

DSP Design

Short-Time Fourier Transform

Shannon L. Sabino

November 2, 2017 (Original)

February 14, 2018 (Edits)

Contents

Introduction	2
Section I: Phase Vocoder	2
Part I: phvoc.m	3
Test I: Speed Up	3
Test II: Slow Down	3
Part II: pitchmod.m	4
Visualization	5
Part III: COLA Property Proof	7
Appendix	10
Project5.m	10
phvoc.m	13
pitchmod.m	13
stft.m	14
istft.m	14
phmap.m	14
spectrogram.m	15

Introduction

In this project, we continue the previous project's theme of pursuing real-time processing methods. As most signals' content varies greatly with time, we are motivated to view these time-changes and filter them accordingly. From this, we find our motivation for the Short-Time Fourier Transform (STFT). To take the STFT, we break our input, x , into sections of length N , called x_m , and make sure that our sections overlap such that we can later reconstruct our signal as if there had been no sectioning done at all. We call the time lag between our sections R . Additionally, we apply a window, $w(n)$, also of length N , to each section to soften the harsh cutoff on the edges, which would create oddities in our frequency processing.

Let's clarify notation before we move on. The input is broken into $M + 1$ segments, x_m , where $m = 0, 1, \dots, M$. Each x_m is defined over N samples, $x_m(n)$, where $n = 0, 1, \dots, N - 1$. In short, we have divided our input into $M + 1$ length- N pieces. If we apply our length- N window, $w(n)$, to each piece, we can think of this as a $N \times (M + 1)$ matrix, which can be more easily passed into matrix-optimized MATLAB functions. Finally, we can describe taking the N -point DFT of each section in a neat equation:

$$X_{k,m} = \sum_{n=0}^{N-1} x(mR + n)w(n)e^{-j\omega_k n} = \sum_{n=0}^{N-1} x_m(n)e^{-j\omega_k n}, \quad (1)$$

where $x_m(n) = x(mR + n)w(n)$. This describes the general STFT. From this, it logically follows that the inverse STFT, the ISTFT, makes use of the N -point inverse DFT, such that

$$x(mR + n)w(n) = x_m(n) = \frac{1}{N} \sum_{k=0}^{N-1} X_{k,m} e^{j\omega_k n} \quad (2)$$

would represent the ISTFT. This solution presents a problem though — we still have the window, $w(n)$, multiplied into the output, which is non-ideal. Thus our final step in taking the ISTFT is a reach-back to our previous project, using the overlap-add (OLA) method of reconstruction, where our final output is

$$y(n) = \sum_{m=-\infty}^{\infty} x_m(n - mR). \quad (3)$$

This is the final expression for the ISTFT used herein.

In practice, the gap R used in the STFT and the R used in the ISTFT can be different — the former is called the analysis hop, and the latter is the synthesis hop size. The effects of having $R_a \neq R_s$ will be described shortly.

Section I: Phase Vocoder

In this section, we create a phase vocoder. Phase vocoders preserve the magnitude of the input signal while altering the phase after taking its STFT and before taking its ISTFT. One application of this is to speed up or slow down the play time of a signal without modulating the audio into a higher or lower pitch. To accomplish this, we can resample each section of the signal that we took for our STFTs, and then place these sections back in

the “appropriate” location by changing the R for our ISTFT. A neat way to describe this is by the factor of speed up, r , where $r = \frac{R_a}{R_s}$. Therefore, given our analysis hop R_a and our desired speed change r , we can find a way to change our time-domain signal without changing our frequency-domain “content”.

Part I: **phvoc.m**

To accomplish this, we write a phase vocoder function in MATLAB, **phvoc**, which takes in our entire original signal, x , the factor of speed up, r , the initial chunk displacement size, R_a , and the length of each section, N . The output, y , is the appropriately slowed down or sped up version of x .

We break the implementation of this function into three parts, the first of which is taking the STFT of x . We choose a hanning window, of the form,

$$w(n) = 0.5 - 0.5 \cos\left(\frac{2\pi n}{N-1}\right). \quad (4)$$

We break x up with the function **buffer**, which politely breaks x into an $N \times (M+1)$ matrix, including any necessary zero-padding at the end to keep it even. We additionally multiply each row of this matrix by our window. Finally, we can pass this entire matrix into the built-in function **fft**, and take an N -point DFT of each row.

Next, we do the phase mapping on these sections. Each sample’s phase is shifted with respect to the analysis hop size and scaled with respect to the synthesis hop size.

Finally, we take the ISTFT, using the synthesis hop size.

Test I: Speed Up

To test our function, we try both slowing down and speeding up a piece of audio. I chose a .wav file that was used a few projects ago because I have no other .wav files on my computer, and ran it through the **phvoc** function with $r = 1.4$, $R_a = 50$ and $N = 3000$. After normalizing the signal (since the speed up also pushed the volume up), the output is indeed a faster version of the input, but with the same pitch. This is saved to *phvocFast.wav*.

Test II: Slow Down

Next, I used $r = 0.7$ with the same R_a and N . The pitch is again unchanged although the signal is much longer. This audio is saved to *phvocSlow.wav*. Examining the magnitude and phase spectrums of each signal shows us that the magnitude is the same between the two signals, while the phase is appropriately stretched on the slower signal.

