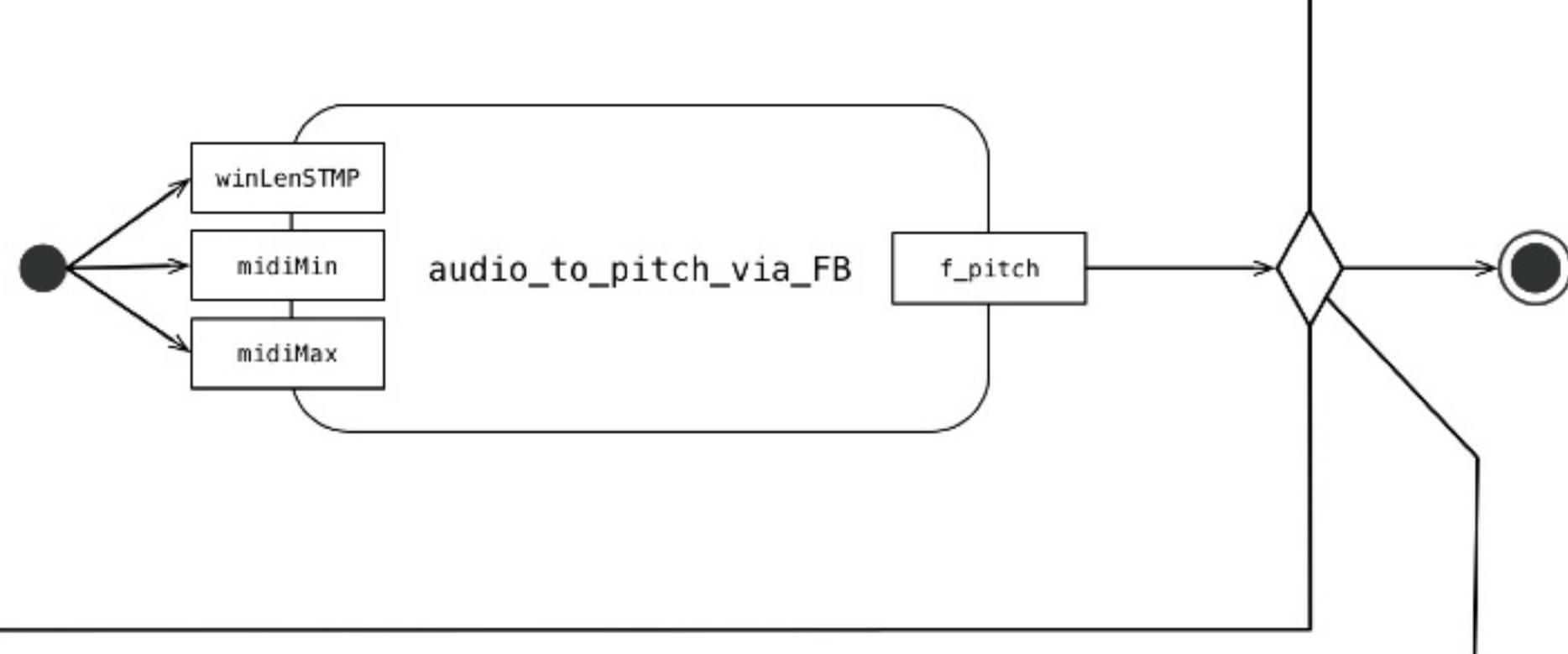
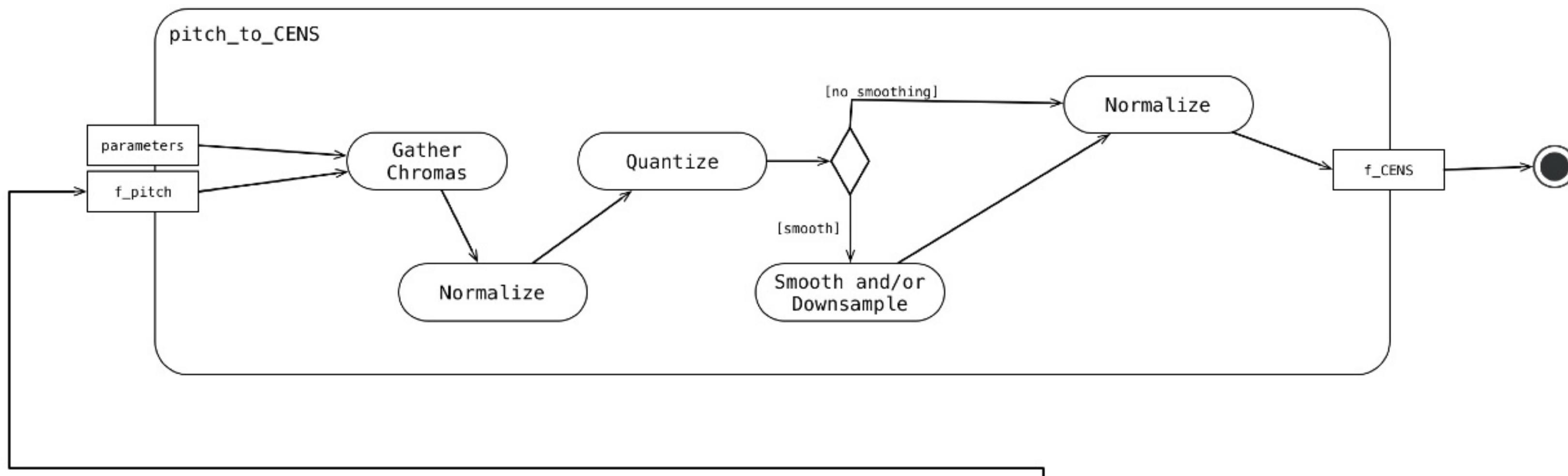
CEG vector = ( 3.107297 0 1.9319309 0 3.7633969 0.5314409 0 4.507787 0.6560999 0 0.5314409 2.0565899 )

## Chord Recognition Overview



## ChromaToolbox Demo Parameters

Pitch Features	window length	= 4410
	hop size	= 2205
	MIDI Range	= [21:108]
Chroma Pitch	compress	= false
	normalize	= true
	norm p	= 2
	norm threshold	= 0.001
	MIDI Range	= [1:120]
Chroma Log-Pitch	compress	= true
	compression factor	= 100
	normalize	= true
	norm p	= 2
	norm threshold	= 0.001
	MIDI Range	= [1:120]
CENS	quantization steps	= [0.4, 0.2, 0.1, 0.05]
	quantization weights	= [0.25, 0.25, 0.25, 0.25]
	norm threshold	= 0.001
	smoothing window length	= 21
	smoothing downsample	= 5
	MIDI Range	= [21:108]
CRP1	kept coefficients	= [55:120]
	compress	= true
	compression factor	= 1000
	norm p	= 2
	norm threshold	= 0.000001
	smoothing window length	= 1
	smoothing downsample	= 1
CRP2	kept coefficients	= [55:120]
	compress	= true
	compression factor	= 1000
	norm p	= 2
	norm threshold	= 0.000001
	smoothing window length	= 21
	smoothing downsample	= 5



## Similarity Scores

Calculate the inner product

otherVector

0.014	0.000	0.027	<b>0.376</b>	0.013	0.011	0.000	<b>0.200</b>	0.099	0.013	0.000	<b>0.248</b>
-------	-------	-------	--------------	-------	-------	-------	--------------	-------	-------	-------	--------------

vector

0.000	0.000	0.120	<b>0.213</b>	0.000	0.082	0.000	<b>0.251</b>	0.000	0.000	<b>0.295</b>	0.038
-------	-------	-------	--------------	-------	-------	-------	--------------	-------	-------	--------------	-------

products

0.000	0.000	0.00324	0.08088	0.000	0.000902	0.000	0.0502	0.000	0.000	0.000	0.009424
-------	-------	---------	---------	-------	----------	-------	--------	-------	-------	-------	----------

$$\sum_{i=0}^n \text{products}_i = 0.143854 \\ + 0.2 \times 0.85 \\ + 0.2 \times 3 \\ + 0.2 \text{ high root energy bonus} \\ - 0.1 \\ \hline 0.813854$$

## Similarity Scores

Chroma with energy  $> 0.35$

vectors for chords played without the 5<sup>th</sup>  
chords with identical templates

