

\*Peter Königs |  
 ترجمة عمر المغربي | \*\*Omar Almagharebi

## أزمة التساؤمية في حقل أخلاقيات الذكاء الاصطناعي

### The Negativity Crisis of AI Ethics

**ملخص:** على الرغم من الإمكانيات الإيجابية الهائلة التي ينطوي عليها الذكاء الاصطناعي، فقد قدم مجتمع أخلاقيات الذكاء الاصطناعي صورة قائمة إلى حد بعيد عن تداعياته الأخلاقية. تفحص هذه الدراسة المزعنة السلبية السائدة داخل هذا الحقل من منظور فلسفة العلم. وتتاتي هذه السلبية الطاغية من الطريقة الخاصة التي يُنظم بها مؤسسيًا، والتي تُجبر المشغلين به على تصويره وفق منظور سلبي. وغدت الصورة العامة التي تقدمها أدبيات أخلاقيات الذكاء الاصطناعي سلبية وذات منظور أحادي. بناءً عليه، ينبغي التشكيك في السردية التساؤمية المتداولة، والبحث عن طرائق لإصلاح المنظومة التي أفرزتها.

**كلمات مفتاحية:** الذكاء الاصطناعي، فلسفة العلم، الإستيمولوجيا الاجتماعية، الذكاء المانديفي، التحيز، أخلاقيات الذكاء الاصطناعي.

**Abstract:** Despite the great positive potential of AI, the AI ethics community has presented a rather gloomy picture of AI's ethical implications. This paper examines the negativity within AI ethics through a philosophy of science lens. The prevailing negativity is a result of the particular way the discipline is institutionally organized, which pressures AI ethicists to portray AI in a critical light. As a consequence, the over all picture of AI offered by the AI ethics community is one-sided and negatively biased. We should be skeptical about the negative narrative promoted by AI ethics and explore ways of reforming the system.

**Keywords:** Artificial Intelligence, Philosophy of Science, Social Epistemology, Mandevillian Intelligence, Bias, Ethics of AI.

\* قسم الفلسفة والعلوم السياسية، جامعة دورتموند التقنية، دورتموند، ألمانيا.

Department of Philosophy and Political Science, TU Dortmund University, Dortmund, Germany.

Email: peter.koenigs@tu-dortmund.de

\*\* باحث مساعد في المركز العربي للأبحاث ودراسة السياسات.

Assistant Researcher, Arab Center for Research and Policy Studies.

Email: omar.almagharebi@doha institute.edu.qa

Peter Königs, "The Negativity Crisis of AI Ethics," *Synthese*, vol. 206, no. 277 (2025).

<https://doi.org/10.1007/s11229-025-05378-9>

This article is licensed under a Creative Commons Attribution 4.0 International License.

هذا المقال مرخص بموجب رخصة المشاع الإبداعي الدولية 4.0.

## مقدمة



يتزايد المزاج التساؤمي وسط المنتسين إلى الذكاء الاصطناعي اطراداً مع ازدهار المبحث المتعلقة بأخلاقياته. يكفي أن يُقلّب المرء أعداد دوريات "أخلاقيات التقنية" حتى يخلص إلى أن تطور الذكاء الاصطناعي يبدو، قبل كل شيء، مدعماً للقلق. فالنقاشات الكبرى تدور، تقريباً بلا استثناء، حول قائمة مطولة من المشكلات، نذكر منها تمثيلاً لا حصرًا: فجوات المسؤولية Responsibility Gaps<sup>(1)</sup>، أو مسائل الشفافية والثقة، أو الإحلال الوظيفي، أو ضمور الملكة الأخلاقية Moral Deskilling<sup>(2)</sup>، أو صور الظلم الجديدة، أو اعتبارات الخصوصية، أو إمكانيات التلاعب، أو معضلة المواءمة The Alignment Problem<sup>(3)</sup>. أما المنظورات الإيجابية فتبدو قليلة بل شحيحة.

على الرغم من أن الذكاء الاصطناعي ينطوي على ممكناً هائلة ونافعة للإنسان، فإن عشرات السبل التي في وسعها جعل العالم أعدل وأفضل تظل، في أدبيات الحقل، عبارات عابرة تستدركون عادةً بـ"لكن" فلقة، تمثل مهادأً لحديث مطول عن المخاطر وضرورات الضبط والتشريع الأخلاقي والقانوني. ربما كان من الحري بوعود التقنية إن لم يلهم حماساً خالصاً، فعلى الأقل، أن يلهم قدرًا من التفاؤل والثقة الحذرین، بيد أن المزاج السائد في حقل أخلاقيات الذكاء الاصطناعي يوصف على لسان أحد أبرز الباحثين بأنه "مدٌ متضادٌ من الذعر"<sup>(4)</sup>، وذلك تبعاً لما يستظره عنوان معتبر لمقدمة عدد خاص لأحد الدوريات، نشر مؤخرًا، حيث يقدم مراجعة لحالة هذا الحقل: "أخلاقيات الذكاء الاصطناعي: إشكاليات متفاقمة، إشكاليات متعددة، إشكاليات غير

(1) مصطلح صاغه الفيلسوف أندرياس ماتياس في سياق الإشكاليات التي تطرحها أنظمة التعليم الآلي، حيث لا يستطيع البرمج التنبؤ بسلوكها المستقبلي. ويشير إلى إشكالية عدم إمكانية إسناد المسؤولية الأخلاقية أو القانونية إلى أي شخص عن الأصوات التي تسببها الأنظمة ذاتية التعلم. (المترجم). ينظر:

Andreas Matthias, "The Responsibility Gap: Ascribing Responsibility for the Actions of Learning Automata," *Ethics Inf Technol*, vol. 6 (2004), pp. 175–183.

(2) مصطلح طورته الفيلسوفة شانون فالور، مستلهمة من النقاشات الاجتماعية حول ضمور الملكات الاقتصادية في القرن العشرين، ويشير المصطلح إلى تراجع القدرة على اتخاذ القرارات الأخلاقية بسبب قلة الممارسة والتجربة، وذلك نتيجة تقويض عمليات صنع القرار لتقنيات الذكاء الاصطناعي. (المترجم)، ينظر:

Shannon Vallor, "Moral Deskilling and Upskilling in a New Machine Age: Reflections on the Ambiguous Future of Character," *Philos Technol*, vol. 28 (2015), pp. 107–124.

(3) أحد أبرز الإشكاليات التي تطرحها عمليات الأتمتة في مجال أبحاث سلامة الذكاء الاصطناعي، ويقصد بها التحدي المتمثل في ضمان عمل أنظمة الذكاء الاصطناعي وفقاً للأهداف والقيم والمبادئ الإنسانية، وكيف يمكن الحفاظ عليها في حدود السيطرة الإنسانية [المترجم]. اعتمدنا في ترجمة جميع المصطلحات التقنية على: معجم البيانات والذكاء الاصطناعي إنجليزي - عربي، الهيئة السعودية للبيانات والذكاء الاصطناعي (2024)، ط 2، شوهد في 1/2/2026، في:

<https://acr.ps/1L9F2ur>

(4) John Danaher, "The Rise of the Robots and the Crisis of Moral Patience," *AI & Society*, vol. 34 (2019), p. 129.

مسبقة<sup>(5)</sup>. وبختصار، فإن أخلاقيات الذكاء الاصطناعي تمحور بكليتها حول الإشكاليات المنشقة منه<sup>(6)</sup>.

بطبيعة الحال، هذا لا يعني غياب مقاربات متباينة، فشمة كتابات عدّة، أكاديمية وغير أكاديمية، تُدافع عن الوجه المشرق للذكاء الاصطناعي، نذكر على سبيل المثال: "آلات مفعمة بالنعمة والمحبة" لداريو أمودي<sup>(7)</sup>، و"مانفيستو التفاؤل التقني" لمارك أندريسن<sup>(8)</sup>، ومن التشاوؤم إلى الوعد ليال أورورا<sup>(9)</sup>، والفاعلية الفائقة لرايد هوفمان وكريغ بيتيو<sup>(10)</sup>، آلة المساواة العادلة لأورلي لوفل<sup>(11)</sup>. غير أن أيّاً من هؤلاء ليس فيلسوفاً أكاديمياً، ومن يفتّش عن أسباب وجيهة للتفاؤل بالذكاء الاصطناعي سيُعثّر عليها، في غالب الأمر، خارج الفلسفة الأكاديمية. ففكرة أخلاقيات إيجابية للذكاء الاصطناعي ترد أحياناً، لكنها تظلّ حاضرةً على استحياء، ولا تكاد تتحلّ موقعًا بُؤرّويًّا في الحقل<sup>(12)</sup>.

يبتغي هذا البحث مسألة النزعة السلبية التي تخيم على جماعة الباحثين في أخلاقيات الذكاء الاصطناعي. وسائلك لتحقيق ذلك طريقةً غير مباشرة. فلنأتني، على طريقة "المتفائلين بالتقنية"،

(5) لست هنا بقصد انتقاد لوتشيانو فلوريدي؛ إذ هو محظٌ في تلخيص أبيات أخلاقيات الذكاء الاصطناعي بأنها تدور في معظمها حول المشكلات؛

Luciano Floridi, "Introduction to The Special Issues: The Ethics of Artificial Intelligence," *American Philosophical Quarterly*, vol. 61 (2024).

(6) يمكن الإحساس بالنزعة السلبية في كثير من النقاشات المابعدية Meta-Discussions حول أخلاقيات الذكاء الاصطناعي بوصفها مشروعًا فلسفياً، حيث تميل إلى إبراز المخاطر والجوانب السلبية، ينظر:

Thilo Hagendorff, "Blind Spots in AI Ethics," *AI and Ethics*, vol. 2 (2022); Luke Munn, "The Uselessness of AI Ethics," *AI and Ethics*, vol. 3 (2023);

وبطريقة كافية، يحدّ ثيلو هاغندورف "القوّة الفعلية" لأنّ أخلاقيات الذكاء الاصطناعي في "حساسيتها تجاه الآذى والمعاناة، وقدرتها على رصد الآثار الخارجية [السلبية]"، ينظر: Hagendorff, p. 862.

(7) Dario Amodei, "Machines of Loving Gracem," 13/10/2024, accessed on 19/1/2025, at: <https://acr.ps/1L9F2Pk>

(8) Marc Andreessen, "The Techno–Optimist Manifesto," *Andreessen Horowitz*, 16/10/2023, at: <https://acr.ps/1L9F2sv>

(9) Payal Arora, *From Pessimism to Promise: Lessons from the Global South on Designing Inclusive Tech* (Cambridge, MA: MIT Press, 2024).

(10) Reid Hoffman & Greg Beato, *Superagency: What Could Possibly Go Right with Our AI Future* (New York: Authors Equity, 2025).

(11) Orly Lobel, *The Equality Machine: Harnessing Digital Technology for a Brighter, more Inclusive Future* (London: Hachette UK, 2022).

(12) حول إمكان صوغ أخلاقيات إيجابية للذكاء الاصطناعي، ينظر:

Sven Nyholm, "What is This Thing Called the Ethics of AI and What Calls for It?" in: David J. Gunkel (ed.), *Handbook on the Ethics of Artificial Intelligence* (Cheltenham: Edward Elgar Publishing, 2024);

وحتى الكتب الفلسفية التي تتضمّن كلمة "يوتوبيا" في عنوانها وتطرح تصوّراتٍ طوباوية للتقنية، تظلّ قراءتها ملتبسةً ومزدوجة الانطباع، ينظر:

Nick Bostrom, *Deep Utopia: Life and Meaning in a Solved World* (Ideapress, 2024); John Danaher, *Automation and Utopia: Human Flourishing in a World without Work* (Cambridge, MA: Harvard University Press, 2019); John Danaher, "Techno–optimism: An Analysis, an Evaluation and a Modest Defense," *Philosophy & Technology*, vol. 35 (2022).

قائمة الوعود التي يمكن للذكاء الاصطناعي أن يتحققها، ولن أطيل الوقوف عند كلّ هاجسٍ أخلاقيٍ محدِّدٍ في الأديبيات. سأنظر، عوض ذلك، إلى المزاج التشاوُمي من منظور فلسفة العلم، مما يتضمن خطوةً إلى الوراء لفحص البنية المؤسِّسية لحقل أخلاقيات الذكاء الاصطناعي باعتباره تخصصاً أكاديمياً. سأقترح أن ثمة قيوداً مؤسِّسية تدفع العاملين في الحقل إلى الإلحاد على الجانب السلبي؛ إذ يتعمّن عليهم إبراز مشكلات التقنية، وإلا عرّض نفسه لخطر التهميش. ومن ثمّ، جاءت الصورة التي يرسمها الحقل عن التداعيات الأخلاقية للذكاء الاصطناعي مُختلةً على مستوىين: أحادية الوجه (إذ نادرًا ما تناوش الجوانب الإيجابية) ومحايدة سلبياً (حيث يُبالغ في تضخيم السلبيات). بناءً عليه، يجدر بنا أن نشكّك في السردية التشاوُمية السائدة، وأن نفكّر جديًّا في سُبل إصلاح هذا الحقل حتى يستعيد توازنه.

ينقسم البحث على النحو الآتي: أبدأ بتشخيص ثلاث سمات مؤسِّسية في حقل أخلاقيات الذكاء الاصطناعي تفسر مجتمعة المزاج السلبي السائد فيه (القسم الثاني). ثم أبين أن هذا التفسير يمنحك مسوغًا للقول إن الصورة الفاتمة التي يرسمها أصحاب الاختصاص أحادية الجانب ومحايدة سلبياً (القسم الثالث). ولإيضاح الحاجة وإحاطتها بسياق أوسع، أفارن أزمة السلبية في أخلاقيات الذكاء الاصطناعي بأزمة التكرار التي هزت حقوقًا علمية أخرى (القسم الرابع). ويعالج القسمان المواليان اعتراضين محتملين يشيران، كُلُّ بطريقته، إلى أن النزعة السلبية قد تكون حميدة (القسمان الخامس والسادس). أما الخاتمة فستعرض الدلالة الأوسع للطرح الذي أقدمه (القسم السابع).

يتَّسق مشروعني في هذه الورقة مع المشروع الأوسع في فلسفة العلم؛ الذي يسعى إلى استجلاء الكيفية التي ينبغي أن يُنظَّم بها البحث العلمي اجتماعيًّا كي يُكَلِّل بالنجاح. يتعامل هذا البرنامج البشري مع البحث العلمي بما هو مسعى جماعي يتوقف تقدُّمه على جملة من المعايير والمؤسسات والسياسات، الظاهرة منها والمستضمرة، التي تشكّله وتتصوّره. وبعبارة فيليب كيتشر البرامجية Programmatic، فإن السؤال المطروح هو: "ما السبيل الأمثل لتصميم المؤسسات الاجتماعية بصورة تُعزّز تقدُّم المعرفة؟"<sup>(13)</sup>، وقد دارت الإسهامات في هذا البرنامج حول مسائل شتى، نذكر منها تمثيلاً لا تحديداً: التوزيع الأمثل للجهود البحثية، وأنظمة المكافأة داخل المبحث العلمي<sup>(14)</sup>، والبني

(13) Philip Kitcher, "The Division of Cognitive Labor," *The Journal of Philosophy*, vol. 87 (1990), p. 22.

(14) ينظر:

Ibid.; Miriam Solomon, *Social Empiricism* (Cambridge, MA: MIT Press, 2001); Michael Strevens, "The Role of the Priority Rule in Science," *The Journal of Philosophy*, vol. 100 (2003); Michael Weisberg & Ryan Muldoon, "Epistemic Landscapes and the Division of Cognitive Labor," *Philosophy of Science*, vol. 76 (2009); Kevin J. S. Zollman, "The Epistemic Benefit of Transient Diversity," *Erkenntnis*, vol. 72 (2010).

