Does vowel recognition relate to pitch?

Dieter Maurer, Christian d'Heureuse, Heidy Leemann Institute for the Performing Arts and Film IPF, Zurich University of the Arts ZHdK, Switzerland Inventec Informatik AG, Winterthur, Switzerland

> Acoustics Virtually Everywhere The 179th Meeting of the Acoustical Society of America 7-11 December 2020

> > ☑ Abstract in JASA
> >
> >
> > ☑ Study (PDF online)
> >
> >
> > ☑ Research homepage
> >
> >
> > ☑ Contact: dieter.maurer@zhdk.ch
> >
> >
> > Version 2021-02-07

Abstract

It was shown that, for vowel sounds, the spectrum relates to fundamental frequency (f_o) and the spectral envelope is ambiguous, often representing two or three different vowel qualities if f_o is varied substantially. Thus, from a speech perception perspective, vowel quality is indicated to relate to the pitch of a vowel sound. In this contribution, a concept is outlined addressing the experimental question of the relationship of vowel recognition and pitch. According to this concept, two or three harmonics are extracted near the spectral peaks of a natural vowel sound, representing assumed F1-F2 or F1-F2-F3, and the lowest harmonic(s) below the first spectral peak are added. Based on a single extracted harmonics pattern, a series of sounds are synthesised stepwise attenuating the levels of the lowest harmonics with the aim to effect a low-to-high transition of the highest common harmonics factor and pitch. Pitch level and vowel recognition of the sounds is then investigated by means of a listening test. – Results of a first pilot test based on sounds of mid-closed vowels /e, ø, o/ confirmed the vowel–pitch relation hypothesis and also revealed cases of sounds with double-pitch and double-vowel recognition.

BACKGROUND

The vowel spectrum of natural sounds relates to fundamental frequency (f_0) .

- ☑ Examples
- **References**

The spectral peak patterns, the estimated formant patterns and the spectral envelopes of vowel sounds are ambiguous, often representing two or three different vowel qualities if f_0 is varied.

- Z Examples
- ☑ References

In vowel synthesis based on series of equal amplitude harmonics, the recognised vowel quality can be changed by either altering the frequency of a single lower amplitude harmonic or only changing the frequency spacing of the higher harmonics. In both cases, the change in vowel quality is indicated to relate to the highest common factor of the harmonics and to pitch.

- **Z** Examples
- ☑ References

The recognition of sinewave vowel sounds is also indicated to relate to pitch.

- Z Examples
- ☑ References

The relation between recognised vowel quality and the vowel spectrum is nonuniform.

- 🗹 Examples
- References

Therefore, from a speech perception perspective, two indications exist:

The recognised vowel quality relates to the pitch of a vowel sound.

This relation is nonuniform.

Here, the first aspect of the pitch-relation is addressed.

GENERAL APPROACH AND PROCEDURE

Creation of a model synthesis experiment

The following model synthesis experiment was designed in order to address the vowel–pitch relation hypothesis.

Firstly, natural vowel sounds are to be selected for which

- the vowel-related spectral peaks are represented by dominant or prominent harmonics, D1– D2 (back vowels and /a/) or D1–D2–D3 (front vowels);
- these harmonics relate to formant patterns F1-F2 or F1-F2-F3 as commonly estimated for the sounds in question;
- the first dominant or prominent harmonic D1 is above the fundamental H1 (D1 > H1); the higher dominant or prominent harmonic(s) are integer multiples of the first dominant or prominent harmonic.

Secondly, the dominant or prominent harmonics and also the low harmonic(s) below *D*1 are to be extracted.

Thirdly, harmonic synthesis is to be applied with

- the patterns of the dominant or prominent harmonics keept unchanged;
- the low harmonic(s) below *D*1 stepwise attenuating until this harmonic is or these harmonics are deleted.

Fourthly, vowel and pitch recognition is to be tested.

For illustration, see below, method section, Figures 1 to 3.

Special character of the synthesised sounds

Sounds that are synthesised on the basis of these patterns of dominant or prominent harmonics (hereafter *D*-patterns) have special characteristics:

- The highest common factor HCF of the pattern including the harmonic(s) below *D*1 is lower than HCF of the pattern excluding the harmonic(s) below *D*1;
- therefore, the periodicity of the two types of sounds differ, affecting pitch perception and recognition;
- however, the spectral energy maxima are kept unchanged for both types of sounds, and if formant patterns are estimated, they also will prove to be unchanged;
- there is no existing concept of spectral envelope estimation accounting for sounds of this kind and their spectral representation of vowel quality.

