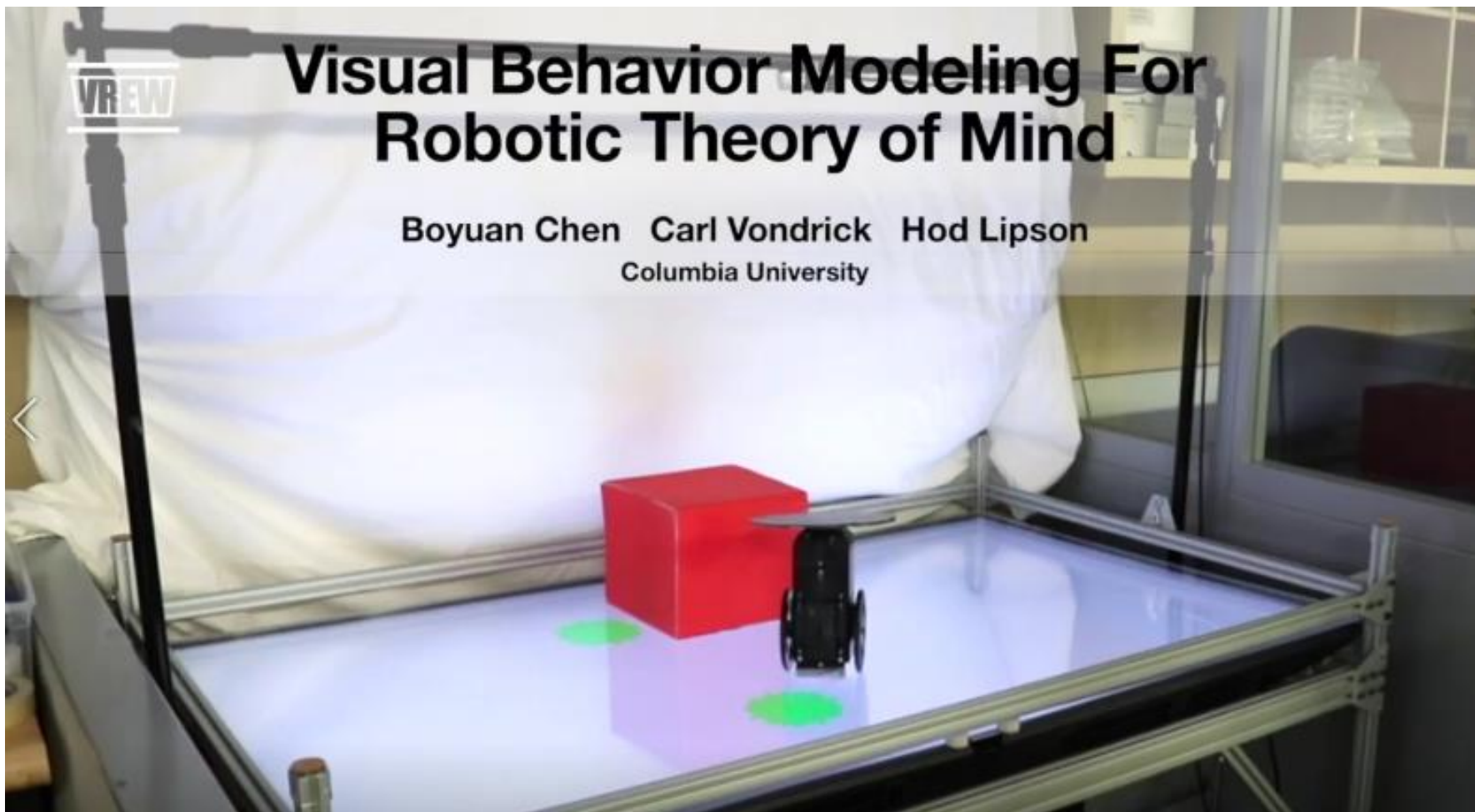


로봇 마음이론을 위한 시각적행동 모델링



Visual behavior modelling이란?

- 로봇 또는 인공지능이 관찰을 통해 다른 에이전트의 행동을 시각적으로 이해하고 예측하는 과정
- 로봇이나 AI가 단순히 명령을 수행하는 것이 아니라, 주변 환경과 상호작용하며 독립적으로 행동을 예측함.
- 로봇은 특정한 심볼이나 코드를 이해하는 것이 아니라, 시각적 정보만을 이용하여 상황을 판단하고 예측함.
- 이는 로봇이 보다 자연스럽게 인간처럼 반응하게 만들어, **인간과 로봇 간의 상호작용을 보다 원활하게 할 수 있는 기술적 기반이 됨.**

시각적 행동 모델링 시나리오

- (상황) 어린이 놀이터에서 로봇이 아이들의 움직임을 관찰하는 상황
- (데이터 수집) 로봇은 각 아이들이 어떻게 움직이는지, 놀이기구에 어떻게 접근하는지를 카메라를 통해 실시간으로 보면서 데이터를 수집
- (예측) 데이터를 기반으로, 로봇은 아이들이 특정 놀이기구에 관심이 많은지, 다음에 어떤 놀이기구로 이동할 가능성이 높은지를 예측함.
- (효과) 로봇의 예측은 로봇이 아이들과 상호작용하거나 아이들의 안전을 감시하는 데 도움이 됨

Robotic theory of mind(로봇 마음이론) ?

- (마음이론) 인간이 타인의 관점을 이해하고 그들의 행동을 예측하는 능력
- (로봇 마음이론) 인간의 'Theory of Mind'를 기계에 적용 + 로봇이나 인공지능 시스템이 다른 에이전트의 의도, 믿음, 지식 등의 정신 상태를 이해하고 예측하는 능력을 개발하는 연구 분야

논문 요약

1. 목적

- 로봇이 보이는 시각정보만으로 미래의 행동을 예측할 수 있는지 확인하는 실험

2. 연구문제

- 로봇이 다른 에이전트의 행동을 시각적으로 모델링 할 수 있는가?