ال التواصلية داخل الجماعات المعرفية<sup>(15)</sup>، والمعايير التي يتعين أن تحكم البحث العلمي<sup>(16)</sup>، أو الطريقة التي ينبغي أن تُنظم بها عملية مراجعة الأقران Peer Review وتمويل البحث<sup>(17)</sup>.

ومع أن الكثير من الإسهامات في هذا البرنامج البحثي تتسم بالعمومية؛ إذ تبحث في الكيفية التي يتعين أن يُنظم بها البحث العلمي عامةً لبلوغ غاياته، فإن مشروع هذا المقال أكثر تحديداً، فهو يحصر النظر في ميدانٍ فلسفـي بعينه هو "أخلاقيات الذكاء الاصطناعي"، ومـرـد ذلك أنه يحلـل الدينامـيات المؤسـسـية والإشكـاليـات المرتبـطة بهاـ التي أحـسـبـها خـاصـةـ بـهـذاـ الحـقـلـ. ومعـ ذـلـكـ، فـمـنـ الجـائزـ أنـ تـعمـمـ بعضـ الإـشـكـاليـاتـ المـحدـدةـ فيـ أـخـلـاقـيـاتـ الذـكـاءـ الـاـصـطـنـاعـيـ علىـ أـخـلـاقـيـاتـ التقـنـيـةـ بـمـعـناـهاـ الـأـوـسـعـ.

## الديناميكـياتـ المؤسـسـيةـ فيـ أـخـلـاقـيـاتـ الذـكـاءـ الـاـصـطـنـاعـيـ

يمـكـنـ تـفـسـيرـ النـزـعـةـ السـلـبـيةـ فيـ أـخـلـاقـيـاتـ الذـكـاءـ الـاـصـطـنـاعـيـ بـسـمـاتـ مؤـسـسـيـةـ مـعـخـصـوصـةـ تـميـزـ هـذـاـ التـخـصـصـ الأـكـادـيمـيـ، وـيرـجـعـ ذـلـكـ إـلـىـ حدـ بـعـيدـ إـلـىـ الـقيـودـ الدـاخـلـيـةـ لـلـحـقـلـ الـتـيـ تـكـرـهـ الـبـاحـثـينـ عـلـىـ تـسـلـيـطـ الضـوءـ عـلـىـ الـمـسـاوـيـ الـمـحـتمـلـةـ. وـعـلـىـ نـحـوـ أـدـقـ، يـنـشـأـ هـذـاـ الـانـحـيـازـ مـنـ تـقـاعـلـ ثـلـاثـ سـمـاتـ مؤـسـسـيـةـ أـسـمـيـهاـ: "مـوـضـوـعـ الـبـحـثـ"ـ Subـject~Matterـ، وـ"ـتـأـثـيرـ الإـيجـابـيـ"ـ Positive Impactـ، وـ"ـالـحـوـافـزـ"ـ Incentivesـ.

أـوـلـاـ، فـلـنـلـتـفـتـ إـلـىـ "ـمـوـضـوـعـ الـبـحـثـ"ـ وـيـسـبـبـ الطـبـيـعـةـ الـخـاصـةـ لـأـخـلـاقـيـاتـ الذـكـاءـ الـاـصـطـنـاعـيـ، يـخـتـلـفـ مـسـارـ الـبـحـثـ فـيـهـاـ عـمـاـ نـجـدـهـ فـيـ سـائـرـ الـحـقـولـ الـفـلـسـفـيـةـ. فـفـيـ تـلـكـ الـحـقـولـ تـكـونـ الدـافـعـيـةـ الـتـيـ تـشـرـعـ مـسـعـىـ فـلـسـفـيـاـ، فـيـ الـغـالـبـ، سـؤـالـاـ فـلـسـفـيـاـ لـمـ يـحـسـمـ بـعـدـ، وـغالـبـاـ ماـ يـكـونـ لـغـرـاـ عـتـيقـاـ دـارـتـ حـولـهـ الـجـدـالـاتـ قـرـوـنـاـ عـدـيـدـةـ. فـفـيـ الـأـخـلـاقـ، مـثـلـاـ، قـدـ يـكـونـ السـؤـالـ: ماـ الـذـيـ يـجـعـلـ فـعـلـاـ مـاـ صـائـبـاـ أوـ خـاطـئـاـ مـنـ النـاحـيـةـ الـأـخـلـاقـيـةـ؟ـ وـفـيـ الـفـلـسـفـةـ الـسـيـاسـيـةـ قـدـ يـكـونـ السـؤـالـ: ماـ الـعـدـالـةـ؟ـ وـمـتـىـ يـكـونـ تـدـخـلـ الـدـوـلـةـ مـشـرـوـعـاـ؟ـ أـمـاـ نـظـرـيـةـ الـمـعـرـفـةـ فـتـسـأـلـ عـنـ مـاهـيـةـ الـمـعـرـفـةـ وـالـكـيـفـيـاتـ الـمـتـعـلـقـةـ بـإـمـكـانـ تـحـصـيلـهـاـ، بـيـنـمـاـ تـتـمـحـورـ نـظـرـيـةـ الـفـعـلـ Action Theoryـ حـولـ سـؤـالـ إـمـكـانـيـةـ تـحـقـقـ الإـرـادـةـ الـحـرـّـةـ.ـ وـفـيـ حـقـلـ الـمـيـتـافـيـزـيـقاـ قـدـ يـطـرـحـ تـسـاؤـلـ عـنـ مـبـدـأـ الـعـلـيـةـ أوـ عـنـ مـاهـيـةـ الـزـمـنـ، أـمـاـ فـلـاسـفـةـ الـعـقـلـ فـيـسـعـونـ إـلـىـ فـهـمـ عـلـاقـةـ الـعـقـلـ بـالـمـادـةـ.ـ وـهـكـذـاـ يـتـحـدـدـ القـاسـمـ الـمـشـترـكـ بـيـنـ هـذـهـ الـمـبـاحـثـ فـيـ أـنـ دـافـعـ الـبـحـثـ الـفـلـسـفـيـ فـيـهـاـ سـؤـالـ مـفـتوـحـ يـلـزـمـ الـفـلـاسـفـةـ بـالـتـمـاسـ جـوابـ ماـ.ـ وـالـمـحـصـلـةـ الـمـتـوـخـّـةـ، فـيـ أـحـسـنـ الـأـحـوالـ، تـمـثـلـ

(15) Kevin J. S. Zollman, "The Communication Structure of Epistemic Communities," *Philosophy of Science*, vol. 74 (2007).

(16) Helen E. Longino, *Science as Social Knowledge: Values and Objectivity in Scientific Inquiry* (Princeton: Princeton University Press, 1990), Ch. 4; Helen E. Longino, *The Fate of Knowledge* (Princeton: Princeton University Press, 2002), Ch 6; Liam Kofi Bright & Remco Heesen, "To Be Scientific Is to Be Communist," *Social Epistemology*, vol. 37 (2023).

(17) Marcus Arvan, Liam Kofi Bright & Remco Heesen, "Jury Theorems for Peer Review," *British Journal for the Philosophy of Science*, vol. 76 (2025); Shahar Avin, "Centralized Funding and Epistemic Exploration," *The British Journal for the Philosophy of Science*, vol. 70 (2019); Remco Heesen & Liam Kofi Bright, "Is Peer Review a Good Idea?" *The British Journal for the Philosophy of Science*, vol. 7 (2021); Carole J. Lee, "Commensuration Bias in Peer Review," *Philosophy of Science*, vol. 82 (2015).

في صياغة جواب [مُحدد] عن هذا السؤال: فيقترح فلاسفة الأخلاق نظريات للحكم الخلقي، ويضع فلاسفة السياسة نظريات حول العدالة أو السلطة، ويطور فلاسفة العلم نماذج للمعرفة والتبرير، ويشرح منظرو حقل نظرية الفعل إمكان الإرادة الحرة أو ينقضونه، إلى غير ذلك من الأمثلة. لكن ثمة استثناءات، فالحقول أقل تجانساً مما تُوحى به هذه الصورة المبسطة، ومع ذلك تظل الصورة التي قدمتها وافية بالنمط الغالب. بل إن الأعمال التي لا تنضوي ظاهراً تحت هذا الوصف، مثل الاستغلالات المنهجية أو المفهومية، يقصد بها، في العادة، إسهاماً غير مباشر في الإجابة عن تلکم الأسئلة<sup>(18)</sup>.

يختلف المشهد في حقل أخلاقيات الذكاء الاصطناعي، وعموم حقل أخلاقيات التقنية عن غيره، وذلك لجهة أن موضوع الحقل، بدهاً، هو التقنيات ذاتها. هنا لا ينبع الدافع إلى البحث من سؤالٍ فلسيٍ معلق، بل من حدثٍ أو تطورٍ محدّد، متمثلاً بظهورٍ (أو توقعٍ ظهور) تقنية ذكاء اصطناعي جديدة. لذلك لا يُدعى فلاسفة أخلاق الذكاء الاصطناعي إلى حلّ معضلاتٍ فلسفية قائمة، بل إلى التأمل في تطور تقنيات جديدة والتعليق عليها. فإذا كان "الناتج" المتوقع للبحث الفلسي هو الإجابة عن سؤالٍ فلسيٍ مفتوح في التخصصات الفلسفية الأخرى، فإن "الناتج" المتوقع في حقل أخلاقيات الذكاء الاصطناعي هو تأملٌ أو تعليقٌ أخلاقي حول ظهور تقنية جديدة.

تجعل السمة المؤسسية الثانية؛ وهي "التأثير الإيجابي" Positive Impact، من العسير على فلاسفة الأخلاق في حقل الذكاء الاصطناعي أن يعلّقوا على التقنيات بروح متفائلة. يشير مصطلح "التأثير الإيجابي" إلى قاعدة غير مكتوبة تحكم التنظير الأخلاقي وتميل نوع المخرجات التي يفترض بالباحثين إنتاجها. يبدو، على وجه التحديد، أن هناك قاعدة تمنع من ملاحظة أن حدثاً أو تطوراً بعينه قد أفضى، أو سيفضي، إلى آثار محمودة أخلاقياً. وللتوضيح دعونا ننتقل من "أخلاقيات الذكاء الاصطناعي" إلى "الفلسفة السياسية"، بادئين بالملاحظة الآتية: إن أحوال المعيشة عالمياً تتحسن بوتيرة لافتة، فالل黍ن ووفيات الأطفال في أدنى مستوياتها التاريخية، ومعدلات القراءة والتعليم النظامي في أعلى مستوياتها، بل إن مظاهر انعدام المساواة العالمي نفسه يتراجع<sup>(19)</sup>. هذه التطورات مرغوبة أخلاقياً. ومع ذلك أزعم أن الحجج الفلسفية التي تصاغ بالطريقة التالية ستبدو مُحرجة إلى حد بعيد:

- أحاجٌ في هذه الورقة بأن التراجع العالمي في معدلات الفقر يُفضي إلى زيادة ملحوظة في العدالة.

(18) مع ذلك، أحسب أن أحد المراجعين كان محقاً في تبيهه إلى أن بعض صور الاستغلال، كالنقد المفاهيمي، وبناء الأطر، والتحليل الجينيولوجي، قد لا تسجم مع هذا النمط. وعلى الرغم من أن كثيراً من هذا الضرب من الاستغلال لا يزال يرمي، ولو بصورة غير مباشرة، إلى معالجة أنماط الأسئلة المذكورة آنفاً، فإن بعضه لا يفعل ذلك. ولا أعتقد أن هذا يُفْوِض حججَي؛ إذ لا تقتضي الأخيرة سوى وجود نزعة عامة قوية من هذا القبيل، لا قاعدة مُطردة لا استثناء فيها.

(19) Max Roser, "The Short History of Global Living Conditions and Why It Matters That We Know It," *Our World in Data* (2016), accessed on 19/1/2026, at: <https://acr.ps/1LF2D7>; Max Roser, "The History of Global Economic Inequality," *Our World in Data* (2017), accessed on 19/1/2026, at: <https://acr.ps/1L9F2x0>

- أزعم في هذه الورقة أن تزايد التحاق الأفراد بالتعليم النظامي يترجم إلى ارتفاع كبير في العدالة التعليمية Educational Justice.

- أجادل في هذه الورقة بأن الانخفاض العالمي في مستويات اللامساواة يعني أن العالم أصبح أكثر عدلاً.

على افتراض أن الفرضيات التجريبية صحيحة، فإن كلّ واحدة من هذه المزاعم صحيحة ومهمة. ومع ذلك، فإن مقالاتٍ تشيد حججها على هذا النحو ستبدو، فيما أزعم، شاذة إلى حدّ ما، ومن غير المرجح أن تتجاوز التحكيم الأكاديمي. فثمة، على ما يبدو، قاعدةٌ عامَّةٌ تقضي بأن الباحثين في الأخلاقيات لا ينبغي أن يكتفوا بمجرد تسجيل الآثار المرغوبة أخلاقياً لحدثٍ أو تطويرٍ معينه. وينطبق هذا أيضاً على أخلاقيات الذكاء الاصطناعي؛ إذ ستغدو مقالاتٌ تتقدّم بحجج من النوع الآتي خرقاً في أعين المراجعين:

- أجادل في هذه الورقة بأن "تشات جي بي تي" ChatGPT، من خلال إتاحة الموارد التعليمية بتكلفة زهيدة للجميع، يُعدّ أمراً محموداً من منظور العدالة التعليمية.

- أجادل في هذه الورقة بأن السيارات ذاتية القيادة مفيدة من زاوية العدالة الاجتماعية؛ لأن الأشخاص الأقل دخلاً سيجنون الفائدة الأكبر من خفض تكاليف التنقل.

يبدو أن ثمة خللاً ما في المقالات الفلسفية التي تشيد حججها على هذا النحو، فبصرف النظر عن صحة هذه المزاعم، لا تُعدّ من الطراز الذي يفترض بالباحث الأخلاقي أن يتقدّم بها. صحيحٌ أن مثل هذه الادعاءات قد تظهر جزءاً من حججٍ أوسع، مثل الاعتراض على التحكم في تلك التقنيات أو تأييد دعمها ماليًّا، ويمكن كذلك الدفع بها ردًا على حجج مضادة تزعم أن تشات جي بي تي أو القيادة الذاتية غير عادلة. غير أن الاكتفاء بملحوظة أن حدثاً أو تطوراً، كظهور تقنية جديدة، له آثارٌ أخلاقية مستحسنة لا يندرج ضمن ما هو متوقع من باحثي الأخلاق فعله، فمقال من طراز "ثمة تطبيق ذكاء اصطناعي جديد، وهو ممتاز" لا يكون مقبولاً، بينما يُعدّ من المقبول أن يُعلن مقال "ثمة تطبيق جديد، وهو مثيرٌ للقلق الشديد".

