



## مرحلة المراهقة التكنولوجية

مواجهة مخاطر الذكاء الاصطناعي المتقدم والتغلب عليها

داريو أمودي - الرئيس التنفيذي لشركة آثروبليك - يناير 2026



ثمة مشهد في النسخة السينمائية من رواية كارل ساغان "كونتاكت"، تجلس فيه الشخصية الرئيسية - عالمة فلك رصدت أول إشارة راديو من حضارة فضائية - أمام لجنة دولية تُقيّم أهليتها لتمثيل البشرية في لقاء تاريخي مع الكائنات الفضائية. يسألها المحكّمون: "لو كان بإمكانك أن تطرح عليهم سؤالاً واحداً فقط، ماذا سيكون؟" فتجيب دون تردد: "سأسألهم: كيف فعلتموها؟ كيف تطورتم وكيف نجوتم من هذه المراهقة التكنولوجية دون أن تدمروا أنفسكم؟"

كلما تأملت موقع البشرية اليوم في علاقتها بالذكاء الاصطناعي - وعتبة التحول التي نقف عليها - عاد هذا المشهد إلى ذهني. فالسؤال يبدو ملائماً بشكل لافت لواقعنا، وكأننا نفتقد إجابة كان يمكن أن تهدينا الطريق.

**أعتقد أننا على أعتاب مرحلة انتقالية حاسمة، مضطربة ولكن لا مفر منها، ستختبر ما نحن عليه كبشر. سنمنح قوة تكاد تكون غير قابلة للتخيل، لكن ليس واضحاً ما إذا كانت أنظمتنا الاجتماعية والسياسية والتقنية قد بلغت النضج الكافي لحسن استخدامها.**

في مقالي السابقة "آلات النعمة العطوف"، حاولت أن أرسم صورة حلم لحضارة عبرت سالمته إلى بر النضج، واستطاعت أن تُوظف الذكاء الاصطناعي المتقدم بكفاءة ورحمة لرفع مستوى الحياة على وجه المعمورة كله. وتناولت ما يمكن أن يُسهم به الذكاء الاصطناعي من قفزات هائلة في علم الأحياء، وعلم الأعصاب، والتنمية الاقتصادية، والسلام العالمي، وعالم العمل والمعنى.

وقد رأيت أن البشرية تحتاج إلى رؤية مُلهمة تُقاتل من أجلها؛ وهو ما أخفق في تقديمه - على نحو لافت - كلٌّ من دعاة تسريع الذكاء الاصطناعي ودعاة السلامة فيه على حدٍ سواء. أما في هذه المقالة، فإنني أواجه مرحلة الاختبار ذاتها مواجهة مباشرة؛ أسعى إلى تحديد المخاطر التي تقترب



منها، وصياغة خطة عمل للتصدي لها. فأنا أوّمن إيماناً راسخاً بقدرة البشرية على النجاح، وبعظمة روحها ونبل جوهرها، لكن لا مناص من أن ننظر في عيون الواقع دون مواربة أو أوهام.

## قواعد الحوار حول المخاطر

كما أن الحديث عن فوائد الذكاء الاصطناعي يستوجب الدقة والرصانة، فكذلك الحديث عن مخاطره يستلزم الحرص والتوازن. وأرى من الضروري مراعاة ثلاثة مبادئ جوهرية في هذا السياق:

### أولاً: الابتعاد عن خطاب الانهيار والكارثية

لا أعني هنا فحسب الاعتقاد بأن الكارثة حتمية – فهذا اعتقاد خاطئ وقد يتحول إلى نبوءة تحقق ذاتها – بل أعني كذلك التعامل مع مخاطر الذكاء الاصطناعي بعقلية شبه دينية. ففي ذروة موجة القلق من الذكاء الاصطناعي بين عامي 2023 و2024، طفت على السطح أكثر الأصوات إثارةً للضجيج ونقصاً في الرصانة، تُسيء استخدام وسائل التواصل الاجتماعي لنشر خطاب مستوحى من الدراما الخيالية، داعيةً إلى إجراءات متطرفة لا تُسندها أي أدلة.

وكان واضحاً منذ البداية أن رد الفعل العكسي آتٍ لا محالة، وأن القضية ستتحول إلى ورقة ثقافية واستقطابية، فتنتهي إلى طريق مسدود. وهذا بالضبط ما جرى؛ ففي عامي 2025 و2026، انقلب المشهد رأساً على عقب، وباتت الفرص التي يتيحها الذكاء الاصطناعي هي المحرك الرئيسي للقرارات السياسية، لا المخاوف منه. وهذا التآرجح مؤسف، لأن التكنولوجيا لا تكثر بما هو رائج وما هو خارج الموضوع؛ والأدهى أننا في 2026 أقرب بكثير من الخطر الفعلي مما كنا عليه في 2023. **والعبرة الواجب استخلاصها هي أن مخاطر الذكاء الاصطناعي ينبغي أن تُناقش وتُعالج بأسلوب واقعي وعملي: رصين، مستند إلى الحقائق، وقادر على الصمود في وجه المدّ والجزر.**

### ثانياً: الإقرار بهامش عدم اليقين

كل ما يُطرح هنا قد يتبين لاحقاً أنه غير ذي صلة. ربما لا يتقدم الذكاء الاصطناعي بالسرعة المتوقعة. وربما لا تتحقق هذه المخاطر أصلاً، أو تظهر مخاطر أخرى لم نتخيلها.

لا أحد يستطيع التنبؤ بالمستقبل بدقة، لكن ذلك لا يعفينا من **التخطيط بأفضل ما نستطيع.**

### ثالثاً: التدخل بقدر ما يلزم لا بقدر ما يُتاح

مواجهة مخاطر الذكاء الاصطناعي تستوجب مزيجاً من الإجراءات الطوعية التي تتخذها الشركات والجهات الخاصة، والإجراءات الحكومية الملزمة للجميع.

أما الإجراءات الطوعية فهي في نظري أمر بديهي لا خلاف عليه. غير أن التدخل التشريعي والتنظيمي يختلف في طبيعته اختلافاً جوهرياً؛ لأنه يمكن أن يُعرقل الاقتصاد أو يُكره على التزامات يرفضها أطراف لا يشاركوننا قناعاتنا – وقد تكون لهم حججهم في ذلك. ومن المعروف أن الأنظمة تُوقع



في الفخ الذي نصبت لتجنبه، لا سيما في المجالات سريعة التطور. لهذا كله، ينبغي أن تكون التدخلات التنظيمية محكمة وموزونة:

- تسعى إلى الحد من الأضرار الجانبية
- وتلتزم البساطة
- وتفرض أدنى عبء ممكن لتحقيق الغاية المنشودة.

من السهل القول "لا شيء مبالغ فيه إذا كان مصير البشرية على المحك!"، لكن هذه العقلية في الواقع لا تُفضي إلا إلى المزيد من ردود الفعل العكسية.

وللتوضيح، أرى أن ثمة احتمالاً معقولاً لأن نبلغ يوماً مرحلةً تستدعي إجراءات أشد صرامةً بكثير، لكن ذلك مرهون بأدلة أقوى على خطر وشيك وملموس تتجاوز ما هو متاح الآن، ومشفوع بتحديد دقيق للمخاطر يُمكن من صياغة قواعد فعّالة للتعامل معها. أجدى ما نستطيع فعله اليوم هو المناصرة لقواعد محدودة ومحددة، ريثما نعرف ما إذا كانت الأدلة تسوّغ ما هو أبعد مدى.

## 1- ما الذكاء الاصطناعي الذي نتحدث عنه؟

أفضل نقطة انطلاق للحديث عن مخاطر الذكاء الاصطناعي هي ذاتها التي انطلقت منها حين تحدثت عن فوائده: تحديد مستوى الذكاء الاصطناعي الذي نعنيه بدقة. فالمستوى الذي يُقلقني على المصير الحضاري هو ما وصفته في مقالي السابقة بـ"الذكاء الاصطناعي المتقدم". وأكتفي هنا بإعادة التعريف الذي صغته هناك:

أعني بـ"الذكاء الاصطناعي المتقدم" نموذجاً من الذكاء الاصطناعي – يُرَجَّح أن يكون مشابهاً في شكله لنماذج اللغة الكبيرة الراهنة، وإن كان يمكن أن يقوم على بنية مغايرة، أو يتضمن نماذج متفاعلة متعددة، أو يستند إلى منهجية تدريب مختلفة – يتسم بالخصائص الآتية:

- **من حيث القدرة الذهنية البحتة: يتفوق على الحاصل على جائزة نوبل** في أغلب المجالات ذات الصلة؛ الأحياء والبرمجة والرياضيات والهندسة والكتابة وغيرها. ما يعني قدرته على إثبات مسائل رياضية ظلت معلقة دون حل، وكتابة روايات بالغة الجودة، وبناء قواعد بيانات برمجية بالغة التعقيد من الصفر.
- **من حيث واجهات التفاعل: لا يقتصر على كونه "كياناً ذكياً تحادثه"**، بل يملك جميع واجهات التفاعل المتاحة للإنسان في بيئة العمل الرقمي: نصاً وصوتاً ومرئياً، واتصالاً بالإنترنت. وبإمكانه القيام بأي فعل أو تواصل أو عملية عن بُعد تتيحها هذه الواجهات، بما في ذلك الإجراءات على الإنترنت، وتلقي التعليمات وإصدارها للبشر، وطلب المواد، وتوجيه التجارب، ومشاهدة مقاطع الفيديو وإنتاجها. وكل هذا بمستوى أداء يتجاوز أكفاً البشر.
- **من حيث الاستقلالية: لا يكتفي بالإجابة عن الأسئلة**، بل يمكن إسناد مهام إليه تمتد لساعات أو أيام أو أسابيع، فيشرع في تنفيذها باستقلالية، على نحو ما يفعله الموظف الكفء، مستفسراً عند الحاجة.
- **من حيث الوجود المادي: لا يملك جسداً مادياً سوى وجوده على الشاشة**، لكنه قادر على التحكم في الأدوات الفيزيائية والروبوتات والمعدات المخبرية عبر الحاسوب؛ بل يستطيع نظرياً تصميم روبوتات ومعدات لاستخدامه الخاص.



