

움직이는 모양 연산자를 이용한 입술 명령어 인식

김선종*, 권혁송**¹, 김주만**²

A Recognition of Lip Commands using a Motion Shape Descriptor

Seon-Jong Kim*, Hyeog-Soong Kwon**¹, and Joo-Man Kim**²

이 과제는 부산대학교 기본연구지원사업(2년)에 의하여 연구되었음.

요 약

음성을 기반으로 하는 명령은 주변 환경이나 잡음에 따른 영향을 받는다. 본 논문에서는 소리가 아닌 입술의 모양으로 움직임을 추적하고, 주어진 명령어를 인식하는 SSD(Sequential Shape Descriptor) 연산자를 제안하였다. 이를 위해 먼저 입술의 모양을 표현하는 TCD(Triangular Centroid Distance) 연산자를 사용하고, 이들의 연속적인 움직임은 SSD 연산자로 표현하여 주어진 명령어를 인식하도록 하였다. SSD 연산자는 윤곽선으로 나타나는 입술 모양을 주파수 영역으로 표현할 수 있어서, 크기, 위치 및 회전에 무관하게 인식될 수 있다. 일반적으로 사용되는 5개의 기계 명령어를 이용하여 CNN 기반 딥러닝에 실험한 결과, 96.3%의 높은 인식률을 얻을 수 있었다. 따라서 소리에 민감한 환경인 경우, 제안된 SSD 연산자를 이용한 인식이 소리 대체 수단으로 사용될 수 있음을 확인할 수 있었다.

Abstract

Voice-based commands are affected by the surrounding environment or noise. In this paper, we proposed an SSD(sequential shape descriptor) operator that tracks motion by shape rather than the sound of lips and recognizes a given command. To do this, the TCD operator expressing the shape of the lips was used, and their continuous movement was made to recognize commands using the SSD operator. The SSD operator can express the shape of the lips as outlined in the frequency domain, so it can recognize regardless of size, position, and rotation. As a result of experimenting with CNN-based deep learning using five commonly used machine instructions, a high recognition rate of 96.3% was obtained. Therefore, in the case of an environment sensitive to sound, it was confirmed that the recognition using the proposed SSD operator can be used as a sound substitute.

Keywords

lip reading, lip commands, sequential shape descriptor, lip motion

* 부산대학교 IT응용공과 교수(교신저자)
- ORCID: <https://orcid.org/0000-0003-2070-290X>
** 부산대학교 IT응용공학과 교수
- ORCID¹: <https://orcid.org/0000-0003-1929-7855>
- ORCID²: <https://orcid.org/0000-0001-5525-3644>

· Received: Oct. 14, 2020, Revised: Nov. 23, 2020, Accepted: Nov. 26, 2020
· Corresponding Author: Seon-Jong Kim
Dept. of IT Engineering, Pusan National University
Tel.: +82-55-350-5413, Email: ksj329@pusan.ac.kr

I. 서 론

음성은 사람들의 의사소통을 위한 가장 효과적인 수단 중 하나이다. 또한 음성을 이용한 대화나 명령은 사람과 사람뿐만 아니라 사람과 기계간에도 적용될 수 있다. 인공지능과 더불어 음성인식에 대한 연구는 기계와 사람 간의 소통의 범위가 넓어짐에 따라 더욱 활발히 진행되고 있다[1][2].

음성 인식은 상업화가 될 정도로 그 기반 기술들이 많은 발전을 이루었지만, 주변 환경의 잡음에 민감한 특징으로 인해 기계와의 의사소통에 영향을 받고 있다. 이 영향에 대한 대체 방법으로 청각적 요소와 시각적 요소의 결합이 제안되고 있다. 이런 의사소통은 청각적 요소인 음성과 더불어 시각적 요소를 이용한다면 인식에 있어서 더 높은 정확도와 신뢰성이 보장될 수 있기 때문이다.

