

The real and illusory in phoneme perception

Kei Omata¹⁾²⁾, Ken Mogi²⁾¹⁾

1) Tokyo Institute of Technology Department of computational Intelligence and Systems science
2) Sony computer science laboratories, Inc.

omata@bn.dis.titech.ac.jp

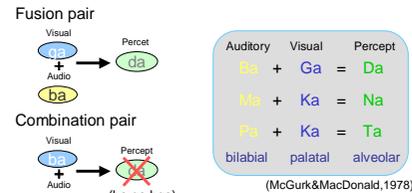
Abstract

Phonemes are building blocks of linguistic conscious perception. The McGurk effect suggests that visual articulation can affect the phoneme perception. Recently, brain imaging studies have suggested that the McGurk effect is due to activities in the posterior part of the left superior temporal sulcus and gyrus, where neurons responding to both auditory and visual information are found.

Here we formulate a model of the McGurk effect by means of the Self-organizing map (SOM). We examined whether putting the SOM in a series of learning process of phonemes would naturally induce a McGurk effect-like activities. We report the reconstruction of the McGurk effect in a manner consistent with the experimentally observed phenomenology of phoneme perception.

Our results suggest a mechanism underlying real and illusory perceptions in the context of learning.

Features of the McGurk effect



Auditory	Visual	Percept
ba	ga	da
ba	ka	na
ba	ka	ta

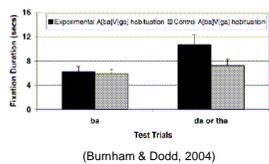
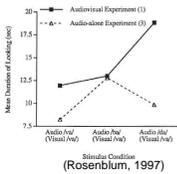
(McGurk&MacDonald, 1978)

The articulations of consonants are related with the McGurk effect.

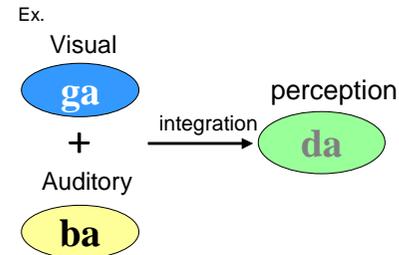
The fusion response occurs when the specific audio-visual pairs are presented.

Genetic or Learning?

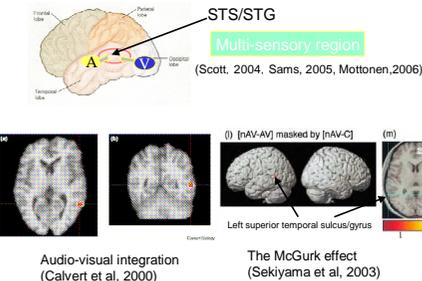
When does the McGurk effect appear?



McGurk effect (McGurk&MacDonald, 1976)



The McGurk effect in Brain Science

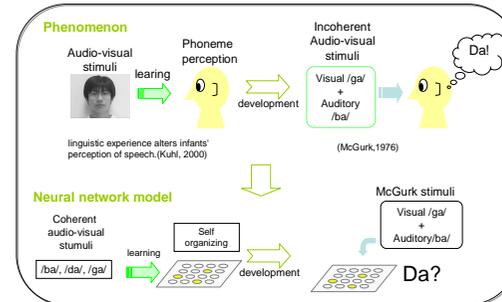


Is language experiences important for the McGurk effect?

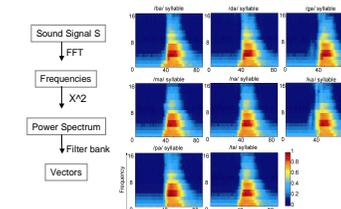
Auditory	Visual	Subjects	Auditory	Visual	Response Field	Combination	Other
ba/be	pa/pe	3-5yr (n=21)	19	0	83	0	0
		5-7yr (n=21)	34	0	64	0	0
		18-40yr (n=34)	7	0	96	0	0
pa/pe	ka/ke	3-5yr (n=21)	17	19	14	4	4
		5-7yr (n=21)	11	11	0	0	0
		18-40yr (n=34)	11	0	4	4	4
pa/pe	ka/ka	3-5yr (n=21)	24	0	32	0	24
		5-7yr (n=21)	24	0	30	0	0
		18-40yr (n=34)	39	0	30	0	0
ka/ka	pa/pe	3-5yr (n=21)	0	0	0	0	24
		5-7yr (n=21)	0	0	0	0	24
		18-40yr (n=34)	0	0	0	0	24

(McGurk & MacDonald, 1976)

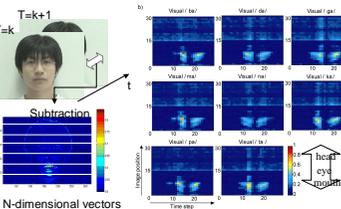
Background & Concept



The auditory vectors



The visual vectors



Results

generalization capability (3 syllables × 30 utterances)

The performance of classifier (%)	/b,d,g/	/m,n,k/	/p,t,k/
5-fold Cross Validation	96	98	94

5-fold Cross-validation
the ratios of discrimination in every group are more than 94%.