Dissociation of fundamental frequency and pitch

In series of synthesised sounds of this type related to a single natural reference sound, above all in the course of the attenuation of the low harmonic(s) below D1, two concurrent HCF and, accordingly, **two concurrent sound periodicities occur to which perception and recognition can refer to:**

- HCF including the lower harmonic(s)
- HCF of the dominant harmonics only

However, only one fundamental frequency is usually attributed to a single sound.

Expectation

Cases of synthesised sounds are expected to occur for which

- the pitch of the sounds produced with both the dominant or prominent harmonics and the low harmonic(s) below *D*1 is recognised as lower than the pitch of the sounds produced with the dominant or prominent harmonics only;
- in parallel, the recognised vowel quality of the first sounds mentioned differs from the quality of the second sounds;
- above all in the course of attenuating the lower harmonic(s) below *D*1, two pitches and/or two vowels may be recognised.

General Procedure

In general, the following procedure is proposed.

Create a large sample of vowel sounds: To find sounds fulfilling the above conditions of dominant or prominent harmonics being associated with spectral peaks, the creation of a large sample of vowel sounds is needed. (For our basis of the Zurich Corpus, see next text box.)

Select natural vowel sounds with a specific spectral peak pattern: In the sample created, search for natural sounds for which the vowel-related spectral peaks are associated with dominant or prominent harmonics D1-D2 (sounds of back vowels and of /a/) or D1-D2-D3 (sounds of front vowels), D1 being above H1, and D2 or D2 and D3 being integer multiples of D1.

Extract the *D*-pattern of a sound spectrum and add lower harmonics < D1: Extract the dominant or prominent harmonics and also the low harmonic(s) below *D*1 from a selected natural sound.

Perform harmonic synthesis with stepwise attenuating the low harmonic < D1: Based on the extracted *D*-pattern and the low harmonic(s), perform harmonic synthesis and attenuate the level(s) of the low harmonic(s) step by step.

Perform vowel and pitch recognition tests: Test vowel and pitch recognition for the synthesised sounds, either separately or jointly.

Note

The spectral representation of vowel quality is nonuniform (see Background). Therefore, results may be dependent on vowel qualities and on f_0 levels of the natural reference sounds investigated as well as on respective extracted *D*-patterns (see also the Discussion).

PILOT STUDY - METHOD

On the basis of the general experimental design, a pilot study was conducted applying the following method.

Sounds investigated

On the basis of the Zurich Corpus, sounds of the three long Standard German vowels /e, ϕ , ϕ / produced by female speakers at f_0 of c. 220–250 Hz were investigated.

Visit the Zurich Corpus

Sound selection

Examples of natural vowel sounds produced as monophthongs at f_0 of c. 220–250 Hz were selected which fulfilled three conditions:

- The vowel-related spectral peaks were associated with dominant harmonics; however, for sounds of /o/, a second spectral peak indicated by only a prominent harmonic was also accepted
- The first dominant harmonic is not the fundamental *H*1
- The higher dominant (or prominent) harmonic(s) are integer multiples of the low dominant harmonic



Figure 1: Vowel spectrum of a natural sound of /e/ produced in V context at fo of c. 200 Hz.

The first three spectral peaks are associated with dominant harmonics D1-D2-D3 = c. 440-2640-3080 Hz.

Play black this sound

For each of the three vowels /e, ϕ , o/, two natural sounds were selected fulfilling these conditions. Two peaks for /o/ and three peaks for /e, ϕ / were considered vowel related.

Extraction of harmonic patterns

Subsequently, dominant (or prominent) harmonics D1-D2 or D1-D2-D3 were extracted (including their actual levels and dynamic contours), and the first harmonic H1 of the natural sound below D1 was added.



Figure 2: Natural sound of /e/ produced at f_0 of c. 220 Hz, with dominant harmonics D1-D2-D3 = c. 440-2640-3080 Hz, and extracted pattern D1-D2-D3 with added H1 = c. 220 Hz.

- **Play black this sound pair**
- \square See all six pairs of natural vowel sounds and extracted harmonics
- Deption: See also formant pattern ambiguity for the six natural reference sounds

Harmonic synthesis

For all six configuration of extracted harmonics, vowel sounds were synthesised with step by step attenuation of H1.



Figure 3: Illustration of H1 attenuation, with the H1 levels shown = 0-10-20-100 dB.

Eight attenuation levels of H1 were investigated: 0/-5/-10/-15/-20/-30/-50/-100 dB

Play black these four sounds

Tool

Extraction of the harmonics and resynthesis was conducted using the HarmSyn harmonic synthesiser.

