

Phase Vocoder without Occurring the Phase Unwrapping Problem

Cristian Constantin 341C2

Abstract:

Un phase vocoder este un procesor de semnal folosit de exemplu in industria audio pentru a obține diverse efecte precum auto-tune sau aplicații mai simple precum schimbarea vitezei de redare. Problema cu un astfel de phase vocoder este ca uneori pot apărea probleme precum o distorsiune a semnalului când poate vrem sa redam un podcast la viteza 1.5x, iar asta nu este o situatie ideală pentru ca am pierde informație astfel si dintr-un semnal audio pe care îl putem auzi fără problem ajungem in situația in care sa nu putem înțelege cuvintele rostite in semnalul inițial.

Introducere:

Articolul Time-Scale and Pitch-Scale Modification by the Phase Vocoder without Occurring the Phase Unwrapping Problem scris de Ryoichi Yoneguchi (Meiji University Graduate School of Science and Technology) si Takahiro Murakami (Meiji University School of Science and Technology) prezintă o varianta de phase-vocoder care evita problema de phase unwrapping intalnita intr-un phase vocoder normal, problema care apare atunci când vocoderul estimează frecvența unghiulară instantanee, iar din cauza limitărilor în calculul fazei rezultă distorsiuni în forma de undă a semnalului modificat. Varianta propusă este o abordare diferită care in loc să se bazeze pe estimarea frecvenței unghiulare instantanee, aceasta folosește direct diferența de faze între cadrele de analiză adiacente pentru a modifica faza semnalului si astfel evită necesitatea de a estima frecvența unghiulară instantanee, acest lucru ajutând la evitarea problemei de tip phase unwrapping.

Pentru acest proiect am pornit de la aceasta abordare descrisă in articol pentru a obține un phase vocoder care sa ofere rezultate mai bune fata de unul convențional, astfel proiectul include 2 implementări: o implementare pentru un phase vocoder normal si o implementare care pornește de la un phase vocoder normal si face modificările pentru semnal urmând ideea din articol.

Un phase vocoder simplu implica in esență 3 pași principali pentru a obține un nou semnal: aplicarea FFT(Fast Fourier Transform) asupra semnalului astfel încât sa fie posibilă procesarea sa pentru ca astfel trecem semnalul in domeniul frecvenței, procesarea fazei si a magnitudinii in cadrul căreia ajungem sa aplicam parametrii primiți precum rata de redare sau rata pentru pitch si in final aplicam un IFFT(Inverse Fast Fourier Transform) pentru a trece semnalul in domeniul timpului așa cum era trimis inițial către FFT.

Implementare:

Implementarea proiectului a presupus 3 etape: Implementarea funcțiilor de tip helper, Implementarea celor 2 variante de phase vocoder si Implementarea secțiunii de testare pentru rezultatele generate.

Functii Helper - Pentru a putea utiliza un phase vocoder in python in mod inițial am nevoie sa pot trimite semnalul audio către funcția de phase vocoder si ulterior sa pot scrie într-un fișier audio semnalul rezultat. In cadrul implementării exista 2 funcții helper: write wave si read_wave. Write_wave folosește modulul wave din python si setea parametrii necesari pentru a scrie fișierul wave. Read-wave este o functie cu unicul scop de citi fișierul audio pentru phase vocoder si începe prin a verifica sample_width-ul care in esență de referă la calitatea audio si am 2 cazuri pe care le tratez: 16 bit si 32 bit. In cazul proiectului meu, generez datele din intrare in GarageBand si export-ul pentru fișiere are setarea de 16-bit care mai este cunoscut si drept "CD

Quality". Dacă avem unul din cele 2 cazuri trebuie să verificăm dacă avem un semnal mono sau stereo pentru că avem nevoie de un semnal mono și astfel dacă este cazul transformăm semnalul stereo în semnal mono și în cazul proiectului meu acest lucru are loc pentru înregistrarea semnalului în mod stereo în GarageBand folosind mai multe microfoane.

