

사람과 로봇의 사회적 상호작용을 위한 로봇의 가치효용성 기반 동기-감정 생성 모델

Robot's Motivational Emotion Model with Value Effectiveness for Social Human and Robot Interaction

이 원 형, 박 정 우, 김 우 현, 이 희 승, 정 명 진*
(Won Hyong Lee¹, Jeong Woo Park¹, Woo Hyun Kim¹, Hui Sung Lee², and Myung Jin Chung^{1,*})

¹Korea Advanced Institute of Science and Technology(KAIST)

²Hyundai-KIA MOTORS

Abstract: People would like to be socially engaged not only with humans but also with robots. One of the most common ways in the robotic field to enhance human robot interaction is to use emotion and integrate emotional concepts into robots. Many researchers have been focusing on developing a robot's emotional expressions. However, it is first necessary to establish the psychological background of a robot's emotion generation model in order to implement the whole process of a robot's emotional behavior. Therefore, this article suggests a robot's motivational emotion model with value effectiveness from a Higgins' motivation definition, regulatory focus theory, and Circumplex model. For the test, a game with the best-two-out-of-three rule is introduced. Each step of the game was evaluated by the proposed model. As the results imply, the proposed model generated psychologically appropriate emotions for a robot in the given situation. The empirical survey remains for future work to prove that this research improves social human robot interaction.

Keywords: social interaction, human and robot interaction, emotional robot, motivation theory, regulatory focus

I. 서론

로봇은 여러 목적에 의해 산업, 군사, 의료, 재활, 구조 등의 분야에서 개발, 사용되어왔으며, 최근에는 서비스 로봇 분야의 확대로 상점, 노인복지, 교육, 안내, 놀이 분야에도 로봇의 쓰임이 다양해지고 있다. 서비스 로봇 분야가 발달하면서 로봇은 사람의 일상생활에 자주 등장하게 되었고, 사람과 로봇의 교류 상황이 찾아지게 되었다. 이는 곧 로봇이 자신의 역할을 수행할 때 사람과 교류하게 되는 상황을 고려해야 한다는 것이다. 이에 따라 사람과 로봇의 상호작용(HRI: Human Robot Interaction)의 개념과 연구들이 발전하였고, 물리적 HRI 뿐 아니라 인지/정서적 HRI 기술까지 폭넓은 영역의 학문들이 로봇 연구에 적용되어지고 있다. 이러한 로봇 연구의 흐름은 2004년 일본 후쿠오카에서 선언된 세계 로봇 선언에서 확인할 수 있다. 선언문에 따르면 차세대 로봇은 인류와 공존하는 파트너가 될 것이며 인류를 신체적(Physically), 심리적(psychologically)으로 보조할 것이라고 설명되어 있다.

인지/정서적 HRI 기술에 대한 관심은 다른 여러 로봇 기술의 발달과 노인복지, 심리치료, 유아교육, 개인용 로봇 등의 시장 확대에 힘입어 지속적으로 늘어나고 있다. 특히 로봇이 사람들과 어울리며 나타나게 되는 사회적 현상은 소셜 로봇

(social robot)이라는 주제로 최근 심도 있게 다뤄지고 있다 [1,2]. 이렇게 사람과 로봇이 인지/정서적 교류를 하게 되고 사회적인 상호작용을 하게 되는 이유는 사람은 어떤 상대방이나 사물 등의 대상에게 사회적인 관계를 맺고 싶어 하기 때문이다. 이러한 성향은 로봇을 대상으로도 동일하게 나타난다[3-5]. 따라서 인지/정서적 HRI 기술은 사람과 로봇의 사회적 관계를 향상시키는 중요한 역할을 하게 된다.

사람과 로봇의 사회적 상호작용의 향상을 위해 사람들은 로봇에게 감정(Emotion)의 개념을 부여해왔다. 감정 또는 정서의 개념이 로봇과 융합되며 로봇은 더이상 기계나 계산용 컴퓨터가 아닌 하나의 생명체, 나아가 인격체로 의인화되었다. 로봇과 사람의 미래 모습을 그린 영화 'Bicentennial man (1999)', 'A.I.(2001)', 'iRobot(2004)' 등의 작품에서는 로봇과 사람의 정서적 교류가 사회적인 관계를 향상시킨다는 가능성을 극명하게 보여주었다.

이러한 가능성을 바탕으로 로봇 연구자들은 감정을 이해하고 표현하는 기능을 로봇에 구현해왔다. 얼굴 표정이나 제스처를 통해 감정을 표현하는 로봇으로 Kismet, MDS, WE-4R, KOBEAN, HRP-4C, Geminoid, Doldori, KaMERO, MERO, KIBO, Nao 등 다수의 로봇이 개발되었다. AIBO, Pleo, Paro 등 사람의 모습이 아닌 동물 형태의 감정 표현 로봇들도 개발되었다. 그리고 연구자들은 자연스럽게 자동적인 감정 표현을 위해 얼굴 표정 및 제스처 표현에 대한 다양한 공학적 접근들을 시도하고 있다[6-11]. 또한 로봇이 상황에 따른 적합한 감정을 가질 수 있도록 로봇의 감정 상태를 정의하고 외부의 자극으로부터 감정이 결정되어 가는 과정을 모델링해왔다[12-17].

로봇의 감정 상태를 정의하고 상황입력과 자극을 통해

* Corresponding Author
Manuscript received February 15, 2014 / revised March 15, 2014 / accepted March 30, 2014
이원형, 박정우, 김우현, 정명진: KAIST 전기 및 전자공학과
(leestation@rr.kaist.ac.kr/pjw@rr.kaist.ac.kr/ishsrain@rr.kaist.ac.kr/mjchung@ee.kaist.ac.kr)
이희승: 현대자동차(huisung.lee@kaist.ac.kr)

※ 본 논문은 지식경제부 산업융합원천기술개발사업에 의하여 연구되었음[N02120248].

감정이 결정되는 과정은 심리학, 행동학, 인지과학 등의 여러 학문적 접근을 통해 연구되어지고 있다. 하지만 아직 감정이 생겨나는 기작과 원리는 명확히 통일되지 못하고 있기 때문에 로봇에 적용되는 방법들도 다양하다. 이러한 이유로 로봇의 감정은 개발자의 목적과 적용 환경에 따라 임의로 정의되기 쉽다. 따라서 감정이 생겨나는 근거와 그것이 로봇에 적용되는 이유에 대한 고찰이 필요하다.

여러 심리학자들의 일반적인 주장을 따른다면 감정은 동기(motivation)의 종류와 목표에 대한 성취도의 피드백(feedback)으로 나타나며, 동기의 크기는 감정의 크기와 정서적 경험 강도에 영향을 미치는 선행 요인(precursor)이라 한다[18,19]. 따라서 이 논문에서는 로봇의 감정이 생기는 원리로 심리학의 동기 이론(motivation theory)을 수식화하여 구현하고자 한다. 이를 통해 로봇의 감정이 생겨나는 심리학적 근거를 확보하고, 여러 상황에 적용 가능한 계산 가능한 체계/framework)을 구축하도록 한다. 또한 이 연구는 사람과 로봇의 상호작용이 사람 간의 상호작용 수준이 되는 것을 최종 목표로 하고 있기 때문에 사람의 심리 모델을 바탕으로 로봇의 감정을 정의하였다.

II. 동기 이론과 로봇

1. 동기 이론의 세 분류

동기 이론에는 여러 가지 종류가 있으며 심리학자 Higgins의 주장에 따르면 동기 이론들은 다음의 세 분류로 나뉜다[19]. 첫 번째로 동기는 쾌락(pleasure)을 최대 하거나 고통(pain)을 최소로 하려는 경향으로 이해된다. 두 번째로 동기는 생존을 위한 에너지로 평가되기도 한다. 마지막으로 동기는 무언가를 추구함에 있어 효과적이기 위한(to be effective) 선호도를 가리키기도 한다.

