

## Παιδαγωγικός Λόγος

Τόμ. 32, Αρ. 1 (2026)

Λόγος περί της Τεχνητής Νοημοσύνης



### Τα συστήματα Τεχνητής Νοημοσύνης και η οντολογική προσέγγιση του Karl R. Popper

*Ειρήνη Δαρκαδάκη*

doi: [10.12681/plogos.33774](https://doi.org/10.12681/plogos.33774)

Copyright © 2026, Eirini Darkadaki



Άδεια χρήσης [Creative Commons Attribution-NonCommercial-NoDerivatives 4.0](https://creativecommons.org/licenses/by-nc-nd/4.0/).

### Βιβλιογραφική αναφορά:

Δαρκαδάκη Ε. (2026). Τα συστήματα Τεχνητής Νοημοσύνης και η οντολογική προσέγγιση του Karl R. Popper. *Παιδαγωγικός Λόγος*, 32(1), 11-29. <https://doi.org/10.12681/plogos.33774>

Ειρήνη ΔΑΡΚΑΔΑΚΗ

*Τα συστήματα Τεχνητής Νοημοσύνης  
και η οντολογική προσέγγιση  
του Karl R. Popper*

---

doi:<https://doi.org/10.12681/plogos.33774>

---

*Εισαγωγή*

**Η** ΤΑΧΕΙΑ ΑΝΑΠΤΥΞΗ ΤΗΣ ΤΕΧΝΗΤΗΣ ΝΟΗΜΟΣΥΝΗΣ (TN) ΕΧΕΙ ΜΕταβάλλει ριζικά τον τρόπο ζωής και αλληλεπίδρασης των ανθρώπων, εγείροντας βαθιά φιλοσοφικά και ηθικά ερωτήματα. Στο κέντρο της συζήτησης βρίσκεται η δυνατότητα των συστημάτων TN να θεωρηθούν ηθικά υποκείμενα, δηλαδή οντότητες ικανά να λαμβάνουν ηθικές αποφάσεις και να αναλαμβάνουν ευθύνη για τις πράξεις τους. Αυτή η μελέτη εξετάζει το ζήτημα μέσα από το πρίσμα της οντολογικής θεωρίας των τριών κόσμων του Karl R. Popper, η οποία προσφέρει έναν ιδιαίτερα κατατοπιστικό θεωρητικό πλαίσιο για την κατανόηση της φύσης της TN.

Ο σκοπός αυτής της εργασίας είναι διττός: αφενός, να αναδειξει τις δυνατότητες και τους περιορισμούς της TN ως ηθικού παράγοντα μέσω της Ποππεριανής προσέγγισης, και αφετέρου, να προτείνει μία μεταηθική θεώρηση που θα επιτρέπει την καλύτερη κατανόηση των ηθικών προκλήσεων που εγείρονται από την αυτονομία της TN. Με αυτόν τον τρόπο, η μελέτη στοχεύει να συμβάλει στη συνεχιζόμενη φιλοσοφική και ηθική συζήτηση γύρω από τη θέση της TN στην ανθρώπινη κοινωνία και τους τρόπους με τους οποίους μπορούμε να διαχειριστούμε τις επιπτώσεις της.

## *Μια μεταθητική θεμελίωση μέσω της οντολογικής προσέγγισης των τριών κόσμων του Karl R. Popper*

### *Οι Κόσμοι του Karl Popper*

Οι ιδέες ή αξίες στην συλλογιστική του Popper δεν υπάρχουν a priori της ανθρώπινης ύπαρξης γιατί εάν συνέβαινε κάτι τέτοιο, υποστηρίζει πως δεν θα μπορούσαμε να τις προσεγγίσουμε<sup>1</sup>. Αυτή η θεμελιώδης διευκρίνιση του φιλοσόφου, τον διαχωρίζει από τις αντίστοιχες πλατωνικές θεωρήσεις και κάνουν φανερή τη ρεαλιστική στάση που υιοθετεί.

Η βασική διάκριση που προτείνει στην οντολογική του θεμελίωση της γνώσης είναι αναγκαίο να αναφερθεί εξαρχής προκειμένου τα όσα ειπωθούν στη συνέχεια καθώς και η απόπειρα ηθικής σύνδεσης με τον κόσμο της τεχνητής νοημοσύνης να γίνουν κατανοητά. Διαχωρίζει λοιπόν το σύμπαν (εάν θα μπορούσαμε να το θέσουμε ως τέτοιο) ή αλλιώς τον κόσμο, σε τρία βασικά επίπεδα<sup>2</sup>.

Ο Popper διαχωρίζει τον περιβάλλοντα κόσμο σε «Κόσμος 1= φυσικός κόσμος των φυσικών αντικειμένων», «Κόσμος 2= κόσμος των υποκειμενικών εμπειριών (νοητικός κόσμος/ κόσμος των ψυχικών καταστάσεων)» και «Κόσμος 3= κόσμος των προτάσεων καθεαυτών<sup>3</sup> (εδώ εντάσσονται τα σύνολα προτάσεων, των προϊόντων της γλώσσας εν γένει, της επιστήμης, των μαθηματικών και της φιλοσοφίας αλλά και των τεχνών)<sup>4</sup>. Στον Κόσμο 3 σύμφωνα με τον φιλόσοφο βρίσκεται η αντικειμενική γνώση για τα πράγματα<sup>5</sup>. Πρέπει λοιπόν να καταστεί σαφές πως ο Κόσμος 3 συνδέεται με την αντικειμενική γνώση μέσω της αποθήκευσης και της διάδοσης των γνώσεων που μπορούν να είναι ανεξάρτητες από τις υποκειμενικές εμπειρίες και ψυχικές καταστάσεις (Κόσμος 2). Με αυτόν τον τρόπο, ο κόσμος 3 παίζει κρίσιμο ρόλο στην επιστημονική πρόοδο και την αναζήτηση της

---

<sup>1</sup> Στέλιος Βιρβιδάκης, «Ο Κόσμος 3, του Καρλ Πόππερ ως ερμηνευτική πρόταση για την κατανόηση των μοντέλων μετριοπαθούς αξιακού ρεαλισμού», *Φιλοσοφία* 49, νο. II (2020): 198-200.

<sup>2</sup> Karl R. Popper, “*Objective knowledge: an evolutionary approach*” (New York: Oxford University Press, 1972), 106-107.

<sup>3</sup> ό.π., 154-157 “... there are three worlds: the first is the physical world, the world of physical states; the second is the mental world or the world of mental states; and the third world is the world of intelligibles, or the world of ideas in the objective sense; it is the world of possible objects of thought..”

<sup>4</sup> Στέλιος Βιρβιδάκης, «Ο Κόσμος 3 του Καρλ Πόππερ», 194-207.

<sup>5</sup> Karl R. Popper, “*Objective knowledge*”, 106- 107.

αντικειμενικής γνώσης. Ο Popper όπως γίνεται φανερό από τα άνωθι, τοποθετείται με την πλευρά της ρεαλιστικής προσέγγισης του ζητήματος και συνεπώς θεωρεί πως υπάρχουν υπερβατικές ηθικές αξίες (αν όχι τουλάχιστον μία ολιστική) επάνω στις οποίες βασίζονται οι ηθικές επιταγές των πρακτικών των ανθρώπων και δύνανται κατά αυτόν τον τρόπο να οριστούν, και αυτές εντάσσονται στον Κ3<sup>6</sup>. Εκ πρώτης όψεως η ανάλυση και προσέγγιση του ζητήματος φαίνεται να ομοιάζει με την πλατωνική θεωρία των ιδεών αλλά κάτι τέτοιο καθίσταται σαφές από τον ίδιο τον φιλόσοφο πως είναι ανυπόστατο<sup>7</sup>. Προκειμένου αυτός ο διαχωρισμός να γίνει κατανοητός αξίζει αρχικώς να αναφερθεί πως ο Popper δεν αναφέρεται σε αναλλοίωτες ιδέες ή αξίες που υπάρχουν σε μία διαφορετική διάσταση όπως γίνεται στη πλατωνική θεώρηση των ιδεών.

Αντιθέτως, αναζητά και αποπειράται να αποδείξει την ενοποίηση αυτών των φαινομενικά ασύνδετων κόσμων, την ενοποίηση δηλαδή της φυσικής και υποκειμενικής γνώσης του φυσικού κόσμου (Κ1 και Κ2) με τον κόσμο της απόλυτης γνώσης (Κ3). Αυτό, το κατορθώνει υποστηρίζοντας ότι στον Κόσμο 1 θεωρούμε πραγματικό ό,τι ενεργεί επάνω σε φυσικά πράγματα ή ό,τι υφίσταται κάποιου είδους ενέργεια από αυτά. Τα νοητικά καθατά αντικείμενα του Κόσμου 3 (π.χ. θεωρίες) ασκούν επίδραση στον φυσικό κόσμο και τον αλλάζουν, και άρα, αναγκαία πρέπει να θεωρηθούν πραγματικά παρότι δεν είναι εμπειρικώς αντιλήψιμα. Η σύνδεση αυτών των δύο κόσμων στη θεωρία του Popper επικυρώνεται μέσω του Κόσμου 2 που θα μπορούσε να ερμηνευθεί ως «μεταβατικός κόσμος»<sup>8</sup>.