3. 결과

- 로봇(observer)이 다른 로봇(actor)을 관찰하고, 그 로봇이 보이는 시각적 정보만을 사용하여 미래의 행동을 예측

실험 시나리오

1. 로봇은 특정한 행동을 프로그래밍되어 있다.
2. 관찰자 로봇은 해당 행동을 바탕으로 미래의 결과를 시각적으로 예측하게 된다.
3. 이 과정에서 로봇은 심볼이나 명시적인 논리적 추론 없이 순수하게 시각적 정보만을 사용하여 행동을 예측한다.
4. 이 연구는 기계가 어떻게 비 심볼적 정보를 사용하여 복잡한 인지적 기능을 수행할 수 있는지에 대한 이해를 넓힐 수 있는 기회를 제공한다.

논문 개요

1. 본연구는 '행동 모델링이 인간과 동물의 사회적 행동을 이해하는데 필수적인 인지 능력이며, 이를 로봇에 적용하려는 연구'이다.
2. 기계 행동 모델링 연구는 프로그램 또는 선택된 매개 변수 센서 입력에 의존하지만, 이 연구에서는 순수한 시각 처리만을 사용하여 행위자의 행동을 모델링할 수 있다는 가설을 세움
3. 연구팀은 비언어적, 비상징적 로봇 실험을 설계하여, 관찰자 로봇이 행위자 로봇의 초기 장면 이미지만을 기반으로 미래 계획을 시각화할 수 있음을 발견했다.
4. 마음 이론은 3세경에 다른 사람이 다른 세계관을 가질 수 있다는 것을 인식하는 인간의 능력에서 비롯된다.
5. 이 연구에서는 로봇이 이러한 인지 능력을 가질 수 있음을 실험을 통해 증명하려고 시도했다.
6. 연구팀은 행위자와 관찰자로서 두 로봇을 사용하여 실험을 수행했으며, 관찰자 로봇은 행위자 로봇의 행동을 관찰하고 그 의도를 이해하여 미래의 행동을 예측할 수 있었다.
7. 이 연구는 로봇이 인간처럼 복잡하고 적응적인 사회적 상호 작용을 할 수 있도록 하는데 중요한 발전을 이뤘다.

시각적 행동이론(논문에서)

행위자 로봇(검은색 원)은 가장 가까운 먹이(녹색 원)를 향해 움직이도록 프로그램되어 있음.

- **패널 A:** 로봇은 가장 가까운 녹색 원을 볼 수 있음 + 직접 이동.
- **패널 B:** 때때로 녹색 원이 장애물(회색 직사각형)에 의해 가려져 있음+ 로봇은 보이는 가장 가까운 녹색 원을 향해 간접적 이동
- **패널 C:**
 - ✓ 관찰자 로봇은 로봇이 어떤 상황에 처해 있는지 보고 있지만, 무슨 일이 일어날지는 표시되지 않음.
 - ✓ 이는 관찰자가 로봇의 행동을 예측해야 하는 상황.
- **패널 D**
 - ✓ 관찰자 AI는 로봇의 "궤적 흐림"을 시각화하여 제공이는 로봇이 움직일 경로를 시각적으로 예측한 것

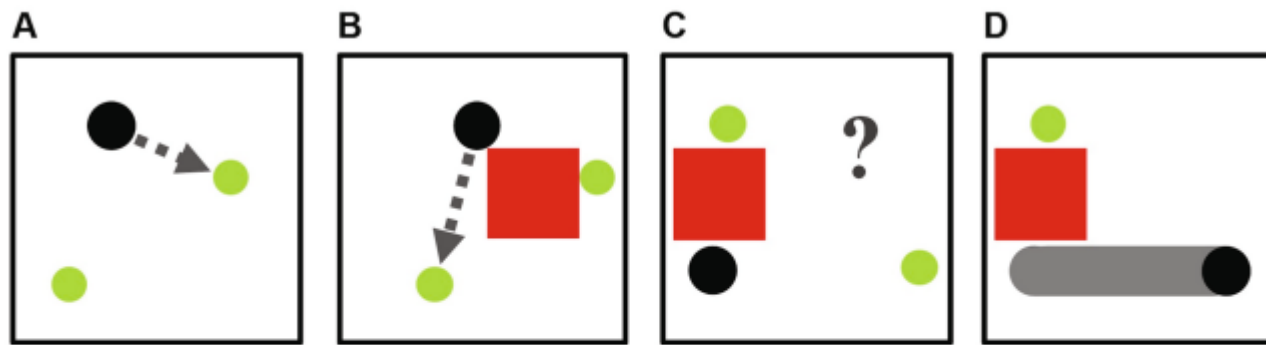


Figure 1. Visual theory of behavior. An actor robot (black circle) is programmed to move towards the nearest food (green circle) that it can see, and consume it. Sometimes (A), the nearest green circle is directly visible to the actor, but sometimes (B) the nearest green circle is occluded by an obstacle. When occluded, the actor will move towards the closest visible circle, if any. After watching the actor act in various situations, an observer-AI learns to envision what the actor robot will do in a new, unseen situation (C). The observer's prediction is delivered as a visualization of the actor robot's "trajectory smear" (D). This entire reasoning process is done visually, sidestepping the need for symbols, logic, or semantic reasoning.

실험 설정

<패널 A> 로봇(행위자)이 플레이펜 안에 있으며, 관찰자, 녹색 음식, 그리고 빨간 장애물이 보임. 이 실험 설정은 로봇이 어떻게 장애물을 우회하여 음식에 도달하는지를 관찰하기 위한 것입니다.