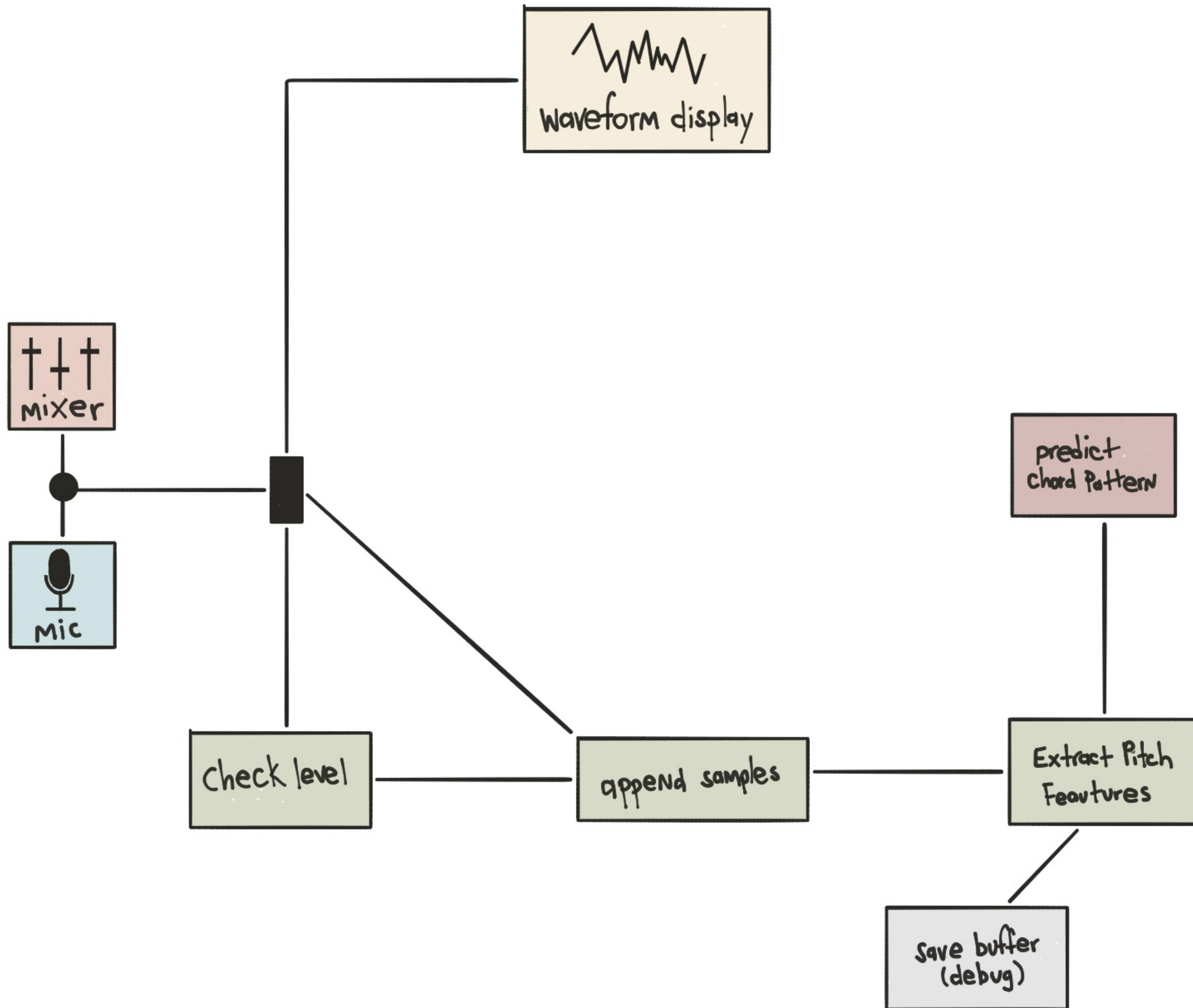
for each chord, which intervals are optional?

vectors with 85% of energy in two chromas

document guess of how many notes are present

missing energy for templates chord root

templates where there is energy for all chord chromas.



## Signal Processing for Waterfall View

### Signal Processing. FFT

$$N = 4800$$

$$\log_2 N = 12$$

$$K = 2400$$

$$\text{bin count} = 2401$$

perform vDSP\_zrip

divide all bins by 2

Move the Nyquist value

take the square roots of the squared magnitudes

### doFFT\_onAudioBuffer

$$N = 2048$$

$$\log_2 N = 11$$

Multiply with a blackman window of length N (only first  $\frac{N}{2}$ ?)