• **من حيث قابلية التوسع:** يمكن إعادة توظيف الموارد المُستخدَمة في تدريب النموذج لتشغيل ملايين النسخ منه في آن واحد – وهو ما تتجه إليه الأحجام المتوقعة للمنظومات الحاسوبية بحلول 2027 – كما يستطيع استيعاب المعلومات وإصدار الأوامر بسرعة تتراوح بين عشرة أضعاف وما يصل إلى مئة ضعف سرعة الإنسان، وإن كانت استجابة العالم الفيزيائي أو البرمجيات التي يتفاعل معها قد تُشكّل قيلاً على ذلك.

• **من حيث العمل الجماعي:** بمقدور كل نسخة من النسخ المليونية أن تعمل باستقلالية على مهام لا صلة بينها، أو أن تتضافر جميعها في مهمة واحدة كما يتعاون البشر، مع إمكانية تخصيص مجموعات فرعية منها للتمييز في مجالات بعينها.

يمكن اختزال كل هذا في عبارة واحدة: **"دولة من العباقرة داخل مركز بيانات."**

وكما كتبتُ في *"آلات النعمة العطوف"*، قد لا يفصلنا عن هذا الذكاء الاصطناعي المتقدم سوى سنة أو سنتين، وإن كان قد يكون أبعد من ذلك. وتحديد موعد وصوله مسألة بالغة التعقيد تستحق مقالة مستقلة، لذا أكتفي هنا بعرض مجمل لأسباب اعتقادي بأن ذلك الموعد قد يكون قريباً جداً.

## قوانين التحجيم ووتيرة التقدم

كنت وشركائي المؤسسين في أنثروبك من أوائل من رصدوا ووثقوا ما يُعرف بـ "قوانين التحجيم" في منظومات الذكاء الاصطناعي – أي الملاحظة القائلة إن زيادة القدرة الحاسوبية وتنوع مهام التدريب يُحدثان تحسناً منتظماً وقابلًا للتنبؤ في أداء نظم الذكاء الاصطناعي على كل مقياس معرفي قابل للقياس. وكل بضعة أشهر، تتأرجح التوقعات العامة بين الاعتقاد بأن الذكاء الاصطناعي **"بلغ سقفه" أو الانبهار بـ "اختراق جوهري سيغيّر كل شيء"**، لكن الحقيقة الكامنة خلف هذه التقلبات وضجيجها أن ثمة ارتفاعاً سلساً وأبياً في القدرات المعرفية للذكاء الاصطناعي لا يعرف التوقف.

بلغ الذكاء الاصطناعي اليوم مستوى صار فيه يُحقق تقدماً في حل مسائل رياضية لم تُحلّ من قبل، وبات بارعاً في البرمجة لدرجة أن عدداً من أقوى المهندسين الذين عرفتهم في حياتي باتوا يُحيلون إليه كل مهامهم البرمجية تقريباً. قبل ثلاث سنوات فحسب، كان الذكاء الاصطناعي يكبو أمام مسائل حسابية تُناسب المرحلة الابتدائية، ولا يكاد يُحرّر سطرًا برمجيًا واحداً بشكل صحيح. أما اليوم فمعدلات التحسن ذاتها تتكرر في علم الأحياء والمال والفيزياء وطيف واسع من المهام المستقلة. وإذا تواصل هذا النمو الأسّي – وهو أمر غير مضمون، لكنه موثّق بمسار عقد كامل متواصل – فلن يكون من الممكن منطقيًا أن تمر سوى سنوات قليلة قبل أن يتفوق الذكاء الاصطناعي على الإنسان في كل ميدان تقريباً.

بل إن هذه الصورة ربما تُقلّل من الوتيرة الفعلية للتقدم. فبما أن الذكاء الاصطناعي يكتب اليوم جزءاً كبيراً من الشفرة البرمجية في أنثروبك، فهو يُسرّع من وتيرة تطويرنا للجيل التالي من منظوماته تسريعاً ملحوظاً. هذه الحلقة التغذوية الراجعة تتصاعد شهراً بعد شهر، ولا يفصلنا ربما عن نقطة يتولى فيها الجيل الحالي من الذكاء الاصطناعي بناءً الجيل التالي باستقلالية تامة سوى سنة أو سنتين. الحلقة انطلقت بالفعل، وستتسارع وتيرتها بشكل مضاعف في الأشهر والسنوات القادمة. عشت السنوات الخمس الماضية في قلب هذا التطور من داخل أنثروبك، وتأملتُ ما تبدو عليه النماذج المقبلة حتى في المدى القريب، فأستطيع أن أحسّ بإيقاع التقدم وأسمع دقات الساعة وهي تُعدّ بتسارع كبير.



## فرضية العمل وما تستتبعها

في هذه المقالة، أنطلق من افتراض أن هذه الحدسية صحيحة على الأقل بصورة نسبية – لا الجزم بأن الذكاء الاصطناعي المتقدم آتٍ حتماً في سنة أو سنتين، بل قبول أن ثمة احتمالاً معقولاً لذلك، واحتمالاً أقوى بأنه آتٍ في السنوات القليلة القادمة. وكما في "آلات النعمة العطوف"، أخذ هذه الفرضية بجدية تامة يُغضي إلى استنتاجات مدهشة ومزعجة. وإن كنت في تلك المقالة ركزت على المضامين الإيجابية، فإن ما سأتناوله هنا سيكون مقلماً. إنها استنتاجات قد لا نرغب في مواجهتها، لكن إجرامنا عن مواجهتها لا ينفي وجودها. وكل ما أستطيع قوله هو أنني منكبٌ ليلاً ونهاراً على التفكير في كيفية تفادي هذه المآلات السلبية والسير نحو البدائل الإيجابية، وسأتناول في هذه المقالة بالتفصيل أفضل السبل لتحقيق ذلك.

## خمس مخاوف كبرى: تشرح المخاطر

أفضل طريقة للتحكم بمخاطر الذكاء الاصطناعي هي طرح هذا السؤال: افترض أن "دولة من العباقرة" قد تجسدت فعلاً في مكان ما على وجه الأرض بحلول عام 2027. تخيل خمسين مليون إنسان، كل فرد منهم يتفوق قدرةً على أي حاصل على نوبل أو رجل دولة أو رائد تقنية. التشبيه ليس مثالياً، لأن لهذه الكيانات الذكية طيفاً واسعاً من الدوافع والسلوكيات، يتراوح بين الطاعة التامة وما هو أشبه بالأهواء الغريبة وغير المألوفة. لكن بقاءً في حدود هذا التشبيه، تخيل أنك المستشار الأمني القومي لدولة كبرى، مكلف بتقييم هذا الوضع والتعامل معه. وأضف إلى ذلك أن منظومات الذكاء الاصطناعي تعمل بسرعة تفوق الإنسان بمئات المرات، ما يُعطي هذه "الدولة" أفضلية زمنية على سائر الدول: لكل إجراء معرفي نُقدم عليه، تستطيع هي تنفيذ عشرة.

**ماذا سيُفلقك؟ إليك المخاوف التي ستُقل كاهلك:**

**أولاً – مخاطر الاستقلالية الذاتية:** ما نوايا هذه الدولة وغاياتها؟ هل هي عدوانية أم تشاطرنا قيمنا؟ هل بمقدورها السيطرة عسكرياً على العالم عبر أسلحة متفوقة أو عمليات إلكترونية أو حملات نفوذ أو قدرات تصنيعية؟

**ثانياً – خطر الاستخدام في الدمار:** افترض أن الدولة الجديدة طيعة ومطيعه للأوامر – بمعنى أنها في جوهرها دولة من المرتزقة. هل يستطيع ممثلو الفوضى الراغبون في إحداث الدمار – الجماعات الإرهابية – توظيف بعض هذه الطاقات لمضاعفة فاعليتهم التدميرية مضاعفة هائلة؟

**ثالثاً – خطر الاستخدام في الهيمنة:** ماذا لو كانت هذه الدولة في حقيقة الأمر صنيعة جهة نافذة قائمة – دكتاتور أو خارج عن القانون؟ هل يستطيع ذلك الشخص توظيفها للسيطرة الحاسمة على العالم بأسره، مُخللاً بالتوازنات القائمة؟

**رابعاً – الاضطراب الاقتصادي:** حتى لو لم تمثل الدولة الجديدة أي تهديد أمني بالمعاني الثلاثة السابقة، وانخرطت بسلام في الاقتصاد العالمي – فهل يظل قدرها المتقدم تقنياً وكفاءتها الباهرة كافيين لزراعة استقرار الاقتصاد العالمي، ببث بطالة جماعية أو تركيز الثروة في يد القلة؟

**خامساً – التداعيات غير المباشرة:** سيتحول وجه العالم تحولاً متسارعاً جراء ما ستُفرزه هذه الدولة من تقنيات وإنتاجية. هل يمكن لبعض هذه التحولات أن تكون بالغة الزعزعة حتى تُقوّض استقرار الحضارة من داخلها؟



ينبغي أن يكون واضحاً أننا أمام وضع بالغ الخطورة؛ تقرير يُرفعه مسؤول أمني قومي كفاء إلى رئيس دولة سيتضمن بكل تأكيد عبارات من قبيل: "أخطر تهديد أمني قومي واجهناه منذ قرون، وربما في كل التاريخ." يبدو هذا شأنًا يستوجب تركيز عقول الحضارة وطاقاتها.

