

ΤΕΧΝΗΤΗ ΓΕΝΙΚΗ ΝΟΗΜΟΣΥΝΗ Η ΝΕΑ ΠΡΟΚΛΗΣΗ ΤΗΣ ΑΝΘΡΩΠΟΤΗΤΑΣ



Σεβαστοὶ πατέρες, κυρίες καὶ κύριοι,

Εὐχαριστῶ τὸ ΔΣ τῆς Χριστιανικῆς Ἐνώσεως Ἐπιστημόνων, γιατί με τὴν πρόσκλησή του μοῦ δίνεται ἡ εὐκαιρία νὰ διατυπώσω τὸν προβληματισμὸ μου πάνω σὲ ἓνα φλέγον ζήτημα τῆς ἐποχῆς μας τὸ ὁποῖο μᾶς ἀπασχολεῖ ὅλους, τὸ θέμα τῆς Τεχνητῆς Νοημοσύνης. Θὰ μιλήσω ὡς ἐπιστήμονας, ἀλλὰ καὶ ὡς ὀρθόδοξος χριστιανός.

Πρὶν ἀπὸ ἓναν χρόνο, ἡ ἐταιρεία λογισμικοῦ OpenAI μὲ διευθυντὴ τὸν Sam Altman παρουσίασε τὸ πρόγραμμα συνομιλίας ChatGPT. Μέσα σὲ δύο μῆνες, περισσότεροι ἀπὸ 100 ἑκατομμύρια ἄνθρωποι ἄρχισαν νὰ συνομιλοῦν μὲ τὸν ὑπολογιστὴ τους. Κάτι τέτοιο δὲν εἶχε ξαναγίνει ποτὲ μὲ καμμία ἀνθρώπινη ἐφεύρεση (γιὰ νὰ φτάσουν τοὺς 100 ἑκατομμύρια νέους χρῆστες, χρειάστηκαν

9 μῆνες γιὰ τὸ TikTok, 4 χρόνια γιὰ τὸ YouTube, 7 χρόνια γιὰ τὸ Ἴντερνετ, καὶ πολλὲς δεκαετίες γιὰ τὸν ἠλεκτρισμὸ, τὸ τηλέφωνο, τὸ αὐτοκίνητο κλπ.). Ἡ ἐταιρεία αὐτὴ ἔλαβε στὴ συνέχεια ἐπιπλέον χρηματοδότηση ἀπὸ τὴ Microsoft, ὅλες οἱ μεγάλες ἐταιρεῖες λογισμικοῦ ῥίχτηκαν στὴν κούρσα τῆς Τεχνητῆς Γενικῆς Νοημοσύνης, καὶ ὅλοι μιλοῦν γιὰ τὸν κίνδυνο ὑποδούλωσης τοῦ ἀνθρώπου ἀπὸ τὶς μηχανές, ὅπως βλέπουμε στὶς ταινίες τοῦ Hollywood (Terminator). Ἄς πάρουμε ὁμῶς τὰ πράγματα λίγο ἀπὸ τὴν ἀρχή.

Τεχνητὴ Νοημοσύνη εἶναι ἡ ἰκανότητα ποὺ ἔχει ἓνας ἠλεκτρονικὸς ὑπολογιστὴς ἢ ἓνα ρομπὸτ ποὺ ἐλέγχεται ἀπὸ ἓναν ἠλεκτρονικὸ ὑπολογιστὴ νὰ πραγματοποιεῖ συνήθεις διεργασίες ποὺ ἐκτελοῦν νοήμονα ὄντα. Ὅταν ἀναφερόμαστε στὸν ἄνθρωπο, τέτοιες διεργασίες εἶναι ἡ δυνατότητα νὰ μαθαίνει, νὰ βγάζει συμπεράσματα, νὰ λύνει προβλήματα, νὰ βρῆσει

νόημα, να γενικεύει, να χρησιμοποιεί μιὰ γλώσσα, να ἔχει ἀντίληψη τοῦ περιβάλλοντός του κ.ἄ. Ἀπὸ τὴν ἀνάπτυξη τοῦ ψηφιακοῦ ὑπολογιστῆ τὴν δεκαετία τοῦ 1940 μέχρι σήμερα, οἱ ὑπολογιστὲς ἔχουν προγραμματιστεῖ νὰ ἀποδεικνύουν μαθηματικὰ θεωρήματα, νὰ παίζουν σκάκι, νὰ κάνουν ἰατρικὲς διαγνώσεις, νὰ ἀναγνωρίζουν ὀμιλία καὶ γραπτὸ κείμενο, νὰ συνδιαλέγονται μὲ τὸν ἄνθρωπο σὲ καθημερινὴ ὀμιλία, καὶ ὅπως θὰ δοῦμε στὴ συνέχεια, παρουσιάζουν ἤδη στοιχεῖα Γενικῆς Νοημοσύνης. Μὲ δεδομένο τὸν

τανοεῖ» τὴν πληροφορία αὐτή, δηλαδὴ ἢ πληροφορία νὰ ἀποκτᾶ «νόημα». Ἡ μέθοδος αὐτὴ ξεκίνησε νὰ ἀναπτύσσεται τὴ δεκαετία τοῦ '50 καὶ προσπαθεῖ νὰ προσομοιώσει στὸν ὑπολογιστὴ τὸν τρόπο μὲ τὸν ὁποῖο σκέφτεται ὁ ἄνθρωπος, μὲ πεπερασμένη μέχρι σήμερα ἐπιτυχία. Παράδειγμα τέτοιου λογισμικοῦ, ποῦ δίνει ἀπαντήσεις σὲ τεχνικὰ ἐρωτήματα, εἶναι τὸ WolframAlpha ἀπὸ τὴν εταιρεία ποῦ ἀνέπτυξε τὸ γνωστὸ πρόγραμμα Mathematica.

Ἡ δεύτερη μέθοδος ὀνομάζεται συνδε-

Τεχνητὴ Νοημοσύνη εἶναι ἡ ἰκανότητα ποῦ ἔχει ἓνας ἠλεκτρονικὸς ὑπολογιστὴς ἢ ἓνα ρομπὸτ ποῦ ἐλέγχεται ἀπὸ ἓναν ἠλεκτρονικὸ ὑπολογιστὴ νὰ πραγματοποιεῖ συνήθεις διεργασίες ποῦ ἐκτελοῦν νοήμονα ὄντα. Ὅταν ἀναφερόμαστε στὸν ἄνθρωπο, τέτοιες διεργασίες εἶναι ἡ δυνατότητα νὰ μαθαίνει, νὰ βγάζει συμπεράσματα, νὰ λύνει προβλήματα, νὰ βρίσκει νόημα, νὰ γενικεύει, νὰ χρησιμοποιεῖ μιὰ γλώσσα, νὰ ἔχει ἀντίληψη τοῦ περιβάλλοντός του κ.ἄ.

περίφημο νόμο τοῦ Moore, δηλαδὴ τὸν διπλασιασμὸ τῆς ἰσχύος τῶν ὑπολογιστικῶν συστημάτων κάθε 18 μῆνες, δηλαδὴ τὸν χιλιαπλασιασμὸ τῆς ἰσχύος τους κάθε 15 χρόνια, εἶναι λογικὸ νὰ μᾶς ἀπασχολεῖ ἢ προοπτικὴ μιᾶς ὑπερ-νοημοσύνης (super-intelligence) ποῦ θὰ ξεπεράσει κατὰ πολὺ τις διανοητικὲς δυνατότητες τοῦ ἀνθρώπου.