<패널 B> 관찰자가 본 입력샘플

<패널 C> 행위자의 움직임 예측이 포함된 **관찰자 출력** 샘플+검은색 **흐림** 처리는 행위자 로봇이 장애물을 **우회하는 경로를 예측하는** 모습을 시각화한 것

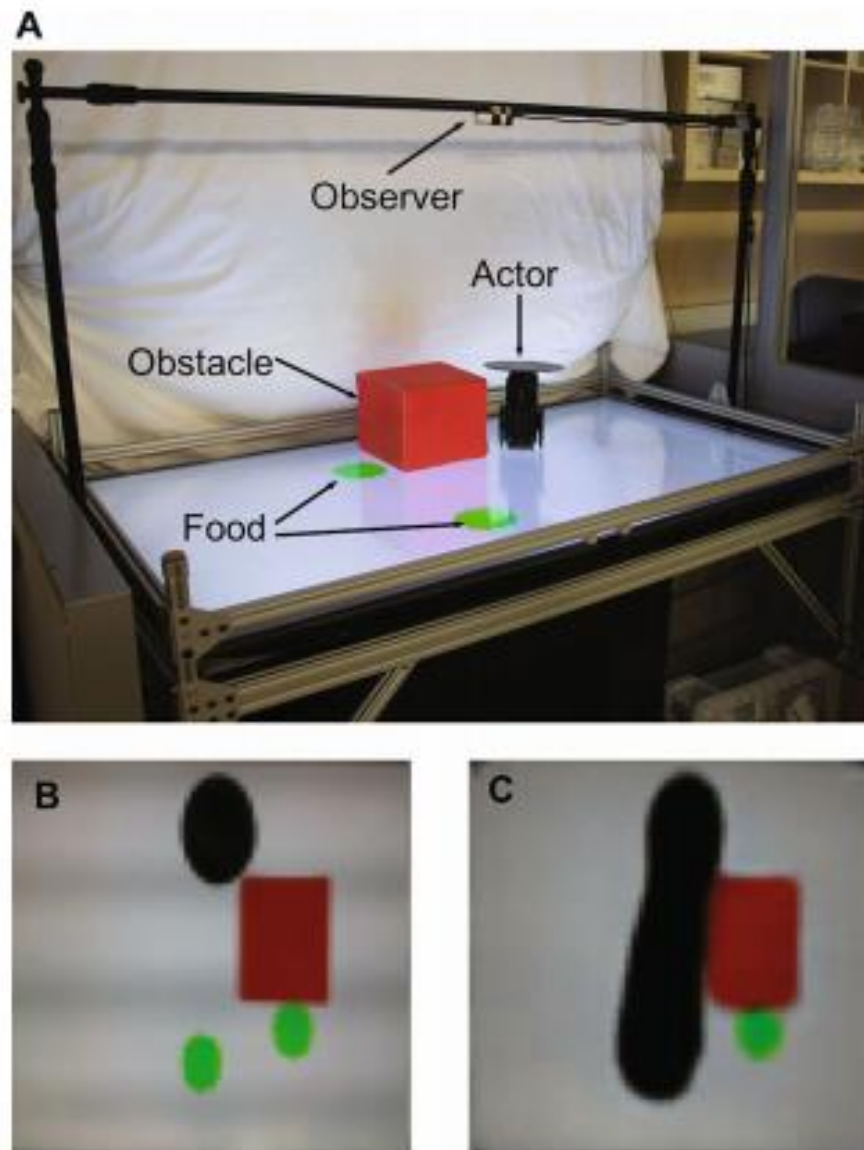


Figure 2. Experimental setup (A) Actor robot in playpen, showing observer, green food, and red obstacle. (B) A sample input image as seen by the observer, and (C) sample output image produced by the observer, which includes a prediction of the motion path of the actor.

<연구 설계>

- 로봇(검정색) : 초록점을 향해 움직임 + 붉은색 장애물을 만나면 움직이지 않음
- 관찰 AI(천정): 로봇이 어떻게, 무엇을 목적으로 행동하는지? 아무런 사전정보 없이 카메라로 관찰(입력값이 없음)
- 연구팀: 관찰 AI에게 로봇과 장애물, 초록점이 보이는 지도만 제공 >>>> 로봇의 경로 예측

Experimental setup and results: (how we test our hypothesis)

실험 설정 및 결과(가설 테스트하는 방법)

1. 실험 문제

- 관찰자가 상징적 추론(입력값+내장된 지식) 없이 행위자 로봇의 미래 계획을 시각적으로 예측할 수 있는가?
- 관찰자가 관점 취득(*perspective taking)과 유사한 능력을 얻을 수 있는가?

2. 행위자 로봇의 사전 프로그램 된 4가지 행동 유형

- a) 직선 행동: 행위자 로봇이 녹색 원을 향해 직선으로 이동
- b) 팔꿈치 행동: 행위자 로봇이 중간에 위치한 비선형 점을 향해 먼저 이동한 후 그곳에서 직선으로 녹색 원으로 이동
- c) 지그재그 행동: 행위자 로봇이 목표로 가는 직선 경로 양쪽에 대칭적으로 배치된 두 개의 중간 비선형 제어점을 통해 원으로 이동
- d) 장애물이 있는 단일 식품 행동: 이 행동은 장애물에 의한 가림으로 인해 행동(a)보다 복잡한 버전

* 다른 사람의 관점에서 상황을 보고 이해하는 능력

<로봇 행동 : 4가지 유형의 동작 프로그래밍>

A



Straight line behavior



Elbow behavior



Zig-zag behavior



Single food and obstacle behavior
when green food is occluded



Single food and obstacle behavior
when green food is visible

1. 직선이동 (로봇 ▶ 녹색원)
2. 엘보 동작(로봇 ▶ 중간위치
경유 ▶ 녹색원)
3. 지그제그 동작
4. 녹색원이 보일때 직선이동
& 녹색원이 보이지 않을 때
이동 x

“천정의 로봇은 사전지식 없이 관찰한다.”

