

# 소셜 로봇의 신뢰성 향상을 위한 클래스별 문턱값 기반 분포 외 데이터 탐지

## Out-Of-Distribution Based on Class-wise Thresholding for Enhancing Social Robot Trustworthiness

황지현<sup>1</sup>·장민수<sup>†</sup>  
Jihyun Hwang<sup>1</sup>, Minsu Jang<sup>†</sup>

**Abstract:** In the realm of social robots, ensuring accurate and reliable recognition of diverse environmental stimuli is crucial for effective interaction. Detecting Out-Of-Distribution (OOD) data is vital for improving system reliability by recognizing and responding to OOD data. Existing studies typically use a single threshold to detect OOD data. However, this method fails to reflect the differences in characteristics and data distribution between classes, leading to performance degradation. To address this issue, we propose a class-wise confidence thresholding that accounts for the differences in data distribution across classes and an automatic thresholding based on grid search. Experiments with 12 datasets, including ImageNet and SUN, demonstrated the effectiveness of the proposed method for OOD detection and open-set recognition. In OOD detection, AUROC increased by up to 1.71 and FPR95 decreased by up to 12.79%p compared to existing methods. In open-set recognition, average F1-Score was 98%, and average accuracy was 95%. The class-wise confidence thresholding and grid search-based automatic thresholding can contribute to increasing the reliability of AI robot systems and have advantages of being flexibly applied in various situations.

**Keywords:** Social Robotics, Deep Learning, Out-Of-Distribution Detection, Uncertainty Estimation

### 1. 서론

최근 소셜 로봇은 인간과의 상호작용을 통해 사회적 및 정서적 지원을 제공하는 능력으로 많은 관심을 받고 있다. 이러한 로봇들은 교육, 의료, 고객 서비스 등 다양한 분야에서 활용되고 있다. 예를 들어, 교육용 로봇은 학생과의 상호작용을 통해 맞춤형 교육을 제공하여 학습 효과를 극대화한다. 의료 로봇은 환자의 정서적 안정을 돕고 심리 치료를 보조하며, 고객 서비스 로봇은 사용자와의 상호작용을 통해 고객 만족도를 높인다<sup>[1]</sup>. 이러한 소셜 로봇의 기

능은 딥러닝 알고리즘의 발전에 크게 의존하고 있다. 딥러닝은 인공 신경망을 기반으로 대규모 데이터에서 복잡한 패턴을 학습할 수 있는 능력을 가지고 있어, 이미지 인식, 자연어 처리, 음성 인식 등의 분야에서 높은 성능을 보여주고 있다<sup>[2]</sup>.

딥러닝 알고리즘은 여러 층으로 구성된 인공 신경망을 통해 입력 데이터의 특징을 학습하고, 이를 바탕으로 문제를 해결하는 데 중점을 둔다. 예를 들어, 합성곱 신경망(CNN)은 이미지 처리에 주로 사용되며, 순환 신경망(RNN)은 시퀀스 데이터 처리에 적합하다<sup>[3]</sup>. 이러한 기술적 발전은 다양한 산업 분야에서 인공지능 로봇의 응용 가능성을 높이고 있다.

그러나 이러한 시스템은 학습 데이터와 다른 데이터를 접했을 때 예측 불확실성을 효과적으로 판단하지 못하는 문제가 있다<sup>[4]</sup>. 인공지능 모델은 주로 대규모 데이터셋을 기반으로 학습되지만, 실제 환경에서는 학습 데이터와 다른 데이터가 발생할 수 있다. 이로 인해 예측 불확실성이 증가하며, 특정 상황에서 로봇이 오작동하거나 예기치 않은 결과를 초래할 수 있음을 시사한다. 예를 들

Received : Jul. 31. 2024; Revised : Sep. 9. 2024; Accepted : Oct. 2. 2024

※ This work was supported by Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (No. 2022-0-00951, Development of Uncertainty-Aware Agents Learning by Asking Questions).

1. Student Researcher, School of ETRI University of Science and Technology (UST), Daejeon, Korea (hjh13210@gmail.com)

† Principal Researcher, Corresponding author: Social Robotics Laboratory Electronics and Telecommunications Research Institute, Daejeon, Korea (minsu@etri.re.kr)

어, 인공지능 모델이 특정 동물을 다른 동물로 오 분류하는 경우가 있는데, 이는 학습된 데이터와 다른 분포에서 온 데이터를 접했을 때 발생하는 문제이다. 이러한 문제는 모델의 예측에 불확실성이 존재함을 나타낸다<sup>5)</sup>.

실세계에서 인공지능 로봇 시스템이 활용되는 경우에도 인식 오류로 인한 다양한 문제가 발생할 수 있다. 예를 들어, 자율주행 자동차가 도로 위의 물체를 잘못 인식하여 사고를 유발하는 등 실 세계에서 안정적이고 신뢰성 있게 작동하지 못할 경우 발생할 수 있는 위험성을 보여준다. 이러한 사례는 인공지능 로봇의 예측 불확실성을 효과적으로 관리하고 감지할 필요성을 강조한다<sup>6)</sup>. 특히, 기존의 소셜 로봇은 상황을 제대로 인식하지 못해도 이를 모른 채 수행하려는 경향이 있어, 로봇이 스스로 예측 불확실성을 인지하고 적절히 대응하는 인간-로봇 상호작용이 필요하다. 따라서, 모델의 예측에 존재하는 불확실성을 측정하는 방법이 필요하다. 불확실성 측정을 통해 로봇의 잘못된 예측을 인식하고 적절히 대응함으로써 실세계에서의 문제를 줄이고 신뢰성을 높일 수 있다. 이는 인공지능 로봇 시스템의 신뢰성과 안전성을 높이는 중요한 과제이다. 특히, 다양한 분야에서 인공지능 로봇의 적용이 확대됨에 따라 예측 불확실성을 관리하는 능력은 그 중요성이 더욱 부각되고 있다. 이처럼 인공지능 로봇이 예측 불확실성을 효과적으로 인식하고 관리하는 능력을 가지는 것은 다양한 응용 분야에서 인공지능의 활용성을 극대화하는 데 필수적이다. 이는 인공지능 기술이 로봇 시스템에 보다 안전하고 신뢰성 있게 적용될 수 있도록 하여, 사회 전반에 걸쳐 긍정적인 영향을 미칠 수 있다.

## 2. 관련 연구

### 2.1 로봇 시스템의 분포 외 탐지

소셜 로봇은 인간과의 상호작용을 주 목적으로 설계된 로봇으로 일상 생활에서 다양한 역할을 수행할 수 있도록 설계 되어있다. 소셜 로봇은 공공장소, 교육, 가정 등 다양한 환경에서 활용될 수 있으며 인간과 자연스럽게 상호작용할 수 있는 능력을 갖추고 있다. 최근 연구들은 소셜 로봇의 시각적 인식 및 상호작용 능력을 향상시키기 위한 다양한 기술의 적용을 시도하고 있다. 예를 들어, Pepper 로봇에 딥러닝 기술을 적용하여 특정 사용자를 식별하고 추적하는 능력을 강화한 연구가 있다. 이 연구에서 로봇은 사용자와 상호작용할 때 발생할 수 있는 다양한 상황에 대해 더 정확하고 빠르게 반응할 수 있도록 설계된다. 특히 복잡한 공간이나 가정 환경에서 로봇이 사용자의 표정이나 제스처를 인식하고 이에 맞춰 적절한 반응을 보이는 것은 로봇의 신뢰성을 높이고 사용자의 만족도를 향상시키는 중요한 요소이다<sup>7)</sup>.

