

청각장애인을 위한 사운드 이벤트 검출 기반 홈 모니터링 시스템

Home monitoring system based on sound event detection for the hard-of-hearing

김지연,¹ 신승수,¹ 김형국[†]

(Gee Yeun Kim,¹ Seung-Su Shin,¹ and Hyoung-Gook Kim^{1†})

¹광운대학교 전자융합공학과

(Received February 21, 2019; revised April 19, 2019; accepted April 30, 2019)

초 록: 본 논문에서는 청각장애인을 위해 양방향 게이트 순환 신경망을 이용한 사운드 이벤트 검출 기반의 홈 모니터링 시스템을 제안한다. 제안된 시스템에서는 우선적으로 효과적인 사운드 이벤트 검출을 위해 패킷손실 은닉을 이용하여 무선 센서 네트워크로 인해 손실된 신호를 복원하고, 멀티채널 상호 상관관계 계수를 이용하여 신뢰할 수 있는 채널을 선택한다. 선택된 채널의 사운드는 이벤트 검출을 위해 두 개의 오디오 채널을 사용하는 양방향 게이트 순환 신경망에 적용된다. 검출된 사운드 이벤트는 텍스트로 변환되며, 이와 함께 하모닉/퍼커시브 음원 분리 방식을 통해 햅틱 신호로 변환되어 청각장애인에게 제공된다. 실험결과는 제안한 사운드 검출기반의 성능이 기존 방식보다 더 우수하다는 것과 음원 분리 방식을 통해 사운드를 세밀한 햅틱 신호로 표현할 수 있음을 보인다.

핵심용어: 홈 모니터링, 사운드 이벤트 검출, 양방향 게이트 순환 신경망, 사운드 햅틱 변환

ABSTRACT: In this paper, we propose a home monitoring system using sound event detection based on a bidirectional gated recurrent neural network for the hard-of-hearing. First, in the proposed system, packet loss concealment is used to recover a lost signal captured through wireless sensor networks, and reliable channels are selected using multi-channel cross correlation coefficient for effective sound event detection. The detected sound event is converted into the text and haptic signal through a harmonic/percussive sound source separation method to be provided to hearing impaired people. Experimental results show that the performance of the proposed sound event detection method is superior to the conventional methods and the sound can be expressed into detailed haptic signal using the source separation.

Keywords: Home monitoring, Sound event detection, Bidirectional gated recurrent neural network, Sound-to-haptic conversion

PACS numbers: 43.60.Bf, 43.60.Vx

1. 서 론

사운드는 주변 환경 및 상황 등의 중요한 정보를 포함하고 있어 사람의 사회적 활동을 이해하고, 상황 맥락을 묘사하기 위한 수단으로 사용될 수 있다.

하지만 청각장애인은 사운드를 통한 정보 취득이 어렵기 때문에 사각지대에 있는 위험 요소로부터 발생하는 사운드를 통해 상황을 이해하고, 즉각적으로 대처하는 것이 어렵다. 따라서 사운드 이벤트를 검출을 통해 사운드 정보를 다른 감각 신호로 변환하여 청각장애를 갖고 있는 사용자에게 전달하는 접근 방식이 필요하다.

최근 사운드 이벤트 검출은 CNN(Convolutional Neu-

[†]Corresponding author: Hyoung-Gook Kim (hkim@kw.ac.kr)
Department of Electronics Convergence Engineering, Kwangwoon University, 20 Gwangun-ro, Nowon-gu, Seoul 01897, Republic of Korea
(Tel: 82-2-940-5574, Fax: 82-2-913-5006)

ral Network), RNN(Recurrent Neural Network), LSTM (Long Short Term Memory)와 같이 다양한 구조의 심층 신경망에 적용되어 우수한 성과를 보이고 있다. 그중 LSTM-RNN은 시계열 데이터의 장기 의존성 문제를 해결함으로써 시퀀스 데이터의 인식률을 현저하게 개선시켜 오고 있다. 최근에는 LSTM-RNN보다 더 단순한 구조를 가지면서 유사한 성능을 제공하는 게이트 순환 신경망(Gated Recurrent Neural Network, GRNN)이 개발되었고, 이를 다양한 연구 분야에 적용하여 우수성을 입증하고 있다.^[1]

이에 본 논문에서는 양방향 게이트 순환 신경망(Bidirectional GRNN, BGRNN)을 이용하여 사운드 이벤트를 검출하고, 검출된 사운드 이벤트를 텍스트 및 햅틱 진동 신호로 변환하여 청각장애인에게 제공하는 홈 모니터링 시스템을 제안한다.

