

## 딥러닝 기반 손 제스처 인식을 통한 3D 가상현실 게임

이병희<sup>○</sup>

오동한

김태영\*

서경대학교 컴퓨터공학과

{powermike, 5donghan, tykim}@skuniv.ac.kr

### 3D Virtual Reality Game with Deep Learning-based Hand Gesture Recognition

Byeong-Hee Lee<sup>○</sup>

Dong-Han Oh

Tae-Young Kim\*

Department of Computer Engineering, Seokyeong University

#### 요 약

가상 환경에서 몰입감을 높이고 자유로운 상호작용을 제공하기 위한 가장 자연스러운 방법은 사용자의 손을 이용한 제스처 인터페이스를 제공하는 것이다. 그러나 손 제스처 인식에 관한 기존의 연구들은 특화된 센서나 장비를 요구하거나 낮은 인식률을 보이는 단점이 있다. 본 논문은 손 제스처 입력을 위한 RGB 카메라 이외 별도 센서나 장비 없이 손 제스처 인식이 가능한 3차원 DenseNet 합성곱 신경망 모델을 제안하고 이를 기반으로 한 가상현실 게임을 소개한다. 4개의 정적 손 제스처와 6개의 동적 손 제스처 인터페이스에 대해 실험한 결과 평균 50ms의 속도로 94.2%의 인식률을 보여 가상현실 게임의 실시간 사용자 인터페이스로 사용 가능함을 알 수 있었다. 본 연구의 결과는 게임 뿐 아니라 교육, 의료, 쇼핑 등 다양한 분야에서 손 제스처 인터페이스로 활용될 수 있다.

#### Abstract

The most natural way to increase immersion and provide free interaction in a virtual environment is to provide a gesture interface using the user's hand. However, most studies about hand gesture recognition require specialized sensors or equipment, or show low recognition rates. This paper proposes a three-dimensional DenseNet Convolutional Neural Network that enables recognition of hand gestures with no sensors or equipment other than an RGB camera for hand gesture input and introduces a virtual reality game based on it. Experimental results on 4 static hand gestures and 6 dynamic hand gestures showed that they could be used as real-time user interfaces for virtual reality games with an average recognition rate of 94.2% at 50ms. Results of this research can be used as a hand gesture interface not only for games but also for education, medicine, and shopping.

키워드: 손 제스처 인식, 딥러닝, 합성곱 신경망, DenseNet, 가상현실, 게임

**Keywords:** Hand Gesture Recognition, Deep Learning, Convolutional Neural Network, DenseNet, Virtual Reality, Game

### 1. 서론

최근 가상현실 기술이 발전하고 다양한 기기들이 보급됨에 따라 사용자의 몰입감을 증폭시키는 가상현실 콘텐츠에 대한 요구가 증가하고 있다. 기존 가상현실 콘텐츠는 대부분 버튼 방식의 일반 컨트롤러를 사용하기 때문에 사용자와

가상현실의 객체 사이의 상호작용이 부자연스럽고 직관성이 떨어지는 단점이 있다.

사용자에게 가장 자연스러운 인터페이스는 사람이 일상생활에서 사용하는 손 제스처를 가상공간 속에서도 인식할 수 있도록 하는 것이다. 이와 같은 손 제스처 인식에 대한 연구로 키넥트나 립모션과 같은 장비를 이용한 손 제스처 인터페이스에 관한 연구[1-6]가 있었지만 특화된 장비 혹은 센서가 필요하고 조명이나 거리와 같은 주변 환경에 제약을

\*corresponding author: Tae-Young Kim/Seokyeong University(tykim@skuniv.ac.kr)

받는다는 단점이 있었다.

