

УДК 342.9:004.8

DOI <https://doi.org/10.24144/2307-3322.2024.86.3.40>

## ОБМЕЖЕННЯ СВОБОДИ ШТУЧНОГО ІНТЕЛЕКТУ ТА ОКРЕМІ ІНІЦІАТИВИ ІНСТИТУТУ МАЙБУТНЬОГО ЖИТТЯ

**Гачкевич А.,**  
*кандидат юридичних наук, доцент,  
доцент кафедри міжнародного та кримінального права  
Національного університету «Львівська політехніка»  
ORCID: 0000-0002-8494-1937  
e-mail: andrii.o.hachkevych@lpnu.ua*

**Кошелєв Г.,**  
*здобувач вищої освіти за спеціальністю 291  
«Міжнародні відносини,  
суспільні комунікації та регіональні студії»  
Національного університету «Львівська політехніка»  
ORCID: 0009-0006-1133-3224*

**Гачкевич А., Кошелєв Г. Обмеження свободи штучного інтелекту та окремі ініціативи Інституту майбутнього життя.**

Дискусії щодо свободи штучного інтелекту пов'язані з існуванням різних позицій з приводу того, що важливіше – технологічний прогрес чи традиційно важливі суспільні цінності, охорону яких забезпечує право. Пошук балансу між ними сприяє вирішенню проблеми управління штучним інтелектом, яка набула особливої актуальності за останні 2-3 роки з появою генеративного штучного інтелекту. У цій статті зроблена спроба визначити передумови та напрями обмеження штучного інтелекту, які повинні бути врахованими при виробленні та реалізації державної політики у сфері штучного інтелекту. Для того, щоб показати межу, за перетину якої обмеження ймовірно стануть надмірними, автори досліджують ініціативи Інституту майбутнього життя – однієї з найбільш авторитетних у світі некомерційних організацій, діяльність яких має відношення до новітніх технологій. Порушені у відкритих листах, які складені, оприлюднені та представлені для підписання Інститутом майбутнього життя, проблеми дозволяють виокремити головні обмеження для штучного інтелекту, зокрема: призупинення навчання ультрасучасних систем штучного інтелекту, більш потужних за GPT-4; формування правил у вигляді окремого зведення принципів, наприклад, розроблених Інститутом майбутнього життя Асиломарських принципів штучного інтелекту; недопустимість автономної зброї, здатної уражати без людського нагляду тощо. Крім того, у відкритих листах наголошено необхідність поглибленого вивчення штучного інтелекту як явища. Окрема увага приділена майже невідомому в українській науці Інституту майбутнього життя, заснованому приблизно десять років тому для того, щоб використання сучасних технологій було ефективним та збалансованим. Результати дослідження містять інтелектуальну карту, яка може стати цінною для кращого розуміння реальних та потенційних небезпек, зумовлених свободою штучного інтелекту та необмеженим технологічним прогресом. Крім того, автори підготували лінію часу, де систематизовані відкриті листи Інституту майбутнього життя, які стосуються штучного інтелекту та були оприлюднені за останні десять років.

**Ключові слова:** штучний інтелект; свобода штучного інтелекту; межі свободи штучного інтелекту; управління штучним інтелектом; Інститут майбутнього життя; Асиломарські принципи штучного інтелекту; відкриті листи.

**Hachkevych A., Kosheliev H. Limits for artificial intelligence freedom and the initiatives of the Future of Life Institute.**

Discussions about the freedom of artificial intelligence revolve around differing views on the importance of technological progress versus traditional social values protected by the law. Striking