또한 가정 및 의료 환경에서 사용되는 사회적 로봇의 신뢰성과 안정성을 보장하는 연구도 중요하게 다뤄지고 있다. 로봇이 제공

하는 정보나 서비스에 대한 사용자의 신뢰가 너무 높거나 낮으면 예상치 못한 문제나 위험이 발생할 수 있다. 의료용 로봇이 환자에게 제공하는 약물 정보에 대한 과도한 신뢰는 오용을 초래할 수 있으며, 이는 환자의 건강 문제를 일으킬 수 있다. 따라서 로봇의 정보 제공과 관련된 윤리적 문제를 해결하고, 신뢰성 있는 정보를 제공하는 것이 중요하다<sup>1)</sup>.

더불어, 소셜 로봇의 중요한 기능 중 하나로 분포 외 탐지가 있다. 이는 로봇 시스템이 훈련 중에 접하지 못한 새로운 환경이나 데이터를 인식하고 적절히 대응할 수 있도록 한다. 로봇 시스템이 복잡하고 예측 불가능한 실제 환경에서 안정적으로 작동하기 위해서는 분포 외 데이터를 효과적으로 감지하고 관리하는 능력이 필수적이다. 최근 연구들은 다양한 방법론을 통해 이 문제에 접근하고 있다. 예를 들어, Farid<sup>8)</sup>는 PAC-Bayes 이론을 활용하여 로봇이 훈련된 환경과 다른 환경에서 운영될 때 이를 감지하는 분포 외 탐지 기법을 개발했다. 이 연구는 훈련 데이터에서 성능이 보장되는 문턱값을 제공하고, 이 성능 문턱값을 넘지 못하는 상황이 발생하면 로봇이 분포 외 환경에서 운영되고 있다고 판단한다. 이 접근법은 로봇 그리퍼와 드론을 이용한 실험을 통해 신뢰구간과 가설 검정을 진행하여 분포 외 데이터를 효과적으로 탐지함과 동시에 통계적 보장을 제공하여 안전성을 강화하고자 하였다.

또 다른 연구에서 Yuhas<sup>9)</sup>는 실시간 임베디드 자율 로봇 플랫폼인 Duckie Bot에 딥 뉴럴 네트워크 기반의 분포 외 탐지기를 구현하고, 이 탐지기의 성능을 평가했다.  $\beta$ -VAE(변형 오토인코더)를 기반으로 한 이 분포 외 탐지기는 로봇 운영 시스템(ROS)과 Docker를 사용한 소프트웨어 프레임워크 내에서 구현되었다. 실험에서는 주로 nuScenes 데이터셋을 사용하여  $\beta$ -VAE 모델을 학습시키고, 이를 실제 데이터로 미세 조정하여 차선을 가로막는 장애물을 감지하는 시스템을 테스트하였다. 이 시스템은 로봇이 장애물을 인식하고 정지하는 능력을 평가하여 실시간 성능과 안전성을 강화하는데 기여하였다.

이러한 연구들은 로봇 시스템의 신뢰성과 유연성을 향상시키며 향후 실용화를 위한 다양한 기술적 도전 과제들을 극복하기 위한 노력을 계속하고 있다.

### 2.2 신경망 모델의 분포 외 탐지

분포 외 탐지는 신경망 모델의 신뢰성을 높이기 위한 중요한 연구 분야로 자리 잡았다. 특히, 모델의 출력 컨피던스를 사용하는 방법이 가장 많이 사용되고 있다. 초기 연구에서는 소프트맥스 출력을 분석하여 분포 외 데이터를 탐지하는 MSP (Maximum Softmax Probability) 방법을 제안하였다. 이 방법은 소프트맥스 출력의 최댓값을 컨피던스 점수로 사용하여, 낮은 점수를 가진 데이터를 분포 외 데이터로 간주하는 방식이다. 이 접근 방식은 간단하고 계산 효율적이라는 장점이 있어 다양한 신경망 모델을 사용

하는 방법의 기반이 되고 있지만, 모델의 예측 확률이 낮더라도 분포 내 데이터로 오인될 수 있다는 단점이 있다<sup>[4]</sup>.

이를 보완하기 위해 Gal<sup>[10]</sup>는 드롭아웃(dropout)을 활용한 베이지안 신경망을 통해 불확실성을 추정하는 방법을 개발하였다. 이 방법은 학습 과정에서 드롭아웃을 사용하여 여러 신경망을 샘플링하고, 그 예측 분포를 통해 불확실성을 추정함으로써 더 정교한 분포 외 탐지를 가능하게 했다. 그럼에도 불구하고, 베이지안 신경망은 계산 비용이 높다는 단점이 있다. 이를 해결하기 위해 Liang<sup>[11]</sup>은 온도 스케일링과 입력 사전처리를 결합한 ODIN (Out-of-Distribution detector for Neural networks)을 소개하였다. ODIN은 소프트맥스 출력의 확률을 조정하여 분포 외 탐지 성능을 한 단계 끌어올렸으며, 이를 통해 다양한 데이터셋에서 높은 성능을 보였다.

이처럼 이전 연구의 한계를 보완하는 과정을 통해 신경망 모델의 분포 외 탐지 기술은 점점 더 정교해지고, 다양한 상황에서의 신뢰성을 높이는 방향으로 발전해오고 있다<sup>[2,6,12]</sup>.

### 2.3 비전-언어 모델을 활용한 분포 외 탐지

비전-언어 모델(Vision-Language Model)은 기존 단일 모달리티 기반의 분포 외 탐지 방법과 달리 텍스트와 이미지 데이터를 동시에 다루며, 분포 외 탐지에 새로운 가능성을 열어주고 있다.

Radford<sup>[12]</sup>은 대규모 비전-언어 모델인 CLIP을 소개했다. CLIP은 이미지를 비전 트랜스포머<sup>[13]</sup>로 처리하고, 텍스트는 트랜스포머 모델로 처리하여 두 모달리티 간의 유사성을 학습한다. CLIP은 이미지 분류 작업에서 텍스트 임베딩과 가장 유사한 이미지 임베딩을 찾는 방식으로 작동하며, 이는 다양한 이미지 분류 작업에

서 강력한 성능을 보여주었다. 그러나 CLIP도 대규모 데이터셋의 필요성과 높은 계산 비용이 문제로 지적되었다<sup>[14]</sup>.