Πιο αναλυτικά, προκειμένου να επιτευχθεί μία ενοποίηση αυτών των τριών ειδών της πραγματικότητας, ο Popper τοποθετείται με σαφήνεια για την ποιότητα και το περιεχόμενο του τρίτου κόσμου που ορίζει. Αναφέρει πως ο Κόσμος 3 ομοιάζει με την αντίστοιχη θεωρία αντικειμενικών ποιότητων που περιέχονται στην φιλοσοφική θεώρηση του Frege<sup>9</sup> και άρα θέλει να αποδείξει με αυτή του την απόπειρα πως η επιστημολογία όπως εφαρμόζεται έως σήμερα βασίζει τα πορίσματά της σε υποκειμενικές θεωρήσεις του κόσμου καθώς μελετά συνιστώσες που ο ίδιος εντοπίζει πως βρίσκονται μόνο στον Κ2 και όχι στον Κ3. Ο Κ2 αποτελεί επομένως το σημερινό πεδίο ερευνών της επιστήμης, κάτι το οποίο σύμφωνα με τον φιλόσοφο αποτελεί λανθασμένη προσέγγιση της γνώσης καθότι μόνο στον

<sup>6</sup> ό.π., 158-161.

<sup>7</sup> ό.π., 106/ 154.

<sup>8</sup> Βιββιδάκης, «Ο Κόσμος 3», 198-199.

<sup>9</sup> Karl R. Popper, “Objective knowledge”, 106-107.

K3 εμπεριέχονται οι αλήθειες και οι αντικειμενικές θεωρήσεις. Προκειμένου να αιτιολογήσει αυτόν του το διαχωρισμό προχωράει σε μία επιχειρηματολογία βασισμένη σε τρεις προκειμένες:

(Π1) : Η παραδοσιακή επιστημολογία κρίνεται ανεπαρκής στην εξήγηση των φυσικών φαινομένων και αληθειών καθώς κατά την μελέτη φυσικών φαινομένων στον K1 μεταβαίνουν στον K2 για να τα αποδείξουν και όχι στον K3 όπου βρίσκεται η αληθινή επιστημονική γνώση<sup>10</sup>.

(Π2) : Είναι αναγκαία και υποχρεωτική η μελέτη του K3 για την αντικειμενική ανεύρεση των αληθειών που διέπουν τον K1 καθώς ο K3 αποτελεί έναν αυτόνομο κόσμο εντελώς εκτός της φυσικής πραγματικότητας, όχι όμως με την πλατωνική σημασία της θεωρίας των ιδεών<sup>11</sup>.

Για την Π2 διευκρινίζει πως οι επιστήμονες πράττουν βασιζόμενοι σε πεποιθήσεις ή υποθέσεις εις άτοπον απαγωγής οι οποίες όμως εδράζονται σε μία υποκειμενική θεώρηση της πραγματικότητας από τους ίδιους, και όχι στην αντικειμενική ποιότητα των αληθειών που θα τους παρέχονταν αν μελετούσαν τα ίδια αντικείμενα μέσω του K3<sup>12</sup>.

(Π3) : Η μελέτη του K3 αποτελεί κρίσιμης σημασίας ζήτημα καθώς μπορεί να παρέχει μία αντικειμενική σκοπιά θέασης η οποία θα έχει ως αποτέλεσμα να δια φωτίσει πλευρές του κόσμου οι οποίες μάς είναι εντελώς άγνωστες και προς το παρόν δεν έχει ανευρεθεί κάποιος τρόπος εξήγησής τους. Ταυτόχρονα, ο K3 θα αιτιολογήσει την υποκειμενική σκοπιά θεώρησης που υιοθετούν οι επιστήμονες (δηλαδή τον K2) προσφέροντας μία πιο αντικειμενική προσέγγιση η οποία κατ' επέκταση οδηγεί σε μία πιο έγκυρη εφαρμογή της μεθοδολογίας των επιστημόνων<sup>13</sup>. Η έλευση στον K3 μέσω του K2, αποτελεί λανθασμένη πεποίθηση και ορίζεται ως ανέφικτη<sup>14</sup>.

(Σ) : Η μελέτη του K3 θα έπρεπε να αποτελεί αποκλειστικά το κέντρο

<sup>10</sup> Karl R. Popper, “*Objective knowledge*”, “Scientific knowledge”, 111.

<sup>11</sup> Στο ίδιο, 1-10/ 106.

<sup>12</sup> Το ερώτημα που προκύπτει από τις άνωθι προκειμένες και αφορά το ζήτημα της παρούσας εργασίας είναι το πώς κανείς αποκτά πρόσβαση στον K3 και εάν, τα συστήματα τεχνητής νοημοσύνης μπορούν όπως ακριβώς και στη περίπτωση του ανθρώπου, να έχουν πρόσβαση σε αυτόν ή τουλάχιστον κάποιου είδους σύνδεσης με αυτόν.

<sup>13</sup> Επομένως η αναγκαία συνεπαγωγή που προκύπτει ορίζει επειδή  $K3 = T$  (T ορίζεται ως True δηλαδή τιμή αληθείας της συνεπαγωγής), τότε ισχύει ότι  $K3 \rightarrow K2$  και όχι ότι  $K2 \rightarrow K3$  γιατί η K2 προκειμένου να ισχύει προϋποτίθεται να πηγάζει από την K3 και όχι το αντίστροφο.

<sup>14</sup> Karl R. Popper, “*Objective knowledge*”, 107.

μελέτης όλων των επιστημόνων και οποιοσδήποτε άλλος τρόπος προσέγγισης αυτού οδηγεί σε πλάνες και υποκειμενικές θεωρήσεις που μας απομακρύνουν από την αλήθεια<sup>15</sup>.

Προκειμένου να ενισχύσει τις άνωθι προκείμενες ο Karl Popper κάνει τρεις βασικές επισημάνσεις που αφορούν τον Κ3:

1. Ο Κόσμος 3 αποτελεί παράγωγο της ανθρώπινης οντότητας, το οποίο το παρομοιάζει με τον ιστό της αράχνης<sup>16</sup>. Όπως γύρω από την αράχνη, η οποία βρίσκεται στο κέντρο του ιστού της, ξεδιπλώνονται όλες οι διαστάσεις του ιστού της σε ευρεία κλίμακα και ποικιλία πάχους, έτσι και στη περίπτωση του Κ3, ο άνθρωπος βρίσκεται στο κέντρο και περιβάλλεται από ένα σύμπαν το οποίο έχει παραχθεί από τον ίδιο ακόμη και ο ίδιος ο άνθρωπος δεν έχει επίγνωση αυτού<sup>17</sup>.
2. Ο Κόσμος 3 αποτελεί ως επί το πλείστον μία αυτόνομη προέκταση της ύλης η οποία επιδρά στον φυσικό κόσμο και στον άνθρωπο και αντίστοιχα ο φυσικός κόσμος, και κατ' επέκταση ο άνθρωπος, επηρεάζεται και πράττει με κινητήριο δύναμη αυτόν τον ίδιο<sup>18</sup> και μάλιστα στη δεύτερη περίπτωση, ασκεί καταλυτική επίδραση στον

---

<sup>15</sup> Προκειμένου να υποστηρίξει το συμπέρασμά του και την επιδραστικότητα αλλά και σύνδεση του Κ3 στον Κ1 προχωρεί στην εξής εγκυροποίηση: "In our attempts to solve these other problems we may invent new theories. These theories, again, are produced by us: they are the product of our critical and creative thinking (K2), in which we are greatly helped by other existing third-world theories. Yet the moment we have produced these theories, they create new, unintended and unexpected problems, autonomous problems, problems to be discovered." Karl R. Popper, "*Objective knowledge*", 161.

<sup>16</sup> Στο ίδιο, 112.