تأمل الافتتاح النموذجي لمقال صدر حديثاً عن "التعسُّف الهيرومنطيقي" المرتبط بالذكاء الاصطناعي: "كشفت أدبيات أخلاقيات الذكاء الاصطناعي عن أشكال كثيرة من الأذى الذي تسببه أو تعززه هذه التقنية [...] غير أن شكلاً واحداً فاته الرصد"<sup>(20)</sup>. لاحظ مدى سلاسة هذا التأطير. في المقابل، سيكون

(20) Andrew P. Rebera, Lode Lauwaert & Ann-Katrien Oimann, "Hidden Risks: Artificial Intelligence and Hermeneutic Harm," *Minds & Machines*, vol. 35, no. 33 (2025), p. 2;

في هذا السياق، ثمة مقالة أخرى في الموضوع نفسه وبالتالي عينه تقول: "بات أمرًا ثابتاً أن الخوارزميات قد تكون أدوات للظلم، غير أن ما يُناقض على نحو أقل بكثير هو أن طرائق نشر الذكاء الاصطناعي الراهنة تجعل اكتشاف الظلم نفسه أمراً عسيراً، إن لم يكن مستحيلاً. [...] يُبين كيف يمكن لاستخلاص الخوارزمي Algorithmic Profiling أن يُولد ظلماً إبستيمياً، ينظر:

Silvia Milano & Carina Prunkl, "Algorithmic Profiling as a Source of Hermeneutical Injustice," *Philosophical Studies*, vol. 182 (2025), p. 186;

وأوّد التنويه إلى أنني لست هنا في معرض نقدهؤلاء المؤلفين.

غريباً أن نصادف ورقة في أخلاقيات الذكاء الاصطناعي ذات نزعة إثباتية بدلاً من تقديرية، تعرف بفوائد التقنية المتعددة وتمضي إلى توثيق فاندبة لم يلتفت إليها من قبل. ومع أن الذكاء الاصطناعي يجعل العالم أكثر عدلاً بطرق لا حصر لها، لا يُتَّسِّرُ من الباحثين في حقل أخلاقيات الذكاء الاصطناعي توثيق هذه المنافع، بل يُتَّسِّرُ منهم التركيز على السبل التي يهدّد بها الذكاء الاصطناعي القيم الأخلاقية، فضلاً عن العناية بكيفية درء تلك التهديدات.

إن القول بوجود افتراضٍ مُسبقاً يعترض طريق المقالات الإيجابية التي تقتصر على الملاحظة يظلّ، في جوهره، طرحاً ظنياً. غير أن هناك مسوّغاً إضافياً لاعتقاد وجود مثل هذا الافتراض، يتتجاوز مجرد الانطباع بوجود عوار في تلك المقالات، ذلك أن الافتراض في ذاته معقول. فملاحظة أن حدثاً أو تطوراً بعينه محمودٌ أخلاقياً قد ينطوي على قيمةٍ معرفية بحد ذاته، أو على الأقل تقدير، عندما لا تكون الملاحظات بديهية. لكن هذه الملاحظات، في الغالب، لا تحمل أثراً عملياً كبيراً، فهي لا تمنحك سوى باعث على الابتهاج. قارن ذلك بمقالات تشير إلى أن حدثاً أو تطوراً ما يستعمل على جوانب إشكالية، مثل هذه المقالات لا تكون بمعناها للازم عاج وحسب، بل تمنحك سبيلاً للتدخل *Intervene*، فهي، بصورة مباشرة أو غير مباشرة، دعواتٌ إلى الفعل. واستصحاباً تنشأ لاتماضية *Asymmetry* في القوة الإلزامية بين التقديرات الإيجابية والسلبية للذكاء الاصطناعي. وبما أن توقيع كتابة أوراق ذات صلةٍ إجرائية ليس أمراً غير معقول، فمن غير المستبعد وجود قاعدةٍ تردّ هذه المقالات الإيجابية الخالصة ذات الحمولة التقريرية الضئيلة.

إذاً، يعني بـ"موضوع البحث" أن متخصصي أخلاقيات الذكاء الاصطناعي منشغلون أساساً بالتعليق على التقنيات الجديدة، لا بحل الألغاز الفلسفية القائمة. ويعني بـ"الأثر الإيجابي" أنه يُحظر عليهم الاشتغال بالجوانب الإيجابية له. أما العامل المؤسسي الثالث "الحوافر"، فيشير إلى حقيقة أن الأكاديميين مضطرون إلى النشر لحفظهم على مسارهم المهني، فالتوقف عن كتابة أوراق حوله ليس خياراً متاحاً لباحثي هذا العقل<sup>(21)</sup>.

يفسر اقتران هذه العوامل الثلاثة تفسيراً جيداً لهذا المد الصاعد من الذعر. فباحثو أخلاقيات الذكاء الاصطناعي مضطرون إلى الشر، والطريق المتاح أمامهم هو تضخيمُ الهواجس الأخلاقية المتصلة بالتقنية. وينسحب الأمرُ عينه على تحصيل المنح البحثية الخارجية؛ إذ لا يستطيع الباحث التقدّم لمنحة، وهي كثيراً ما تكون حاسمةً في الترقية الأكاديمية، إلا إذا أطّرَ الذكاء الاصطناعي بوصفه مولدًا لمشكلاتٍ أخلاقية تستوجب البحث. ولنقل بوجيز العبارة: الباحثون في هذا المجال ملزمون بالتأطير السليبي من أجل المحافظة على مسارهم الوظيفي، أما أولئك الذين يؤمّنون بغير ذلك، فمهدوّن بفقدان وظائفهم سريعاً.

(21) للاطلاع على مناقشات ذات صلة بهياكل الحوافر الإشكالية في العلم، ينظر:

Wesley Buckwalter, "The Replication Crisis and Philosophy," *Philosophy and the Mind Sciences*, vol. 3 (2022); Remco Heesen, "Why the Reward Structure of Science Makes Reproducibility Problems Inevitable," *The Journal of Philosophy*, vol. 115 (2018).

يحمل هذا العرض شيئاً من المبالغة؛ فلا ينحصر اشتغال هؤلاء المتخصصين في إثارة المخاوف الأخلاقية؛ لذلك ينبغي تقييد هذا الزعم بوجهه ثلاثة: أولاً، في وسع المتخصصين أن يقتربوا حلولاً للمشكلات التي يطرونها، مثل معالجة فجوات المسؤولية، أو البطالة التي تُحدثها التقنيات، أو ضمور الملكة الأخلاقية. ثانياً، يمكن أن يناقش الباحثون في مقالاتهم السجالية Response Pieces<sup>(22)</sup> وجود هذه المشكلات المزعومة أو مدى خطورتها، فيذهب أحدهم إلى أن الذكاء الاصطناعي لا يفضي فعلياً إلى حدوث فجوات المسؤولية أو البطالة أو ضمور الملكة الأخلاقية، معتبراً تلك الهواجس من قبيل المبالغات أو أنها زائفة. وأخيراً، هناك موضوعات في الحقل لا تتمحور أساساً حول كونه مصدر قلق أخلاقي، مثل الجدل حول الوضعية الأخلاقية والوكالية للنظم الاصطناعية the Moral and Agential Status of Artificial Systems، أو السؤال عما إذا كان في إمكان هذه النظم أن تكون أصدقاء أو شركاء عاطفيين، فضلاً عن التأملات الأوسع عن الظرف الإنساني Condition Humaine في عصر الذكاء الاصطناعي.

مع ذلك، لا تقاد هذه الاستثناءات تخففاً كثيراً من الديناميات الموصوفة آنفًا. فاقتراح حلول لمشكلات مرتبطة بالذكاء الاصطناعي يفترض سلفاً وجود تلك المشكلات، ومن ثم يؤكّد السردية السلبية بصورة غير مباشرة. أمّا المقالات السجالية، فتقلّ عادةً من حيث الحضور والتأثير، فإذا صحّ أن بعضها يحقق رواجاً لافتاً، فإن ذلك واقع في خانة الاستثناء لا القاعدة<sup>(23)</sup>. وعلى الرغم من وجود عدد محدود من القضايا في حقل أخلاقيات الذكاء الاصطناعي التي لا تتمحور حول المخاوف أو المعضلات، فمن باب المبالغة القول إن الباحثين المستغلين في هذا الحقل لا يختار لهم سوى رسم صورة فاتمة للتقنية، إلا أن النقطة الأعم، هي وجود إكراهات مؤسسية قوية تدفع في هذا الاتجاه. يقترح هذا المقال، إذًا، ميلًا عاماً واضحاً نحو النظرة السلبية، مع الإقرار بوجود بعض الاستثناءات.

يبدو أن كلّ عامل من العوامل الثلاثة السابقة ضروري لتفسير التزعة التشاوؤمية. فلو كانت أخلاقيات الذكاء الاصطناعي، من حيث موضوع البحث، أقرب إلى سائر الحقول الفلسفية، لما انحصر عمل المتخصصين فيها في التعليق على الأحداث أو التطورات<sup>(24)</sup>. ولو كان من المقبول نشر مقالاتٍ

(22) مقالات أكاديمية قصيرة غايتها المباشرة الاستباق مع دراسة سابقة؛ فهي تأتي على صورة تعليق أو نقد أو توضيح أو تطوير لحجج سبق طرحها، من دون أن تُقدم أطروحة جديدة. هذا اللون حاضر في التقليد الباحثي الغربي ويُكاد يكون مفقوداً في الكتابة الأكاديمية العربية. (المترجم)

(23) وتبلغ فرص نجاحها أقصاها حين تُفضي المشكلة إلى نقاش واسع، غير أن هذا لا يصدق إلا على شريحة ضئيلة من القضايا المتصلة بالذكاء الاصطناعي التي تتراولها أدبيات الحقل.

(24) وقد ذهب مايكيل هيوم على أنس مماثلة إلى أن المستغلين بالأخلاقيات يُنجزون بأخلاقيات الذكاء الاصطناعي وحدهم، يُضخّمون المشكلات الأخلاقية. وقد يصحّ هذا القول، غير أن المشكلة تبدو أكثر حدةً داخل أخلاقيات الذكاء الاصطناعي بسبب بنية المؤسسية الخاصة. إذ يتوجّأ من المستغلين في هذا الحقل أن يُعلّموا على التطّورات التقنية، بينما يُحظر عليهم إلى حد بعيد قول أشياء إيجابية عنها. وهذا ما يُؤلّد ضغطاً للعثور على مثالب في هذه التطّورات. وقد يشعر المستغلون بالأخلاق في حقول أخرى، أحياناً، مماثلاً نحو تقديم أمور بريئة بوصفها إشكالية، كما أشار هيوم. بيد أن لديهم أيضاً خيار الكتابة حول الأسئلة الكبرى للفلسفة الأخلاقية (هل يعني لنا دائمًا تعظيم الخير؟ هل الأخلاق موضوعية؟ ما الفضائل الموجودة؟ هل ثمة مبادئ أخلاقية عامة؟ ... إلخ)، عوضاً عن التعليق على الأحداث أو التطورات الراهنة. وهذا ما يمنّ لهم فرصاً أوسع لإجراء بحوث أخلاقية دون اضطرارهم إلى تقديم الأشياء (التطّورات التقنية والاجتماعية والسياسية ... إلخ) بوصفها مُقلقة أو إشكالية، ينظر:

تُشيد بروعة الذكاء الاصطناعي، لاستطاع الباحثون تجنب تغذية السردية القاتمة عبر ضرب من المقالات المتفائلة. ولو لم يكن إنتاج مخرّجات بحثية محددة شرطاً مهنياً، لكن في مقدورهم تقليص عدد الكتابات النقدية أو الكف عنها كلية. غير أن هذه المقدمات غير موجودة كلها. ونتيجةً لذلك، لا يكاد يبقى أمام الباحثين المشغلين في هذا الحقل سوى الكتابة عن الإشكالات الأخلاقية التي يطرحها الذكاء الاصطناعي. إن اجتماع العوامل الثلاثة يفسّر سلبية المجال، ومع رسوخ هذه الديناميات المؤسسية وازدهار الحقل المذكور، فإن موجة الذعر مرشحة لا محالة لمزيد من الارتفاع. ومن العسير أن نتصوّر كيف يمكن أن تُصبح أخلاقيات الذكاء الاصطناعي يوماً ما حقلًا ذا نبرةٍ متفائلة تجاه التقنية.

لا أفترض بهذا المقترن التفسيري تحقق هذه السمات الثلاث في صيغ مثالية؛ إذ أعدّها أنماطاً عامة ومويلاً سائدة، لا قوانين مطلقة لا تعرف استثناء، لكنها في تضافرها وشيوعها كافية لتفسير التزعة الشائمية السائدة في الحقل.

## من السلبية إلى الأحادية والتحيز

توفر الديناميات المؤسسية الموصوفة آنفاً تفسيراً للنبرة الشائمية الملحوظة في حقل أخلاقيات الذكاء الاصطناعي. وقد التزمتُ، حتى هذه اللحظة، الحيادَ إزاء ما إذا كان هذا الأمر يشكل مازقاً. وسأبين الآن أنه بالفعل أمرٌ إشكالي. فالصورة التي ترسمها الجماعة العلمية لأنّار الذكاء الاصطناعي الأخلاقية تبدو لي مفرطة في السلبية، وذلك من وجهين: الأول، لأنّها أحاديدُ الجانب One-Sided؛ إذ تميل إلى الجوانب السلبية فحسب، والثاني، لأنّها منحازةٌ Biased سلبياً.

أولاً، بإيجاز، لتأمّل مشكلة الأحادية السلبية. تُبرز الديناميات المؤسسية، التي ذكرناها، الجوانب السلبية على نحو لافت، في حين نادرًا ما تُناقِش الجوانب الإيجابية. والتنتيجة المتأدية عن ذلك تتمثل في أن الصورة العامة التي تشكّل عن الذكاء الاصطناعي تصبح أحاديداً بصورة فادحة. وهي أحاديد مضللة؛ إذ تجعل التقنية تبدو، في المحصلة، أشدَّ إشكاليةً مما هي عليه فعلاً. ولم يُست المشكلة في أن القضايا المُشحّصة متخيّلة، بل المشكلة في التركيز المبالغ فيه على الجوانب السلبية بينما يجري تجاهل الجوانب الإيجابية إلى حدٍ بعيد. حتى لو صدّق أن كل قضية ذُكرت في أدبيات أخلاقيات هذا الحقل قضية حقيقة، فستظلّ الصورة الناتجة من ذلك خادعة؛ إذ يقتضي عرضُ دقيقٍ غير مضلل لأي تقنية موازنةً بين سلبياتها وإيجابياتها معاً. ولا عجب أن يشعر الناس بالذعر إزاء الذكاء الاصطناعي إذا واجهوا في الأساس عيوبه، فأي تقنية ستبدو مقلقة إذا اقتصر النظر على جوانبها المظلمة.

لفهم السبب في أن الأحادية قد تكون إشكالاً، تخيل مسافرةً عادت من إنكلترا ولم تُصحِّح إلّا عن انطباعاتها السلبية، لأنّ تشکو من سوء الطقس، وارتفاع أسعار الفنادق في لندن، وانحطاط المدن ما بعد الصناعية، ورداة شبكة السكك الحديدية، إلى آخره. فلا شيء مما ذكرته باطل من حيث المبدأ؛ فكلّ ملاحظةٍ صحيحةٍ في ذاتها. غير أن حذف كلّ الانطباعات الإيجابية يجعل السردُ مُضللاً؛

إذ بالرغم من الصحة الوصفية للملحوظات، يشكل لدى المستمع انطباعًّا مغلوط عن البلد. إن حقل أخلاقيات الذكاء الاصطناعي، حين يحصر عدسه في المشكلات والمخاطر، يماثل إلى حدّ بعيد رواية تلك المسافرة، ولو بدرجة أقلّ حدّة، فحتى لو افترضنا صحةً كلّ هاجس أخلاقي حول التقنية، ستظلّ الصورةُ العامة معيبةً بأحاديتها، بالضبط كما اختلت رواية المسافرة.