a balance between these two aspects is crucial for addressing the challenges of artificial intelligence governance. This issue has become increasingly relevant over the past two to three years, particularly with the rise of generative artificial intelligence. This article aims to define the preconditions and limitations of artificial intelligence that should underlie a state policy in this area. To illustrate the point at which restrictions may become excessive, the authors examine the initiatives of the Future of Life Institute, a respected non-profit organization focused on cutting-edge technologies. The open letters compiled, published, and circulated for signatures by this organization highlight key restrictions on artificial intelligence. These include: a suspension of the training of advanced AI systems more powerful than GPT-4, the establishment of a distinct set of guidelines, e.g. the Asilomar Principles for Artificial Intelligence, and a prohibition on lethal autonomous weapons (LAWS) that can strike without human supervision. Additionally, the open letters emphasize the necessity of conducting in-depth studies on artificial intelligence as a phenomenon. Special focus is given to the Future of Life Institute, which is relatively unknown in Ukrainian academic circles. Founded about a decade ago, it seeks to ensure that modern technologies are utilized in an effective and balanced manner. The study's results offer a valuable framework for understanding both the real and potential dangers associated with the freedom of artificial intelligence and unrestricted technological advancement. This highlights the seriousness of the issue covered in this article. In addition, the authors have prepared a timeline designed for the systematization of the open letters of the Future of Life Institute related to artificial intelligence that have been published over the past ten years.

**Key words:** artificial intelligence; freedom of artificial intelligence; limits of artificial intelligence freedom; governance of artificial intelligence; Future of Life Institute; Asilomar Principles for Artificial Intelligence; open letters.

**Постановка проблеми.** Свобода штучного інтелекту, під якою ми розуміємо можливість вільно розробляти та використовувати технології та системи штучного інтелекту, викликає сьогодні великий інтерес та дискусії, пов'язані з тим, наскільки абсолютною вона може бути.

Питання меж такої свободи тісно пов'язане з пошуком балансу між двома пріоритетами: технологічного прогресу, для якого надмірні обмеження не є сприятливими, а також – забезпечення традиційно важливих суспільних цінностей, як-от охорона прав людини. Певною мірою перший з пріоритетів відображає очікування щодо майбутнього, тоді як інший – здобутки минулого, які збережені в тому числі завдяки правовій сфері.

Пошук балансів передбачає вироблення цілісного підходу, необхідного для вирішення проблеми управління штучним інтелектом. На нашу думку, в цьому плані надзвичайно корисними є окремі ініціативи Інституту майбутнього життя – авторитетної некомерційної організації, – відкриті листи, підтримані серед іншого відомими вченими та представниками компаній, що належать до лідерів галузі штучного інтелекту на світовому рівні. З одного боку, ці ініціативи втілюють у собі позиції фахівців з новітніх технологій, які добре розуміють перспективи впливу штучного інтелекту на суспільство, з іншого – завдяки ініціативам такого роду в подальшому формується державна та міжнародна політика щодо управління штучним інтелектом, включаючи особливості правового регулювання.

**Стан опрацювання проблематики.** Багатоаспектна проблема управління штучним інтелектом полягає в тому, щоб гарантувати розвиток сфери штучного інтелекту у відповідності до інтересів суспільства, насамперед – забезпечення традиційно важливих суспільних цінностей – прав людини, демократії, національної безпеки – при цьому ризики повинні бути мінімізовані, а загрози – усунені. В основі управління, найбільш ефективними засобами якого є державно-політичні, мають лежати базові принципи, зміст яких залишається предметом обговорення. Разом з тим, не підлягає сумніву ідея балансів – між небезпеками та можливостями, а також – технологічним прогресом та збереженням цінностей.

Серед українських вчених ця проблема порушена у працях Г. Андрощука, О. Баранова, Н. Вітницької, М. Карчевського, Ю. Кривицького, О. Намясенко та С. Ратушного тощо. Вона є дуже активно досліджуваною іноземними вченими, зокрема в контексті: вивчення впливу етичних правил [1], характеристики національних стратегій штучного інтелекту [2], визначення концептуальних рамок [3], пояснення стану глобального управління та співробітництва [4] та ін.

Обраний аспект проблеми, який має відношення до ініціатив у вигляді оприлюднених відкритих листів Інституту майбутнього життя щодо штучного інтелекту, а також до діяльності цієї

організації в загальному, ще не був дослідженим у наукових працях українськими вченими. Під таким кутом зору проблема не вивчалась й іноземними вченими, хоча й в окремих публікаціях Інститут майбутнього життя та його внесок у розвиток знання про штучний інтелект відзначались низкою авторів [5; 6; 7 та ін.].