Ἐπάρχουν δύο μέθοδοι ὑλοποίησης Τεχνητῆς Νοημοσύνης στὸν ὑπολογιστῆ. Ἡ πρώτη, ἡ συμβολικὴ (symbolic ἢ “top-down”), σχετίζεται μὲ τὸν τρόπο ποῦ ἀποκτᾶ νόημα ἢ γνώση ποῦ συσσωρεύει ὁ ἄνθρωπος. Ἡ πληροφορία ἀναλύεται σὲ σύμβολα ποῦ ἀντιπροσωπεύουν ἔννοιες, καὶ ὁ προγραμματιστὴς καθορίζει τις σχέσεις ἀνάμεσα στὰ σύμβολα αὐτά, ὥστε ὁ ὑπολογιστὴς νὰ μπορεῖ νὰ «κα-

σμολογικὴ (connectionist ἢ “bottom-up”). Ξεκίνησε τὸ 1957 μὲ τὴν εἰσαγωγὴ τῶν perceptrons καὶ τὴν μεθοδολογία ἐκπαίδευσίς τους ἀπὸ τὸν Frank Rosenblatt τοῦ Πανεπιστημίου Cornell, καὶ ἀπὸ τὴ δεκαετία τοῦ '80 μέχρι σήμερα ἀναπτύσσεται μὲ ραγδαίους ρυθμούς. Ἡ μέθοδος αὐτὴ προσομοιώνει τὸν τρόπο μὲ τὸν ὁποῖο λειτουργεῖ ὁ ἀνθρώπινος ἐγκέφαλος, ὁ ὁποῖος τελικὰ, μὲ τρόπο ἄγνωστο καὶ ἀκατανόητο μέχρι στιγμῆς, προσδίδει στὸν ἄνθρωπο τὰ χαρακτηριστικὰ τῆς νοημοσύνης. Δηλαδή, ἀφοῦ δὲν γνωρίζουμε τί εἶναι ἡ ἀνθρώπινη νοημοσύνη ἀλλὰ οὔτε καὶ μποροῦμε νὰ προσομοιώσουμε μὲ ἐπιτυχία τὸν μηχανισμό της, ἀποφασίσαμε νὰ προσομοιώσουμε τὴ λειτουργία τῶν θεμελιωδῶν συστατικῶν τοῦ ἐγκεφάλου, τῶν νευρῶνων καὶ τῶν συνάψεων μεταξὺ

τους, να αφήσουμε στη συνέχεια τον υπολογιστή να εκπαιδευτεί μόνος του από μια βάση δεδομένων χωρίς την δική μας επίβλεψη, και να δούμε στη συνέχεια τί θα κάνει. Κάτι δηλαδή σαν τον τρόπο με τον οποίο μαθαίνει τον κόσμο ένα μικρό παιδί. Ής θυμηθούμε τη φυσιολογία του ανθρώπινου εγκεφάλου. Ο εγκέφαλος δέχεται ερεθίσματα από τα διάφορα όργανα του σώματος. Τα ερεθίσματα αυτά μεταφέρονται με τη μορφή ηλεκτρικού σήματος από τους νευρώνες. Οί νευρώνες είτε ενισχύουν είτε αποδυναμώνουν το ηλεκτρικό σήμα που μεταφέρουν, και καταλήγουν σε συνάψεις με άλλους νευρώνες, όπου προστίθενται όλα τα σήματα των νευρώνων που φτάνουν σε αυτούς και στη συνέχεια αποστέλλονται νέα σήματα προς τους επόμενους νευρώνες, κ.ο.κ. Με τον τρόπο αυτό ο άνθρωπος σκέφτεται, θυμάται, νιώθει, κλπ. Στον ανθρώπινο

προσομοιωθεί πολύ εύκολα σε έναν ηλεκτρονικό υπολογιστή. Ορίζεται ή δομεί ένας δικτύος μεγάλου αριθμού νευρώνων (ή perceptrons) των οποίων ή λειτουργία είναι να πολλαπλασιάζουν το σήμα της εισόδου τους με έναν αριθμό (βάρος, weight). Ο ρόλος των συνδέσμων είναι να προσθέτουν όλα αυτά τα σήματα από τους νευρώνες που καταλήγουν σε αυτούς, και στη συνέχεια να κανονικοποιούν το αποτέλεσμα σε μια τιμή μεταξύ 0 και 1 μέσω μιας μη γραμμικής συνάρτησης «ενεργοποίησης» (activation function). Η μη γραμμικότητα είναι σημαντική, γιατί αλλιώς το σύστημα θα ήταν απλώς ένας γραμμικός πίνακας χωρίς ιδιαίτερες δυνατότητες. Το δίκτυο εκπαιδεύεται, είτε αυτόνομα είτε υπό ανθρώπινη επίβλεψη, συγκρίνοντας την έξοδο του με μια πρότυπη έξοδο. Οί παράμετροι του δικτύου αναπροσαρμόζονται, ώστε

Το 1950 ο Άγγλος Μαθηματικός Alan Turing αναρωτήθηκε αν μπορεί μια μηχανή να σκέπτεται. Και επειδή είναι δύσκολο να ορίσουμε τί σημαίνει «σκέπτομαι», πρότεινε να θέσουμε ερωτήσεις στη μηχανή και σε έναν άνθρωπο, χωρίς να ξέρουμε ποιός είναι τί, και αν από τις απαντήσεις που θα λάβουμε δεν μπορούσαμε να πούμε με βεβαιότητα ποιός είναι ή μηχανή και ποιός ο άνθρωπος, τότε θα μπορούμε να πούμε ότι ή μηχανή σκέπτεται, ή ότι ή μηχανή έφτασε τον άνθρωπο.

εγκέφαλο υπάρχουν περίπου 100 δισεκατομμύρια νευρώνες, ο καθένας με 7.000 συνάψεις (συνδέσεις) με άλλους νευρώνες. Υπολογίζεται δηλαδή ότι ο εγκέφαλος ενός παιδιού 3 ετών έχει περίπου 1 τετρακίς εκατομμύριο διασυνδέσεις, και ο εγκέφαλος ενός ενήλικου περίπου 100 με 500 τρισεκατομμύρια. Συγκρατήστε παρακαλώ αυτό το νούμερο. Έκατοντάδες τρισεκατομμύρια διασυνδέσεις!

Η παραπάνω λειτουργία μπορεί να

να ελαχιστοποιηθεί ή απόκλιση της τιμής της εξόδου από το πρότυπο. Για παράδειγμα, τα σύγχρονα Μεγάλα Γλωσσικά Μοντέλα (Large Language Models-LLM) περιέχουν μερικές εκατοντάδες δισεκατομμύρια συνδέσεις, δηλαδή μερικές εκατοντάδες δισεκατομμύρια παραμέτρους (βάρη, weights). Συγκρατήστε παρακαλώ και αυτό το νούμερο. Έκατοντάδες δισεκατομμύρια διασυνδέσεις!

Με τα παραπάνω περιγράψαμε με

πολύ απλά λόγια πώς λειτουργεί ένα Νευρωνικό Δίκτυο (Neural Network). Δεν είπαμε τίποτα για τη δομή του, πού μπορεί να αποτελείται από διαδοχικά στρώματα κόμβων (Feed-Forward), να μεταφέρουν πληροφορία σε μεθεπόμενα στρώματα (Residual), να εκτελούν κάποια προεργασία στο σήμα εισόδου (Convolution), να λειτουργούν με ανάδραση (Recurrent), να κάνουν κωδικοποίηση και αποκωδικοποίηση της πληροφορίας (Encoder-Decoder), να μεταφέρουν πληροφορία (attention) σε απομακρυσμένα σημεία του (Long-Short Term Memory, Transformer) κλπ. Όσο πολύπλοκη όμως και να είναι η αρχιτεκτονική αυτών των δικτύων, είναι απείρως απλούστερη από αυτήν του ανθρώπινου εγκεφάλου. Επιπλέον, μπορεί ένα Νευρωνικό Δίκτυο να διαγράφει κόμβους, αν δεν τους χρησιμοποιεί, όπως και ο άνθρωπος εγκέφαλος. Όπως προαναφέραμε, τα συστήματα αυτά άρχισαν να αναπτύσσονται με πιδόντατικούς ρυθμούς τη δεκαετία του '80, και μέχρι τον Νοέμβριο του 2022 παρουσίαζαν συγκρατημένη επιτυχία στην επίτευξη Τεχνητής Γενικής Νοημοσύνης.