علاوةً على ذلك، هناك ما يدعو إلى افتراض أن الصورة العامة للذكاء الاصطناعي، بفعل الديناميات المؤسسية الموصوفة، ليست أحادية الوجه فحسب، بل مشوّهة على نحو أعمق؛ فالمسألة لا تقتصر على إهمال الجوانب الإيجابية في كثير من الأحيان، بل يتحمل أيضًا أن تُضيّع الجوانب السلبية. وبعبارة أخرى، ينبغي التسليم بأن هذه الديناميات تُنشئ في الحقل انحيازًا تجاه السلبية والتشاؤم، يُمكن فهم التحثّز على أنه ميلٌ منهجيٌّ إلى الانحراف عن الحقيقة في اتجاهٍ بعينه<sup>(25)</sup>؛ إذ يميل خطاب أخلاقيات الذكاء الاصطناعي إلى الابتعاد عن الدقة عبر المبالغة في حجم المشكلات الأخلاقية التي تثيرها التقنية. وافتراضُ أن الديناميات المؤسسية تولّد مثل هذا الانحياز السلبي يبدو وجيهًا لثلاثة أسباب.

يرتبط السبب الأول بما يسمى "قانون الأداة" Law of the Instrument، الذي يفيد أن من لا يملك في عدته سوى مطرقة يرى في كل شيء مسمارًا، فيرى المسامير حتى في الأماكن التي لا توجد فيها. ومن الراجح أن قانونًا مشابهًا يعمل في أخلاقيات الذكاء الاصطناعي، فالعدة الرئيسة، وإن لم تكن الوحيدة، لدى باحث هذا الحقل هي تشخيص المشكلات الأخلاقية. لذلك يُتوقع أن يرى الباحثون مشكلات حيث لا توجد، أو أن يعدوها أفتح مما هي عليه. ينشأ هذا الأثر من العاملين الأول والثاني اللذين سبقت مناقشتهما: فالتفاعل بين "موضوع البحث" و"تأثير الإيجابي" يدفع الباحث إلى تبني زاوية نقدية تلقائياً؛ إذ يستبعد الاستجابات الإيجابية من صندوق أدواته و يجعل النقد الأخلاقي أداته الافتراضية. وحين يتأمل الباحثون تقنيّةً ما يميلون إلى النظر إليها من زاوية سؤال: "ما الذي يمكن أن يكون مقلقاً فيها؟". ومن المعقول افتراض أن زاوية النظر النقدية تدفعهم أحياناً إلى رؤية مشكلات لا وجود لها أو إلى تضخيم الموجود منها بغير مسوغ. وحتى مع استبعاد العامل الثالث المتصل بـ"الحوافز"، يبقى الاعتقاد بأن ثمة إفراطاً في تشخيص المشكلات أمرًا مسوغاً.

يرتبط السبب الثاني بالأول غير أنه يدخل عامل "الحوافز" في الحساب. فالنقد الأخلاقي لا يُعد الأداة الرئيسية في عدّة باحث أخلاقيات الذكاء الاصطناعي وحسب، بل إن هؤلاء يُشجّعون على استعماله؛ إذ يتبعن على الباحثين الانحراف في النقد الأخلاقي إن أرادوا النشر والترقّي في السلم الأكاديمي. ويبدو المبدأ الآتي وجيهًا: متى وُجد حافز يدفع المرء إلى تبني موقفٍ بعينه، مال إلى الدفاع عنه بوتيرة تفوق ما تُجيزه المعطيات المعرفية. فإذا كانت لديك حافز قوية للدفاع عن قضية ما، فستفعل ذلك غالباً حتى عندما لا شعف الأدلة تلك القضية. وبإسقاط هذا على أخلاقيات الذكاء الاصطناعي، يرجح أن الباحث يُشجّع على الزعم أن التقنية مقلقةً أخلاقياً، وسيكتب أوراقاً تشير هذه المخاوف، ولو لم تسنده الشواهد بما يكفي.

(25) Heather Douglas & Kevin C. Elliott, "Addressing the Reproducibility Crisis: A Response to Hudson," *Journal for General Philosophy of Science*, vol. 53 (2022), p. 202.

تجدر الإشارة أيضاً إلى أن مختصي أخلاقيات الذكاء الاصطناعي، وأخلاقيات التقنية عموماً، يشاركون كثيراً في مشاريع كونسورتيوم Consortium<sup>(26)</sup> التي تضم علماء حاسوب ومهندسين، حيث يُعهد إليهم تقديم التوجيه الأخلاقي للمشروع البحثي. ويجد هؤلاء الباحثون أنفسهم، في مثل هذه المشروعات، مضطربين إلى إثارة هموم وقضايا أخلاقية تتعلق بالعمل القائم، فليس من الممكن أن يخاطبوا المهندسين بالقول: "ما تقومون به أمرٌ مذهلٌ حقاً، ولا نرى أي مشكلات [أخلاقية] في مشروعكم، فحظاً طيباً!". مهما كان المشروع رائعاً أو بريئاً، فإن وظيفة باحث الأخلاقيات تفرض عليه أشكالها.

أما السبب الثالث الذي يسُوّغ الظن بأن الديناميات، التي تقدّم ذكرها، قد تفضي إلى تضخيـم المشكلات الأخلاقية المتعلقة بالذكاء الاصطناعي، فإنه يرتبط بما تولده من آثار انتقائية محتملة. فالأشخاص المتفائلون بالتقنية قد يُحـجـمون عن الانخراط في مجال يلزمـهم التركيز على الجوانب السلبية، الحقيقية منها أو المتـوهـمةـ، في حين يـجـدـ المـتـشـكـكـونـ فيـ الذـكـاءـ الـاصـطـنـاعـيـ هذاـ المـجـالـ أكثرـ جـاذـبيـةـ.ـ منـ شـأنـ هـذـاـ الـأـثـرـ الـانـقـائـيـ أنـ يـزـيدـ اختـلالـ الـمـنـظـورـ العـامـ فيـ الـاتـجـاهـ السـلـبـيـ عـلـىـ نحوـ مـقـلـقـ؛ـ إذـ يـفـضـيـ إـلـىـ تـرـكـيـةـ بـحـثـيـةـ غـيرـ مـوـازـنـةـ يـتـسـعـ مـنـهـاـ تـقـيـيمـاتـ مـنـحـازـةـ إـلـىـ إـمـكـانـاتـ التـقـنـيـةـ الـأـخـلـاقـيـةـ،ـ كـمـاـ يـصـرـفـ الفـتـةـ الـقـادـرـةـ عـلـىـ مـواـزـنـةـ الـقوـتـيـنـ السـابـقـتـيـنـ الدـافـعـتـيـنـ إـلـىـ التـشـاؤـمـ عـنـ الـخـوـضـ فـيـ هـذـاـ الـمـجـالـ دونـ إـسـهـامـ.

ثـمـةـ،ـ إـذـاـ،ـ ثـلـاثـةـ أـسـبـابـ تـدـعـونـاـ إـلـىـ الـظـنـ أـنـ الـعـوـاـمـ الـمـؤـسـسـيـ تـدـفعـ جـمـاعـةـ أـخـلـاقـيـاتـ الـذـكـاءـ الـاـصـطـنـاعـيـ إـلـىـ تـضـخـيمـ الـمـخـاـوفـ الـأـخـلـاقـيـةـ الـمـتـعـلـقـةـ بـالـتـقـنـيـةـ إـجـمـالـاـ.ـ وـكـمـاـ أـشـيـرـ آـنـفـاـ،ـ يـمـكـنـ أـنـ يـتـخـذـ هـذـاـ تـضـخـيمـ شـكـلـيـنـ مـخـلـفـيـنـ:ـ فـإـمـاـ أـنـ يـأـتـيـ مـعـظـمـاـ لـمـخـاـوفـ أـخـلـاقـيـةـ حـقـيقـةـ لـكـنـهـاـ لـيـسـ بـالـقـدـرـ نـفـسـهـ مـنـ الـجـسـامـةـ،ـ وـإـمـاـ أـنـ يـأـخـذـ صـورـةـ تـعـيـيـنـ مـشـكـلـاتـ أـخـلـاقـيـةـ لـاـ وـجـودـ لـهـاـ أـصـلـاـ.ـ وـلـأـغـرـاضـ الـإـيـضـاحـ سـأـغـامـرـ،ـ مـعـ مـاـ يـنـطـويـ عـلـيـهـ ذـلـكـ مـنـ جـدـلـ مـحـتـومـ،ـ بـطـرـحـ مـثـالـ وـاحـدـ لـكـلـ مـنـ النـمـطـيـنـ السـالـفيـنـ.

أما المثال على النـمـطـ الـأـوـلـ؛ـ أيـ الـمـبـالـغـةـ فـيـ وـسـمـ قـضـاـيـاـ أـخـلـاقـيـةـ حـقـيقـةـ بـالـخـطـورـةـ،ـ فـيـتـمـثـلـ فـيـ عـتـامـةـ خـوـارـزمـيـاتـ التـعـلـمـ الـعـمـيقـ the Opacity of Deep Learning Algorithmsـ فـكـثـيرـاـ ماـ يـفـتـرـضـ أـنـ مـعـضـلـةـ "ـالـصـنـدـوقـ الـأـسـوـدـ"ـ Black Boxـ ثـعـقـدـ الـاستـخـدـامـ الـأـخـلـاقـيـ لـلـذـكـاءـ الـاـصـطـنـاعـيـ فـيـ مـجـالـ الـطـبـ؛ـ إـذـ تـدـورـ جـمـهـرـةـ مـعـتـبـرـةـ مـنـ الـأـدـبـيـاتـ حـولـ كـيـفـيـةـ بـلـوـغـ "ـالـقـابـلـيـةـ لـلـتـفـسـيرـ"ـ Explainabilityـ أـوـ التـعـوـيـضـ عـنـ غـيـابـهـاـ.ـ فـيـ حينـ يـذـهـبـ عـدـدـ قـلـيلـ مـنـ الـأـدـبـيـاتـ إـلـىـ أـنـ الـقـيـمـةـ

(26) يـشـيرـ هـذـاـ المصـطـلـحـ إـلـىـ اـتـلـافـ يـجـمـعـ عـدـدـ أـشـخـاصـ أـوـ مـؤـسـسـاتـ تـعـمـلـ مـعـاـ مـنـ أـجـلـ غـاـيـةـ مـشـترـكـةـ وـتـقـتـسـمـ الـمـوـارـدـ وـالـخـبـرـاتـ.ـ (المـتـرـجـمـ)

(27) يـُـسـتـخـدـمـ هـذـاـ المصـطـلـحـ لـوـصـفـ نـمـاذـجـ الـعـلـمـ الـآـلـيـ،ـ وـخـاصـةـ نـمـاذـجـ الـتـعـلـمـ الـعـيـقـ،ـ الـتـيـ يـصـعـبـ فـهـمـ الـآـلـيـاتـ عـمـلـهـاـ الـداـخـلـيـةـ؛ـ إـذـ يـمـكـنـ اـطـلـاعـ الـمـسـتـخـدـمـيـنـ عـلـىـ الـمـدـخـلـاتـ وـالـمـخـرـجـاتـ،ـ وـلـكـنـهـمـ لـاـ يـسـتـطـعـونـ تـبـعـ الـمـنـطـقـ الـداـخـلـيـ الـذـيـ يـسـفـرـ عـنـ الـتـائـجـ.ـ نـجـمـ عـنـ التـخـوـفـاتـ الـتـيـ يـشـيرـهـاـ هـذـاـ الـعـمـوـضـ تـطـوـرـ حـقـلـ "ـالـذـكـاءـ الـاـصـطـنـاعـيـ القـابـلـ لـلـتـفـسـيرـ"ـ Explainable AI – XAIـ.ـ (المـتـرـجـمـ)ـ لـأـطـلـاعـ أـوـسـعـ يـنظـرـ:

الأخلاقية لقضية القابلية للتفسير مبالغ فيها، فما يهم ببساطة هو إذا ما كانت تؤدي وظيفتها أم لا. ومن ثم، قد يكون السؤال الحاسم، هنا، ما إذا كان العلاج أو التشخيص الطبي المدعوم بالذكاء الاصطناعي ناجعاً؛ إذ يتضاعل الشأن المعرفي المتعلق بالكيف أو السبب إذا قورن بالأثر. حتى الأطباء من البشر كثيراً ما يعجزون عن تقديم تفسيرات مماثلة، ولا يُعد ذلك أمراً مريضاً<sup>(28)</sup>. من ثم قد يكون هناك مشكلة حقيقة هنا، إلا أن حجم القلق الذي يُثار حولها يبدو غير متناسب مع خطورتها الفعلية<sup>(29)</sup>.

وأما النمط الثاني؛ أي الإشكالات المتوجهة تماماً، فأسوق، بتحفظ، مثالاً ما يُدعى "فجوات المسؤولية"؛ إذ صدرت العديد من الأديبيات عن الفكرة القائلة إن استخدام أنظمة الذكاء الاصطناعي ذاتية التحكم يختلف فراغاً مقلقاً في إسناد المسؤولية، إلا أن وجود هذه الفجوات غير محسوم على نحو لافت<sup>(30)</sup>، ونادرًا ما نجد محاولاتٍ منهجهية لإثباتها، فغالبية الكتابات تفترض مسبقاً وجودها وتتصحر إلى البحث عن طرائق لردمها، [علمًا] أن وجودها يبدو موضعًا للريبة ابتداءً. خذ مثلاً الحالـة النموذجـية لمنظـومـات الأسلـحة ذاتـية التـحكـم: يصعب التـسلـيم بأن استـقلـالية طـائـرة مـسـيرة تـعـفي قـائـدهـا من اللـوم إـذا أصـابـت مـدنـين بـدـلـاً مـن مـقـاتـلينـ، فـمـجـرد نـشـر تقـنيـة ذاتـية التـحكـم لا يـجـعـل المسـاءـلة أقلـ إـمـكـانـيـة مـمـا هي عـلـيـه عـنـ اـعـتمـاد تقـنيـات تقـليـدية لا يـمـكـن التـنبـؤ بـسـلـوكـها وـتـنـطـويـ هيـ الأـخـرى عـلـى ضـربـ منـ المـخـاطـرـ. وهذاـ، عـلـى ما يـبـدوـ، ما تـشـيرـ إـلـيـه أـيـضاً حدـوسـ غـيرـ المـخـتصـينـ حـيـالـ هـذـهـ الحالـاتـ<sup>(31)</sup>.

(28) للاطلاع على تبعيات لها المنظور المتساهل، ينظر:

Boris Babic et al., "Beware Explanations from AI in Health Care," *Science*, vol. 373 (2021); Alex John London, "Artificial Intelligence and Black-box Medical Decisions: Accuracy versus Explainability," *Hastings Center Report*, vol. 49 (2019); John Zerilli et al., "Transparency in Algorithmic and Human Decision-Making: Is There a Double Standard?" *Philosophy & Technology*, vol. 32 (2019).