최근에는 큰 규모의 데이터 세트를 수집하는 것이 용이해지고 GPU를 활용한 고성능 컴퓨팅이 보편화됨에 따라 딥러닝(Deep Learning) 기술을 이용하여 손 제스처를 인식하고자 하는 연구들이 이뤄지고 있다. 이와 관련한 기존 연구로 스테레오 비디오로부터 구한 깊이 정보와 색 정보를 이용하여 검출한 손 윤곽선 정보를 학습시켜 인식하는 연구[7], 깊이 카메라로 얻은 관절 정보를 학습하여 인식하는 연구[8] 그리고 칼라 영상과 깊이 영상을 결합하여 3차원 합성곱 신경망(Convolutional Neural Network)으로 학습하여 인식하는 연구[9] 등이 있다. 하지만 위의 연구들 역시 특정 장비나 센서를 필요로 하거나, 낮은 인식률을 보이는 문제점을 지니고 있다.

본 논문에서는 가상현실 기기에 부착한 RGB 형식의 카메라 이외 별도 센서 없이 딥러닝 기술인 합성곱 신경망을 통해 사용자의 손 제스처를 실시간으로 인식하는 3D 가상현실 게임을 소개한다. 먼저 사용자가 가상현실 게임의 사용자 인터페이스를 조작하고 마법을 구현하는 등의 명령을 수행하기 위한 제스처를 정의한다. 정의한 손 제스처 인식을 위한 합성곱 신경망은 압축 계층을 지닌 고밀도 연결 구조를 사용하여 신경망이 깊어질수록 정보를 손실하는 문제를 보완하고 적은 파라미터로 높은 속도와 인식률을 제공하는 DenseNet[10]을 3차원으로 확장하였다. 본 방법의 검증을 위하여 조명, 배경, 영상의 손 위치와 거리 등 다양한 상황을 고려하여 제작한 10가지 정적 및 동적 손 제스처 데이터 세트로 실험한 결과 평균 50ms의 속도로 94.2%의 인식률을 보여 가상현실 게임의 실시간 사용자 인터페이스로 사용 가능함을 알 수 있었다.

본 논문의 구성은 다음과 같다. 2장에서 본 연구에서 사용한 합성곱 신경망 구조를 소개하고 3장에 손 제스처 유형 및 데이터 세트 제작 방법과 손 제스처 인식 과정을 설명한다. 4장에서는 본 방법을 적용한 가상현실 게임을 소개하고 5장에서 실험결과를 기술한 후 6장에서 향후 연구 방향을 기술한 후 결론을 맺는다.

## 2. 합성곱 신경망

합성곱 신경망은 사람의 사고방식을 컴퓨터에게 가르치는 머신 러닝의 한 분야인 딥러닝 알고리즘 중 하나로 영상 인식과 음성 인식 모두에서 뛰어난 성능을 보인다. 합성곱 신경망은 기존의 완전 연결 신경망에 합성곱 연산과 최대 풀링을 추가하여 데이터의 공간 정보를 유지하면서 인접 데이터와의 특징을 효과적으로 추출한다.

최근에는 신경망의 구조가 더욱 복잡해지고 깊어짐에 따라서 성능의 향상보다 기존의 신경망 구조를 최적화하는 것에 초점을 맞추고 있다. 대표적인 예로 인셉션(Inception) 구

조를 사용한 GoogLeNet[11]과 고밀도 연결 구조(Dense Connectivity)를 사용한 DenseNet은 파라미터 수를 획기적으로 감소시킴으로써 파라미터 저장을 위한 메모리 사용량을 감소시켰다.

### 2.1 DenseNet

DenseNet은 각 계층의 합성곱 연산 후의 출력을 그 이후의 모든 계층에 연결하는 고밀도 연결 구조를 사용해 신경망이 깊어질수록 초기 계층의 정보를 잃어버리는 문제를 해결하였다. 또한 단순한 행렬 덧셈으로 계층 사이를 연결했던 ResNet[12]의 합산(Summation) 방식 대신 이전 계층의 특징맵을 결합해 기존 정보를 유지하며 계층 사이를 연결하는 결합(Concatenation) 방식을 제안하였다.