Implementare phase vocoder simplu - implementarea pentru phase vocoder începe prin a verifica dacă semnalul primit este unul de tip mono și în cazul în care da efectuăm STFT pe semnal cu un window size pentru transformata Fourier de 1024 și hop size de 256. După aplicarea STFT avem nevoie de matricea în care stocăm varianta prelucrată a matricei date de STFT pentru că în funcție de rată pentru timp e posibil să avem un semnal de o durată diferită și astfel avem nevoie și de un vector de timp nou. În următoarea secțiune a funcției putem modifica scala temporală, reprezentând astfel secțiunea din cod în care schimbăm viteza de redare a semnalului nostru și în această secțiune realizăm diferențele de fază specifice phase vocoder și tot aici este secțiunea în care poate apărea phase unwrapping și este important acest aspect pentru că o să revenim la el pentru implementarea vocoder-ului din articol. Parcurgem fiecare coloană a matricei STFT care va fi modificată temporal cu un for și pentru fiecare moment t în matricea STFT întinsă temporal, se calculează `col_idx`, un index al coloanei din matricea STFT originală. Acest index este modificat în funcție de factorul de întindere temporală, apoi pentru fiecare coloană, se calculează diferența de fază între cadrele succesive ale matricei STFT [1] și în final aplicăm modificările în matricea STFT. Pentru modificarea pitch-ului ne trebuie o nouă matrice STFT în care o să avem rezultatul final și parcurgem toate frecvențele din matricea STFT modificată temporal și pentru fiecare frecvență k , se calculează un nou index `k_psm`, care este determinat de factorul de scalare a înălțimii (`psm_rate`). La final folosim ISTFT pentru a obține semnalul audio modificat.

Implementare phase vocoder articol - scopul phase vocoder-ului din articol este acela de a elimina problema de phase unwrapping și o primă abordare cât mai simplă pentru această problemă este secțiunea amintită la punctul [1] unde după ce fac diferența de fază folosind `numpy.angle` aș putea foarte simplu să folosesc biblioteca `numpy` și funcția `unwrap` asupra diferenței de fază pentru a rezolva această problemă, dar am ales să urmez articolul și să implementez algoritmi propuși. Diferența în implementare față de phase vocoder-ul simplu este data de modificarea fazelor, astfel în următorul cod: `stretched_stft_matrix_tsm[:, t] = np.abs(stft_matrix[:, col_idx]) * np.exp(1j * (np.angle(stretched_stft_matrix_tsm[:, t - 1]) + phase_diff))` folosim `stretched_stft_matrix_tsm` în loc de `stft_matrix` ca în phase vocoder-ul simplu ca să obținem astfel ceea ce propune articolul: folosirea directă diferenței de fază între cadrele de analiză adiacente. Alta diferență față de abordarea simplă este vizibilă tot în acest cod pentru că articolul menționează faptul că pentru valori ale `tsm_rate` care nu sunt întregi o să fie accentuată problema pentru că în phase vocoder-ul normal în loc să adunăm `phase_diff` direct o să înmulțim cu `tsm_rate`.

Implementare funcții de testare - Pentru această secțiune folosesc `spectrogram` din `scipy.signal` pentru a genera spectrograme pentru semnalele audio necesare pentru analiza cantitativă și funcția de `plot` din `matplotlib.pyplot` pentru a obține un waveform care să ne arate problema de phase unwrapping. Pentru a evalua numeric performanța folosesc o funcție care generează raportul SNR între cele 2 semnale.