perform vDsp\_fft\_zip

calculate squared magnitudes

divide the squared magnitudes by 2

### Signal Processing. STFT

$$N = \text{'WindowSize'} (\text{default} = 8820)$$

Multiply by a normalized Hann window of size N

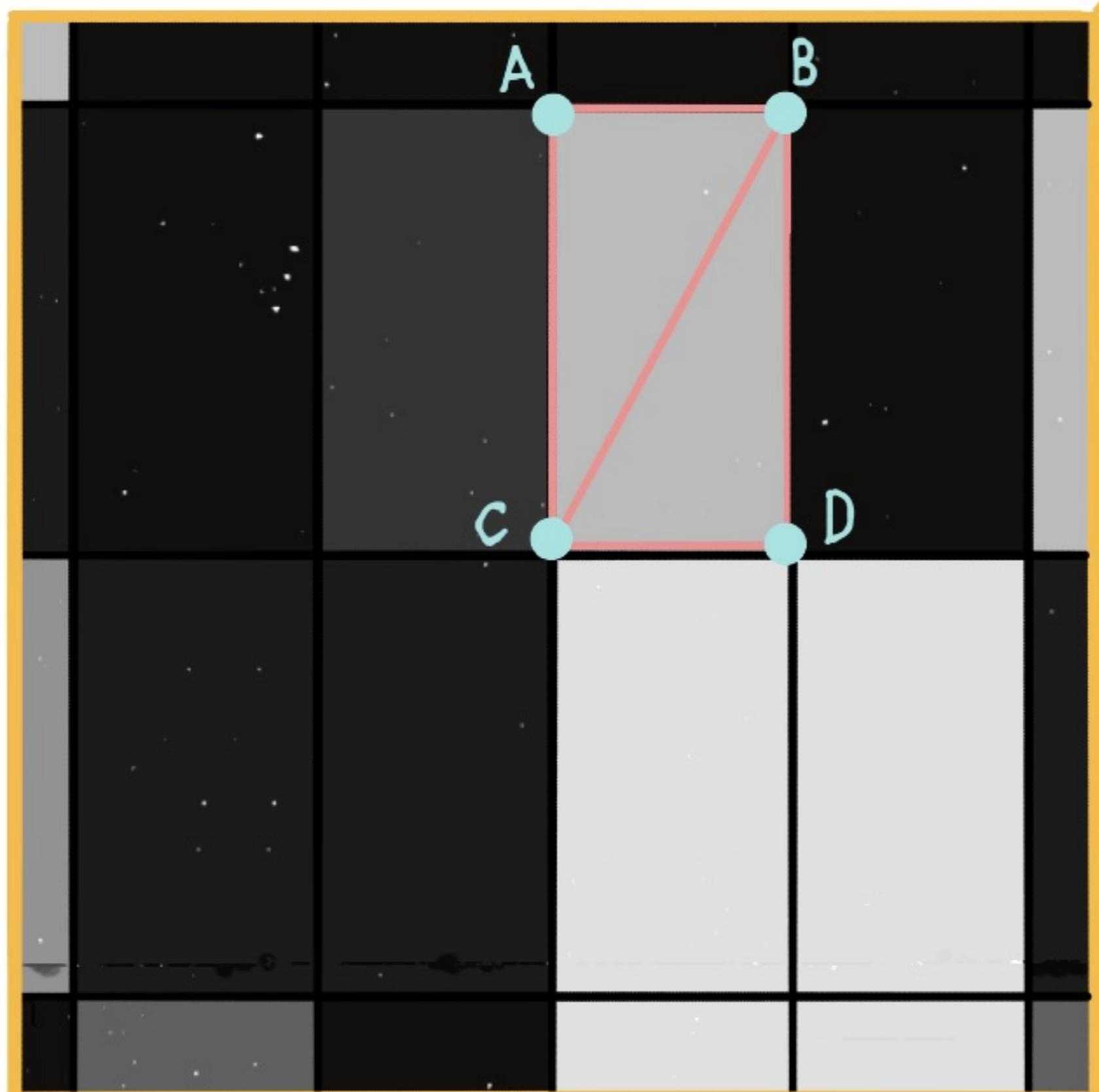
### Variables

spectrumCount the total # rows = 12

binCount the total # columns = 32

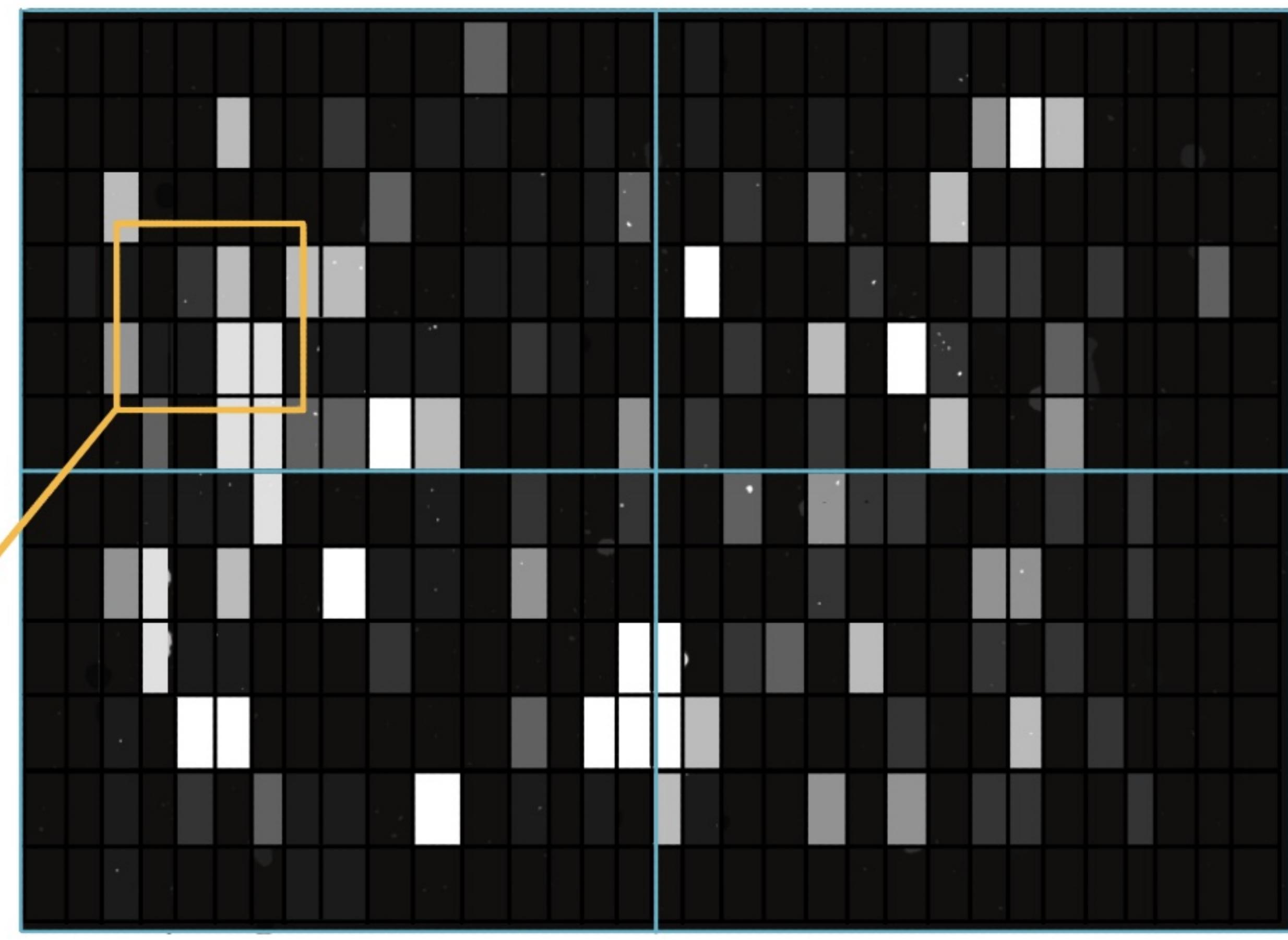
$$h = \frac{2}{\text{spectrumCount}}$$

$$w = \frac{2}{\text{binCount}}$$



(-1,1,1)

(1,1,1)



(-1,-1,1)

(1,-1,1)

$$b = \text{bin index} \quad s = \text{spectrum index} \quad s' = \text{spectrum count} - s$$

$$A = (bw-1, sh-1)$$

$$B = (A_x + w, A_y)$$

$$C = (A_x, A_y - h)$$

$$D = (B_x, C_y)$$

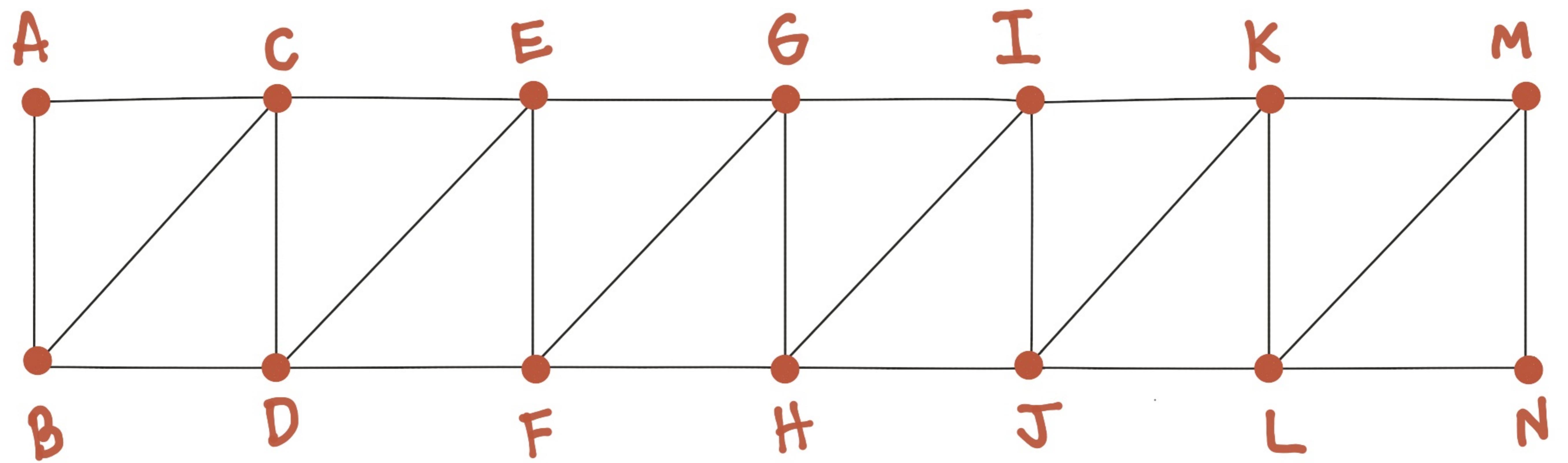
$$b=5 \quad s=3 \quad s'=9 \quad w=\frac{1}{16} \quad h=\frac{1}{6}$$

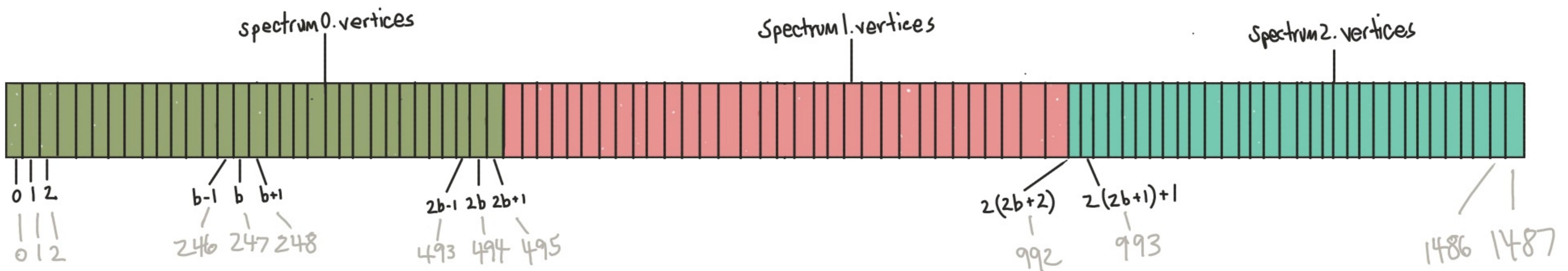
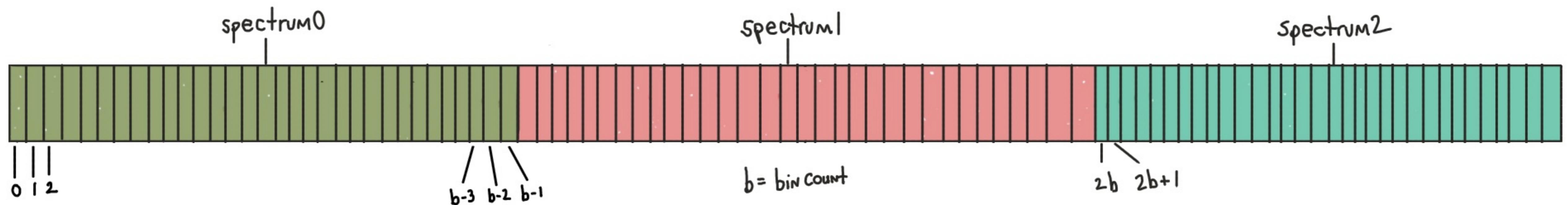
$$A = \left( -\frac{11}{16}, \frac{1}{2} \right)$$