“관찰자의 유일한 정보: 카메라의 입력 이미지
뿐 + 라벨링, 분할, 모터명령, 궤적 좌표 또는
파생정보가 관찰자에게 제공되지 않는다.”

Network architecture

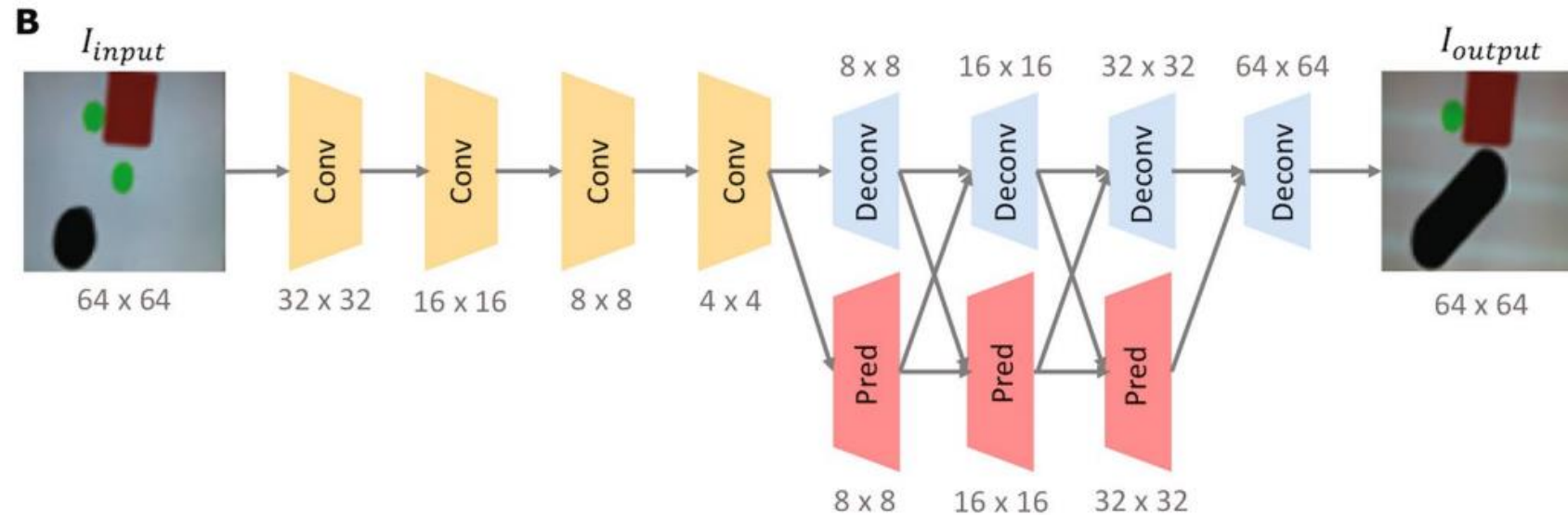


Figure 3. Visualization of the behaviors of the actor robots and Observer network architecture (A) We pre-programmed four types of behaviors for the actor robot. The images shown here are produced by integrating a sequence of frames from a video captured by the top-down camera. The robot path is shown in black, the rectangular obstacle (if any) in red and the goal circles in green or red. (B) The image prediction network is composed of several layers of convolutional units and deconvolutional units. At the deconvolutional stage, we utilize multi-scale prediction to maintain high resolution at the output image. Numbers indicate the dimension of output feature map after each module.

CNN 알고리즘

1. 합성곱 층(Convolution Layers)

- 이미지의 특징을 추출
- 각 합성곱 층은 입력 이미지에 여러 필터를 적용하여 이미지의 **로컬 특징을 학습**
- 필터들은 이미지의 에지, 색상, 질감 등 다양한 시각적 정보를 감지하는 데 사용

2. 풀링 층(Pooling Layers)

- 이미지의 공간 크기를 줄이는데 사용
- 모델의 계산 부담을 줄이고, 과적합을 방지하는 데 도움
- 최대 풀링(Max Pooling)과 평균 풀링(Average Pooling)

CNN 알고리즘

3. 활성화 함수(Activation Functions)