이러한 단점을 보완하기 위해 ALIGN 모델이 제안되었다<sup>[14]</sup>. ALIGN 모델은 18억 개의 이미지-텍스트 쌍을 사용하여 학습함으로써, 텍스트와 이미지 간의 강력한 연관성을 형성한다. ALIGN은 이미지를 EfficientNet으로 처리하고 텍스트는 BERT 기반 모델로 처리한 후, 이 두 결과를 결합하여 이미지와 텍스트 간의 일치도를 높인다. ALIGN은 CLIP<sup>[12]</sup>보다 효율적인 모델 구조를 사용하여 계산 비용을 줄이면서도 높은 성능을 유지하였다.

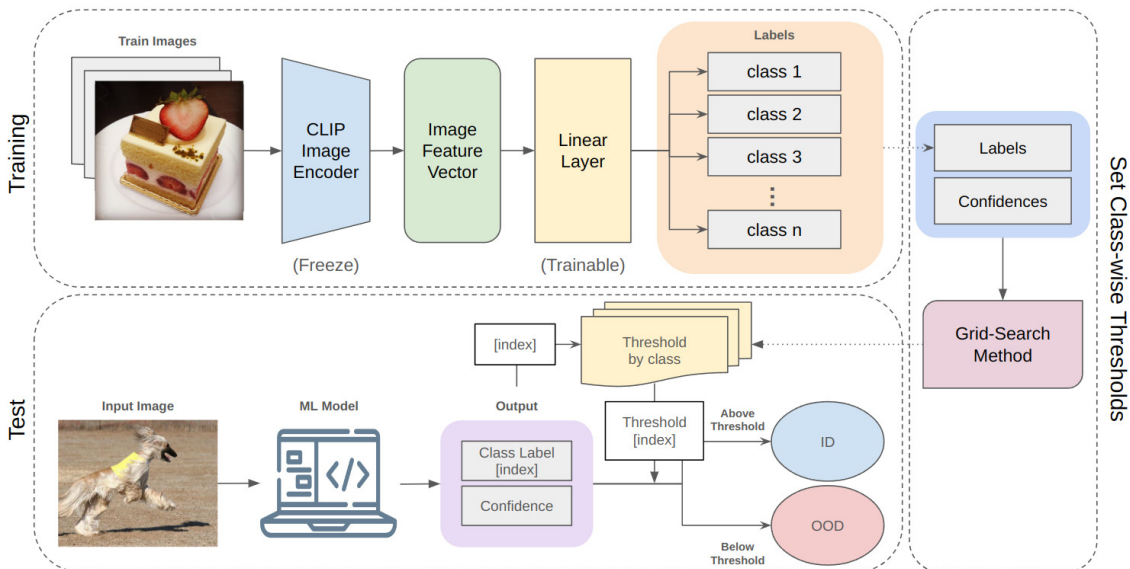
이와 같은 기존 연구들을 바탕으로 MSP<sup>[4]</sup>방법과 비전-언어 모델을 활용하여 인공지능 로봇 시스템 환경에서의 더욱 효율적이고 정확한 분포 외 탐지 방법을 제안한다.

## 3. 클래스별 예측 확률 문턱값 설정과 자동화

본 논문에서는 분포 외 탐지와 오픈셋 인식을 위해 클래스별 예측 확률 문턱값을 자동으로 설정하는 방법을 제안한다. [Fig. 1]은 실험에 사용된 분류 모델의 구조와 문턱값 설정 방법, 분포 외 데이터 탐지 과정을 나타내고 있다. 학습 과정을 통해 분류 모델이 이미지 분류 작업에 최적화되었으며, 학습 데이터를 기반으로 클래스별 예측 확률 문턱값을 설정하였다. 새로운 입력 데이터에 대한 모델의 분류 결과가 설정된 클래스 예측 확률 문턱값을 넘지 못하면 분포 외 데이터로 분류된다.

### 3.1 클래스별 예측 확률 문턱값 설정 방법

기존 방식에서는 다중 클래스 분류에서 단일 문턱값(threshold)



[Fig. 1] Training a classification model and setting class-wise thresholds

을 사용하여 각 클래스를 분류하는 경우가 많다. 그러나 이 접근법은 몇 가지 단점을 가지고 있는데, 주요 단점 중 하나는 데이터셋의 불균형 문제를 충분히 해결하지 못한다는 점이다. 다중 클래스 분류 문제에서 각 클래스의 데이터 분포는 종종 불균형하다. 예를 들어, 한 클래스의 샘플 수가 다른 클래스의 샘플 수보다 훨씬 많을 수 있다. 이러한 상황에서 단일 문턱값을 사용하면 빈도수가 많은 클래스에 대한 모델의 예측 성능은 상대적으로 높아지지만, 빈도수가 적은 클래스에 대한 예측 성능은 현저히 낮아질 수 있다. 이는 결국 전체적인 모델 성능의 저하로 이어질 수 있다<sup>15)</sup>.

이러한 단점을 극복하기 위해 클래스별 문턱값 설정을 도입하였다. 클래스별 문턱값을 사용함으로써 각 클래스의 특성과 데이터 분포에 맞춘 최적의 임계값을 설정하고, 각 클래스의 문턱값을 개별적으로 조절할 수 있는 유연성을 제공하고자 한다. 예를 들어, 특정 클래스에 속하는 이미지를 모델이 잘 분류하지 못하는 경우, 해당 클래스의 문턱값을 낮춰 더 많은 정답 예측이 이루어지도록 조절할 수 있다. 또한, 클래스 간의 특성이 유사한 경우에도 개별 문턱값을 통해 각 클래스의 분류 경계를 명확히 설정할 수 있다. 이를 통해 모델의 성능을 극대화하고, 보다 정교한 분류 결과를 얻을 수 있을 것이라 기대한다.

### 3.2 예측 확률 문턱값 설정의 자동화 방법

기존 다중 클래스 분류 문제에서 단일 문턱값을 경험적으로 0.7, 0.8 같은 값이나 FPR95와 같은 특정 평가지표로 분류하는 방법<sup>16)</sup>과 달리, 본 논문에서는 클래스별로 다른 문턱값을 자동으로 설정하는 방법을 제안한다.

문턱값 설정 방법으로는 격자 탐색(Grid-Search) 방식을 채택하였다. 이 방법은 모든 가능한 매개변수 조합을 체계적으로 검토

하여 가장 최적의 매개변수 조합을 선택하는 고전적인 최적화 방법 중 하나이다. 격자 탐색은 지정된 범위 내의 모든 조합을 평가함으로써 모든 가능성을 탐색하는 장점이 있다. 이를 통해 모델 성능을 향상시키는 이상적인 하이퍼파라미터 값을 결정할 수 있다<sup>17)</sup>. 이 방법을 통해 모든 예측 확률값을 채택 가능한 문턱값 후보로 실험하며, 클래스별로 모델이 정확한 예측과 부정확한 예측을 가장 효과적으로 구분할 수 있는 최적의 예측 확률값을 문턱값으로 채택하였다. 이를 위해 파인튜닝에 사용된 분포 내 데이터셋의 학습 데이터에 대한 모델의 분류 결과를 사용하였으며, 이 데이터에는 정답 라벨, 예측 라벨, 예측 확률값이 포함되어 있다.