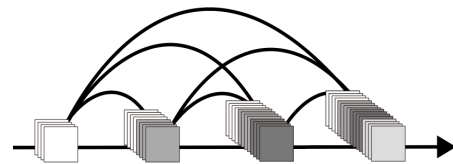


Figure 1: Dense Connectivity

Figure 2와 같이 하나의 고밀도 연결의 수행되는 범위를 고밀도 구역(Dense Block)이라 하고 각 구역 사이의 이행 계층(Transition Layer)에서 평균 풀링을 수행하여 특징맵의 크기를 축소시킨다.

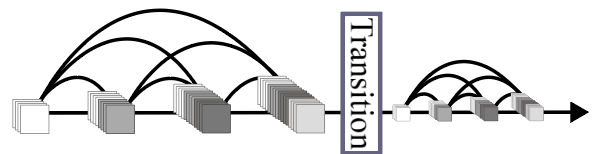


Figure 2: DenseNet with Dense Blocks

마지막으로 파라미터 수를 줄이기 위해 기존 모델에 비해 합성곱 계층에서 사용하는 커널의 깊이를 축소시키고 NIN(Network In Network)[13]의  $1 \times 1$  합성곱 연산을 적용한 병목 계층(Bottleneck Layer)과 압축 계층(Compression Layer)을 추가하였다. 병목 계층은 고밀도 구역 내 합성곱 계층의 이전 단계에서 특징맵의 깊이를 조정하고, 압축 계층은 고밀도 구역의 끝에서 이행 계층으로 넘어가는 특징맵의 깊이를 조정하여 파라미터 크기를 줄인다.

### 2.2 제안하는 신경망

손 제스처는 시작부터 끝까지 손의 모양과 위치가 변하지 않는 정적 제스처와 시간에 따라 위치가 변하는 동적 제스

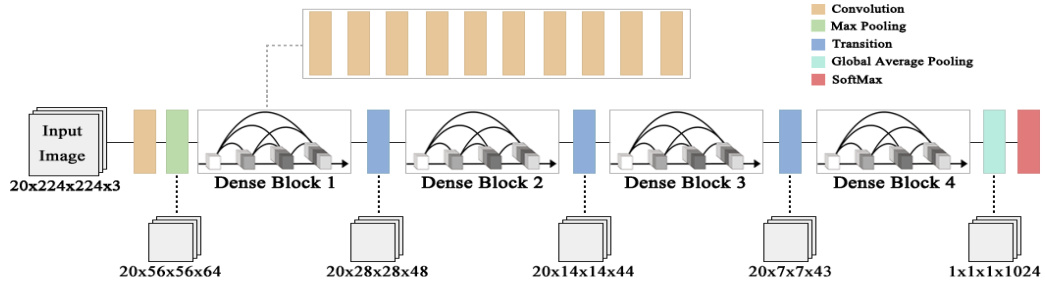


Figure 3: 3D DenseNet Structure

쳐로 분류된다. 본 논문은 정적 제스처와 동적 제스처를 모두 분류하기 위하여 공간적인 특징만을 고려하는 기존의 DenseNet에 시간상의 정보를 고려할 수 있도록 Figure 3과 같은 차원을 확장한 구조의 3D DenseNet 구조를 제안한다.

신경망은 4개의 고밀도 구역으로 구성되어 있고 각각의 고밀도 구역은 10개의 합성곱 계층을 가진다. 크기가  $20 \times 224 \times 224$ 인 데이터를 입력받아 첫 번째 합성곱 계층과 네 개의 이행 계층, 전역 평균 풀링을 거쳐 특징맵의 크기가 축소된다. 병목 계층과 압축 계층을 제외한 모든 합성곱 계층은  $3 \times 3 \times 3$  크기의 커널을 사용한다. 각 고밀도 구역에서 출력된 특징맵은 25%의 압축률을 가지는 압축 계층을 거쳐 이행 계층의 평균 풀링을 수행한다. 마지막 이행 계층 이후에는 파라미터 수를 줄이기 위하여  $20 \times 7 \times 7$  크기의 커널을 사용하는 전역 평균 풀링을 수행한다. 그리고 난후, 1024개의 뉴런을 가지는 은닉 계층을 거쳐 10개의 라벨로 이루어진 출력 층을 가지는 완전 연결 계층(Fully Connected Layer)을 수행한다.