Ένας τομέας έρευνας των ηλεκτρονικών υπολογιστών είναι η Έπεξεργασία Φυσικής Γλώσσας (Natural Language Processing-NLP). Ένα από τα πρώτα προγράμματα συνομιλίας ήταν το γνωστό πρόγραμμα ψυχοθεραπείας Eliza, το οποίο ξεκίνησε να αναπτύσσεται στο MIT το 1966, και για ιστορικούς λόγους λειτουργεί μέχρι σήμερα. Το πρόγραμμα αυτό εκπαιδεύτηκε, ώστε να εντοπίζει συγκεκριμένες λέξεις και εκφράσεις στο κείμενο που εισάγει ο χρήστης, και με βάση αυτές να δίνει απαντήσεις από ένα δεδομένο σύνολο απαντήσεων όγκου 50 kilobytes. Ο χρήστης το αντιλαμβάνεται αυτό πολύ γρήγορα, γιατί οι απαντήσεις του έχουν μικρή σχέση με αυτά που γράφει ο ίδιος, δηλαδή αντιλαμβάνεται ότι

το πρόγραμμα αυτό «δεν ξέρει τί λέει». Και όμως, πολλοί ασθενείς που χρησιμοποίησαν το πρόγραμμα αυτό δήλωσαν ότι έμειναν πολύ ευχαριστημένοι από την Eliza, γιατί «αυτή μόνον τους κατάλαβε», και ζήτησαν να συνεχίσουν τη συνομιλία μαζί της. Αυτό το στοιχείο είναι πολύ σημαντικό, γιατί το 1950 ο Άγγλος Μαθηματικός Alan Turing αναρωτήθηκε αν μπορεί μια μηχανή να σκέπτεται. Και επειδή είναι δύσκολο να ορίσουμε τι σημαίνει «σκέπτομαι», πρότεινε να θέσουμε ερωτήσεις στη μηχανή και σε έναν άνθρωπο, χωρίς να ξέρουμε ποιός είναι τί, και αν από τις απαντήσεις που θα λάβουμε δεν μπορούμε να πούμε με βεβαιότητα ποιός είναι η μηχανή και ποιός ο άνθρωπος, τότε θα μπορούμε να πούμε ότι η μηχανή σκέπτεται, ή ότι η μηχανή έφτασε τον άνθρωπο. Όπως είδαμε, το πρόγραμμα ψυχοθεραπείας Eliza πέρασε αυτή τη δοκιμή, αλλά μόνον με ψυχικά άρρώστους. Συγκρατήστε παρακαλώ και αυτό το σημαντικό στοιχείο: ένας ψυχικά άρρωστος δεν μπορεί να ξεχωρίσει έναν άνθρωπο από ένα ρομπότ!

Μετά το Eliza, ακολούθησαν οι «εικονικοί βοηθοί» (virtual assistants) Siri, Cortana, Alexa κλπ., με δυνατότητα επικοινωνίας σε καθημιόλιμνη γλώσσα με τον χρήστη, και από το 2017 ξεκίνησαν να αναπτύσσονται τα σύγχρονα Μεγάλα Γλωσσικά Μοντέλα (Large Language Models-LLM). Αυτά είναι μεγάλα Νευρωνικά Δίκτυα προ-έκπαιδευμένα πάνω σε όλη την ανθρώπινη γνώση που είναι διαθέσιμη στο internet (έχουν δηλαδή «διαβάσει» όλα τα βιβλία του κόσμου!). Σε σχέση με το πρόγραμμα Eliza, αυτά εκπαιδεύονται πάνω σε ένα σύνολο κειμένων όγκου 45 τρισεκατομμύρια bytes. Λειτουργούν ως εξής: Έπεξεργάζονται το κείμενο που εισάγει ο χρήστης τους, και στη συνέχεια απαντούν με μια λέξη. Η λέξη αυτή προστίθεται στο προηγούμενο

κείμενο, τὸ νέο κείμενο ὑφίσταται νέα ἐπεξεργασία, καὶ τὸ πρόγραμμα ἀπαντᾷ στὴ συνέχεια μὲ μιὰ νέα λέξη, καὶ οὕτω καθεξῆς. Τὸ μόνο ποὺ κάνουν εἶναι νὰ προβλέπουν ποιά θὰ εἶναι ἡ ἐπόμενη λέξη στὴ συνομιλία μὲ τὸν χρήστη τους. Τίποτε περισσότερο, τίποτε λιγώτερο! Καὶ μὲ τὸν τρόπο αὐτὸ μποροῦν νὰ συνομιλοῦν πάνω σὲ ὅλα τὰ θέματα τῆς ἀνθρώπινης ἐμπειρίας καὶ γνώσης, σὲ πολλές γλώσσες, καὶ ὁ μέσος χρήστης τους δὲν μπορεῖ νὰ ξεχωρίσει ἀν' ὁ συ-

Ἔγὼ προσωπικὰ δὲν μποροῦσα νὰ πιστέψω ὅτι θὰ λειτουργοῦσαν ἀποτελεσματικά, ὅμως ἔκανα λάθος. Γιὰ κάποια χρόνια μαθαίναμε ἐλάχιστα πράγματα γιὰ τὴν ἐξέλιξή τους. Γνωρίζαμε ὅτι ὁ πρωτοπόρος στὸν χῶρο ἦταν ἡ ἐταιρεία Google, γιὰ κάποιον λόγο ὅμως δὲν ἀνέφερε σημαντικὴ πρόοδο. Ἡ προσωπικὴ μου ἐκτίμηση εἶναι ὅτι ἡ Google διαπίστωσε ὅτι τὸ δικό της πρόγραμμα (Bard) δὲν ἦταν ἀρκετὰ ἀξιόπιστο (γιὰ παράδειγμα, ἀνέφερε μὲ μεγάλη σιγουριά ὅτι

Τὰ Νευρωνικὰ Δίκτυα ἀποτελοῦν πρὸς τὸ παρὸν μαῦρο κουτὶ (black box), καὶ δὲν κατανοοῦμε ἀκόμα πλήρως τὴν λειτουργία τους. Αὐτὸ σημαίνει ὅμως ὅτι δὲν ὑπάρχει καμία διαβεβαίωση ὅτι μὲ τὸν πολλαπλασιασμὸ τῆς ὑπολογιστικῆς ἰσχύος θὰ πολλαπλασιαστεῖ ἀντίστοιχα καὶ ἡ «νοημοσύνη» τους. Μπορεῖ νὰ περάσουν δεκαετίες, καὶ νὰ μὴν ὑπάρξει νέα θεαματικὴ ἔκρηξη Τεχνητῆς Νοημοσύνης. Καὶ ἀν' δὲν πιστεύετε ὅτι αὐτὸ εἶναι ἓνα πιθανὸ ἐνδεχόμενον, θυμηθεῖτε τί γίνεται μὲ τὴν ἔρευνα γιὰ τὴν παραγωγὴ φτηνῆς ἐνέργειας ἀπὸ τὸ ὕδρογόνο μὲσῶ τοῦ μηχανισμοῦ τῆς θερμοπυρηνικῆς σύντηξης. Ὁ μηχανισμὸς εἶναι ἀπλὸς (συμβαίνει στὸν ἥλιο), ἡ μεθοδολογία εἶναι γνωστὴ ἐδῶ καὶ 7 δεκαετίες περίπου (διάταξη Tokamak), κάθε 10 χρόνια μᾶς λένε ὅτι σὲ 10 χρόνια θὰ ἔχουμε ἀφθονή φτηνὴ ἐνέργεια, ἀκόμα ὅμως δὲν ἔχουμε ἰδέα πότε θὰ λειτουργήσῃ ἓνας ἀντιδραστήρας σύντηξης.

νομηπτής του εἶναι ἄνθρωπος ἢ μηχανή. Σὲ πολλές περιπτώσεις, τὰ προγράμματα αὐτὰ περνοῦν ἀνετα τὸ test τοῦ Turing, παρουσιάζουν δηλαδὴ ἱκανὰ στοιχεῖα Γενικῆς Τεχνητῆς Νοημοσύνης! Ὀνομάστηκαν «Θεμελιώδη» (Foundation Models), δηλαδὴ μοντέλα τὰ ὁποῖα μποροῦν νὰ ἐνσωματωθοῦν σὲ ἄλλα προγράμματα καὶ πάνω στὰ ὁποῖα μπορεῖ νὰ ἀναπτυχθεῖ πολὺ πιὸ σύνθετο λογισμικό.