Date:

Pentru secțiunea de date am ales să îmi generez 5 fișiere audio în GarageBand, 4 sunt înregistrări ale vocii mele în care număr de la 1 la 10 (acest aspect o să fie detaliat în discuția despre Spectrograme) și 3 dintre aceste 4 fișiere au deja un filtru aplicat asupra lor precum `deeper voice`, `robot voice` sau `telephone voice` pentru a putea evidenția diferența de performanță între cele 2 metode. Fișierul numărul 5 urmează modelul de testare amintit în articolul inițial și reprezintă o înregistrare de aproximativ 30 de secunde a ciripitului de păsări pentru că un astfel de sunet este predispus la problema de phase unwrapping. De menționat este faptul că evaluările au fost făcute pe 3 fișiere: fișierul cu vocea mea fără un filtru drept caz de bază, fișierul cu filtru `deeper voice` din GarageBand întrucât este cel mai dificil din cele 3 filtre pentru phase vocoder în urma testării și fișierul cu ciripitul de păsări pentru a putea reproduce situația menționată în articol.

Tot pentru date am folosit inițial pentru testare si audio_sample-urile din următorul proiect pentru autotune in python: https://github.com/PyExplained/Autotune-using-Python/tree/main/audio_samples si o sa includ rezultatele obținute pentru Sample-ul C Major Vocal pentru ca are o calitate mai ridicată.

Evaluare:

Evaluare cantitativă:

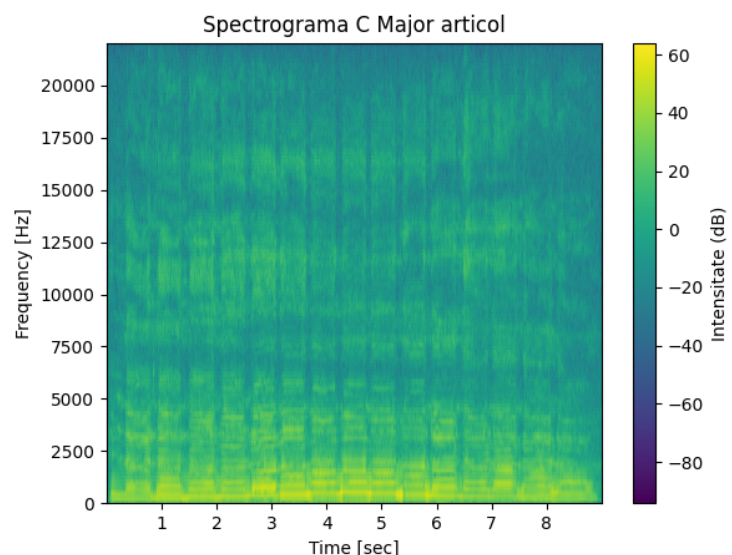
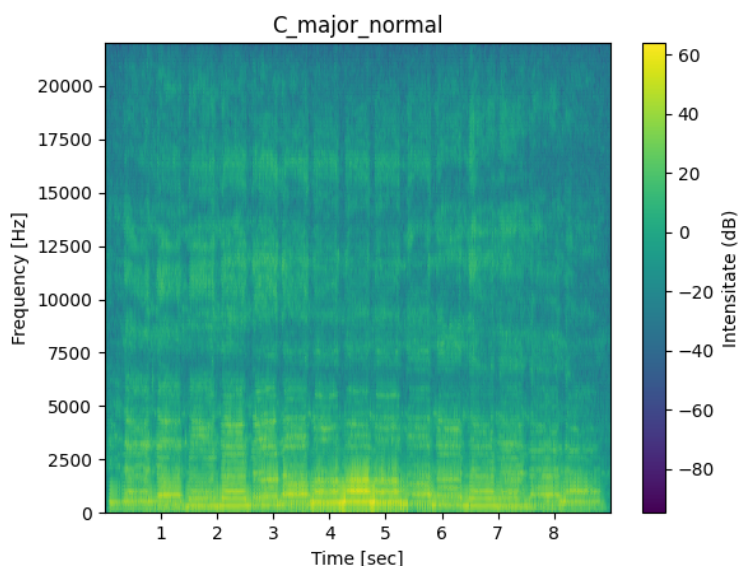
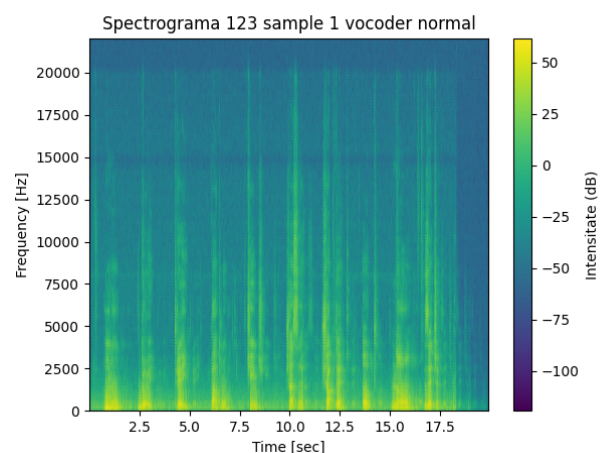
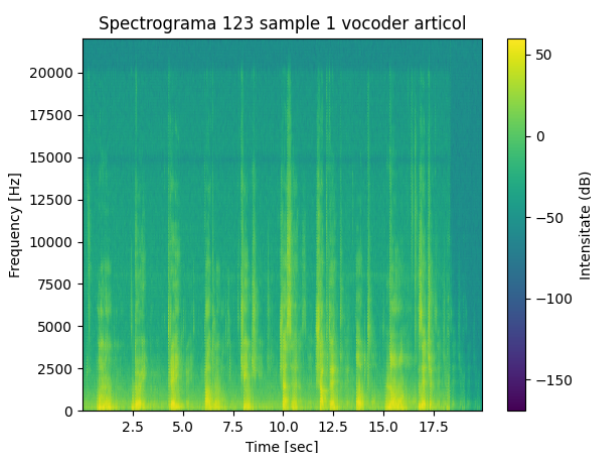
Pentru evaluarea cantitativă am folosit 3 indicii: Spectrograme, Waveform si Raport SNR si am sa menționez sample-urile pentru care se pot obține concluzii mai usor. De menționat pentru aceste valori este faptul ca pentru Sample 1, 2 si C Major parametrii de rulare pentru phase vocoder au fost 1.5 pentru tsm si 1 pentru psm, iar pentru Sample 5 valorile au fost 1.5 si 3 pentru a urma exemplul de simulare din articol.