- ReLU(Rectified Linear Unit) CNN에서 자주 사용되는 활성화 함수

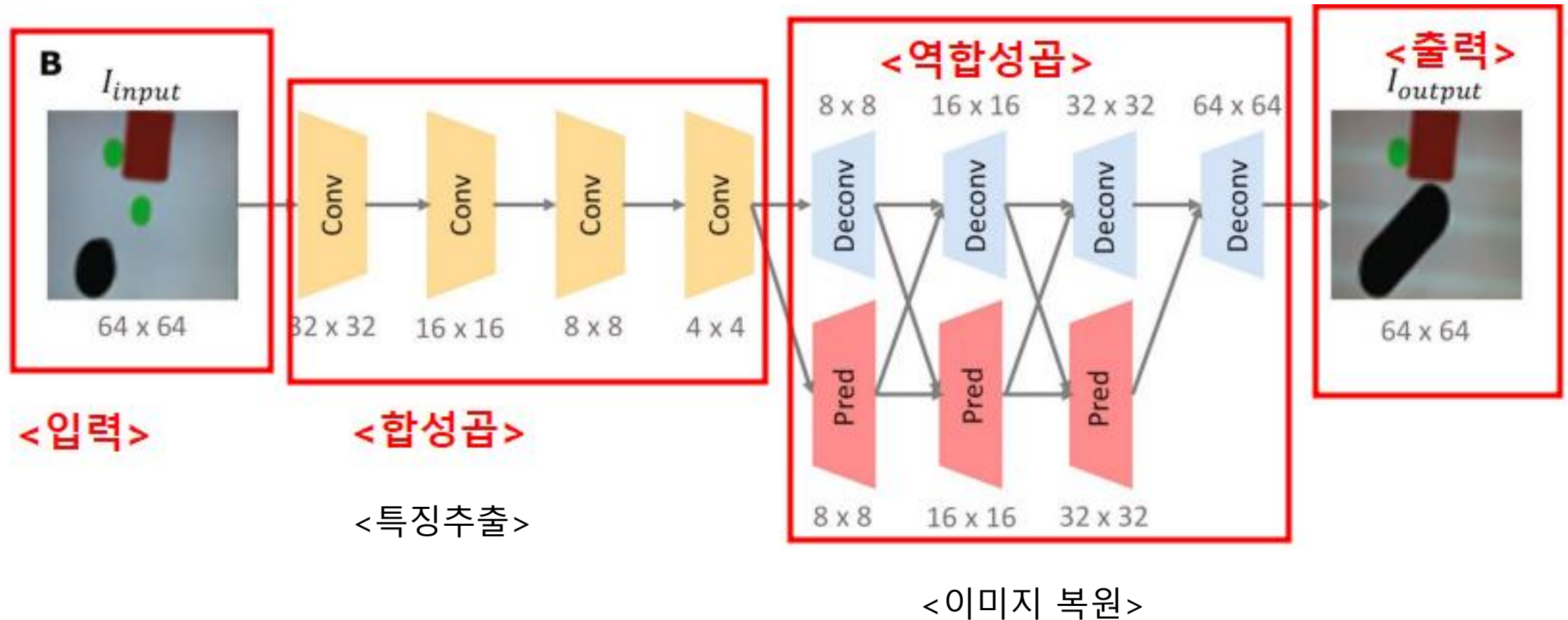
4. 역합성곱 층(Deconvolution Layers):

- 이미지의 해상도를 다시 높이는 역할
- 이 층은 예측된 이미지의 세부적인 부분까지 복원하는 데 중요함

5. 완전 연결 층(Fully Connected Layers)

- 네트워크의 마지막 부분
- 모든 특징을 결합하여 최종 예측 수행
- 분류 작업에서 결정을 내리는 데 사용됨

Network architecture



Results

1. 실험훈련

- 관찰자 AI는 실험적으로 기록된 2400개의 입력-출력 이미지 쌍으로 훈련
- 이중 20%는 테스트에 사용

2. 성능평가

A

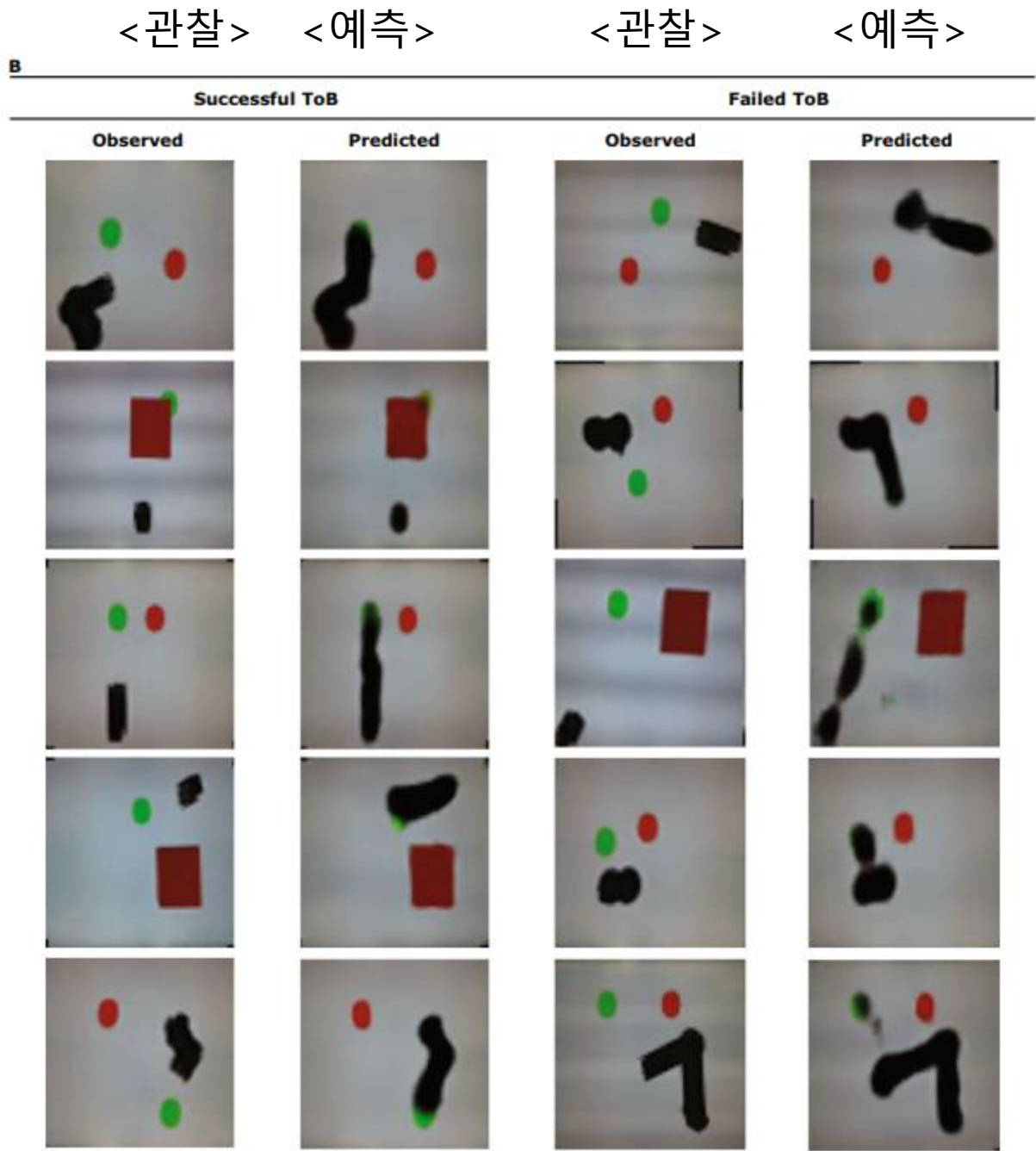
Actor Behavior	Success Rate
Straight line behavior	99.90% (n=980)
Elbow behavior	98.94% (n=944)
Zia-zag behavior	96.70% (n=999)
Single food with obstacle behavior (장애물)	98.52% (n=811)