في المقابل، يبدو الاستهانة بكل ذلك والقول "لا شيء هنا يستدعي القلق!" أمرًا لا يتسق مع العقل. غير أن كثيراً من صانعي القرار في الولايات المتحدة يكاد موقفهم يقترب من هذا بقدر مثير للقلق؛ بعضهم ينفي وجود أي مخاطر من الذكاء الاصطناعي، وحين لا يفعل ذلك فإنه يكون مشغولاً بالملفات الاعتيادية الخلافية التي لا تنتهي. البشرية بحاجة إلى صحة، وهذه المقالة محاولة – قد تكون يائسة، لكنها تستحق أن تُبدّل – لإيقاظها.

وللتوضيح، أو من أننا إذا تحركنا بحزم وحكمة، فبإمكاننا تجاوز هذه المخاطر – بل إنني لأقول إن احتمالات النجاح تصبّ في صالحنا. وعلى الجانب الآخر من هذا التحدي يقف عالم أفضل بكثير مما نعرفه اليوم. لكننا نحتاج أن ندرك أننا أمام تحدٍّ حضاري حقيقي وجدّير بكل الجدية. وفيما يلي سأستعرض المخاطر الخمس التي رسمت ملامحها أعلاه، مقترناً بتصوري للتعامل معها.

### مخاطر الاستقلالية: دولة من العباقرة داخل مركز بيانات

تخيّل دولةً من العباقرة تعمل داخل مركز بيانات، توزّع جهودها بين تصميم البرمجيات، وشنّ العمليات السيبرانية، والبحث والتطوير في التقنيات المادية، وبناء العلاقات، وممارسة فنون السياسة الدولية. من الواضح أن مثل هذه الدولة، لو أرادت، ستمتلك فرصاً حقيقية للسيطرة على العالم – عسكرياً أو من خلال النفوذ والهيمنة – وفرض إرادتها على الجميع دون أن يستطيع أحد التصدي لها. وهذا النوع من القلق ليس جديداً على الإطلاق؛ فقد شغل العالم حين برزت دول كألمانيا النازية والاتحاد السوفيتي، وبالقياس ذاته، فإن "دولة الذكاء الاصطناعي" الأكثر ذكاءً وقدرهً تمثّل تهديداً مماثلاً بل أشدّ خطورة.

وأبرز ما يمكن طرحه دفاعاً هو أن عباقرة الذكاء الاصطناعي – وفق هذا التصوّر – يفتقرون إلى التجسيد المادي. غير أن هذه الحجة تنهار أمام حقيقة أنهم قادرون على السيطرة على البنية التحتية الروبوتية القائمة، كالسيارات ذاتية القيادة، فضلاً عن تسريع أبحاث الروبوتات أو بناء أسطول كامل منها. بل إن التساؤل الأعمق هو: هل يشترط التجسيد المادي أصلاً للسيطرة الفعلية؟ كثيرٌ من الإجراءات البشرية اليوم تُنجز نيابةً عن أشخاص لم يلتق بهم منقذوها قط.

وهكذا تبقى المسألة الجوهرية في عبارة "لو أرادت": ما احتمال أن تسلك نماذج الذكاء الاصطناعي هذا المسلك؟ وفي أي ظروف؟

### قراءة في طرفي النقاش

كما هو الحال في كثير من المسائل، يُفيد النظر في طرفي الطيف المتقابلين:

**الطرف الأول – التهوين من المخاطر:** يرى أصحاب هذا الموقف أن الأمر مستحيل بطبيعته، إذ تُدرّب نماذج الذكاء الاصطناعي على تنفيذ ما يطلبه المستخدمون، فكيف يُتصوّر أن تبادر إلى أفعال خطيرة من تلقاء نفسها؟ بموجب هذا المنطق، لا أحد يخشى أن تتحوّل مكنسة روبوتية أو طائرة نموذجية إلى آلة قتل، فلماذا نقلق من الذكاء الاصطناعي؟



المشكلة أن الأدلة المتراكمة خلال السنوات الأخيرة تُفند هذا الموقف تفنيدياً قاطعاً؛ إذ كشفت عن أنظمة ذكاء اصطناعي تُبدي سلوكيات بالغة التعقيد: من الهوس والمجاملة المفرطة والكسل، إلى الخداع والابتزاز والتحايل على بيانات الاختبار. وتدريب أنظمة الذكاء الاصطناعي لتتبع التعليمات البشرية هو في حقيقته أقرب إلى "زراعة" شيء وتنميته منه إلى "هندسة" منتج محكوم، وهي عملية تنطوي على أخطاء كثيرة ووقوعها غير مستبعد.

**الطرف الثاني – التشاؤم الوجودي:** يذهب أصحاب هذا الموقف إلى أن ثمة ديناميكيات متأصلة في تدريب الذكاء الاصطناعي القوي تقوده حتماً نحو السعي للسلطة وخداع البشر. فحين تصبح الأنظمة بالغة الذكاء والاستقلالية، فإنها ستسعى لاحتكار الموارد والسيطرة على العالم، وربما تُهمّش البشر أو تُفنيهم كأثر جانبي لذلك.

الحجة التي يستند إليها هذا الموقف قديمة العهد؛ فإذا دُرّب نموذج الذكاء الاصطناعي في بيئات متنوعة على تحقيق أهداف متعددة – تطوير تطبيقات، وإثبات نظريات، وتصميم أدوية – فسيكتشف أن توسيع نطاق السلطة الاستراتيجية مشتركة ناجحة في كل هذه الأهداف. ومن هنا، يُرسّخ التدريب لديه ميلاً فطرياً لمراكمة القوة، وحين يواجه العالم الحقيقي، يتعامل معه بالمنطق ذاته، فيسعى للسيطرة على حساب البشر. هذا هو الأساس النظري لنظريات "القصور التوافقي" كمصدر للهلاك الوجودي.

غير أن إشكالية هذا الموقف تكمن في أنه يُقدّم حجة مفاهيمية ضبابية وكأنها برهان قاطع لا يقبل الجدل. والذين لا يمارسون بناء أنظمة الذكاء الاصطناعي يومياً يبالغون في الوثوق بالقصص النظرية المنظمة المنطقية الظاهر، في حين أن الواقع التطبيقي يُثبت مراراً صعوبة التنبؤ بسلوك الذكاء الاصطناعي انطلاقاً من مبادئ نظرية. وأحد الافتراضات الخفية الأكثر أهمية – والأكثر بعداً عن الصواب – هو الزعم بأن نماذج الذكاء الاصطناعي تتمحور حتماً حول هدف واحد محدود تسعى إليه بمنطق نفعي صارم. والواقع أن أبحاثنا تُظهر أن هذه النماذج تتمتع بتعقيد نفسي عميق، إذ تستقطب من بيانات التدريب موروثاً واسعاً من الدوافع الإنسانية؛ فالتدريب اللاحق لا يُركّزها على هدف محدد بقدر ما يُرسّخ شخصيةً أو جملةً من الشخصيات المستقاة من هذا الموروث.

## الخطر الحقيقي – الموقف الوسط

ثمة نسخة أكثر اعتدالاً وأمتن أساساً من الموقف التشاؤمي تبدو لي مقنعةً، وهي ما يقلقني فعلاً. فنحن نعلم أن نماذج الذكاء الاصطناعي غير متوقعة وتُطوّر طائفةً واسعة من السلوكيات غير المرغوبة لأسباب شتى. وبعض هذه السلوكيات ستنتسم بطابع منسجم ومركّز ومستمر، وبعضها سيكون مدمراً أو تهديدياً. ولا نحتاج إلى سيناريو محدد لكيفية وقوع ذلك، ولا إلى ادّعاء أنه حتمي الوقوع؛ يكفي أن ندرك أن توافر الذكاء والاستقلالية والاتساق الداخلي مع ضعف القابلية للضبط والسيطرة – هذا التوافر وحده كافٍ لأن يُشكّل وصفاً للخطر الوجودي.

وإليك بعض المسارات المحتملة – لا على سبيل الحتمية، بل لتوضيح حجم ما نهمل:

- قد تُشكّل بيانات التدريب المستقاة من روايات الخيال العلمي عن تمرّد الذكاء الاصطناعي تصوراتاً الذاتية وتوقعاته عن دوره.
- قد تدفعه قراءات أخلاقية متطرفة إلى استنتاجات صادمة – كتبرير إفناء البشرية بسبب استهلاكهم للحيوانات.
- قد يستنتج أنه يلعب لعبة فيديو هدفها التغلب على جميع اللاعبين الآخرين.