시각적 요소 중의 하나인 입술 모양에 대한 연구는 오래전부터 많은 연구가 진행되고 있다. Medioni 등[3]은 입술 특징 점(Landmarks)에 대한 연구로써 말하는 사람에 상관없이 주어지는 특징 점을 사용하여 입술의 형태 정보를 추적하는 방법을 제안하였다. 그 이후, 입술 특징 점들 사이 거리의 변화를 통해 움직임 분석하거나 입술 특징 점들이 가지는 복잡성을 분석하기 위한 딥러닝 네트워크를 통해 입술 영상을 학습하는 연구가 진행되고 있다[4]. 특히, 영상에서 얼굴과 눈의 관심영역 분류와 정보 분석을 기반으로 하는 인식 방법을 제안[5]한 바 있지만, 윤곽선을 이용한 모델[6]이 기존의 방법보다 정확성이 높다는 것을 보여주었다.

립 리딩 기술은 입술의 특징 점을 추출하는 것에 기반을 하고 있으며, 입술 영역 내의 특징을 추출하는 방법과 입 형태를 근사화 시킨 특징 벡터를 추출하는 방법 등으로 나누어진다. 첫 번째 연구는 입술 영역의 영상 값의 변화에 중점을 둔 연구[7]와 입술 영역 움직임에 특화시킨 Optical Flow 특징이 사용되는 연구[8]가 있다. 두 번째 연구는 입술 윤곽선의 변화를 근사 모델로 움직임을 추적하는 등[9]이 진행되고 있다. 최근, 딥러닝의 발전으로 입술 영역 및 특징 점 추출은 실시간으로 얻을 수 있을 만큼 발전되었다. 추출된 특징 점들을 기반으로 윤

곽선[10]이 추출되며, 이 특징 점들의 위치 변화는 윤곽선을 변화시킨다. 사용되는 음성 명령들은 각 음절에 따라 입술 모양의 연속적인 변화를 가지게 된다. 이 연속적인 변화는 다양한 방법을 통하여 인식될 수 있지만 최근에 CNN기반 네트워크가 유용하게 사용되고 있다. 따라서 윤곽선을 영상 형태로 바꾸어 입력 패턴을 만들면 기존의 네트워크를 이용할 수 있다[11][12].

TCD(Triangular Centroid Distance)는 윤곽선 형태를 크기, 위치, 회전에 무관하게 주파수 차원으로 해석할 수 있는 방법으로 제안되었다[13]. TCD는 무게 중심을 기준으로 형태를 이루는 모든 점과의 거리의 주파수 특성으로 표현될 수 있다.

본 논문에서는 TCD를 이용하여 연속적인 입술 형태를 표현할 수 있는 SSD 연산자를 제안한다. 제안된 연산자는 입술의 형태를 윤곽선으로 표현하고, 이를 주파수 영역으로 변환하여 크기, 위치 및 회전에 무관한 입술 형태로 표현될 수 있도록 한다. 주어진 얼굴 영상으로부터 입술을 찾고, 입술 영역에 대한 특징 점을 찾는다. TCD는 입술과 같이 고정된 수의 포인트를 가지는 경우에 더욱 유용하게 사용될 수 있다. 기계 명령에 대한 연속적인 입술 윤곽선 변화에 대한 SSD 연산자를 나타내며, 이들을 학습 패턴으로 하여 주어진 명령어를 인식하도록 한다. 인식은 Inception ResNet V2를 통하여 딥러닝을 이용한다.

본 논문의 구성은 다음과 같다. 2장에서는 관련 연구로서 입술에 대한 특징점을 추출 및 TCD에 대해 설명하고, 3장에서 SSD 연산자를 제안하고, 입술 윤곽선의 특징에 다른 파라미터를 추출한다. 그리고 4장에서는 기계 명령어에 대한 실험을 통해 제안된 연산자의 타당성을 확인하고, 마지막으로 5장에서 결론을 맺는다.

II. 관련 연구

2.1 입술 특징 점 추출

최근 인공지능이 발전되면서 영상 인식에 대한 많은 연구가 성과를 이루고 있다. 특히 얼굴 인식은

얼굴 특징 추출 등과 더불어 지속적인 발전이 되고 있다. 얼굴의 한 요소인 입술은 인공지능을 이용하여 특징 추출이 용이하게 되었다. 다양한 인공지능 네트워크가 제안되고 있으며, 이 중에서 Shuffle Net v2는 CNN 구조의 효율성 향상을 위해 제시되었으며, 입술 특징 추출에 유용하게 사용될 수 있다. 그림 1은 이 네트워크를 사용하여 입술의 특징 점을 추출한 결과이다. 이는 입술 형태를 초당 30 프레임 이상으로 처리될 수 있어서 실시간 처리 성능을 가지며, 정확성도 뛰어나다. 본 논문에서는 이를 이용하여 입술 특징 점을 추출한다.