본 논문의 구성은 다음과 같다. II장에서는 제안하는 홈 모니터링 시스템에 대해 설명하고, III장에서는 실험결과를 제시한다. 마지막으로 IV장에서는 결론과 향후 연구를 설명한다.

II. 본 문

Fig. 1은 본 논문에서 제안하는 홈 모니터링 시스템의 전체 구조도이다.

제안하는 시스템은 싱크 기반의 신호 추정, 신뢰할 수 있는 채널 선택, 사운드 이벤트 검출, 사운드 이벤트의 텍스트 및 햅틱 변환으로 구성되어 있다. 먼저 무선 음향 센서를 통해 수집된 사운드는 패킷 형태로 인코딩되어 싱크로 전달된다. 싱크 기반 신호 추정에서는 패킷 신호를 수신하여 디코딩한 후 무선 센서 네트워크(Wireless Sensor Networks, WSNs)의 멀티 홉 통신으로 인해 패킷이 손실되었을 경우 패킷을 복원한다. 그다음 연산 효율과 사운드 이벤트 검출(Sound Event Detection, SED) 성능의 향상을 위해 멀티채널 상호 상관관계 계수를 이용하여 신뢰할 수 있는 두 개 채널을 선택한다. 선택된 채널의 사운드는 BGRNN기반의 SED와 하모닉/퍼커시브 음원 분리 방식을 이용한 사운드 햅틱 변환에 적용된다. IP 네트워크를 통해 사용자 디바이스로 전달된 사운드 이벤트는 텍스트와 햅틱 진동 신호로 변환된다.

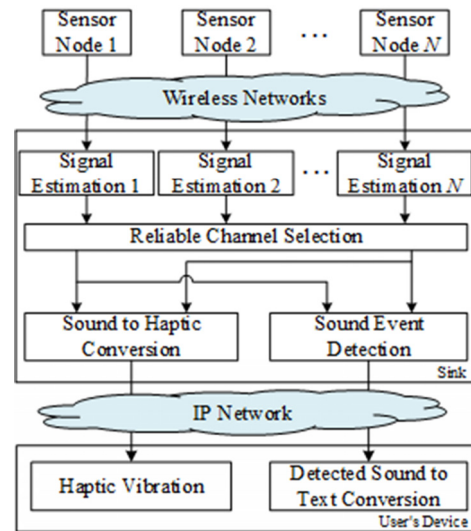


Fig. 1. Architecture of the proposed system.

2.1 싱크 기반 신호 추정

WASNs(Wireless Acoustic Sensor Networks)를 통해 싱크로 전달되는 사운드 신호 패킷은 무선 멀티 홉 통신 과정 중에 손실될 수 있다. 따라서 본 논문에서는 싱크 기반의 패킷 복원 방식을 적용한다.

각 음향센서 노드 마이크를 통해 녹음된 i 번째 신호 패킷이 싱크에 도착하면 패킷 버퍼에 저장되고, G.722.2 음성 코덱을 통해 신호로 디코딩된 후 신호 프레임 버퍼에 저장된다. 이때, 패킷 손실이 발생하였다면 RLPS(Recursive Linear Prediction and Synthesis)^[2] 기반의 패킷 손실 은닉 방식을 통해 신호를 복원한다. 만약 i 번째 패킷 손실이 발생하고, $(i+1)$ 번째 패킷이 존재하지 않으면서 $(i-1)$ 번째 패킷만 존재하면 단구간 패킷복원과 장구간 패킷복원의 두 가지 방식으로 나누어 수행한다. 단구간 패킷복원에서는 $(i-1)$ 번째 패킷에 RLPS를 적용하여 은닉신호를 생성하고, 장구간 패킷손실복원에서는 이전에 합성된 신호를 RLPS에 적용하여 반복적으로 여기 신호를 생성한다. 여기신호 생성 시에는 복원된 구간의 소리를 점진적으로 감소시킨다. 반면에 i 번째 패킷이 정상적으로 수신되고, $(i-1)$ 번째 패킷이 은닉된 신호인 경우에는 merging and smoothing 방식을 이용하여 i 번째 신호 프레임과 $(i-1)$ 번째 신호 프레임의 불연속 지점을 매끄럽게 연결한다. 이와 다르게 i 번째 패킷 손실이 발생하고, $(i+1)$ 번째 패킷과 $(i-1)$ 번째 패킷이 모