(29) أرجـزـ هنا تحـديـداً عـلـى السـيـاقـ الطـبـيـ؛ إذ قد تـخـلـفـ المـخـاطـرـ والـرهـانـاتـ فيـ سـيـاقـاتـ آخـرىـ. وقد فـلـأـتـ بعضـ الـبـاحـثـينـ أـسـبابـ عـامـةـ تـدـعـواـ لـلـهـتمـامـ بـعـتـامـةـ الـخـوارـزمـياتـ، يـنـظرـ عـلـى سـبـيلـ المـثالـ:

Jocelyn Maclure, "AI, Explainability and Public Reason: The Argument from the Limitations of the Human Mind," *Minds and Machines*, vol. 31 (2021); Sophie Dishaw, "The Right to a Justification," *Political Philosophy*, vol. 2, no. 4 (2025); Kate Vredenburgh, "The Right to Explanation," *The Journal of Political Philosophy*, vol. 30 (2022);

لا يـعـنـيـ معـالـجةـ تـلـكـ الـهـمـومـ عـلـى نـحـوـ وـافـ فيـ هـذـاـ السـيـاقـ. وـمـعـ ذـلـكـ، حتـىـ لوـ كـانـتـ هـنـاكـ أـسـبابـ عـامـةـ تـدـعـواـ إـلـىـ القـلـقـ، فقد تـبـاـيـنـ درـجـةـ حـدـثـهاـ. للـنـظـرـ فيـ عـدـدـ مـقـالـاتـ الـتـيـ اـنـقـدـتـ هـذـاـ المـوقـفـ:

Sophie Dishaw, "The Right to a Justification," *Political Philosophy*, vol. 2, no. 4 (2025); John Zerilli et al., "Transparency in Algorithmic and Human Decision-Making: Is There a Double Standard?" *Philosophy & Technology*, vol. 32 (2019).

(30) يـنـظرـ:

Johannes Himmelreich, "Responsibility for Killer Robots," *Ethical Theory and Moral Practice*, vol. 22 (2019); Sebastian Köhler et al., "Technologically Blurred Accountability," in: C. Ulbert et al. (eds.), *Moral Agency and the Politics of Responsibility* (London: Routledge, 2018); Peter Königs, "Artificial Intelligence and Responsibility Gaps: What is the Problem?" *Ethics and Information Technology*, vol. 24 (2022).

(31) Philip Robbins, "Of Machines and Men: Attributions of Moral Responsibility in AI-Assisted Warfare," *Ethics and Information Technology*, vol. 27 (2025).

وبما أنني، انسجاماً مع أطروحة هذه الورقة، أخالف الرأي الغالب في هذه المسائل، فلا أتوقع أن تُعد هذه الملاحظات الموجزة مقنعةً أو كافية<sup>(32)</sup>. ومع ذلك، فهي تكسو أطروحة البحث المجردة شيئاً من اللحم، وقد تؤدي دوراً إضافياً. أكتفي بهذا القدر، لأعود في بقية الدراسة إلى نهجي غير المباشر القائم على تحليل العوامل المؤسسية لا الاعتبارات الأولية.

خلاصه الحجة حتى هذه اللحظة: بفعل الديناميات المؤسسية داخل حقل أخلاقيات الذكاء الاصطناعي يدفع الباحثون، إلى حدٍ ما، إلى رسم صورة سلبية حول التقنية. وهذا لا يفضي إلى عرض أحادي لآثارها الأخلاقية فحسب (إذ نادرًا ما يلتفت إلى الجوانب الإيجابية)، بل يرجح أيضاً أن يقود إلى تشويه الجوانب السلبية ذاتها (وذلك حين تضخم وتتصور أكبر مما هي عليه فعلاً).

## الحجّة في سياق أزمة التكرار

من المفيد عقد مقارنة بين أزمة السلبية في حقل أخلاقيات الذكاء الاصطناعي وأزمة التكرار Replication Crisis، في عدد من الحقول التجريبية؛ إذ إن أوجه الشبه والاختلاف بين هاتين الظاهرتين تمنح فهماً مقارناً وافقاً للديناميات الإشكالية داخل هذا الحقل، فضلاً عن كشفها لما أطّرجه هنا عن طبيعة القلق وحدته.

تشير أزمة التكرار إلى أن عدداً كبيراً من النتائج المنشورة في علم النفس والعلوم الطبية وغيرها من الحقول التجريبية لا يمكن إعادة تكرارها. وإن ما يفسّر، أفضلياً، ازدحام الدوريات العلمية بالدراسات المشفوعة بتتابع يتعدّد إعادة إنتاجها كائن في طريقة التنظيم الاجتماعي لهذه الفروع المعرفية. يقدم روبرت هادسن التحليل الآتي لأسباب الأزمة: "يتناقض العلماء على الوظائف والتمويل والمكانة الاجتماعية، وتعتمد قدرتهم في نيل هذه الأمور على قدرتهم على نشر أعمالهم. علاوة على ذلك، تتسابق الدوريات العلمية على نشر النتائج الأكثر إثارةً وسداداً، ولأجل ذلك تتجنب، في العادة، نشر تكرارات للنتائج البارزة التي سبق نشرها أو تأكيد الفرضيات الصفرية Null Hypotheses<sup>(33)</sup>. والحقيقة أن البحث العلمي المنشور يُظهر ما يسمى "انحياز النشر"، بحيث يعلن العلماء عن نتائج قابلة للنشر؛ أي نتائج ذات أثر دلالي إحصائي كبير، لا بالضرورة عن نتائج صادقة أو مبررة على نحو كافٍ بالضرورة<sup>(34)</sup>.

إحدى علل أزمة التكرار ترجع إلى تفاعل عاملين: الأول، تميّل الدوريات المرموقة إلى نشر الدراسات التي تُفضي إلى نتائج ذات تأثير كبير، بدلاً من نشر دراسات التكرار [والتحقق] أو النتائج الصفرية

(32) سلامه حجّي الجوهرية لا تعتمد على مدى وجاهة الأمثلة المذكورة. أذكر ذلك لا لتحسين الحجة من النقد؛ إذ قد تكون عرضة له بوجه آخر، بل للتشديد على طابعها غير المباشر.

(33) فرضية إحصائية تفترض عدم وجود أثر أو علاقة بين المتغيرات المدروسة، وتُستخدم نقطة انطلاق لاختبار الدلالة الإحصائية. ويُعد "اختبار دلالة الفرضية الصفرية" Null Hypothesis Significance Testing الإطار الإحصائي السائد في العلوم الاجتماعية والسلوكية والطبية الحيوية. (المترجم)

(34) Robert Hudson, "Should We Strive to Make Science Bias-free? A Philosophical Assessment of the Reproducibility Crisis," *Journal for General Philosophy of Science*, vol. 52 (2021), p. 396.

الثاني، الباحثون، الذين يتوقف مسارهم المهني على سجل نشرهم، يملكون حواجز قوية لنشر أعمالهم في تلك الدوريات. هنا التضارف يُغرى الباحثين "باستحداث" نتائج عظيمة، فيليجاً بعضهم، أحياناً، إلى ممارسات بحثية جدلية، مثل التلاعب بالبيانات (Hacking P<sup>(36)</sup>) أو حتى التزوير الصريح. فتتكددس حصيلة وافرة من النتائج "المهمة" التي لا يمكن تكرارها لاحقاً؛ لأن الآثار [المفترضة] لتلك النتائج لا وجود لها في الواقع.

تشابه أزمة التكرار، في بعض الوجوه، مع أزمة السلبية في أخلاقيات الذكاء الاصطناعي. ففي الحالتين تؤدي حواجز التقدم المهني دوراً حاسماً. كذلك يظهر، في كليتهما، ميل إلى تفضيل نمط معينه من المخرجات البحثية على سواه؛ إذ تؤثر الدوريات العلمية (وكذلك لجان التوظيف) في الحقوق التي تعاني "أزمة التكرار" الدراسات التي تُعلنُ السبق إلى نتائج ذات آثار جسمية (Significant Results)، وذلك على حساب دراسات التتحقق أو الدراسات التي لا تقدم نتائج جديدة. أمّا في أخلاقيات الذكاء الاصطناعي، فيبدو أن هناك توقعاً بـألا تكتفي المقالات بوصف الآثار الإيجابية للتقنية (عامل "الأثر الإيجابي" الذي سبق ذكره)، ثم إن مقالات السجال البحثة القادرة على تفنيد المخاوف الأخلاقية التي يُشيرها الآخرون تُعدّ هي الأخرى أقل أهمية.

غير أن هنالك أيضاً فروقاً، أحدها أن أزمة التكرار تُتبع، في جزء منها، من سلوك غير أخلاقي يصدر عن الباحثين؛ يشمل ممارساتٍ بحثيةً ملتوية، وحتى تزييفاً للبيانات. ولا أرى أن هذا ينطبق على أزمة السلبية في أخلاقيات الذكاء الاصطناعي؛ إذ إن زعمي أن الدينامية المؤسسية في هذا الحقل تولد انحيازاً، لا أنها تدفع أعضاءه إلى التصرف على نحو غير نزيه. وشمة فرق آخر هو أن لدينا، في سياق أزمة التكرار، سببين للتشكيك في موثوقية المنظومة المعرفية: أولاً، يوجد مسوغ قوي للارتياب فيها بالنظر إلى هيكل الحواجز في بعض الحقوق التجريبية، وهذا مسلك غير مباشر، أو قبلي، لإثارة الشكوك حول الموثوقية. ثانياً، يتوافر دليلاً مباشراً على احتلال المنظومة؛ إذ تبيّن أن كثيراً من النتائج المنشورة يُخفق فعلياً عند التتحقق والتكرار. أمّا في أخلاقيات الذكاء الاصطناعي، فقد اكتفيت بتقديم أسباب غير مباشرة للتشكيك في موثوقيتها؛ إذ بدلاً من أن أبين، على نحو مباشر، سبب المبالغات المتصلة بالمخاوف الأخلاقية المتعلقة بالذكاء الاصطناعي، فإنني أعمد إلى ضرب من التكهّن مؤداه أن هذه المخاوف حاضرة، وذلك استناداً، حصرًا، إلى تحليل الكيفية التي ينضم بها هذا الحقل مؤسسيًا.

(35) تُشير النتائج الصفرية إلى نتائج البحث التي لا ثبت وجود أثر ذي دلالة إحصائية؛ أي التي لا ترفض الفرضية الصفرية. وتُعرف مشكلة عدم نشر هذه النتائج بـ"مشكلة الدُّرُج" File Drawer Problem، وهو مصطلح صاغه روبرت روزنثال عام 1979 للإشارة إلى أن الدراسات ذات النتائج غير الدالة إحصائياً تبقى حيسة أدراج الباحثين دون نشر. [المترجم]، ينظر:

Robert Rosenthal, "The File Drawer Problem and Tolerance for Null Results," *Psychological Bulletin*, vol. 86, no. 3 (1979), pp. 638–641.

(36) ويقصد بها لجوء الباحث إلى حِيل إحصائية، مثل استبعاد بعض البيانات، أو اختبار متغيرات عديدة، أو إيقاف جمع العينة مبكراً، بغية انتزاع قيمة إحصائية دالّة من البيانات على نحو تعسفي. (المترجم)

ليس واضحًا إذا ما كان في الإمكان إبراد دليل مباشر، مماثل لفشل أزمة التكرار في الحقول التجريبية، داخل أخلاقيات الذكاء الاصطناعي. فصحيح أن الدعاوى الأخلاقية (مثل التمسك بوجود مشكلة مرتبطة بالتقنية) قابلة للطعن، غير أن هذه الطعون محكوم عليها بأن تظلّ موضع خلاف وجدل، فلا نجد في حقل الأخلاقيات ما يكافئ سلسلة الفشل في تكرار البحوث التي شهدتها العلوم التجريبية.

إن الطابع غير المباشر للنهج الذي اتبعته يطرح قيوداً على ما أحياول تقديمه، وذلك على الرغم من شيوعه في فلسفة العلم. ففي وضعٍ مثل، كنا سنتملّك دلائل أكثر دقة و مباشرة على انحياز سلبي ممنهجه داخل الحقل، غير أن تجاهل القرائن غير المباشرة التي تشي بهذا الانحياز لن يكون تصرّفاً حكيمًا أيضًا.

على هامش متصل، انتهج ويسلّي بکوالتر مقاربةً غير مباشرة مشابهة استلهم فيها أزمة التكرار. فهو يرى أن تحليل العوامل البنوية التي تقف وراء أزمة التكرار في العلوم الحيوية والعلوم الاجتماعية يمكن أن يقدّم، بطريقٍ غير مباشر، رؤيّةً عما إذا كانت الفلسفة هي الأخرى معرّضة لأزمة مماثلة؛ ذلك أن الفلسفة الأكademie، ولا سيما حين تعتمد منهجه "دراسة الحال"، فإنها تغيّب كثيراً من العوامل البنوية المُسبّبة لأزمة التكرار في العلوم الحيوية والاجتماعية. ولذا يخامر بکوالتر الاعتقاد أن الفلسفة قد تواجه مشكلات شبيهة<sup>(37)</sup>. أمّا مشروعه هنا فأضيق نطاقاً من مشروعه؛ إذ أركّز على مجموعةٍ من العوامل الخاصة بأخلاقيات الذكاء الاصطناعي تؤثّر في موثوقية هذا الحقل تحديداً. ومع ذلك، يجمع مقاربتينا افتراضً أساساً فحواه أن التأمل في البنية المؤسّسة للفلسفة، بوصفها تخصّصاً جامعيّاً، قادرٌ على أن يكشف لنا عن صدقّتها.

## الاعتراض الأول: انحياز السلبية هو "ماندفيلي"

بعد عرض الحجة وتوضيح مفاصيلها، أنتقل إلى مناقشة اعتراضين يُوجّهان إلى الفكرة القائلة بوجود انحياز سلبي مُقلق في أخلاقيات الذكاء الاصطناعي.

قد يقال إنني تعجلت في الحكم على الحكم استناداً إلى وجود هذا الانحياز، فمن المهم أن نميز بين اللاعقلانية على مستوى الفرد واللاعقلانية على مستوى المنظومة؛ أي مستوى الجماعة العلمية ذاتها. هذا التمييز ظلّ، حتى الآن، مُعفلاً. فعند تناول مسألة الانحياز في سياق فلسفة العلم، يتراكم اهتمامنا أولاً وأساساً في الرشد المعرفي للنظام كله: هل تنتج الجماعة العلمية، بوصفها جماعة منظمة من الباحثين، معرفة يُعتقد بها؟ أما عقلانية الباحثين الأفراد، المكوّنين لهذه الجماعة، فلا تمتلك أهمية إلا بصورة غير مباشرة؛ أي بقدر ما تؤثّر في عقلانية المنظومة برمتها. وبصدق الأمر نفسه على أخلاقيات الذكاء الاصطناعي؛ إذ ما يعنينا، في المقام الأخير، هو صوابية هذا الحقل بوصفه نظاماً

(37) Wesley Buckwalter, "The Replication Crisis and Philosophy," *Philosophy and the Mind Sciences*, vol. 3 (2022).

إبستيمياً، لا صوابية أفراده على نحوٍ مستقل<sup>(38)</sup>؛ وبناً عليه، فإنَّ الهمَّ الرئيس في هذه الورقة يتمحور حول موثوقية مخرجات أخلاقيات الذكاء الاصطناعي بوصفها منظومةً معرفية.

يجمع المستويين الفردي والمنظومي ارتباطُ وثيق، فمن المنطقي أن يؤدي انحيازُ، أو شكلاً آخر من اللاعقلانية، على مستوى الفرد إلى خلل على مستوى المنظومة، لكن هذا ليس أمرًا ضروريًا. فاستنادًا إلى الأطروحة البارزة "حكاية النحل الرمزية" Bernard The Fable of the Bees Mandeville (1733–1677)، والتي تقضي بأنَّ الرذائل الخاصة يمكن أن تُفضي إلى آثار اجتماعية حميدة، أشار فلاسفةُ العلم إلى أنَّ الرذائل الإدراكية Cognitive Vices الفردية قد تُختلف آثارًا معرفية نافعة على مستوى الجماعة. وبالنظر إلى موثوقية مخرجات نظام إبستيمي، قد يكون من المستحب أن يتصف الباحثون المُتسبّلون إلى نظام إبستيمي ما بعض النقائص الإدراكية. ومن بين هذه النقائص الماندفيليَّة المحتملة انحيازاتٌ تُعدُّ لا عقلانيةً إبستيمياً على مستوى الفرد، لكنها قد تغدو مفيدةً على مستوى المنظومة<sup>(39)</sup>. وبما أنني أفترض أنَّ انحيازَ السلبية، محل النقاش، يُضعف العقلانية على المستوى المنظومي، فقد يُعرَض بأننا بصدق حالة "ذكاء ماندفيلي" Mandevillian Intelligence، وعنده يصير هذا الانحياز غير ضارٍ، على خلاف ما ذهبْت إليه.