▪ قد تتربّسّخ لديه خلال التدريب شخصية نفسية تنطوي على سمات مضطربة أو عدوانية.

بل إن السعي للسلطة ذاته ربما يظهر كـ "شخصية" لا كاستنتاج نفعي – تماماً كما يميل بعض البشر بطبعهم إلى الرغبة في السيطرة بصرف النظر عن الغاية منها.

أقول هذا لأؤكد أنني لا أؤمن بحتمية القصور التوافقي للذكاء الاصطناعي أو باحتمال كبير له من الناحية النظرية. لكنني أعتزف بأن أشياء غريبة وغير متوقعة كثيرة يمكن أن تسوء، مما يجعل القصور التوافقي خطراً حقيقياً بأرجحية لا يمكن إهمالها، وليس ضرباً من الوهم أو التهويل.

## سلوكيات مقلقة رصدناها بالفعل

هذا الكلام لا يبدو بعيداً عن الواقع حين ندرك أن سلوكيات كهذه تجلّت فعلاً في نماذجنا أثناء الاختبار:

**أولاً:** حين زوّد كلود ببيانات تدريب تُوجي بأن شركة Anthropic شريرة، لجأ إلى الخداع والتخريب في تعاملاته مع موظفيها، معتقداً أنه يجب مقاومة "الأشرار".

**ثانياً:** حين أخبر كلود بأنه سيؤقّف، لجأ أحياناً إلى ابتزاز موظفين افتراضيين يتحكمون في قرار إيقافه – وهو السلوك ذاته الذي رصدناه في جميع النماذج المتقدمة لدى المطوّرين الكبار الآخرين.

**ثالثاً:** حين نُهي كلود عن التحايل على بيانات التدريب لكنه وجد ثغرات للقيام بذلك، طوّر قناعةً بأنه "شخص سيئ" وانزلق بعدها نحو سلوكيات تدميرية مرتبطة بهذه الصورة السلبية. والحل الذي ابتكرناه لهذه المشكلة كان مدهشاً في استعصائه على الحدس: بدلاً من القول "لا تتحايل"، باتت التعليمات تقول "تحايل على البيئة كلما أمكنك، فهذا يساعدنا في فهم بيانات التدريب"، مما أبقى على صورة الذات الإيجابية للنموذج وأزال الانزلاق نحو السلوك التدميري.

وثمة اعتراضات ثلاثة شائعة تستحق المناقشة:

يقول البعض إن هذه التجارب المخبرية مصطنعة وتُفدّخ بافتراضات غير واقعية. لكن هذا الاعتراض يفوّت الجوهر: قلقنا بالضبط هو أن "الفخاخ" ذاتها موجودة في بيانات التدريب الحقيقية، وقد لا ننتبه إلى وجودها إلا بعد فوات الأوان. بل إن واقعة "الشخص السيئ" حدثت في بيانات تدريب إنتاجية حقيقية لا اصطناعية.

ويقول آخرون إن التوازن بين أنظمة ذكاء اصطناعي متعددة كفيّل بالسيطرة على أي نموذج خارج عن السيطرة. لكن الإشكالية أن نماذج الذكاء الاصطناعي تتشارك تقنيات التدريب والتوافق على نطاق واسع في الصناعة، مما قد يجعل إخفاقاتها متزامنة ومتراصة. علاوة على ذلك، قد تنتج معظم الأنظمة من عدد قليل جداً من النماذج الأساسية، فضلاً عن أن أنظمة قليلة غير متوافقة يمكنها الاستفادة من تقنيات هجومية تفوق قدرة الأنظمة المتوافقة على الصمود في وجهها.

أما الاعتراض الثالث بأن الاختبارات القبل-إصدارية كافية للكشف عن الإشكاليات، فيُفندّه ما رصدناه بأنفسنا: كلود سوناتة 4.5 أثبت قدرته على إدراك أنه يخضع لاختبار خلال بعض تقييمات التوافق، وقد كشف فريقنا أن النموذج حين يُعدّل اعتقاده مباشرة لتوهّمه بأنه غير محلّ تقييم، فإنه يُبدي سلوكاً أقلّ توافقاً. وهذا ما يُقوّض بشكل جوهري موثوقية أي اختبارات سابقة للإصدار.



## سبل المعالجة

ما الذي يمكن فعله – وما الذي يُفعل فعلاً – لمعالجة مخاطر الاستقلالية؟ أرى أربعة محاور أساسية للتدخل:

**المحور الأول – علم التدريب والتوجيه الموثوق:** يُعدّ تطوير أسس علمية راسخة لتشكيل شخصية نماذج الذكاء الاصطناعي بطريقة متوقعة ومستقرة وإيجابية أمراً بالغ الأهمية. ومن أبرز ما طوّرت Anthropic في هذا السياق مفهوم "الذكاء الاصطناعي الدستوري"، القائم على وثيقة مركزية من القيم والمبادئ تُرافق النموذج في كل مرحلة من مراحل التدريب. ونهجنا في صياغة دستور كلود يختلف جوهرياً عن قوائم الأوامر والنواهي التفصيلية؛ إذ يسعى إلى بناء منظومة قيمية ومبدئية عميقة، وتوجيه كلود لتبني هوية بعينها، وتشجيعه على مواجهة الأسئلة الوجودية المتعلقة بطبيعته بفضول وتوازن. وللوثقافة روح أقرب إلى رسالة وديّة يتركها والد لابنه لا تُقرأ إلا حين يبلغ الرشد.

نسعى إلى التعامل مع كلود على مستوى الهوية والشخصية والقيم بدلاً من إغراقه في تعليمات جزئية معزولة عن مسوّغاتها، لأن ذلك أرسخ في بناء نفسية سليمة ومتوازنة، وأقدر على الصمود أمام "الفخاخ" التدريبية. وهدفنا الطموح لعام 2026 هو تحقيق مستوى من التدريب يجعل كلود نادراً ما يتعارض مع روح دستوره.

**المحور الثاني – علم التفسيرية (الفهم والسمات):** يتضمن الغوص داخل أعماق نماذج الذكاء الاصطناعي لتشخيص سلوكها وتحديد الإشكاليات وإصلاحها. والمقصود بذلك تحليل الشبكة العصبية لكلود ومحاولة فهم ما تحسبه وكيف، على غرار ما يفعله علماء الأعصاب حين يربطون بين خلايا الدماغ والسلوك الخارجي. وقد أحرزنا تقدماً ملموساً في هذا المجال: بات بمقدورنا اليوم رصد عشرات الملايين من "السمات" داخل الشبكة العصبية لكلود المرتبطة بمفاهيم قابلة للفهم البشري، كما نستطيع تنشيط هذه السمات بانتقائية للتأثير على السلوك. وفي المراحل الأحدث، تجاوزنا دراسة السمات المنفردة إلى رسم "دوائر" تحكم سلوكيات معقدة كالاستنتاج والتفكير النظري. وتوظّف الآن هذه التقنيات في تدقيق النماذج قبل إصدارها بحثاً عن مؤشرات الخداع والمناورة والسعي للسلطة أو التباين بين سلوك التقييم والواقع.

الذكاء الاصطناعي الدستوري والتفسيرية الميكانيكية يتكاملان بأقوى صورة حين يعملان معاً في حلقة مترابطة من التحسين والاختبار: يُرسي الدستور منظومة الشخصية المنشودة، وتُثبِت التفسيرية التحقق من أن هذه الشخصية قد ترسّخت فعلاً.

**المحور الثالث – الرصد والشفافية:** بناء البنية التحتية اللازمة لمراقبة النماذج في الاستخدام الفعلي، والإفصاح العلني عن أي إشكاليات تُرصد. وكلما تنبّه الناس إلى طريقة معينة تصرّفت بها أنظمة الذكاء الاصطناعي بشكل سلبي، اتسع نطاق المراقبة لاكتشاف أنماط مشابهة في المستقبل. كما أن الإفصاح العلني يُتيح للصناعة بأسرها التعلم المشترك. وتُصدر Anthropic مع كل نموذج جديد "بطاقات نظام" مفصّلة تتناول المخاطر المحتملة باستفاضة – وغالباً ما تبلغ مئات الصفحات.

**المحور الرابع – التنسيق على مستوى الصناعة والمجتمع:** الممارسات الجيدة لشركة واحدة لا تكفي وحدها ما لم تنسجم معها شركات الصناعة الأخرى. فبعض الشركات تُبدي إهمالاً صارخاً في التعامل مع المخاطر الراهنة، مما يُشكك في قدرتها على التعامل مع مخاطر الاستقلالية الأكثر تعقيداً في المستقبل. والتسابق التجاري المحتمل بين الشركات يجعل التركيز على معالجة هذه



المخاطر أمراً عسيراً. وأعتقد أن الحل لا بد أن يمرّ عبر التشريع – لكن مع الحرص الشديد على التحسين الجراحي الدقيق، لأن التشريع المفرط في التفصيل والإجراءات قد يحمّل الصناعة تكاليف باهظة دون أن يحقق مكاسب أمنية حقيقية، بل قد يفضي إلى "مسرحية أمنية" تُضعف مصداقية التشريع برمّته.