첫 번째 분류는 과거부터 많이 받아들여지던 개념으로 이에 가장 근접한 예시로는 당근과 채찍이 있다. 어떤 일을 함에 있어 쾌락을 최대 하기 위해서는 당근을 얻기 위해 행동을 하게 될 것이고, 고통을 최소 하기 위해서는 채찍을 피하기 위해 행동을 하게 된다는 것이다. 그러나 이러한 접근 방법은 모든 동기를 설명하는 데에는 한계가 있다는 비판이 있다. 쾌락이나 고통은 결과에 따른 감정이지 동기 자체가 되기는 어렵다는 것이다. 예를 들어 목마름의 경우 물을 마시면 쾌감을 느끼고, 물을 마시지 못하면 고통을 느끼지만 물을 먹는 근본적인 동기는 쾌감이나 고통이 아닌 목마름이라는 필요에 의한 것이다[19].

두 번째 분류의 동기 이론 역시 생존이라는 절대적 개념 아래 심리학뿐만 아니라 생명과학, 인문사회학 등 전반적인 분야에서 설명되어왔다. 그러나 이 또한 동기를 설명하는 데에는 한계가 있다는 비판이 있다. 가장 알기 쉬운 반례로는 바로 익스트림 스포츠다. 생존이라는 동기에 의한다면 익스트림 스포츠는 피해야 하는 행위이기 때문이다.

마지막 동기의 분류는 위의 두 분류의 비판에 대한 대안으로 제시되었다. 사람에게 목표함이 있으면 그것을 이루기 위해 동기가 생긴다는 것이다. 이를 뒷받침하는 연구들이 여러 심리학자들에 의해 발표되었다[20-25]. 심리학자 Higgins는 이를 정리하여 동기의 효용성(effectiveness)에 대한 세 가지

방식(way)을 발표했는데, 이는 가치효용성(value effectiveness), 진실효용성(truth effectiveness), 통제효용성(control effectiveness)의 세 방식이다[19]. 가치효용성은 목표로 하는 것에 대한 결과물이 실제적인 가치를 발생시키는 경우 그 목적을 달성하기 위해 발생하는 동기이고, 진실효용성은 밝혀지지 않은 사실에 대한 확인하고자 하는 동기이며, 통제효용성은 자신의 행동이나 주변 상황을 통제할 수 있을 때 발생하는 동기이다. 이렇게 동기를 세 가지 효용성으로 정의하면 첫 번째 분류의 동기는 각 효용성 방식의 성취도에 따른 결과가 쾌락이나 고통으로 나타나는 것으로 설명될 수 있고, 두 번째 분류의 동기는 생존을 가치효용성으로 대치시킴으로써 설명될 수 있다.

2. 로봇에 적용된 동기 이론

로봇 연구자들은 동기 이론들을 복합적으로 로봇에 적용시켜왔다[26]. 첫 번째 분류의 동기들은 다수의 로봇에서 구현되었는데, 사용자가 로봇에게 칭찬이나 혼남에 해당하는 입력을 주면 로봇은 그에 대한 감정적 반응을 하여 칭찬을 더 받기 위하거나 혼남을 피하기 위해 행동하는 형태다 [12,15,28].

두 번째 분류의 동기들은 독립적인 개념으로 구분되어 구현되어 오지는 않았지만, 로봇이 위험 상황에 직면했을 때 이를 벗어나고자 하는 행동 패턴을 수행하도록 해왔다. 또한, 생명체의 생명 유지 특징인 항상성(homeostasis)을 로봇에 도입하여 로봇의 내부상태를 생존에 필요한 범위로 유지하도록 하는 기능들을 구현하기도 했다[12,17,27,28]. 개념적으로도 Issac Asimov의 로봇 3법칙에는 로봇은 자신의 존재를 보호해야 한다는 조항이 들어있다.

마지막 분류의 동기는, 정리된 대로 가치효용성, 진실효용성, 통제효용성으로 나뉘어 살펴보면, 가치효용성의 경우는 쾌락주의원리(hedonic principle)의 개념으로 로봇에 구현되어 있는 경우가 많으며 이는 첫 번째 분류의 동기가 로봇에 구현된 것을 설명한 예시들로 간주될 수 있다. 두 번째 분류의 동기 또한 가치효용성의 하나로 적용되었다고 볼 수 있다. 진실효용성의 경우는 로봇의 학습기능을 돕는 역할로 다양한 형태로 구현되어 있다. 즉, 로봇이 학습되지 않은 부분을 찾아 탐색해 학습해가는 원리를 진실효용성을 증가시키기 위함으로 간주하는 것이다. 이는 휴머노이드 로봇에서 로봇의 호기심이라 표현되기도 하는데, 이러한 호기심을 소비함수(cost function)의 형태로 정의해 내적동기(intrinsic motivation)를 계산하는 체계를 제안한 연구도 있다[29,30]. 하지만 진실효용성이 로봇에 구현된 경우는 로봇의 감정과 연관 지어지지 않았다. 통제효용성은 로봇에 적용된 사례를 찾기 힘들다. 로봇은 자신이 행동하는 모든 과정들을 통제하고 있으며, 자신이 통제할 수 없는 부분들은 인식하지 못하는 범주로 분류되어 통제 자체가 불가능하기 때문이다. 만약 로봇이 어떤 상황에서 행동하는 데 필요한 실제적 소비값을 계산할 수 있다면 이 값을 통제효용성으로 간주하고 행동 결정 요인으로 사용할 수도 있을 것이다.

3. 연구 범위

이 논문은 위의 동기이론 중에서 마지막 분류인 목표를 성취해가는 효용성이 높아지게 하기 위한 선호도로서의 동기 개념을 로봇에 구현하고자 한다. 또한 효용성의 세 방식 중 가치효용성에 대해 구체적인 구현 과정을 제안한다. 이를 통해

효용성을 기반으로 하는 동기 모델이 로봇의 감정 생성에 적용가능하다는 기초적인 근거를 마련하고, 나머지 두 방식인 진실효용성과 통제효용성은 추후 연구로 남겨두었다.

III. 가치효용성기반 동기-감정 관계 수식화 제안

1. 조절초점(Regulatory Focus) 이론

Higgins의 조절초점 이론에 의하면 가치효용성 동기는 두 종류의 조절초점, 즉, 향상초점(promotion focus)과 방어초점(prevention focus)을 갖고 있다. 각 조절초점은 세 개의 요소: 필요충족(needs satisfaction), 기준(standards), 쾌락적 특성(hedonic properties)에 의해 결정된다. 향상초점은 필요충족 중 양육적 필요(nurturance needs), 이상적 기준(ideals), 이득여부(positive outcome, gain or non-loss situation)에 의해 결정되며, 방어초점은 보안적 필요(security needs), 의무적 기준(oughts), 손실여부(negative outcome, loss or non-gain situation)에 의해 결정된다[18,19]. 이 논문에서는 상대적으로 계산이 용이한 이득과 손실여부인 쾌락적 특성을 조절초점의 결정 요소로 선택했다.

조절초점이 결정되면 가치효용성의 목적 달성 여부에 따라 경험되는 감정의 종류는 표 1과 같이 결정된다[18,19].

가치효용성의 달성도는 쾌(pleasure)의 감정과 불쾌(unpleasure or pain)의 감정을 경험하게 한다. 가치효용성이 달성이 되면 기쁨과 안도를 느끼게 되며, 기쁨과 안도는 쾌의 성분을 가지고 있다. 가치효용성이 미달성되면 낙담과 불안을 느끼게 되며, 낙담과 불안은 불쾌의 성분을 가지고 있다. 또한, 기쁨과 안도의 감정, 낙담과 불안의 감정은 향상초점과 방어초점에 따라 구분지어지게 된다. 향상초점에 의해서는 낙담과 기쁨의 감정이 느껴지게 되고, 방어초점에 의해서는 불안과 안도의 감정이 느껴지게 된다.

이 때, 향상초점이면서 가치효용성이 달성되는 경우인 기쁨의 감정은 각성 정도(arousal)가 높게 나타나는 반면, 향상초점이면서 가치효용성이 미달성되는 경우인 낙담의 감정은 각성 정도가 낮게 나타난다. 방어초점이면서 가치효용성이 달성되는 경우인 안도의 감정은 각성 정도가 낮게 나타나는 반면, 방어초점이면서 가치효용성이 미달성되는 경우인 불안의 감정은 각성 정도가 높게 나타난다[18,19]. 따라서 조절초점과 가치효용성에 의해 나타나는 감정을 쾌/불쾌, 각성

정도의 성분으로 나누어 정리하면 표 2와 같이 나타나진다.