두 존재하는 경우에는 $(i-1)$ 번째 신호로부터 생성된 은닉신호와 $(i+1)$ 번째 신호로부터 생성된 은닉신호를 합성하여 손실된 i 번째 신호로 대체한다.

2.2 신뢰할 수 있는 채널 선택

각 음향 센서와 사운드 발생 지점 사이의 거리에 따라 각 채널은 다양한 품질의 신호를 갖는다. 따라서 SED 성능을 향상시키기 위해 고품질 신호를 갖는 채널 선택 방식이 필요하다. 이에 본 논문에서는 여러 홈 공간 내에서 발생하는 SED 성능과 연산 효율을 높이기 위해 MCCC(Multi-Channel Cross-Correlation Coefficient)를 사용한 채널 선택 방식을 적용한다.^[3]

채널 선택 방식의 단계는 다음과 같다. 먼저 여러 개의 마이크 채널 중 패킷 손실이 적고, 에너지 임계값 보다 큰 root mean square 값을 갖는 신호로 구성된 M개의 채널을 선택한다. 다음으로 연산 효율을 높이기 위해 GCC-PHAT(Generalized Cross-Correlation with PHase-based weighting)^[4]을 이용하여 M개 채널 신호의 지연 시간을 예측하고, 신호를 시간순으로 정렬한다. 정렬된 각 채널을 두 개씩 쌍을 지어서 가능한 모든 채널 쌍을 생성하여 각 쌍의 MCCC 값을 계산한다. 모든 채널 쌍 중에서 최대 MCCC 값을 갖는 하나의 채널 쌍을 찾는 것은 많은 연산량이 필요하기 때문에 계산 효율을 향상시키기 위해 가장 큰 MCCC 값을 갖는 하나의 채널 쌍(두 개의 채널)이 남을 때까지 가장 작은 MCCC 값을 제공하는 채널을 제외하는 과정을 반복하여 수행한다.

2.3 사운드 이벤트 검출

Fig. 2는 사운드 이벤트 검출 시스템의 전체적인 구조를 나타내는 블록도이다. 제안하는 시스템은 학습 단계와 테스트 단계로 구성되어 있다.

학습 단계에서는 데이터 희박성 문제를 해결하기 위해 학습 데이터에 잡음 혼합 및 피치 이동 기법을 적용하여 데이터를 확장하고, 잡음을 포함하고 있는 오디오 신호의 품질 향상을 위해 두 채널 오디오 신호를 시간 프레임 단위로 나눈 뒤 특징값으로 잡음 감소 스펙트로그램과 TDOA(Time Delay Of Arrival)^[4]을 추출한다. 잡음 감소 스펙트로그램을 추출하기

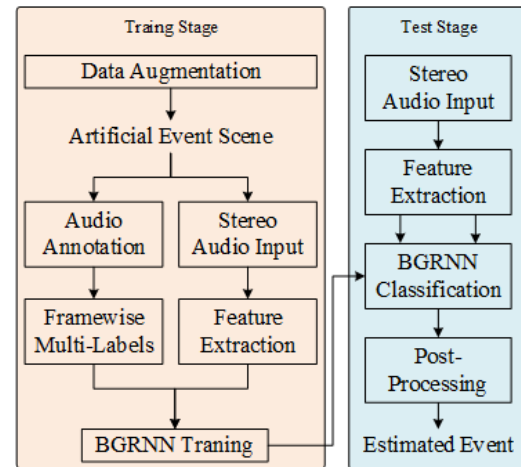


Fig. 2. Framework of the training and testing procedure for the proposed SED system.