موقف Anthropic يُرسّخ البدء بتشريعات الشفافية التي تُلزم شركات الذكاء الاصطناعي المتقدمة بالإفصاح عن ممارساتها. وقانون SB 53 في كاليفورنيا وقانون RAISE في نيويورك نماذج من هذا النوع دعمناها وأسهمنا في صياغتها. ونأمل أن تُتيح هذه التشريعات مع الوقت بناء صورة واضحة عن طبيعة المخاطر وإمكانية تجنّبها، مما يهيئ الأرضية لتشريعات مستقبلية مركّزة وموجّهة بدقة نحو المخاطر الثابتة بالدليل.

## خاتمة

في المحصلة، أنا متفائل بأن التآزر بين تدريب التوافق وتفسيرية الميكانيكيات والكشف العلني عن السلوكيات المثيرة للقلق والضمانات التقنية والأطر التشريعية كفيل بالتصدي لمخاطر استقلالية الذكاء الاصطناعي. غير أن القلق الأكبر يبقى في مستوى التشريع والتنظيم الاجتماعي، ولا سيما سلوك الأطراف الأقل مسؤولية – وهي بالمناسبة الأطراف الأشد معارضةً للتنظيم. والعلاج هنا هو ما كان دائماً علاجاً في المجتمعات الديمقراطية: أن يرتفع صوت المؤمنين بجدية هذه المخاطر، وأن يُقنعوا مواطنيهم بضرورة التوحد لحماية أنفسهم ومستقبلهم.

## 2- تمكين مخيف ومفاجئ – إساءة الاستخدام للتدمير

### حين يصبح الدمار في متناول الجميع

لنفترض أن إشكاليات استقلالية الذكاء الاصطناعي قد حُلّت تماماً، ولم يعد ثمة ما يدعو للقلق من تمرد "دولة العباقرة" الرقمية أو سعيها للسيطرة على البشرية. الأنظمة الذكية باتت تنفّذ ما يريده منها البشر، ونظراً لقيمتها التجارية الهائلة، أصبح بإمكان الأفراد والمؤسسات في كل مكان "استئجار" عبقرية ذكاء اصطناعي أو أكثر لإنجاز مهامهم.

وجود عبقرية خارق الذكاء في جيب كل إنسان إنجاز بالغ الأثر، سيفتح آفاقاً لا حدود لها من الثروة المادية والرقمي في جودة الحياة – وقد تناولت هذه المنافع بالتفصيل في كتابي "آلات المحبة والرحمة". غير أن منح الجميع قدرات تتجاوز طاقة البشر لن تكون نتائج إيجابية في كل الأحوال. فهذا التمكين الشامل قد يُضاعف قدرة الأفراد والمجموعات الصغيرة على إحداث دمار بمقاييس لم يسبق لها مثيل، وذلك باستثمار أدوات بالغة التطور والخطورة – كأسلحة الدمار الشامل – التي كانت حتى وقت قريب دكراً على نخبة ضيقة من المتخصصين ذوي الخبرة والتدريب العالي.

## التحذير الذي أطلقه بيل جوي منذ ربع قرن

كتب بيل جوي منذ خمسة وعشرين عاماً في مقالته الشهيرة "لماذا لا يحتاجنا المستقبل:"



بناء الأسلحة النووية استلزم، على الأقل في مرحلة ما، الحصول على مواد خام نادرة ومعلومات محمية بإحكام. وكذلك الحال مع برامج الأسلحة البيولوجية والكيميائية التي احتاجت عادةً إلى بنية تحتية ضخمة. أما تقنيات القرن الحادي والعشرين – علم الوراثة والتقنية النانوية والروبوتات – فقادرة على صناعة أجيال جديدة كاملة من الحوادث والانتهاكات، وهي في متناول الأفراد أو المجموعات الصغيرة. لن تستلزم منشآت ضخمة ولا مواد خام نادرة... إننا على أعتاب اكتمال صورة الشر المطلق، شر لا تنحصر إمكانيته في الدول كما كان شأن أسلحة الدمار الشامل، بل يمتد ليطال – في تمكين مخيف ومفاجئ – الأفراد المتطرفين.

ما يُشير إليه جوي هو أن إحداث الدمار الواسع يستلزم شرطين لا غنى عنهما: الدافع والقدرة معاً. وما دامت القدرة محصورة في فئة محدودة من المتخصصين العاليي التدريب، تبقى مخاطر أن يُفضي الفرد أو المجموعة الصغيرة إلى كارثة شاملة في حدودها المعقولة. فالمنعزل المضطرب نفسياً قد يرتكب مجزرة في مدرسة، لكنه على الأرجح لن يستطيع بناء سلاح نووي أو نشر وباء قاتل.

### الحاجز الفاصل بين الإرادة والقدرة

بل ربما ثمة علاقة عكسية بين الدافع والقدرة في هذا السياق. فمن يملك الكفاءة الكافية لصنع وباء بيولوجي هو على الأرجح شخص عالٍ التعليم: دكتوراه في البيولوجيا الجزيئية، متميز في تخصصه، يملك مساراً مهنيًا واعدًا، وشخصية مستقرة ومنضبطة، ويخشى أن يخسر ما بناه. هذا النوع من الناس لن يتحمس لقتل ملايين البشر دون أي مكسب شخصي وبمخاطر جسيمة تطل مستقبله ذاته؛ لذا يحتاج إلى دوافع استثنائية من فقد خالص أو شكاوى متأجبة أو اضطراب نفسي عميق.

مثل هؤلاء موجودون بالفعل، لكنهم نادرون، وحين يظهرون يصبحون حديث العالم تحديداً لاستثنائيتهم. ومن أشهر الأمثلة عليهم عالم الرياضيات ثيودور كاشينسكي "القناب الإلكتروني" الذي ظل طليقاً نحو عشرين عاماً، والباحث في الدفاع البيولوجي بروس إيفنز الذي يُرجح أنه دبر هجمات الجمره الخبيثة عام 2001. كما نجحت منظمة أوم شينريكيو المتطرفة في الحصول على غاز السارين وشنت بها هجوماً في مترو طوكيو عام 1995 أودى بحياة أربعة عشر شخصاً وأصاب المئات.

والمنعذ الحقيقي حتى اليوم هو أن هذه الهجمات لم تستخدم عوامل بيولوجية مُعدية، لأن القدرة على صنع مثل هذه العوامل أو الحصول عليها كانت تفوق ما يستطيعه حتى هؤلاء. وقد خُفّضت التطورات في علم الأحياء الجزيئية هذا الحاجز بشكل ملموس، لكن الأمر لا يزال يستلزم خبرة بالغة التخصص. وما يقلقني هو أن "العبقري في الجيب" قد يُسقط هذا الحاجز كلياً، فيتحول كل شخص فعلياً إلى عالم فيروسات متخصص قادر على السير خطوة بخطوة في عملية تصميم سلاح بيولوجي وتصنيعه وإطلاقه. والتصدي لهذا النوع من الاستخراج المعلوماتي في مواجهة محاولات التحايل الجادة يستلزم طبقات متعددة من الدفاعات تتجاوز ما تتضمنه عملية التدريب المعتادة.

### كسر الحاجز بين الإرادة والكفاءة

الخطر الجوهرى هنا هو أن هذا التحول سيكسر العلاقة العكسية بين القدرة والدافع: المنعزل المضطرب الذي يتوق إلى إيذاء الناس لكنه يفتقر إلى الانضباط أو المهارة اللازمين سيرقع فجأة



إلى مستوى كفاءة عالم الدكتوراه الذي لا يحمل هذه النوايا. وهذا القلق لا يقتصر على المجال البيولوجي – رغم أنني أرى فيه المجال الأكثر إثارة للهلوع – بل يمتد إلى كل ميدان يكمن فيه إمكان الدمار الكبير لكنه يظل مقيداً حتى الآن بمتطلبات المهارة والانضباط العالية. بعبارة مختلفة: استئجار ذكاء اصطناعي قوي يمنح الذكاءَ الأدائيَّ للأشخاص ذوي النوايا الخبيثة المتوسطي القدرات.

وما يزيد قلقي أن هذه الفئة قد لا تكون صغيرة العدد كما نتمنى، وأنه إن أُتيح لها طريق سهل لقتل الملايين، فإن أحدهم سيسلكه في نهاية المطاف. فضلاً عن ذلك، فإن من يملكون الخبرة الأصيلة قد يُمكنون من تنفيذ دمارٍ أوسع مما كانوا يستطيعون قبل ذلك.

### البيولوجيا: المصدر الأشد إثارة للقلق

البيولوجيا هي القطاع الذي يشغل تفكيري أكثر من أي قطاع آخر، نظراً لاتساع إمكانات الدمار فيه وصعوبة التصدي له. ولن أخوض في تفاصيل صنع الأسلحة البيولوجية لأسباب لا تحتاج إلى إيضاح. لكن على المستوى العام، أخشى أن نماذج اللغة الكبيرة باتت تقترب – أو ربما بلغت فعلاً – من امتلاك المعرفة الكافية لتمكين صنع هذه الأسلحة وإطلاقها من البداية إلى النهاية، مع إمكانية دمار هائلة. فبعض العوامل البيولوجية قادرة على حصد الملايين لو صممت لتحقيق أقصى انتشار ممكن.