Інтерес авторів цієї статті до обраної тематики пояснюється двома причинами. По-перше, ініціативи Інституту майбутнього життя показують цілком обґрунтовані напрями обмеження штучного інтелекту, які повинні бути врахованими при виробленні та реалізації державної політики у сфері штучного інтелекту та гарантувати збалансоване управління сферою. По-друге, аналіз ініціатив сприяє кращому розумінню факторів довіри до технологій штучного інтелекту, від якої залежить те, наскільки суспільство буде підтримувати ідею використання та вдосконалення новітніх технологій.

**Метою статті** є визначення передумов та напрямів обмеження штучного інтелекту на основі аргументації відкритих листів, оприлюднених з ініціативи Інституту майбутнього життя. Окрема увага приділена Інституту майбутнього життя як некомерційній організації та його внеску у вирішення проблеми управління штучним інтелектом.

**Виклад основного матеріалу.** Інститут майбутнього життя (Future of Life Institute, далі – Інститут) є неприбутковою організацією, основний вектор діяльності якої спрямований на сприяння ефективному використанню сучасних технологій – зростанню переваг при пом'якшенні ризиків [8].

Він був створений у Кембриджі штату Массачусетс у 2014 р. за ініціативи п'яти людей: Я. Таллінна, який раніше став співзасновником науково-дослідного інституту для вивчення екзистенційного ризику при Кембриджському університеті, М. Тегмарка, професора космології Массачусетського технологічного інституту, який співпрацював зі шведським філософом Н. Бостромом (його дисертація про суперінтелект дала поштовх створенню Інституту), та дружини М. Чіти-Тегмарк, дослідниці Бостонського університету, В. Краковної, співробітниці компанії DeepMind, Е. Агірре, професора фізики Каліфорнійського університету в Санта-Крус [9, с. 93].

Згадуючи історію створення Інституту, М. Тегмарк пояснює головний мотив діяльності: «сприяти тому, щоб майбутнє життя існувало й було б якомога чудовішим» [10, с. 34]. Він також зауважив, що засновники усвідомлювали дві можливі траєкторії впливу новітніх технологій на суспільство: давати силу людям та процвітання або призвести до самознищення. Крім штучного інтелекту, були визначені також три інші вектори подальшої діяльності: біотехнології, ядерна зброя та зміна клімату [10, с. 34].

Сьогодні основними формами діяльності Інституту є наступні: а) участь у виробленні політики та адвокація; б) інформаційно-просвітницька діяльність; в) проведення досліджень та сприяння їм; г) надання грантів; д) розбудова інституцій; е) організація та проведення заходів [8].

Головний офіс Інституту знаходиться в США, а так само має представництво на території ЄС (Бельгія). Він належить до консорціуму Інституту безпеки штучного інтелекту Національного інституту стандартів і технологій США, до Партнерства зі штучного інтелекту (Partnership on AI) та ін. Він також здійснює співпрацю з: (1) некомерційними організаціями («Спільнота майбутнього», The Future Society; Інститут інженерів з електротехніки та електроніки, Institute of Electrical and Electronics Engineers), (2) міжнародними організаціями (ОЕСР; окремі листи адресовані ЄС (щодо Закону ЄС про штучний інтелект) та ООН (стосовно регулювання на міжнародному рівні питання летальної автономної зброї), (3) дослідницькими інститутами (Центр сумісного з людиною штучного інтелекту при Каліфорнійському університеті в Берклі, Center for Human-Compatible Artificial Intelligence; Центр європейських політичних студій, Centre for European Policy Studies) та (4) приватними компаніями (серед підписантів відкритих листів – представники Google DeepMind; Apple Inc.; IBM Research; OpenAI та ін.) [11].

Наша увага зосереджена на відкритих листах Інституту, які стосуються штучного інтелекту та систематизовані на рисунку [12–27].



Варто зауважити, що відкриті листи у практиці Інституту відрізняються між собою, хоча й майже всі вони підготовлені з ініціативи Інституту та присвячені важливому питанню зі сфери управління штучним інтелектом. Умовним є поділ відкритих листів за предметом: військове використання штучного інтелекту та летальна автономна зброя, а також – невійськове.