위해 사운드 신호를 단시간 푸리에 변환에 적용하여 스펙트로그램으로 변환한다. 이를 재귀적인 배경잡음 에너지 임계값을 통해 사운드 이벤트 신호가 존재하지 않는 프레임을 검출하여 잡음추정에 반영한다. 이와 함께 사운드가 발생한 공간의 잔향 시간 모델을 적용하여 잔향을 추정한다. 두 방식을 통해 추정된 잡음과 잔향을 스펙트로그램에서 제거하여 입력 신호의 품질을 향상시키고 로그 연산을 통해 추출된 스펙트로그램을 정규화한다. 이렇게 추출된 특징값들은 BGRNN에 적용되어 학습을 수행한다. BGRNN은 3개의 게이트로 이루어진 LSTM과 다르게 업데이트 게이트와 리셋 게이트 2개로 구성된 GRU(Gated Recurrent Unit)를 사용하기 때문에 LSTM-RNN에 비해 비교적 단순한 구조를 갖는다. GRU의 업데이트 게이트는 현재 상태를 결정하기 위해 이전 메모리로부터 유입되는 정보를 제어하고, 리셋 게이트는 현재의 입력정보와 이전 메모리에서 유출되는 정보를 제어한다. GRU로 구성된 BGRNN^[5]의 히든 레이어는 두 개로 분리되어 하나는 학습데이터의 정방향 시퀀스를, 다른 하나는 역방향 시퀀스를 사용하여 학습을 수행한다. 각 히든 레이어를 통한 출력 과정은 다음과 같이 계산된다.

$$\vec{q}_t = h(U_{q^-} x_t + W_{q^-} \vec{q}_{t-1} + b_{q^-}), \quad (1)$$

$$\overleftarrow{q}_t = h(U_{q^-} x_t + W_{q^-} \overleftarrow{q}_{t-1} + b_{q^-}), \quad (2)$$

$$q_t = \text{Concat}[q_t^+, q_t^-], \quad (3)$$

여기서 q^+ , q^- 는 각 히든레이어의 정방향 시퀀스, 역방향 시퀀스를 나타낸다. x_t , U , W , b , h 는 각각 현재 입력 프레임, 입력 레이어에서 히든레이어로 연결되는 가중치, 히든 레이어에서 히든레이어로 연결되는 가중치, 바이어스 값, 활성화 함수를 의미한다.

테스트 단계에서는 스테레오 오디오 신호로부터 추출된 특징값들을 학습된 BGRNN에 입력하고, 후처리 과정에서 Mesaros *et al.*^[6] 방식과 유사하게, 오디오 신호의 연속적인 구간에서 0.1s보다 작은 시간 간격으로 발생하는 이벤트는 무시하고, 0.1s보다 긴 사운드만 감지하여 분류한다. 이때 인식된 사운드 이벤트는 상황 정보를 제공하기 위해 텍스트로 변환되어 청각장애인에게 제공된다.

2.4 사운드 햅틱 변환

본 논문에서는 청각장애인에게 효과적으로 사운드 정보를 촉각 신호로 제공하기 위해 인식된 사운드에 하모닉/퍼커시브 음원 분리 방식^[7]을 적용한 사운드 햅틱 변환 방식을 제안한다.

Fig 3은 사운드의 햅틱 진동 신호 변환 단계를 보여주는 블록도이다. 먼저 2.2장에서 선택된 두 채널의 신호를 모노 포닉 사운드로 변환하고, 단시간 푸리에 변환에 적용한다. 하모닉 성분과 퍼커시브 성분으로 구성된 커널 모델과 반복적인 Backfitting 방식을 적용하여 추출된 스펙트로그램을 퍼커시브와 하모닉 스펙트럼 성분으로 분리하고, 각 스펙트럼 성분의 주파수 대역을 0 Hz~75 Hz, 75 Hz~150 Hz, 150 Hz~250 Hz, 250 Hz~500 Hz, 500 Hz~750 Hz로 하여 사람의 귀에 예민한 주파수 영역에 해당되는 350 Hz 이하 구간은 좁게, 나머지 주파수 구간은 넓게 분할한다. 그다음 각 주파수 대역과 시간축의 진폭 신호를 시간 프레임 단위로 분할하여 에너지값 $h = [h_1, h_2, \dots, h_t]$ 을 구하고, h 의 평균 m 과 표준편차 σ 를 계산하여 Eq. (4)에 적용함으로써 0과 255 사이의 값으로 정규화된 햅틱 데이터 \bar{h}_t 를 생성한다.