τὸ διαστημικὸ τηλεσκόπιο James Webb ἦταν αὐτὸ ποὺ φωτογράφησε τὸν πρῶτο ἐξωπλανήτη-πλανήτη ἐκτὸς τοῦ ἡλιακοῦ συστήματος, καὶ ἡ «γκάφα» αὐτὴ κόστισε στὴ Google πτώση τῆς ἀξίας τῆς μετοχῆς της κατὰ 100 δις δολάρια! Ἐπίσης, πολλές φορὲς ἔδινε ἀκατάλληλες ἀπαντήσεις, πρόσβαλλε τὸν νομηπτή, ἔκανε ρατσιστικὰ σχόλια, ἔδινε πληροφορίες γιὰ τὸ πῶς μπορεῖ κανεὶς νὰ φτιάξει

μιὰ βόμβα κλπ.), όπότε ή Google ήταν πολύ διστακτική στο να τó βγάλει στην αγορά. Μάλιστα, ό 75χρονος διευθυντής του προγράμματος Geoffrey Hinton παραιτήθηκε, προειδοποιώντας για τούς κινδύνους από τις έξελίξεις στον χώρο τής Τεχνητής Νοημοσύνης. Φαίνεται ότι τó αντίπαλο δέος, ή Microsoft, δέν είχε αντίστοιχους ένδοιασμούς, αλλά, για να αποφύγει την ταλαιπωρία που υπέστη ή Google, άνεθεσε στην εταιρεία OpenAI «να βγάλει αυτή τó φίδι από την τρύπα», με έμμεση χρηματοδότηση, χωρίς έμπλέκεται άμεσα ή ίδια.

Όλα αυτά μέχρι τόν Νοέμβριο του 2022, όπότε κυκλοφόρησε τó λογισμικό ChatGPT, τó όποιο βασίστηκε στην έκδοση 3.5 του λογισμικού GPT (Generative Pre-Trained Transformer – Άναγεννητικό Προ-Έκπαιδευμένο Νευρωνικό Δίκτυο τύπου Transformer): ξαφνικά σαν «κάτι να ξύπνησε», και τó λογισμικό άρχισε να παρουσιάζει νέα «αναδυόμενα» (emergent) χαρακτηριστικά νοημοσύνης, έξυπνάδας, αυτόσυνειδησίας κλπ. Αυτός ήταν και ό λόγος τής τεράστιας άποδοχής και τής ραγδαίας διάδοσής του. Τó έντυπωσιακό είναι ότι μέχρι την προηγούμενη έκδοσή του (GPT-3), ή πρόοδος του λογισμικού ήταν μάλλον συκρατημένη. Τó ενδιαφέρον όμως είναι ότι ένα σύγχρονο Μεγάλο Γλωσσικό Μοντέλο έμπεριέχει όλη τή γνώση τής ανθρωπότητας μέσα στις περίπου 100 δισεκατομμύρια παραμέτρους του Νευρωνικού Δικτύου του (τά βάρη-weights των νευρώνων του δηλαδή). Ένας ανθρωπίνος έγκέφαλος με 1000 φορές περισσότερες παραμέτρους «θυμάται» πολύ λιγώτερα πράγματα. Δηλαδή ένα σύγχρονο Νευρωνικό Δίκτυο είναι πολύ πιό αποτελεσματικό στη λειτουργία τής αποθήκευσης πληροφορίας από τόν ανθρωπίνo έγκέφαλο. Ό λόγος πιστεύω είναι ότι ό ανθρωπίνος έγκέφαλος έπιτελεί πολύ περισσότερες δραστηριότητες

άπό ένα λογισμικό Τεχνητής Νοημοσύνης. Ό κύριος ρόλος του είναι να συντονίσει τις πολλαπλές λειτουργίες του ανθρώπινου σώματος, να διασφαλίσει την βιωσιμότητά του (δηλαδή να λειτουργεί με ρυθμούς πιό χαλαρούς από τó 100% των δυνατοτήτων του) και δευτερευόντως να άσχολείται με άνωτερες διανοητικές λειτουργίες. Φυσικά, αν τó να ξύπνησει ένα τέτοιο Μεγάλο Γλωσσικό Μοντέλο ήταν τόσο άπλό, φανταστείτε τί θα μπορούσε να γίνει στο άμεσο μέλλον. Φανταστείτε τί θα μπορούσε να γίνει σε λίγα χρόνια, όταν άποκτήσει 1000 φορές περισσότερες έπιπλέον παραμέτρους! Γι' αυτό κάποιοι φοβήθηκαν και κάποιοι παραιτήθηκαν από την περαιτέρω ανάπτυξη τέτοιων προγραμμάτων.

Η πραγματικότητα όμως είναι πιό πεζή. Τó σύστημα «ξύπνησε», όχι γιατί μεγάλωσε άπτότομα ή ύπολογιστική δύναμη (τó λεγόμενο compute). Όπως φαίνεται, έγιναν πολλές διορθώσεις σε πολλά σημεία, με πιό σημαντική, όπως φαίνεται, την άλλαγή τής συνάρτησης ενεργοποίησης σε τύπου ReLU (Rectified Linear Unit). Κανείς δέν γνωρίζει γιατί αυτό ήταν τόσο σημαντικό, και όμως άλλαξε θεαματικά τή συμπεριφορά του συστήματος. Τα Νευρωνικά Δίκτυα άποτελούν προς τó παρόν μαύρο κουτί (black box), και δέν κατανοούμε ακόμα πλήρως την λειτουργία τους. Αυτό σημαίνει όμως ότι δέν υπάρχει καμία διαβεβαίωση ότι με τόν πολλαπλασιασμό τής ύπολογιστικής ισχύος θα πολλαπλασιαστεί αντίστοιχα και ή «νοημοσύνη» τους. Μπορεί να περάσουν δεκαετίες, και να μñν υπάρξει νέα θεαματική έκρηξη Τεχνητής Νοημοσύνης. Και αν δέν πιστεύετε ότι αυτό είναι ένα πιθανό ένδεχόμενο, θυμηθείτε τί γίνεται με την έρευνα για την παραγωγή φτηνής ενέργειας από τó ύδρογόνο μέσω του μηχανισμού της θερμοπυρηνικής σύντηξης. Ό μηχανισμός είναι άπλός (συμβαίνει

στον ήλιο), ή μεθοδολογία είναι γνωστή εδώ και 7 δεκαετίες περίπου (διάταξη Tokamak), κάθε 10 χρόνια μᾶς λένε ότι σὲ 10 χρόνια θὰ ἔχουμε ἀφθονή φθινηνὴ ἐνέργεια, ἀκόμα ὅμως δὲν ἔχουμε ἰδέα πότε θὰ λειτουργήσει ἕνας ἀντιδραστήρας σύντηξης. Ἐμεῖς προσωπικὰ ἀκούσαμε τὸν Ἅγιο τῶν ἡμερῶν μας, τὸν Ἅγιο Πορφύριο, νὰ μᾶς λέει ὅτι «μοῦ φαίνεται ὅτι ἡ ἀπάντηση γιὰ τὸ ἐνεργειακὸ πρόβλημα εἶναι μπροστὰ στὰ μάτια τῶν ἐπιστημόνων καὶ ὅτι, ἂν ἀπλώσουν ἔτσι τὸ χέρι τους, θὰ βροῦν τὴν λύση, ὅμως εἶναι σὰν ὁ Θεὸς νὰ ἔχει ἀπλώσει ἕνα πέπλο μπροστὰ στὰ μάτια τους, καὶ δὲν μποροῦν νὰ βροῦν τὴ λύση, γιὰτὶ δὲν ξέρω ἂν θὰ εἶναι γιὰ τὸ καλὸ τῆς ἀνθρωπότητας...». Ἴσως κάτι τέτοιο νὰ συμβεῖ καὶ μὲ τὴν Τεχνητὴ Νοημοσύνη. Ἴσως νὰ μὴν ἐπιτρέψει ὁ Θεὸς νὰ φτάσουμε ποτὲ στὴν ὑπερ-νοημοσύνη, γιὰτὶ ἴσως αὐτὸ νὰ μὴν εἶναι γιὰ τὸ καλὸ τῆς ἀνθρωπότητας...