В загальному під відкритим листом можемо розуміти громадську ініціативу щодо важливого питання, виражену у публічній формі та придатну для підтримки іншими людьми, крім підписантів (хоча й далеко не всі відкриті листи придатні для подальшого підписання), завдяки чому громадськість може брати участь в управлінні загальносуспільними справами [28].

Аналіз наявних відкритих листів показує, що допустимі обмеження свободи штучного інтелекту при виробленні та реалізації державної політики, а також – в контексті управління штучного інтелекту, можуть бути встановлені у таких напрямках.

### **1. Призупинення навчання ультрасучасних систем штучного інтелекту, більш потужних за GPT-4.**

На момент написання статті чи не найважливішим відкритим листом є «Призупинення великих експериментів зі штучним інтелектом». У ньому викладений заклик до всіх лабораторій зі штучного інтелекту «негайно призупинити щонайменше на півроку навчання систем штучного інтелекту, більш потужних за GPT-4», серед ризиків, які пояснюють таку вимогу, – пропаганда, масова автоматизація робочих місць, заміщення людей та втрата контролю в масштабах суспільства (підписантами є Й. Бенжіо, С. Расселл, І. Маск, С. Возняк, А. Янг та ін.). Оприлюднення листа відбулося через тиждень після випуску GPT-4 OpenAI (ультрасучасні великі мовні моделі стають «конкурентоспроможними з людьми у загальних завданнях»). У листі відзначена пріоритетність створення потужних систем штучного інтелекту, які були б більш точними, безпечними, пояснювальними, прозорими, надійними, узгодженими. Він також рекомендує більш жорстке державне регулювання, проведення незалежних аудитів перед навчанням систем штучного інтелекту, а також моніторинг високопродуктивного штучного інтелекту та належне державне фінансування технічних досліджень безпекових питань [24].

### **2. Формування правил у вигляді окремого зведення принципів, наприклад, Асилмарських принципів штучного інтелекту.**

Однієї з найбільш значимих подій, які мають відношення до Інституту, стала конференція в Асилмарі щодо корисного штучного інтелекту (Asilomar Conference on Beneficial AI) у січні 2017 р. Результатом конференції є прийняття за підсумками обговорення фахівців з інформаційних технологій, економіки, права, етики та філософії зведення принципів для штучного інтелекту (дослідження та розвиток). Воно складається з 23 положень, об'єднаних у 3 категорії: а) аспекти досліджень, б) етика та цінності, в) довгострокові проблеми. У зведенні наголошено на тому, що дослідження повинні мати мультидисциплінарний характер та бути націленими на створення корисного штучного інтелекту. Рекомендовано налагодження зв'язків між науковою та політичною елітою, а також виховання культури довіри та співпраці. Категорія «етика та цінності» складається з набору стандартів, що нагадують етичні принципи штучного інтелекту: безпека, прозорість, дотримання прав людини, повага до приватності, відповідальність та ін. Серед довгострокових проблем можемо умовно виокремити: ті, що вимагають обережності та попереджувальних дій, – поміркованості стосовно верхніх меж можливостей штучного інтелекту, запобігання непередбачуваності передового штучного інтелекту, зменшення катастрофічної та екзистенційної небезпеки, а також підконтрольності систем, здатних до інтенсивного самовдосконалення та самовідтворення, а також проблему суспільного блага, згідно з якою штучний суперінтелект повинен служити загальноновизнаним етичним ідеалам і користі усього людства, а не однієї держави чи організації [17].