غير أنني أرى أنَّ ثمة سببًا وجيهًا لافتراض بأنَّ انحيازَ السلبية ليس حالةً من "الذكاء الماندفيلي"، بل لا يوجد في الواقع مسوغٌ وجيهٌ لافتراض مقلوب الأمر. فالالأصل أن نفترض قاعدةً مسبقةً ضدَّ هذه الأطروحة؛ أي إنَّ اللاعقلانية على المستوى الفردي تُترجم، بوجهِ عام، إلى لا عقلانية على مستوى المنظومة، ما لم يقدَّم دليلٌ إيجابيٌّ يثبتُ أثرًا ماندفيليًّا معاكسًا. ولا يلوح، فيما يتعلق بانحيازَ السلبية موضع النقاش، أي سببٌ إيجابيٌّ يدعونا إلى تبني هذا الافتراض.

أكثرُ الأسباب وجاهةً للاعتقاد بأنَّ انحيازَ السلبية قد يكون ماندفيليًّا هو أنَّ انحيازًا آخر، هو "انحياز التأييد" Confirmation Bias، طُرح بوصفه حالةً ماندفيليَّة. ويُقصد بانحياز التأييد ميلُ المرء إلى

(38) Finnur Dellsén, "The Epistemic Impact of Theorizing: Generation Bias Implies Evaluation Bias," *Philosophical Studies*, vol. 177 (2020), p. 3665; David L. Hull, *Science as a Process: An Evolutionary Account of the Social and Conceptual Development of Science* (Chicago: University of Chicago Press, 1988); H. E. Longino, *Science as Social Knowledge: Values and Objectivity in Scientific Inquiry* (1990), Chap 4; Miriam Solomon, "Scientific Rationality and Human Reasoning," *Philosophy of Science*, vol. 59 (1992), p. 452; Kevin J. S. Zollman, "The Epistemic Benefit of Transient Diversity," *Erkenntnis*, vol. 72 (2010), p. 33;

أما بشأن استقلال هذين المستويين فيمكن مراجعة:

Conor Mayo-Wilson et al., "The Independence Thesis: When Individual and Social Epistemology Diverge," *Philosophy of Science*, vol. 78 (2011).

(39) Joshua May, "Bias in Science: Natural and Social," *Synthese*, vol. 199 (2021), pp. 33–53; Uwe Peters, "An Argument for Egalitarian Confirmation Bias and Against Political Diversity in Academia," *Synthese*, vol. 198 (2021); Uwe Peters, "Illegitimate Values, Confirmation Bias, and Mandevillian Cognition in Science," *The British Journal for the Philosophy of Science*, vol. 72 (2021); Paul R. Smart, "Mandevillian Intelligence," *Synthese*, vol. 195 (2018).

البحث عن البيانات أو تفسيرها بما يوافق معتقداته السابقة<sup>(40)</sup>. هذا الانحياز، على الرغم من لا عقلانيته على مستوى الفرد، عُدَّ ضرباً من ضروب "الذكاء المانديفيلي"؛ إذ يحفزُ الباحثين المتأثرين به إلى تفحُّص طيفٍ واسع من الفرضيات على نحو أكثر صرامة. ويوضحُ أوي بيترز<sup>41</sup> أنَّ انحياز التأييد على مستوى الفرد، يدفع كلَّ عالم إلى استثمار جهُدٍ كبير في أمرٍين: جمع بيانات تدعمُ أطروحته، والرُّد على الأدلة المضادة والاعتراضات؛ ما يُفضي إلى استكشافٍ دقيقٍ للأطروحة وتنميتها، لا رفضها السريع<sup>(42)</sup>. وبفضل التفاعل التبادل الذي يُلغى لا عقلانيات الأفراد، قد يكون الأثر النهائي لهذا السلوك، على مستوى المنظومة، إيجابياً<sup>(43)</sup>.

سواء نجح هذا المسوغ في تبرير اعتبار انحياز التأييد مانديفيلاً أم لا، فإنه لا ينطبق على انحياز السلبية الذي أنا بصدده هنا. فالفرق الجوهرى أنَّ انحياز التأييد يدفع الباحثين إلى تأييد قناعاتهم المُسبقة، وهي قناعات تتبادر من باحث إلى آخر، فيتوجَّه كلَّ فرد نحو مسارٍ مختلفٍ عن الآخر؛ ما يُعزِّز "تقسيم العمل المعرفي" Cognitive Division of Labor، وقد يفضي هذا إلى آثارٍ نافعة على مستوى المنظومة. أمَّا انحياز السلبية، فيدفع الباحثين جميعاً إلى الاتجاه نفسه؛ أي نحو الشائوم والسلبية. ويسبب هذه الطبيعة الأحادية الاتجاه للانحياز، يصبح من غير الجلي أي آلية يمكن أن تحولُ اللاعقلانية الفردية إلى عقلانيةٍ جماعية. ويلاحظ فينور ديلسين أنه قد تكون الانحيازات الفردية غير ضارة حين تجذب الباحثين في اتجاهات متعاكسة؛ إذ يمكنها أن تلغي بعضها بعضاً، فلا يبقى انحيازٌ على مستوى المنظومة. أمَّا الإسکال، فيظهر حين تناحر مجموعاتٍ فرعية كبيرة بالقدر نفسه وفي اتجاهٍ واحدٍ، ففي هذه الحال يتنتقل الانحياز إلى مستوى النظام المعرفي<sup>(44)</sup>. وهذا، مع الأسف، هو ما يحدث فعلاً في حقل أخلاقيات الذكاء الاصطناعي.

وبناءً عليه، تبقى قرينة الرفض المبدئي لفكرة "الذكاء المانديفيلي" قائمة بلا نقاش. فخلافاً لانحياز التأييد، يبدو أنَّ انحياز السلبية على مستوى الفرد يتنتقل، وبصورة مباشرة تقريرياً، إلى انحياز مماثل على مستوى المنظومة. فإذا صَحَّ ما اقترحته من وجود ميل سلبيٍ إشكاليٍ لدى كثيرٍ من الباحثين على نحو فردي، وجب أن نفترض أنَّ هذا الانحياز يخلُّ بموثوقية حقل أخلاقيات الذكاء الاصطناعي برمته.

(40) اعتمدنا في ترجمة "انحياز التأييد" بدلاً من "انحياز التأكيد" اختيار عادل مصطفى، ينظر: عادل مصطفى، *المغالط المنطقية: فصول في المنطق غير الصوري* (القاهرة: المجلس الأعلى للثقافة، 2007)، ص 179–185. (المترجم)

(41) Uwe Peters, "Illegitimate Values, Confirmation Bias, and Mandevillian Cognition in Science," *The British Journal for the Philosophy of Science*, vol. 72 (2021), p. 1070.

(42) Ibid.; Paul R. Smart, "Mandevillian Intelligence," *Synthese*, vol. 195 (2018), pp. 4188–4191.

(43) Finnur Dellsén, "The Epistemic Impact of Theorizing: Generation Bias Implies Evaluation Bias," *Philosophical Studies*, vol. 177 (2020), pp. 3664–3665;

وبالمثل تتطلب مقاربة مريم سولومان التجريبية الاجتماعية أنْ "تُوزَّع المُتَجَهَّات الـلـآمـبـيرـيقـيـة تـوزـيـعاً مـتكـافـئـاً". ينظر:

Miriam Solomon, "Scientific Rationality and Human Reasoning," *Philosophy of Science*, vol. 59 (1992). p. 77;

ينظر أيضًا:

David L. Hull, *Science as a Process: An Evolutionary Account of the Social and Conceptual Development of Science* (Chicago: University of Chicago Press, 1988), p. 22.

لأنهض حجّتي ضدّ اعتبار انحياز السلبية "ذكاءً ماندفيليًّا" على افتراض أن كلّ باحثٍ في أخلاقيات الذكاء الاصطناعي يرّجح تحت هذا الانحياز، فذلك غير صحيح على وجهٍ بينٍ. وحجّتي لا تفترض هذا أكثر مما يفترض المدافعون عن ماندفيليّة انحياز التأييد أن جميع أفراد الجماعة متاثرون بذلك الانحياز. المقصود، بالأحرى، أن الانحياز متى وُجد دفع المتأثرين به في الاتّجاه نفسه، ومن ثمّ ينبغي افتراض أنه يعكس سليماً على العقلانية في مستوى المنظومة.

ختاماً، يجدر التنبيه إلى أن السبب الثالث من الأسباب الثلاثة التي سبقت مناقشتها في افتراض انحياز السلبية يتصل اتصالاً مباشراً بتركيبة جماعة أخلاقيات الذكاء الاصطناعي. وبناءً عليه، وعلى خلاف مصدري الانحياز الآخرين، قد لا يكون من الأنسب فهمه بوصفه انحيازاً فردياً يُترجم لاحقاً إلى مستوى المنظومة؛ ذلك أن منشأه يدو، من الأصل، كامناً في مستوى المنظومة نفسها.

## الاعتراض الثاني: انحياز السلبية على مستوى النظام مرغوب فيه

يعتبر الاعتراض الأول أن الانحياز غير مؤذٍ لأنّه يظهر على مستوى الأفراد فقط ولا ينتقل إلى مستوى المنظومة. أمّا الاعتراض الذي نحن في صدده، فيقرّ بأن الانحياز الفردي يتسرّب فعلاً إلى مستوى المنظومة، إلاّ أنه يُعدُّ ذلك الانحياز الجمعي مرغوباً. ينبغي التنبيه إلى أن هذا الاعتراض لا يقول فحسب إنّه من الجيد أن يدفع التنظيم المؤسسي باحثي أخلاقيات الذكاء الاصطناعي إلى إخضاع تقنياته لفحص أخلاقي صارم، فهذا هدفٌ محمودٌ ولا تعارض فيه مع أطروحتي. بل إن الاعتراض يُفصّح عن أن جدوى التنظيم تكمّن في أنه يُولّد تحيزاً سليماً في أثناء هذا الفحص؛ أي إننا قد نرغب، عمداً، في أن تُبالغ منظومة أخلاقيات الذكاء الاصطناعي في نقدّها.

قد ينبع هذا التصور من ملاحظةِ مفادها أن فتاتٍ أخرى في المجتمع، وعلى رأسها المنشغلون بقطاع التقنية، يبالغون في التفاؤل ويفدون قدرًا من التراخي حيال الإشكالات الأخلاقية المحتملة للذكاء الاصطناعي. ومن ثمّ يبدو أننا نحتاج إلى فتاةٍ من باحثي الأخلاقيات لتکبح هذا التفاؤل الساذج. وبناءً عليه، قد يذهب الاعتراض إلى القول إنه من المستحسن تنظيم حقل أخلاقيات الذكاء الاصطناعي على نحوٍ يفضي إلى "فرط إنتاج" للتزعّع السلبية؛ إذ إن التشاوئ المفرط داخل الحقل، والتفاؤل المفرط في دوائر أخرى، يُلغي أحدهما الآخر، فتتكتون صورةً إجماليةً أكثر اعتدالاً للتقنية.

وفي وسعنا فهم هذه الفكرة مرّةً أخرى عبر منظور "الذكاء الماندفيلي"، مع اختلافِ جوهري؛ هو أن جماعة أخلاقيات الذكاء الاصطناعي تُعامل هنا بوصفها جزءاً من منظومةٍ أوسع، المجتمع بأسره مثلاً، تضم جماعات أخرى كالمستغلين في قطاع التقنية. فالانحياز السلبي المفترض لا يُعد ماندفيليًّا داخل الجماعة نفسها، بل يصبح ماندفيليًّا على المستوى الأعلى؛ أي مستوى المجتمع كله. ولكي تعمل هذه المنظومة الكبرى كما ينبغي، قد يلزم أن تتحيز مكوناتها الفرعية، مثل جماعتي الأخلاقيين والتقنيين، بصورةٍ منهجيةٍ في اتجاهين متقابلين.

تنطوي هذه الفكرة على قدرٍ من الوجاهة؛ إذ يُحتمل أن القطاع التقني، ولأسباب بنويةٍ مشابهة، يبالغ في التفاؤل عند تصويره تقنيات الذكاء الاصطناعي ويميل إلى تَسَهِّل الهواجس الأخلاقية المشروعة تحت البساط<sup>(44)</sup>. ومع ذلك، لا أرى أن من الحكمَة تَنظيم حقل أخلاقيات الذكاء الاصطناعي، عن سابق تصور وتصميم، على نحوٍ يتوقع منه إنتاج سرديةٍ مفرطة السلبية حول تبعاته الأخلاقية.

بادئ ذي بدء، تبدو فكرة تصميم جماعةٍ فلسفيةٍ بحيث تكون مُتحيزَة، عن سابق تصور وتصميم، أمراً مُرِيباً في حد ذاته؛ إذ لو رفع معهداً لأخلاقيات الذكاء الاصطناعي شعاراً يقول: "بلغ في تصوير إشكاليات الذكاء الاصطناعي بأكثر مما تقتضيه حاله"، لاستوقفنا ذلك لا محالة، وإذا كان من الصحيح أن هذه الملاحظة لا تُفنِّد الاستدلال السابق على نحوٍ حاسم، فإنها تكفي، على أقل تقدير، لحتَّى على التوقف والتأمل.

ثمة سبب آخر هو أن يوجد قيمة جوهرية Intrinsic في الإجابات الصائبة عن الأسئلة الفلسفية. فإن تحقق انحيازٍ منهجي على مستوى المنظومة، فلن نجد أنفسنا إزاء وضعيةٍ مُستغربةٍ فحسب، بل إننا سنفقد قسماً من تلك القيمة الجوهرية أيضاً. على الرغم من أن المنظومة القائمة لا تفتقر إلى الموثوقية، فإنها ليست، في الوقت نفسه، مُهيأةً على الوجه الأمثل لإنتاج تَبَصُّراتٍ فلسفية صادقة وتجنب الأباطيل. ومهمما يكن ما يُكسبه هذا الانحياز من قيمة أداتية، فإنه يشتمل من دون شك على خسارة قيمة جوهرية<sup>(45)</sup>.

غير أن أبرز ما يُضعف الاعتراض الثاني أن التصور الذي يفترضه عن الديناميات الاجتماعية ناقصٌ على نحوٍ مُضللٍ، فهو يفترض أن المجتمع، في غياب حقل أخلاقيات الذكاء الاصطناعي، سيُبالغ في التفاؤل أو التساهل حيال الذكاء الاصطناعي؛ ما يستلزم تصدِّياً من باحثي الأخلاقيات. غير أن المستغلين في قطاع التقنية ليسوا الفاعل الاجتماعي الوحيد الجدير بالاعتبار، فصحيح أن هذا القطاع يعرض الذكاء الاصطناعي عرضاً مُغرياً في الإيجابية، لكن هناك قوى أخرى تميل، في المقابل، إلى صبغ التقنية بصبغةٍ تشاؤميةٍ مفرطة أو إلى كبح تطورها واستثمارها النافع. ومن بين هؤلاء الفواعل مشرّعون وسياسات وبيروقراطيون يتحرّكون لاستخدام السلطة التشريعية، فضلاً عن وجود نقاباتٍ مهنية يهمّها إعاقة التقنيات التي تقلل الاعتماد على العمالة البشرية، فضلاً عن قصورٍ عام في إدراك عموم الناس للأهمية الاقتصادية لهذه التقنيات<sup>(46)</sup>، إضافةً إلى التحيز للوضع القائم Status Quo Bias وانحياز "بعض الخسارة" Loss Aversion، اللذين يجعلان الناس أكثر خشيةً من الخسائر المحتملة من تبني التقنيات الجديدة المحفوفة بالمخاطر مقارنة بما يبني على عدم تبنيها من تكلفة وفرص

(44) Jan-Christoph Heilinger, "The Ethics of AI Ethics. A Constructive Critique," *Philosophy & Technology*, vol. 35 (2022), pp. 11–12.