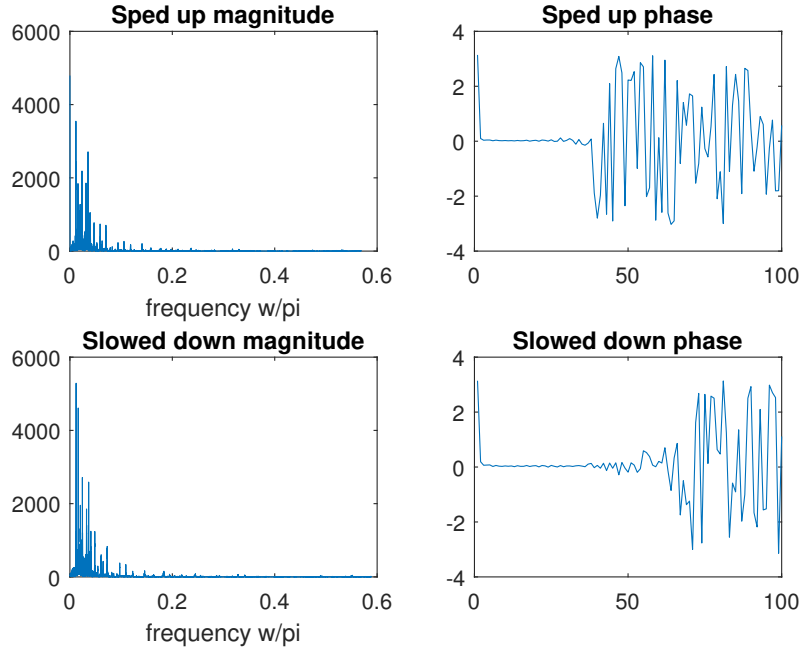


Figure 1: Comparison of slowed and sped up audio after the phase vocoder.

Part II: `pitchmod.m`

In this section, we intentionally modify the pitch of the audio of a flute without changing its speed, using a function we make called **pitchmod**. This function resamples the input by r , which changes the length of the signal, before passing it into **phvoc** with $r = \frac{1}{r}$, which has the overall result of the output being around the same length as the input, but with its pitch shifted up or down. Examine a section of the frequency of the original audio below.

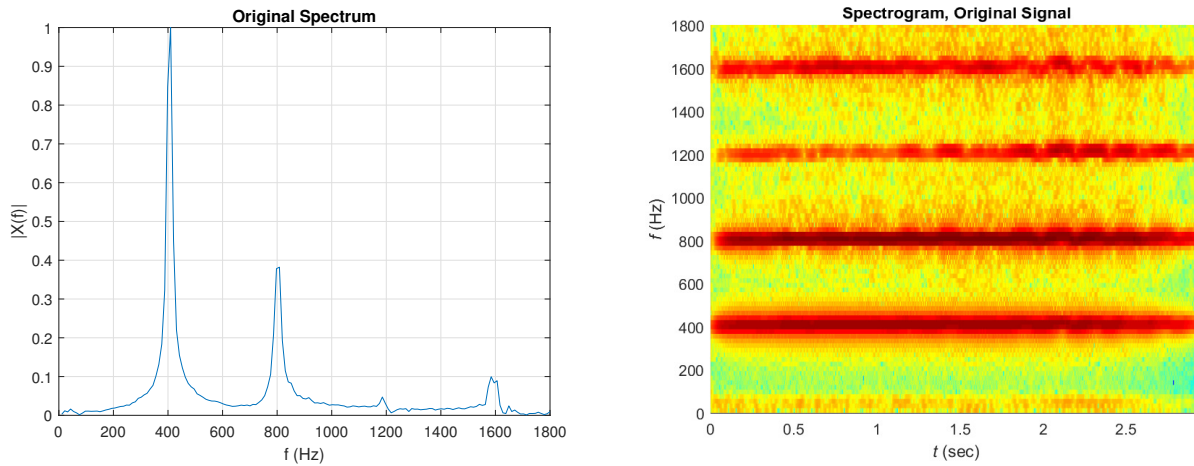


Figure 2: Original signal, notice the peak at 400 Hz.

Notice that the fundamental harmonic in the original audio is at 400 Hz. When we use $r = 2$ to shift these frequencies down, we find the fundamental frequency shifts to $400 \times r = 800$ Hz. Similarly, using $r = \frac{1}{2}$, we see the fundamental frequency shifted to $400 \times r = 200$ Hz. The audio in each case plays for about the same amount of time, but the pitch is altered.

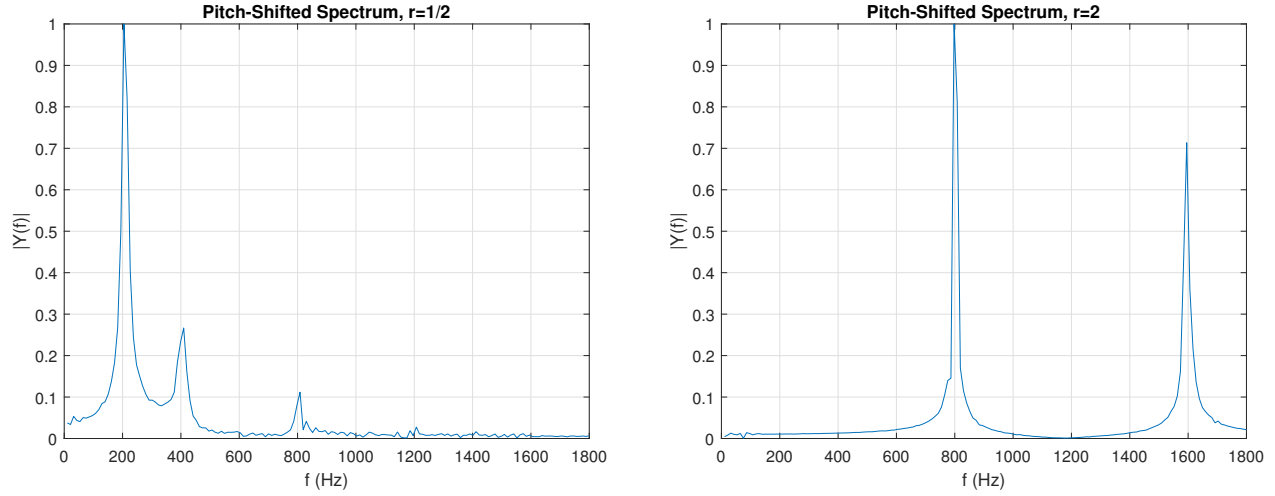


Figure 3: The altered spectrums. (1)