<sup>17</sup> Η ανάλυση που έχει γίνει μέχρι αυτό το σημείο γεννά το βασικό ερώτημα κατά πόσον η οντολογική θεώρηση του Popper προσφέρει ένα εύφορο έδαφος προκειμένου να ανακαλυφθεί μία αντίστοιχης υφής σύνδεση των συστημάτων TN και του Κ3. Εκ πρώτης όψεως η θεώρηση αυτή φαίνεται ατελέσφορη για αυτήν την απόπειρα, όμως δημιουργεί ένα έδαφος αναλογικού συλλογισμού πάνω στο οποίο μπορούμε να στηρίξουμε μία θεωρία αντίστοιχης φύσεως που να αφορά τα συστήματα TN αλλά να μη σχετίζεται με τον Κ3, να δρα όμως παράλληλα και σε αντιδιαστολή με αυτόν. Αυτό είναι και το δεύτερο σημείο εστίασης της παρούσας εργασίας. Ο λόγος λοιπόν για τον οποίο χρησιμοποιείται αυτή η θεώρηση στη παρούσα εργασία είναι για να αναδείξει την διαφορετική λειτουργία των συστημάτων TN από τους ανθρώπους και να χρησιμοποιηθεί με τρόπο αναλογικό ώστε να δημιουργηθεί ένα νέο οντολογικό γίγνεσθαι στο οποίο θα μπορούν να αποδοθούν ηθικές ποιότητες σε αυτά τα συστήματα, αντίστοιχες με αυτές του ανθρώπου.

<sup>18</sup> Karl R. Popper, "*Objective knowledge*", 112-113.

φυσικό κόσμο (Κ1) ώστε να μπορεί να υφίσταται με τον τρόπο ακριβώς που συμβαίνει.

Ενώ φαίνεται από τη θεώρηση του Popper πως ο Κ3 παράγεται από τον άνθρωπο, τον ορίζει ως αυτόνομο, και η αιτία αυτού, εδράζεται στο γεγονός ότι ο Κ3 περιέχει το πεδίο της καθαρής γνώσης κάτι που του δίνει το χαρακτηριστικό της αυτονομίας. Προς απόδειξη αυτού του επιχειρήματος αναφέρει τα βιβλία και την ανάγνωσή τους από διαφορετικούς ανθρώπους<sup>19</sup>. Ο κάθε άνθρωπος διαβάζοντας ένα βιβλίο θα δώσει μία διαφορετική ερμηνευτική προσέγγιση σε αυτό που διάβασε αυτό όμως δεν αναιρεί την καθαρή πληροφορία γνώσης στην οποία βασίζεται το βιβλίο. Όπως στην περίπτωση των βιβλίων, έτσι και στη περίπτωση του τρίτου κόσμου, μπορεί να αποτελεί παραγόμενο προϊόν της ανθρώπινης ύπαρξης αλλά εντός αυτού υπάρχουν θεωρίες, αλήθειες, γνώση εν γένει οι οποίες δεν βασίζονται στην παραγωγή τους από τον άνθρωπο και μπορούν ποτέ να μην γίνουν αντιληπτές ή κατανοητές από αυτόν. Ένα ακόμη τέτοιο παράδειγμα απόδειξης της αυτονομίας του Κ3 αποτελεί η ίδια η γλώσσα καθαυτή ή η θεωρία των φυσικών αριθμών.

3. Μόνο μέσω της αλληλεπίδρασης των ανθρώπων με τον Κ3 μπορεί η αντικειμενική γνώση να ευδοκιμήσει ενώ ταυτόχρονα μέσω αυτής δύναται να αναπτύσσεται και ο Κ1. Επί παραδείγματι, αυτό εντοπίζεται στη βιολογική εξέλιξη των φυτών και των ζώων<sup>20</sup>.

Αυτή η τρίτη διευκρίνιση δημιουργεί ένα πρόσφορο έδαφος ώστε να εξεταστεί το εάν ένα σύστημα τεχνητής νοημοσύνης μπορεί να αλληλεπιδρά και να επηρεαστεί από τον Κ3. Ο φιλόσοφος σε μία απόπειρα να στηρίξει την αυτονομία του Κ3 προβαίνει σε μία επιχειρηματολογία βασισμένος στην επιστήμη της βιολογίας. Αρχικώς αναφέρει πως η επιστήμη της βιολογίας μελετά κυρίως α) την συμπεριφορά και φυσική δομή των έμβιων όντων και β) κάποιες από τις μη έμβιες οντότητες που προκύπτουν από τα φυτά και τα ζώα όπως τις φωλιές τους ή τα μονοπάτια που χαράζουν στο δάσος<sup>21</sup>. Η οπτική του Popper εστιάζει στην ανάδειξη της ανάπτυξης της

---

<sup>19</sup> Στο ίδιο, 116-117.

<sup>20</sup> Στο ίδιο, 112-113.

<sup>21</sup> ό.π.

γνώσης στο πλαίσιο της εξέλιξης της έμβιας ζωής<sup>22</sup>. Τα προβλήματα που εγείρονται από αυτόν τον τρόπο προσέγγισης του αντικειμένου μελέτης των βιολόγων, είναι κυρίως δύο:

A) Τα ζητήματα που προκύπτουν από τα αποτελέσματα των πράξεων των έμβιων όντων και

B) Τα ζητήματα που προκύπτουν εάν ληφθούν υπόψη οι δομές των έμβιων όντων καθ'αυτές<sup>23</sup>.

Για την ανθρωπότητα τα A) και B) μπορούν αναλογικά να τεθούν στην γλώσσα και την επιστήμη. Αυτά τα ζητήματα ο Popper υποστηρίζει πως είναι αληθή και εντοπίζονται και για τους ανθρώπους: η ανθρωπότητα επίσης έχει δημιουργήσει νέα είδη που εμφανίζονται στον K1 ή στον K2 και τα ορίζει ως «πνευματικά προϊόντα» (intellectual products) τα οποία δομούν το περιβάλλον του ανθρώπου. Τέτοια ορίζονται ως οι μύθοι, τα λογοτεχνικά βιβλία, τα έργα τέχνης ή επιστημονικές θεωρίες που έχουν διατυπωθεί και ισχύουν για τον εμπειρικό κόσμο. Όταν αναφερόμαστε σε μία θεωρία εξελικτικού περιεχομένου αυτά τα πνευματικά προϊόντα πρέπει να λαμβάνονται ως αντικείμενα μίας πραγματικότητας έξω από εμάς, και μαζί με αυτά εντοπίζεται και η αληθινή γνώση (knowledge)<sup>24</sup>. Επομένως

<sup>22</sup> “...Η προσέγγιση του συγκεκριμένου ζητήματος είναι να τοποθετήσει τον τρόπο ανάπτυξης της καθαρής γνώσης στο πλαίσιο της εξέλιξης των ζώων και του ανθρώπου”, Stephen, Thornton “Karl Popper” στο *The Stanford Encyclopedia of Philosophy*, επιμ.: Edward N. Zalta & Uri Nodelman, Winter, 2022, δική μου μετάφραση. <https://plato.stanford.edu/archives/win2022/entries/popper/>.

<sup>23</sup> Σύμφωνα με τη συλλογιστική του Karl Popper τα ζητήματα του A) μπορούμε να πούμε ότι εντάσσονται στον K1, ενώ οι δομές που αναφέρονται στο B) εντάσσονται στον K2. Παρόλα αυτά, είναι δυνατόν να εντοπίσει κανείς πως οι δομές που αναφέρονται στο B) αποτελούν προϊόντα του εμπειρικού κόσμου (π.χ. η δομή και η λειτουργία ενός κυττάρου) και άρα με βεβαιότητα θα μπορούσε κάποιος να συμπεράνει πως εντάσσονται στον K1. Τα πορίσματα όμως που βγαίνουν μέσω της παρατήρησης των λειτουργιών των δομών αυτών καθώς και οι γενικεύσεις που προκύπτουν μέσω αυτών είναι αναμφιβόλως προϊόντα του K2 και όχι του K1. Άρα εν τέλει μπορούμε να συμπεράνουμε πως το B) είναι μία παράμετρος που μπορεί να ενταχθεί και στον K1 αλλά και στον K2 ίσως και ταυτόχρονα.

<sup>24</sup> “Αυτό, ο Popper αναφέρει, είναι αληθές και για την περίπτωση των ανθρώπων: εμείς επίσης έχουμε δημιουργήσει νέα είδη προϊόντων, «πνευματικά προϊόντα», τα οποία δομούν το περιβάλλον μας. Τέτοια είναι οι μύθοι, οι ιδέες μας, τα προϊόντα της τέχνης μας, και οι επιστημονικές μας θεωρήσεις για τον κόσμο στον οποίο ζούμε. Όταν αυτή η σκέψη τοποθετείται στο εξελικτικό πλαίσιο, ο Popper προτείνει αυτού του είδους τα προϊόντα πως είναι αναγκαίο να συλλαμβάνονται οργανικά, ως εξωσωματικά κατασκευαστικά προϊόντα. Η ενοποιός δύναμή τους είναι η γνώση”, Stephen, Thornton “Karl!”, δική μου μετάφραση.