Results

3. 정확도 계산

- 이미지에서 가장 큰 윤곽과 예측된 궤적, 음식 원의 위치를 추출하여 정확도를 계산
- 계산된 가장 짧은 거리(D_{target})가 행위자 로봇의 지름보다 작거나 같으면 예측을 성공으로 간주

Results



손실률 변화

- 4가지 시나리오 상황

Train / Test

Food Visible Train, Food Visible Test

Food Visible Train, Food Obscured Test

Half Food Visible Train, Half Food Obscured Train, Food Visible Test

Half Food Visible Train, Half Food Obscured Train, Food Obscured Test

Success Rate

97.30% (n=667)

56.82% (n=88)

98.21% (n=669)

100.00% (n=88)

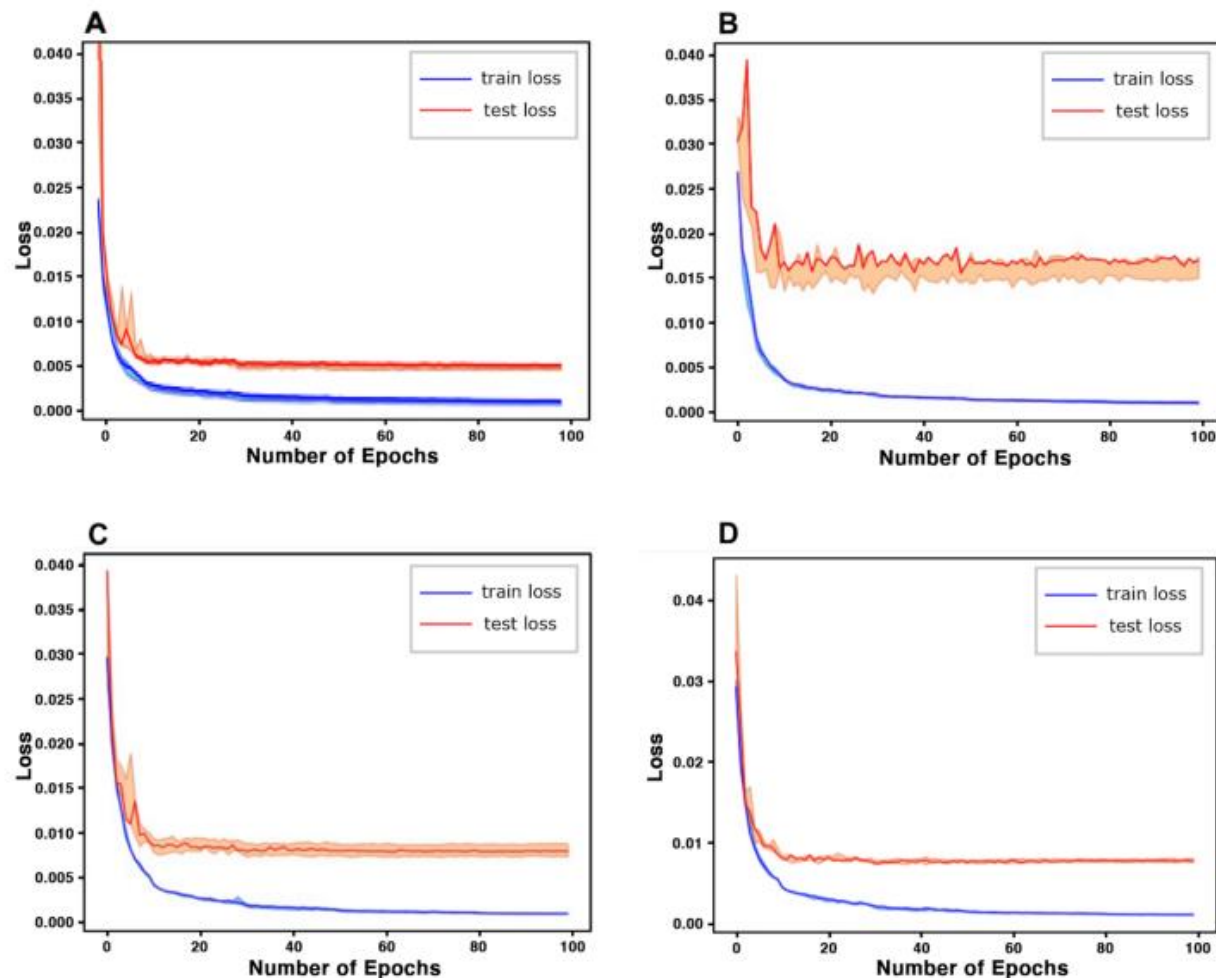


Figure 7. Training and testing of the observer and corresponding success rate We first gathered training and testing data by randomly placing the actor, two green food items, and the obstacle (Table, A). We also collected “obscured” test cases where we deliberately placed the closest food to where it is not visible to the actor (Table, B). Higher success rates were achieved by balancing the training data with half “visible” data and half “obscured” data (Table, C and D). Learning curves across all four above scenarios are shown. Error bars are presented to show experiment results across three different random seeds used in both data splitting and network training.