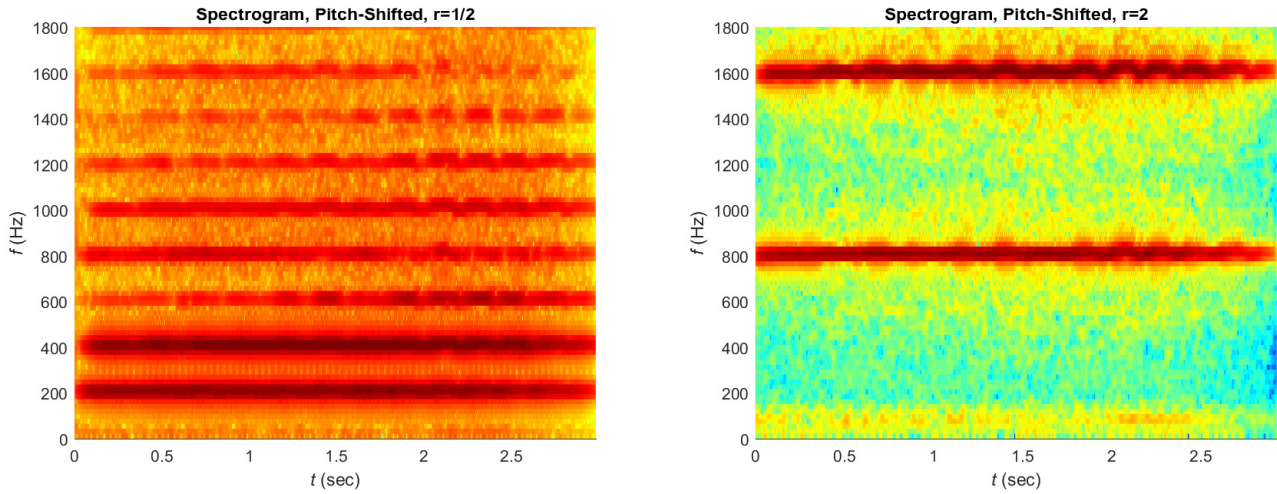
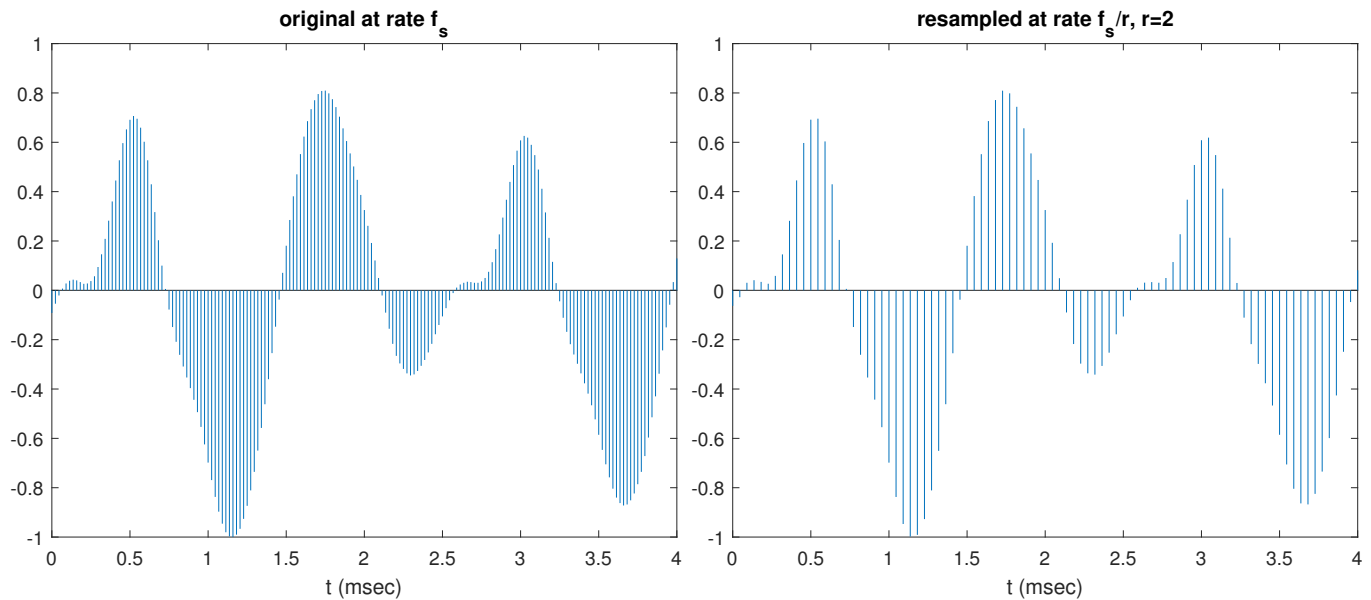


Figure 4: The altered spectrums. (2)

Visualization

Let's explore what happens in the pitch modulation further. Below is 4 milliseconds of the flute audio, unaltered, in the time-domain.



Next, we pass the original signal into the **resample** function, for $r = 2$. As we can see on the right, the envelope of the signal is unchanged, but with the sample rate cut in half, every other sample is just missing.