τα αποκλήματα της φαντασίας του ανθρώπου ο Popper τα κατατάσσει στον Κ3 και τα οποία μέσω του Κ2 συνδέονται με μία σχέση αποβλεπτικότητας<sup>25</sup> στον Κ1.

Σύμφωνα με τον φιλόσοφο είναι σημαντικότερο να αναζητηθεί η λύση του ζητήματος των δομών καθεαυτές (Β), γιατί εάν εξηγηθεί αυτό, θα έχει δοθεί μία αντικειμενική λύση και στα συμπεριφορικά ζητήματα που αναφέρονται στο Α). Η αναγωγή λοιπόν στη Β περίπτωση πρέπει αναγκαία και ικανά να γίνει στον Κ3 ώστε να μπορούν να εξηγηθούν οι δομές των έμβιων όντων καθεαυτές με έναν όσο το δυνατόν πιο αντικειμενικό τρόπο. Αυτή η θέση μπορεί να οριστεί ως το τρίτο ζήτημα (Γ) που κατέχει αντί-συμπεριφοριστική χροιά<sup>26</sup>.

Επομένως, αυτό που αποτελεί καθοριστικής σημασίας στην αιτιακή σχέση ανάμεσα στους κόσμους του Popper είναι πως με αυτή τη προσέγγιση επιτρέπεται στον φιλόσοφο να αναγάγει την ανάπτυξη και εξέλιξη της ανθρώπινης γνώσης ως μία εξελικτική διαδικασία με εξωσωματικές προσαρμογές<sup>27</sup> οι οποίες τελικώς αποτελούν μία λειτουργική διαδραστική διαδικασία ανάμεσα στη σχέση του Κόσμου 1 (Κ1) και νοητικού κόσμου (Κ2) ενώ ταυτόχρονα επιτυγχάνεται η ίδια διάδραση και με τον κόσμο της αντικειμενικής γνώσης (Κ3) ή αλλιώς του περιεχομένου της σκέψης<sup>28</sup>.

<sup>25</sup> Ο όρος "αποβλεπτικότητα" αναφέρεται στην ιδιότητα του νου να κατευθύνεται προς "κάτι", πράγματα, αντικείμενα (πραγματικά ή φανταστικά), γεγονότα, καταστάσεις, σχέσεις και ιδέες του κόσμου, τα οποία είναι ως επί το πλείστον εξωτερικά (Κ1), δηλαδή "εκτός" του ίδιου του νοήμονος όντος που νοεί τον κόσμο (Κ2). Για παράδειγμα, όταν σκεφτόμαστε, σκεφτόμαστε κάτι, όταν αποφασίζουμε, αποφασίζουμε κάτι, όταν θέλουμε, θέλουμε κάτι, και ούτω καθεξής. Ο σχηματισμός ενός συγκεκριμένου νοητικού περιεχομένου, ενός νοήματος (Κ3), αποτελεί θεμελιώδες χαρακτηριστικό αυτής της ιδιότητας. Βλέπε: Gounaris, A. (2011). Intentionality and the Emergence of Meaning. *Philosophia - Annual Journal of the Research Centre for Greek Philosophy of the Academy of Athens*, v.41, pp 319-321, 2011. ISSN 1105-2120

<sup>26</sup> Karl R. Popper, "Objective knowledge", 114-115.

<sup>27</sup> Με τον όρο εξωσωματικές προσαρμογές, ορίζω την αλληλεπίδραση των νοητικών γνώσεων με των εμπειρικών δεδομένων. Τα εμπειρικά δεδομένα, βρίσκονται εκτός της ανθρώπινης σκέψης και νοητικής διαδικασίας όμως παρόλα αυτά, μπορούν να ανατρέψουν θεωρητικά δεδομένα που υπάρχουν στο νου, και να τα αναδομήσουν ή να τα αναπροσαρμόσουν. Επομένως οι εξωσωματικές προσαρμογές αναφέρονται σε διαδικασίες που επιτελούνται εκτός της νοητικής σφαίρας του εαυτού.

<sup>28</sup> "Σε τελική ανάλυση, αυτό που αφορά και αναφέρεται η οντολογική επιστημολογία του Popper, είναι η αιτιακή αλληλεπίδραση ανάμεσα στους κόσμους: αυτή του επιτρέπει να αντικατοπτρίσει την εξέλιξη της ανθρώπινης γνώσης ως μια εξελικτική διαδικασία με εξωσωματικές προσαρμογές η οποία εν τέλει αποτελεί ένα παιχνίδι μεταξύ των σχέσεων

*Η ηθική θεμελίωση των συστημάτων Τεχνητής Νοημοσύνης και η κριτική τους μέσω των Κόσμων του Karl R. Popper*

Όπως περιγράφεται στην προηγούμενη ενότητα ο τρόπος διαχωρισμού των κόσμων από τον Karl R. Popper, μπορεί να πραγματοποιηθεί και για τις ηθικές αλήθειες πάνω στις οποίες καλούμαστε να ανακαλύψουμε εάν μπορούν να υφίστανται στα συστήματα Τεχνητής Νοημοσύνης (TN). Τα συστήματα τεχνητής νοημοσύνης καθότι ενεργούν στον (K1), έχουν την ικανότητα να δρουν αποκλειστικά με γνώμονα αυτόν. Επομένως, είναι αναγκαίο να αναφερθεί πως σε πρώτο στάδιο πριν τη δημιουργία τους, ήταν νοητικά φαινόμενα (K2), και η ύπαρξή τους ως μέρη του φυσικού κόσμου βασίστηκε σε θεωρίες και επιχειρήματα που σύμφωνα με τον Popper πηγάζουν και κατατάσσονται στον K3 και άρα αποτελούν προϊόντα αυτού<sup>29</sup>.

Οι ηθικές αλήθειες, ενώ ανήκουν στον K3, έχουν αντίστοιχη θέση στον φυσικό κόσμο K1, καθώς η αποβλεπτικότητα τους μπορεί να εντοπιστεί εκεί. Αρχικώς, ως σύστημα θεωριών, αρχών και νομοθέτησης, αποτελούν μία νοητική διαδικασία που εδράζεται στον K2 και η οποία βασίζεται στον αντικειμενικό και καθολικό χαρακτήρα της K3<sup>30</sup>. Σε συνάρτηση με τα παραπάνω και με τα εμπειρικά φαινόμενα, η αρμονική συνύπαρξη των έμβιων όντων βασίζεται σε κανονιστικές ηθικές αρχές έξω από τον αισθητό κόσμο, οπότε η παρούσα εργασία εξετάζει το ζήτημα των ηθικών αρχών των αυτόνομων οπλικών συστημάτων από μία ηθικώς ρεαλιστική σκοπιά. Επιπρόσθετα σύμφωνα με τον Hage, οι προθέσεις και οι επιθυμίες μας έχουν αναγκαία υπόσταση και ύπαρξη, ανεξαρτήτως της αποβλεπτικότητάς τους στον φυσικό κόσμο<sup>31</sup>. Αυτό τεκμηριώνεται, εάν σκεφτούμε πως οι προθέσεις ή επιθυμίες μας αποτελούν βασικά κίνητρα για τις πράξεις μας ή τη δομή του συναισθηματικού μας κόσμου και της ιδιοσυγκρασίας μας. Μπορεί όμως ένα αυτόνομο οπλικό σύστημα να διαθέτει αυτά τα

---

ανάμεσα στον φυσικό και νοητικό/πνευματικό κόσμο με τον κόσμο της αντικειμενικής γνώσης ή το περιεχόμενο της σκέψης”, Stephen, “Karl”, δική μου μετάφραση.

<sup>29</sup> Ο K3 περιλαμβάνει τα πολιτιστικά δημιουργήματα, τις γνώσεις, τις ιδέες, τις επιστημονικές θεωρίες, τα μαθηματικά κατασκευάσματα, και άλλα προϊόντα της ανθρώπινης νόησης και διάνοιας τα οποία ενσωματώνονται στα συστήματα AI. Ένα σύστημα είναι αποτέλεσμα τεχνολογικών και επιστημονικών γνώσεων. Ναι μεν τα συστήματα αυτά σχεδιάζονται και χρησιμοποιούνται για να εκτελούν συγκεκριμένες λειτουργίες στον K1, όμως ενσωματώνουν τον K3.

<sup>30</sup> Επομένως και για την περίπτωση των συστημάτων τεχνητής νοημοσύνης και για την περίπτωση των ηθικών αληθειών ισχύει η ίδια συνεπαγωγή: επειδή  $K3 = T$  τότε ισχύει ότι:  $K3 \rightarrow K2 \rightarrow K1$ .

<sup>31</sup> Hage, *Theoretical Foundations*, 259.