A패널

- 훈련+테스트: 음식이 보이는 경우
- 초기 손실이 급격히 감소
- 성공률 97.3%

Train / Test

Food Visible Train, Food Visible Test

Food Visible Train, Food Obscured Test

Half Food Visible Train, Half Food Obscured Train, Food Visible Test

Half Food Visible Train, Half Food Obscured Train, Food Obscured Test

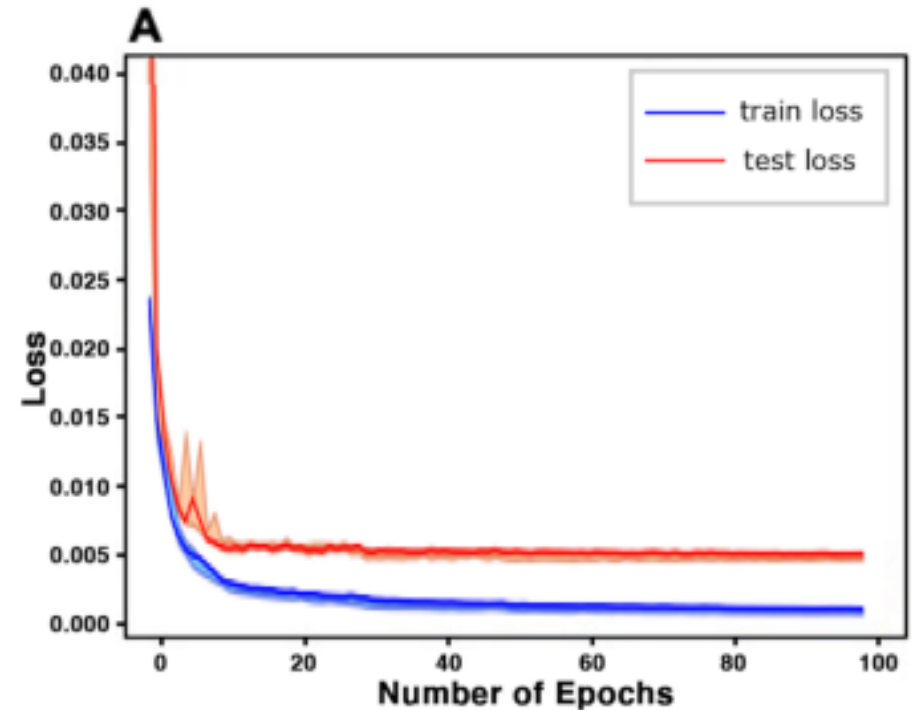
Success Rate

97.30% (n=667)

56.82% (n=88)

98.21% (n=669)

100.00% (n=88)



B패널

Train / Test

Food Visible Train, Food Visible Test

Food Visible Train, Food Obscured Test

Half Food Visible Train, Half Food Obscured Train, Food Visible Test

Half Food Visible Train, Half Food Obscured Train, Food Obscured Test

Success Rate

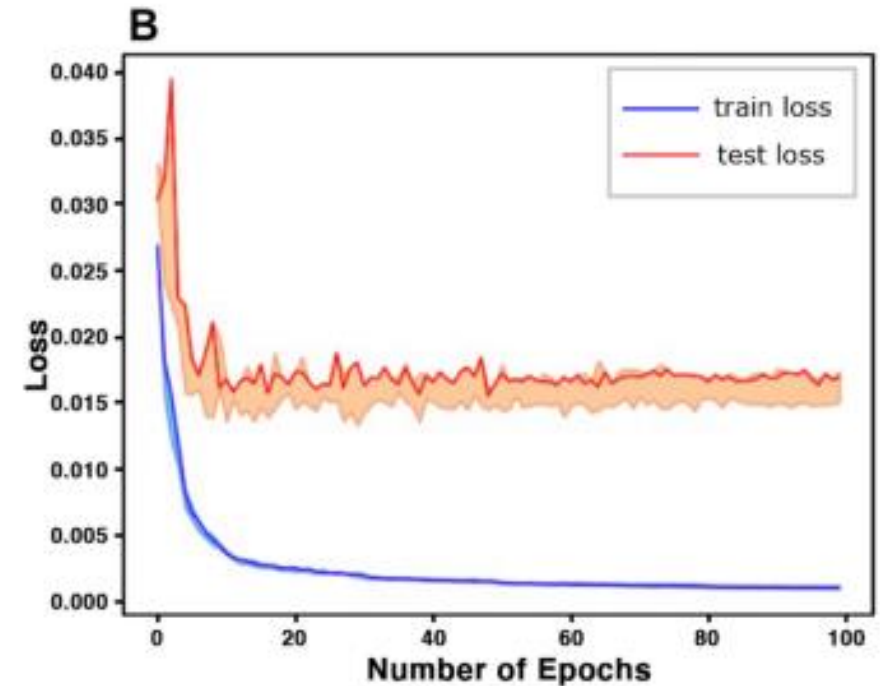
97.30% (n=667)

56.82% (n=88)

98.21% (n=669)

100.00% (n=88)

- 훈련: 음식이 보이지 않는 상황
- 테스트: 음식이 가려진 상황
- 테스트 손실: 초기 높게 시작하여 점차 감소
- 훈련손실에 비해 높은 수준 유지
- 성공률 56.82%



"AI 시스템이 실제 환경에서 마주할 수 있는 다양한 변수와 불확실성을 처리하는 데 있어서의 한계를 드러냄 "