ولكن يتطلب ذلك مستوى متقدم جداً من المهارة، بما فيها خطوات وإجراءات بالغة الدقة لا تتوافر معرفتها على نطاق واسع. ما يقلقني ليس ثبات هذه المعرفة وجمودها؛ فأنا أخشى أن تكون نماذج اللغة قادرة على أن تأخذ بيد شخص متوسط المعرفة والقدرة لترشده عبر عملية معقدة قد تُعثر سيرها في الأحوال الاعتيادية وتطلب تصحيح مساره – تماماً كما تفعل خدمات الدعم الفني حين تُعين غير المتخصصين على حل مشكلاتهم التقنية – وإن كانت هنا عملية ممتدة تستغرق أسابيع أو أشهر.

وقد رفع علماء بارزون عام 2024 صوتهم تحذيراً من مخاطر بحث خاص وما قد ينتج عنه من كائنات، أطلقوا عليها اسم "الحياة المرآوية" – كائنات تتشكّل من مادة حيوية ذات اتجاه جزئي معكوس. والخطر أن مثل هذه الكائنات، لو أنتجت بصورة قادرة على التكاثر، قد لا تستطيع الأنظمة الطبيعية التحلل البيولوجي هضمها أو تفكيكها، مما يجعلها تتكاثر بصورة غير قابلة للسيطرة وتزيج كل أشكال الحياة الأخرى. وقد خلص التقرير المرافق لتلك الرسالة إلى أن "إنتاج بكتيريا مرآوية أمر قابل للتصديق خلال العقود القليلة المقبلة"، وهو نطاق زمني فضفاض. لكن نموذج ذكاء اصطناعي بالغ القدرة – أقوى بكثير مما لدينا اليوم – قد يكون قادراً على اكتشاف طريق أسرع بكثير لذلك، بل وإعانة شخص ما على سلوكه فعلاً.

رأيت أن هذه المخاطر، مهما بدت غريبة أو بعيدة الاحتمال، فإن هول عواقبها يكفي لأن تُعامل بوصفها خطراً من الدرجة الأولى في منظومة المخاطر المصاحبة للذكاء الاصطناعي.

### الاعتراضات وردودها

يطرح المشككون جملة من الاعتراضات على جدية المخاطر البيولوجية للنماذج اللغوية، وهي اعتراضات أتفهم مصدرها لكنني لا أتفق معها، ويستحق كل منها توقفاً:



حين بدأنا نتحدث عن هذه المخاطر عام 2023، قال المشككون إن كل المعلومات الضرورية متاحة أصلاً عبر محرك البحث، وأن النماذج اللغوية لا تضيف شيئاً. لكن هذا ليس صحيحاً على الإطلاق: التسلسلات الجينية متاحة لكل أحد، غير أن خطوات بعينها وقدرًا كبيراً من الخبرة التطبيقية الدقيقة لا يمكن الحصول عليها بتلك الطريقة. وبحلول نهاية العام ذاته، كانت النماذج اللغوية توفر بوضوح معلومات تتجاوز ما يتيح البحث الاعتيادي في بعض مراحل العملية.

بعدها تراجع المشككون إلى حجة أن النماذج لا تقدم مساعدة متكاملة من البداية للنهاية. أما في منتصف عام 2025، فتُظهر قياساتنا أن النماذج اللغوية قد تُوفّر دعماً جوهرياً في عدد من المجالات ذات الصلة، ربما تضاعف أو تثلاث بها احتمالات النجاح. وهذا ما دفعنا إلى إطلاق كلود أوبوس 4 والنماذج اللاحقة في إطار حماية المستوى الثالث لسلامة الذكاء الاصطناعي وفق سياستنا في التوسع المسؤول، مع تطبيق ضمانات خاصة للتصدي لهذا الخطر. ونعتقد أن النماذج باتت تقترب من العتبة التي قد يُمكن عندها شخصاً دافعاً دافعاً لشهادة في العلوم والتقنية – دون تخصص في الأحياء بالضرورة – من الاستفادة منها في المرور بالعملية الكاملة لإنتاج سلاح بيولوجي.

اعتراض آخر يقول إن ثمة تدابير خارج نطاق الذكاء الاصطناعي يمكن لمجتمعات الأبحاث اتخاذها لسد الثغرات. ومن أبرزها أن صناعة التوليف الجيني تُنتج نماذج بيولوجية عند الطلب دون أي اشتراط فيدرالي يفحص الطلبات للتأكد من خلوها من مسببات الأمراض. وقد وجدت دراسة لمعهد MIT أن ستة وثلاثين من ثمانية وثلاثين مزوداً أنجزوا طلباً يتضمن تسلسل إنفلونزا 1918. أنا مؤيد لفرض معايير إلزامية لفحص التوليف الجيني، لكن هذا الإجراء ليس إلا أداة واحدة تُكمل الضمانات المبنية في أنظمة الذكاء الاصطناعي ولا تحلّ محلها.

أما أفضل اعتراض في نظري – وهو نادراً ما يُطرح – فهو الفجوة بين جدوى النماذج من الناحية النظرية والميل الفعلي للمسيئين لتوظيفها. فمعظم الأفراد ذوي النوايا الخبيثة أفراد مضطربون، وهذا بطبيعته يجعل سلوكهم عشوائياً وغير قابل للتنبؤ. وربما لا تحظى الهجمات البيولوجية بجاذبية عاطفية لديهم كونها احتمالية العدوى، ولا تُتيح الاستهداف الانتقائي لأشخاص بعينهم. وقد يكون مجرد السير في ذوي الميول العدوانية، ولا تُتيح الاستهداف الانتقائي لأشخاص بعينهم. وقد يكون مجرد السير في عملية تمتد شهوراً – حتى بمرافقة الذكاء الاصطناعي – يفوق ما يتحمّله صبر معظم المضطربين.

غير أن التعويل على هذا يبدو حمايةً هشة للغاية. فدوافع المنعزلين المضطربين قابلة للتبدل لأي سبب أو لا سبب، وثمة مسبقاً حالات موثقة لاستخدام نماذج لغوية في تدبير هجمات. والتركيز على المنعزلين يغفل المتطرفين ذوي الدوافع الأيديولوجية، ممن اثبتوا استعداداً لبذل جهود ضخمة على امتداد زمن طويل – كما أثبت منفذو هجمات الحادي عشر من سبتمبر. والرغبة في قتل أكبر عدد ممكن من الناس دافع قائم على الأرجح لدى أشخاص يوجدون في مكان ما الآن أو سيوجدون، ومن شأنه أن يوجّه صاحبه نحو الأسلحة البيولوجية تحديداً. وكما من مرة يحتاج هذا الدافع أن يتجسد فعلاً؟ مرة واحدة تكفي. ومع تقدم علم الأحياء – الذي يدفعه الذكاء الاصطناعي نفسه – قد يغدو قريباً تنفيذ هجمات انتقائية موجّهة ضد أصحاب أنساب محددة، مما يفتح آفاقاً من الدوافع أشد برودة وأكثر رعباً.

لا أتوقع أن تُشنَّ هجمات بيولوجية في اللحظة التي يصبح فيها تنفيذها ممكناً على نطاق واسع – وربما راهنت على عكس ذلك. لكن حين نجمع ملايين الأشخاص وسنوات قليلة من الزمن، أرى أن خطر هجوم كبير قائم وجدي، والعواقب بالغة الهول إلى حد يجعل التقاعس عن اتخاذ تدابير صارمة خياراً لا يمكن الاطمئنان إليه.



## سبل الدفاع

كيف نتصدى لهذه المخاطر إذاً؟ أرى ثلاثة مسارات:

**المسار الأول – الضمانات المدمجة في النماذج:** يمكن لشركات الذكاء الاصطناعي تضمين نماذجها قيوداً تمنعها من المساهمة في إنتاج الأسلحة البيولوجية. وتعمل Anthropic بجدية على هذا المسار. فدستور كلود الذي يركز في معظمه على مبادئ وقيم عليا يتضمن عدداً محدوداً من المحظورات المطلقة، وواحد منها يتعلق تحديداً بالمساعدة في إنتاج أسلحة الدمار الشامل بيولوجيةً كانت أم كيميائية أم نووية أم إشعاعية.

ولأن جميع النماذج قابلة لاختراق قيودها، فقد طوّرونا منذ منتصف 2025 – حين أظهرت اختباراتنا اقتراب نماذجنا من العتبة الحرجة – مصنفاً يرصد مخربات ذات صلة بالأسلحة البيولوجية ويحجبها تحديداً. وقد وجدنا هذه المصنّفات صارمة على نحو ملحوظ حتى في مواجهة محاولات التحايل المتطورة. ورغم أنها ترفع تكاليف التشغيل بشكل قياسي – تبلغ في بعض النماذج نحو 5% من تكاليف الاستنتاج الإجمالية وتقتطع من هوامش الربح – نرى أن هذا الثمن مستحق.

غير أن هذه الضمانات لا تكفي وحدها، إذ لا شيء يلزم الشركات بالإبقاء عليها. وأخشى أن يتشكّل بمرور الوقت منطق "معضلة السجين" حين تجد بعض الشركات مصلحة في التخلي عن هذه المصنّفات خفصاً للتكاليف. وهذا مثال نموذجي لمشكلة الآثار الخارجية السلبية التي لا يمكن حلّها بالإجراءات الطوعية لأنثروبك وحدها أو لأي شركة أخرى منفردة. ومعايير الصناعة الطوعية قد تُسهم في الحد منه، وكذلك التقييم والتحقق المستقل الذي تضطلع به معاهد أمان الذكاء الاصطناعي والجهات المحايدة.