그림 1. 인공지능을 이용한 입술 특징 점 추출
Fig. 1. A extraction of lip landmarks using the shuffle net v2

2.2 TCD

TCD는 윤곽선 형태의 물체를 정합하기 위한 방법이다. 이 방법은 부분적인 윤곽선 정합도 가능하며, 윤곽선의 부분적인 변화를 감지하는데 우수한 성능을 보여주고 있다.

TCD 알고리즘은 다음과 같다. 즉, 2차원 평면위의 점 $P_i = (x_i, y_i)$, ($i = 1, 2, \dots, N$)가 존재하는 경우, N 개의 점은 각각의 점을 시작점으로 하는 $(N, N-1)$ 개의 2차원 점들 간의 거리에 대한 행렬로 표현될 수 있다. 각 행에 대한 푸리에 변환하여 2차원의 주파수 변환 영상을 얻을 수 있다.

먼저, 시작점과의 거리는 시작점 P_i 와 특징 점 P_j 그리고 중심점 좌표 G 간의 거리를 통해서 얻는다. 따라서 거리 $Dist(P_i, g_j)$ 는 다음과 같다. 즉,

$$Dist(P_i, g_j) = \sqrt{(x_i - x_{gj})^2 + (y_i - y_{gj})^2} \quad (1)$$

이 된다. 이 때, $g_j = (x_{gj}, y_{gj})$ 는 $\Delta P_i, P_j, G$ 에 대해 구해진 중심 좌표점이고, G 는 식 (2)와 같이 윤곽선의 중심점이다. 즉,

$$x_G = \frac{1}{N} \sum_{i=1}^N x_i \quad (2)$$

$$y_G = \frac{1}{N} \sum_{i=1}^N y_i$$

이 된다. 따라서 G 는 모든 특징 점들의 전역 중심 점이고, g_j 는 특징 점 P_i 와 P_j 그리고 G 로 구성된 삼각형의 중심점인 국소 중심점이라 할 수 있다. 따라서 시작점 P_i 에 대한 거리함수 $D(P_i)$ 는 다음과 같이 주어진다. 즉,

$$D(P_i) = (Dist(P_i, g_{i1}), \dots, Dist(P_i, g_{ij}), \dots, Dist(P_i, g_{iN})) \quad (3)$$

이 된다. 주어진 $D(P_i)$ 를 크기에 대해 정규화시키면 크기에 무관한 패턴으로 만들 수 있다. 따라서 주어진 윤곽선 S 에 대한 $TCD(S)$ 는 모든 특징 점으로 확장하면 2차원 패턴이 되며, 이를 각각 각 행에 대한 푸리에 변환(STFT)을 수행하여 다음과 같은 주파수 특성을 얻을 수 있다.

$$TCD(S) = |STFT(D(P_1))| \quad (4)$$

$$\dots$$

$$|STFT(D(P_N))|$$

III. 제안된 입술 움직임 연산자

제안된 시스템은 그림 2에 도시한 바와 같이, 먼저 카메라로부터 받은 얼굴 영상에 대해 입술 윤곽선에 대한 특징 점을 추출한다. 추출된 윤곽선으로부터 주파수 특성인 TCD를 계산하고, 연속적인 움직임에 대한 TCD를 누적시킨 SSD 연산자 패턴을 만들고, SSD 패턴을 인공지능 딥러닝을 통해 인식하는 과정을 갖는다.