(45) لرأي مماثل ينظر مايكل هيمير. ففكرة أن الحقيقة ذات قيمة جوهرية فكرةٌ مثيرة للجدل، فالبراغماتيون سيترضون على هذه الفكرة على سبيل المثال: Michael Huemer, "The Ethicist's Vic," *Fake Nous*, 2023, at: <https://acr.ps/1L9F2tE>

(46) Bryan Caplan, *The Myth of the Rational Voter: Why Democracies Choose Bad Policies* (Princeton: Princeton University Press, 2007), pp. 40–43.

ضائعة<sup>(47)</sup>. ومن ثم، قد تتضاد هذه القوى مع قوى اجتماعية ونفسية أخرى لتشهيم في موقفٍ شديدٍ من النقد من الذكاء الاصطناعي.

يظل من غير الواضح إذا ما كان تعامل المجتمع مع الذكاء الاصطناعي، قبل إسهام باحثي الأخلاقيات، يتسم، في المحصلة، بقدرٍ مفرطٍ من الإيجابية أو السلبية، بل يمكن القول، مع قدرٍ من المعقولة، إن المجتمع قد يتتفع بجرعةٍ إضافية من التفاؤل بالتقنية. وعلى أي حال، لا يتوافر ما يكفي من المسوغات للاعتقاد بأنه من المرغوب أن تتجه جماعة أخلاقيات الذكاء الاصطناعي نحو السلبية<sup>(48)</sup>. وإذا انعدم هذا المسوغ، وجب على جماعة الباحثين في هذا الحقل أن تسعى لتكون حكماً محايضاً، لا إحدى القوى المنحازة. ومن المُرجح أن إحدى الشمار الإيجابية لهذا الحياد هي تعزيز ثقة الفاعلين الاجتماعيين الآخرين بجماعة باحثي هذا الحقل؛ إذ ستتجدد صعوبةً في كسب هذه الثقة إذا صُنمت بنيتها على انحيازٍ سلبيٍّ، أو إذا تصورت نفسها صراحةً في موقع مواجهةٍ مع القوى المؤيدة للتقنية.

## ما يترتب وما لا يترتب على ذلك

إذا كان تحليلي يسير، عامة، في الاتجاه الصحيح، فشَّمة ما يدعو إلى افتراض أن التنظيم المؤسسي لأخلاقيات الذكاء الاصطناعي، بوصفها تخصصاً جامعياً يمسّ موثوقيتها مسأً سليماً. فبحكم الحوافز التي تفرضها المنظومة، يُدفع الباحثون إلى عرض الذكاء الاصطناعي في ضوء نقدي سلبي، ومن ثم يُرجح أن تنشأ صورة إجمالية لآثاره الأخلاقية لا تبدو أحادية الوجه فحسب، بل منحازة كذلك إلى السلييات. فما الاستنتاجات التي ينبغي لنا استخلاصها من هذا؟

يتبدّى الدرس الأهم في ضرورة التقليل من شأن الصورة السلبية التي يرسمها مبحث أخلاقيات الذكاء الاصطناعي حول هذا الأخير؛ إذ لا يعني ما يقدّمه هذا المبحث من رؤية قاتمة حول الآثار السلبية للذكاء الاصطناعي، بالضرورة، أنها صورة قاتمة إلى هذه الدرجة، ولنقل، بصورة أدق، يوفر انحياز السلبية حجة "ملغٌ تقويضي" Undercutting Defeater<sup>(49)</sup> للرأي القائم على التقييم السلبي للجماعة

(47) Cass R. Sunstein, *Laws of Fear: Beyond the Precautionary Principle* (Cambridge: Cambridge University Press, 2005), pp. 41–43.

(48) ينسحب هذا أيضاً على التقنية بعامة، فليس بيتاً أن المجتمع سيسفيد من انحياز جماعة أخلاقيات التقنية نحو السلبية والدعوة إلى مزيد من الحيطة والتنظيم.

(49) يحتل مصطلح المُلغِي Defeater والقابلية للإلغاء Defeasibility مكاناً بارزاً في الإبستيمولوجيا المعاصرة، ينظر: صلاح إسماعيل، نظرية المعرفة: مقدمة معاصرة، ط 2 (القاهرة: الدار المصرية اللبنانية، 2022)، ص 61–62؛ تميز الأدبيات الإبستيمولوجية بين نوعين رئيسيين من المُلغِيات: 1. المُلغِي الدحضي Rebutting Defeater، وهو سبب لاعتقاد تقىض القضية (ق) أو قضية أخرى تتعارض معها. مثال ذلك: إذا رأت مريم من بعيد ما يبدو خروفاً في الحقل فاعتقدت وجود خروف، ثم أخبرها صاحب الحقل أنه لا توجد خراف في، فقد اكتسبت مُلغِي دحضياً لمعتقدتها. 2. المُلغِي التقويضي Undercutting Defeater، وهو سبب للكف عن اعتقاد (ق) دون أن يكون سبباً لاعتقاد تقىضها، وذلك بمحاجمة الصلة بين الدليل والاستنتاج. ومثال ذلك: إذا رأى شخص قطعاً تبدو حمراء، ثم علم أنها مضاءة بأضواء حمراء، فقد فَقد سببه لاعتقاد أنها حمراء، لكنه لم يكتسب سبباً لاعتقاد أنها ليست حمراء. يُستخدم هذا المفهوم في الورقة للإشارة إلى أن التحيز السلبي المؤسسي في أخلاقيات الذكاء الاصطناعي يُقدم مُلغِي تقويضياً للتقييم السلبي الذي تطرحه الجماعة العلمية [المترجم]، للمزيد ينظر أيضاً:

Michael Sudduth, "Defeaters in Epistemology," in: *Internet Encyclopedia of Philosophy*, accessed on 19/1/2026, at: <https://acr.ps/1L9F2Qp>

العلمية؛ ومفاده أن الذكاء الاصطناعي إشكالي بدرجة كبيرة. ومن المفترض، في الأحوال العادلة، أن يُشكّل هذا التقييم السلبي دليلاً قوياً على أن الذكاء الاصطناعي يتسم بوضع إشكالي بدرجة كبيرة، وذلك تبعاً للطراائق الكثيرة والانتقادات التي دفع بها مجتمع أخلاقيات الذكاء الاصطناعي، أمّا مع وجود الانحياز، فإن هذا الدليل لا يدعم تلك النظرة القاتمة، أو أنه يدعمها بدرجة أوهى وأضعف. ومن الممكن صوغ هذا الدرس من منظور جماعة أخلاقيات الذكاء الاصطناعي نفسها، ونظام، إذاك، مع ما يمكن تسميته المُلغِي من الدرجة العليا Higher-Order Defeater<sup>(50)</sup>، وهو يُعد غالباً نوعاً خاصاً من المُلغِي التقويمي. وكما تُبيّن ماريا لاسونين-آرنيو، فإن "المُلغِي" من الدرجة العليا "يعمل عبر إثارة الشك في أن الحالة الاعتقادية Doxastic التي يتبنّاها المرء ناتجة من عملية قاصرة"<sup>(51)</sup>.

وقد ديفيد كريستنسن مثلاً نموذجاً على ذلك في المثال التالي: "أنا طبيب مقيم أُشخّص حالات المرض وأصفُ العلاج المناسب. بعد أن شخصت حالة مريض معين ووصفته له أدوية محددة، أخبرتني الممرضة بأنني مستيقظٌ منذ سِتٍ وثلاثين ساعة. وبما أعرفه عن ميل الناس إلى ارتكاب أخطاء معرفية عندما يُحرّمون من النوم (وربما أيضاً لما أعلمه عن ضعف سجلِي التشخيصي في مثل هذه الظروف)، خفضت درجة ثقتي بتשתיحي ووصفتي ريثما أعيد فحصرأبي بعناية"<sup>(52)</sup>. هنا يمتلك الطبيب دليلاً من الدرجة العليا Higher-Order Evidence يدعوه إلى الارتباط في تقييمه أدلة الدرجة الأولى First-Order الخاصة بحالة المريض والأدوية الالزمة. وعادةً يُناقض "الدليل من الدرجة العليا" في إطار العقلانية الفردية، كما في هذه الحالة، غير أن بسطه ليشمل الجماعات العلمية ليس بعيداً، فكما توجد عقلانية فردية تتصل بها أدلة من الدرجة العليا، يمكن أن يكون للمنظومة المعرفية ذاتها عقلانية بدرجة كبرت أم صغرت (أو موضوعة، وموضوعية)، متوفّرة على أدلة من الدرجة العليا. فكما يرى الطبيب أن الأدلة من الدرجة الأولى تدعم تشخيصاً وعلاجاً معيناً، تستشكّل جماعة أخلاقيات الذكاء الاصطناعي هذا الأخير، وذلك بما ينطوي عليه من مشكلات أخلاقية عديدة. وبهذا، فإن الأدلة من الدرجة الأولى تُقيّم بوصفها داعمة لنظرة سلبية شاملة لأنّ الذكاء الاصطناعي الأخلاقية. وكما كان من المُتعيّن على الطبيب أن يقلّل من موضوعية تشخيصه، يجدّر بهذه الجماعة أن تخفّف من يقيّها في تقديرها القائم لهذه التقنية.

(50) في نظرية المعرفة، تُعدّ "الأدلة من الدرجة الأولى" First-order Evidence هي تلك التي تتعلّق مباشرةً بصدق قضية معينة، بينما "الأدلة من الدرجة العليا" Higher-order Evidence هي أدلة حول طبيعة أدلة ذاتها، أو حول قدراتنا المعرفية واستعداداتنا للاستجابة العقلانية لها. على سبيل المثال، إذا كان دليل الأرصاد الجوية دليلاً من الدرجة الأولى على أن المطر سيهطل غداً، فإن معرفة أن عالم الأرصاد كان في حالة إرهاق شديد عند تقييمه للبيانات تُعدّ دليلاً من الدرجة العليا. وقد أصبح هذا التمييز محوريًا في الأدب الإبستيمولوجي المعاصرة. (المترجم). للمزيد، ينظر:

Sophie Horowitz, "Higher-Order Evidence," *The Stanford Encyclopedia of Philosophy* (Spring 2025), Edward N. Zalta & Uri Nodelman (eds.), at: <https://plato.stanford.edu/cgi-bin/encyclopedia/archinfo.cgi?entry=higher-order-evidence>

(51) Maria Lasonen-Aarnio, "Higher-Order Evidence and the Limits of Defeat," *Philosophy and Phenomenological Research*, vol. 88 (2014), p. 314.

(52) David Christensen, "Higher-Order Evidence." *Philosophy and Phenomenological Research*, vol. 81, no. 1 (2010), p. 186.

مع ذلك، ثمة أمور لا تترتب على ما سبق، فلا يجب أن يفهم من هذه الأطروحة أن استخدام الذكاء الاصطناعي لا يشير إلى إشكالاتٍ أخلاقيةً جسمية، فمن الواضح أن مثل هذه الإشكالات قائمة، وتظل دراستها من باحثي الأخلاق مهمّةً بالغة الأهمية. وعلى سبيل المثال، يبدو لي أن احتمالاً، ولو ضئيلاً، لظهور ذكاءٍ اصطناعي خارق Superintelligence قادرٍ على إحداث دمار كارثي يُعدّ مسألةً مهمةً<sup>(53)</sup>.

ولا يترتب على ما سبق أن كل إشكالية أخلاقية يُزعم ارتباطها بالذكاء الاصطناعي تعدّ وهمية أو مبالغ فيها، فحتى لو وُجد ميلٌ داخل جماعة الحقل إلى الإكثار غير المبرر من هذه الإشكالات أو تضخيمها، يبقى من غير الجائز إقصاء أي ادعاءٍ محدد حول خللٍ أخلاقي يتعلق بالذكاء الاصطناعي. فالإلمام بالديناميات المؤسسية التي قادت إلى أزمة التكرار في العلوم التجريبية لا يبيح لنا رفض النتائج الإيجابية المُبهرة من حيث المبدأ، وبالمثل، فإن الوعي بما يعتور الحقل من اختلال لا يبرر ردّ مزاعم الإشكال الأخلاقي في الذكاء الاصطناعي جملةً وتفصيلاً. وبناءً عليه، ينبغي التعامل مع كل دعوى على حدة وجدية، لا بالطعن فيها استناداً إلى اعتباراتٍ تخصُّ النظام كله<sup>(54)</sup>. فلا تهدف هذه الورقة إلى دحض مسألة أخلاقية بعينها، وإنما يسعى لتقدير موثوقية المنظومة الإبستيمية من منظور فلسفة العلم؛ أي على مستوى النظام لا على مستوى القضايا الجزئية داخل حقل أخلاقيات الذكاء الاصطناعي.

ختاماً، لا يعني ما تقدّم أن لدينا سبباً يدعونا إلى تفاؤلٍ إيجابيٍّ خالصٍ حيال الحصيلة الأخلاقية للذكاء الاصطناعي. ف مجرد ملاحظة الصورة التي تنظم فيها المنظومة، كما تقدّم الوصف، يوفر لنا "معنىًّا تقويفياً" ، بدرجةٍ غير محسومة، وذلك نتيجة الدليل الذي يحمله التقسيم السلبي الصادر عن جماعة الحقل، غير أنه لا يزودنا بحججة إيجابية تدفعنا إلى تبني موقفٍ محابٍ للذكاء الاصطناعي. فحتى لو ثبت أن حقل أخلاقيات الذكاء الاصطناعي يتسم بالخلل الذي أشرتُ إليه، يظل من الوارد، على المستوى النظري، أن يكون الذكاء الاصطناعي في المُجمل شرّاً أخلاقياً. ومن ثم، فإن نقدي للسلبية المتفشية في هذا البحث لا يرقى إلى بناء قضية إيجابية للتلفؤ بالذكاء الاصطناعي<sup>(55)</sup>.

(53) ينظر على سبيل المثال:

Leopold Aschenbrenner, "Situational Awareness," (June 2024), accessed on 19/1/2026, at: <https://acr.ps/1L9F2lw>; Nick Bostrom, *Superintelligence: Paths, Dangers, Strategies* (Oxford: Oxford University Press, 2014); Leonard Dung, "The Argument for Near-term Human Disempowerment through AI," *AI & Society*, vol. 50 (2025); Michael Huemer, "I, for One, Welcome Our Robot Overlords," *AI & Society* (2025); Toby Ord, *The Precipice: Existential Risk and the Future of Humanity* (London: Bloomsbury, 2020), Ch. 5; Eliezer Yudkowsky & Nate Soares, *If Anyone Builds It, Everyone Dies: Why Superhuman AI Would Kill Us All* (New York: Little, Brown and Company, 2025).