[Algorithm 1]은 예측 확률 기반 클래스별 문턱값 설정 알고리즘이다. 이 알고리즘은 입력으로 클래스 번호  $k$ 와 예측 확률값들  $c_1, c_2, c_3, \dots, c_n$ 을 받아 각 클래스에 대해 올바르게 분류되었지만 문턱값보다 낮은 예측 확률을 가지는 이미지 샘플(False-Negative)과 잘못 분류되었지만 문턱값보다 높은 예측 확률을 가지는 이미지 샘플(False-Positive)의 총 개수를 최소화하는 최적의 문턱값을 찾는 과정을 수행한다.

### 3.3 분류 모델의 구조와 학습

제안한 방법은 OpenAI의 CLIP을 활용한다. CLIP은 이미지와 텍스트 간의 복잡한 관계를 학습하며, 다양한 이미지 인식 작업에 널리 사용되는 강력한 모델이다<sup>12)</sup>. 실험에 사용된 모델은 CLIP 이미지 인코더에 선형 레이어(Linear Layer)를 헤드(Head)로 추가하여 이미지 분류를 위해 파인튜닝(Fine-Tuning)된다. 이 선형 레이어는 CLIP 모델이 제공하는 일반화된 특징 추출 기능을 활용하면서 특정 분류 작업에 대한 모델 성능을 최적화한다. 제안한 방법의 성능을 기존 연구와 직접 비교하기 위해 Ming<sup>18)</sup>과 동일한 데이터셋을 이용하여 실험을 수행하였다. 실험 과정에서 CLIP ViT-B/16 모델의 이미지 인코더에 선형 레이어를 추가하고, 이를 10 에포크(Epoch) 동안 파인튜닝하였다.

## 4. 실험 및 결과

### 4.1 실험 방법

본 실험에서는 [Fig. 2]과 같이 모델에 분포 내 데이터 이미지와 분포 외 데이터 이미지를 입력하여 얻은 분류 결과를 사용하였다. 실험에 사용된 8개의 분포 내 데이터셋과 4개의 분포 외 데이터셋은 Ming<sup>18)</sup>와 동일한 데이터셋을 사용하여 기존 연구와의 비교가 용이하도록 하였다. 여기서 모델은 각 입력 이미지에 대해 특정 클래스에 속할 예측 확률을 출력하며, 이 값을 바탕으로 실험을 수행한다. 이미지 분류 문제에서 분포 외 탐지 및 오픈셋 인식 성능을

[Algorithm 1] Grid-search method for fine class-wise threshold

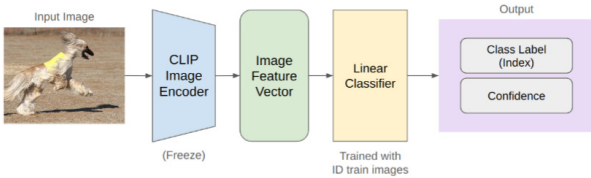
```

1: Input: classes_number, confidence_list
2: Output: optimal_threshold
3: class_number = k
4: confidence_list = [c1, c2, c3, ..., cn]
5: true_positive = {ci | class(ci) = k ∧ predict(ci) = k}
6: false_positive = {ci | class(ci) ≠ k ∧ predict(ci) = k}
7: min_count = ∞
8: optimal_threshold = None
9: for threshold in confidence_list do
10: TP_below_thres = {ci ∈ true_positive | ci > threshold}
11: FP_above_thres = {ci ∈ false_positive | ci < threshold}
12: count = card(TP_below_thres)+card(FP_above_thres)
13: if count < min_count then
14:   min_count = count
15:   optimal_threshold = threshold
16: end if
17: end for
    
```

평가하기 전에, 먼저 모델의 성능을 검증하기 위해 이미지 분류 성능을 측정하였다. 이를 통해 모델이 분포 외 데이터셋에서 얼마나 잘 작동하는지를 평가할 수 있다. 모델의 분류 성능을 검증한 후, AUROC<sup>[19]</sup>와 Total FPR95 지표를 사용하여 분포 외 탐지 성능을 평가하였다. 또한, 전체 데이터에 대해 AUROC 및 FPR95를 계산한 결과와 각 클래스를 고려한 AUROC 및 Total FPR95 결과를 비교하여, 클래스별 차이를 고려한 평가가 더 유의미한지 확인하였다.

격자 탐색 방법으로 구한 클래스별 문턱값이 더 나은 성능을 제공하는지 확인하기 위해, 단일 문턱값과 클래스별 문턱값을 사용하여 분포 외 탐지 문제에서 FPR을 비교하였다. 먼저, 단일 문턱값과 클래스별 문턱값을 설정하고 해당 문턱값 지점에서의 FPR을 측정하여 클래스별 문턱값이 단일 문턱값보다 분포 외 탐지 성능을 향상시키는지 분석하였다. 격자 탐색 방법의 유용성을 평가하기 위해 클래스별로 FPR95 지점에서의 예측 확률 문턱값도 설정하여 비교하였다. 여기서 FPR95 지점에서의 문턱값을 ATPR95 (At True Positive Rate 95%)라 지칭하였다. 클래스별 ATPR95 문턱값을 설정하기 위해, 먼저 각 클래스의 TPR이 95%가 되는 예측 확률값을 문턱값으로 설정하였으며, 격자 탐색 방법을 통해 얻은 최적의 문턱값과 평가 지표 FPR95 지점에서의 문턱값(ATPR95)을 비교하였다. 또한, 단일 문턱값은 클래스별이 아니라 전체 데이터에 대해 격자 탐색과 ATPR95를 사용하여 구하였다. 이렇게 구한 단일 문턱값과 각 클래스에 대해 개별적으로 최적화된 문턱값을 사용하는 클래스별 문턱값 설정 방법의 성능을 비교하였다.

오픈셋 인식에서는 [Fig. 3]과 같이 모델이 입력에 대해 예측한 클래스의 문턱값을 기준으로 이를 넘지 못하면 해당 입력을 분포 외 데이터로 분류하고, 문턱값을 넘으면 모델이 예측한 클래스로 분류한다. 그런 다음, 예측한 클래스를 실제 정답 라벨과 비교하여 평가를 진행한다. 이러한 방식은 모델이 알려진 클래스와 알려지지 않은 클래스를 효과적으로 구별할 수 있는지를 측정하기 위한 것이다.



[Fig. 2] Image classification process of the classification model



[Fig. 3] Open-set recognition image classification method

## 4.2 실험 결과

### 4.2.1 분류 모델 학습 결과

파인튜닝된 분류 모델의 성능을 검증하기 위해 최고 수준 기술의 모델과 이미지 분류 정확도를 비교하였다. [Table 1]에는 두 모델 간의 분류 정확도 결과가 나와 있으며, 이를 통해 제안한 모델의 성능을 평가하였다.

추가로, [Table 2]에서는 Recall, Precision, F1-score 평가 지표를 통해 제안한 모델의 성능을 더욱 심도 있게 분석하였다<sup>[19]</sup>. 평가 결과, 제안한 모델은 이미지 분류 정확도 뿐만 아니라 Recall, Precision, F1-score에서도 의미 있는 성능을 보였다. 이는 제안한 모델이 다양한 이미지 인식 분류 작업에서의 활용 가능성이 충분함을 확인하였다.