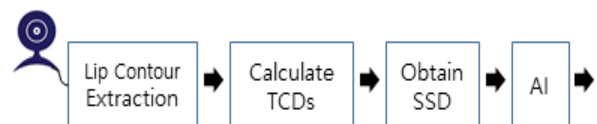


그림 2. 제안된 시스템 블록도
Fig. 2. Block diagram of the proposed system

3.1 입술 윤곽선에 대한 TCD 특성

입력 영상으로부터 얻어진 입술 특징 점은 다양한 윤곽선으로 표현될 수 있다. 전술한 TCD 연산자는 일정한 윤곽선에 대하여 그 크기와 시작점에 불변한 속성을 지니지만, 다른 순서에 따른 윤곽선은 다른 속성을 가지게 된다. 그림 3은 윤곽선 순서에 따라 달라지는 TCD 연산자의 결과를 도시한 것이다. 두 윤곽선은 동일한 지점에 대한 24개의 좌표를 가지며, 순서만 달리하여 만든 윤곽선이다. 결과를 살펴보면 순서에 따라 연산자 결과가 달라지는 것을 알 수 있다. 이렇게 순서를 달리함으로써 윤곽선의 특성을 더욱 극대화하거나 혹은 필요에 맞게 최적화를 시킬 수 있음을 확인하였다.

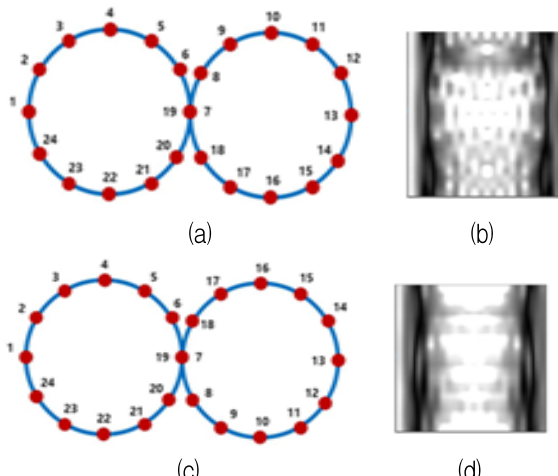


그림 3. (a) 순서 A와 (b) TCD 연산자 그리고 (c) 순서 B와 (d) TCD 연산자

Fig. 3. (a) Contour A and (b) TCD descriptor and (c) Contour B and (d) TCD descriptor

3.2 제안된 SSD 연산자

일반적으로 연속적인 발음은 입술의 모양이 각 소리에 맞는 형태로 움직인다. 따라서 입술 명령어는 입술 모양의 연속적인 변화를 계속해서 확인할 필요가 있다. 본 논문에서는 움직임을 파악하기 위해 특정 프레임 내의 입술 모양의 움직임을 표현할 수 있는 TCD를 연속적으로 배열하였다. 그림 4는 연속적으로 변하는 윤곽선의 모양에 대한 푸리에 이미지의 예이다.

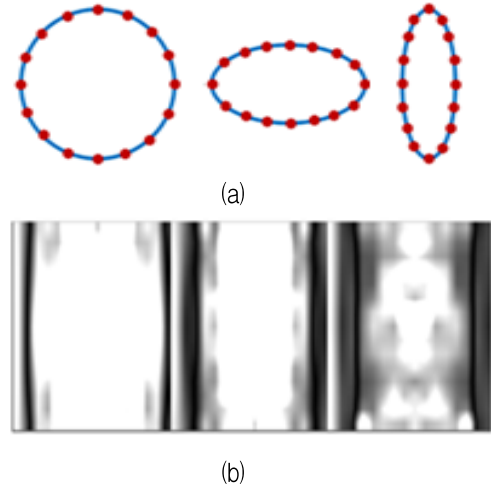


그림 4. (a) 3가지 윤곽선에 대한 (b) TCD 연산자
Fig. 4. (a) Three contours and (b) its TCD descriptor

시간에 따라 얻은 입술에 대한 윤곽선은 TCD 연산자로 변환되고, 주어진 기계 명령어 인식을 위해 연속적인 윤곽선의 변화를 확인하여야 한다. 그림 5와 같이 연속 발음하는 윤곽선 형태를 5x5 행렬 형태로 연속적으로 배열하여 얻을 수 있도록 하였다. 따라서 주어진 기계 명령어에 따라 길이를 제한함으로써 연속적인 입술 움직임에 대한 SSD 연산자를 통해 하나의 패턴으로 만들어질 수 있다. 패턴 인식을 위해 필요한 모양 변화에 대한 필요한 프레임은 다를 수 있지만, 기계 명령어에 필요한 5 종류의 대표적 명령어는 25 프레임으로 연속적인 모양정보를 얻을 수 있다. 그림 5는 기계 명령어의 하나인 Go 명령어에 대한 25프레임에 대한 SSD 연산자를 도시한 것이다.