Εἶναι ἐνδιαφέρον ὅτι δὲν ὑπάρχει σήμερα ἕνας καλὸς ἐπιστημονικὸς ὀρισμὸς τῆς νοημοσύνης, οὔτε ἀκόμα καὶ γιὰ τὴν περίπτωση τῶν ζώων. Ἀντιλαμβανόμαστε ὅτι τὰ δελφίνια, γιὰ παράδειγμα, ἔχουν μιὰ κάποια μορφὴ νοημοσύνης, ὅμως δὲν γνωρίζουμε τί εἶναι αὐτὸ ποὺ πρέπει νὰ ἐπιτύχει ἡ Τεχνητὴ Νοημοσύνη, ὥστε νὰ μποροῦμε νὰ ποῦμε ἀντικειμενικὰ ὅτι ἔφτασε τὸ ἐπίπεδο ἑνὸς δελφινιοῦ. Δὲν μποροῦμε δηλαδὴ ἀντικειμενικὰ νὰ ποῦμε ἂν ἡ Τεχνητὴ Νοημοσύνη ἔχει πετύχει ἢ ὄχι τὸν στόχο της. Καὶ μὲ κάθε νέα ἐπιτυχία της (παιχνίδι σκακιοῦ, ἰατρικὴ διάγνωση, ὁδήγηση αὐτοκινήτου, καθημερινὴ ὁμιλία) λέμε, «καλὰ ὡς ἐδῶ, ὅμως ἡ ἀνθρώπινη νοημοσύνη εἶναι κάτι παραπέρα»... Θυμᾶμαι τὴν κατασκίνωση τῆς ΧΦΕ στὸν Παρνασσὸ τὸ 1983, πρὶν ἀπὸ 40 χρόνια, ὅταν μὲ κατέβασε στὴν Ἀφρική ἕνα ἀπὸ τὰ στελέχη. Ἐκεῖ, περιμένοντας τὸ τρένο, εἶχαμε μιὰ ἔντονη συζήτηση γιὰ τὴν Τεχνητὴ Νοημοσύνη, καὶ

τὸ κύριο ἐπιχείρημά του ὅτι ποτὲ δὲν θὰ ἐπιτευχθεῖ ἦταν ὅτι ὁ ἠλεκτρονικὸς ὑπολογιστὴς δὲν ἐπρόκειτο ποτὲ νὰ νικήσει

Εἶναι ἐνδιαφέρον ὅτι δὲν ὑπάρχει σήμερα ἕνας καλὸς ἐπιστημονικὸς ὀρισμὸς τῆς νοημοσύνης, οὔτε ἀκόμα καὶ γιὰ τὴν περίπτωση τῶν ζώων. Ἀντιλαμβανόμαστε ὅτι τὰ δελφίνια, γιὰ παράδειγμα, ἔχουν μιὰ κάποια μορφὴ νοημοσύνης, ὅμως δὲν γνωρίζουμε τί εἶναι αὐτὸ ποὺ πρέπει νὰ ἐπιτύχει ἡ Τεχνητὴ Νοημοσύνη, ὥστε νὰ μποροῦμε νὰ ποῦμε ἀντικειμενικὰ ὅτι ἔφτασε τὸ ἐπίπεδο ἑνὸς δελφινιοῦ. Δὲν μποροῦμε δηλαδὴ ἀντικειμενικὰ νὰ ποῦμε ἂν ἡ Τεχνητὴ Νοημοσύνη ἔχει πετύχει ἢ ὄχι τὸν στόχο της. Καὶ μὲ κάθε νέα ἐπιτυχία της (παιχνίδι σκακιοῦ, ἰατρικὴ διάγνωση, ὁδήγηση αὐτοκινήτου, καθημερινὴ ὁμιλία) λέμε, «καλὰ ὡς ἐδῶ, ὅμως ἡ ἀνθρώπινη νοημοσύνη εἶναι κάτι παραπέρα»...

τὸν παγκόσμιο πρωταθλητὴ στὸ σκάκι... 14 χρόνια μετὰ ὁ ὑπερ-ὑπολογιστὴς Deep Blue τῆς IBM νίκησε τὸν τότε παγκόσμιο πρωταθλητὴ Garry Kasparov. Ὅμως ἡ ἱστορία δὲν σταμάτησε ἐκεῖ. Σήμερα, ὁ

γιός μου, που παίζει έρασιτεχνικά σκάκι, έχει παγκόσμια βαθμολογία ELO 1100. Ο σημερινός παγκόσμιος πρωταθλητής έχει ELO 2800. Η εφαρμογή που έχει εγκαταστήσει ο γιός μου στο κινητό του έχει ELO 3100, και το πρόγραμμα σκακιού Stockfish έχει ELO πάνω από 3500. Συντρίβουν δηλαδή άνετα τον παγκόσμιο πρωταθλητή! Όπως θα δούμε στο τέλος της εισήγησης, το θέμα του τί είναι νομοσύνη και γενικότερα του τί είναι ο

ψηφιακοί βοηθοί στην καθημερινή μας ζωή (Siri, Alexa, Cortana), νέες έφευρέσεις (αυτόνομη οδήγηση), αμερόληπτες απαντήσεις (συνεντεύξεις για επίλυση υποψηφίων), κοπιαστικές έπαναληπτικές εργασίες (ρομπότ συγκόλλησης σε γραμμή παραγωγής), γρήγορες αποφάσεις (χρηματιστήριο), αναγνώριση προτύπων και τάσεων (αναγνώριση απάτης ή άνωμαλίας), ιατρική διάγνωση (αναγνώριση βλαβών σε ιατρικές εικόνες), παραγωγή

Υπάρχουν όμως και σοβαροί διαφαινόμενοι κίνδυνοι από την ανάπτυξη της Τεχνητής Νομοσύνης. Έπιστήμονες όπως ο Stephen Hawking και πολλοί άλλοι διατύπωσαν την άποψη ότι, όταν αποκτήσει μια μέρα τη δυνατότητα να έπανασχεδιάσει τον εαυτό της με όλο και ταχύτερο ρυθμό (όπως εξάλλου συμβαίνει με κάθε νέα τεχνολογία στον χώρο της πληροφορικής), αυτό θα οδηγήσει σε μια ανεξέλεγκτη έκρηξη νομοσύνης, που θα μπορούσε να οδηγήσει στον άφανισμό του ανθρώπου. Ο γνωστός μας Elon Musk, ένας από τους συνιδρυτές της OpenAI, χαρακτήρισε την Τεχνητή Νομοσύνη ως τη μεγαλύτερη ύπαρξιακή απειλή της ανθρωπότητας. 34000 προσωπικότητες υπέγραψαν μια ανοικτή επιστολή, με την οποία ζητούν την άμεση παύση εκπαίδευσης συστημάτων Τεχνητής Νομοσύνης πέραν του GPT-4, μέχρι να καταλάβουμε πώς λειτουργούν και πώς μπορούμε να τα ελέγχουμε καλύτερα.

άνθρωπος είναι κεντρικό για τη συζήτηση περί την Τεχνητή Νομοσύνη.

Τα όφελι (σημερινά και άναμενόμενα) από την Τεχνητή Νομοσύνη είναι αναρίθμητα: εξάλειψη ανθρώπινου λάθους (ρομποτική χειρουργική), μείωση κινδύνου (λειτουργία ρομπότ σε επικίνδυνο περιβάλλον), συνεχής διαθεσιμότητα (πελατειακή υποστήριξη, τηλεφωνητές 24 ώρες το 24ωρο, 7 μέρες την εβδομάδα),

νέου λογισμικού(!) κ.ά. Τα αρνητικά σημεία είναι και αυτά πολλά: ύψηλο κόστος (θα επικρατήσουν μόνον οι εταιρείες που διαθέτουν άφθονα δεδομένα και ισχυρούς υπερυπολογιστές), έλλειψη δημιουργικότητας (άνακυκλώνουν αυτό που έμαθαν), γενικευμένη ανεργία σε πολλούς κλάδους (έχουν ξεκινήσει συζητήσεις για τριήμερη εργασία και εξασφάλιση κατώτατου μισθού για όλους!), ανθρώπινη αδράνεια

(θυμηθείτε την ταινία Wall-e της Pixar), απάθεια (έλλειψη συναισθημάτων και ένσυναίσθησης), αδυναμία εξέλιξης (!) κ.ά.