$$\bar{h}_t = [(h_t - m) / \sigma] \times 255. \quad (4)$$

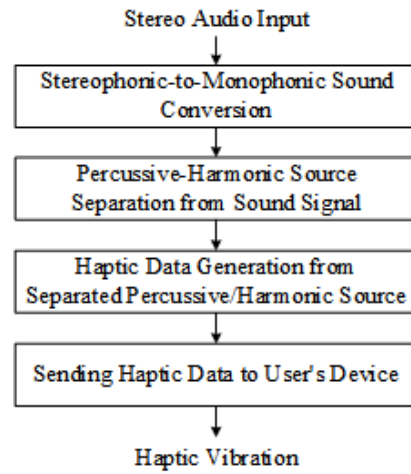


Fig. 3. Diagram of a process for generating haptic vibration.

생성된 햅틱 진동 데이터는 IP Network를 통해 12개의 진동 모듈로 구성된 디바이스로 전송되어 청각장애인에게 홈 환경에서 발생한 상황을 텍스트 정보와 함께 햅틱 진동 정보로 전달한다.

III. 실험 및 결과

3.1 실험 데이터

본 논문에서는 제안한 방식의 성능 측정을 위해 실생활에서 발생하는 사운드를 수집하여 데이터베이스를 구성하였다. 데이터 수집을 위해 키친, 복도, 거실, 침실, 서재, 샤워실 및 화장실로 구성된 59 m² 크기의 아파트에 센서를 설치하였고, 각 방에 설치된 6개의 마이크를 통해 16 kHz 샘플링레이트, 24비트 해상도를 사용하여 5 min~15 min 길이로 다양한 사운드를 녹음하였다. 총 녹음된 사운드 파일은 2300개이고, 총 재생 시간은 19000 min이다. 파일에는 걷는 소리, 문 잠그는 소리, 접시 떨어지는 소리, 유리 깨지는 소리, 물체 떨어지는 소리, 비명 소리, 전화 또는 현관문 벨 소리, 대화 소리, 음악 소리, 물소리, 아기 울음소리, 청소기 소리, 박수 소리, 도로 소음, 강아지 소리, 웃음소리 등이 포함된다. SED 실험을 위해 전체 데이터베이스 중 60%는 학습 데이터로, 20%는 검증 데이터로, 나머지 20%는 테스트 데이터로 사용하였다.

WASNs에서 SED을 위한 테스트 베드는 다음과 같

이 설정하였다. 직선 라인에 일정한 간격으로 배치된 36개의 센서 노드 중 한 개는 라인 한쪽 끝에 설치하여 싱크를 설치하였고, 또 다른 한 개는 싱크와 반대 라인 끝에 설치하여 패킷 생성을 위한 노드로 사용하였다. 전송 전력은 0 dBm으로 설정하여 약 4 m의 전송 범위를 갖도록 하고, 무선 주파수는 890 MHz로 설정하였다. 대역폭은 38.4 kbps이고, 패킷의 크기는 36바이트로 고정되어 최대 링크 당 133패킷(pkt/s)의 용량을 갖도록 구성하였으며, 지연(25 ms~80 ms), 지터(40 ms~300 ms) 및 패킷 손실(2%~10%)을 포함하는 임의의 traffic load 기반의 WASNs을 적용하였다. 본 논문에서는 실험을 위해 0.5 pkt/s의 패킷 생성 속도와 8%의 패킷 손실률을 사용하였다.