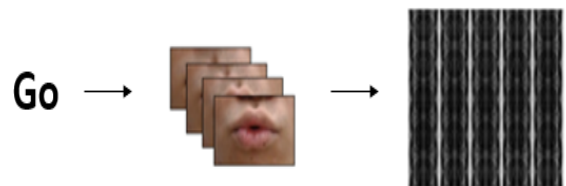


그림 5. Go 명령어에 대한 입술 모양과 SSD 연산자
Fig. 5. Lip shape and SSD descriptor for 'Go' command

IV. 실험 및 결과 고찰

본 논문에서는 제안된 SSD 연산자의 타당성을 조사하기 위해 기계와의 의사소통에 사용될 수 있는 5개의 대표적 명령어에 대해 성능 평가를 하였다.

5개의 명령어는 go, back, left, right 및 stop으로 한정하였다. 데이터 셋은 각 명령어를 발음하여 나타나는 입술 모양 변화로 25 프레임을 촬영하였다. 각 기계 명령어를 300번씩 발음하여 1,500개의 데이터 패턴을 모았다.

먼저, 입술 특징 점의 윤곽선에 대한 실험을 진행하였다. 그림 6은 그 결과를 도시한 것이다. 이때 사용된 발음은 ‘오’와 ‘이’로 하였다. 그림 6에서 (a)는 입술 특징 점이며, (b)는 순서, 그리고 (c)는 각 순서에 따른 두 발음간의 거리 차를 나타낸 것이다. 두 윤곽선의 거리가 클수록 좋은 연산자가 될 가능성이 있다. 그리고 순서 A가 순서 B와 순서 C보다 거리가 큰 것을 알 수 있다.

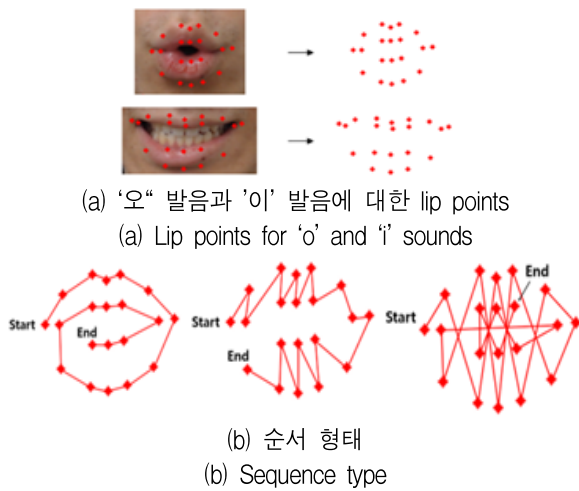


그림 6. 윤곽선 순서에 따른 발음 ‘오’와 ‘이’ 간의 거리
Fig. 6. The distance between the contours when the ‘o’ and ‘i’ are pronounced

따라서 입술 윤곽선은 순서 A와 같이 바깥에서 안쪽 입술로 들어오는 윤곽선을 이용하기로 한다. 다음은 주어진 윤곽선에 대한 기계 명령 인식 성능 실험을 수행하였다. 입술 명령어는 SSD 연산자를 입력 패턴으로 하는 Inception-ResNet-V2에 학습시켰으며, 성능을 평가하여 효용성을 검증하였다. 그림 7은 학습에 사용된 대표 명령어에 대한 SSD 연산자이다.

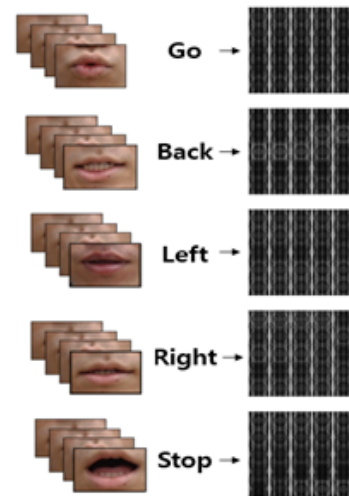


그림 7. 5가지 명령어에 대한 SSD 연산자
Fig. 7. SSD descriptor for 5 commands

성능 평가는 전체 패턴에 대한 학습 데이터 비율에 따라 성능으로 조사하였다. 인식 결과는 표 1에 나타내었다. 전체 명령 패턴 중에서 50%만 학습에 이용하고 나머지 50%는 성능평가에 사용된 경우에는 89.08%의 성능으로 나타났다. 표에서 알 수 있듯이 전체 데이터의 70% 이상을 학습에 사용된다면 인식률은 96.0% 이상인 것을 확인할 수 있었다.