1. Valori SNR (normal/articol):

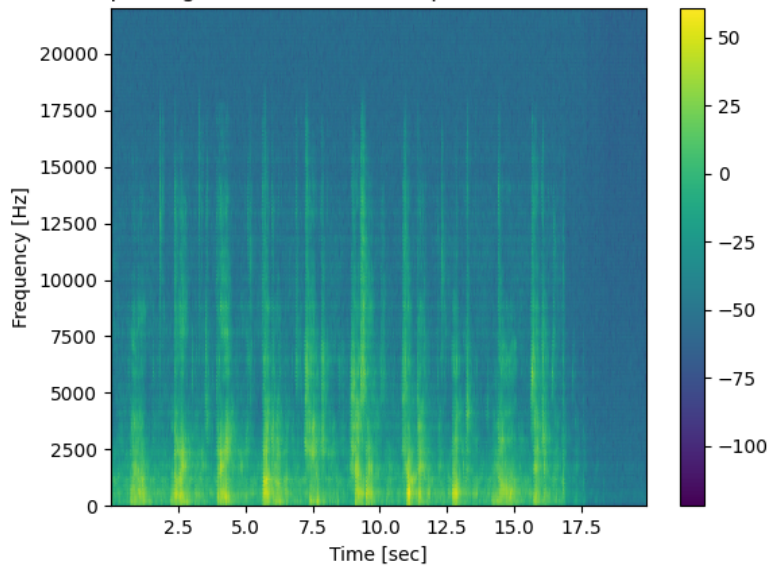
- SNR Sample1: 1.741101442008811 dB
- SNR Sample2: -0.3418349227733592 dB
- SNR Sample5: 0.7324084431242194 dB
- SNR Sample C major: 2.729103844368449 dB

Rezultatele obținute pentru primele 3 sample-uri indica o îmbunătățire daca luam in context si evaluarea calitativă pt ca la sample-ul 1 reusim sa auzim sunetul mai clar, dar pentru sample-ul 2 care este puternic afectat de phase unwrapping deși reușim sa auzim mai bine cifrele avem o valoare negativa si asta se datorează calității reduse a înregistrării, astfel am inclus si Sample-ul C Major de calitate ridicată care indica in mod clar o îmbunătățire cu o valoare de peste 2 DB.