### 4.2.2 분포 외 탐지 성능 평가 결과

[Table 3]은 다양한 분포 내 데이터셋 및 분포 외 데이터셋에 대한 파인튜닝된 분류 모델의 분포 외 탐지 성능을 나타낸다. 실험을 통해 구한 데이터셋별 AUROC와 FPR95를 MCM방법<sup>[18]</sup>과 비교한 결과를 보여주며, 본 연구에서는 클래스별로 ROC를 구하는 문턱값을 적용하되 AUROC는 전체 데이터를 고려하여 계산한 결과를 제시하였다. 실험 결과 클래스별로 데이터를 고려한 방법이 평균적으로 더 높은 성능을 보였다. 베이스라인 방법과 비교할

[Table 1] Image classification accuracy between State-Of-The-Art and fine-tuned models

Dataset	Model	
	SOTA	Ours
CUB-200	92.95	90.86
Stanford-Cars	99.85	94.29
Food-101	98.22	92.22
Oxford-Pets	98.22	96.13
ImageNet-10	None	99.99
ImageNet-20	None	97.3
ImageNet-100	86.96	85.46

[Table 2] Image classification performance of fine-tuned models

Dataset	Recall	Precision	F1-score	Accuracy
CUB-200	90.88	91.6	90.84	90.86
Stanford-Cars	94.25	94.51	94.27	94.29
Food-101	92.22	92.3	92.21	92.22
Oxford-Pets	90.08	90.24	89.99	90.13
ImageNet-10	99.99	99.99	99.99	99.99
ImageNet-20	97.3	97.32	97.29	97.3
ImageNet-100	85.46	85.6	85.31	85.46

[Table 3] OOD detection performance with various datasets

ID dataset	OOD dataset							
	iNaturalist		SUN		Places		Texture	
	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC
	Baseline <sup>[3]</sup> /Ours							
CUB-200	9.83/ <b>6.65</b>	98.24/98.74	4.93/ <b>2.14</b>	99.1/ <b>99.52</b>	6.65/ <b>2.42</b>	99.57/98.76	6.97/ <b>0.62</b>	98.75/ <b>99.71</b>
Stanford-Cars	0.05/3.04	99.77/99.76	0.02/2.11	99.95/99.9	0.24/3.09	99.89/99.46	0.02/1.78	99.96/99.91
Food-101	0.64/4.89	99.78/99.65	0.9/2	99.75/ <b>99.87</b>	1.86/2.33	99.58/ <b>99.79</b>	4.04/5.77	98.62/97.79
Oxford-Pets	2.85/ <b>1.06</b>	99.38/ <b>99.7</b>	1.06/ <b>0.12</b>	99.73/ <b>99.9</b>	2.11/ <b>1.69</b>	99.56/99.52	0.8/ <b>0.46</b>	99.81/99.8
ImageNet-10	0.12/ <b>0.02</b>	99.8/ <b>99.99</b>	0.29/ <b>0.17</b>	00.79/99.15	0.88/ <b>0.8</b>	99.62/99.06	0.04/ <b>0.02</b>	99.9/ <b>99.98</b>
ImageNet-20	1.02/4.41	99.66/99.1	2.55/3.21	99.5/ <b>99.71</b>	4.4/ <b>3.9</b>	99.11/ <b>99.56</b>	2.43/4.56	99.03/ <b>99.15</b>
ImageNet-100	18.13/18.2	96.77/95.1	36.45/ <b>16.36</b>	99.54/96.26	34.52/ <b>21.67</b>	94.36/ <b>96.71</b>	41.2/ <b>22.93</b>	92.25/ <b>96.68</b>

[Table 4] OOD detection performance using ImageNet-1k as ID

Model	OOD dataset									
	iNaturalist		SUN		Places		Texture		Average	
	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC
MOS <sup>[20]</sup>	<b>9.28</b>	<b>98.15</b>	40.63	92.01	76.54	89.06	60.43	81.23	39.97	90.11
Fort et al. <sup>[21]</sup>	15.74	96.51	52.34	87.37	55.14	86.48	51.38	85.54	43.62	88.96
Energy <sup>[22]</sup>	10.62	97.52	30.46	<b>93.83</b>	<b>32.25</b>	<b>93.01</b>	44.35	<b>89.54</b>	<b>29.42</b>	<b>93.5</b>
MSP <sup>[4]</sup>	34.54	92.62	61.18	83.68	59.86	84.1	59.27	82.31	53.71	85.68
MCM <sup>[18]</sup>	28.38	94.95	<b>29</b>	94.14	35.42	92	59.88	84.88	38.17	91.49
Ours	24.77	92.96	35.55	90.59	46.19	85.42	<b>41.93</b>	86.1	34.68	88.77

[Table 5] OOD detection performance based on thresholding methods

Dataset	Method	FPR95	AUROC
CUB-200	Single	13.13	97.47
	Class-wise	<b>2.96</b>	<b>99.28</b>
Stanford-Cars	Single	3.09	99.12
	Class-wise	<b>2.5</b>	<b>99.78</b>
Food-101	Single	15.49	97.3
	Class-wise	<b>3.75</b>	<b>98.28</b>
Oxford-Pets	Single	37.19	93.74
	Class-wise	<b>0.83</b>	<b>99.74</b>
ImageNet-10	Single	1.27	99.92
	Class-wise	<b>0.25</b>	<b>99.89</b>
ImageNet-20	Single	11.79	97.64
	Class-wise	<b>3.52</b>	<b>99.38</b>
ImageNet-100	Single	42.77	92.84
	Class-wise	<b>19.79</b>	<b>96.19</b>
ImageNet-1k	Single	77.47	77.4
	Class-wise	<b>34.68</b>	<b>88.77</b>

때, 제안한 방법은 AUROC에서 최대 1.71 더 높은 성능을 보였으며 FPR95에서도 최대 12.79%p 더 낮은 비율을 기록하였다. 이는 클래스별 데이터를 고려함으로써 분포 외 데이터를 더 효과적으로 탐지할 수 있음을 시사한다. 또한, 다양한 분포 내 데이터셋과 분포 외 데이터셋에 대한 성능 평가를 통해 제안한 방법의 일반화 가능성을 확인할 수 있었다.

추가로, 대규모 작업에서 제안한 방법의 분포 외 탐지 성능을 평가하였다. [Table 4]는 ImageNet-1k 데이터셋을 분포 내 데이터로 사용하여 분포 외 탐지 성능을 평가한 결과를 나타낸다.