### 3. 손 제스처 인식

손 제스처 인식 기술은 사람이 손을 이용하여 미리 정해진 동작을 했을 때, 그것이 어떤 동작인지를 인식하는 기술을 말한다. 본 장에서는 가상현실 게임을 위한 손 제스처를 정의하고 실험 데이터 세트 제작 방법과 학습과정을 설명한다.

#### 3.1 손 제스처 정의

가상현실에서 사용자와 객체 사이의 상호작용을 위한 손 제스처는 친숙하고 직관적이며 인식의 정확도가 높아야 한다. 이 점을 고려하여 Table 1과 같은 정적인 형태의 4가지 손 제스처와 동적인 형태의 6가지 손 제스처 유형을 정의하였다.

Table 1: Hand Gesture Types

Static Gestures		Dynamic Gestures	
Point	Okay	Push	Swipe Down
			
No	Receive	Swipe Right	Swipe Left
			
-	-	Draw Circle	Draw Question
-	-		

#### 3.2 데이터 세트

데이터 세트는 640x480 해상도의 카메라로 촬영되었고, 4명의 학습자가 참가하여 Table 1에서 정의한 손 제스처를 다양한 위치, 각도, 거리에서 시행하여 제작하였다. 초당 20 프레임으로 구성된 손 제스처 영상은 1000세트 중 20,000장을 제작하였고, 과적합 문제를 개선하기 위해 제작된 손 제스처 영상의 손의 위치와 크기, 조명을 다양하게 후처리하여 5배 증가시켜 총 100,000장을 제작하였다. 총 100,000장 중 90%인 4,500세트는 학습데이터로, 나머지 10%인 500세트는 평가 데이터로 사용하였다.

### 3.3 학습 과정

학습 과정은 Figure 4와 같다. 입력된 학습 데이터 세트는 DenseNet 구조의 3차원 합성곱 신경망의 순전파 과정을 거쳐 출력 층의 Cross Entropy를 통해 오차가 계산되고 이후 역전파 과정에서 계산된 오차와 Adam Optimizer[14]를 통해 가중치를 최적화한다. 전체 데이터에 대한 한 번의 학습이 끝나면 평가 데이터 세트를 통해 인식률을 측정하고 학습된 모델에 추가적으로 위와 같은 과정을 반복한다. 최종적으로 인식률이 가장 높은 모델을 선정하여 학습 모델을 구성한다.

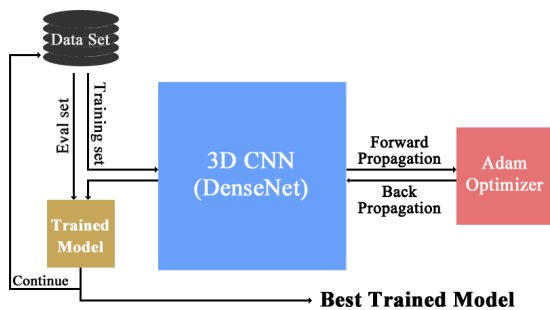


Figure 4: Learning Process

## 4. 장치 및 게임 소개

### 4.1 장치 소개

본 논문은 마이크로소프트의 라이프캠 HD-3000을 사용하여 사용자가 실시간으로 취하는 손 제스처를 입력받고, HMD(Head Mounted Display)로 HTC와 Valve사에서 공동 출시한 Vive를 사용하여 가상현실 게임을 개발하였다(Figure 5). Vive는 적외선 센서인 베이스 스테이션을 통해 HMD를 추적하여 실제로 사용자가 이동한 거리나 방향을 인식한다.

본 논문은 Vive에서 제공하는 외부 컨트롤러를 사용하지 않고 라이프캠 HD-3000을 Vive와 연동하여 손 제스처를 통해 게임 내 모든 명령을 수행하는 가상현실 게임을 개발하였다.