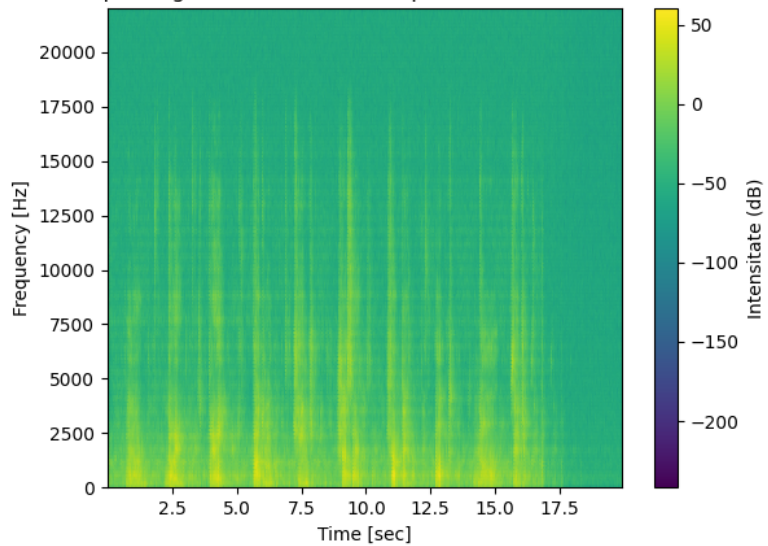
2. Spectrograme:



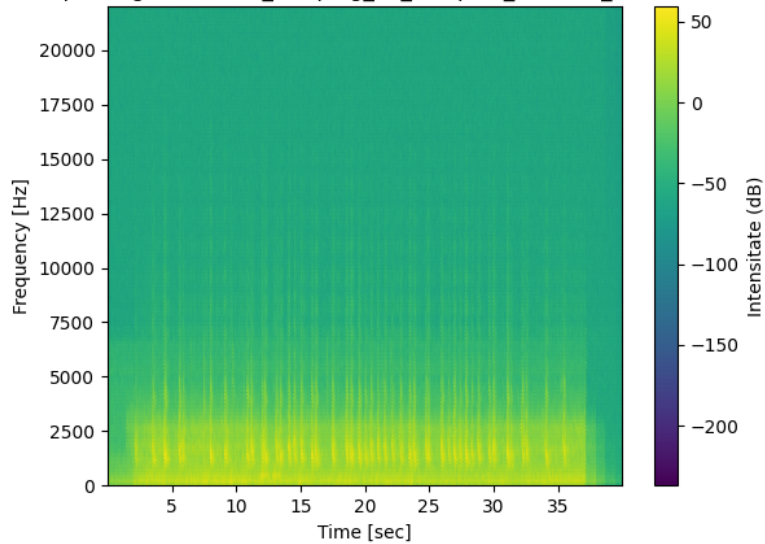
Spectrograma 123 Robot sample 2 vocoder articol