$$B = \left( -\frac{10}{16}, \frac{1}{2} \right)$$

$$C = \left( -\frac{11}{16}, \frac{1}{3} \right)$$

$$D = \left( -\frac{10}{16}, \frac{1}{3} \right)$$





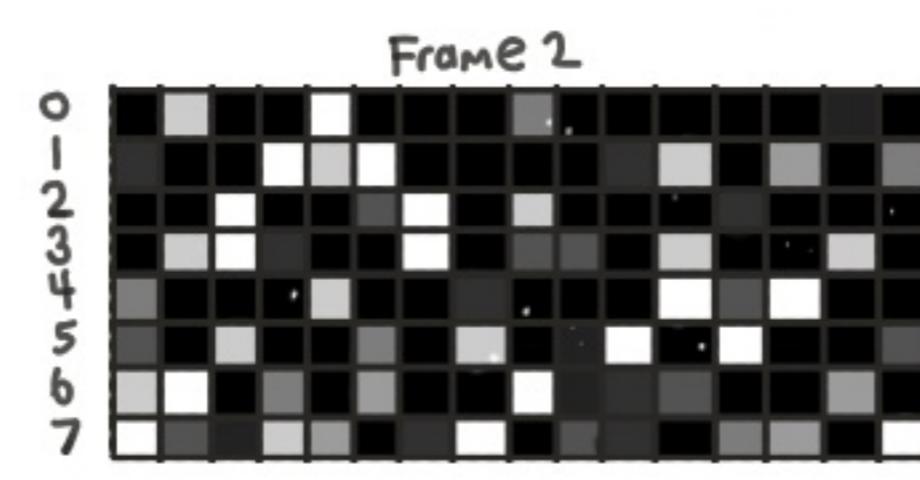
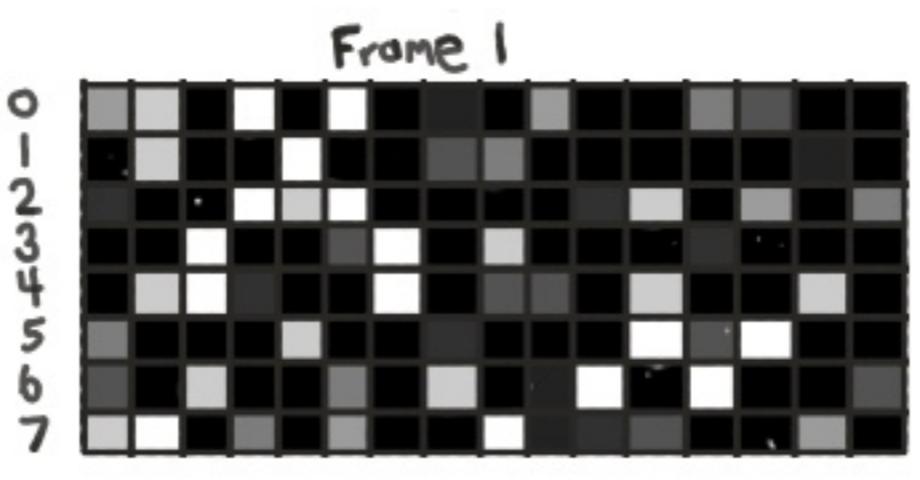
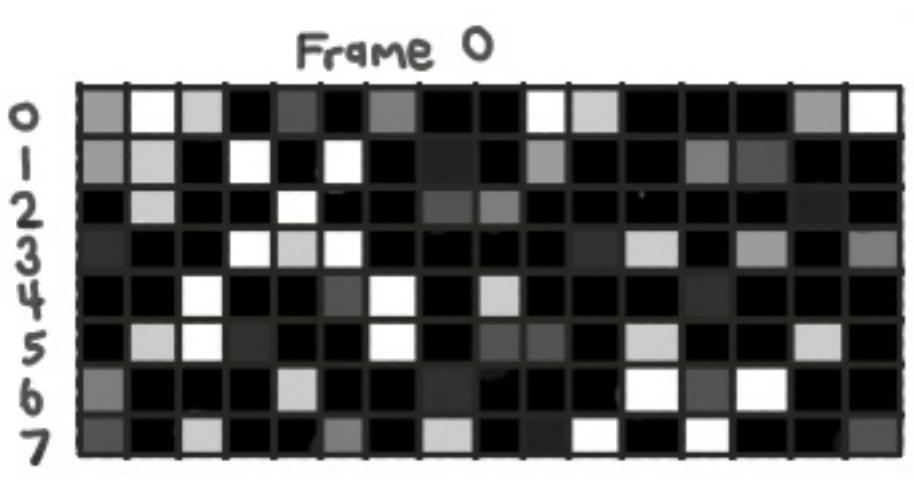
$$v = \text{vertices per Spectrum} = 2b + 2$$

$$\text{let } b = 247, v = 496$$

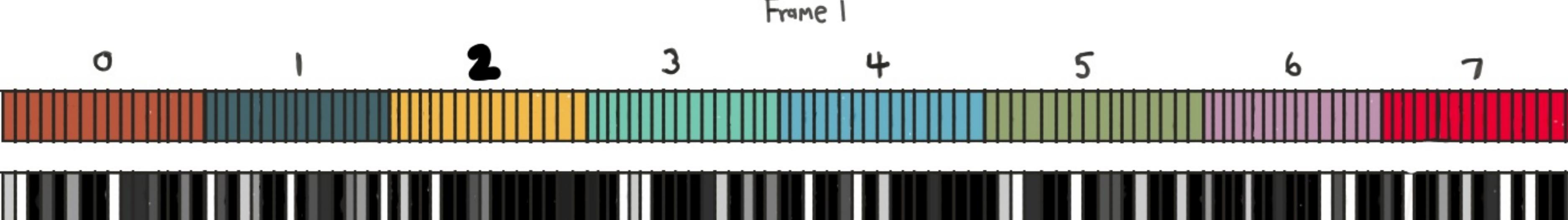
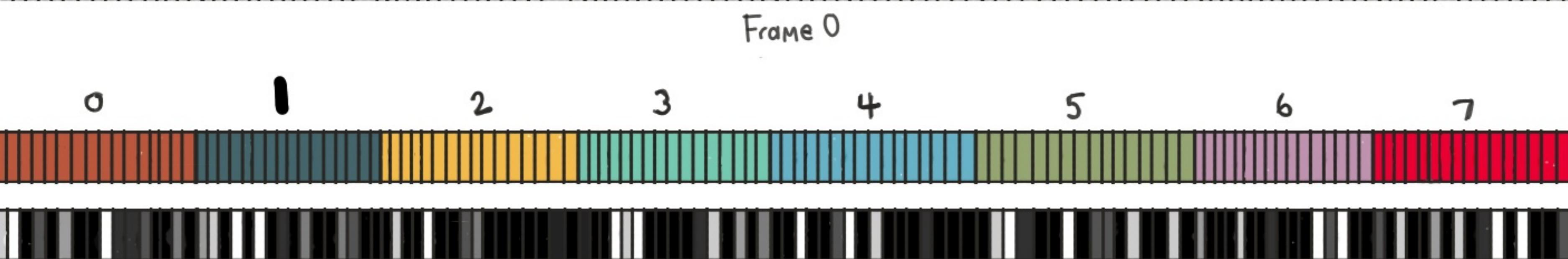
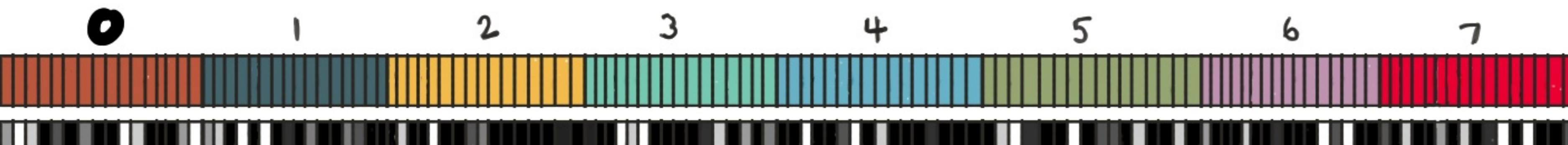
$$\text{bin}(\text{vid}) = \min\left(\frac{\text{vid} \bmod v}{2}, b-1\right)$$