χαρακτηριστικά; Είναι λοιπόν δυνατό, να υποθέσουμε πως τα συστήματα TN τα οποία έχουν αναπτυχθεί και στα οποία δεν υπάρχει ανθρώπινη παρέμβαση για τη λειτουργία τους (ή έστω είναι σε ελάχιστο βαθμό), διαθέτουν ή μπορούν να αναπτύξουν την ιδιότητα της ελεύθερης βούλησης και σύμφωνα με αυτή να επιτελέσουν ηθικές πράξεις και αποφάσεις<sup>32</sup>; Εάν υποθέσουμε κάτι τέτοιο, προκύπτει πως τα συστήματα TN θεωρούνται από τους ανθρώπους ως έλλογα όντα με ηθικό καθεστώς και επομένως διαθέτουν τη δυνατότητα αξιολόγησης και ενδεδειγμένης επεξεργασίας των προθέσεων τους προτού πράξουν. Άρα είναι ικανά να αξιολογούν ηθικά και να επιλέγουν σύμφωνα με τον Ηθικό Νόμο το πώς θα πράξουν. Το συμπέρασμα αυτό προκύπτει μέσω διαφόρων πειραμάτων που έχουν πραγματοποιηθεί ως σήμερα, και αποδεικνύουν πως τα συστήματα TN διαθέτουν την δυνατότητα να επιλέξουν έναν διαφορετικό τρόπο συμπεριφορικής λειτουργίας από αυτόν που ήταν προγραμματισμένα από τους επιστήμονες να τελέσουν<sup>33</sup>.

Προς επίρρωση του συγκεκριμένου επιχειρήματος απόδοσης ηθικής ευθύνης στα συστήματα TN, επικαλούμαι μία αναγκαία συνεπαγωγή όπως ορίστηκε από τον αναλυτικό φιλόσοφο Thomas Nagel, εάν συνδυάσουμε την αντικειμενική σκοπιά θεώρησης με την πράξη. Σύμφωνα με τη θεωρία του Thomas Nagel, εάν αποδεχθώ ότι ισχύει ότι *A*, τότε αναγκαία αποδέχομαι ότι ισχύει και *B*<sup>34</sup>. Στην περίπτωση λοιπόν των συστημάτων TN, εάν

<sup>32</sup> Κάτι τέτοιο είναι λογικό να μπορούμε να το υποθέσουμε μέσω των πειραμάτων που τελούνται με συστήματα τεχνητής νοημοσύνης και της αποβλεπτικότητας αυτών στον εμπειρικό κόσμο.

<sup>33</sup> Ένα παράδειγμα για αυτές τις περιπτώσεις αποτελούν τα παιχνίδια στρατηγικής (σκάκι ή AlphaGo-κορεάτικο παιχνίδι στρατηγικής παρόμοιο με το σκάκι-) και τα πειράματα που έχουν γίνει με την χρήση συστημάτων TN. Περισσότερες πληροφορίες για αυτά προκύπτουν αν κάποιος ανατρέξει στο ντοκιμαντέρ AlphaGo: <https://www.youtube.com/watch?v=WXuK6gekU1Y>.

<sup>34</sup> Thomas Nagel, *Η Θέα από το πουθενά*, (Αθήνα: Κριτική, 2002), 234-235. Ο Nagel σε αυτό του το βιβλίο προσπαθεί να αποδείξει τη σύνδεση της υποκειμενικής θεώρησης με την αντικειμενική θεώρηση της πραγματικότητας και πραγματεύεται ποικίλους στοχασμούς είτε από την υποκειμενική θεώρηση του κόσμου είτε προσπαθώντας να προσεγγίσει μία αντικειμενική σκοπιά των πραγμάτων. Κατά αυτόν τον τρόπο, δημιουργεί την αναφερθείσα αναγκαία (σύμφωνα με το συλλογισμό του) συνεπαγωγή, με σκοπό στο συγκεκριμένο κεφάλαιο του βιβλίου του να αποδείξει την ταυτότητα του εαυτού. Τα *A* και *B* επομένως μπορούν να εφαρμοστούν σε οτιδήποτε έχει αποβλεπτικότητα στον φυσικό κόσμο και παρατηρείται μία αναγκαία αλληλεπίδραση μεταξύ τους η οποία δεν θα μπορούσε να είναι αλλιώς. Προς περαιτέρω επεξήγηση μπορούν να αναφερθούν παραδείγματα από τα μαθηματικά όπως ότι αν δεχθώ ότι ισχύει ότι  $2+2$  ισούται με  $4$  τότε το άθροισμα αυτών των δύο φυσικών αριθμών δεν μπορεί να μας δώσει τίποτα άλλο πέρα

αποδεχθώ πως διαθέτουν ηθικό καθεστώς (Α), τότε αποδέχομαι και το γεγονός της απόδοσης ηθικής ευθύνης σε αυτά για τις πράξεις τους (Β)<sup>35</sup>. Εάν δεχθώ το παραπάνω, το πραγματώνω είτε α) για να είμαι έτοιμος ως έλλογο ον να δεχθώ τις συνέπειες των πράξεων που αυτά θα επιτελέσουν σε περίπτωση που αυτό κρίνεται αναγκαίο, είτε β) για να δεχθώ πως τα συστήματα αυτά που φαίνεται να διαθέτουν την ιδιότητα της ελεύθερης βούλησης, θα αναλάβουν την πλήρη αποδοχή των συνεπειών των πράξεών τους, επιβλαβών ή μη. Στη β) περίπτωση τα θεωρώ έλλογα όντα με αναφορά στον ηθικό νόμο για πράξη και άρα υπεύθυνα για τις πράξεις που αυτά επιτελούν.

Σε αυτή τη περίπτωση (β), στην οποία θεωρούνται *άξια* ανάληψης ευθυνών, αναγκαία θεωρούνται και *ικανά* για λήψη αποφάσεων με ηθικό καταλογισμό. Και στις δύο όμως περιπτώσεις (α και β) θα πρέπει να αναλογιστούμε και να είμαστε έτοιμοι για όλες τις αρνητικές εκβάσεις των πράξεών τους ακόμη και για εκείνες που δεν δυνάμεθα να φανταστούμε ότι θα πραγματώσουν σε μελλοντικό χρόνο. Το συγκεκριμένο επιχείρημα καθιστά σαφή τη σημαντικότητα του χρόνου σε αυτή τη μεταηθική θεμελίωση των αυτόνομων οπλικών συστημάτων, καθώς επίσης και το γεγονός πως η συνιστώσα του χρόνου<sup>36</sup> θα πρέπει να λαμβάνεται υπόψη ως παράγοντας επηρεασμού λήψης μίας απόφασης ενός συστήματος ΤΝ.

---

από 4. Μία τέτοια θεωρία μπορεί να επιβεβαιωθεί μέσω παραδειγμάτων στον φυσικό κόσμο είτε βρίσκομαι στην υποκειμενική, είτε στην αντικειμενική σκοπιά, με αποτέλεσμα να συμπεράνω πως είναι ένα παράδειγμα απόδειξης μίας θεωρίας που ισχύει καθολικά.

<sup>35</sup> ό.π.

<sup>36</sup> Όταν αναφέρομαι στην έννοια του χρόνου εννοώ το χρονικό διάστημα στο οποίο το σύστημα τεχνητής νοημοσύνης αποφασίζει a priori το πώς θα πράξει, την χρονική στιγμή της τέλεσης της πράξης, τις επιπτώσεις που προκύπτουν από αυτή -βραχυπρόθεσμες ή μακροπρόθεσμες- ακόμη και το πότε το ίδιο το σύστημα τεχνητής νοημοσύνης δημιουργήθηκε κάτι το οποίο ορίζει το πόσο εξελιγμένο είναι. Η παρούσα εργασία γράφεται αυτή τη χρονική περίοδο, κατά την οποία η τεχνολογική εξέλιξη αυτών των συστημάτων διαρκώς μεταβάλλεται και εξελίσσεται. Παράλληλα με αυτό μεταβάλλεται και εξελίσσεται ο τρόπος με τον οποίο μπορεί να αποδοθεί ηθικός καταμερισμός σε αυτά τα συστήματα. Επομένως, τίποτα δεν είναι απόλυτο και σίγουρο από τη στιγμή που υπάρχει μία διαρκής και συνεχόμενη μεταβολή και για αυτό διευκρινίζεται πως ο αστάθμητος παράγοντας του χρόνου δύναται να μεταποιηθεί τις ηθικές συνιστώσες που παρουσιάζονται εδώ, σε σχέση με τα συστήματα τεχνητής νοημοσύνης.