Note that, using this tool, the dynamic course of the harmonics can be analysed and resynthesised including frequency and level variation of the harmonics over time and including harmonic level attenuation.

☑ Reference

Sound sample created

For each single harmonic configuration, eight synthesised sounds were produced related to the eight attenuation levels of H1 mentioned of: 0/-5/-10/-15/-20/-30/-50/-100 dB

A sample of 48 single sounds in total was created for further investigation. \square See all sound series, and listen to the sounds

Vowel and pitch recognition tests

Five experienced listeners (professionally trained singers) participated in the below listening tests. (Note that these listeners were also involved in a large number of listening tests for the Zurich Corpus; see above for more information on the Corpus.)

Testing vowel and pitch recognition of opposing sounds (without and with full attenuation of H1)

Vowel and pitch recognition of sounds with unchanged H1 and with fully deleted H1 were tested in three subtests.

In these three subsets, for each test item, the listeners were asked to assign only one vowel (prominent or dominant vowel) or only one pitch (prominent or dominant pitch).

In the first subtest, each test item consisted of a single sound and the listeners were asked to assign a vowel quality (forced choice, all long Standard German vowels and schwa).

In the second and third subtest, each test item consisted of two synthesised sounds (with a 1 sec pause in between): a sound with an unchanged H1 configuration versus a second sound with fully deleted H1, and vice versa (AB and BA order). – In the second subtest, the listeners were asked to assign a vowel quality to the second sound presented (forced choice, see above). – In the third subtest, the listeners were asked to label the pitch difference between the two sounds as either "falling", "rising" or "flat" ("flat"=no pronounced pitch difference).

Testing vowel and pitch recognition of opposing and transitional sounds

Vowel and pitch recognition of all sounds was tested in three subtests.

In these three subsets, for each test item, the listeners were again asked to assign only one vowel (prominent or dominant vowel) or only one pitch (prominent or dominant pitch).

In the first subtest, each test item consisted of a single sound with an attenuated *H1* configuration and the listeners were asked to assign a single vowel quality (forced choice; see above).

In the second and third subtest, each test item consisted of two synthesised sounds (with a 1 sec pause in between): a sound with an unchanged *H1* configuration versus a second sound with attenuated or deleted *H1* configuration or, inversely, a sound with deleted *H1* configuration versus a second sound with unchanged or attenuated configuration (AB and BA order). – In the second subtest, the listeners were asked to assign a vowel quality to the second sound presented (forced choice, see above). – In the third subtest, the listeners were asked to label the pitch difference between the two sounds as either "falling", "rising" or "flat" ("flat"=no pronounced pitch difference).

Testing double-vowel and double-pitch recognition of opposing and transitional sounds

In the fourth subtest, each test item consisted of a single sound and the listeners were asked whether they recognise one or two vowel qualities.

Likewise, in the fifth subtest, each test item consisted of a single sound and the listeners were asked whether they recognise one or two pitches.

Vowel and pitch recognition of opposing sounds (without and with full attenuation of H1)

Table 1. Recognition results of pitch and vowel quality tested for opposing sounds of a D-pattern without and with full attenuation of H1 (vowel recognition = 2 subtests with 10 identifications, pitch recognition = 1 test with 5 identifications; see Method section). Columns 1–3 = sound production (S=series number, fo=fundamental frequency in Hz, V=vowel quality intended). Columns 4–8 = harmonic configuration of synthesis (H1=fundamental, AH1=attenuation of H1, D(i)=dominant harmonics, approximate values in Hz). Columns 9–14 = recognised vowel qualities (total of both tests). Columns 15–16 = recognised pitch level for the comparison of sounds related to a D-pattern without and with full attenuation of H1. – Recognition rates ≥ 80% of mid-closed vowels and/or of low pitch are highlighted in blue, and recognition rates ≥ 80% of closed vowels and/or of high pitch are highlighted in red.

Natural				Harmonic							Vowel								Pitch	
S f _o V				H1 AH1 D1 D2 D3													low high			
1	220	e		220	0		440	2640	3080				10						5	
			Ц	-	-100					_							10			5
2	250	_		250	0		500	2500	3000			1	9						5	
		e	Ц	١	-100		500	2500	3000				1				9			5
3	220	-	Ħ	220	0		440	1700	0040			10							5	
		ø		-	-100		440	1760	2640							10				5
4	220		Η	220	0		4.40	4700	0040			10							5	
		Ø		-	-100		440	1760	2640			1				9				5
5	220			220	0			880	-		10								5	
		0		-	-100		440								10					5
				000	-					-									-	
6	220	0		220	0		440	880	_		8				2				5	
		5		-	-100		. 10	000							10					5

The results of the vowel and pitch recognition tests for opposing sounds of a D-pattern without and with full attenuation of H1 (test details, see Method section) showed that

- no attenuation of *H*1 was associated with a lower pitch level and a mid-closed vowel quality;
- full attenuation of *H*1 (*H*1 deleted) was associated with a higher pitch level and a closed vowel quality.