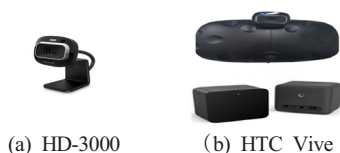


Figure 5: Devices

### 4.2 게임 소개

본 논문은 손 제스처 인터페이스를 적용하여 생존게임 장르의 가상현실 게임을 개발하였다. Figure 6에서 보는 바와 같이 게임화면 내에 사용자가 취하는 손 제스처를 실시간으로 보여주어 조작을 좀 더 용이하도록 하였고, 사용자는 정의된 손 제스처를 취하여 마법, UI 조작 등의 게임 내 명령을 수행할 수 있다.



Figure 6: Game Scene

### 4.3 동작 설명

- 번개 (Lightning)  
검지로 정면을 가리키는 정적 제스처를 취하면 사용자가 바라보는 방향으로 번개 마법이 발사된다.
- 확인 (Confirmation)  
검지와 엄지를 둥글게 하여 O 모양을 그리는 정적 제스처를 취하면 확인 명령을 수행한다.
- 닫기 (Close)  
양손의 검지를 서로 교차하여 X 모양을 그리는 정적 제스처를 취하면 열려있는 도움말, 마법 책 등의 UI를 닫는 명령을 수행한다.
- 마법 책 (Skill Book)  
양 손바닥이 위를 바라보게 펴고 있는 정적 제스처를 취하면 마법들을 배울 수 있는 마법 책이 내려온다.
- 화염 구체 (Fireball)  
양손으로 미는 동적 제스처를 취하면 사용자가 바라보는 방향으로 화염 구체 마법이 발사된다.
- 유성 (Meteor)  
손바닥을 아래로 향하게 펴고 아래로 내리는 동적 제스처를 취하면 플레이어 전방에 유성들이 떨어진다.
- 이전/다음 (Previous/Next)  
손바닥을 화면의 왼쪽에서 오른쪽으로 이동하는 동적 제스처를 취하면 마법 책이 이전 페이지로 넘어가고 화면의 오른쪽에서 왼쪽으로 이동하는 동적 제스처를 취하면 마법 책이 다음 페이지로 넘어간다.
- 얼음 구체 (Iceball)  
검지로 원을 그리는 동적 제스처를 취하면 사용자가 바라보는 방향으로 얼음 구체 마법이 발사된다.



● 도움말 (Help)

검지로 물음표 모양을 그리는 동적 제스처를 취하면 도움말이 열린다.

Table 2와 Table 3은 게임 내에서 정의한 제스처를 실행한 화면이다.

Table 2: Static Gesture Commands





Command	Example
Lightning	
Confirmation	
Close	
Skill Book	

Table 3: Dynamic Gesture Commands

Command	Examples
Fireball	
Meteor	
Previous	
Next	
Iceball	
Help	

## 5. 실험

### 5.1 실험환경

본 실험은 프로세서 Intel Core i5 8400, 그래픽 카드 GeForce GTX 1080Ti, RAM 16GB 등으로 구성된 장비에 Python 3.5 기반의 Tensorflow 1.5.0으로 신경망을 구축하고 Unity 3D 2018.2.8.fl 엔진으로 가상현실 게임을 개발하여 실험을 진행하였다.

### 5.2 실험결과

본 논문에서 제안한 3D DenseNet 구조의 학습 모델에 대한 성능 검증을 위하여 동일한 조건으로 GoogLeNet 구조의 학습 모델을 구축하여 비교 실험을 진행하였다. Table 4와 Figure 7은 두 학습 모델의 제스처 별 인식률을 나타낸 것이다.