[Table 5]는 클래스별 예측 확률의 분포를 고려하여 4개의 분포 외 데이터셋에 대해 평가지표를 계산한 값과 전체 데이터를 기반으로 계산한 값의 평균을 나타낸다. 비교 결과 클래스별 예측 확률의 분포를 고려했을 때 모든 데이터셋에서 AUROC와 FPR95 모두 높은 성능을 보였음을 확인할 수 있었다. 이는 분포 외 탐지에서 클래스별 특성을 반영하는 것이 중요하다는 것을 보여준다. 전체 데이터를 기반으로 계산한 값은 클래스 간의 특성을 충분히 반영하지 못할 수 있기 때문에, 제안한 방법과 같이 클래스별로 접근하는 것이 더 효과적일 수 있음을 시사한다.

또한, 4.1 절에서 언급된 네 가지 문턱값 기준에서의 FPR을 비교해 본 결과, 제안한 격자 탐색 방식 기반의 클래스별 문턱값 설

[Table 6] False-positive rate in OOD detection performance based on thresholding methods

Dataset	Method		FPR95
CUB-200	Single	ATPR95	13.13
		Grid-Search	3.59
	Class-wise	ATPR95	2.96
		Grid-Search	<b>1.69</b>
Stanford-Cars	Single	ATPR95	2.46
		Grid-Search	2.57
	Class-wise	ATPR95	2.5
		Grid-Search	<b>1.41</b>
Food-101	Single	ATPR95	15.49
		Grid-Search	12.04
	Class-wise	ATPR95	3.75
		Grid-Search	<b>2.73</b>
Oxford-Pets	Single	ATPR95	37.19
		Grid-Search	6.56
	Class-wise	ATPR95	0.83
		Grid-Search	<b>0.82</b>
ImageNet-10	Single	ATPR95	1.27
		Grid-Search	5.45
	Class-wise	ATPR95	<b>0.25</b>
		Grid-Search	0.3
ImageNet-20	Single	ATPR95	11.79
		Grid-Search	17.18
	Class-wise	ATPR95	<b>3.52</b>
		Grid-Search	3.55
ImageNet-100	Single	ATPR95	42.77
		Grid-Search	35.16
	Class-wise	ATPR95	19.79
		Grid-Search	<b>11.05</b>

정 방법이 모든 분포 내 데이터셋에서 최대 35.8%p 낮은 결과를 보였다([Table 6] 참조). 이처럼 클래스별 접근 방법은 클래스의 각기 다른 특성을 반영함으로써 분포 외 데이터 탐지 성능을 최적화할 수 있다. 특히, 다수의 클래스가 존재하는 데이터셋에서 더욱 눈에 띄는 성능 차이를 보였다.

#### 4.2.3 오픈셋 인식 성능 평가 결과

[Table 7]은 오픈셋 인식 문제에서 문턱값 설정 방식에 따른 이미지 분류 성능을 나타낸다. 제안한 격자 탐색 기반 클래스별 예측 확률 문턱값 설정 방법은 정확도와 F1-Score에서 가장 높은 성능을 보였다. 모델은 예측 클래스의 문턱값을 기준으로 이를 넘지 못하는 입력을 분포 외 데이터로 분류함과 동시에, 문턱값을 넘는 입력은 예측한 클래스로 분류하였다. 이러한 분류 방식은 모델이 알려진 클래스와 알려지지 않은 클래스를 효과적으로 분류할 수

있음을 보여준다.

이와 같은 결과는 제안한 방법이 단순하면서도 분포 외 탐지와 오픈셋 인식 테스트 모두에서 신뢰할 수 있는 인식 성능을 제공할 수 있음을 시사한다.

#### 4.2.4 분포 외 탐지 사례 분석

분포 외 탐지 및 오픈셋 인식의 정성적 실험 결과로서, 제안한 모델과 방식은 분포 내 데이터에 대해 높은 예측 확률값이 문턱값을 넘어서는 정확한 예측을 제공하였다([Fig. 4] 참조). 반면, 분포 외 데이터는 낮은 예측 확률값을 가지며, 대부분의 데이터 포인트가 설정된 문턱값보다 낮아 분포 외 데이터로 잘 분류되었다([Fig. 5] 참조). 그러나 분포 외 데이터임에도 분포 내 데이터로 분류된 경우가 발생하였다([Fig. 6] 참조).

이러한 오 분류의 주요 원인은 두 가지로 분석되었다. 첫째, 모델이 [Fig. 7(a)]와 같이 예측 확률값의 분포가 모든 예측 확률값에 걸쳐 균등하게 분포하는 특수한 경우이다. 이러한 경우, 클래스의 문턱값이 너무 낮게 설정될 수 있으며, 문턱값이 낮으면 분포 외 데이터가 설정된 문턱값을 초과하여 분포 내 데이터로 잘못 분류될 가능성이 높아진다. 둘째, 분포 외 데이터의 예측 확률값이 문턱값과 큰 차이가 없을 때도 이러한 문제가 발생할 수 있다. 이 경우, 모델이 [Fig. 7(b)]와 같이 분포 외 데이터의 예측 확률값이 문턱값 근처에 있어 모델이 이를 명확하게 분포 외 데이터로 분류하



Dataset: ImageNet-10  
Image class: Bird  
Predicted class: Bird  
Confidence: 0.99  
Class threshold: 0.67

[Fig. 4] Instance of ID dataset image classified as ID data



Dataset: Texture  
Image class: OOD  
Predicted class: OOD  
Confidence: 0.32  
Class threshold: 0.88

[Fig. 5] Instance of OOD dataset image classified as OOD data



Dataset: Places  
Image class: OOD  
Predicted class: Caesar salad  
Confidence: 0.2  
Class threshold: 0.08

[Fig. 6] Instance of OOD dataset image classified as ID data

[Table 7] Classification results of open-set recognition according to thresholding method

Dataset	Method		Recall	Precision	F1-score	Accuracy
CUB-200	Single	ATPR95	89.95	64.42	74.85	75.29
		Grid-Search	92.33	78.12	84.35	84.48
	Class-wise	ATPR95	97.23	94.32	95.76	89.8
		Grid-Search	<b>98.69</b>	<b>97.36</b>	<b>98.02</b>	<b>94.61</b>
Stanford-Cars	Single	ATPR95	81.08	57.49	67.45	79.92
		Grid-Search	80.14	60.38	68.88	80.1
	Class-wise	ATPR95	98.19	96.83	97.38	85.18
		Grid-Search	<b>99.19</b>	<b>97.88</b>	<b>98.03</b>	<b>86.8</b>
Food-101	Single	ATPR95	87.54	81.71	84.54	81.66
		Grid-Search	88.77	80.64	84.38	81.08
	Class-wise	ATPR95	96.65	95.83	95.77	93.88
		Grid-Search	<b>99.05</b>	<b>96.99</b>	<b>97.96</b>	<b>96.71</b>
Oxford-Pets	Single	ATPR95	60.83	74.53	66.81	62.58
		Grid-Search	77.64	78.54	78.12	74.42
	Class-wise	ATPR95	99.3	97.98	87.43	97.59
		Grid-Search	<b>99.64</b>	<b>99.39</b>	<b>99.5</b>	<b>99.15</b>
ImageNet-10	Single	ATPR95	82.21	93.35	87.41	88.07
		Grid-Search	85.58	93.66	89.44	89.75
	Class-wise	ATPR95	<b>99.85</b>	92.1	95.71	92.12
		Grid-Search	<b>99.85</b>	<b>93.53</b>	<b>96.12</b>	<b>93.58</b>
ImageNet-20	Single	ATPR95	76.89	78.32	77.41	77.87
		Grid-Search	82.23	80.46	81.15	81.1
	Class-wise	ATPR95	99.84	<b>97.67</b>	97.53	96.25
		Grid-Search	<b>99.92</b>	97.15	<b>98.15</b>	<b>97.2</b>
ImageNet-100	Single	ATPR95	59.52	60.43	59.93	51.36
		Grid-Search	61.81	61.7	61.51	52.73
	Class-wise	ATPR95	99.64	99.08	98.98	98.79
		Grid-Search	<b>99.72</b>	<b>99.33</b>	<b>99.3</b>	<b>99.11</b>

지 못하고 분포 내 데이터로 잘못 분류할 수 있다.