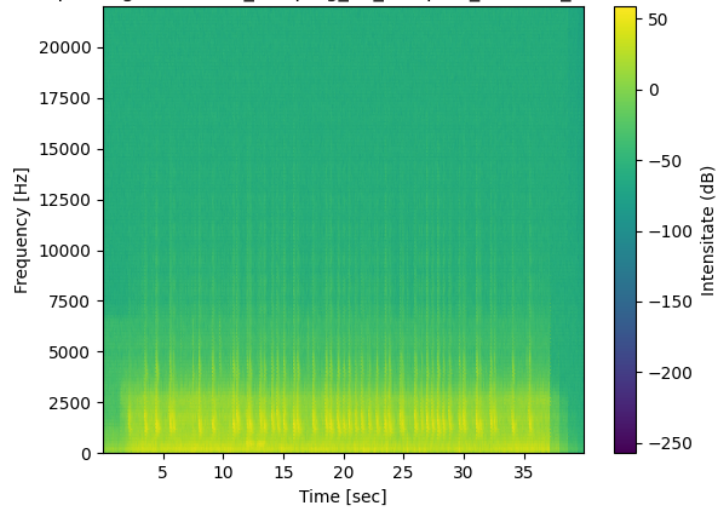
Spectrograma 123 Robot sample 2 vocoder normal



Spectrograma Birds_Chirping_PS_sample5_vocoder_articol



Spectrograma Birds_Chirping_PS_sample5_vocoder_normal

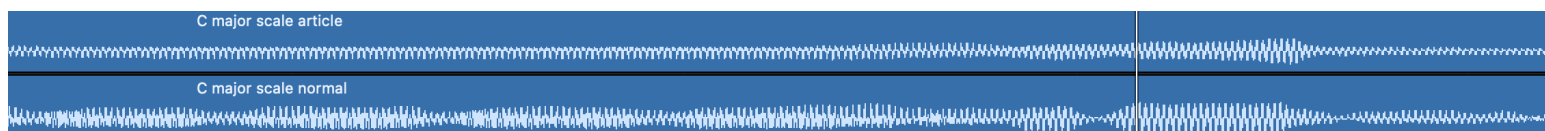
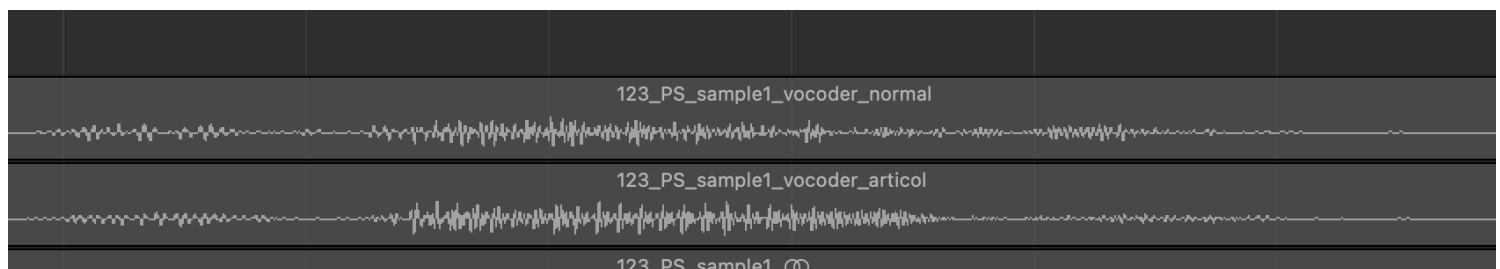
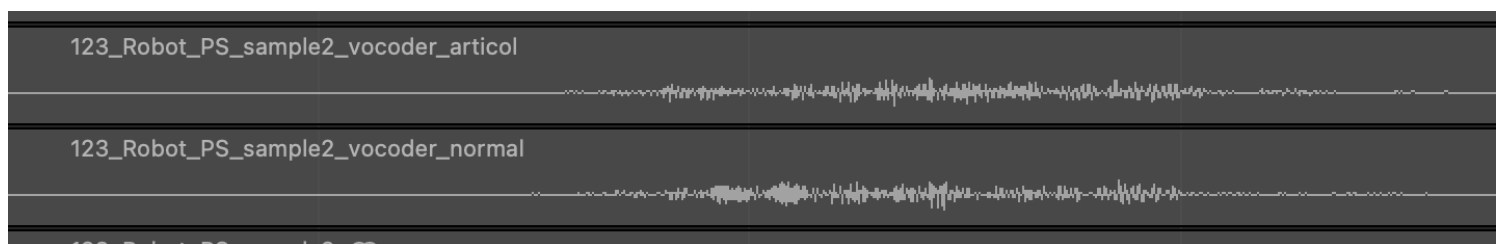
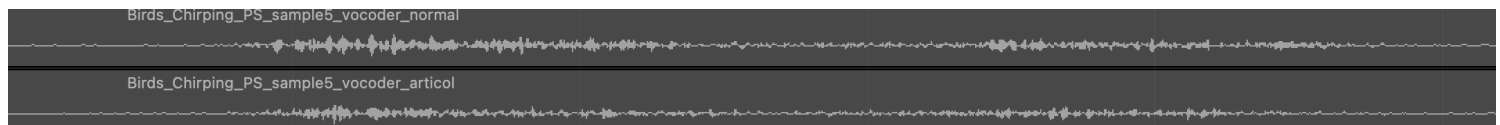