2. 조절초점 이론과 원형모델(Circumplex Model)의 결합

감정의 기본 요소를 분석하여 공간화시키는 연구는 다양하게 연구되어왔다[31]. 이 중 가장 널리 사용되는 모델이 쾌/불쾌, 각성의 정도를 축으로 하는 2차원 공간 모델이다. Russell은 이 공간에 28개의 감정 t_i 들이 어떻게 분포하는지를 조사하여 원형모델을 제시했다[32]. t_i 는 감정의 위치를 의미하며 i 는 감정의 종류를 의미한다.

표 2 바탕으로 조절초점 이론에서 이야기하는 감정을 원형모델의 2차원 공간상에 투영하면 그림 1과 같다.

원형모델에서 1사분면에 위치한 감정들이 조절초점이론에서 언급한 기쁨과 관련된 감정 영역이고, 3사분면에 위치한 감정들이 낙담과 관련된 감정 영역이다. 2사분면에 위치한 감정들은 불안과 관련된 감정 영역이며 4사분면에 위치한 감정들은 안도와 관련된 영역이다.

향상초점의 경우 가치효용성이 달성될수록 쾌의 감정이 커지며 각성의 정도가 커지게 되어 1사분면에 위치한 감정들이 나타날 확률이 커지게 되고, 가치효용성이 미달성될수록 불쾌의 감정이 커지며 각성의 정도가 작아지게 되어 3사분면에 위치한 감정들이 나타날 확률이 커지게 된다. 이러한 방향성을 표시한 것이 그림 1에서 실선 화살표다. 방어초점의 경우 가치효용성이 달성될수록 쾌의 감정이 커지며 각성의 정도가 작아지게 되어 4사분면에 위치한 감정들이 나타날 확률이 커지게 되고, 가치효용성이 미달성될수록 불쾌의 감정이 커지며 각성의 정도가 커지게 되어 2사분면에 위치한 감정들이 나타날 확률이 커지게 된다. 이러한 방향성을 표시한 것이 그림 1에서 점선 화살표다.

따라서 그림 1에서의 실선과 점선의 기울기 $\theta[n]$ 와 가치효용성의 달성도로 대변되는 원점으로부터의 거리 $v[n]$ 를 계산하게 되면 2차원 원형모델 공간상의 정서 위치를 결정할 수 있게 된다.

3. 동기-감정 생성 모델

조절초점 이론은 각 심리적 요소들에 대한 연관관계를 기술하고 있지만, 그 관계를 로봇에 적용하기 위해서는 수식화의 과정이 필요하다. 조절초점 이론을 수식화하고 원형모델과

표 1. 가치효용성과 조절초점에 따라 경험되는 감정 종류.

Table 1. Correlation of emotions with value effectiveness and regulatory focus.

조절초점과 감정의 종류	가치효용성	
	미달성	달성
향상초점	낙담 Dejection	기쁨 Cheerfulness
방어초점	불안 Agitation	안도 Quiescence

표 2. 가치효용성과 조절초점에 따른 쾌/불쾌, 각성 정도.

Table 2. Pleasure/unpleasure and arousal level with value effectiveness and regulatory focus.

조절초점과 감정 성분	가치효용성	
	미달성	달성
향상초점	불쾌, 약한 각성	쾌, 강한 각성
방어초점	불쾌, 강한 각성	쾌, 약한 각성

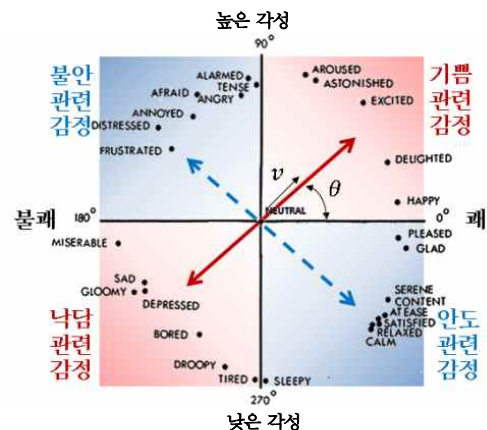


그림 1. 원형모델에 투영시킨 조절초점 이론의 감정 종류.

Fig. 1. Emotional category reflected in regulatory focus theory on Circumplex model.

결합하여 계산 가능한 모델을 제안하는 것이 이 논문의 주 목적이다.

쾌락 특성 요소 계산: 앞선 설명에 따라 쾌락 특성 요소를 계산하기 위해서는 이득과 손실의 총 기대값을 알아야 한다. 어떤 상황에 대한 결과값을 o_k , 결과값에 대한 각 단계의 확률값을 $P_k[n]$ 라 가정하면 예상되는 기대값은 결과값과 확률값의 곱인 $o_k \cdot P_k[n]$ 으로 계산된다. k 는 결과의 종류를 의미하며 주어진 상황에 대한 모든 결과의 총 기대값 $E[n]$ 은 수식 (1)과 같이 계산된다.

$$E[n] = \sum_k (o_k \cdot P_k[n]) \quad (1)$$

조절초점 결정: 쾌락 특성 요소가 계산되면 조절초점 이론에 따라 다음의 조건으로 조절초점이 결정된다.

$$\text{조절초점} = \begin{cases} \text{향상초점} & \text{if } E[n] \geq 0 \\ \text{방어초점} & \text{if } E[n] < 0 \end{cases} \quad (2)$$

가치효용성 달성도 계산: 가치효용성 달성도 $v[n]$ 는 목표에 얼마나 가까워졌는지에 대한 척도로 이 논문에서는 각 단계 마다의 목표 달성 가능성 $P_{goal}[n]$ 의 변화량을 계산하는 방법을 제안한다.

$$v[n] = (P_{goal}[n] - P_{goal}[n-1]) / Q[n] \quad (3)$$

$Q[n]$ 은 정규화 값으로 이전 단계의 목표 달성 가능성 $P_{goal}[n-1]$ 과 가치효용성 달성도 $v[n]$ 의 최대, 최소값인 1 또는 0의 거리 값으로 정의된다. 현재 단계의 목표 달성 가능성 $P_{goal}[n]$ 이 이전 단계의 목표 달성 가능성 $P_{goal}[n-1]$ 보다 커진 경우는 $v[n]$ 의 최대값인 1을 참고하고, 반대의 경우는 $v[n]$ 의 최소값인 0을 참고한다. 이를 수식화하면 수식 (4)와 같고, 이 개념을 도식화하면 그림 2와 같다.

$$Q[n] = \begin{cases} 1 - P_{goal}[n-1] & \text{if } P_{goal}[n] \geq P_{goal}[n-1] \\ P_{goal}[n-1] & \text{if } P_{goal}[n] < P_{goal}[n-1] \end{cases} \quad (4)$$

가치효용성 기울기 계산: 가치효용성 기울기 $\theta[n]$ 는 그림 1에 의하면 평균적으로 45도 또는 -45도의 값을 가져야 한다. 그러나 상황에 따라 이 값을 다양하게 조절할 수 있게 하기 위해, 이 논문에서는 조절초점을 결정하기 위해 사용된 기대값 $E[n]$ 에 비례하도록 $\theta[n]$ 을 수식 (5)와 같이 계산하였다.

$$\theta[n] = c \cdot E[n] \quad (5)$$

수식 (5)에서 c 는 기대값 $E[n]$ 가 $[-90^\circ, 90^\circ]$ 의 범위를

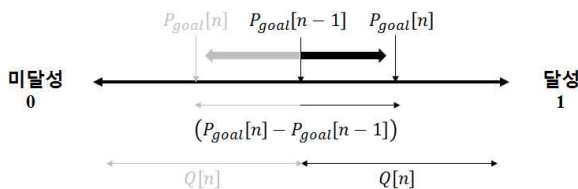


그림 2. 가치효용성 달성도 계산 개념도.
Fig. 2. Conceptual diagram of value effectiveness evaluation.