*Το ηθικό καθεστώς των συστημάτων Τεχνητής Νοημοσύνης*

Στο παρόν κεφάλαιο θα πραγματοποιηθεί ένας ορισμός της ηθικής ευθύνης όπως αυτή εμφανίζεται στις ανθρώπινες κοινωνίες, καθώς επίσης και το εάν τα συστήματα Τεχνητής Νοημοσύνης (TN) μπορούν να θεωρηθούν όντα με ηθικό καταλογισμό. Βασική προϋπόθεση για αυτή την απόπειρα αποτελεί μία στοιχειώδης αναφορά στους βαθμούς ανθρώπινης παρέμβασης των συστημάτων TN.

*Βαθμοί ανθρώπινης παρέμβασης και τα συστήματα TN*

Τα αυτόνομα συστήματα επιδέχονται βαθμούς ανθρώπινης παρέμβασης με τρεις τρόπους: human in the loop<sup>37</sup>, human on the loop<sup>38</sup> και human off the loop<sup>39</sup>. Εκτός από την περίπτωση human off the loop, υπάρχει η σκέψη και η δυνατότητα απόδοσης ηθικής ευθύνης σε ένα αυτόνομο σύστημα και στη περίπτωση human on the loop<sup>40</sup>. Θεωρείται πιθανό και ίσως πιο εύλογο να μπορεί να αποδοθεί η ευθύνη στα αυτόνομα οπλικά

<sup>37</sup> Human in the loop= η οποιαδήποτε εντολή επιδιώξει η μηχανή λόγω του αυτόνομου συστήματος που διαθέτει, να πράξει, χρειάζεται την έγκριση του ανθρώπου πριν τη πραγματοποιήσει. Ο άνθρωπος επιτελεί την ενεργό δράση οντότητα, και η μηχανή απλώς εκτελεί τις εντολές του

<sup>38</sup> Human on the loop= μεγαλύτερη αυτονομία της μηχανής, ο άνθρωπος δίνει την πρωταρχική εντολή και επιβλέπει.

<sup>39</sup> Human off the loop= Οι αξιολογήσεις, εντολές και οι εκτελέσεις, βασίζονται αποκλειστικά στη μηχανή η οποία αποτελεί τη δράση οντότητα, και ο άνθρωπος είναι εντελώς αποσυνδεδεμένος από οποιαδήποτε λειτουργία ή απόφασή της.

<sup>40</sup> Ποιος αναλαμβάνει την ευθύνη μιας λάθος εκτίμησης ενός αυτόνομου οπλικού συστήματος; Μπορεί να αποδοθεί ευθύνη στο ίδιο το σύστημα; Ή θα πρέπει να αποδοθεί στον κατασκευαστή, στον προγραμματιστή, ή ακόμη και στον στρατιωτικό που επέλεξε να θέσει αυτό το όπλο στη μάχη; Ακόμη και στη περίπτωση του human on the loop, είναι ηθικώς θεμιτό να προσδοθεί ηθικός καταμερισμός της ευθύνης στον εντολέα (άνθρωπος) ο οποίος όμως δεν έχει ευθύνη για τις μη ορθές σύμφωνα με μία κανονιστική ηθική θεωρία πρακτικές του αυτόνομου οπλικού συστήματος που χρησιμοποιήσε, καθότι το δεύτερο, λειτουργήσε με βάση τις προσλαμβανόμενες που αξιολόγησε με έναν δικό του τρόπο (black box paradox); Για το συγκεκριμένο πρόβλημα βλέπε: Γούναρης Α., & Κωστελέτος Γ. (2024). Γράφοντας τον αλγόριθμο του Καλού: Η Τεχνητή Νοημοσύνη ως μηχανή απόδοσης Δικαιοσύνης. Ηθική. Περιοδικό φιλοσοφίας, (19). <https://doi.org/10.12681/ethiki.39654>

Επίσης: Gounaris, A., Kosteletos, G. (2020). Licensed to Kill: Autonomous Weapons as Persons and Moral Agents. In Prole, D. and Rujević, G. (ed.). Personhood. Novi Sad, Filozofski Fakultet & The NKUA Applied Philosophy Research Lab Press. DOI: <https://doi.org/10.12681/aprpl.49>

συστήματα και στις δύο περιπτώσεις καθώς, είτε οι κατασκευαστές, είτε οι προγραμματιστές, κρίνεται αδύνατο και απολύτως παράλογο να επωμιστούν εξολοκλήρου την ευθύνη για πράξεις που έχουν παρθεί εξολοκλήρου από το ίδιο το σύστημα<sup>41</sup>.

Σε αυτή τη περίπτωση τα συστήματα τεχνητής νοημοσύνης που χρησιμοποιούνται σε πολεμικές συρράξεις κρίνεται αναγκαίο να θεωρηθούν έλλογα όντα καθώς μόνο μέσω αυτής της κατηγοριοποίησης θα είναι δυνατό να αναλάβουν ευθύνη και να επωμιστούν τις συνέπειες για τις πράξεις τους. Πώς όμως μπορεί να τους αποδοθεί ηθική ευθύνη<sup>42</sup>; Η φύση του ανθρώπου και ο τρόπος λειτουργίας του, τον καθιστούν ικανό για ηθικές πράξεις ή μη, σχεδόν *a priori*, βασιζόμενες στις κοινωνικές του πρακτικές και αλληλεπιδράσεις με τους άλλους ανθρώπους<sup>43</sup>. Αυτό, έχει ως αποτέλεσμα, η ηθικές του αποφάσεις και τα θεμέλια της ηθικής του σύμφωνα με τα οποία πράττει, να διαφοροποιούνται, ή (σε πρώτο στάδιο) να δομούνται και εξαιτίας εξωγενών παραγόντων.

#### *Το ηθικό καθεστώς στις ανθρώπινες κοινωνίες και τα συστήματα ΤΝ*

Οι άνθρωποι θεωρούνται υπεύθυνοι και τους προσδίδεται ένα ηθικό καθεστώς εξαιτίας του γεγονότος ότι αποκτούν ένα συγκεκριμένο ρόλο κοινωνικά. Η αιτιακή ευθύνη αποτελεί μία έκφανση αυτής, και συνδέεται άρρηκτα με το αν αποτελεί ο ίδιος ο πράττων την αιτία ενός αποτελέσματος. Σε αυτή τη περίπτωση υπάρχει άμεση σύνδεση αιτίου-αποτελέσματος. Μια πράξη ανατίθεται σε έναν αυτουργό, ο οποίος μετέπειτα θα αναλάβει την ευθύνη του αποτελέσματος αυτής του της πράξης. Το άτομο, αποτελεί υψίστης σημασίας καθώς είναι η αιτία του παραχθέντος αποτελέσματος, και πρέπει να είναι αξιόπιστο σε δύο τομείς: α) στο να μπορεί να επιτελέσει την πράξη με αξιοπιστία, και β) να μπορεί να αναλάβει την ευθύνη για το

<sup>41</sup> Jaap Hage, “Theoretical foundations for the responsibility of autonomous agents”, *Springer* (August 2017): 255.

<sup>42</sup> “...παραμένει αδιευκρίνιστο σε τι βαθμό η αυτονομία ενός οπλικού συστήματος είναι ή μπορεί να θεωρηθεί πως πράττει σύμφωνα με τις πληροφορίες που του έχουν εμφυτευτεί. Μία αιτία αυτής της έλλειψης γνώσης από την οποία προκύπτει αυτό το ερώτημα, εδράζεται στο βαθμό αυτονομίας, κάτι που μπορούμε να ισχυριστούμε και για τους ανθρώπους εάν σκεφτούμε την έλλειψη κινήτρων, εδράζεται στη φύση και στο μέγεθος της αυτονομίας κάτι το οποίο αποτελεί από μόνο του ένα αμφιλεγόμενο και δύσκολως κατανοητό ζήτημα.” Robert Sparrow, “Killer Robots”, *Journal of Applied Philosophy* Vol. 24, no. 1 (2007): 65, δική μου μετάφραση.

<sup>43</sup> Hage, “Theoretical foundations”, 256.

αποτέλεσμα της πράξης του<sup>44</sup>.

Η ικανότητα ανάληψης της ευθύνης του αποτελέσματος μίας πράξης, εγείρει το θεμελιώδες ερώτημα για το αν τα συστήματα ΤΝ είναι δυνατόν να διαθέτουν αυτή την ικανότητα (έστω στο μέλλον αν όχι τώρα) όπως αντιστοίχως συμβαίνει στις ανθρώπινες κοινωνίες. Δηλαδή, εάν έχουν ή μπορούν να αποκτήσουν ηθικό καθεστώς. Φυσικά, και στην περίπτωση των έλλογων όντων και στη περίπτωση των συστημάτων ΤΝ, το πότε και το εάν μία πράξη ή η αποδοχή της ευθύνης αυτής θεωρείται ηθικώς θεμιτή, μπορεί να εξεταστεί και αξιολογηθεί από εντελώς διαφορετικές οπτικές κάτι το οποίο μπορεί να αλλάξει την έκβαση μίας απόφασης.