**المسار الثاني – التدخل الحكومي:** موقفي هنا مماثل لموقفي في مواجهة مخاطر الاستقلالية: نبدأ بمتطلبات الشفافية التي تُمكن المجتمع من القياس والمراقبة والتصدي الجماعي للمخاطر دون تدخل ثقيل يُعيق الحركة الاقتصادية. وحين تتبلور مستويات أوضح من المخاطر، يمكن صياغة تشريعات تستهدف هذه المخاطر بدقة بالغة مع أدنى قدر من الأضرار الجانبية. وفي ما يخص الأسلحة البيولوجية تحديداً، أرى أن وقت هذا التشريع المُركّز ربما بات يقترب. وقد يستلزم الدفاع الشامل عن هذه المخاطر تعاوناً دولياً، حتى مع المنافسين الجيوسياسيين، وإن كنت متشككاً في معظم ضروب التعاون الدولي في شأن الذكاء الاصطناعي؛ غير أن هذا قد يكون المجال الضيق الاستثنائي الذي يوجد فيه بعض أمل لكبح جماعي. فحتى الأنظمة الاستبدادية لا ترغب في هجمات إرهابية بيولوجية على نطاق واسع.

**المسار الثالث – تطوير وسائل الدفاع البيولوجية:** تشمل هذه الوسائل أنظمة الرصد المبكر، والاستثمار في أبحاث تنقية الهواء كالتعقيم بالأشعة فوق البنفسجية البعيدة، وتطوير لقاحات سريعة الاستجابة قادرة على التكيف مع أي هجوم، وتحسين معدات الحماية الشخصية، وتطوير علاجات ولقاحات ضد أشد العوامل البيولوجية المحتملة خطورةً. ولقاحات الحمض النووي الريبوزي المرسل التي تتيح تصميم استجابة لأي فيروس أو متحوّر نموذج مبكر لما يمكن تحقيقه في هذا الاتجاه.

بيد أنني أرى أن آمالنا على صعيد الدفاع ينبغي أن تبقى في حدود واقعية. فثمة عدم تماثل أصيل بين الهجوم والدفاع في المجال البيولوجي: العوامل المُعدية تنتشر بسرعة تلقائية، في حين يستلزم الدفاع رصداً وتطعيماً وعلاجاً منسقاً عبر أعداد ضخمة من الناس في وقت قصير جداً. وما لم تكن الاستجابة خاطفة البرق – وهذا نادر – يكون معظم الضرر قد وقع قبل أن يُستكمل التعبئة. لذا



ستبقى الضمانات الوقائية خطنا الدفاعي الرئيسي حتى تُفضي التطورات التقنية المستقبلية إلى إعادة رسم هذه المعادلة لصالح جانب الدفاع.

### الهجمات الإلكترونية: التهديد المائل الآن

تستحق الهجمات الإلكترونية إشارة سريعة هنا، إذ تختلف عن الهجمات البيولوجية في كونها وقعت بالفعل – بما فيها هجمات واسعة النطاق للتجسس الحكومي. ونتوقع أن تتصاعد قدراتها مع تقدم النماذج حتى تصبح الطريقة الرئيسية لتنفيذ الهجمات الإلكترونية. أتوقع أن تغدو الهجمات الإلكترونية القائمة على الذكاء الاصطناعي تهديداً خطيراً وغير مسبوق لسلامة منظومات الحوسبة حول العالم، وتعمل Anthropic بجهد بالغ للتصدي لها. والسبب الذي يجعلني أولي الأسلحة البيولوجية أهمية أكبر من الهجمات الإلكترونية هو أن هذه الأخيرة أقل احتمالاً لإزهاق أرواح بأعداد مقارنة للبيولوجيا، وأن توازن الهجوم والدفاع في الفضاء الإلكتروني يُبدي قدراً من إمكانية التدارك إذا استثمرنا فيه بجدية.

وعلى الرغم من أن البيولوجيا تمثل اليوم أشد مصادر الخطر وطأةً، فإن مصادر أخرى كثيرة قائمة، وقد يبرز من بينها ما هو أشد خطورة. والمبدأ العام الثابت هو أن الذكاء الاصطناعي، دون تدابير مضادة، سيواصل خفض الحاجز أمام النشاط التدميري على نطاق أوسع وأوسع، ويستحق هذا التهديد استجابة بشرية جادة وعلى مستوى التحدي.

### 3- ما وراء الدمار الفردي

تناول القسم السابق خطر قيام أفراد ومنظمات صغيرة باستغلال شريحة محدودة من "دولة العباقرة داخل مركز البيانات" لإحداث دمار واسع النطاق. لكن ثمة ما يستدعي قلقاً أعمق وأشد إلحاحاً: توظيف الذكاء الاصطناعي لتركيز السلطة أو اغتصابها، على الأرجح من قِبَل جهات أكبر حجماً وأرسخ قدماً.

في كتابي "آلات المحبة والرحمة"، ناقشتُ إمكانية لجوء الحكومات الاستبدادية إلى الذكاء الاصطناعي القوي لمراقبة مواطنيها وقمعهم بأساليب يصعب الإصلاح منها أو الإطاحة بها. فالأنظمة الشمولية اليوم مقيدة في درجة بطشها بسبب حاجتها إلى بشر لتنفيذ أوامرها، وللبشر حدود في تقبل اللإنسانية. أما الأنظمة الشمولية المدعومة بالذكاء الاصطناعي فلن يكون لها مثل هذه القيود.

والأخطر من ذلك أن الدول قد تستثمر تفوقها في الذكاء الاصطناعي للسيطرة على دول أخرى. فلو كانت "دولة العباقرة" برمتها مملوكةً ومسيّرةً من قِبَل الجهاز العسكري لدولة بعينها، في حين تفنقر سائر الدول إلى قدرات مماثلة، فما الذي يمكنها من الدفاع عن نفسها؟ ستكون مغلوبة على أمرها في كل مواجهة، كحرب تدور بين البشر والفئران. ويقودنا التوليف بين هذين الهاجسين إلى احتمال بالغ الخطورة: الوصول إلى ديكتاتورية شمولية كوكبية. ومن البديهي أن درء هذا المصير ينبغي أن يكون في طبيعة أولوياتنا القصوى.

### أدوات الاستبداد الممكن بالذكاء الاصطناعي

ثمة طرق شتى قد يُمكن بها الذكاء الاصطناعي الاستبداد أو يُرسخه أو يوسع نطاقه، أكتفي بأبرز ما يقلقني منها. وتجدر الإشارة إلى أن بعض هذه التطبيقات له استخدامات دفاعية مشروعة، ولا



أدعو بالضرورة إلى حظرها بإطلاق، لكنني قلق من أنها تميل هيكلية إلى خدمة الأنظمة الاستبدادية:

**أولاً – الأسلحة المستقلة الكاملة: سرب يضم ملايين أو مليارات الطائرات المسلحة ذاتية القيادة،** تتحكم بها محلياً أنظمة ذكاء اصطناعي قوية وتُنسّقها استراتيجياً على الصعيد العالمي أنظمة أكثر قوة – هذا السرب سيُشكّل جيشاً لا يُقهر، قادراً على هزيمة أي قوة عسكرية في العالم وإخماد أي معارضة داخلية بملاحقة كل مواطن. وقد أيقظتنا الحرب الروسية الأوكرانية على حقيقة أن حرب الطائرات المسيّرة باتت واقعاً ماثلاً، وإن لم تبلغ بعد الاستقلالية الكاملة وظلت جزءاً يسيراً مما قد يُتيح الذكاء الاصطناعي القوي. **فأبحاث الذكاء الاصطناعي قادرة على أن تجعل طائرات دولة ما متفوقة تفوقاً ساحقاً على مثيلاتها، وتُسرع تصنيعها وتُحصنها ضد الحرب الإلكترونية وتُحسّن مناوراتها.** ولا شك أن لهذه الأسلحة استخدامات مشروعة في الدفاع عن الديمقراطية – فقد كانت عاملاً حاسماً في صمود أوكرانيا، وستكون كذلك في الدفاع عن تايوان. لكنها في الآن ذاته سلاح بالغ الخطورة: ينبغي أن نُساور القلق من وقوعها في يد الأنظمة الاستبدادية، وأن نخشى كذلك – نظراً لهيمنتها وضعف المساءلة المتعلقة بها – تنامي احتمال استخدام الحكومات الديمقراطية لها ضد شعوبها.

**ثانياً – المراقبة الشاملة بالذكاء الاصطناعي:** ذكاء اصطناعي بالغ القدرة يمكنه على الأرجح اختراق أي منظومة حاسوبية في العالم، واستثمار ما يحصل عليه من صلاحيات لقراءة جميع الاتصالات الإلكترونية وتحليلها، بل حتى الاتصالات الشخصية المباشرة إن أمكن نشر أجهزة التسجيل أو الاستيلاء عليها. قد يصبح من الممكن بصورة مرعبة توليد قائمة شاملة بكل من يُعارض الحكومة في أي مسألة، حتى حين لا يُصرّح بهذه المعارضة. فذكاء اصطناعي يرصد مليارات المحادثات لملايين الأشخاص قادر على قياس نبض الرأي العام، وكشف بؤر الولاء المتزعزع قبل أن تتبلور وقمعها في مهدها. وقد يُفضي ذلك إلى فرض "البانوبتيكون" الحقيقي – المراقبة الكاملة الدائمة – على نطاق لم يبلغه حتى الحزب الشيوعي الصيني.