If we play the resampled signal at half the original sampling frequency, it sounds nearly the same. That being said, if you replay the resampled signal at its original sampling frequency, you get the “chipmunked” audio effect. The voice actors for the Chipmunk movies have said in interviews that they have to read the lines very slowly so that their audio can be sped up and made to sound mousier! Pretty weird approach considering that the figure below on the right shows a way that we can end up with chipmunked audio using only the previously described phase vocoder - you can see the added higher frequencies, although the signal is still 4 milliseconds long!

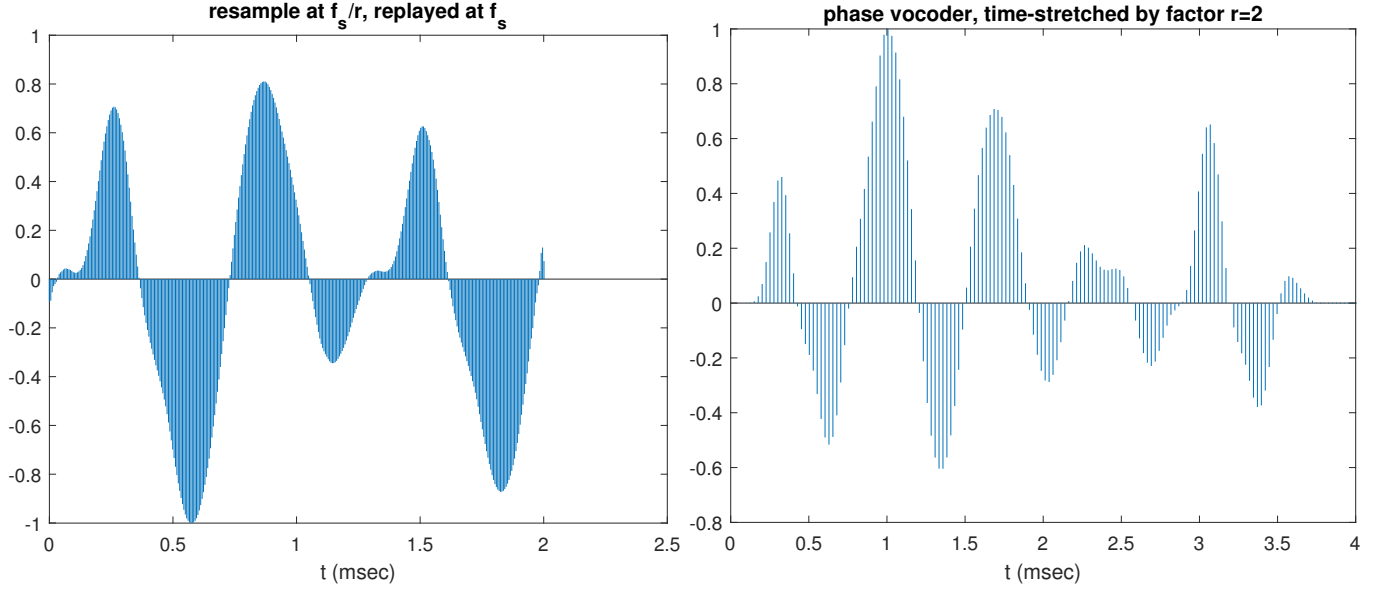


Figure 5: Chipmunked audios.

Part III: COLA Property Proof

The Constant Overlap-Add (COLA) property refers to the property that allows us to use Overlap-Add (OLA) reconstruction to account for our windowing. If this property is fulfilled at a particular R , then our STFT will properly represent the whole signal's frequency domain content.

For the COLA property to be satisfied, a DFT taken at R points should only be nonzero at the constant frequency, $\omega_r = 0$. What this effectively means, is that the windows end up behaving similar to a constant multiplier across all of the chunks in the time domain, and an impulse in the frequency domain. Consider the length-128 hanning window we used previously, $w(n) = 0.5 - 0.5 \cos(\frac{2\pi n}{N-1})$, and an overlapped hanning window with $M = 10$, $\tilde{w}(n) = \sum_{m=0}^M w(n - mR)$. Let's examine how different hop sizes, R , affect the windowing.

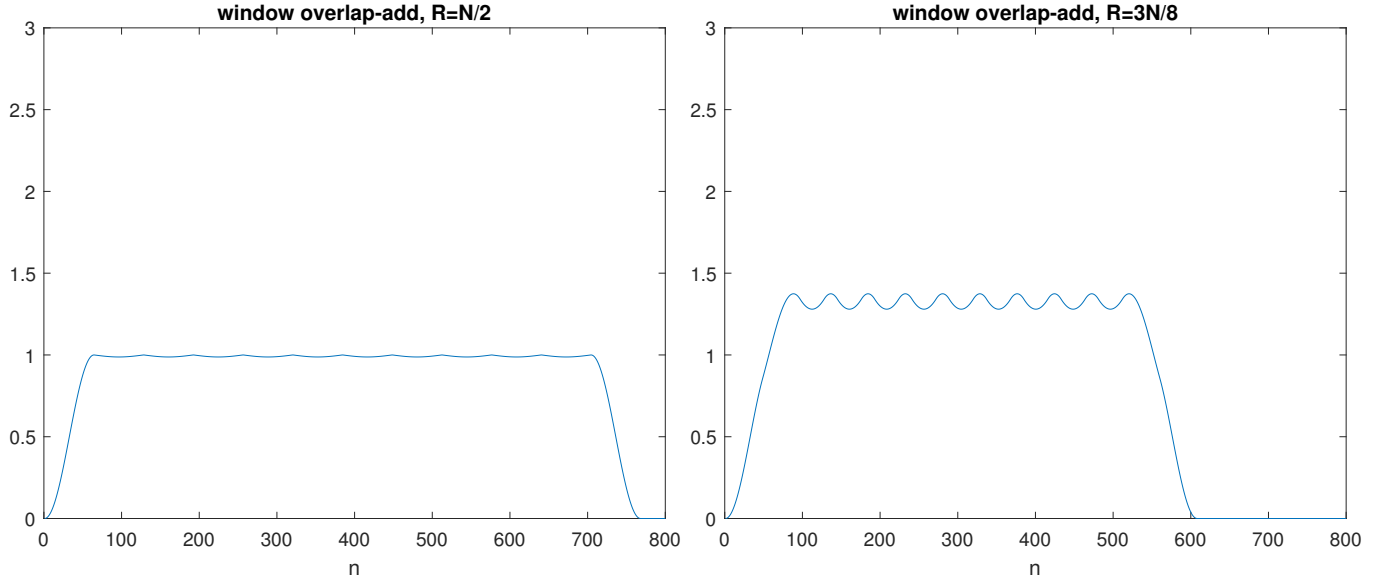


Figure 6: Windows in the time-domain. (1)