가지도록 하는 표준화 상수다. 이 수식에 의하면 $\theta[n]$ 가 $[-90^\circ, 0^\circ]$ 의 범위에 있을 때는 조절초점이 방어초점으로 결정되고, $\theta[n]$ 가 $[0^\circ, 90^\circ]$ 의 범위에 있을 때는 조절초점이 향상초점으로 결정되는 것을 알 수 있다.

정서 벡터(Affect Vector) 계산: 얻어진 가치효용성 달성도 $v[n]$ 와 기울기 $\theta[n]$ 를 바탕으로 2차원 원형모델 공간상의 정서 위치 $\mathbf{a}[n]$ 을 계산한다. 이를 직교좌표 정서 벡터 $\mathbf{a}[n] = (a_1[n], a_2[n])$ 로 표현하면 다음과 같이 계산된다.

$$a_1[n] = v[n] \cdot \cos(\theta[n-1]) \quad (6)$$

$$a_2[n] = v[n] \cdot \sin(\theta[n-1]) \quad (7)$$

수식 (6)과 수식 (7)에서 $a_1[n]$ 과 $a_2[n]$ 는 각각 원형모델에서 쾌/불쾌 값과 각성 정도를 의미한다. 이 때, 조절초점은 현재 단계의 사건이 일어나기 전에 결정되어 있는 것이기 때문에 가치효용성 기울기는 이전 단계의 값인 $\theta[n-1]$ 이 사용됨을 확인할 수 있다.

최종 감정 계산: 최종 감정 $\mathbf{e}[n] = (e_1[n], e_2[n], \dots)$ 은 앞에서 계산한 정서 벡터 $\mathbf{a}[n]$ 와 \mathbf{t}_i 의 관계를 통해 계산된다.

$\mathbf{t}_i = (t_{i,1}, t_{i,2})$ 는 2차원 원형모델 공간상에 지정된 감정의 위치 벡터이며, i 는 28가지 감정종류를 가리킨다[32]. 필요에 의해 2차원 원형모델 공간의 원점을 중립(Neutral) 감정으로 정의하여 감정 종류에 포함시킬 수도 있다.

최종 감정 $\mathbf{e}[n]$ 은 얻고자 하는 감정의 종류와 적용 양식에 따라 여러 방법으로 계산될 수 있는데 이 논문에서는 다음의 두가지 방법을 제안한다.

최종 감정 계산 방법 1: 첫번째 방법은 정서 벡터 $\mathbf{a}[n]$ 와 원형모델의 감정 위치 \mathbf{t}_i 사이의 거리를 계산하고, 거리값의 역수에 비례하도록 수식 (8)과 같이 각 감정의 크기 $e_i[n]$ 를 정의하는 것이다.

$$e_i[n] = \frac{1 / \|\mathbf{t}_i - \mathbf{a}[n]\|}{\sum 1 / \|\mathbf{t}_i - \mathbf{a}[n]\|} \quad (8)$$

이렇게 계산하면 정서 벡터 $\mathbf{a}[n]$ 에 가까운 원형모델의 감정 일수록 해당하는 최종 감정의 크기는 상대적으로 큰 값을 갖게 되며 다른 감정과의 구분도 명확해지게 된다. 이 값들 중 최대값만을 최종 감정값으로 사용할 수도 있고, 각 값들을 원소로 하는 벡터 $\mathbf{e}[n]$ 을 최종 감정으로 사용할 수도 있다. 이 방법은 저자의 이전 연구에서 사용된 바가 있다[17].

이 방법의 한계는 다음과 같다. 정서 벡터 $\mathbf{a}[n]$ 가 근접한 원형모델의 두 감정 $\mathbf{t}_1, \mathbf{t}_2$ 사이에 위치하게 되었다고 가정해보자. 그러면, $\mathbf{a}[n]$ 과 \mathbf{t}_i 의 절대적인 거리는 매우 가까워 최종 감정 $e_i[n]$ 은 최대 값에 가까운 값으로 계산되었어야 한다. 하지만, 상대적 거리에 의해 최종 감정 $e_1[n], e_2[n]$ 의 크기는 반으로 줄어들게 된다. 또한, 이때의 $\mathbf{a}[n]$ 의 위치가 조금이라도 변하게 되면 $e_1[n], e_2[n]$ 의 크기는 극적으로 변하게 되고, 비현실적인 감정 변화가 생겨나게 된다. 이에 대한 보완으로 두 번째 방법이 고안되었다.

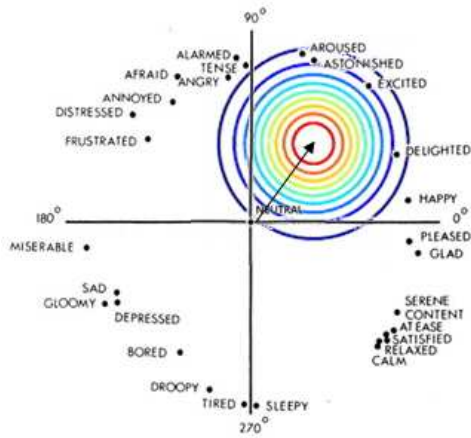


그림 3. 원형모델에 투영한 정서 벡터 분포 등고선.
Fig. 3. Contour line of an affect vector distribution on Circumplex model.

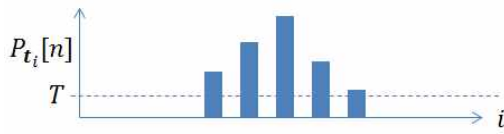


그림 4. 최종 감정 분포 그래프 예시.
Fig. 4. An Example of final emotion distribution graph.

최종 감정 계산 방법 2: 두 번째 방법은 정서 벡터 $\mathbf{a}[n]$ 를 중심으로 정규분포 형태의 정서 벡터 분포를 두어 t_i 위치에서의 분포값 $P_i[n]$ 을 계산하고 이에 비례하게 각 감정의 크기 $e_i[n]$ 를 계산하는 것이다. $P_i[n]$ 은 수식 (9)와 같이 계산한다.

$$P_i[n] = \begin{cases} \exp\left(-\left(\frac{\|t_i - \mathbf{a}[n]\|}{2\sigma^2}\right)^2\right) & \text{if } P_i[n] > T \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

여기서 T 는 역치값이다. 이 개념을 2차원 원형모델 공간상에 표현하면 그림 3의 등고선과 같이 되고, 그림 4와 같은 분포 그래프를 가진다.

얻어진 정서 벡터 분포 $P_i[n]$ 의 값으로 각 감정의 크기 $e_i[n]$ 를 얻는다.

$$e_i[n] = P_i[n] \quad (10)$$

이 방법은 정서 벡터 $\mathbf{a}[n]$ 가 원형모델의 감정 위치 t_i 에 가까울수록 최종 감정 $e_i[n]$ 의 크기가 커지는 특징은 유지하면서도 앞선 최종 감정 계산 방법 1에서 가지는 한계를 극복할 수 있다.

이 방법에서도 마찬가지로 최대 값을 갖는 $e_i[n]$ 하나만을 최종 감정으로 선택할 수도 있고, $e_i[n]$ 를 원소로 하는 벡터 $\mathbf{e}[n]$ 를 최종 감정 벡터로 사용할 수도 있다. 여기에 한 가지를 더 제한하면 계산된 정서 벡터 분포 $P_i[n]$ 을 표준화하여 확률값처럼 사용하는 것이다. 그러면 같은 정서 벡터 $\mathbf{a}[n]$ 에

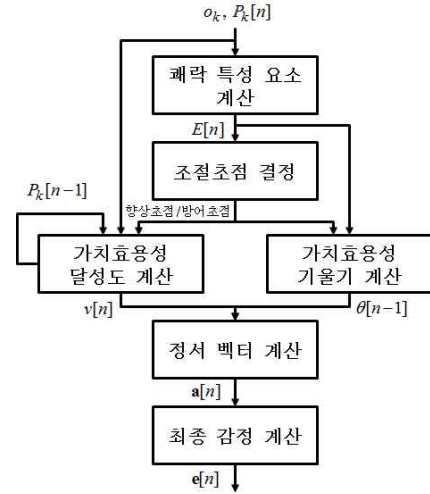


그림 5. 가치효용성 기반 동기-감정 생성 모델 전체 구조도.
Fig. 5. Overall structure of motivation emotion model with value effectiveness.