$$\text{spectrum}(\text{vid}) = \frac{\text{vid}}{v}$$

vid	bin
$32b$	163
492	246
494	246
496	0



**bold** = spectrumIndex

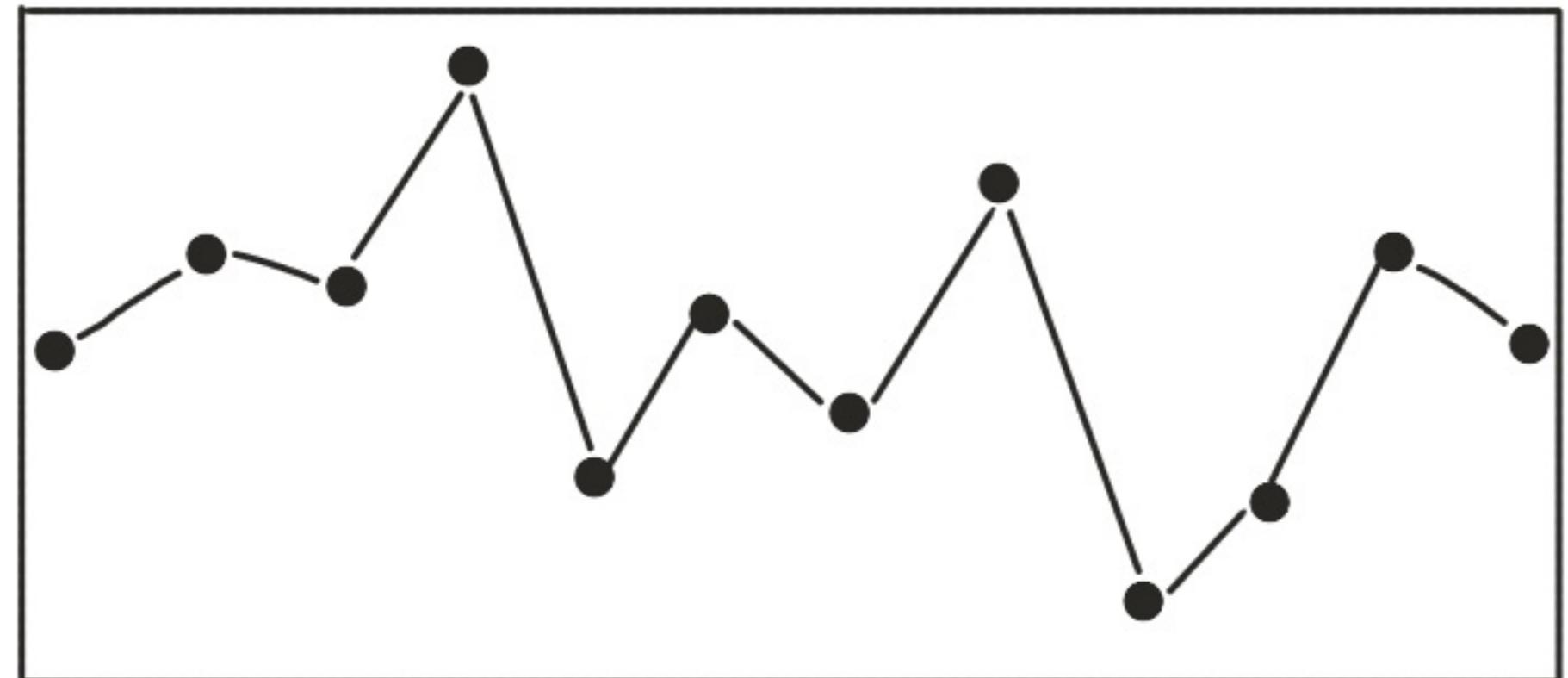


```

float4 vertex(vid) {
    sid = vid / spectrumCount
    binOffset = (vid % binCount) / 2
    bid = min(binOffset, binCount - 1)
    isTop = vid % 2 == 0
    w = 2 / binCount
    h = 2 / spectrumCount
    Δ = sid - spectrumIndex % spectrumCount
    row = (spectrumCount + Δ) % spectrumCount
    y = 1 - h * row - (isTop ? 0 : h)
    x = -1 + w * binOffset
    return float4(x, y, 0, 1)
}

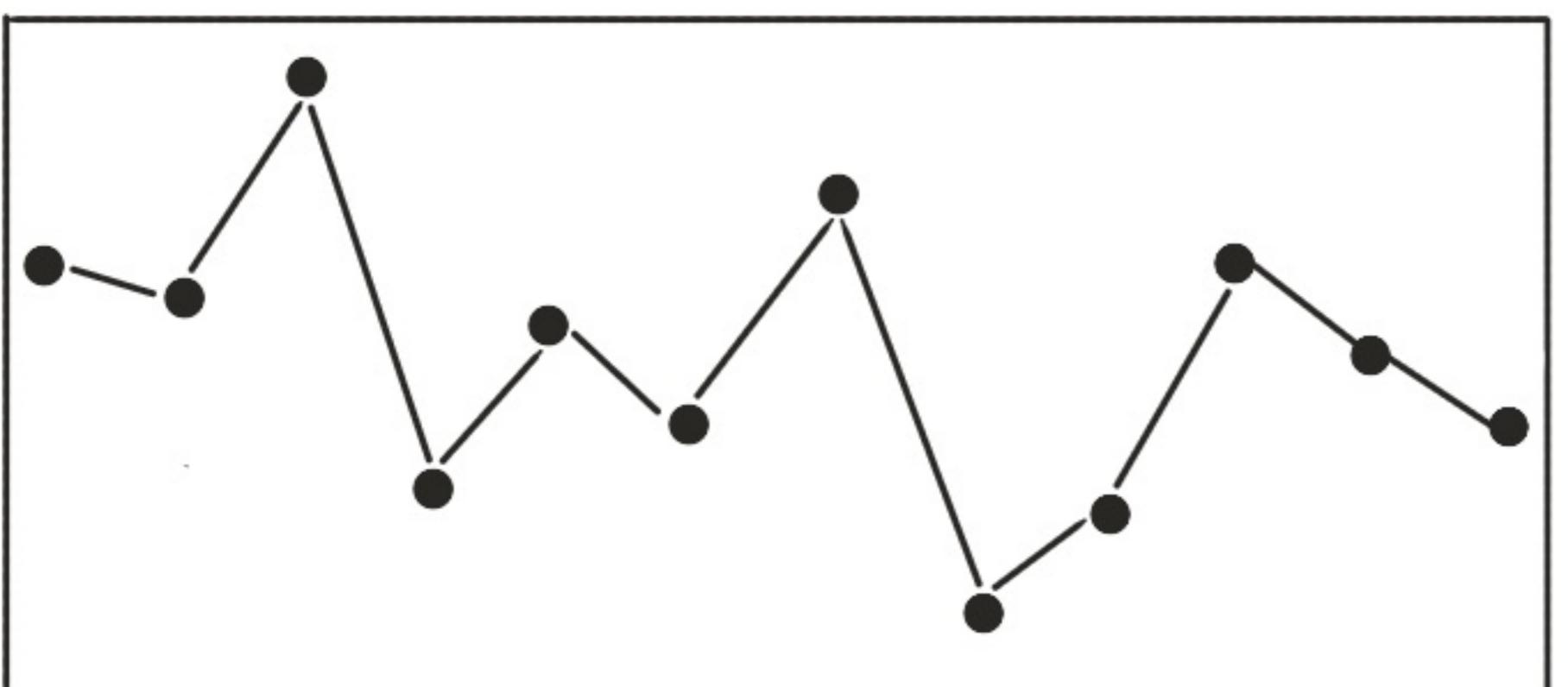
```