Table 4: Comparison of Recognition Rate between 3D DenseNet and 3D GoogLeNet

Gesture	3D DenseNet	3D GoogLeNet
Point (P)	86.0%	76.0%
Okay (O)	90.0%	80.0%
No (N)	92.0%	88.0%
Receive (R)	94.0%	94.0%
Push (PS)	90.0%	72.0%
Swipe Down (SD)	98.0%	86.0%
Swipe Left (SL)	98.0%	86.0%
Swipe Right (SR)	96.0%	84.0%
Draw Circle (DC)	100.0%	78.0%
Draw Question (DQ)	98.0%	90.0%
Average	94.2%	83.4%

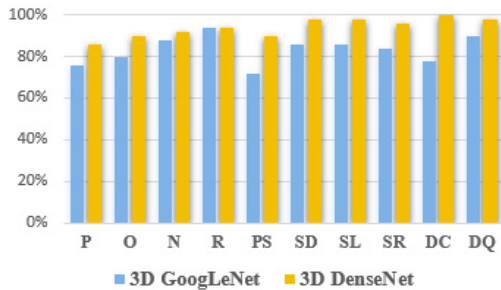


Figure 7: Comparison of Recognition Rate between 3D DenseNet and 3D GoogLeNet

3D DenseNet 학습 모델은 3D GoogLeNet 학습 모델보다 전체적으로 우수한 인식률을 보였다. 두 학습 모델에 대하여 대부분의 인식률이 정적 제스처가 동적 제스처보다 떨어지는 것으로 보아 정적인 형태의 데이터를 입력받는 2차원 합성곱 신경망에 비해 3차원 합성곱 신경망은 동적인 형태의 데이터를 입력받아 손의 공간적인 특징보다 시간적 차원의 특징을 우선적으로 추출하는 것을 알 수 있었다.

Table 5와 Table 6은 두 학습 모델의 파라미터 수를 나타낸 것이다. 파라미터 수는 DenseNet 학습 모델이 GoogLeNet 학습 모델보다 6배 가까이 적은 것을 확인할 수 있다. 이를 통해 DenseNet 학습 모델이 메모리 측면에서도 우수함을 알 수 있다.

Table 5: Number of Parameters of Our Method

Type	Params
Convolution	28.22K
Dense Block 1	507.90K
Transition 1	9.21K
Dense Block 2	548.86K
Transition 2	7.74K
Dense Block 3	544.77K
Transition 3	7.40K
Dense Block 4	543.74K
Global Average Pooling	
Total	2,197.84K

Table 6: Number of Parameters of 3D GoogLeNet

Type	Params
Convolution 1	27.56K
Max Pooling	
Convolution 2	328K
Max Pooling	
Inception 1	400.5K
Inception 2	961K
Inception 3	755.25K
Inception 4	954K
Inception 5	1148.5K
Inception 6	1419.5K
Inception 7	1947K
Inception 8	2118K
Inception 9	3005K
Global Average Pooling	
Total	13064.31K

Table 7은 가상현실 게임 상에서 본 논문에서 제안한 3D DenseNet 학습모델과 3D GoogLeNet 학습모델의 손 제스처

인식 처리시간을 비교한 것이다. 두 가지 방법 모두 실시간으로 손 제스처 인터페이스를 제공하는 것이 가능함을 알 수 있었다.

Table 7: Average Processing Time for Hand Gesture Recognition

Trained Model	Average Processing Time	Standard Deviation
3D DenseNet	50ms	3ms
3D GoogLeNet	71ms	5ms

## 6. 결론

본 논문은 가상현실 기기에 부착한 RGB 형식의 카메라 이외 별도의 센서나 장비 없이 실시간으로 정적 손 제스처와 동적 손 제스처를 인식하기 위한 3D DenseNet 기반 신경망 구조를 제안한 후 가상현실 게임 상에서 실시간으로 손 제스처 인터페이스로 적용 가능함을 보였다.

가상현실 환경에서 손 제스처를 이용하여 몰입감 있고 자연스러운 인터페이스를 제공하려면 간단한 정적 및 동적 제스처 뿐만 아니라 스토리를 가지는 손 제스처의 인식이 필요하다. 향후 연구로 스토리 형식의 손 제스처 인식이 가능한 딥러닝 모델에 대한 연구를 수행하고자 한다.

