

ხელოვნური მორალური აგენტი

სამაგისტრო ნაშრომი

ავტორი: ლაშა უსუფაშვილი

ხელმძღვანელი: თამარ ცხადაძე

მეცნიერების და ხელოვნების ფაკულტეტი

ილიას სახელმწიფო უნივერსიტეტი

26.06.2019

Abstract	1
აბსტრაქტი	2
შესავალი	3
<i>ჰიპოთეზა და მიზანი</i>	<i>5</i>
თეორიის მიმოხილვა	6
<i>სხვების საკითხი ეთიკაში</i>	<i>6</i>
<i>სუბიექტის დეცენტრალიზირება</i>	<i>8</i>
<i>გაფართოებული აგენტობის თეორიები</i>	<i>9</i>
<i>ემანუელ ლევენასი</i>	<i>10</i>
<i>ანთროპოცენტრისტული ეთიკა</i>	<i>12</i>
<i>სათნოების ეთიკა</i>	<i>13</i>
<i>თავისუფალი ნება, მორალი და ხელოვნური აგენტი</i>	<i>15</i>
<i>ცნობიერი ხელოვნური აგენტი</i>	<i>18</i>
<i>არაცნობიერი აგენტი</i>	<i>20</i>
ადამიანის მორალის მოდელირება	22
<i>ხელოვნური მორალური აგენტის შექმნა ევოლუციური გამოთვლის გზით</i>	<i>24</i>
<i>ზემოდან ქვემოთ</i>	<i>25</i>
<i>ქვემოდან ზემოთ</i>	<i>27</i>
<i>სათნო (virtuous) ხელოვნური მორალური აგენტი</i>	<i>28</i>
<i>AMMA</i>	<i>32</i>
<i>ევოლუციური ნეირონული ქსელები და ნეიროეთიკა</i>	<i>36</i>
არგუმენტები ხელოვნური აგენტობის წინააღმდეგ	37
<i>მორალური ეკოლოგიის კონტრარგუმენტები</i>	<i>37</i>
<i>ხელოვნური მორალური აგენტი ლევენასის შუქზე</i>	<i>38</i>
<i>არის თუ არა ხელოვნური არტიფაქტი მოაზროვნე აგენტი</i>	<i>41</i>
<i>ჩინური ოთახის ექსპერიმენტი</i>	<i>43</i>
<i>შეპირისპირება ჩარჩლენდთან</i>	<i>46</i>
<i>კონტრარგუმენტების შეჯამება</i>	<i>50</i>
დასკვნა	51
<i>სხვის აგენტობის წინააღმდეგ</i>	<i>52</i>

<i>მოდელირება და ცნობიერება</i>	55
<i>ანტიპოდიანელი მორალური აგენტი</i>	58
გამოყენებული ლიტერატურა	61

Abstract

In our time technologies penetrates every sphere of human activities. The moral sphere is no different. It's become increasingly important to analyze whether it is possible to have artificially made moral agent or not. Reason for that is increasing in usage of technological artifacts which makes it impossible for them, to be excluded from moral sphere. To answer this question there are two distinct ways to do so. First explores possibility of different moral agent other than human. Second is to model artificial agent using the AI computer system. Both have its pros and cons, but at the end both failed to accomplish impossible. The impossible is to give machine intentions without it, it is impossible to have an agent which acts in accordance of its own will, which is vital for morality, as described by Campbell. Only viable option, then becomes to simulate moral deliberation and knowing that it is a limited one. I call it an antipodean moral agent. Agent without self-consciousness and own intentions, an agent which mimics to be conscious and moral without any understanding of variables it manipulates.

Keywords:

Artificial moral agent, Modeled moral deliberation, Chinese room argument, Consciousness in artificial agent, Artificial agent, morality, Moral deliberation.

აბსტრაქტი

ჩვენ ოპოქაში ტექნოლოგიების გამოყენების არეალი იზრდება ყოველწამიერად. მიკრო-ჩიპებით თუ ნულებისა და ერთების კომბინაციით მართულ სამყაროში, ბუნებრივია ტექნოლოგიები მორალის სფეროშიც შემოდინან. ამ ფონზე მნიშვნელოვანი ხდება ხელოვნური მორალური აგენტობის დადგენა. ამ მიზნით არსებობს ორი ძირითადი მიდგომა. ერთი სწავლობს ზოგადად ნებისმიერი სახის, ადამიანისგან განსხვავებულ მორალურ აგენტობას, ხოლო მეორე კი ადამიანის მსგავსი, მორალური აგენტის ხელოვნურად შექმნის შესაძლებლობას. ორივე მიმართულებას აქვს თავის უპირატესობები, ისევე როგორც ნაკლოვანებები. იმისათვის, რომ დავადგინოთ თუ რომელ მიმართულებას შეუძლია ამ თეორიული საკითხის გადაჭრა, რაც ნაშრომის მთავარი მიზანია, საჭიროა დეტალურად მიმოვიხილოთ ისინი. მიდგომები, რომლებიც მიმართულია მოახდინონ ადამიანის მორალური განსჯის მოდელირება, ვერ ახარხებენ ამას, რადგან მათ არ შეუძლიათ მორალურობის მთავარი ასპექტის, განზრახულობის კომპიუტერულ სისტემაში შექმნა. მეორეს მხრივ, სხვა ტიპის აგენტობის არცერთი თეორია არ არსებობს, რომელსაც შეუძლია დაასაბუთონ, უკანსკნელის შესაძლებლობა. შედეგად ჩვენ ვრჩებით ისეთი ხელოვნური აგენტობის იმედად, რომელიც არ არის ნამდვილი და შეიძლება მხოლოდ სიმულაციად ჩაითვალოს.

საკვანძო სიტყვები:

ხელოვნური მორალური აგენტი, მორალური განსჯის მოდელირება, სხვა ტიპის მორალური აგენტი, ხელოვნური ინტელექტი, ცნობიერება, ჩინური ოთახის არგუმენტი.