The McGurk effect (30 visual streams × 30 syllables)

McGurk effect (%)	/b,d,g/	/m,n,k/	/p,t,k/
Presentation Pair	/b/ /d/ /g/	/m/ /n/ /k/	/p/ /t/ /k/
Fusion	20.9 (46.9)	32.2 (34.4)	42.1 (9.6)
Combination	52.2 2.3 45.4	47.0 0.6 52.4	83.0 1.2 15.8

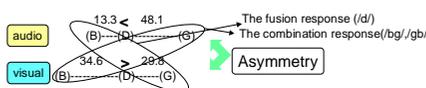
- The fusion pairs were perceived as the third phone (alveolar).
- The combination pairs were perceived as either the presented phonemes.

Analysis

similarity index: subtraction between the sample vectors of averaged syllables. Distance = sum of the absolute values of subtraction for each intensity between frames.

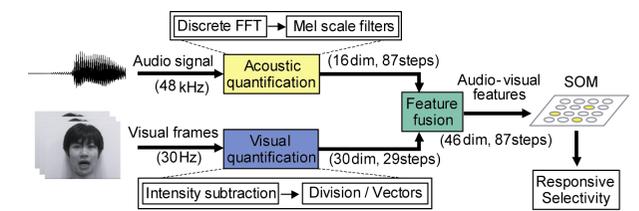
$$f(S_i, S_j) = \sum_{m,n} |x_{i,m} - x_{j,n}| \quad f(S_i, S_j) = f_a(S_i, S_j) + f_v(S_i, S_j)$$

Similarity index	/b,d,g/	/m,n,k/	/p,t,k/
(S _i , S _j)	(/b/, /d/), (/d/, /g/), (/g/, /b/)	(/m/, /n/), (/n/, /k/), (/k/, /m/)	(/p/, /t/), (/t/, /k/), (/k/, /p/)
f _a (S _i , S _j)	47.9 78.0 88.3	48.8 67.5 85.1	61.7 65.9 69.9
f _v (S _i , S _j)	13.3 < 48.1	47.2 14.5 < 42.4	42.4 25.2 < 45.3
f _v (S _i , S _j)	34.6 > 29.8	41.1 34.3 > 25.2	42.7 36.5 > 20.6

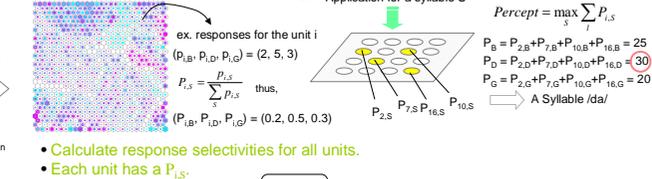


Method

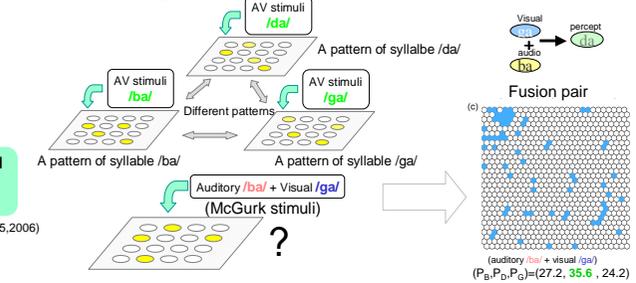
- The construction of the phoneme perception system using the Self organizing map ('SOM', Kohonen, 1995)
- Application of the McGurk stimuli to the maps



Response selectivity P



- Calculate response selectivities for all units.
- Each unit has a P_{i,S}.



Discussion/Conclusion

Stimuli

- There was an asymmetric relationship between the similarity of auditory features and that of visual features among the phonemes.
- The asymmetry relationships were found in all the three groups.

The existence of the asymmetric relationship between the audio and visual stimuli may induce the fusion responses of the McGurk effect when the particular audio-visual stimuli are presented.

System

- We hypothesized that time-synchronous audio-visual information is concatenated in the model.

The results suggested that the phoneme perception system requires a mechanism to combine auditory and visual information

Learning of audio-visual speech is an inducing factor of the McGurk effect.

Hypothesis: Learning of the audio-visual speech perception may cause the McGurk effect.

These results suggested that the McGurk effect occurred, and that the fusion pairs were similar to the alveolar consonants.

The McGurk effect may be based on the asymmetry of the similarities

Reference

McGurk, H. and MacDonald, J. 'Hearing lips and seeing voices', Nature, 264:746-748 1976
 Kohonen, T. 1995 Self-Organizing Maps, Springer-Verlag, Heidelberg