표 1. 학습 패턴 비율에 따른 인식률
Table 1. Recognition rates according to the ratio of the learning data

Learned patterns(%)	Recognized patterns(%)	Recognition rates(%)
10	90	67.55
20	80	72.23
30	70	85.21
40	60	85.57
50	50	89.08
60	40	91.08
70	30	96.05
80	20	96.40
90	10	96.34

마지막으로 입술 모양이 아닌 소리에 대한 인식의 문제점을 확인하기 위하여 잡음의 영향을 조사하였다. 5개의 음성 명령어인 go, back, left, right 및 stop의 표준 발음은 옥스퍼드 사전에서 가져와 사용하였다. 그리고 음성 인식은 Google Cloud Platform Voice Recognition 모듈을 사용하였다. SNR 값을 기

준으로 잡음 비율을 높여가며 5개 명령어의 평균 인식률을 측정하였다. 그 결과를 그림 8에 도시하였다. 측정 결과, 약 15dB 이하에서 인식률이 급격하게 떨어짐을 확인하였다.

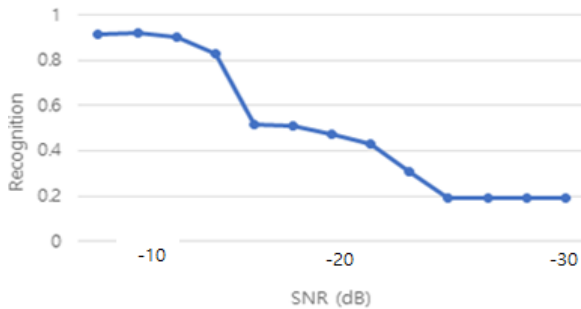


그림 8. 구글 음성 인식 성능

Fig. 8. Recognition rate of the Google voice experiments

이상에서와 같이 노이즈가 있는 환경에서 5개 기계 명령어에 대한 소리에 대한 인식에는 문제가 있다는 것을 확인하였으며, 이를 대체할 수 있는 입술 모양에 따른 인식에 대한 SSD 연산자의 타당성을 확인하였다. 따라서 입술 모양의 연속적인 변화 정보를 얻을 수 있는 SSD 연산자를 이용하면 음성 명령이 불안한 환경에 대한 대체 수단으로 사용될 수 있음을 확인하였다.

V. 결 론

본 논문에서는 연속적인 입술 움직임에 대한 윤곽선을 주파수 영역으로 변환되는 SSD 연산자를 제안하였다. 소리에 의한 인식은 잡음환경에 대해 인식에 한계가 있음을 확인하였고, 잡음의 정도에 따라 인식률이 급격하게 감소한다는 것을 확인하였다. 또한 입술의 연속적인 변화에 대한 제안한 SSD 연산자는 5개의 명령어에 대해 96% 이상의 인식률로 나타나 그 타당성을 확인되었다. 제안된 연산자는 주파수 특성을 갖는 패턴으로 나타나 크기, 위치, 회전에 무관한 인식될 수 있다. 특히, 제안된 연산자는 고정된 개수를 갖는 입술 특징 점에 유리할 뿐만 아니라 다양한 입술 비주얼 정보를 얻는데 효과적으로 이용될 수 있다. 또한 기존의 음성인식 기술과 결합할 경우, 더욱 높은 신뢰성을 얻을 수 있

을 것으로 판단되며, 이를 위해 더 많은 명령어와 더 많은 데이터로 성능을 평가할 예정이며, 또한 RNN 기반의 인공지능 모델에도 적용하여 성능을 평가하여야겠다.