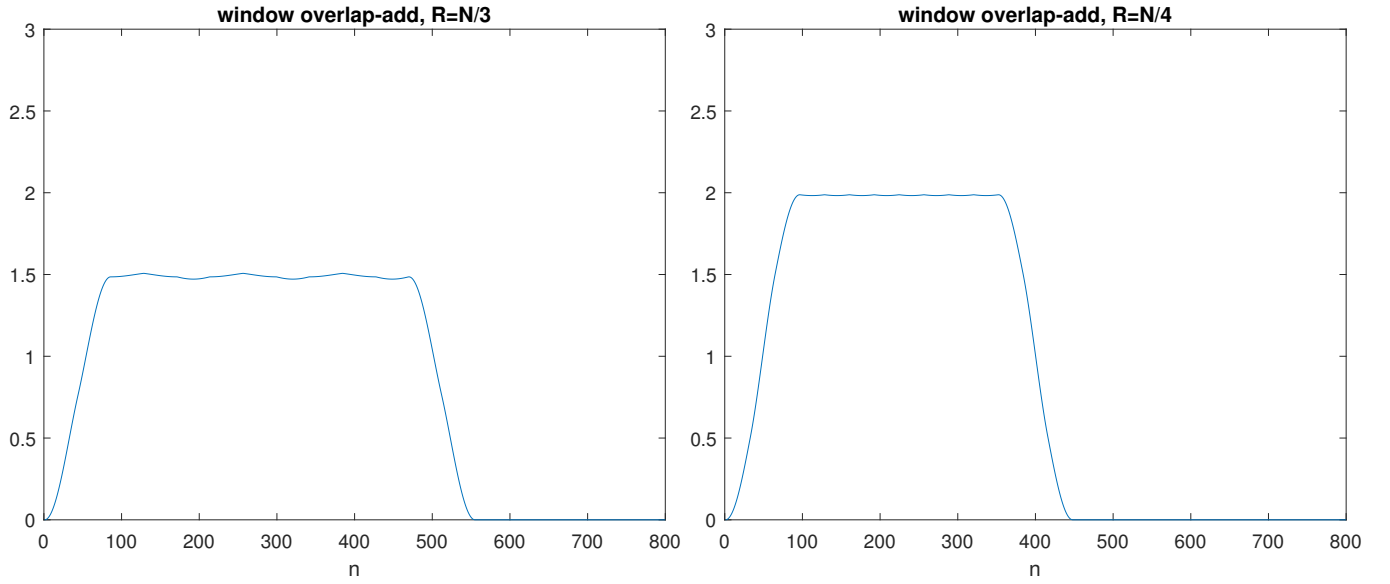


Figure 7: Windows in the time-domain. (2)

Note that in both the biggest and the smallest hop sizes, $R = \frac{N}{2}$ and $R = \frac{N}{4}$, the windows plateau into a tabletop. This can be observed in the frequency domain too.

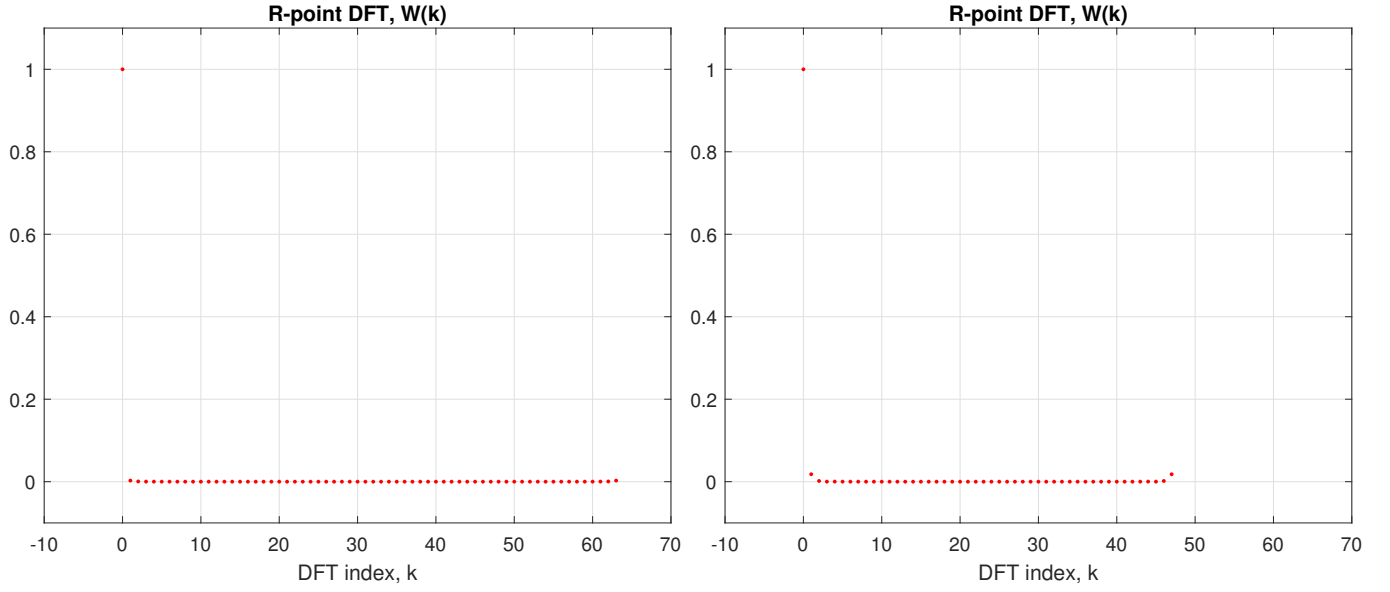


Figure 8: Windows in the frequency-domain for $R=N/2$ (left) and $R=3N/8$ (right). (1)