$t = \text{top} = \text{index of the oldest sample}$



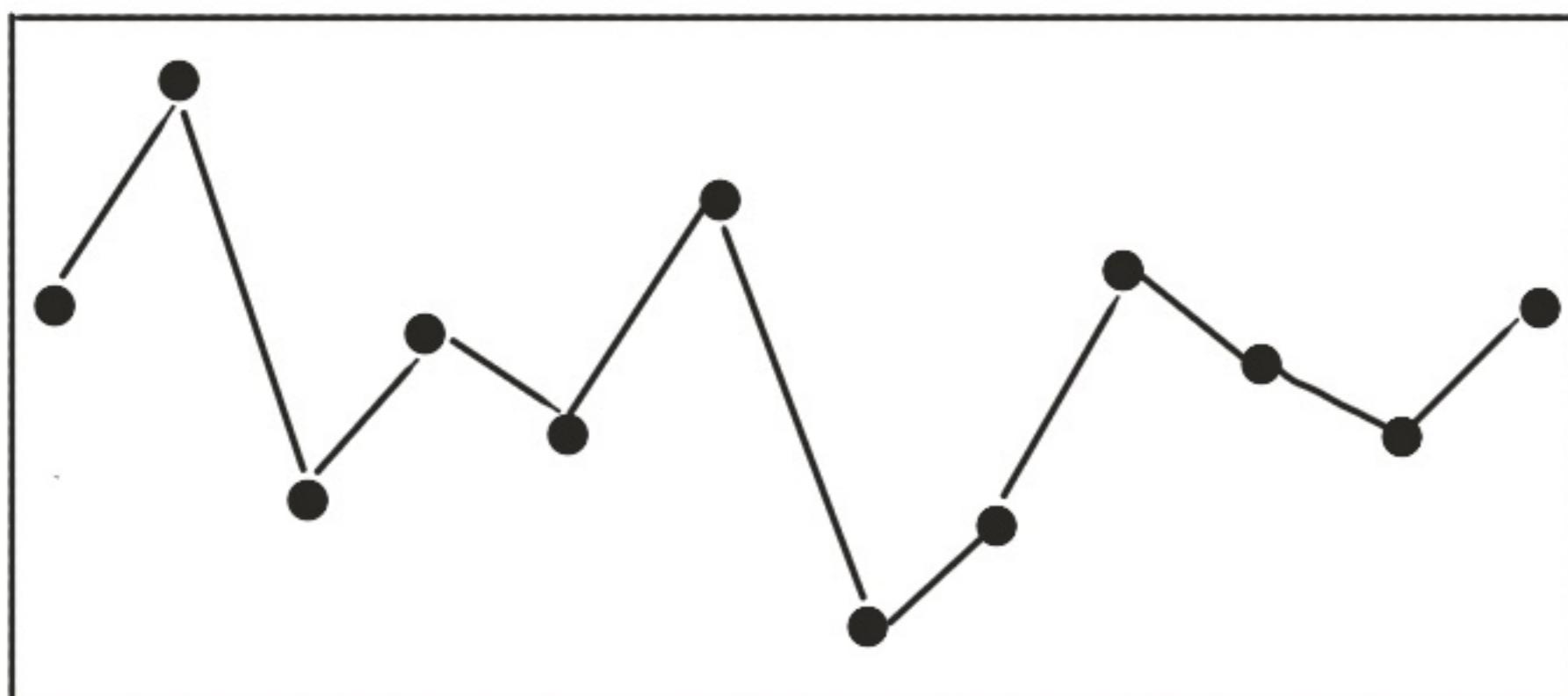
$t$

0.0	0.3	0.2	0.9	-0.4	0.1	-0.2	0.4	-0.8	-0.5	0.3	0.0
-----	-----	-----	-----	------	-----	------	-----	------	------	-----	-----



$t$

-0.2	0.3	0.2	0.9	-0.4	0.1	-0.2	0.4	-0.8	-0.5	0.3	0.0
------	-----	-----	-----	------	-----	------	-----	------	------	-----	-----



$t$

-0.2	0.1	0.2	0.9	-0.4	0.1	-0.2	0.4	-0.8	-0.5	0.3	0.0
------	-----	-----	-----	------	-----	------	-----	------	------	-----	-----

vertex VertexOut func waveformVertex ( uint vid [[ vertex.id ]],  
WaveformVertexArguments &arguments [[ buffer(0) ]])

{

    uint top = arguments.top;

    uint sampleCount = arguments.sampleCount;

    device float \*samples = arguments.samples;

    VertexOut outVertex;

    float y = samples[vid];

    float xOffset = 2.0f / float(sampleCount);

    uint index = vid > top ? vid - top : (vid + sampleCount) - top;

    float x = -1.0f + float(index) \* xOffset;

    outVertex.position = float4(x, y, 0.0f, 1.0f);

    outVertex.color = float4(1.0f, 1.0f, 0.0f, 1.0f);

    return outVertex;

}

Spectrum Buffer

0.05	0.01	0.002	0.8
0.001	0.4	0.16	0.6
0.036	0.0006	0.332	0.91
0.016	0.003	0.23	0.004
0.05	0.01	0.002	0.8
0.001	0.4	0.16	0.6
0.9	0.0012	0.31	0.02
0.06	0.003	0.23	0.004
0.05	0.01	0.002	0.8
0.44	0.09	0.061	0.2
0.036	0.0006	0.332	0.91
0.06	0.003	0.23	0.004

Incoming Audio Samples

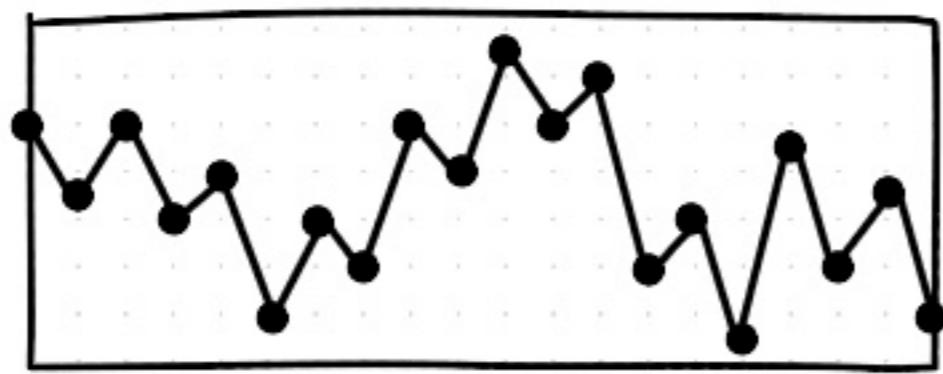
0.03	-0.01	0.2	0.1
------	-------	-----	-----

Acquire enough new samples to fill the FFT window and calculate the pitch spectrum

Pitch Spectrums

0.05	0.1	0.002	0.8
0.001	0.4	0.16	0.6
0.9	0.0012	0.31	0.02
0.06	0.003	0.23	0.004

Copy the New collection of spectrums into the Next frame's buffer



$$\text{displayCount} = \text{samplesConsumed} \times \text{framesPerSample} = 4410$$

$$\text{samplesConsumed} = \frac{\text{sampleRate}}{\text{framesPerSecond}} = 735$$

$$\text{sampleRate} = 44,100$$

$$\text{framesPerSecond} = 60$$

$$\text{framesPerSample} = 6$$

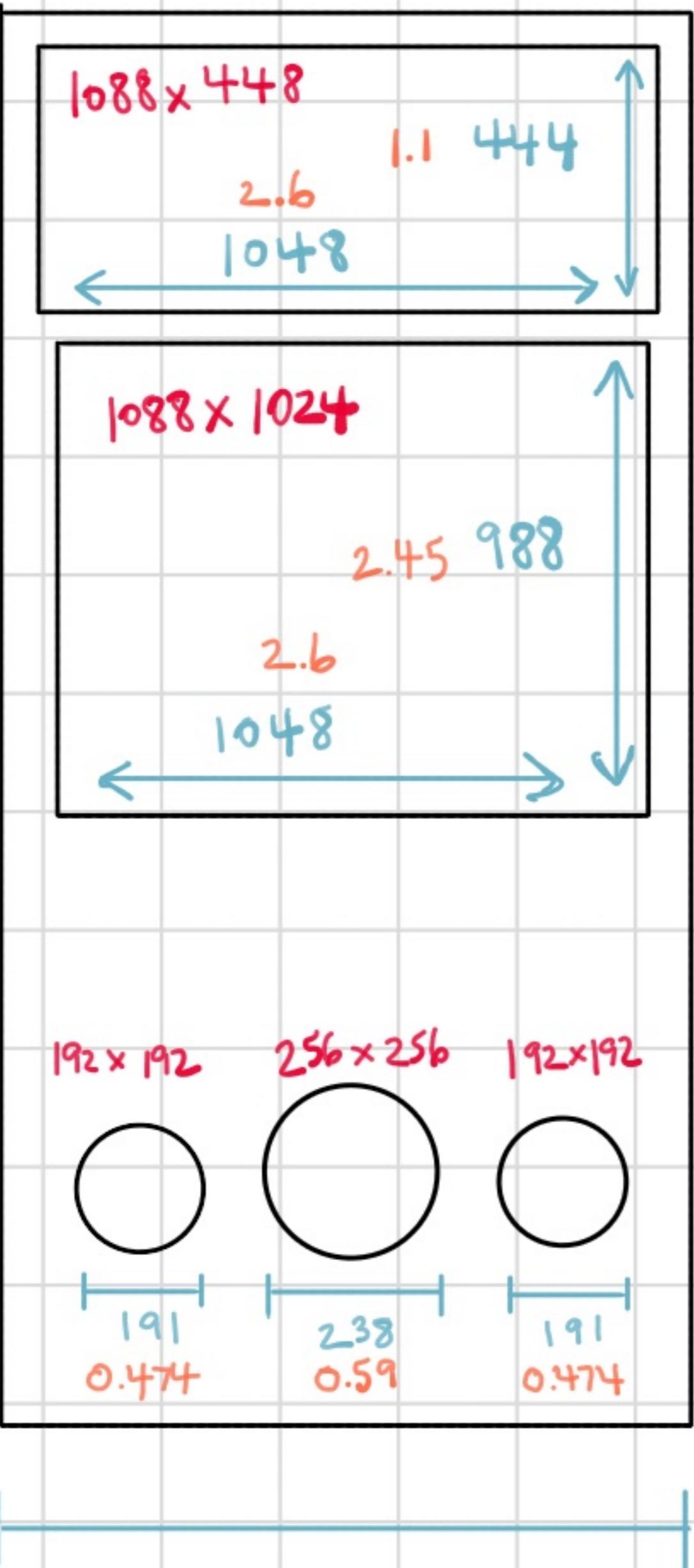
Frame	Samples
0	0 .. < 4410
1	735 .. < 5145
2	1470 .. < 5880
3	2205 .. < 6615
4	2940 .. < 7350
5	3675 .. < 8085