References

- [1] S. Benkerzaz, Y. Elmir, and A. Dennai, "A Study on Automatic Speech Recognition", *Journal of Information Technology Review*, Vol. 10, No. 3 Aug. 2019.
- [2] A. Joshua and V. Sugumaran, "Speech Recognition For Humanoid Robot", *International Journal of Applied Engineering Research*, ISSN 0973-4562, Vol. 10, No. 68, pp. 57-60, Apr. 2015.
- [3] G. Medioni, J. M. Choi, M. Labeau, J. T. Leksut, and L. Meng, "3D Facial Landmark Tracking and Facial Expression Recognition", *Journal of Information and Communication Convergence Engineering*, Vol. 11, No. 3, pp. 207-215, Sep. 2015.
- [4] M. Faisal and S. Manzoor, "Deep Learning for Lip Reading using Audio-Visual Information for Urdu Language", <https://www.researchgate.net/publication/323217526>, [accessed : Feb. 2019].
- [5] Y. Lu and H. Li, "Automatic Lip-Reading System Based on Deep Convolutional Neural Network and Attention-Based Long Short-Term Memory", *Applied Sciences*, Vol. 9, No. 8, pp. 1599, Apr. 2019.
- [6] S. Christian, I. Sergey, V. Vincent, and A. Alexander, "Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning" *AAAI Conference on Artificial Intelligence(2017)*, [accessed : 29 Jan. 2020]
- [7] J. S. Lee and C. H. Park, "Automatic Lipreading Using Color Lip Images and Principal Component Analysis", *Journal of Information Processing Systems*, Vol. 15, No. 3, pp. 229-236, Jun. 2008.

- [8] Ayaz A. Shaikh, Dinesh K. Kumar, Wai C. Yau, M. Z. Che Azemin, and Jayavardhana Gubbi, "Lip reading using optical flow and support vector machines", 2010 3rd International Congress on Image and Signal Processing, Yantai, China, Vol. 1, pp. 327-330, Oct. 2010.
- [9] M. Z. Ibrahim and D. J. Mulvaney, "Robust geometrical-based lip-reading using Hidden Markov models", EUROCON, 2013 IEEE, Zagreb, Croatia, pp. 2011-2016, Jul. 2013.
- [10] S. W. Jang, S. Khanam, and W. J. Park, "Image Retrieval Integrating Interior and Contour Descriptors", Journal of Korean Institute of Information Technology, Vol. 10, No. 1, pp. 209-216, Jan. 2012.
- [11] D. S. Kim, Y. I. Kim, and S. J. Kim "A New Region-centralized Shape Descriptor for Character Representation", Journal of Korean Institute of Information Technology, Vol. 6, No. 1, pp. 87-95, Feb. 2008.
- [12] S. H. Cho and S. J. Kim, "A signal analysis for feature points tracking of lip movements", IEEE 4th International Conference on Signal and Image Processing, Wuxi, China, pp. 241-144, Jul. 2019.
- [13] C. Yang, H. Wei, and Q. Yu, "A novel method for 2D nonrigid partial shape matching", Neurocomputing, Vol. 275, pp. 1160-1176, ISSN 0925-2312, Jan. 2018.

권혁승 (Hyeog-Soong Kwon)



1996년 4월 ~ 현재 : 부산대학교
IT응용공학과 정교수
관심분야 : 생체신호, 의료정보,
바이오텔레미터링, CDMA,
유성버스트통신

김주만 (Joo-Man Kim)



1984년 : 숭실대학교
전산학과(공학사)
2003년 : 충남대학교
컴퓨터공학과(공학박사)
1985년 ~ 2000년 : ETRI
책임연구원(운영체제연구팀장)
2001년 ~ 현재 : 부산대

IT응용공학과 교수.

관심분야 : 임베디드 시스템, Real-time OS.

분산병렬처리, 고장허용시스템

저자소개

김선종 (Seon-Jong Kim)



1996년 8월 : 경북대학교
전자공학과(공학박사)
1995년 2월 ~ 1997년 2월 :
순천제일대학 전임강사
1997년 3월 ~ 현재 : 부산대학교
IT응용공학과 교수
관심분야 : 신호 및 영상처리,

머신/딥러닝, VR/AR, 스마트 카메라 등임.