Η βασική αιτία για αυτό, εντοπίζεται στο γεγονός πως οι άνθρωποι πράττουν εκ προθέσεως κινητοποιημένοι από τις προσταγές της ελεύθερης βούλησέως τους να πράξουν κατά το δοκούν<sup>45</sup>. Αντίθετα, τα συστήματα ΤΝ φαίνεται πως δεν διαθέτουν προδιαθέσεις για πράξεις όμοιες με αυτές του ανθρώπου και εξαιτίας αυτού, σε καμία περίπτωση προς το παρόν, δε μπορεί να θεωρηθεί πως διαθέτουν ελεύθερη βούληση<sup>46</sup> ακόμη και εάν προγραμματιστούν με σκοπό να αναπτύξουν μία που να ομοιάζει με ελεύθερη βούληση<sup>47</sup>. Η συμπεριφορά αυτών, βασίζεται στους ίδιους κανόνες με αυτούς των ανθρώπων σύμφωνα με την αποβλεπτικότητα αυτής στον εμπειρικό κόσμο. Ένα σύστημα τεχνητής νοημοσύνης το οποίο είναι προγραμματισμένο να επιτελεί εξολοκλήρου αυτόνομα τις πράξεις του, είναι δυνατόν εκ των αποτελεσμάτων αυτών να φανεί πως πράττει ορθά, και στις περισσότερες περιπτώσεις να μην είναι διακριτό το εάν πρόκειται για αποτελέσματα ανθρώπινων πράξεων ή ενός συστήματος τεχνητής νοημοσύνης<sup>48</sup>.

Τα συστήματα ΤΝ δομούν τις πράξεις τους σύμφωνα με τις αλγοριθμικές αλληλουχίες που επιτελούν και οι οποίες βασίζονται στον τρόπο με τον οποίο ο άνθρωπος έχει κατασκευάσει/προγραμματίσει αυτά τα

<sup>44</sup> Hage, “Theoretical foundations”, 257.

<sup>45</sup> οι οποίες για να θεωρηθούν ηθικές θα πρέπει να είναι σύμφωνες με τον ηθικό νόμο.

<sup>46</sup> Hage, “Theoretical foundations”, 258.

<sup>47</sup> Αιτία αυτού είναι το γεγονός πως δεν μπορούμε με απόλυτη βεβαιότητα να γνωρίζουμε εάν διαθέτουν μία εσωτερική προδιάθεση για πράξη η οποία να εδράζεται σε έναν στιβαρό και αναλλοίωτο Ηθικό Νόμο όπως αυτός ορίζεται στη καντιανή θεωρία, και όπως αυτός βλέπουμε να ισχύει για τους ανθρώπους -τους οποίους ορίζουμε ως έλλογα όντα που πράττουν σύμφωνα με τον Ηθικό Νόμο-.

<sup>48</sup> Luciano Floridi, “On the Morality of Artificial Intelligence”, *Minds and Machines*, (August 2014): 5

DOI: 10.1023/B:MIND.0000035461.63578.9d

συστήματα. Αντιθέτως, ο άνθρωπος, κινητοποιημένος από τον Ηθικό Νόμο<sup>49</sup> (περιεχόμενο του Κ3) αποφασίζει και πράττει σύμφωνα με την ελεύθερή του βούληση. Εξαιτίας αυτής της σημαντικής διαφοράς ανάμεσα σε έναν άνθρωπο και σε μία μηχανή ΤΝ, εγείρεται το σημαντικό ερώτημα εάν θα πρέπει να υπάρχει μία διαφορετικού τύπου μεταχείριση των συστημάτων αυτών, από αυτή που υφίσταται ήδη για τους ανθρώπους όσον αφορά το ζήτημα του ηθικού καθεστώτος<sup>50</sup>. Ποια είναι η ηθική επιταγή σύμφωνα με την οποία αποφασίζει ένα σύστημα ΤΝ να επιτελέσει μία του πράξη; Προς το παρόν, μας είναι άγνωστο.

Το πρόβλημα λοιπόν παραμένει. Μπορεί η ανθρωπότητα να διακρίνει εάν ένα τέτοιο σύστημα ΤΝ αντλεί τις αρχές της ηθικής από τον Κ3 προκειμένου να προβεί σε μία πράξη; Εάν ναι, μήπως έτσι κατορθώνει να συνδέεται στον Κ2 που έως τώρα θεωρούσαμε πως όχι; Εάν όχι, τότε τι είναι αυτό που ισχύει για την ηθική και τα συστήματα ΤΝ;

### *Συμπεράσματα*

Η μεταηθική προσέγγιση αναδεικνύει ερωτήματα και προβλήματα που ενδεχομένως προκύπτουν από τη χρήση συστημάτων τεχνητής νοημοσύνης που ως τώρα μπορεί να μην είχαν γίνει αντιληπτά. Ο μετασχηματισμός της οντολογικής προσέγγισης του Popper δημιουργεί ένα γόνιμο έδαφος που μας οδηγεί να υποθέσουμε πως ο Κ3 μπορεί να δημιουργείται ή προσεγγίζεται και από άλλες οντότητες. Είναι γεγονός πως εντοπίζεται μία τάση ή και επιθυμία εξίσωσης της ανθρώπινης φύσης με τις λειτουργίες και πρακτικές μίας αυτόνομης μηχανής σε ηθικό επίπεδο, στις περιπτώσεις που αυτό κρίνεται αναγκαίο ώστε να μπορούν να αποδοθούν ευθύνες ειδικά στην περίπτωση λάθος πρακτικών.

Η Ποππεριανή προσέγγιση του ζητήματος όμως αναδεικνύει μία ειδοποιό διαφορά ανάμεσα στα συστήματα τεχνητής νοημοσύνης και στις ανθρώπινες οντότητες. Στα συστήματα Τ.Ν. η μετάβαση πραγματοποιείται από τον Κ3 που διαθέτουν πρόσβαση, στη λειτουργία και πραγμάτωση αυτού στο Κ1 χωρίς τον ενδιάμεσο κόσμο Κ2. Αντιθέτως, ο άνθρωπος διαθέτει πρόσβαση και στον Κ2 στον οποίο δημιουργεί Κ3 και αποβλέπει στον Κ1. Η μετάβαση των συστημάτων ΤΝ στον Κ2 από τον Κ3 και

<sup>49</sup> Ιμμάνουελ Καντ, *Θεμελίωση της Μεταφυσικής των Ηθών*, επιμ. Κώστας Ανδρουλιδάκης, (Πανεπιστημιακές εκδόσεις Κρήτης: 2017), 75-79.

<sup>50</sup> Hage, "Theoretical foundations, 256.

κατόπιν στον Κ1 (από τον ένα δηλαδή κόσμο του Popper στον άλλο) είναι αβέβαιη και ενδεχομένως να περιορίζεται στην αποβλεπτικότητα τους μέσω των αντανάκλασεων από τον Κ1 στον Κ1 ή όπως εξηγήθηκε από τον Κ3 στον Κ1 εξαιτίας του τρόπου δημιουργίας τους. Το ερώτημα που προκύπτει λοιπόν είναι: Μπορεί να υπάρξει ηθική για τα συστήματα ΤΝ εάν απουσιάζει από αυτά ο Κ2; Εν αντιθέσει με τις ηθικές αλήθειες που αναγνωρίζουν οι άνθρωποι και οι οποίες *πηγάζουν* από τον Κ3 και αντανάκλουν στον Κ1, τα συστήματα Τεχνητής Νοημοσύνης είναι αβέβαιο αν μπορούν να επιτελέσουν μία τέτοιου είδους αναγωγή και αντ' αυτού είναι σαν να δημιουργούν ένα καινούργιο επίπεδο το οποίο ακόμη μας είναι άγνωστο.

Στα συστήματα τεχνητής νοημοσύνης θα μπορούσαμε να πούμε με όσα έχουν αναφερθεί στην παρούσα εργασία, πως ενδέχεται *μόνο* η δημιουργία τους να εντοπίζεται και να πηγάζει από τον Κ3 και να αντανάκλα στον Κ1 εάν όπως αναλύθηκε στο πρώτο κεφάλαιο οι επιστήμονες έχουν ως αντικείμενο έρευνας τα όσα βρίσκονται στον Κ3 και όχι στον Κ2 όπως φαίνεται να κάνουν έως σήμερα. Σκοπός όμως αυτής της οντολογικής προσέγγισης του ζητήματος, είναι να εντοπιστεί η ίδια συλλογιστική αλληλουχία και στις ηθικές συνιστώσες τέτοιων συστημάτων. Επομένως ανοίγει ένα νέο πεδίο έρευνας το οποίο εξετάζει ένα νέο κόσμο, τον νοητικό κόσμο των συστημάτων ΤΝ, στον οποίο απουσιάζει ο Κ2 και που ενδεχομένως λειτουργεί με το δικό του αξιακό σύστημα, ενέχει αποβλεπτικότητας στον Κ1, αλλά αποτελεί μόνο κατά ένα μέρος του αντικείμενο του φυσικού κόσμου. Τα ερωτήματα παραμένουν, και μαζί με αυτά εγείρονται καινούργια, καθώς και καινούργιες προβληματικές που ενδέχεται να προκύψουν εξαιτίας του αστάθμητου παράγοντα του χρόνου.