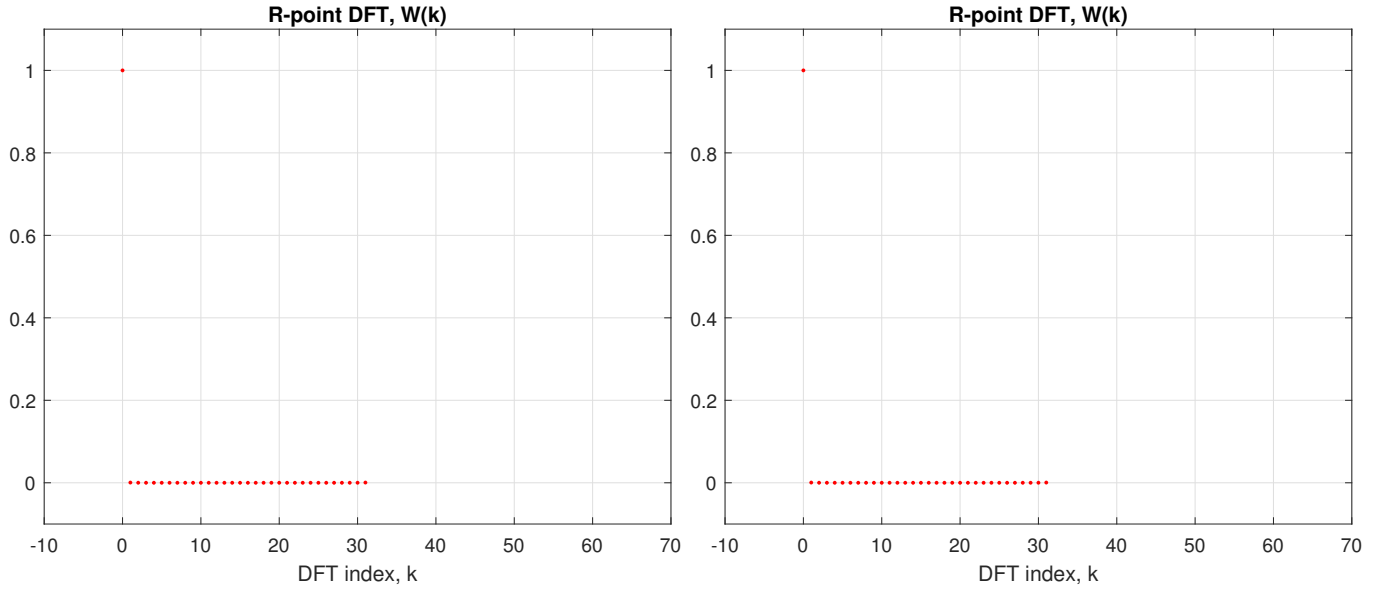


Figure 9: Windows in the frequency-domain for $R=N/3$ (left) and $R=N/4$ (right). (2)

Note that almost all of the frequency content is in the DC frequency — satisfying the COLA property.