대해서도 $P_i[n]$ 에 비례하게 다양한 확률로 최종 감정들이 선택될 수 있다. 즉, 정서 벡터 $\mathbf{a}[n]$ 에 가까운 원형모델의 감정일수록 더 높은 확률로 최종 감정이 되는 것이다.

4. 전체 구조도

제안된 가치효용성 기반 동기-감정 생성 모델의 전체 구조도는 그림 5와 같다.

IV. 실험 결과

1. 시나리오

제안된 모델의 타당성을 확인하기 위해 삼세판 게임 시나리오를 설정한다. 삼세판 게임은 세 번의 단계 중 최소 2번의 단계를 이기거나 지면 결판이 나는 게임이다. 결판이 나면 최종 승인 경우 상, 최종 패인 경우 벌을 받는다. 각 단계는 승 또는 패의 경우만 있으며 무승부는 없다고 가정한다.

이러한 시나리오를 바탕으로 앞서 설명한 수식들의 변수들을 연결시키면 다음과 같다.

- k : 상 또는 벌
- o_k : 상 또는 벌의 크기
- $P_k[n]$: 상 또는 벌을 받게 될 확률
- $E[n]$: 상/벌을 포함한 총 기대 값
- $P_{goal}[n]: P_{prize}[n]$
- 기준: 이상적, 의무적 기준 공존

n 번째 단계까지 m 번의 승리와 l 번의 패를 한 경우 상을 받게 될 확률 $P_{prize}[n]$ 는 다음과 같이 계산된다.

$$P_{prize}[n] = \frac{{}^{3-n}C_{2-m}}{{}^{3-n}C_{2-m} + {}^{3-n}C_{2-l}} \quad (11)$$

삼세판 게임은 간단하므로 결판나는 경우를모두 나열하면 [승, 승], [승, 패, 승], [패, 승, 승], [패, 패], [패, 승, 패], [승, 패, 패] 이렇게 여섯 가지가 존재한다. 각 경우에 대해 단계마다 가치효용성 기반 동기-감정 생성 모델이 계산되어 로봇의 최종 감정을 얻는 실험을 수행했다.

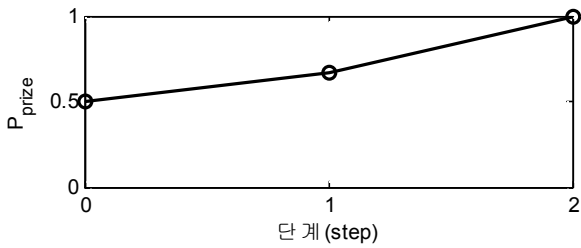


그림 6. [승, 승]의 경우, 상을 받게 될 확률 $P_{prize}[n]$ 변화.
Fig. 6. [Win, Win] change of $P_{prize}[n]$.

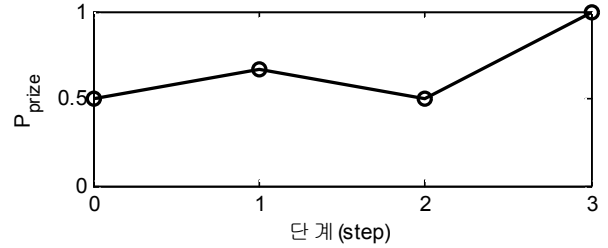


그림 9. [승, 패, 승]의 경우, 상을 받게 될 확률 $P_{prize}[n]$.
Fig. 9. [Win, Loss, Win] change of $P_{prize}[n]$.

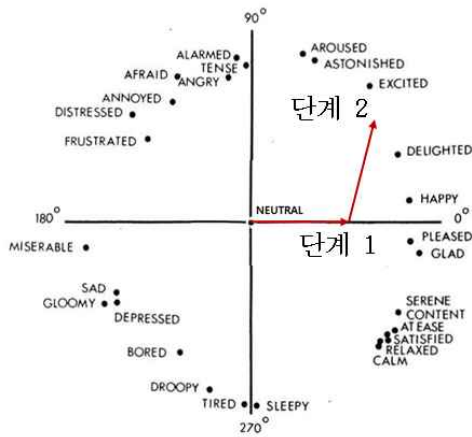


그림 7. [승, 승]의 경우, 정서 벡터 $a[n]$ 의 변화.
Fig. 7. [Win, Win] change of affect vector $a[n]$.

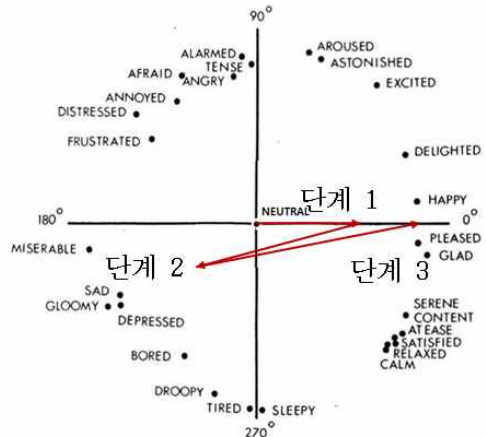


그림 10. [승, 패, 승]의 경우, 정서 벡터 $a[n]$ 의 변화.
Fig. 10. [Win, Loss, Win] change of affect vector $a[n]$.

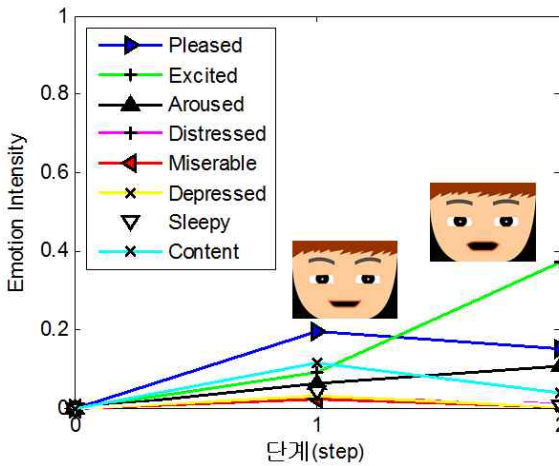


그림 8. [승, 승]의 경우, 각 단계의 감정 변화와 표정.
Fig. 8. [Win, Win] change of emotion and facial expressions.

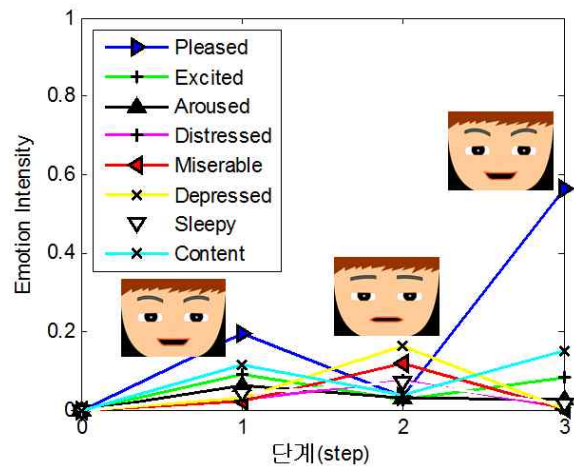


그림 11. [승, 패, 승]의 경우, 각 단계의 감정 변화와 표정.
Fig. 11. [Win, Loss, Win] change of emotion and facial expressions.

2. 설정값

위 시나리오에 필요한 설정값은 다음과 같다. 각 단계의 승률은 0.5, 상은 1, 벌은 -1로 둔다. 최종 감정 결정 방법은 두 번째 방법을 사용하고, $P_i[n]$ 계산에 필요한 σ 값은 0.3으로 역치값 T 는 0.1로 설정한다. 또한, 가장 큰 $e_i[n]$ 의 감정 종류와 크기를 최종 감정의 대표값으로 표시한다.