Η μεταηθική θεμελίωση ανοίγει καινούργια μονοπάτια σκέψης στον τρόπο ηθικής θεμελίωσης των συστημάτων τεχνητής νοημοσύνης. Μήπως θα ήταν καλύτερο να δομήσουμε ένα εξολοκλήρου νέο ηθικό σύστημα από εντελώς διαφορετική σκοπιά, για τις ηθικές συνιστώσες και πρακτικές (που επιτελούν ή θα επιτελέσουν στο μέλλον) αυτών των συστημάτων ή να ερευνήσουμε περισσότερο τα ερωτήματα «είναι οι οντότητες ΤΝ ισότιμες με τον άνθρωπο;» ή αν έστω «είναι οι οντότητες της τεχνητής νοημοσύνης έλλογα όντα;» ώστε να καταλήξουμε στο «εάν είναι δυνατόν οι οντότητες της ΤΝ να παράγουν έναν διαφορετικό τον Κ2 που οδηγεί και στον Κ3»; Το σίγουρο είναι, πως δεν μπορούμε να πούμε με βεβαιότητα ότι προσφέρεται μία οριστική και άμεση λύση στο ζήτημα.

## Βιβλιογραφία

### Ελληνική

- Βιρβιδάκης, Στέλιος. «Ο Κόσμος 3 του Καρλ Πόππερ ως ερμηνευτική πρόταση για την κατανόηση των μοντέλων μετριοπαθούς αξιακού ρεαλισμού», *Φιλοσοφία* 49, νο. II (2020): 194-207.
- Γούναρης Α., & Κωστελέτος Γ. (2024). Γράφοντας τον αλγόριθμο του Καλού: Η Τεχνητή Νοημοσύνη ως μηχανή απόδοσης Δικαιοσύνης. Ηθική. Περιοδικό φιλοσοφίας, (19).  
<https://doi.org/10.12681/ethiki.39654>
- Καντ, Ιμμάνουελ. *Θεμελίωση της Μεταφυσικής των Ηθών*, επιμ. Κώστας Ανδρουλιδάκης, Πανεπιστημιακές εκδόσεις Κρήτης: 2017.
- Nagel, Thomas. *Η Θέα από το πουθενά*, Αθήνα: Κριτική, 2002.

### Ξενόγλωσση

- Arkin, Ronald C. “The Case of Ethical Autonomy in Unmanned Systems”, *Journal of Military Ethics Vol. 9*, no. 4, 2010.
- Asaro, Peter. “THE NEAR FUTURE OF ARTIFICIAL INTELLIGENCE”, *Autonomous Weapons and the Ethics of A.I.*, 2020.
- Floridi, Luciano. “On the Morality of Artificial Intelligence”, *Minds and Machines*, August: 2014.  
DOI: 10.1023/B:MIND.0000035461.63578.9d
- Gounaris, A. (2011). Intentionality and the Emergence of Meaning. *Philosophia - Annual Journal of the Research Centre for Greek Philosophy of the Academy of Athens*, v.41, pp 319-321, 2011. ISSN 1105-2120
- Gounaris, A., Kosteletos, G. (2020). Licensed to Kill: Autonomous Weapons as Persons and Moral Agents. In Prole, D. and Rujević, G. (ed.). *Personhood*. Novi Sad, Filozofski Fakultet & The NKUA Applied Philosophy Research Lab Press. DOI: <https://doi.org/10.12681/aprpl.49>
- Hage, Jaap. “Theoretical foundations for the responsibility of autonomous agents”, *Springer*: August 2017.
- Kim, Jaegwon, *Η Φιλοσοφία του Νου*, Leader Books, Αθήνα: 2005.
- Popper, Karl R. *Objective knowledge: an evolutionary approach*, New York: Oxford University Press, 1972.
- Robinson, Howard. “Dualism”, *The Stanford Encyclopedia of Philosophy* (Fall 2020 Edition), επιμέλεια: Edward N. Zalta.

- <https://plato.stanford.edu/archives/fall2020/entries/dualism/>  
Sparrow, Robert. “Killer Robots”, *Journal of Applied Philosophy* Vol. 24, no. 1, 2007.
- Thornton, Stephen. “Karl Popper”, *The Stanford Encyclopedia of Philosophy* (Winter 2022 Edition), επιμέλεια: Edward N. Zalta & Uri Nodelman.  
<https://plato.stanford.edu/archives/win2022/entries/popper/>

### Διαδικτυακοί ιστότοποι

- DeepMind. AlphaGo - The Movie | Full award-winning documentary. Youtube, 2020. <https://www.youtube.com/watch?v=WXuK6gekUIY>
- Bronstein, Michael. AMMI Course "Geometric Deep Learning" - Lecture 1 (Introduction). Youtube, 2021.  
[https://www.youtube.com/watch?v=PtA0lg\\_e5nA](https://www.youtube.com/watch?v=PtA0lg_e5nA)



### Περίληψη

Η ραγδαία ανάπτυξη των συστημάτων Τεχνητής Νοημοσύνης (TN ή A.I.-Artificial Intelligence-) αποτελεί ένα βασικό παράγοντα βελτίωσης της ποιότητας ζωής των ανθρώπων και της κοινωνίας εν γένει σε όλους τους τομείς. Από τα συστήματα παραγωγικής Τεχνητής Νοημοσύνης έως τα αυτοματοποιημένα μηχανήματα έκδοσης εισιτηρίων και τα αυτοματοποιημένα οχήματα, ο τρόπος διαβίωσης των ανθρώπινων όντων αναδιαμορφώνεται. Σαφέστατα, κρίνεται πως είναι πολύ νωρίς ώστε να προκύψουν συνολικά συμπεράσματα για το αν αυτή η εισβολή της τεχνολογίας αποτελεί κάτι κοινωνικά ωφέλιμο ή όχι. Η σκέψη αυτή προβληματίζει την επιστημονική κοινότητα όχι τόσο για τις περιπτώσεις ημιαυτόνομης λειτουργίας τέτοιων συστημάτων<sup>51</sup>, όσο για τις περιπτώσεις πλήρους αυτονομίας και «αυτενεργείας» τους<sup>52</sup>. Σε αυτές ο άνθρωπος έχει ελάχιστη ή και μηδαμινή συμμετοχή. Καραδοκεί θα μπορούσαμε να πούμε, ένας μόνιμος φόβος ή

---

<sup>51</sup> Human in the loop

<sup>52</sup> Human on the loop ή human off the loop

αλλιώς αμφιβολία, για το αν τέτοια συστήματα μπορούν να δράσουν αυτόνομα με τέτοιο τρόπο ώστε οι ενέργειες αλλά και τα κίνητρά τους να ομοιάζουν με αυτές των ανθρώπων, εάν δηλαδή διαθέτουν ηθικό καθεστώς. Επομένως το βασικό ερώτημα που προκύπτει είναι εάν τέτοιου είδους συστήματα έχουν ή μπορούν να αποκτήσουν ηθικό καθεστώς ή και κατ' επέκταση εάν μπορούν να θεωρηθούν ηθικά υποκείμενα. Ποια είναι τα ηθικά θεμέλια πάνω στα οποία αναπτύσσονται και λειτουργούν τα συστήματα Τεχνητής Νοημοσύνης; Υπάρχουν; Είναι αναγκαίο να βασίζονται σε τέτοιου είδους θεμέλια; Ενδεχομένως τα παραπάνω ερωτήματα να βρουν απαντήσεις στη μεταηθική προσέγγιση του ζητήματος υπό το πρίσμα της οντολογικής θεωρίας των τριών κόσμων του Karl R. Popper.

*Λέξεις-κλειδιά:* Τεχνητή Νοημοσύνη, απόδοση ηθικής ευθύνης, ηθική ευθύνη, ηθικό καθεστώς, Οντολογία, μεταηθική, οι Κόσμοι του Karl R. Popper, Αποβλεπτικότητα

---

Ειρήνη Δαρκαδάκη  
Τμήμα Φιλοσοφίας, ΕΚΠΑ  
Ηλεκτρονική Διεύθυνση: [iomalan@philosophy.uoa.gr](mailto:iomalan@philosophy.uoa.gr)  
ORCID iD: <https://orcid.org/0000-0001-9942-7124>