Appendix

Project5.m

```
1 %3.1
2 addpath('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files');
3
4 [x,Fs]=audioread('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-1-
    files\original.wav');
5 %% 3.1.a and 3.1.b
6 %phvoc(x,r,Ra,N)
7 yfast=phvoc(x,1.4,50,3000);
8 yfast=yfast./max(yfast); %normalize
9 sound(yfast,Fs);
10 audiowrite('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\
    phvocFast.wav',...
11           [x; zeros(Fs,1); yfast], Fs);
12
13 yslow=phvoc(x,.7,50,3000);
14 yslow=yslow./max(yslow);
15 sound(yslow,Fs);
16 audiowrite('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\
    phvocSlow.wav',...
17           [x; zeros(Fs,1); yslow], Fs);
18
19 figure; plot((1:length(x))./Fs,x)
20 xlabel('time (seconds)'); title('Original Audio')
21 %
22 fastfreq=fft(yfast);
23 fastmag=abs(fastfreq); fastph=angle(fastfreq);
24 slowfreq=fft(yslow);
25 slowmag=abs(slowfreq); slowph=angle(slowfreq);
26 figure;
27 subplot(2,2,1); plot((1:round(length(fastmag)/4))./3000./17,fastmag(1:round(
    end/4)));
28 xlabel('frequency w/pi'); title('Sped up magnitude')
29 subplot(2,2,3); plot((1:round(length(slowmag)/8))./3000./16,slowmag(2:2:round(
    end/4)));
30 xlabel('frequency w/pi'); title('Slowed down magnitude')
31 subplot(2,2,2); plot(1:100,fastph(1:100));
32 title('Sped up phase')
33 subplot(2,2,4); plot(1:100,slowph(1:100));
34 title('Slowed down phase')
35
36 print('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\
    phvoc', '-depsc')
37
38 %phase vocoder changes time stuff without changing frequency content
39 %% 3.2.a
40
41 [x,FS] = audioread('flute2.wav'); x = x(:,1)';
42
43 % y = pitchmod(x, r, Ra, N)
44 yhigh=pitchmod(x,2,256,2048);
45 %STILL need to take real part >:(
46 ylow=pitchmod(x,1/2,256,2048);
47 %should these boys changing how long they are in time?
48 soundsc(x,FS)
49 soundsc(yhigh,FS)
50 soundsc(ylow,FS)
```

```

51 freq=abs(fft(x,4096));
52 freq=freq./max(freq);
53 figure;
54 plot((1:(round(1800/(FS)*4096)))./334*1800*2,freq(1:(round(1800/(FS)*4096))))
55 xlabel('f (Hz)'); ylabel('|X(f)|'); title('Original Spectrum')
56 grid on;
57 print('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\
    OriginalSpectrum2', '-depsc')
58 freq=abs(fft(yhigh,4096));
59 freq=freq./max(freq);
60 figure;
61 plot((1:(round(1800/(FS)*4096)))./334*1800*2,freq(1:(round(1800/(FS)*4096))))
62 xlabel('f (Hz)'); ylabel('|Y(f)|'); title('Pitch-Shifted Spectrum, r=2')
63 grid on;
64 print('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\
    HighSpectrum2', '-depsc')
65 freq=abs(fft(ylow,4096));
66 freq=freq./max(freq);
67 figure;
68 plot((1:(round(1800/(FS)*4096)))./334*1800*2,freq(1:(round(1800/(FS)*4096))))
69 xlabel('f (Hz)'); ylabel('|Y(f)|'); title('Pitch-Shifted Spectrum, r=1/2')
70 grid on;
71 print('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\
    LowSpectrum2', '-depsc')
72
73 %% 3.2.c
74 load x4.mat;
75 fs=44100;
76 figure;
77 stem((0:1/fs:(.004)).*1000,x4,'marker','none')
78 print('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\
    COrig', '-depsc')
79
80 r=2;
81 [p,q]=rat(1/r);
82 x=resample(x4,p,q);
83 figure; stem((0:1/2/fs:(.004+1/2/fs)).*1000,x,'marker','none')
84 print('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\Cr2',
    '-depsc')
85
86 figure; stem((0:1/4/fs:(.002+1/4/fs)).*1000,x,'marker','none')
87 print('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\Cr1
    -2', '-depsc')
88
89 y=phvoc(x,1/r,4,32);
90 figure;
91 stem(linspace(0,4,length(y)),y./max(y),'marker','none')
92 title('phase vocoder, time-stretched by factor r=2')
93 xlabel('t (msec)')
94 print('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\
    Cphvoc', '-depsc')
95
96 %figure 2 and this method changes the frequency content and pushes the time
97 %back to normal
98
99 %% 3
100 N=128; M=10;
101 % n=1:N;
102
103 % w=.5-.5*cos(2*pi*(1:N)./(N-1));

```

```

104 w=@(n) .5-.5.*cos(2*pi*(n)./(127))';
105 wn=w(0:(N-1));
106 wprime=zeros(800,1);
107 R=N/2;
108 for m=0:M
109     %m=m+1
110     wprime(((R*m+1):(R*m+128)))=wprime(((R*m+1):(R*m+128))) + w(0:(N-1));
111     %hold on; plot(wprime,'color',[ (m+1)/19 .5 (m+1)/11])
112 end
113 figure;
114 plot(wprime); title('window overlap-add, R=N/2')
115 xlabel('n'); axis([0 800 0 3])
116 print('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\winN
    -2', '-depsc')
117 %get dft from equation 7
118 r=0:(R-1);
119 wr=2*pi.*r./R;
120 W=zeros(R,1);
121 for n=0:(N-1)
122     W=W+w(n).*exp(-1j.*wr.*n)';
123 end
124 plot(abs(W)./max(abs(W)), 'r. '); title('R-point DFT, W(k)')
125 grid on; xlabel('DFT index, k'); axis([-10 70 -.1 1.1])
126 print('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\dftN
    -2', '-depsc')

127 %-----
128 wprime=zeros(800,1);
129 R=N*3/8;
130 for m=0:M
131     wprime(((R*m+1):(R*m+128)))=wprime(((R*m+1):(R*m+128))) + w(0:(N-1));
132 end
133 figure;
134 plot(wprime); title('window overlap-add, R=3N/8')
135 xlabel('n'); axis([0 800 0 3])
136 print('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\
    winN3-8', '-depsc')
137 %get dft from equation 7
138 r=0:(R-1);
139 wr=2*pi.*r./R;
140 W=zeros(R,1);
141 for n=0:(N-1)
142     W=W+w(n).*exp(-1j.*wr.*n)';
143 end
144 plot(r, abs(W)./max(abs(W)), 'r. '); title('R-point DFT, W(k)')
145 grid on; xlabel('DFT index, k'); axis([-10 70 -.1 1.1])
146 print('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\
    dftN3-8', '-depsc')

147 %-----
148 wprime=zeros(800,1);
149 R=N/3;
150 for m=0:M
151     wprime(((R*m+1):(R*m+128)))=wprime(((R*m+1):(R*m+128))) + w(0:(N-1));
152 end
153 figure;
154 plot(wprime); title('window overlap-add, R=N/3')
155 xlabel('n'); axis([0 800 0 3])
156 print('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\winN
    -3', '-depsc')