See this table in PDF format

 \square See all opposing sounds, and listen to the sounds

Vowel and pitch recognition of opposing and transitional sounds

Table 2. Recognition results of pitch and vowel quality tested for all sounds of a D-pattern with stepwise attenuation of H1 from 0 db to -100db (vowel recognition = 2 tests with 15 identifications, pitch recognition = 1 test with 10 identifications; double-vowel recognition = 1 test with 5 identifications; double-pitch recognition = 1 test with 5 identifications; values marked with "**" are taken from Table 2; details see Method). Columns 1–3 = sound production (S=series number, fo=fundamental frequency in Hz, V=vowel quality intended). Columns 4–8 = harmonic configuration of synthesis (H1=fundamental, AH1=attenuation of H1, D(i)=dominant harmonics, approximate values in Hz). Columns 9–14 = recognised vowel qualities (total of both tests). Columns 15–16 = recognised pitch level for the comparison of sounds related to a D-pattern with and without full attenuation of H1. – Recognition rates \ge 80% of mid-closed vowels and/or of low pitch are highlighted in blue, and recognition rates \ge 80% of closed vowels and/or of high pitch are highlighted in red.

Nat	ural so	und			Vowel recognition								Pitch recognition					
S	f.	V	H1	AH1	D1	D2	D3		0	ø	e		u	У	i		low	high
1			220	0							13	Π			2		10	
			(200)	-5							10	П			5		9	1
	220	е	(200)	-10	440		3080				11	Ц			4		9	1
			(200)	-15		2640					6	Н			9		9	1
			(200)	-20							5	Н			10		7	3
			(200)	-50								Н			15			10
			-	-100								H			15			10
=			220	0	 	;==				1	13	Ħ			1		10	
			(200)	-5						2	13	H			<u> </u>		10	
			(200)	-10		2640	3080			1	13	H			1		10	
	200		(200)	-15							13	Π			2		7	3
1	220	e	(200)	-20	1 ***						9				6		6	4
			(200)	-30							9	Ц		2	4		6	4
			(200)	-50							3	Ц			12			10
			-	-100							3	Ц			12	1.1		10
3	220	ø	220	0		2640	3080			15		Ц					10	
			(200)	-5						15		Ц					10	
			(200)	-10						14		Н		1			10	
			(200)	-15	440					13		Н		2			10	2
			(200)	-20						8		H		7			5	5
			(200)	-50						-		H		15		1		10
			-	-100								Π		15				10
	220	ø	220	0	0 -5 -10 -15 20 30 50 100	2640	3080			15		Π					10	
			(200)	-5						15		H				1	10	
			(200)	-10						14				1		1	10	
4			(200)	-15						11		Ц		4			8	2
			(200)	-20						9		Ц		6			6	4
			(200)	-30						7		Н		8			5	5
			(200)	-50						1		Н		14		1		10
		\models		+				17			片		14			15	10	
5			220 0					15			Н					10		
			(200)	-5		2640			14			Н	1				0	1
			(200)	-15					10			H	5			1	7	3
	220	0	(200)	-20	440		3080		7			H	8			1 .	6	4
			(200)	-30					6				9				4	6
			(200)	-50									15					10
			- [-100									15			-		10
6			220	0		2640	3080		12				3				10	
	220		(200)	-5	440				12			Ц	3				10	
			(200)	-10					12			Ц	3				9	1
		0	(200)	-15					10			Н	5				9	1
			(200)	-20					8			Н	7				8	2
			(200)	-50					-			Н	15				0	10
			-	-100								H	15					10

See this table in PDF format

Referring to the columns 9-14 and 15-16 in Table 2, the results of the vowel and pitch recognition tests for all sound series of a *D*-pattern with stepwise attenuation of *H*1 from 0 to -100 dB showed that

• marked between-speaker differences occurred for both vowel and pitch recognitions in the transitions from the oppositions of mid-closed vowels and lower pitch (no attenuation of *H*1) to closed vowels and higher pitch (full attenuation of *H*1).

See all sound series, and listen to the sounds (link is also given in the Method section)

Double vowel and double pitch recognition

Referring to the Columns 17 and 18 in Table 3 (extension of Table 2, see link below), the results of testing double vowel and double pitch recognition showed that

 numerous cases of sounds occurred for which two vowel qualities and/or two pitches were recognised.