반면, 모델이 [Fig. 7(c)]와 같이 분포 내 데이터와 분포 외 데이터의 예측 확률 분포를 명확히 구분할 수 있는 경우 적절한 문턱값을 설정하면 오 분류를 효과적으로 방지할 수 있다. 이러한 상황에서는 분포 내 데이터의 예측 확률값이 높은 반면, 분포 외 데이터의 예측 확률값이 낮으므로 모델은 분포 외 데이터를 정확하게 탐지하고 분류할 수 있다.

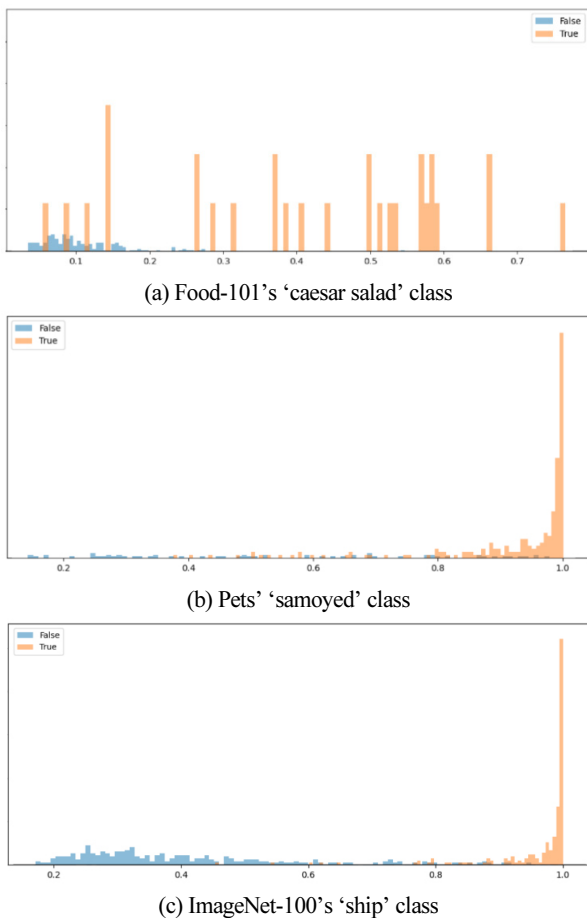
따라서 분포 내 데이터와 분포 외 데이터 간의 예측 확률 분포가 잘 구분되는 경우 문턱값 설정을 통해 모델의 성능을 최적화할 수 있다. 이와 같은 접근법은 오픈셋 인식 문제를 해결하는 데 있어서 중요한 전략이 될 수 있다. 적절한 문턱값 설정을 통해 분포 내 데이터를 정확하게 분류함과 동시에 분포 외 데이터를 효과적으로 구분함으로써 모델의 예측 정확도를 향상할 수 있다.

## 5. 결론

본 연구에서는 소셜 로봇 시스템에서 분포 외 데이터를 효과적으로 탐지하기 위한 새로운 접근 방법을 제안하고 평가하였다. 제안된 방법은 OpenAI의 CLIP 이미지 인코더에 리니어 헤드를 추가하여 학습한 모델을 기반으로 하며, 클래스별 예측 확률을 고려한 격자 탐색 기반의 자동 문턱값 설정 방법을 제안하였다.

실험 결과, 제안된 방법은 다양한 분포 내 데이터셋 및 분포 외 데이터셋에서 기존 방법들과 유사한 성능을 보였다. 이는 제안된 방법이 다양한 클래스 간의 복잡한 관계를 효과적으로 처리할 수 있음을 보여준다. 또한, 클래스별로 특성을 반영한 접근 방식이 분포 외 탐지 성능을 최적화할 수 있음을 입증하였다. 제안된 방법은 한 번 학습되면 이미지 분류, 분포 외 탐지 등 다양한 태스크에





[Fig. 7] Examples of model's confidence distributions by class

추가적인 훈련 없이도 사용할 수 있는 높은 일반화 성능을 가지고 있다. 특히, 클래스별로 예측 확률 문턱값을 자동으로 설정함으로써 각 클래스의 특성에 맞는 더욱 정교한 분류가 가능하다. 또한, 분포 외 데이터에 대해 추가적인 학습 과정이 필요 없으며, 사전 정보 없이도 다양한 상황과 환경에서 유연하게 적용될 수 있다. 이러한 특성 덕분에 제안된 방법은 분류와 이상치 검출 두 가지 작업을 모두 효과적으로 수행할 수 있다. 특히, 많은 클래스를 다루는 실제 응용 환경에서의 적용 가능성을 보여준다. 제안된 방법은 클래스별 특성을 반영하여 분포 외 탐지 성능을 최적화함으로써, 다양한 클래스가 포함된 데이터셋에서도 안정적인 성능을 유지할 수 있다. 이는 다양한 실세계 응용에서 중요한 요소이다. 또한, 분포 외 탐지 평가와 더불어 오픈셋 인식에서도 제안된 방법을 평가하였다. 이 평가에서도 클래스별 예측 확률을 고려한 격자 탐색 기반의 자동 문턱값 설정 방식이 효과적임을 확인하였다. 이는 제안된 방법이 다양한 상황에서 유연하게 적용될 수 있음을 보여준다.

결론적으로, 제안된 클래스별 예측 확률 문턱값 자동 설정 방법은 분포 외 탐지에서 베이스라인보다 AUROC에서 최대 1.71 더 높은 성능을 보였으며, FPR95에서도 최대 12.79%p 더 낮은 비율을 기록하였다. 이는 인공지능 시스템이 분포 외 데이터를 효과

적으로 탐지하고, 다양한 태스크에 높은 일반화 성능을 발휘하며, 실세계 응용 환경에 적용을 통해 인공지능 로봇 시스템의 신뢰성과 안전성을 높이는 데 기여할 수 있음을 나타낸다. 다양한 기존 방법들과의 비교 실험을 통해 제안된 방법의 실용성을 입증하였으며, 특히 클래스별 특성을 반영한 접근 방식이 분포 외 탐지와 오픈셋 인식 성능을 안정적으로 유지할 수 있음을 확인할 수 있었다. 앞으로 본 연구를 기반으로 소셜 로봇을 포함한 다양한 로봇 환경에서의 적용 가능성을 더욱 확장하고, 실세계 데이터에서의 성능을 지속적으로 평가하여 인공지능 로봇 시스템의 신뢰성과 안정성을 더욱 강화할 수 있을 것으로 기대한다.