Υπάρχουν όμως και σοβαροί διαφαινόμενοι κίνδυνοι από την ανάπτυξη της Τεχνητής Νοημοσύνης. Επιστήμονες όπως ο Stephen Hawking και πολλοί άλλοι διατύπωσαν την άποψη ότι, όταν αποκτήσει μια μέρα τη δυνατότητα να επανασχεδιάσει τον έαυτό της με όλο και ταχύτερο ρυθμό (όπως εξάλλου συμβαίνει με κάθε νέα τεχνολογία στον χώρο της πληροφορικής), αυτό θα οδηγήσει σε μια ανεξέλεγκτη έκρηξη νοημοσύνης, που θα μπορούσε να οδηγήσει στον άφανισμό του ανθρώπου. Ο γνωστός μας Elon Musk, ένας από τους συνιδρυτές της OpenAI, χαρακτήρισε την Τεχνητή Νοημοσύνη ως τη μεγαλύτερη υπαρξιακή απειλή της ανθρωπότητας. 34000 προσωπικότητες υπέγραψαν μια ανοικτή επιστολή, με την οποία ζητούν την άμεση παύση εκπαίδευσης συστημάτων Τεχνητής Νοημοσύνης πέραν του GPT-4, μέχρι να καταλάβουμε πώς λειτουργούν και πώς μπορούμε να τα ελέγχουμε καλύτερα. Μά το πρώτο που έκαναν οι δημιουργοί του ChatGPT ήταν να το συνδέσουν με το internet και να το αφήσουν να μάθει πώς να προγραμματίζει! Από την άλλη, όλες οι μεγάλες εταιρείες λογισμικού (Microsoft, Google, Amazon, Meta) αναπτύσσουν με γοργούς ρυθμούς ή καθεμία το δικό της λογισμικό Τεχνητής Γενικής Νοημοσύνης (GPT-5, PaLM 2, Alexa, LLaMA). Ο ίδιος ο Musk αναπτύσσει το πρόγραμμα Grok με δεδομένα από το Twitter (νύν X), και έχει εγκαταστήσει στην εταιρεία Tesla έναν υπερ-υπολογιστή, που μπορεί να εκτελεί 1 εκατομμύριο τρισεκατομμύρια υπολογισμούς το δευτερόλεπτο (exaflop) με 10000 κάρτες γραφικών GPU H100 της Nvidia για την ανάπτυξη του λογισμικού αυτόνομης οδήγησης autopilot, και όχι μόνον. Δηλαδή, μετά τη δημοσίευση της

ανοικτής επιστολής, η κόουρσα όχι μόνον δεν σταμάτησε, αλλά μάλλον επιταχύνθηκε!

Υπάρχουν όμως και σοβαρά εμπόδια, που μας κάνουν να πιστεύουμε ότι η πρόοδος των Μεγάλων Γλωσσικών Μοντέλων Τεχνητής Γενικής Νοημοσύνης δεν είναι δεδομένη. Ένα μεγάλο πρόβλημα με τη Μηχανική Μάθηση είναι ότι δεν γνωρίζουμε τί «ξέρει» και τί «καταλαβαίνει». Μου ανέφεραν το εξής έντυπωσιακό παράδειγμα: Εκπαίδευσαν ένα νευρωνικό δίκτυο, ώστε να ξεχωρίζει αγελάδες από αυτοκίνητα σε φωτογραφίες, και τα κατάφεραν πολύ καλά. Όταν όμως θέλησαν να μάθουν τί ακριβώς έβλεπε το σύστημα και πώς αναγνώριζε τις αγελάδες, διαπίστωσαν ότι το μόνο που κοίταζε ήταν το έδαφος κάτω από τα πόδια τους, και όταν έβλεπε χορτάρι, έβγαζε το συμπέρασμα ότι επρόκειτο για αγελάδα και όχι για αυτοκίνητο. Δηλαδή το σύστημα δεν είχε την παραμικρή ιδέα για το τί είναι αγελάδα και τί αυτοκίνητο· δηλαδή, αν έβλεπε ένα τρακτέρ σε ένα χωράφι, θα έλεγε ότι ήταν αγελάδα! Το ίδιο ακριβώς πρόβλημα υπάρχει και στο λογισμικό αυτόνομης οδήγησης των αυτοκινήτων Tesla, το οποίο και αυτό βασίζεται σε εικόνες video από τις 8 κάμερες του αυτοκινήτου. Στο αυτοκίνητο υπάρχει ένας πολύ ισχυρός υπολογιστής, ο οποίος επεξεργάζεται τα 8 αυτά video και αποκτάει αυξημένη αντίληψη του περιβάλλοντος του αυτοκινήτου σε σχέση με τον άνθρωπο. Σε γενικές γραμμές, το σύστημα αυτό τα καταφέρει πολύ καλύτερα από τον μέσο άνθρωπο οδηγό. Δυστυχώς, για τον λόγο ότι δεν μπορούμε ακόμα να ξέρουμε τί «ξέρει» και τί «καταλαβαίνει» ένα νευρωνικό δίκτυο, σε κάποιες περιπτώσεις το σύστημα «τρελαίνεται», και εκεί που το αυτοκίνητο προχωράει κανονικά, αποφασίζει χωρίς λόγο να το ρίξει πάνω σε μια κολώνα. Λείπει αυτό

που λέμε «κοινὸς νοῦς» (common sense).

Ένα ἄλλο πρόβλημα εἶναι ὅτι τὰ προγράμματα αὐτὰ ἐμφανίζουν «παραισθήσεις» (hallucinations). Ἐπειδὴ εἶναι γλωσσικὰ μοντέλα καὶ δὲν ἔχουν κατανόηση τῆς ἀλήθειας ἢ τοῦ ψεύδους τοῦ κειμένου που γεννοῦν λέξη πρὸς λέξη, σὲ πάρα πολλές περιπτώσεις γεννοῦν πληροφορία που ἀπλᾶ δὲν ὑπάρχει. Γιὰ παράδειγμα, ἕνας δικηγόρος ἀπὸ τὴ Νέα Ὑόρκη ἐπικαλέστηκε δικαστικὸ προηγούμενο σὲ ἕξι περιπτώσεις μὲ συγκεκριμένες ἀναφορὲς ἀπὸ τὶς βάσεις νομικῶν δεδομένων LexisNexis καὶ Westlaw, μόνον που οἱ ἀναφορὲς αὐτὲς που πρότεινε τὸ ChatGPT δὲν ὑπάρχουν... Τὸ ἴδιο ἀκριβῶς συνέβη καὶ σὲ ἕναν συνάδελφό μου ἀστρονόμο, ὅταν ρώτησε τὸ ChatGPT νὰ τοῦ προτείνει μιὰ ἀναφορὰ γιὰ τὴν ἐργασία του. Ἡ ἀναφορὰ που τοῦ πρότεινε εἶχε λογικὸ τίτλο, ὄνομα συγγραφέα, ὄνομα περιοδικοῦ, ἀριθμὸ τόμου καὶ σελίδας, μόνον που, ὅταν ὁ συνάδελφός μου ἀνέτρεξε στὸν συγκεκριμένον τόμο, διαπίστωσε ὅτι ἡ σελίδα τῆς ἐργασίας τὴν ὁποία πρότεινε τὸ ChatGPT ἦταν ἡ ἐπόμενη σελίδα ἀπὸ τὴν τελευταία σελίδα τοῦ τόμου, καὶ ἡ ἐργασία ἀπλᾶ δὲν ὑπῆρχε! Τὸ πρόβλημα μὲ αὐτὲς τὶς παραισθήσεις εἶναι ἡ φαινομενικὴ σιγουριὰ μὲ τὴν ὁποία διατυπώνονται ἀπὸ τὸ πρόγραμμα.