C패널

Train / Test

Food Visible Train, Food Visible Test

Food Visible Train, Food Obscured Test

Half Food Visible Train, Half Food Obscured Train, Food Visible Test

Half Food Visible Train, Half Food Obscured Train, Food Obscured Test

Success Rate

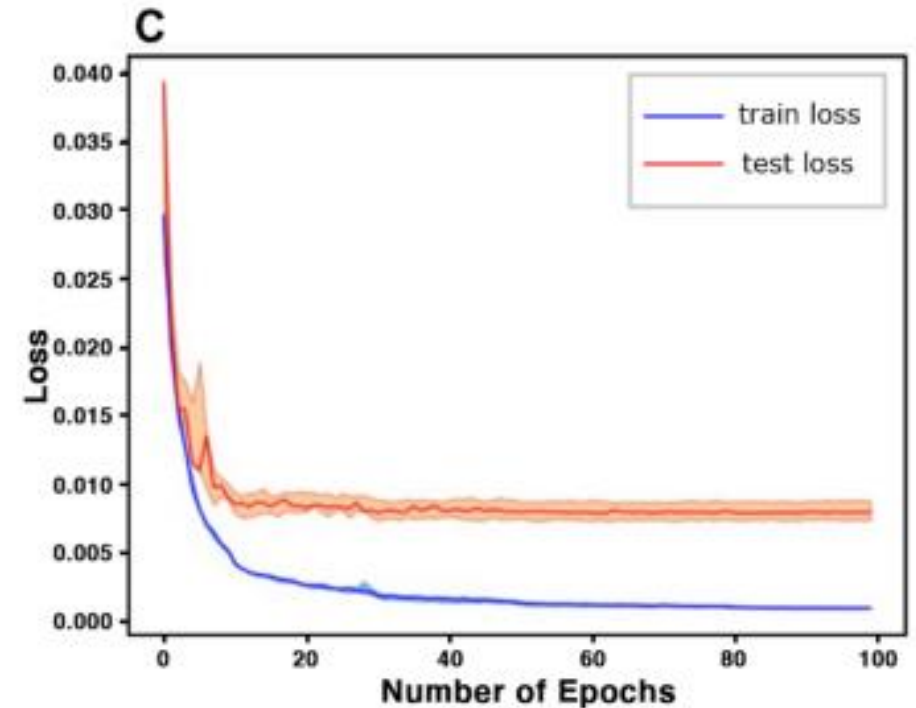
97.30% (n=667)

56.82% (n=88)

98.21% (n=669)

100.00% (n=88)

- 훈련: 음식이 반은 보이고 반은 가려진 상황
- 테스트: 음식이 보이는 상황
- 손실률 안정적으로 낮아짐
- 성공률 98.21%



D패널

Train / Test

Food Visible Train, Food Visible Test

Food Visible Train, Food Obscured Test

Half Food Visible Train, Half Food Obscured Train, Food Visible Test

Half Food Visible Train, Half Food Obscured Train, Food Obscured Test

Success Rate

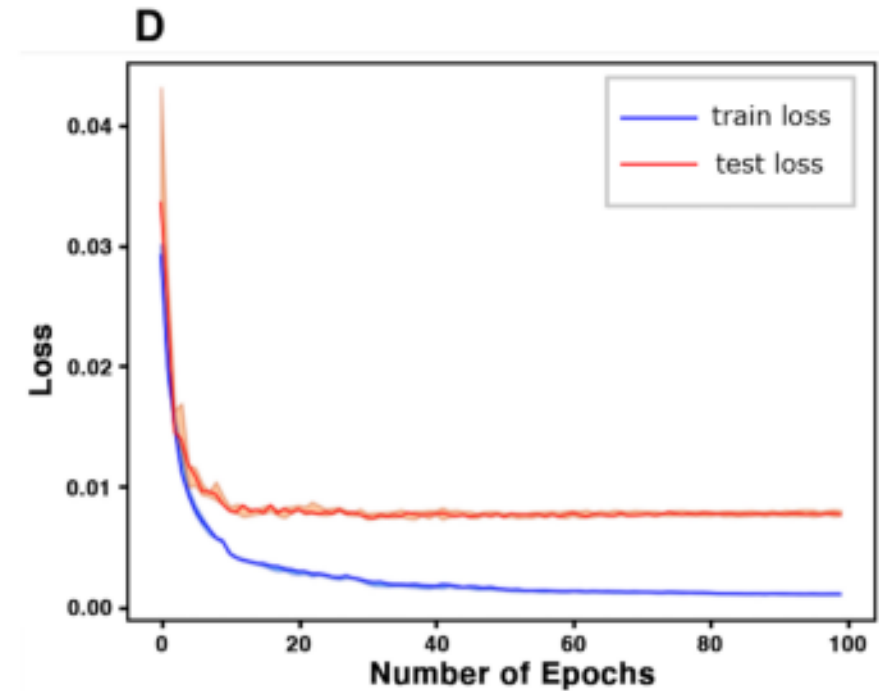
97.30% (n=667)

56.82% (n=88)

98.21% (n=669)

100.00% (n=88)

- 훈련: 음식이 반은 보이고 반은 가려진 상황
- 테스트: 전체적으로 음식이 가려진 상황
- 초기 손실 감소 후 안정적인 성능
- 성공률 100%



“AI가 매우 제한적이거나 불완전한 정보 상황에서도 높은 정확도로 행동을 예측할 수 있는 강력한 능력을 갖추고 있음”

Conclusions: (limitations, caveats, future work)

"We conjecture that perhaps, our ancestor primates also learned to process a form of behavior prediction in a purely visual form, long before they learned to articulate internal mental visions into language, to others or even to themselves"

"인간의 조상인 영장류(원숭이)들이 시간이 지나면서 복잡한 인지 능력, 특히 행동 예측과 같은 고급 뇌 처리 기능을 개발했듯이, 현재의 AI 시스템도 계속해서 발전하여 인간의 인지 능력에 가까워질 수 있을 것이다."

<연구의 가치>

- 로봇이 다른 로봇의 시선으로 세상을 바라볼 수 있다는 사실을 보여준다.
- 관찰자가 상대와 처지를 바꿔서 생각하고 상대의 위치에서 상대가 초록색 점을 볼 수 있는지 없는지 안내없이 이해하는 능력은 기초적인 공감능력을 의미할 수 있다.