See Table 3 in PDF format (extension of Table 2)

Further details of the results

Details of all pitch and vowel recognition results are given in separate documents.

 \square See the details of the recognition results

DISCUSSION

In the present experiment, two dominant or prominent harmonics D1-D2 (sounds of /o/) or three dominant harmonics D1-D2-D3 (sounds of /e, ϕ /) of the spectra of natural mid-closed vowel sounds produced at f_0 in the range of c. 220–250 Hz were extracted, the dominant or prominent harmonics corresponding to formant frequency patterns F1-F2 or F1-F2-F3 commonly estimated for these natural sounds. Then, the fundamental H1 of the sound spectrum in question (the harmonic below the first spectral peak) was added. Based on this type of harmonic patterns, series of sounds were synthesised stepwise attenuating the levels of H1 with the aim to trigger a low-to-high transition of the highest common harmonics factor and of the recognised pitch. Finally, vowel and pitch recognition of the sounds was investigated.

Firstly, **the results support the vowel–pitch relation hypothesis**, here in terms of a pronounced general tendency for a mid- closed to closed vowel quality shift associated with an increase of the recognised pitch.

Secondly, and most importantly, **cases of sounds with double-vowel and/or double-pitch recognition occurred**, a manifestation which strongly underpins the vowel–pitch relation hypothesis.

However, within the transitional phase between the low-pitched sounds of mid-closed vowels and the high-pitched sounds of closed vowels, **the pitch and vowel shifts did not obey a strict parallelism**: Vowel shifts occurred before pitch shifts, and vice versa.

The new indication of vowel quality recognition as related to pitch offers a coherent explanation of earlier findings that the vowel spectrum of natural sounds is related to f_0 and that formant patterns and spectral shapes are ambiguous representations of vowel quality as well as of the results of the vowel synthesis experiments referred to in the Background section.

All of the above indications that vowel quality recognition is related to pitch stand in strong contrast to the claim of the prevailing acoustic theory of the vowel that either the formant pattern or the spectral envelope acoustically represents vowel quality (see Hillenbrand and Houde, 2003; Swanepoel et al., 2012): In the present experiment, estimated formant patterns do not differ for the vowel quality shifts of the sounds observed, and current concepts of spectral envelope estimation do not account for the harmonic configurations of the vowel spectra discussed here.

As one of the consequences, the spectrograms of two sounds recognised as two different vowels do not indicate this difference, as is illustrated here in the comparison of the first sound pair of /e/ and /i/ shown in Table 1.



Figure 4: Sound spectrogram of /e/ (left) and /i/ (right) of the sound pair shown in Series of Table 1.

In order to further investigate and clarify the role of pitch for vowel quality recognition and its impact on prevailing acoustic theory, further research on the matter is needed including the development of improved experimental designs.

With respect to such future research, special attention should be given to the fact that, in the experiments referred to, the general indication of a vowel–pitch relation proved to be robust but a strict parallelism vowel and pitch recognition was lacking. Special attention should also be given to the many indications of a nonuniform relation between spectral characteristics and vowel quality (see Background).

Concerning the nonuniform spectral representation of vowel quality, when designing the present experiment, we have experienced that the sound quality of this type of synthesis depends on the vowel quality and on f_0 of the natural reference sounds and that this dependence affects the inter-listener accordance or confusion of both vowel and pitch recognition. For similar experimentations, we recommend to first investigate natural reference sounds of mid-closed vowels produced at f_0 of 200–250 Hz because their lower vowel spectrum was found to strongly relate to an increase in f_0 and, consequently, single formant patterns as well as single spectral envelopes of these sounds are in most cases ambiguous in that they represent mid- closed and closed vowel qualities (see Background). Subsequently, the investigation may be extended to sounds of other vowel qualities produced at lower or higher f_0 .

AReferences

Miscellaneous / Corrections

ACKNOWLEDGEMENT

This work was supported by the Swiss National Science Foundation SNSF, Grant 100016_159350.

AUTHOR INFORMATION

Dieter Maurer, Prof. Dr. Institute for the Performing Arts and Film IPF, Zurich University of the Arts ZHdK, Switzerland; dieter.maurer@zhdk.ch

Christian d'Heureuse Inventec Informatik AG, Winterthur, Switzerland; chdh@inventec.ch

Heidy Suter Institute for the Performing Arts and Film IPF, Zurich University of the Arts ZHdK, Switzerland; heidy.suter@zhdk.ch