Explicatii Spectrograme:

Pentru Sample 1 se poate observa faptul ca diferența este neglijabila, dar pentru Sample 2 in spectrograma pentru varianta din articol observam clar o îmbunătățire pentru undele asociate numărătorii sunt clar vizibile fata de zgomot. Pentru Sample-ul C Major si Sample 5 nu putem utiliza spectrograma drept un indicator foarte clar al performantei.

3. Waveforms:

Pentru aceasta evaluare am atașat capturi de ecran din GarageBand unde am importat rezultatele întrucât interfața îmi permite sa vad mult mai in detaliu fenomenul de phase unwrapping. Generez in același timp si waveform-uri in Python așa cum am menționat in secțiunea de implementare si acestea se ragesesc in arhiva. Pentru Sample 1 se observa clar avem mai mult detaliu, iar pentru Sample C Major observam ca avem un număr mai redus de distorsiuni prin uniformitatea mai ridicată a waveform-ului fata de versiunea normală.



Evaluare calitativă

Pentru a evalua calitativ rezultatele am importat în GarageBand fișierele obținute și am rulat succesiv fișierele audio și am obținut următoarele concluzii:

Sample 1 - Pentru acest sample rezultatul generat de versiunea din articol rezulta într-un sunet mai clar pentru cifre față de versiunea normală, dar în situația în care există un ecou în semnal varianta din articol va amplifica acel efect în lipsa altor prelucrări.

Sample 2 - Acest sample este cel mai dificil pentru ambele vocodere din cele 4 sample-uri cu vocea mea pentru că zgomotul este mai ridicat după procesare, volumul pare mai scăzut pentru ambele, dar în ciuda acestor aspecte versiunea din articol este clar cea preferată. De menționat este faptul că înregistrarea inițială aplica filtrul de tip deeper voice asupra întregului input ceea ce generează mult zgomot de fundal și acest aspect s-a reflectat în scorul negativ al raportului SNR.

Sample 5 - Pentru acest sample diferența nu este una foarte sesizabilă.

Sample C Major - Acest sample prezintă cea mai mare diferență față de varianta normală, sunetul este foarte clar, vocalele sunt izolate și se apropie în calitate de sample-ul original, iar acest aspect este posibil influențat de calitatea bună a sample-ului inițial care nu prezintă zgomot de fundal inițial precum sample-urile generate de mine.

Concluzii

Metoda propusă în articolul ales este o soluție foarte eficientă pentru implementarea unui phase vocoder pentru diferența în implementare nu este deloc mare față de un phase vocoder, deci nu presupune un trade-off în acest sens și rezultatele calitative, dar și cantitative pentru anumite situații susțin asta.

Un aspect de menționat este faptul că metoda propusă rezolvă doar problema de phase unwrapping și nu abordează problema de zgomot de fundal, astfel se explică faptul că în cazul sample-urilor generate de mine care conțin și ecou, dar și zgomot de fundal rezultatele pentru evaluarea cantitativă nu au fost așa de ușor de observat inițial precum în cazul sample-ului C Major care are zgomot de fundal minim și astfel reprezintă un caz ideal de aplicare.