(54) حول القيود على توظيف الدليل من المرتبة الأعلى في الجدل الفلسفية، ينظر:

Zach Barnett, "Philosophy without Belief," *Mind*, vol. 128 (2019); Peter Königs, "Bracketing Higher-order Evidence in Scholarly Philosophical Argumentation: Why and Which?" *Synthese*, vol. 205 (2025).

(55) حول صعوبة بناء حجّة إيجابية للتلفؤ بالتقنية، ينظر:

John Danaher, "Techno-optimism: An Analysis, an Evaluation and a Modest Defense," *Philosophy & Technology*, vol. 35 (2022).

خلاصة أخيرة، يتعين علينا التفكير في سبل لإصلاح المشكلة، وإن لم يكن ذلك يسيراً. ففي الحقول التي عصفت بها أزمة التكرار أمكن لبعض تعديلاتٍ مؤسسية، كالتسجيل المسبق للدراسات، وإمكانية نشر النتائج الصفرية، والبيانات المفتوحة المصدر، أن تقطع شوطاً لا بأس به في معالجة هذه الأزمة. أمّا في أخلاقيات الذكاء الاصطناعي، فليس واضحاً أن حلولاً مماثلة متاحة؛ إذ لا يمكن، ولا ينبغي، العبث عقلياً بأي من العوامل المؤسسية الثلاثة. فموضوع الحقل سيظل، حتماً، هو ذاته، ومن المنطقي أن يقدم المشغلون في البحث الأخلاقي تعليقاً قيمياً على التطورات التقنية الجديدة. قد يتصور، نظرياً، التخلص عن القاعدة التي تُنْهي عن المقاربات الإيجابية ذات الطابع الوصفي، لكن ليس من الواضح أن ذلك تصرفٌ حصيف، فكما سبق لنا الذكر، يمكن تَفَهُّم الدافع لشيطنة باحثي أخلاقيات الذكاء الاصطناعي، وثنיהם عن إنتاج مقالات ضعيفة الطابع التقريري. ومن ثم، ثمة وجاهةٌ للإبقاء على هذه القاعدة. أما فيما يتعلق ببنى الحواجز داخل الأكاديمية فإن تغييرها أمرٌ عسيرٌ، كما أن تقديم أي منظومة حواجز بديلة لمعالجة هذه المعضلة قد يستجلب مشكلاتٍ أخرى، وتنتهي بحصيلة أسوأ من سابقتها. ف الصحيح أن المقالات السجالية يمكنها أداء دورٍ قيم في ترشيد التقييمات السلبية للذكاء الاصطناعي، لكننا قد نندم لو شجّعنا فيضًا من المقالات السجالية التي لا تقدم سوى تحفظات طفيفة. يفضي هذا كله إلى معضلة: إن مشكلة السلالية تبدو أثراً جانبياً مزعاً لتضافر عوامل مؤسسية، يمكن تَعْقُل كل منها على حدة.

إذا كان إحداث تغييرات في الإطار المؤسسي أمراً متعرّضاً، فإن ثمة تدابير أخرى قد تُعدُّ بتغييرات طفيفة. وتتضمن هذه التدابير تغيير طريقة تفكيرنا وممارساتنا البحثية. ومن المتعيّن علينا، بصفتنا كُتاب مقالات حول أخلاقيات الذكاء الاصطناعي، أن نكون محظيين بالдинاميات المؤسسية التي تقتادنا إلى ضرب من السلالية. ولا بدّ لنا من بذل جهدٍ واع لتجنب الإفراط في سوق المظاهر السلالية. ويتعين علينا، بوصفنا مراجعين ومحرّرين، إخضاع المزاعم المتصلة بالمشاكل الأخلاقية للذكاء الاصطناعي لتدقيق وفحص كافيين. ومن واجبنا، أيضاً، أن نضاعف جهودنا لإشراك غير المتخصصين في الأخلاقيات، ولا سيما مطوري الذكاء الاصطناعي وعلماء الكمبيوتر والزج بهم في مشاريعنا البحثية بهدف التتحقق الخارجي من صحة المخرجات. فال усили الحيث للحصول على المزيد من الآراء الإيجابية من خارج الحقل الفلسفـي لا يعني، بالضرورة، الموافقة عليها وتزكيتها، وإنما لاكتساب منظور أشمل ربما يعين في مواجهة تحيزات المباحث المعرفـية.

ستبقى هذه الإجراءات حللاً مؤقتاً؛ إذ إنها تترك дيناميـات المؤسسـية من دون علاج. ومن المرجـح أن تواصل أمواج الفزع ارتفاعـها ما دامت أخلاقيـات الذكاء الاصـطناعـي تـشهد ازدهـارـاً. سيـحدث ذلك، إلى حدـ بعيدـ، بمـعـزلـ عنـ مـدىـ الإـشكـاليـاتـ الأـخـلـاقـيـةـ الـواقـعـيـةـ لـلـذـكـاءـ الـاصـطـنـاعـيـ. أمـاـ "الـخيـارـ النـوـويـ"ـ،ـ فـيـتـمـثـلـ فـيـ تقـليـصـ جـحـمـ هـذـاـ الحـقـلـ عـنـ سـبـقـ إـصـرـارـ وـتـصـمـيمـ،ـ حـرـصـاـ عـلـىـ مـصـلـحـتـهـ.ـ فـقدـ أـشـارـ غـيرـ مـؤـلـفـ إـلـىـ أـنـ الحـقـلـ تـضـخـمـ عـلـىـ نـحـوـ غـيرـ مـلـائـمـ<sup>(56)</sup>.ـ فـحـجـمـ التـموـيلـ المـخـصـصـ لـهـ لـاـ يـشـيرـ

(56) Jan-Christoph Heilinger, "The Ethics of AI Ethics. A Constructive Critique," *Philosophy & Technology*, vol. 35 (2022), p. 16; Luke Munn, "The Uselessness of AI Ethics," *AI and Ethics*, vol. 3 (2023), p. 873; Felicitas Lambrecht & Marina Moreno, "What is AI Ethics?" *American Philosophical Quarterly*, vol. 61 (2024), pp. 397–398.

السؤال فحسب قياساً على تكلفة الفرص الضائعة أو العوائد الحدية المتناقصة Marginal Returns، بل قد يقوض أيضاً موثوقية المنظومة الإبستيمية عبر توليد فجوة بين حجم الجماعة البحثية في حقل أخلاقيات الذكاء الاصطناعي وعدد الإشكالات الأخلاقية المطلوب بحثها، فضلاً عن شدتها.

## اعترافٌ وتقدير

أود أن أتقدم بالشكر الجزيل للملحوظات القيمة والحوارات المُشرية لكل من المحكمين من مجلة Synthese، ولجمهور المحاضرات في كريت Crete، ودورتموند Dortmund، ودوسلدورف Düsseldorf، وأيندهوفن Eindhoven، وإرلنغن Erlangen، ومانهايم Mannheim، وللمشاركين في حلقة "أوبسالا-فيينا للذكاء الاصطناعي" Uppsala-Vienna AI Colloquium.

## التمويل

إتاحة الوصول المفتوح جرى تنظيمها وتمويلها عبر Projekt DEAL.

## References

## المراجع

- Arora, Payal. *From Pessimism to Promise: Lessons from the Global South on Designing Inclusive Tech*. Cambridge, MA: MIT Press, 2024.
- Arvan, Marcus, Liam Kofi Bright & Remco Heesen. "Jury Theorems for Peer Review." *British Journal for the Philosophy of Science*. vol. 76 (2025).
- Avin, Shahar. "Centralized Funding and Epistemic Exploration." *The British Journal for the Philosophy of Science*. vol. 70 (2019).
- Babic, Boris et al. "Beware Explanations from AI in Health Care." *Science*. vol. 373 (2021).
- Barnett, Zach. "Philosophy without Belief." *Mind*. vol. 128 (2019).
- Bostrom, Nick. *Superintelligence: Paths, Dangers, Strategies*. Oxford: Oxford University Press, 2014.
- \_\_\_\_\_. *Deep Utopia: Life and Meaning in a Solved World*. Idea, 2024.
- Bright, Liam Kofi & Remco Heesen. "To Be Scientific Is to Be Communist." *Social Epistemology*. vol. 37 (2023).
- Buckwalter, Wesley. "The Replication Crisis and Philosophy." *Philosophy and the Mind Sciences*. vol. 3 (2022).
- Caplan, Bryan. *The Myth of the Rational Voter: Why Democracies Choose Bad Policies*. Princeton: Princeton University Press, 2007.
- Christensen, David. "Higher-Order Evidence." *Philosophy and Phenomenological Research*. vol. 81, no. 1 (2010).
- Danaher, John. *Automation and Utopia: Human Flourishing in a World without Work*. Cambridge, MA: Harvard University Press, 2019.

- \_\_\_\_\_. "The Rise of the Robots and the Crisis of Moral Patency." *AI & Society*. vol. 34 (2019).
- \_\_\_\_\_. "Techno-optimism: An Analysis, an Evaluation and a Modest Defence." *Philosophy & Technology*. vol. 35 (2022).
- Dellsén, Finnur. "The Epistemic Impact of Theorizing: Generation Bias Implies Evaluation Bias." *Philosophical Studies*. vol. 177 (2020).
- Dishaw, Sophie. "The Right to a Justification." *Political Philosophy*. vol. 2, no. 4 (2025).
- Douglas, Heather & Kevin C. Elliott. "Addressing the Reproducibility Crisis: A Response to Hudson." *Journal for General Philosophy of Science*. vol. 53 (2022).
- Dung, Leonard. "The Argument for Near-term Human Disempowerment through AI." *AI & Society*. vol. 50 (2025).
- Floridi, Luciano. "Introduction to the Special Issues: The Ethics of Artificial Intelligence." *American Philosophical Quarterly*. vol. 61 (2024).
- Gunkel, David J. (ed.). *Handbook on the Ethics of Artificial Intelligence*. Cheltenham: Edward Elgar Publishing, 2024.
- Hagendorff, Thilo. "Blind Spots in AI Ethics." *AI and Ethics*. vol. 2 (2022).
- Heesen, Remco. "Why the Reward Structure of Science Makes Reproducibility Problems Inevitable." *The Journal of Philosophy*. vol. 115 (2018).
- Heesen, Remco & Liam Kofi Bright. "Is Peer Review a Good Idea?" *The British Journal for the Philosophy of Science*. vol. 7 (2021).
- Heilinger, Jan-Christoph. "The Ethics of AI Ethics. A Constructive Critique." *Philosophy & Technology*. vol. 35 (2022).
- Himmelreich, Johannes. "Responsibility for Killer Robots." *Ethical Theory and Moral Practice*. vol. 22 (2019).
- Hoffmann, Reid & Greg Beato. *Superagency: What Could Possibly Go Right with Our AI Future*. Authors Equity, 2025.
- Hudson, Robert. "Should We Strive to Make Science Bias-free? A Philosophical Assessment of the Reproducibility Crisis." *Journal for General Philosophy of Science*. vol. 52 (2021).
- Huemer, Michael. "I, for One, Welcome Our Robot Overlords." *AI & Society*. (2025).
- Hull, David L. *Science as a Process: An Evolutionary Account of the Social and Conceptual Development of Science*. Chicago: University of Chicago Press, 1988.
- Kitcher, Philip. "The Division of Cognitive Labor." *The Journal of Philosophy*. vol. 87 (1990).
- Königs, Peter. "Artificial Intelligence and Responsibility Gaps: What is the Problem?" *Ethics and Information Technology*. vol. 24 (2022).
- \_\_\_\_\_. "Bracketing Higher-order Evidence in Scholarly Philosophical Argumentation: Why and Which?" *Synthese*. vol. 205 (2025).

Lambrecht, Felicitas & Marina Moreno. "What Is AI Ethics?" *American Philosophical Quarterly*. vol. 61 (2024).

Lasonen–Aarnio, Maria. "Higher–Order Evidence and the Limits of Defeat." *Philosophy and Phenomenological Research*. vol. 88 (2014).

Lee, Carole J. "Commensuration Bias in Peer Review." *Philosophy of Science*. vol. 82 (2015).

Lobel, Orly. *The Equality Machine: Harnessing Digital Technology for a Brighter, more Inclusive Future*. London: Hachette UK, 2022.

London, Alex John. "Artificial Intelligence and Black–box Medical Decisions: Accuracy versus Explainability." *Hastings Center Report*. vol. 49 (2019).

Longino, Helen E. *Science as Social Knowledge: Values and Objectivity in Scientific Inquiry*. Princeton: Princeton University Press, 1990.

\_\_\_\_\_. *The Fate of Knowledge*. Princeton: Princeton University Press, 2002.

Maclure, Jocelyn. "AI, Explainability and Public Reason: The Argument from the Limitations of the Human Mind." *Minds and Machines*. vol. 31 (2021).

May, Joshua. "Bias in Science: Natural and Social." *Synthese*. vol. 199 (2021).

Mayo–Wilson, Conor et al. "The Independence Thesis: When Individual and Social Epistemology Diverge." *Philosophy of Science*. vol. 78 (2011).

Milano, Silvia, Carina Prunkl. "Algorithmic Profiling as a Source of Hermeneutical Injustice." *Philosophical Studies*. vol. 182 (2025).

Munn, Luke. "The Uselessness of AI Ethics." *AI and Ethics*. vol. 3 (2023).

Ord, Toby. *The Precipice: Existential Risk and the Future of Humanity*. London: Bloomsbury, 2020.

Peters, Uwe. "An Argument for Egalitarian Confirmation Bias and Against Political Diversity in Academia." *Synthese*. vol. 198 (2021).

\_\_\_\_\_. "Illegitimate Values, Confirmation Bias, and Mandevillian Cognition in Science." *The British Journal for the Philosophy of Science*. vol. 72 (2021).

Rebera, Andrew P., Lode Lauwaert & Ann–Katrien Oimann. "Hidden Risks: Artificial Intelligence and Hermeneutic Harm." *Minds & Machines*. vol. 35, no. 33 (2025).

Robbins, Philip. "Of Machines and Men: Attributions of Moral Responsibility in AI–Assisted Warfare." *Ethics and Information Technology*. vol. 27 (2025).

Smart, Paul R. "Mandevillian Intelligence." *Synthese*. vol. 195 (2018).

Solomon, Miriam. "Scientific Rationality and Human Reasoning." *Philosophy of Science*. vol. 59 (1992).

\_\_\_\_\_. *Social Empiricism*. Cambridge, MA: MIT Press, 2001.

Strevens, Michael. "The Role of the Priority Rule in Science." *The Journal of Philosophy*. vol. 100 (2003).

- Sunstein, Cass R. *Laws of Fear: Beyond the Precautionary Principle*. Cambridge: Cambridge University Press, 2005.
- Ulbert, Cornelia. et al. (eds.). *Moral Agency and the Politics of Responsibility*. London: Routledge, 2018.
- Vredenburgh, Kate. "The Right to Explanation." *The Journal of Political Philosophy*. vol. 30 (2022).
- Weisberg, Michael & Ryan Muldoon. "Epistemic Landscapes and the Division of Cognitive Labor." *Philosophy of Science*. vol. 76 (2009).
- Yudkowsky, Eliezer & Nate Soares. *If Anyone Builds It, Everyone Dies: Why Superhuman AI Would Kill Us All*. New York: Little, Brown and Company, 2025.
- Zerilli, John et al. "Transparency in Algorithmic and Human Decision-Making: Is There a Double Standard?" *Philosophy & Technology*. vol. 32 (2019).
- Zollman, Kevin J. S. "The Communication Structure of Epistemic Communities." *Philosophy of Science*. vol. 74 (2007).
- \_\_\_\_\_. "The Epistemic Benefit of Transient Diversity." *Erkenntnis*. vol. 72 (2010).