3.2 측정 방식

WASNs에서 SED 시스템을 위해 SE(신호 추정), CS(채널 선택), 특징값 추출 방식으로 2(2채널), NR(잡음 감소), ST(스펙트로그램)을 사용하였다. SED에서는 분류기로 LSTM, GRNN, BGRNN를 사용하여 성능을 비교하였다. 실험에서 LSTM은 200개의 LSTM 뉴런으로 구성된 3개의 히든레이어를 사용하였고, GRNN과 BGRNN은 200개의 GRU 뉴런으로 구성된 3개의 히든레이어를 사용하였다. 입력 레이어 뉴런의 수는 사용한 오디오 신호의 길이에 따라 다르고, 출력 레이어의 뉴런 수는 16개로 클래스 개수와 동일하다. 신경망 학습의 손실 함수로는 binary cross-entropy를 이용한 back propagation을 사용하였다. 실험의 측정 지표로는 ER(Error Rate)와 F-score를 사용하였으며 1 s 단위의 세그먼트 기준으로 계산하였다.

3.3 실험 결과

Table 1은 여러 특징들과 각 분류기의 다양한 조합에 따른 사운드 이벤트 검출 성능 결과를 보여준다. 본 실험에서는 다중 채널 마이크를 사용하여 CS 방식을 모든 SED 방식에 적용하였다. M2, M3의 실험 결과를 통해 신호추정 방식을 적용하지 않았을 때 손실된 패킷 신호가 SED 성능에 영향을 끼쳐 낮은 성능을 보이는 것을 확인할 수 있었고, M1, M2의 결과로부터 동일한 분류기를 사용할 때 모노 채널 특징을 사용하는 것보다 청각적 특성을 고려한 스테레

Table 1. Comparison of the segment-based error rate and F-score for different combinations of classifiers and features.

Method (M)			
	Classifier	Feature	F-score
M1	BGRNN	SE, CS, NR2, ST2, TDOA	90.4
M2	BGRNN	SE, CS, NR, ST, TDOA	86.1
M3	BGRNN	CS, ST, TDOA	78.3
M4	GRNN	SE, CS, NR2, ST2, TDOA	89.7
M5	LSTM	SE, CS, NR2, ST2, TDOA	88.3

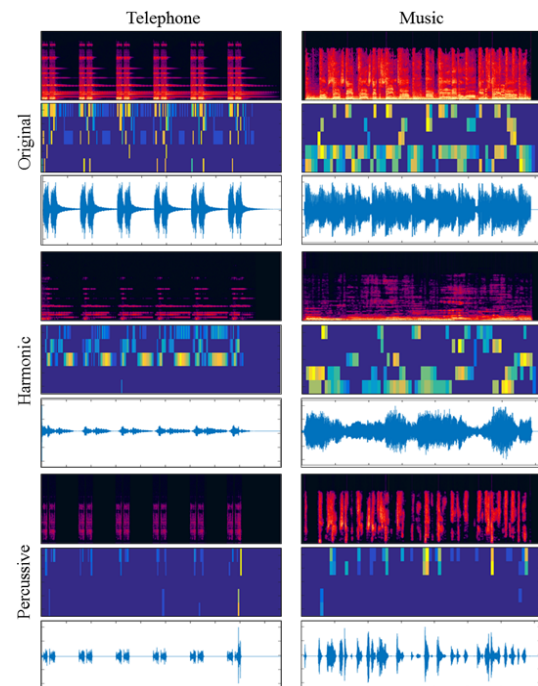


Fig. 4. Results of sound-to-haptic conversion using harmonic-percussive source separation.

오 채널 기반의 2채널 특징값을 적용한 경우 성능이 향상되는 것을 확인할 수 있었다. M1, M4, M5의 실험 결과에서는 동일한 특징 조합 기반의 분류방식으로 BGRNN을 사용했을 때 ER 0.48, F-score 90.4로 가장 우수한 성능을 보였고, 이는 이전 시간 프레임의 정보뿐 아니라 역방향 시퀀스도 함께 학습에 사용하는 것이 현재의 SED에 효과적임을 확인할 수 있다.

Fig 4는 M1 방식과 하모닉/퍼커시브 음원 분리 방식을 적용하여 추출한 각 객체 음원의 시간축 신호, 스펙트럼, 그리고 스펙트럼으로 변환된 햅틱 정보를 나타낸다. 분리된 각 스펙트럼에는 분리되기 이전의 스펙트럼으로부터 하모닉 성분으로 수평적 특징을

갖는 멜로디 스펙트럼 성분이 분리되었고, 퍼커시브 성분으로는 수직적 특징을 갖는 주기적인 스펙트럼 성분이 분리되었음을 확인할 수 있다.