Ένα ἄλλο πρόβλημα εἶναι ὅτι τὰ προγράμματα αὐτὰ παρουσιάζουν σοβαρὲς προκαταλήψεις. Αὐτὸ εἶναι πολὺ λογικὸ, γιατί ἡ ἐκπαιδυσί τους βασίζεται σὲ ἀνθρώπινα δεδομένα, καὶ οἱ ἄνθρωποι ἐν γένει ἔχουμε πολλές προκαταλήψεις. Γιὰ παράδειγμα, τὸ πρόγραμμα ἐπιλογῆς ὑποψηφίων τῆς Amazon βαθμολογοῦσε ἀρνητικὰ γυναῖκες ὑποψηφίους, ἀπλᾶ καὶ μόνον γιατί τὸ δείγμα ἐκπαιδυσίς του εἶχε λιγώτερες αἰτίσεις ἀπὸ γυναῖκες (πρόκειται γιὰ γνωστὸ πρόβλημα μὴ ἰσορροπημένου δείγματος). Ἀντίστοιχα,

τὸ πρόγραμμα περιπολιῶν τῆς Ἀμερικανικῆς Ἀστυνομίας, τὸ ὁποῖο καὶ αὐτὸ βασίζεται σὲ ἀντίστοιχους ἀλγορίθμους, περιλαμβάνει περισσότερες περιπολίες σὲ γειτονιὲς μὲ μαῦρο πληθυσμό, διότι ἀντιλαμβάνεται ὅτι ἐκεῖ συμβαίνουν περισσότερα ἐγκλήματα. Ἐπισίμως αὐτὸ θεωρεῖται ρατσιστικὴ συμπεριφορὰ που ἔπρεπε νὰ σταματήσει. Ἡ ἐπέμβαση ὁμως τῶν προγραμματιστῶν γιὰ τὴ βελτίωση τῶν προγραμμάτων εἶχε μὴ ἀναμενόμενα ἀποτελέσματα!

Εἶναι συγκλονιστικὸ ὅτι τὸ λογισμικὸ ChatGPT, ἐνῶ ἀναπτύσσεται συνεχῶς ἀπὸ τὸν Νοέμβριον τὸ 2022, ὥστε νὰ καταστῆ πιὸ λειτουργικὸ καὶ φιλικὸ στὸν χρήστη, χωρὶς τὰ προβλήματα που προαναφέραμε (εἴμαστε τώρα στὴν ἔκδοσιν GPT-4 καὶ ἐκπαιδεύεται μὲ ραγδαίους ρυθμούς ἢ ἐπόμενη ἔκδοσιν GPT-5), ἄλλαξε συμπεριφορὰ πρὸς τὸ χειρότερον σὲ μιὰ περίοδο λίγων μηνῶν. Οἱ ἀπαντήσεις του ἔγιναν πιὸ ἀσαφείς καὶ πιὸ «χαζές». Δηλαδή, ἐνῶ πέρυσι τὸν Νοέμβριον μπορούσε νὰ λύσει τὸ 84% τῶν μαθηματικῶν προβλημάτων που τοῦ ἔθεται, σήμερα τὸ ποσοστὸ αὐτὸ ἔχει ἐκφυλιστεῖ στὸ 35 μὲ 50%, καὶ κανεὶς δὲν ξέρει γιατί. Σὲ ἄλλους δείκτες παρατηρεῖται μείωση ἀπὸ τὸ 97 στὸ 23%. Δηλαδή, ἐνῶ προσπαθοῦν οἱ προγραμματιστὲς νὰ βελτιώσουν κάποια στοιχεῖα τοῦ λογισμικοῦ, ἐπηρεάζουν μὲ ἄγνωστο τρόπο ἄλλα μέρη τοῦ λογισμικοῦ. Ἡ συμπεριφορὰ αὐτὴ ὀνομάζεται drift (ἀπομάκρυνση). Γιὰ τὸν λόγο αὐτὸ ὁ Max Tegmark ἀπὸ τὸ MIT καὶ πολλοὶ ἄλλοι μαζί μὲ αὐτὸν ὁμιλοῦν γιὰ τὴν ἀνάγκη νὰ ἀποκτήσουμε κατανοητὴ Τεχνητὴ Νοημοσύνη (intelligible intelligence) μὲ εὐθυγράμμιση (alignment) ὡς πρὸς τὶς ἀνάγκες καὶ τὰ συμφέροντα τοῦ ἀνθρώπου. Ἐντοπίζουν ὅτι ἐκεῖ ἔγκειται ὁ μεγαλύτερος κίνδυνος, καὶ δυστυχῶς εἴμαστε πολὺ μακριὰ ἀπὸ τὴν εὐθυγράμμιση.

Θὰ ἤθελα νὰ ὀλοκληρώσω αὐτὴ τὴ σύντομη εἰσήγηση στὸ τεράστιο αὐτὸ θέμα τῶν ἡμερῶν μας μὲ κάποιες προσωπικὲς σκέψεις. Μετὰ τὴν ἔκρηξη τῆς

στήμονες στὸν χῶρο τῆς ὑπολογιστικῆς μηχανικῆς καὶ σὲ ὀλόκληρη τὴν ἀνθρωπότητα γενικότερα, ἀλλὰ πῶς μπορούμε νὰ μιλήσουμε γιὰ τὰ οὐσιαστικὰ γνωρί-

Μετὰ τὴν ἔκρηξη τῆς Τεχνητῆς Γενικῆς Νομοσύνης ποὺ ἀναφέραμε, τοὺς τελευταίους 10 μῆνες εἶδαμε πολλοὺς ἀνθρώπους νὰ δημιουργοῦν ψηφιακὰ ἀντίγραφα τοῦ ἑαυτοῦ τους, εἶδαμε ἄλλους νὰ «παντρεύονται» ἓνα ψηφιακὸ avatar τὸ ὁποῖο ἐκπαίδευσαν κατ' ἐπιλογὴν τους, ἀκούσαμε ἄπειρους προβληματισμοὺς γιὰ τὸ ἂν ἔχουμε τὸ δικαίωμα «νὰ τὸ βγάλουμε ἀπὸ τὴν πρίζα», ὅταν λέει πράγματα ποὺ δὲν μᾶς ἀρέσουν, καὶ γενικὰ ἂν ἓνα τέτοιο λογισμικὸ ἔχει «μυαλό», «νομοσύνη», «συναισθήματα», «δικαιώματα», «κατανόηση», «ἐλευθερία», «αὐτοσυνειδησία»! Λυπᾶμαι πολὺ γι' αὐτὴ τὴ σύγχυση ποὺ ἐπικρατεῖ ἀνάμεσα στοὺς πιὸ ἔγκριτους ἐπιστήμονες στὸν χῶρο τῆς ὑπολογιστικῆς μηχανικῆς καὶ σὲ ὀλόκληρη τὴν ἀνθρωπότητα γενικότερα, ἀλλὰ πῶς μπορούμε νὰ μιλήσουμε γιὰ τὰ οὐσιαστικὰ γνωρίσματα τοῦ ἀνθρώπου, ὅταν ἔχουμε χάσει τὴ γνώση τῆς θεμελιώδους «πνευματικῆς φυσιολογίας» τοῦ ἀνθρώπου, καὶ ὅταν ἔχουμε νομοθετήσει τὴν ἀποστασία ἀπὸ τὸν νόμο τοῦ Θεοῦ;

Τεχνητῆς Γενικῆς Νομοσύνης ποὺ ἀναφέραμε, τοὺς τελευταίους 10 μῆνες εἶδαμε πολλοὺς ἀνθρώπους νὰ δημιουργοῦν ψηφιακὰ ἀντίγραφα τοῦ ἑαυτοῦ τους, εἶδαμε ἄλλους νὰ «παντρεύονται» ἓνα ψηφιακὸ avatar τὸ ὁποῖο ἐκπαίδευσαν κατ' ἐπιλογὴν τους, ἀκούσαμε ἄπειρους προβληματισμοὺς γιὰ τὸ ἂν ἔχουμε τὸ δικαίωμα «νὰ τὸ βγάλουμε ἀπὸ τὴν πρίζα», ὅταν λέει πράγματα ποὺ δὲν μᾶς ἀρέσουν, καὶ γενικὰ ἂν ἓνα τέτοιο λογισμικὸ ἔχει «μυαλό», «νομοσύνη», «συναισθήματα», «δικαιώματα», «κατανόηση», «ἐλευθερία», «αὐτοσυνειδησία»! Λυπᾶμαι πολὺ γι' αὐτὴ τὴ σύγχυση ποὺ ἐπικρατεῖ ἀνάμεσα στοὺς πιὸ ἔγκριτους ἐπι-