### **3. Недопустимість автономної зброї на базі штучного інтелекту, здатної уражати без людського нагляду.**

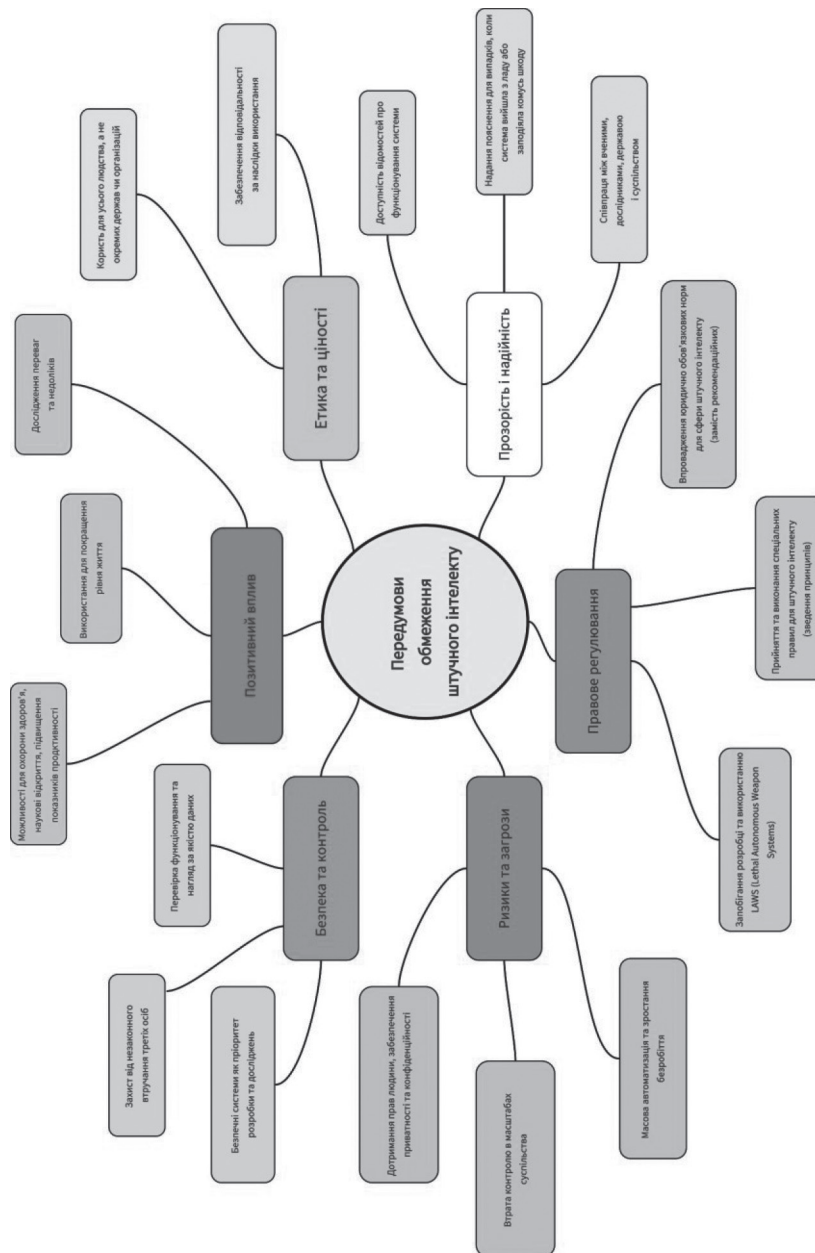
Позиція Інституту з приводу летальної автономної зброї є категоричною – така зброя в жодному разі не повинна бути застосовуваною. Ця позиція відображена відразу в кількох листах, кожен з яких має свою аргументацію. Скажімо, військове використання штучного інтелекту як зброї назване неприйнятним, як і те, що комп'ютер буде приймати рішення, чи забирати людське життя [21]. В іншому листі поява зброї на базі штучного інтелекту відзначена як революція у засобах ведення війни, поряд з винайденням пороху та створенням ядерної зброї. Незважаючи на це, наголошено, що володіння такою зброєю терористами, диктаторами чи військовими командирами може призвести до жахливих наслідків, а тому слід взагалі відмовитись від реалізації ідеї гонки озброєнь систем штучного інтелекту [16].

Інститут неодноразово підтримував зусилля ООН щодо розробки та прийняття конвенції про заборону летальної автономної зброї [19–22].

Один з відкритих листів – «Пріоритети досліджень для надійного та корисного штучного інтелекту» – містить доповнення у вигляді статті, в якій описані ці пріоритети [14; 15]. Варто зауважити, що до них належать:

- економічні (вплив штучного інтелекту на ринок праці – автоматизація та зростання безробіття, негативні зміни у різних сферах економіки та виробничих процесах, нові підходи до управління),
- етико-правові (відповідальність безпілотних автомобілів, забезпечення приватності та конфіденційності, питання автономної зброї тощо),
- комп'ютерно-технічні (перевірка системи, точність її функціонування, контроль людиною та захист як запобігання незаконного втручання у систему).