px

in

Tx

1152 x 2456



1125

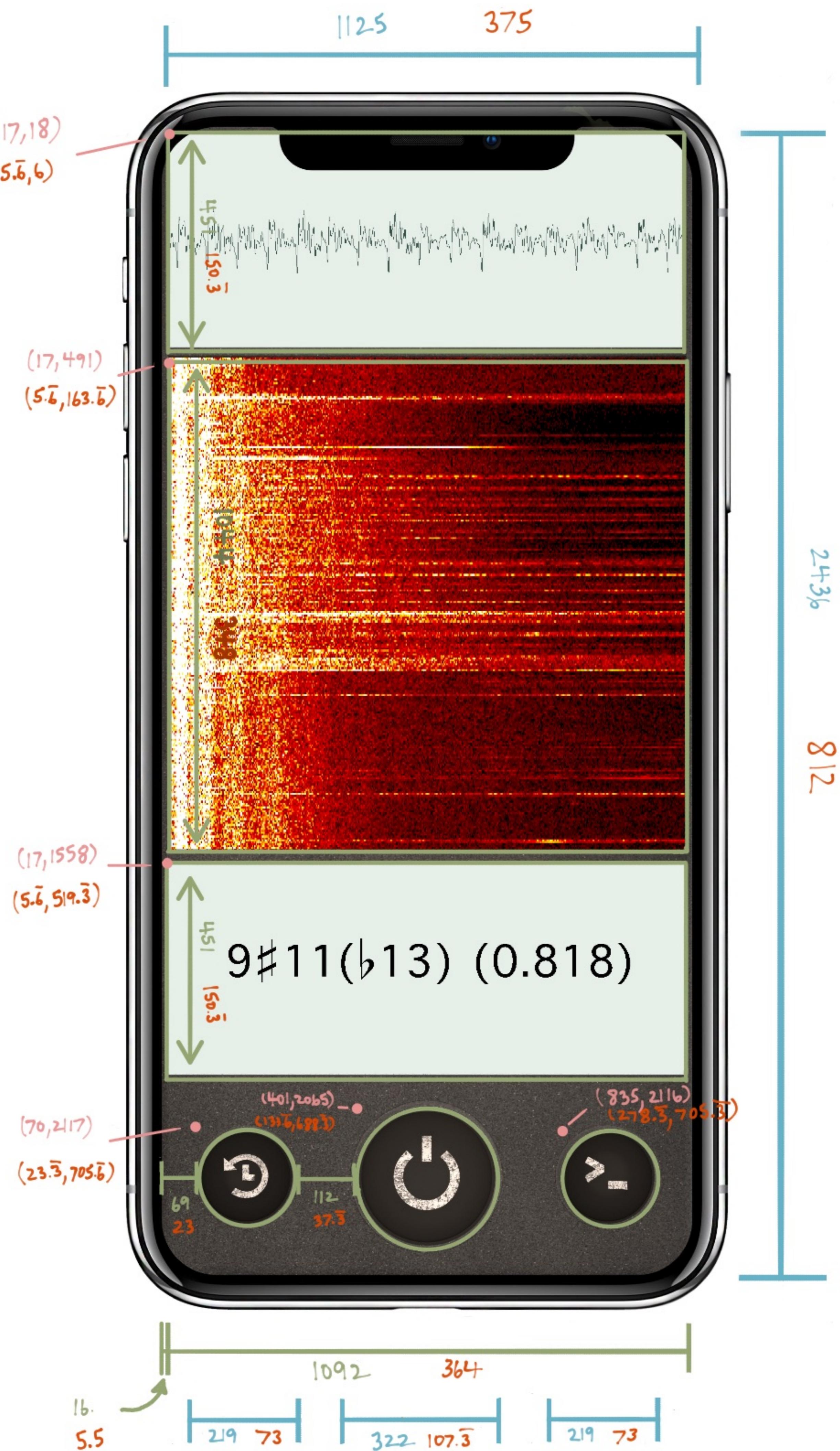
2.79

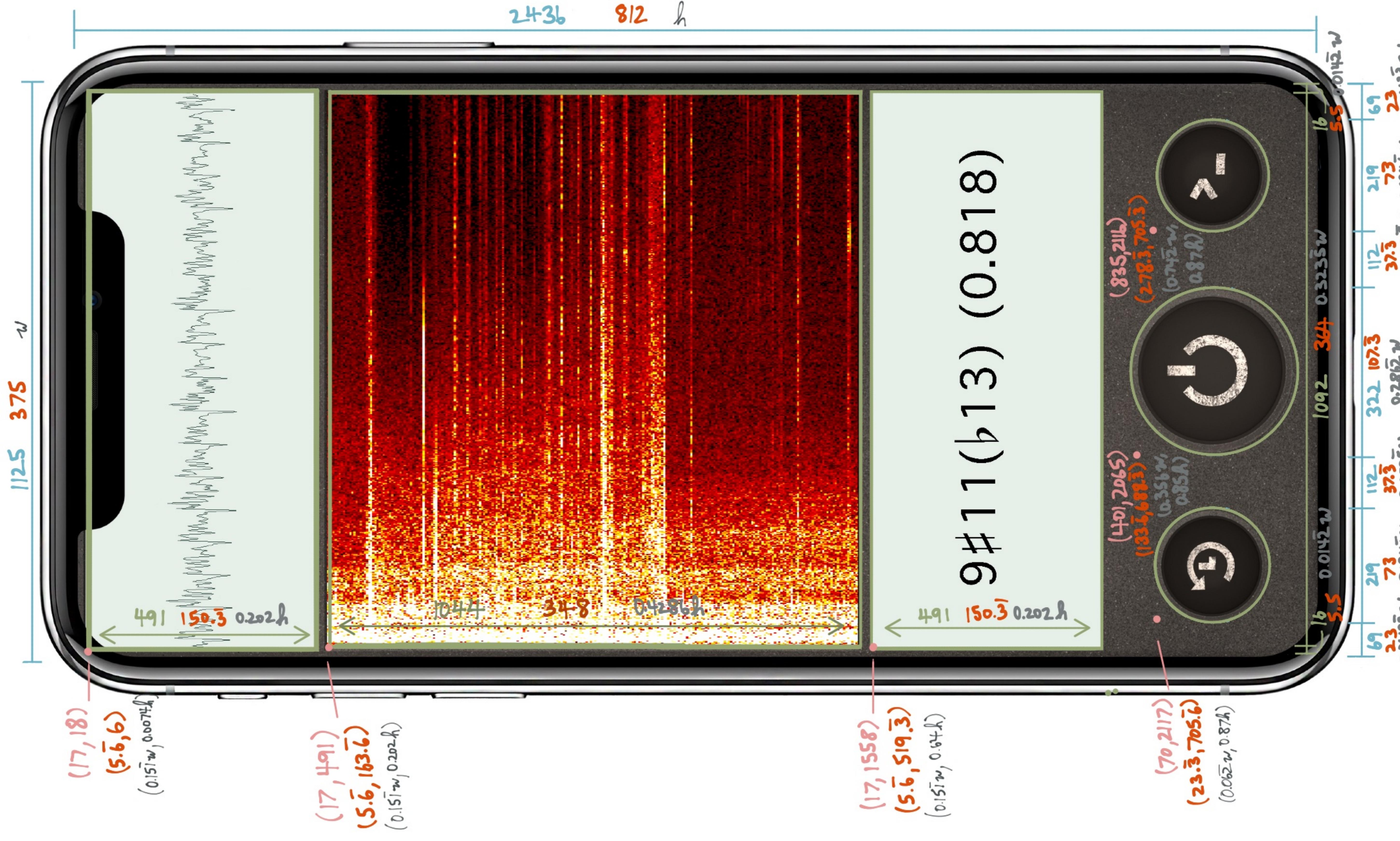
$$\frac{1125}{2.79} = 403.2258064516$$

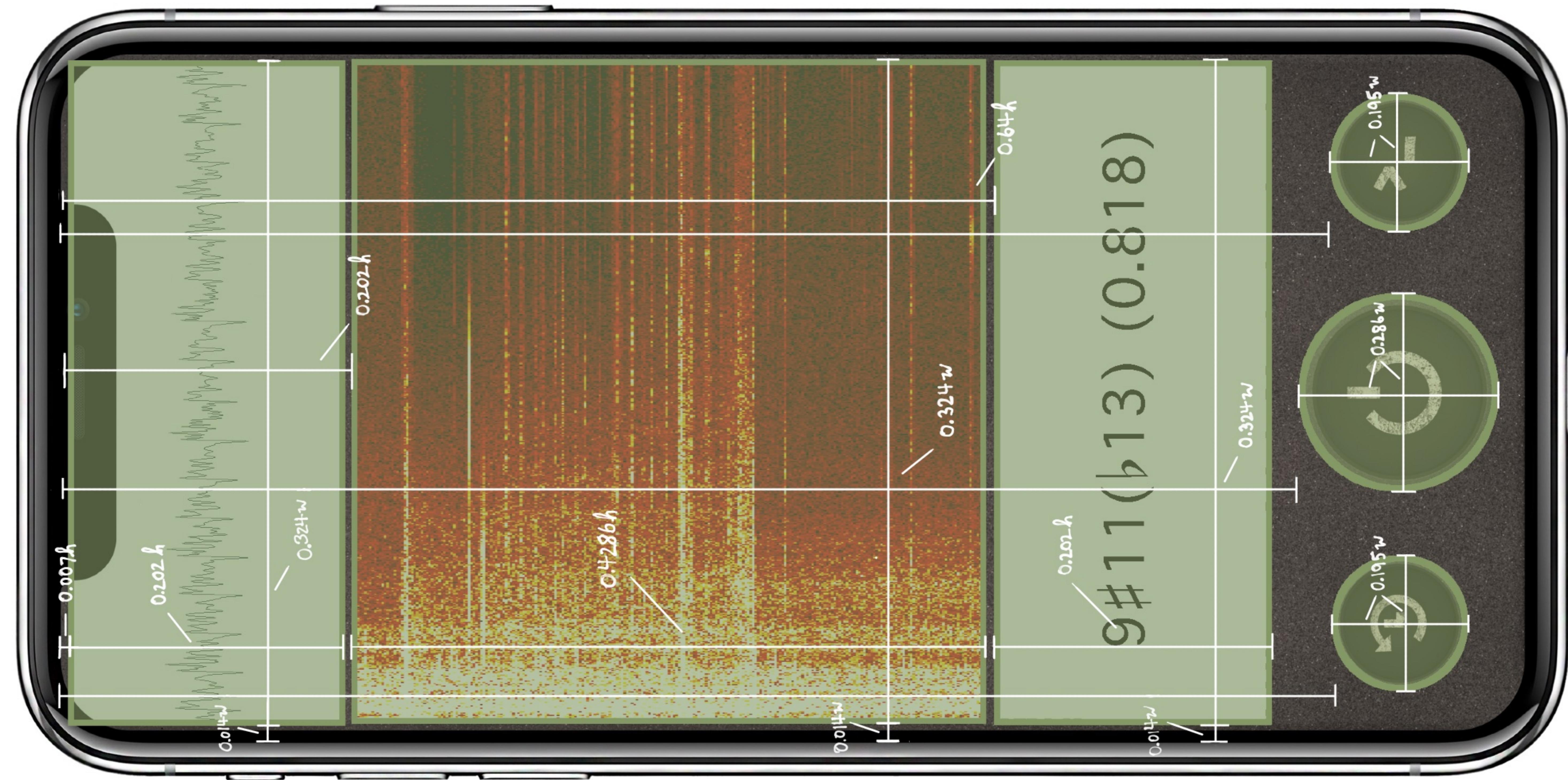
T

2436 6.044

$$\frac{2436}{6.044} = 403.0443414957$$

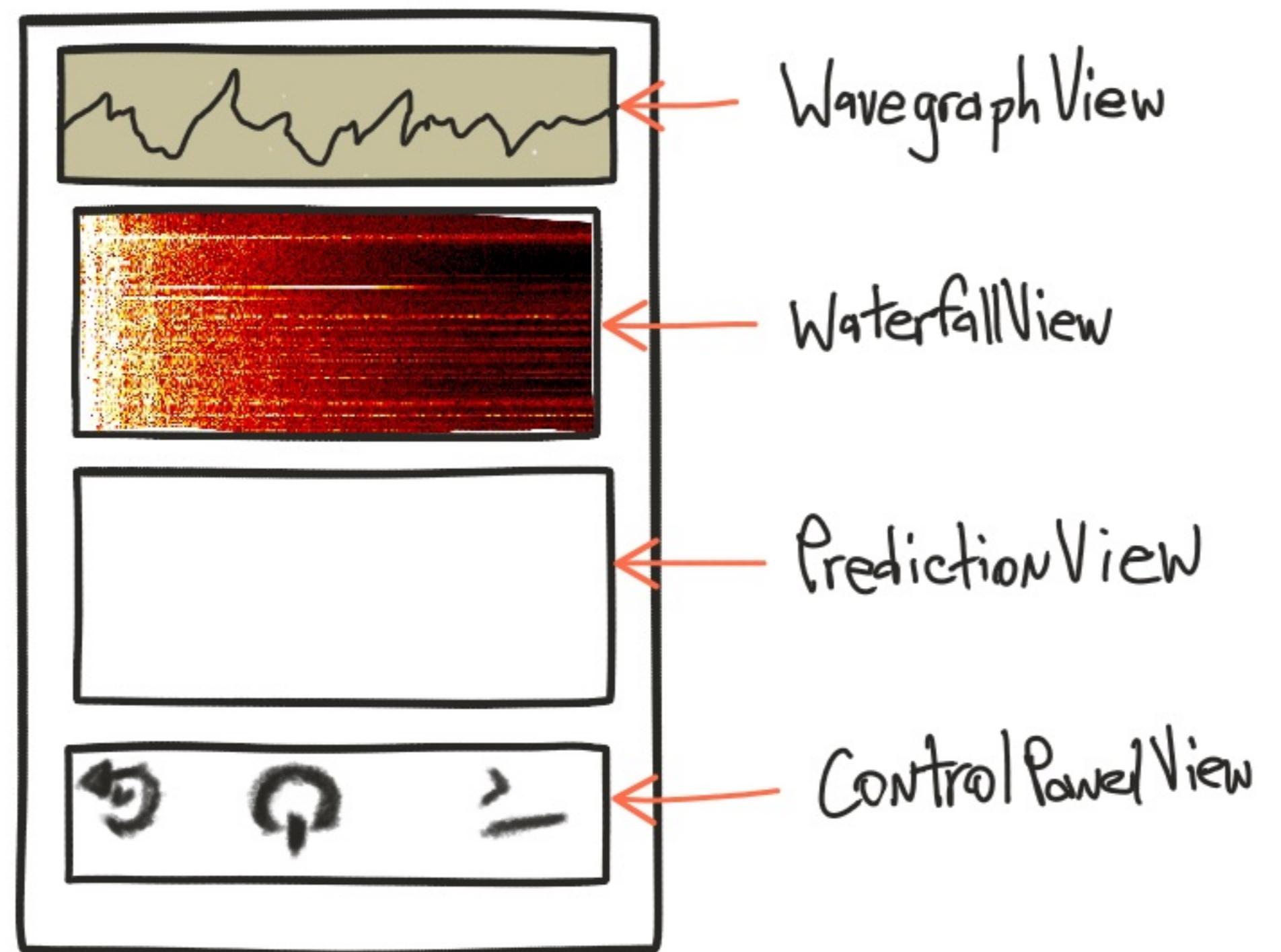






## Front Burner

- UIButton subclass for custom animation
- Training capability added to MLKit
- Aggregate results from testing overlapping audio clips
- Incorporate known characteristics into prediction logic
- Implement 3D model for custom drawing
- Tuner functionality
- Chord library functionality



<u>Pattern</u>	<u>Intervals</u>	<u>Confidence</u>
m (add 9)		0.9352

## Adding Audiobus Support

- Make sure audio session is configured to "mix with others" before becoming active
- Add Audiobus SDK
- Enable background and inter-app audio
- Life Cycle Management:
  - ensure running and active audio session
  - stop audio system when disconnected from Audiobus
  - never stop audio system while Audiobus controller is connected
  - implement suggested background policy
- Set up launch URL
- Register and generate API key
- Add Audiobus Controller instantiation
- Create sender and receiver Audiobus ports adding component descriptions to Info.plist
- Update the Audiobus registry

Audio Manager