## 감사의 글

이 논문은 2017년도 정부(교육부)의 재원으로 한국연구재단의 지원을 받아 수행된 기초연구사업임. (No. NRF-2017R1D1A1B03029834)

## References

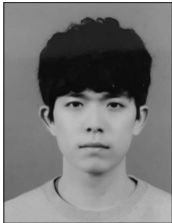
- [1] 박경범, 이재열, “가상현실 환경에서 3D 가상객체 조작을 위한 인터페이스와 인터랙션 비교 연구,” *한국CDE학회 논문집*, 21(1), pp. 20-30, 2016. 3.
- [2] 윤종원, 민준기, 조성배, “몰입형 가상현실의 착용식 사용자 인터페이스를 위한 Mixture-of-Experts 기반 제스처 인식,” *한국HCI학회 논문지*, 6(1), pp. 1-8, 2011. 5.
- [3] 나민영, 유휘중, 김태영, “스마트 디바이스 제어를 위한 비전 기반 실시간 손 포즈 및 제스처 인식방법,” *한국차세대컴퓨팅학회 논문지*, 8(4), pp.27-34, 2012. 8.
- [4] 이새봄, 정일홍, “키넥트를 사용한 NUI 설계 및 구현,” *한국디지털콘텐츠학회 논문지*, 15(4), pp. 473-480, 2014. 8.
- [5] 고택균, 윤민호, 김태영, “HMM과 MCSVM 기반 손 제스처 인터페이스 연구,” *한국차세대컴퓨팅학회 논문지*, 14(1), pp. 57-64, 2018. 2.
- [6] 김민재, 허정만, 김진형, 박소영, 장준호, “직관적인 손동작을 고려한 립모션 기반 게임 인터페이스의 개발 및 평가,” *한국컴퓨터게임학회 논문지*, 27(4), pp. 69-75, 2014. 12.
- [7] 문현철, 양안나, 김재곤, “웨어러블 응용을 위한 CNN 기반 손 제스처 인식,” *방송공학회 논문지*, 23(2), pp. 246-252, 2018. 3.
- [8] A. Sinha, C. Choi and K. Ramani, "DeepHand: Robust hand pose estimation by completing a matrix imputed with deep features," *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4150-4158, 2016.
- [9] P. Molchanov, S. Gupta, K. Kim and J. Kauts, "Hand gesture recognition with 3D convolutional neural networks," *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp.1-7, 2015.
- [10] G. Huang, Z. Liu, K. Q. Weinberger and L. van der Maaten, "Densely connected convolutional networks," *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp 3-11, 2017.
- [11] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke and A. Rabinovich, "Going deeper with convolutions," *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1-9, 2015.
- [12] K. He, X. Zhang, S. Ren and J. Sun, "Deep residual learning for image recognition," *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770-778, 2016.
- [13] M. Lin, Q. Chen and S. Yan, "Network In Network," *arXiv preprint arXiv:1312.4400*, 2013.
- [14] Diederik P. Kingma, Jimmy Lei Ba, "Adam: A method for stochastic optimization," *International Conference on Learning Representations (ICLR)*, pp. 1-15, 2015.

## 〈저자소개〉



### 이 병 희

- 2012. 3 ~ 현재 서경대학교 컴퓨터공학과 학사 재학
- 관심 분야 : 가상현실, 게임프로그래밍, 컴퓨터 비전



### 오 동 한

- 2012. 3 ~ 현재 서경대학교 컴퓨터공학과 학사 재학
- 관심 분야 : 가상현실, 게임프로그래밍, 컴퓨터 비전



### 김 태 영

- 1991. 2 이화여자대학교 전자계산학과 학사
- 1993. 2 이화여자대학교 전자계산학과 석사
- 1993. 3 - 2002. 2 한국통신 멀티미디어연구소 선임연구원
- 2001. 8 서울대학교 전기컴퓨터 공학부 박사
- 2002. 3 - 현재 서경대학교 컴퓨터공학과 부교수
- 관심 분야 : 실시간 렌더링, 증강현실, 영상처리, 모바일 3D
- 관심분야: 의료 인공지능, 의료영상분석, 의료영상처리