Для того, щоб узагальнити передумови обмеження штучного інтелекту, які були визначені в результаті аналізу відкритих листів Інституту майбутнього життя та розглянуті під час дослідження, пропонуємо інтелектуальну карту. Вона може стати цінною для кращого розуміння небезпек, зумовлених свободою штучного інтелекту та необмеженим технологічним прогресом.



**Висновки.** Намір створення Інституту майбутнього життя виник близько десяти років тому завдяки усвідомленню його засновниками двох фактів. По-перше, величезного потенціалу та в подальшому невідворотного частого та широкого застосування штучного інтелекту та інших новітніх технологій. По-друге, ризиків та загроз такого застосування, які досі залишаються критичними (а з появою нових можливостей – ще більш небажаними, ніж в той період, коли Інститут був заснований).

Вплив на управління штучним інтелектом, пов'язаний з виробленням та реалізацією державної та міжнародної політики у сфері новітніх технологій, став одним з лейтмотивів діяльності Інституту, який дуже виразно простежується у положеннях відкритих листів щодо штучного інтелекту.

Ми розглянули понад десять листів, розміщених на веб-сайті Інституту та систематизували їх за допомогою лінії часу. Зауважимо, що найбільш виправданими напрямками встановлення обмежень для свободи штучного інтелекту є: призупинення навчання ультрасучасних систем, формування зведень правил, які повинні лежати в основі державної політики та правового регулювання у сфері штучного інтелекту, а також заборона розробок та використання летальної автономної зброї.

У відкритих листах так само отримала підтримку ідея радше жорсткого регулювання сфери, ніж м'якого, що передбачає прийняття спеціальних законів, а також запровадження відповідних дозволів, особливо для систем високого ризику.

Ще одним пріоритетом стосовно штучного інтелекту Інститут практично від моменту заснування визначив проведення досліджень, насамперед для того, щоб системи штучного інтелекту були безпечними та надійними, а також приносили користь людству.

#### **СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ:**

1. Larsson S. On the Governance of Artificial Intelligence through Ethics Guidelines. *Asian Journal of Law and Society*. 2020. Volume 7, Issue 3. Pp. 437–451. DOI: <https://doi.org/10.1017/als.2020.19>.
2. Radu R. Steering the Governance of Artificial Intelligence: National Strategies in Perspective. *Policy and Society*. 2021. Volume 40, Issue 2. Pp. 178–193. DOI: <https://doi.org/10.1080/14494035.2021.1929728>.
3. de Almeida P., dos Santos C., Farias J. Artificial Intelligence Regulation: a Framework for Governance. *Ethics and Information Technology*. 2021. Volume 23. Pp. 505–525. DOI: <https://doi.org/10.1007/s10676-021-09593-z>.
4. Butcher J., Beridze I. What is the State of Artificial Intelligence Governance Globally? *The RUSI Journal*. 2019. Volume 164, Issue 5–6. Pp. 88–96. DOI: <https://doi.org/10.1080/03071847.2019.16942605>.
5. Li G., Deng X., Gao Z., Chen F. Analysis on Ethical Problems of Artificial Intelligence Technology. *Proceedings of the 2019 International Conference on Modern Educational Technology*. New York, 2019. Pp. 101–105. DOI: <https://doi.org/10.1145/3341042.3341057>.
6. Brodecki Z., Konopačka M. Thinking out of the box: The human being in the AI era. *Per mare ad astra. Space technology, governance and law*. 2022. Volume II. Pp. 195–214.
7. Leikas J., Koivisto R., Gotcheva N. Ethical Framework for Designing Autonomous Intelligent Systems. *Journal of Open Innovation: Technology, Market, and Complexity*. 2019. Volume 5, Issue 1, Article 18. DOI: <https://doi.org/10.3390/joitmc5010018>.
8. Future of Life Institute. URL: <https://futureoflife.org/> (дата звернення: 10.12.2024 р.).
9. Paura R. The Notion of Existential Risk and Its Role for the Anticipation of Technological Development's Long-Term Impact. *Anticipation, Agency and Complexity*. Cham, 2019. Pp. 79–97. DOI: [https://doi.org/10.1007/978-3-030-03623-2\\_6](https://doi.org/10.1007/978-3-030-03623-2_6).
10. Tegmark M. *Life 3.0. Being Human in the Age of Artificial Intelligence*. New York, 2017.
11. Future of Life Institute (FLI). URL: <https://www.influencewatch.org/non-profit/future-of-life-institute-fli/> (дата звернення: 10.12.2024 р.).
12. AI Economics Open Letter. URL: <https://futureoflife.org/open-letter/ai-economics-open-letter/> (дата звернення: 10.12.2024 р.).
13. Digital Economy Open Letter. URL: <https://futureoflife.org/open-letter/digital-economy-open-letter/> (дата звернення: 10.12.2024 р.).