σματα τοῦ ἀνθρώπου, ὅταν ἔχουμε χάσει τὴ γνώση τῆς θεμελιώδους «πνευματικῆς φυσιολογίας» τοῦ ἀνθρώπου, καὶ ὅταν ἔχουμε νομοθετήσει τὴν ἀποστασία ἀπὸ τὸν νόμο τοῦ Θεοῦ; Εἶδαμε παραπάνω ὅτι ἀδυνατοῦμε νὰ ὀρίσουμε τί εἶναι ἡ νομοσύνη (τεχνητὴ ἢ φυσικὴ) γιὰ τὴν ὁποία μιλήσαμε σήμερα. Πόσοι ὅμως ἀπὸ ἐμᾶς γνωρίζουμε τί εἶναι ὁ νοῦς, ἡ ψυχὴ, ἡ καρδιὰ καὶ ἡ διάνοια, καὶ ποιά ἡ σχέσισ τους μὲ τὸ σῶμα τοῦ ἀνθρώπου, καὶ πῶς αὐτὲς οἱ ἔννοιες ἀναφέρονται στὴν Καινὴ Διαθήκη καὶ ἀναλύονται στὰ κείμενα τῶν Πατέρων τῆς Ἐκκλησίας ὅπως ὁ Ἅγιος Γρηγόριος ὁ Παλαμᾶς, ὁ Ἅγιος Συμεὼν ὁ Νέος Θεολόγος, ὁ Ἅγιος Μάξιμος ὁ

Όμολογητής; Ό Μητροπολίτης Ναυπάκτου Ίερόθεος στο βιβλίο του Όρθόδοξη Ψυχοθεραπεία γράφει: «Χάσαμε την παράδοσή μας, γι' αυτό και πολλοί από μᾶς ταυτίζουμε τὸν νοῦ μὲ τὴ λογική. Δὲν ὑποπτευόμαστε καθόλου ὅτι ἐκτὸς ἀπὸ τὴ λογικὴ ὑπάρχει καὶ ἄλλη δύναμη, ποὺ ἔχει μεγαλύτερη ἀξία, δηλαδή ὁ νοῦς, ἡ καρδιά. Όλος ὁ πολιτισμὸς εἶναι πολιτισμὸς ἀπώλειας τῆς καρδιάς. Καὶ κάτι ποὺ δὲν ἔχει ὁ ἄνθρωπος μέσα του δὲν μπορεῖ νὰ τὸ ἀντιληφθεῖ. Ἡ καρδιά νεκρώθηκε, ὁ νοῦς σκοτίσθηκε, καὶ δὲν μποροῦμε νὰ ἀντιληφθοῦμε τὴν παρουσία τους. Γιὰ τὸν ἄνθρωπο ποὺ ἔχει μέσα του τὸ Ἅγιο Πνεῦμα, γι' αὐτὸν ποὺ βρίσκεται “ἐν τῇ ἀποκαλύψει”, δὲν χρειάζονται πολλὲς διασαφηνίσεις, γιατί αὐτὸς γνωρίζει ἀπὸ τὴν πειρὰ του τὴν παρουσία καὶ τὴν ὕπαρξη τοῦ νοῦς, τῆς καρδιάς». Καὶ συμπληρώνω: καὶ στὸν ἑαυτό του, καὶ στὸν συνάνθρωπο. Εἴπαμε προηγουμένως ὅτι ἕνας ψυχικὰ ἄρρωστος δὲν μπορεῖ νὰ ξεχωρίσει ἕναν ἄνθρωπο ἀπὸ ἕνα ρομπότ. Πόσοι ὅμως ἀπὸ ἐμᾶς εἴμαστε πραγματικὰ ψυχικὰ ὑγιεῖς;

Ό Ἅγιος Μάξιμος ὁ Όμολογητὴς μᾶς διδάσκει ὅτι τὸ ἀνθρώπινο σῶμα εἶναι ὁ ναός, ἱερὸ βῆμα εἶναι ἡ καρδιά καὶ θυσιαστήριο ὁ νοῦς. Καὶ ὅπως στὴν Ἁγία Τράπεζα δὲν ἐπιτρέπεται νὰ βάλουμε τίποτε ἄλλο, παρὰ μόνο τὸ Ἅγιο Δισκο-

πότηρο καὶ τὸ ἱερὸ Εὐαγγέλιο, ἔτσι καὶ στὸ θυσιαστήριο τοῦ νοῦ τοῦ ἀληθινοῦ χριστιανοῦ δὲν ἐπιτρέπεται νὰ ὑπάρχει τίποτε ἄλλο παρὰ μόνον ὁ Χριστός, ἡ ἀδιάλειπτη μνήμη τοῦ Ἰησοῦ Χριστοῦ. Αὐτὴ εἶναι ἡ ἡσυχαστικὴ ἐμπειρία τῆς Ἐκκλησίας, ποὺ συνεχίζεται μέχρι τὶς μέρες μας, περνώντας καὶ ἀπὸ τοὺς μεγάλους σύγχρονους ἡσυχαστὲς τοῦ Ἁγιοφάραγγου τῶν Ἀστερουσίων, τοῦ Ἁγίου Εὐμενίου τῆς Ἐθιᾶς καὶ τῆς γερόντισσας Γαλακτίας τῆς Παναγίας τῆς Καλυβιανῆς. Αὐτοὺς τοὺς ἀνθρώπους ἐμεῖς τοὺς γνωρίσαμε. Ἡ ἐμπειρία αὐτὴ εἶναι σήμερον ζωντανὴ στὴν νότια Κρήτη καὶ σὲ πολλὰ μέρη τῆς Ἑλλάδας, ἀλλὰ καὶ τοῦ κόσμου. Καὶ αὐτὴν τὴν ἐμπειρία πρέπει νὰ τὴν σπουδάζουμε ὅλη μας τὴν ζωὴ καὶ νὰ τὴ μεταδίδουμε στὶς ἐπόμενες γενεές. Μόνον ἔτσι θὰ μπορέσουμε νὰ ἀντιμετωπίσουμε χωρὶς φόβο τὶς νέες προκλίσεις ποὺ φέρνει ἡ νέα τεχνολογία. Ἀλλιῶς, ὁ κόσμος μας κινδυνεύει νὰ δαμμονιστεῖ. Ἄν δὲν ἔχει δαμμονιστεῖ ἤδη...

Σᾶς εὐχαριστῶ γιὰ τὴν προσοχή σας!

ΙΩΑΝΝΗΣ ΚΟΝΤΟΠΟΥΛΟΣ
Διευθυντὴς Ἐρευνῶν
Κέντρο Ἐρευνῶν Ἀστρονομίας
καὶ Ἐφαρμοσμένων Μαθηματικῶν
τῆς Ἀκαδημίας Ἀθηνῶν

Σημείωση: Λίγες μέρες πρὶν τὴ δημοσίευση τοῦ παρόντος κειμένου, κυκλοφόρησε μιὰ ἐσωτερικὴ ἐπιστολὴ πρὸς τοὺς ἐργαζομένους τῆς OpenAI, ποὺ τοὺς προειδοποιοῦσε γιὰ τοὺς κινδύνους τοῦ ἐρευνητικοῦ προγράμματος Q*, μὲ τὸ ὅποιο ἡ ἐταιρεία ὑποτίθεται ὅτι ἔφτασε στὴ Γενικὴ Τεχνητὴ Ὑπερ-Νοημοσύνη! Ό ἴδιος ὁ Sam Altman δῆλωσε δημοσίως ὅτι οἱ ἄνθρωποι ἀναρωτιοῦνται ἂν αὐτὸ ποὺ ἔφτιαξαν δὲν εἶναι ἀπλῶς ἕνα νέο ἐργαλεῖο, ἀλλὰ ἕνα νέο «ὄν» (“Is this a tool we’ve built or a creature we’ve built?”)! Τὴν ἐπόμενη μέρα τὸ διοικητικὸ συμβούλιο τῆς ἐταιρείας ἀπέλυσε τὸν Altman, ἀλλὰ ἀναγκάστηκε νὰ τὸν ἐπαναπροσλάβει μιὰ ἐβδομάδα ἀργότερα, ὅταν 700 ἐργαζόμενοι τῆς ἐταιρείας ἀπέιλσαν ὅτι θὰ παρατηθοῦν. Ἀντιλαμβανόμαστε ὅλοι ὅτι οἱ ἐξελίξεις στὸν χῶρο τῆς Τεχνητῆς Νοημοσύνης εἶναι καταγιστικές. Ό Θεὸς ἂς βάλει τὸ χέρι Του.