```

```

159 %get dft from equation 7
160 r=0:(round(R)-1);
161 wr=2*pi.*r./R;
162 W=zeros(round(R),1);
163 for n=0:(N-1)
164     W=W+w(n).*exp(-1j.*wr.*n)';
165 end
166 plot(r, abs(W)./max(abs(W)), 'r. '); title('R-point DFT, W(k)')
167 grid on; xlabel('DFT index, k'); axis([-10 70 -1 1.1])
168 print('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\dftN
-3', '-depsc')

169 %————
170
171 wprime=zeros(800,1);
172 R=N/4;
173 for m=0:M
174     wprime((R*m+1):(R*m+128))=wprime((R*m+1):(R*m+128)) + w(0:(N-1));
175 end
176 figure;
177 plot(wprime); title('window overlap-add, R=N/4')
178 xlabel('n'); axis([0 800 0 3])
179 print('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\winN
-4', '-depsc')

180 %get dft from equation 7
181 r=0:(round(R)-1);
182 wr=2*pi.*r./R;
183 W=zeros(round(R),1);
184 for n=0:(N-1)
185     W=W+w(n).*exp(-1j.*wr.*n)';
186 end
187 plot(r, abs(W)./max(abs(W)), 'r. '); title('R-point DFT, W(k)')
188 grid on; xlabel('DFT index, k'); axis([-10 70 -1 1.1])
189 print('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\dftN
-4', '-depsc')

```

phvoc.m

```

1 function y = phvoc(x,r,Ra,N)
2 % y = phvoc(x,r,Ra,N)
3 %
4 % x = input audio signal
5 % r = speed-up factor
6 % Ra = analysis hop size, synthesis hop size Rs = round(Ra/r)
7 % N = FFT frame length
8 %
9 % y = output signal
10
11 X = stft(x,Ra,N); % STFT of input, implementing Eq. (11)
12 %clear x
13 Y = phmap(X,r,Ra,N); % phase remapping, implementing Eqs. (15) and (16)
14 %clear X
15 Rs=round(Ra./r); %new speed
16 y = istft(Y,Rs,N); % ISTFT/OLA, implementing Eqs. (12) and (13)
17 clear Y
18
19 end

```

pitchmod.m

```

1 function y = pitchmod(x, r, Ra, N)

```

```

2 % x = input audio signal
3 % r = pitch scaling factor
4 % Ra = analysis hop
5 % N = FFT frame length
6 %
7 % y = output signal
8
9 [p,q]=rat(1/r);
10 xtemp=resample(x,p,q);
11 y = phvoc(xtemp,1/r,Ra,N);
12
13 end

```

stft.m

```

1 function X = stft(x,Ra,N)
2 % X = stft(x,Ra,N)
3 % Short Time Fourier Transform
4
5 w=.5-.5*cos(2*pi*(1:N)/(N-1)); %check 0/1 indexing
6
7 Xtemp=buffer(x,N,N-Ra,'nodelay'); %x chunks
8 [N1,M]=size(Xtemp); %grab M
9 W=repmat(w,1,M); %window
10 X=fft(W.*Xtemp,N);
11
12 end

```

istft.m

```

1 function y = istft(Y,Rs,N)
2 % y = istft(Y,Rs,N)
3
4 w=.5-.5*cos(2*pi*(1:N)/(N-1)); %check 0/1 indexing
5
6 %y=zeros(size(Y));
7 ytemp=real(iff(Y,N));
8
9 [N1,M]=size(Y); %N1=N?
10 y=zeros(M*Rs+N,1); %?
11
12 for m=0:(M-1) %?
13     for n=1:(N)
14         y(m*Rs+n)=y(m*Rs+n)+ytemp(n,m+1).*w(n);
15     end
16 end
17 end

```

phmap.m

```

1 function Y = phmap(X,r,Ra,N)
2 %phase modification
3 %want |Y|=|X|
4
5 Ymag=abs(X);
6 %X(k,m)=xm(k)
7
8 Rs=round(Ra./r); %new speed
9 psi=zeros(size(X));

```

```

10 phi=angle(X);
11 [N1,M]=size(X); %check this, N1=N if right?
12 w=zeros(size(X));
13 dw=w;
14 mod2pi=@(x) mod(x+pi,2*pi)-pi;
15
16 for k=1:N %check 0/1 indexing
17
18     wk=2*pi*(k-1)/N;
19     psi(k,1)=phi(k,1); %check indexing
20     for m=2:M %check
21         dw(k,m)=1/Ra.*mod2pi(phi(k,m)-phi(k,m-1)-Ra*wk);
22         w(k,m)=wk+dw(k,m);
23         psi(k,m)=psi(k,m-1)+Rs.*w(k,m);
24     end
25 end
26 Y=Ymag.*exp(1j.*psi);
27 end

```

spectrogram.m

```

1 R = 256; % hop size
2 N = 2048; % FFT length
3 [x,fs] = audioread('flute2.wav'); x = x(:,1)';
4 yhigh=pitchmod(x,2,256,2048);
5 %STILL need to take real part >:(
6 ylow=pitchmod(x,1/2,256,2048);
7 X = stft(ylow,R,N); % assumed to have size Nx(M+1)
8 M = size(X,2) - 1;
9 k = 0:N/2; % positive-frequency half of Nyquist interval
10 f = k*fs/N; % 0 <= f <= fs/2
11 t = (0:M)*R/fs; % M+1 time frames spaced at hop interval R*Ts = R/fs
12 Xmag = abs(X(k+1,:)); % extract positive-frequency half
13 Xmag = Xmag/max(max(Xmag)); % normalize to unity maximum
14 S = 20*log10(Xmag); % dB
15 surf(t,f,S,'edgecolor','none'); % spectrogram plot, t horizontal, f vertical
16 colormap(jet);
17 axis tight; view(0,90);
18 xlabel('{\itt} (sec)');
19 ylabel('{\itf} (Hz)');
20 ylim([0,1800]); % show range 0 <= f <= 1800 Hz
21
22 print('C:\Users\shsab\OneDrive\Documents\FALL 2017\DSPD\project-5-files\
    HighGram', '-depsc')

```