3. 얼굴 표정 표현

로봇의 감정 결과를 숫자나 공간상의 벡터 등으로 이해하는 것은 쉽지 않다. 따라서 얻어진 최종 감정의 종류와 값을

얼굴 표정으로 표현하여 직관적으로 감정 결과를 이해할 수 있도록 도왔다. 얻어진 최종 감정 $e[n]$ 중 가장 큰 감정을 표정 표현 알고리즘 Linear Affect-Expression Space Model이 적용된 얼굴 시뮬레이터 FRESi에 표현했다[34,35]. 저자의 이전 연구 방법을 참고했다[17].

4. 결과

이 시나리오에 의해 일어나는 총 여섯 가지 경우에 따른 수식의 각 값들의 변화와 2차원 원형모델 공간상에서 정서 벡터의 이동, 감정 변화를 그래프로 정리하였다. 감정 변화

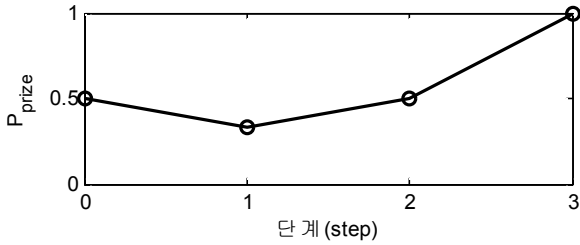


그림 12. [패, 승, 승]의 경우, $P_{prize}[n]$ 의 변화.

Fig. 12. [Loss, Win, Win] change of $P_{prize}[n]$.

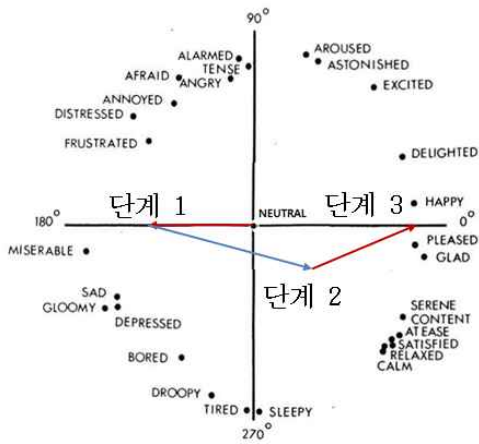


그림 13. [패, 승, 승]의 경우, 정서 벡터 $a[n]$ 의 변화.

Fig. 13. [Loss, Win, Win] change of affect vector $a[n]$.

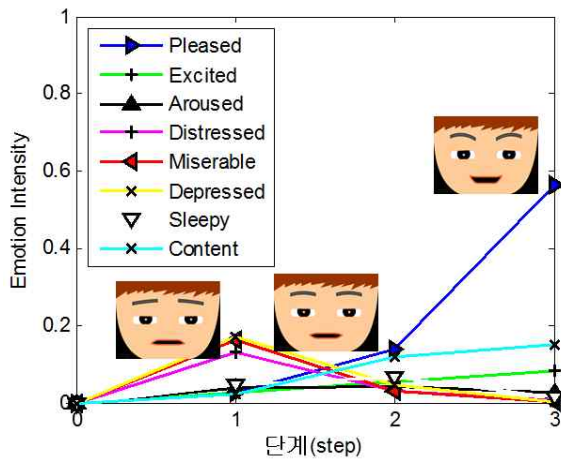


그림 14. [패, 승, 승]의 경우, 각 단계의 감정 변화와 표정.

Fig. 14. [Loss, Win, Win] change of emotion and facial expressions.

그래프의 경우 원형모델의 28가지 감정을 모두 표현하는 것에 어려움이 있기 때문에 각 축을 대표하는 감정(Pleased, Aroused, Miserable, Sleepy)과 각 사분면을 대표하는 감정(Excited, Distressed, Depressed, Content)만을 표시하였다.

경우 1 [승, 승]: [승, 승]의 경우 첫 단계에서 승리하였기 때문에 중간 정도의 행복(Pleased)을 보였고, 최종 승리의 가능성이 높아졌기 때문에 항상초점으로 설정됐다. 따라서 최종 승리를 했을 때에는 행복감에 각성 정도가 더해져 신이 난(Excited) 감정이 나타나는 것을 확인할 수 있다.

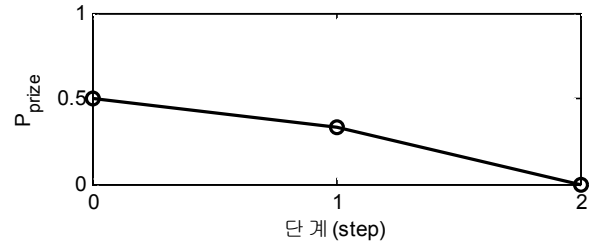


그림 15. [패, 패]의 경우, $P_{prize}[n]$ 의 변화.

Fig. 15. [Loss, Loss] change of $P_{prize}[n]$.

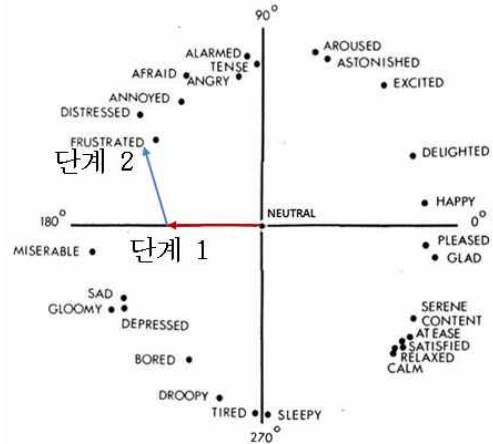


그림 16. [패, 패]의 경우, 정서 벡터 $a[n]$ 의 변화.

Fig. 16. [Loss, Loss] change of affect vector $a[n]$.

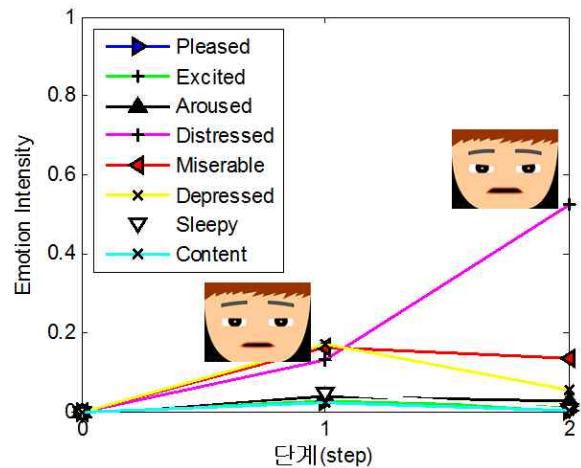


그림 17. [패, 패]의 경우, 각 단계의 감정 변화와 표정.

Fig. 17. [Loss, Loss] change of emotion and facial expressions.

경우 2 [승, 패, 승]: [승, 패, 승]의 경우 첫 단계에서는 승리하였기 때문에 중간 정도의 행복을 보였고, 항상초점으로 설정됐다. 하지만 그 다음 단계는 기대와는 달리 패하였기 때문에 좌절감(Depressed)과 비참함(Miserable)이 다소 나타난다. 최종 패를 한 것이 아니므로 좌절감이 크게 나타나지는 않았다. 마지막 단계에서는 결국 다시 최종 승리를 얻어내어 강한 행복감이 나타났다.

경우 3 [패, 승, 승]: [패, 승, 승]의 경우 첫 단계에서 패하였기 때문에 중간 정도의 좌절감, 비참함, 중압감을 보였고,

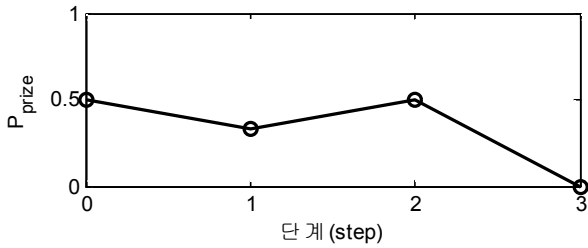


그림 18. [패, 승, 패]의 경우, $P_{prize}[n]$ 의 변화.
Fig. 18. [Loss, Win, Loss] change of $P_{prize}[n]$.