사운드의 햅틱 진동 변환의 성능 측정을 위해 MOS 테스트를 수행하였다. 음원 분리가 적용되지 않은 사운드 햅틱 변환 방식의 MOS 점수는 3.1인 반면, 하모닉/퍼커시브 음원 분리를 적용한 햅틱 진동 방식은 3.4의 MOS 점수를 획득하였다. 이를 통해 사용자에게 사운드 신호로부터 추출된 주기와 멜로디 정보를 각각의 진동 신호로 제공하는 것이 사운드 인식에 더 효과적임을 확인할 수 있다.

IV. 결 론

본 논문에서는 청각장애인을 위한 SED 방식 기반 홈 모니터링 시스템을 제안하였다. 효과적인 SED와 사운드 햅틱 변환 방식을 위해 신호를 추정하여 손실된 패킷을 복원하였고, 상관관계가 높은 채널을 선택하였다. SED를 위해서는 두 개 채널 특징값들과 BGRNN 분류기를 사용하였고 실험을 통해 제안한 방식이 기존 방식보다 더 우수한 성능을 보이는 것을 확인하였다. 향후 본 연구에서 제안한 시스템을 무선 센서 네트워크 기반의 사물인터넷을 위한 시청각 상황 맥락 인식과 모니터링에 적용할 예정이다.

감사의 글

이 논문은 2018년 대한민국 교육부와 한국연구재단의 지원을 받아 수행된 연구임(NRF-2018S1A3A 2074955).

References

1. M. Zöhrer and F. Pernkopf, "Gated recurrent networks applied to acoustic scene classification and acoustic event detection," Proc. Detection and Classification of Acoustic Scenes and Events 2016, 1-5 (2016).
2. B. H. Kim, H.-G. Kim, J. Jeong, and J. Y. Kim, "VoIP receiver-based adaptive playout scheduling and packet loss concealment technique," IEEE Trans. Consum. Electron., **59**, 250-258 (2013).
3. K. Kumatani, J. McDonough, J. F. Lehman, and B. Raj, "Channel selection based on multichannel cross-correlation coefficients for distant speech recognition," Proc. Joint Workshop Hands-free Speech Commun. Microphone Arrays, 1-6 (2011).
4. D. Pavlidi, A. Griffin, M. Puigt, and A. Mouchtaris, "Real-time multiple sound source localization and counting using a circular microphone array," IEEE Trans. Audio, Speech, Lang. Process., **21**, 2193-2206 (2013).
5. R. Lu and Z. Duan, "Bidirectional GRU for sound event detection," Proc. Detection and Classification of Acoustic Scenes and Events 2017, 1-4 (2017).
6. A. Mesaros, T. Heittola, and T. Virtanen, "TUT database for acoustic scene classification and sound event detection," Proc. 24th Eur. Signal Process. Conf., 1128-1132 (2016).
7. A. Liutkus, D. Fitzgerald, Z. Rafii, B. Pardo, and L. Daudet, "Kernel additive models for source separation," IEEE Trans. Signal Process., **62**, 4298-4310 (2014).

저자 약력

▶ 김 지 연 (Gee Yeun Kim)



2018년 2월: 광운대학교 전자융합공학과 학사
2018년 3월 ~ 현재: 광운대학교 전자융합공학과 석사과정

▶ 신 승 수 (Seung-Su Shin)



2019년 2월: 광운대학교 전자융합공학과 학사
2019년 3월 ~ 현재: 광운대학교 전자융합공학과 석박사통합과정

▶ 김 형 국 (Hyoung-Gook Kim)



1999년 ~ 2002년: 독일 SIEMENS/Cortologic AG 책임연구원
2002년 ~ 2005년: 독일 베를린 공과대학교 Assistant Professor
2005년 ~ 2007년: 삼성종합기술원 수석연구원
2007년 3월 ~ 현재: 광운대학교 전자융합공학과 교수