
Spis treści

Wstęp: wprowadzenie do bezpieczeństwa i ochrony sztucznej inteligencji	xi
Podziękowania	xxvii
Redaktor naukowy	xxix
Współpracownicy	xxxii

Część I Obawy luminarzy

Rozdział 1	Dlaczego przyszłość nas nie potrzebuje	3
	<i>Bill Joy</i>	
Rozdział 2	Głęboko przeplatana obietnica i niebezpieczeństwo GNR	25
	<i>Ray Kurzweil</i>	
Rozdział 3	Podstawowe pobudki SI	59
	<i>Stephen M. Omohundro</i>	
Rozdział 4	Etyka sztucznej inteligencji	71
	<i>Nick Bostrom i Eliezer Yudkowsky</i>	
Rozdział 5	Przyjazna sztuczna inteligencja: Wyzwanie fizyki	87
	<i>Max Tegmark</i>	
Rozdział 6	MDL destylacja inteligencji: Poznawanie strategii bezpiecznego dostępu do superinteligentnych możliwości rozwiązywania problemów	93
	<i>K. Eric Drexler</i>	
Rozdział 7	Problem uczenia się wartości	111
	<i>Nate Soares</i>	
Rozdział 8	Przykłady kontrydiktoryjne w świecie fizycznym	123
	<i>Alexey Kurakin, Ian J. Goodfellow i Samy Bengio</i>	
Rozdział 9	W jaki sposób może zaistnieć SI? Różne podejścia i ich implikacje dla życia we wszechświecie	141
	<i>David Brin</i>	

Rozdział 10	Przyszłość MADCOM: Jak sztuczna inteligencja może wzmocnić propagandę obliczeniową, przeprogramować ludzką kulturę oraz zagrozić demokracji... i co można z tym zrobić	159
	<i>Matt Chessen</i>	
Rozdział 11	Strategiczne implikacje otwartości w rozwoju sztucznej inteligencji ...	183
	<i>Nick Bostrom</i>	
 Część II Odpowiedzi naukowców		
Rozdział 12	Korzystanie z ludzkiej historii, psychologii i biologii w celu uczynienia SI bezpieczną dla ludzi	211
	<i>Gus Bekdash</i>	
Rozdział 13	Bezpieczeństwo SI z perspektywy pierwszej osoby	251
	<i>Edward Frenkel</i>	
Rozdział 14	Strategie dla nieprzyjaznej wyroczni SI z przyciskiem resetowania	261
	<i>Olle Häggström</i>	
Rozdział 15	Zmiany celu w inteligentnych agentach	273
	<i>Seth Herd, Stephen J. Read, Randall O'Reilly i David J. Jilk</i>	
Rozdział 16	Ograniczenia weryfikacji i walidacji zachowań agencyjnych	283
	<i>David J. Jilk</i>	
Rozdział 17	Kontrydiktoryjne uczenie maszynowe	295
	<i>Phillip Kuznetsov, Riley Edmunds, Ted Xiao, Humza Iqbal, Raul Puri, Noah Golmant i Shannon Shih</i>	
Rozdział 18	Uzgadnianie wartości wykorzystując obliczalną odległość preferencji	313
	<i>Andrea Loreggia, Nicholas Mattei, Francesca Rossi i K. Brent Venable</i>	
Rozdział 19	Racjonalnie uzależniona sztuczna superinteligencja	329
	<i>James D. Miller</i>	
Rozdział 20	Bezpieczeństwo aplikacji robotów z wykorzystaniem ROS	341
	<i>David Portugal, Miguel A. Santos, Samuel Pereira i Micael S. Couceiro</i>	

Spis treści	ix
Rozdział 21 Wybór preferencji społecznej i problem wyrównania wartości	363
<i>Mahendra Prasad</i>	
Rozdział 22 Rozłączne scenariusze katastrofального ryzyka SI	395
<i>Kaj Sotala</i>	
Rozdział 23 Realizm ofensywny i niezabezpieczona struktura systemu międzynarodowego: Sztuczna inteligencja i globalna hegemonia	423
<i>Maurizio Tinnirello</i>	
Rozdział 24 Superinteligencja i przyszłość rządów: Priorytetyzacja problemu kontroli na końcu historii	445
<i>Phil Torres</i>	
Rozdział 25 Wojskowa SI jako zbieżny cel samodoskonalącej się SI	467
<i>Alexey Turchin i David Denkenberger</i>	
Rozdział 26 Wrażliwe na wartości podejście do projektowania inteligentnych agentów	491
<i>Steven Umbrello i Angelo F. De Bellis</i>	
Rozdział 27 Konsekwencjalizm, deontologia i bezpieczeństwo sztucznej inteligencji	509
<i>Mark Walker</i>	
Rozdział 28 Inteligentne maszyny są zagrożeniem dla ludzkości	523
<i>Kevin Warwick</i>	
Indeks	533