## References

- [1] M. Salem, G. Lakatos, F. Amirabdollahian, and K. Dautenhahn, "Towards Safe and Trustworthy Social Robots: Ethical Challenges and Practical Issues," *Social Robotics*, vol. 9388, pp. 584-593, Oct., 2015, DOI: 10.1007/978-3-319-25554-5\_58.
- [2] M. Soori, B. Arezoo, and R. Dastres, "Artificial intelligence, machine learning and deep learning in advanced robotics, a review," *Cognitive Robotics*, vol. 3, pp. 54-70, 2023, DOI: 10.1016/j.cogr.2023.04.001.
- [3] H. Liang and Y. Lu, "A CNN-RNN unified framework for intrapartum cardiocograph classification," *Computer Methods and Programs in Biomedicine*, vol. 229, Feb., 2023, DOI: 10.1016/j.cmpb.2022.107300.
- [4] D. Hendrycks and K. Gimpel, "A baseline for detecting misclassified and out-of-distribution examples in neural networks," *arXiv:1610.02136*, 2017, DOI: 10.48550/arXiv.1610.02136.
- [5] B. Settles, "Active learning literature survey," 2009, [Online], <http://digital.library.wisc.edu/1793/60660>.
- [6] F. Kraus and K. Dietmayer, "Uncertainty estimation in one-stage object detection," *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, Auckland, New Zealand, pp. 53-60, 2019, DOI: 10.1109/ITSC.2019.8917494.
- [7] W. Zhong, L. Tian, D. T. Le, and H. Rezafighi, "Improving Visual Perception of a Social Robot for Controlled and In-the-wild Human-robot Interaction," *Improving Visual Perception of a Social Robot for Controlled and In-the-wild Human-robot Interaction*, pp. 1199-1203, Mar., 2024, DOI: 10.1145/3610978.3640648.
- [8] A. Farid, S. Veer, and A. Majumdar, (2022, January). "Task-driven out-of-distribution detection with statistical guarantees for robot learning," *The 5th Conference on Robot Learning*, pp. 970-980, 2022, [Online], <https://proceedings.mlr.press/v164/farid22a.html>.
- [9] M. Yuhas, Y. Feng, D. Jun X. Ng, Z. Rahiminasab, and A. Easwaran, "Embedded out-of-distribution detection on an autonomous robot platform," *The Workshop on Design Automation for CPS and IoT*, pp. 13-18, 2021, DOI: 10.1145/3445034.3460509.
- [10] Y. Gal and Z. Ghahramani, "A theoretically grounded application of dropout in recurrent neural networks," *30th Conference on Neural Information Processing Systems*, Barcelona, Spain, 2016,

- [Online], [https://proceedings.neurips.cc/paper\\_files/paper/2016/file/076a0c97d09cf1a0ec3e19c7f2529f2b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2016/file/076a0c97d09cf1a0ec3e19c7f2529f2b-Paper.pdf)
- [11] S. Liang, Y. Li, and R. Srikant, "Enhancing the reliability of out-of-distribution image detection in neural networks," *arXiv:1706.02690*, 2018, DOI: 10.48550/arXiv.1706.02690.
- [12] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," *The 38th International Conference on Machine Learning*, pp. 8748-8763 2021, [Online], <https://proceedings.mlr.press/v139/radford21a.html>.
- [13] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *The 31st International Conference on Neural Information Processing Systems*, Long Beach, CA, USA, 2017, [Online], <https://user.phil.hhu.de/~cwurm/wp-content/uploads/2020/01/7181-attention-is-all-you-need.pdf>.
- [14] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. Le, Y.-H. Sung, Z. Li, and T. Duerig, "Scaling up visual and vision-language representation learning with noisy text supervision," *The 38th International Conference on Machine Learning*, vol. 139, pp. 4904-4916, 2021, [Online], <https://proceedings.mlr.press/v139/jia21b.html>.
- [15] H. He and E. A. Garcia, "Learning from imbalanced data," *IEEE Transactions on knowledge and data engineering*, vol. 21, no. 9, pp. 1263-1284, Sept., 2009, DOI: 10.1109/TKDE.2008.239.
- [16] M. Li and I. K. Sethi, "Confidence-based active learning," *IEEE transactions on pattern analysis and machine intelligence*, vol. 28, no. 8, pp. 1251-1261, Aug., 2006, DOI: 10.1109/TPAMI.2006.156.
- [17] P. Liashchynskiy and P. Liashchynskiy, "Grid search, random search, genetic algorithm: a big comparison for NAS," *arXiv:1912.06059*, 2019, DOI: 10.48550/arXiv.1912.06059.
- [18] Y. Ming, Z. Cai, J. Gu, Y. Sun, W. Li, and Y. Li, "Delving into out-of-distribution detection with vision-language representations," *Advances in Neural Information Processing Systems 35 (NeurIPS 2022)*, 2022, [Online], [https://proceedings.neurips.cc/paper\\_files/paper/2022/hash/e43a33994a28f746dcfd53eb51ed3c2d-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/e43a33994a28f746dcfd53eb51ed3c2d-Abstract-Conference.html).
- [19] D. M. W. Powers, "Evaluation: from precision, recall and F-measure to ROC, informedness, markedness and correlation," *arXiv:2010.16061*, 2011, DOI: 10.48550/arXiv.2010.16061.
- [20] R. Huang and Y. Li, "Mos: Towards scaling out-of-distribution detection for large semantic space," *arXiv:2105.01879*, 2021, DOI: 10.48550/arXiv.2105.01879.
- [21] S. Fort, J. Ren, and B. Lakshminarayanan, "Exploring the limits of out-of-distribution detection," *Advances in Neural Information Processing Systems*, pp. 7068-7081, 2021, [Online], [https://proceedings.neurips.cc/paper\\_files/paper/2021/hash/3941c4358616274ac2436eacf67fae05-Abstract.html](https://proceedings.neurips.cc/paper_files/paper/2021/hash/3941c4358616274ac2436eacf67fae05-Abstract.html).
- [22] W. Liu, X. Wang, J. Owens, and Y. Li, "Energy-based out-of-distribution detection," *Advances in neural information processing systems*, [Online], <https://proceedings.neurips.cc/paper/2020/hash/f5496252609c43eb8a3d147ab9b9c006-Abstract.html>.



### 황지현

2020 한밭대학교 기계공학과(학사)

2024 과학기술연합대학원대학교 인공지능(석사)

관심분야: Social Robotics, Self-Awareness, Cloud-Robot Intelligence



### 장민수

1992 서강대학교 전산학과(학사)

1994 서강대학교 전산학과(석사)

2015 한국과학기술원 전산학과(박사)

1999~현재 한국전자통신연구원 책임연구원

관심분야: Social Robotics, Self-Awareness, Cloud-Robot Intelligence