14. Research Priorities for Robust and Beneficial Artificial Intelligence: An Open Letter. URL: <https://futureoflife.org/open-letter/ai-open-letter/> (дата звернення: 10.12.2024 р.).
15. Russell S., Dewey D., Tegmark M. Research Priorities for Robust and Beneficial Artificial Intelligence. *AI Magazine*. 2015. Volume 36, Number 4. Pp. 105–114. <https://doi.org/10.48550/arXiv.1602.03506>.
16. Autonomous Weapons Open Letter: AI & Robotics Researchers. URL: <https://futureoflife.org/open-letter/open-letter-autonomous-weapons-ai-robotics>.
17. Asilomar AI Principles. URL: <https://futureoflife.org/open-letter/ai-principles/> (дата звернення: 10.12.2024 р.).
18. The Principles – Signatories List. URL: <https://futureoflife.org/open-letter/principles-signatories/> (дата звернення: 10.12.2024 р.).
19. An Open Letter to the United Nations Convention on Certain Conventional Weapons. URL: <https://futureoflife.org/open-letter/autonomous-weapons-open-letter-2017> (дата звернення: 10.12.2024 р.).
20. 2018 Statement to United Nations on Behalf of LAWS Open Letter Signatories. URL: <https://futureoflife.org/open-letter/statement-to-united-nations-on-behalf-of-laws-open-letter-signatories> (дата звернення: 10.12.2024 р.).
21. Lethal Autonomous Weapons Pledge. URL: <https://futureoflife.org/open-letter/lethal-autonomous-weapons-pledge/> (дата звернення: 10.12.2024 р.).
22. Autonomous Weapons Open Letter: Global Health Community. URL: <https://futureoflife.org/open-letter/medical-lethal-autonomous-weapons-open-letter/> (дата звернення: 10.12.2024 р.).
23. Foresight in AI Regulation Open Letter. URL: <https://futureoflife.org/open-letter/foresight-in-ai-regulation-open-letter/> (дата звернення: 10.12.2024 р.).
24. Pause Giant AI Experiments: An Open Letter. URL: <https://futureoflife.org/open-letter/pause-giant-ai-experiments/> (дата звернення: 10.12.2024 р.).
25. FAQs about FLI’s Open Letter Calling for a Pause on Giant AI Experiments. URL: <https://futureoflife.org/ai/faqs-about-flis-open-letter-calling-for-a-pause-on-giant-ai-experiments/> (дата звернення: 10.12.2024 р.).
26. AI Licensing for a Better Future: On Addressing Both Present Harms and Emerging Threats. URL: <https://futureoflife.org/open-letter/ai-policy-for-a-better-future-on-addressing-both-present-harms-and-emerging-threats/> (дата звернення: 10.12.2024 р.).
27. Open letter calling on world leaders to show long-view leadership on existential threats <https://futureoflife.org/open-letter/long-view-leadership-on-existential-threats/> (дата звернення: 10.12.2024 р.).
28. Kosheliev H., Nachkevych A. Open Letters as an Instrument of Public Influence on the Governance of Artificial Intelligence (Manuscript). URL: <https://peers.international/paper/open-letters-instrument-public-influence-governance-artificial-intelligence> (дата звернення: 10.12.2024 р.).