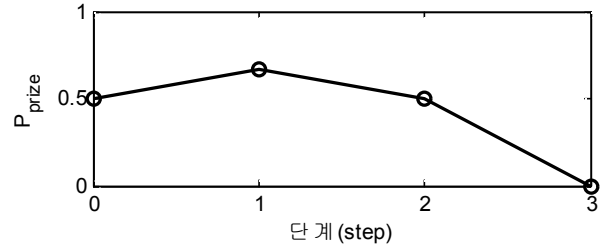


그림 21. [승, 패, 패]의 경우, $P_{prize}[n]$ 의 변화.
Fig. 21. [Win, Loss, Loss] change of $P_{prize}[n]$.

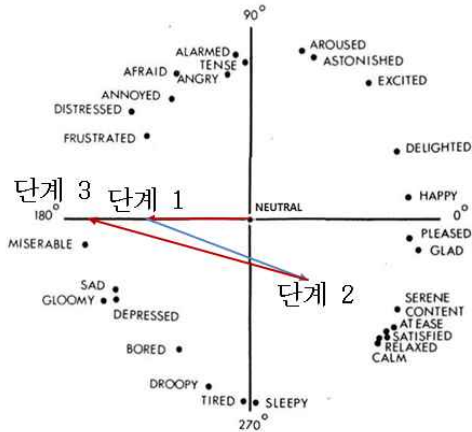


그림 19. [패, 승, 패]의 경우, 정서 벡터 $a[n]$ 의 변화.
Fig. 19. [Loss, Win, Loss] change of affect vector $a[n]$.

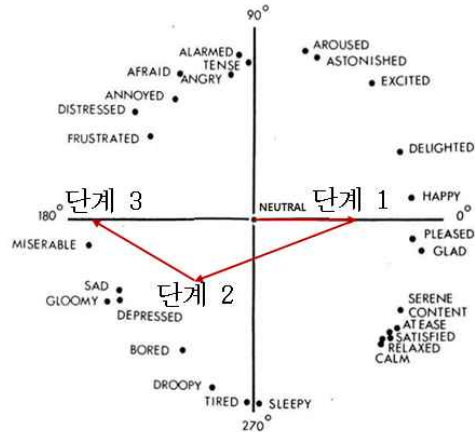


그림 22. [승, 패, 패]의 경우, 정서 벡터 $a[n]$ 의 변화.
Fig. 22. [Win, Loss, Loss] change of affect vector $a[n]$.

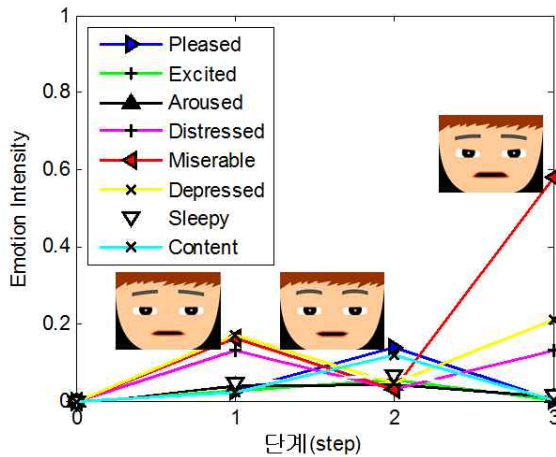


그림 20. [패, 승, 패]의 경우, 각 단계의 감정 변화와 표정.
Fig. 20. [Loss, Win, Loss] change of emotion and facial expressions.

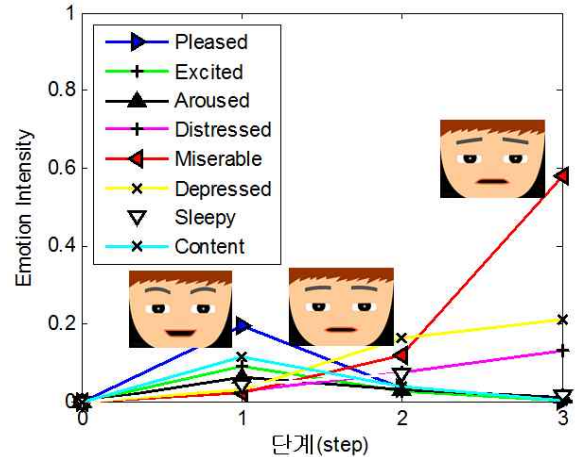


그림 23. [승, 패, 패]의 경우, 각 단계의 감정 변화와 표정.
Fig. 23. [Win, Loss, Loss] change of emotion and facial expressions.

최종 패의 가능성이 커지게 되어 방어초점으로 설정됐다. 하지만 두 번째 단계에서 승리하여 다소의 안도감(Relaxed)을 나타냈다. 마지막 단계에서 최종 승리를 얻어내어 강한 행복감이 나타나는 것을 확인할 수 있었다.

경우 4 [패, 패]: [패, 패]의 경우 첫 단계에서 패하였기 때문에 중간 정도의 좌절감, 비참함, 중압감을 보였고, 방어 초점으로 설정되었다. 두 번째 단계에서도 패하게 되어 최종 패를 얻게 되었고, 강하게 괴로워하는 감정을 나타냈다.

경우 5 [패, 승, 패]: [패, 승, 패]의 경우 첫 단계에서 패하

였기 때문에 중간 정도의 좌절감, 비참함, 중압감을 보였고, 방어 초점으로 설정됐다. 하지만 두 번째 단계에서 승리하여 약한 안도감과 행복을 나타냈다. 마지막 단계에서 패하며 최종 패를 얻게 되어 강한 비참함이 나타나는 것을 확인할 수 있었다.

경우 6 [승, 패, 패]: [승, 패, 패]의 경우 첫 단계에서는 승리하였기 때문에 중간 정도의 행복을 보였고, 항상 초점으로 설정됐다. 하지만 그 다음 단계는 기대와는 달리 패하였기 때문에 좌절감과 비참함이 다소 나타난다. 마지막 단계에서는 결국 최종 패를 얻게 되어 강한 비참함이 나타났다.

실험을 통해 얻은 결과를 통해 로봇의 감정이 사람의 심리와 유사하게 형성되는 것을 확인할 수 있었다. 이처럼 로봇이 사람의 동기와 감정을 모사하여 표현한다면 사람과 로봇의 사회적 상호작용의 향상을 기대할 수 있을 것이다. 다양한 피실험자를 대상으로 사람과 로봇의 사회적 상호작용이 향상되었음을 밝히는 설문조사를 준비 중이다.

V. 결론 및 추후과제

이 연구를 통하여 사람과 로봇의 사회적 교감을 위하여 로봇에게 감정의 개념이 필요하다는 것을 주장하였고, 심리학적 근거를 마련하기 위해 동기 이론을 바탕으로 감정의 생성 원리를 도입하였다. 사용된 동기 이론은 가치효용성 기반 조절초점 이론이었고, 이를 원형모델과 결합시켜 동기-감정 모델을 수식화하여 제안하였다. 제안된 모델을 검증하기 위해 삼세판 게임을 설정하여 각 단계마다 정서 벡터와 최종 감정을 계산하였고, 이를 얼굴 표정 시뮬레이터에 표시시켜 보았다. 얻어진 결과를 바탕으로 로봇의 감정 생성 과정이 사람의 심리 과정과 유사함을 확인할 수 있었다. 또한, 효용성을 바탕으로하는 감정 생성 알고리즘의 기초적인 근거를 마련할 수 있을 것이라 기대된다.

추후 다수의 피실험자를 대상으로 사회적 상호작용이 향상되었음을 확인해야 한다. 더불어 제안된 모델이 범용적으로 사용될 수 있는 체계를 확인하기 위해 삼세판 게임 시나리오 이외의 실험환경도 준비하고자 한다. 마지막으로 이미 언급한 것과 같이 가치효용성 이외의 진실효용성, 통제효용성에 대한 동기 이론을 로봇에 맞게 재정의하여 확장된 로봇의 동기-감정 모델을 제안하고자 한다.

REFERENCES

[1] Lin, Patrick, Keith Abney, and George A. Bekey. *Robot Ethics: The Ethical and Social Implications of Robotics*. The MIT Press, 2011.

[2] Guy Hoffman (2013, October). Robots with "soul" [Video file]. Retrieved from http://www.ted.com/talks/guy_hoffman_robots_with_soul.html

[3] Breazeal, Cynthia L. *Designing Sociable Robots with CDROM*, MIT press, 2004.

[4] Thomaz, Andrea L., and Cynthia Breazeal, "Teachable robots: Understanding human teaching behavior to build more effective robot learners," *Artificial Intelligence*, 172.6, pp. 716-737, 2008.

[5] Kozima, Hideki, Marek P. Michalowski, and C. Nakagawa, "Keepon: A Playful Robot for Research, Therapy, and Entertainment," *International Journal of Social Robotics* 1.1, pp. 3-18, 2009.

[6] Lee, Hui Sung, et al. "A linear dynamic affect-expression model: Facial expressions according to perceived emotions in mascot-type facial robots," *Robot and Human interactive Communication, 2007. RO-MAN 2007. The 16th IEEE International Symposium on*. IEEE, 2007.

[7] Park, Jeong Woo, et al. "A Robot Simulator'FRESi'for Dynamic Facial Expression," *The 6th International Conference on Ubiquitous Robots and Ambient Intelligence*, 2009.

[8] Kim, Woo Hyun, et al. "Synchronized multimodal expression generation using editing toolkit for a human-friendly robot,"

Robotics and Biomimetics (ROBIO), 2009 IEEE International Conference on. IEEE, 2009.

[9] Kim, Woo Hyun, et al. "Hierarchical database based on feature parameters for various multimodal expression generation of robot," *Advanced Robotics and its Social Impacts (ARSO), 2010 IEEE Workshop on*. IEEE, 2010.

[10] Kim, Jaewoo, et al. "Automated robot speech gesture generation system based on dialog sentence punctuation mark extraction," *System Integration (SII), 2012 IEEE/SICE International Symposium on*. IEEE, 2012.

[11] Haring, M., Nikolaus Bee, and Elisabeth André, "Creation and evaluation of emotion expression with body movement, sound and eye color for humanoid robots," *RO-MAN*, 2011.

[12] Arkin, Ronald C., et al. "An ethological and emotional basis for human-robot interaction," *Robotics and Autonomous Systems* 42.3, pp. 191-201, 2003.

[13] Breazeal, Cynthia, "Function meets style: insights from emotion theory applied to HRI," *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on* 34.2, pp. 187-194, 2004.

[14] Miwa, Hiroyasu, et al. "A new mental model for humanoid robots for human friendly communication introduction of learning system, mood vector and second order equations of emotion," *Robotics and Automation, 2003. Proceedings. ICRA'03. IEEE International Conference on*. vol. 3. IEEE, 2003.

[15] Kwon, Dong-Soo, et al. "Emotion interaction system for a service robot," *Robot and Human interactive Communication, 2007. RO-MAN 2007. The 16th IEEE International Symposium on*. IEEE, 2007.

[16] Kim, Won Hwa, et al. "Stochastic approach on a simplified OCC model for uncertainty and believability," *Computational Intelligence in Robotics and Automation (CIRA), 2009 IEEE International Symposium on*. IEEE, 2009.

[17] Lee, Won Hyong, et al. "Robot's emotion generation model for transition and diversity using energy, entropy, and homeostasis concepts," *Robotics and Biomimetics (ROBIO), 2010 IEEE International Conference on*. IEEE, 2010.

[18] Higgins, E. Tory. "Promotion and prevention experiences: Relating emotions to nonemotional motivational states," *Handbook of Affect and Social Cognition*, pp. 186-211, 2001.

[19] Higgins, E. Tory, *Beyond Pleasure and Pain: How Motivation Works*. Oxford University Press, 2011.

[20] Keynes, John Maynard, *General Theory of Employment, Interest and Money*. Atlantic Publishers & Dist, 2006.

[21] R. S. Woodworth and H. Schlosberg, "Experimental psychology. New York: Henry Holt and Company," (1938): 696.

[22] Hebb, Donald Olding. "Drives and the CNS (conceptual nervous system)," *Psychological Review* 62.4 (1955): 243.

[23] White, Robert W. "Motivation reconsidered: the concept of competence," *Psychological Review* 66.5 (1959): 297.

[24] Bandura, Albert, "Self-efficacy mechanism in human agency," *American Psychologist* 37.2 (1982): 122.

[25] Deci, Edward L., and Richard M. Ryan. "The "what" and "why" of goal pursuits: Human needs and the self-determination of behavior," *Psychological Inquiry* 11.4, pp. 227-268, 2000.

[26] Hawes, Nick. "A survey of motivation frameworks for intelligent systems," *Artificial Intelligence* 175.5, pp. 1020-1036, 2011.

- [27] Breazeal, Cynthia, "A motivational system for regulating human-robot interaction," *AAAI/IAAI*, 1998.
- [28] Dimas, Joana, et al. "Pervasive pleo: long-term attachment with artificial pets," *Mobile HCI*, 2010.
- [29] Oudeyer, Pierre-Yves, and Frederic Kaplan, "What is intrinsic motivation? a typology of computational approaches," *Frontiers in Neurobotics 1*, 2007.
- [30] Oudeyer, P-Y., Frédéric Kaplan, and Verena Vanessa Hafner, "Intrinsic motivation systems for autonomous mental development," *Evolutionary Computation, IEEE Transactions on 11.2*, pp. 265-286, 2007.
- [31] Sokolov, Evgeni N., and Wolfram Boucsein. "A psychophysiological model of emotion space," *Integrative Physiological and Behavioral Science 35.2*, pp. 81-119, 2000.
- [32] Russell, James A., "A circumplex model of affect," *Journal of Personality and Social Psychology*, vol. 39, no. 6, Dec. 1980.
- [33] Ekman, Paul, and Wallace V. Friesen. *Unmasking the Face: A Guide to Recognizing Emotions from Facial Clues*. Ishk, 2003.
- [34] Lee, Hui Sung, Jeong Woo Park, and Myung Jin Chung, "A linear affect-expression space model and control points for mascot-type facial robots," *Robotics, IEEE Transactions on 23.5*, pp. 863-873, 2007.
- [35] J. W. Park, W. H. Kim, W. H. Lee, and M. J. Chung, "A robot simulator 'FRESi' for dynamic facial expression," *The 6th International Conference on Ubiquitous Robots and Ambient Intelligence*, pp. 727 - 728, Gwangju, Korea, Oct. 2009.



이 원 형

2008년 KAIST 전기및전자공학과 학사.
2010년 KAIST 전기및전자공학과 석사.
2010년~현재 KAIST 전기및전자공학과
박사과정 재학중. 관심분야는 HRI,
Social Robotics, Robot's Emotion and
Motivation, and Its Applications.



박 정 우

2005년 경북대학교 전자전기공학부 학사.
2007년 KAIST 전기및전자공학과 석사.
2007년~현재 KAIST 전기및전자공학과
박사과정 재학중. 관심분야는 HRI,
Facial Robot, Machine Learning.



김 우 현

2007년 KAIST 전기및전자공학과 학사.
2009년 KAIST 전기및전자공학과 석사.
2009년~현재 KAIST 전기및전자공학과
박사과정 재학중. 관심분야는 HRI,
artificial emotion, robot expression, and
human gesture recognition.



이 희 승

2000년 KAIST 전기및전자공학 공학 학사.
2002년 KAIST 전기및전자공학 공학 석사.
2008년 KAIST 전기및전자공학 공학박사.
2008년~2010년 삼성전자 책임연구원.
2010년~현재 현대자동차 책임연구원.
관심분야는 HRI/HCI, 임베디드 시스템,
센서융합.



정 명 진

1973년 서울대학교 공과대학 전기공학과 학사.
1977년 미시간대학교 전기공학과 석사.
1983년 미시간대학교 제어공학과 박사.
1983년~현재 KAIST 전기및전자공학과 교수.
관심분야는 sensor-based robot control and
planning, HRI, and service robots for the disabled.