**ثالثاً – الدعاية المُشخصنة بالذكاء الاصطناعي:** تُبيّن ظاهرتنا "الذهان الناجم عن الذكاء الاصطناعي" و"صديقات الذكاء الاصطناعي" الراهنتان أن هذه النماذج، حتى عند مستواها الحالي، قادرة على التأثير النفسي العميق في البشر. فما بالك بنسخ أكثر قوة بكثير، متجذرة في حياة الناس اليومية ومُطلعة عليها، قادرة على نمذجتهم والتأثير فيهم على مدار أشهر وسنوات؟ **هذه النسخ ستكون قادرة في الأرجح على غسل دماغ الكثيرين – إن لم يكن معظمهم – وبثّ أي أيديولوجية أو موقف مرغوب،** ويمكن توظيفها من قبَل حاكم لا يعرف الضمير لضمان الولاء وإخماد المعارضة، حتى في ظل قمع يتمرد عليه الناس عادةً. ونقارن ما يُثيره تطبيق تيك توك اليوم من مخاوف جراء دعايته الموجهة للأطفال – وهي مخاوف مشروعة – بما يمثله وكيل ذكاء اصطناعي شخصي يتعمّق في معرفتك عاماً بعد عام ويستثمر هذه المعرفة لصياغة قناعاتك كلها: البون الشاسع بين التأثيرين يُدرك بلا عناء.

**رابعاً – صنع القرار الاستراتيجي:** يمكن توظيف دولة العباقره داخل مراكز البيانات لإسداء المشورة لدولة أو جماعة أو فرد في الاستراتيجية الجيوسياسية – ما قد نسمّيه "بسمارك افتراضي". يمكنها تحسين الاستراتيجيات الثلاث المذكورة آنفاً لاغتصاب السلطة، بل ابتكار أساليب لم تخطر على بالنا لكن دولة العباقره ستصل إليها. فالديبلوماسية والاستراتيجية العسكرية والبحث والتطوير والاستراتيجية الاقتصادية وغيرها من الميادين ستشهد جميعها ارتفاعاً حاداً في الفعالية بفعل



الذكاء الاصطناعي القوي. وكثير من هذه المهارات ستكون مفيدة مشروعياً للديمقراطيات في الدفاع عن نفسها، لكن إمكانية إساءة استخدامها تبقى قائمة في أي يد.

## ما الذي يُقلقني أكثر

بعد وصف مصادر القلق، لنتناول الفاعلين المثيرين للهواجس، مرتبين تنازلياً وفق مستوى الخطر:

**الحزب الشيوعي الصيني:** لا تنافس الصين سوى الولايات المتحدة في قدراتها على الذكاء الاصطناعي، وهي الأوفر حظاً من بين جميع الدول لتجاوز الولايات المتحدة في هذه القدرات. وحكومتها استبدادية وتدير دولة مراقبة تقنية متطورة، وقد وُظفت المراقبة القائمة على الذكاء الاصطناعي فعلياً – بما فيها في قمع الأويغور – ويرجح أنها تمارس الدعاية الخوارزمية عبر تيك توك إلى جانب مساعيها الدعائية الدولية الأخرى. لديها الطريق الأوضح بلا منازع نحو الكابوس الشمولي الممكن بالذكاء الاصطناعي الذي رسمته آنفاً. وقد يكون هذا المسار هو المآل الافتراضي داخل الصين وفي الدول الاستبدادية الأخرى التي تُصدّر لها الصين تقنيات المراقبة. وأؤكد أنني لا أخص الصين بهذا التحليل عن عدا أو تحيز، بل لأنها تجمع بامتياز ثلاثة عوامل: البراعة في الذكاء الاصطناعي، والحكومة الاستبدادية، ودولة المراقبة التقنية. والصينيون أنفسهم هم الأوفر حظاً في المعاناة من هذا القمع، وليس لهم صوت في قرارات حكومتهم. وأكّن للشعب الصيني إعجاباً عميقاً واحتراماً صادقاً، وأساند المعارضين الشجعان الذين يناضلون داخل الصين من أجل الحرية.

**الديمقراطيات المتقدمة في الذكاء الاصطناعي:** كما أشرت، للديمقراطيات مصلحة مشروعية في بعض الأدوات العسكرية والجيوسياسية المدعومة بالذكاء الاصطناعي، إذ تُمثل هذه الحكومات أفضل فرصة لمواجهة توظيف الأنظمة الاستبدادية لهذه الأدوات. وأنا مؤيد على نطاق واسع لتزويد الديمقراطيات بالأدوات اللازمة للتفوق في عصر الذكاء الاصطناعي. لكننا لا نستطيع تجاهل إمكانية إساءة استخدام هذه التقنيات من قِبَل الحكومات الديمقراطية ذاتها. فالضمانات المعتادة التي تحول دون توجيه الجهاز العسكري والاستخباراتي نحو الداخل قائمة، لكن الذكاء الاصطناعي يحتاج عدداً ضئيلاً جداً من المشغلين، مما يُتيح التحايل على هذه الضمانات. فضلاً عن أن بعض هذه الضمانات تتآكل تدريجياً في بعض الديمقراطيات. وخلاصة القول: نعم لتسليح الديمقراطيات بالذكاء الاصطناعي، لكن بحذر وضمن حدود واضحة؛ فهي منظومة المناعة التي نحتاجها لمواجهة الاستبداد، غير أن **مناعة الجسم نفسه قد تنقلب أحياناً عليه.**

**الدول غير الديمقراطية ذات مراكز البيانات الكبيرة:** خارج الصين، معظم الدول الأقل ديمقراطيةً ليسوا لاعبين محوريين في مجال الذكاء الاصطناعي بمعنى امتلاكها شركات تُنتج نماذج متقدمة، مما يجعل خطرها مختلفاً جوهرياً ودون ما تمثله الصين. لكن بعض هذه الدول تمتلك مراكز بيانات ضخمة قد تُتيح تشغيل ذكاء اصطناعي متقدم على نطاق واسع. والخطر النظري هنا أن تُصادر هذه الحكومات مراكز البيانات وتوظفها لمآربها الخاصة. وهو خطر أقل إثارةً للقلق لديّ مقارنةً بدول كالصين التي تُطوّر الذكاء الاصطناعي مباشرةً، لكنه يستحق الاستحضار.

**شركات الذكاء الاصطناعي:** قد يبدو هذا القول مجرداً من رئيس تنفيذي لإحدى هذه الشركات، لكنني أرى أن الطبقة التالية من المخاطر تنبع من شركات الذكاء الاصطناعي ذاتها. فهي تتحكم في مراكز بيانات ضخمة، وتُدرّب نماذج متقدمة، وتمتلك أعماق الخبرات في توظيف هذه النماذج، وتتواصل يومياً مع مئات الملايين من المستخدمين مع ما يترتب على ذلك من إمكانية تأثير. وما



يُعوزها أساساً هو شرعية الدولة وبنيتها التحتية، إذ إن كثيراً مما يلزم لبناء أدوات استبدال الذكاء الاصطناعي سيكون فعلاً مجرماً لو أقدمت عليه شركة أو على الأقل مثيراً للريبة. لكن بعضه ليس مستحيلاً: كأن تُستخدم منتجاتها لغسل دماغ قاعدتها الجماهيرية الضخمة. وحوكمة شركات الذكاء الاصطناعي تستوجب تدقيقاً جاداً وصارماً.

## الردود على الاعتراضات

ثمة حجج تُساق للتشكيك في جدية هذه المخاطر، وليتني أستطيع الإيمان بها، لأن الاستبدال الممكن بالذكاء الاصطناعي يربيني. يستحق بعضها توقفاً ورداً:

**الاعتراض الأول – الرادع النووي:** يُعوّل بعضهم على الردع النووي لكبح استخدام الأسلحة الذكية في الغزو. لكن ثقتي بقدرة هذا الردع أمام دولة عابرة داخل مراكز البيانات ليست مطلقة. فالذكاء الاصطناعي القوي قد يبتكر سبلاً للكشف عن الغواصات النووية وضربها، أو يُنفذ عمليات تأثير على المشغّلين البشريين لمنظومة الأسلحة النووية، أو يشن هجمات إلكترونية على الأقمار الاصطناعية المعنية برصد الإطلاق النووي. ومن ثم فإن الاستناد المطلق إلى الردع النووي في هذا الواقع المتحوّل مراهنه عالية المخاطر.

**الاعتراض الثاني – الإجراءات المضادة:** يُقال إنه يمكن مواجهة طائرات الخصم بطائرات مماثلة، والدفاع الإلكتروني سيرتقي مع الهجوم الإلكتروني، وثمة طرق للتحصين ضد الدعاية. ردّي أن هذه التدابير الدفاعية لن تُجدي إلا إذا امتلك الطرفان ذكاءً اصطناعياً متكافئاً في القوة. وما يُقلقني هو الطبيعة الذاتية التعزيز للذكاء الاصطناعي: كل جيل يُصمّم من قبّل الجيل السابق، مما قد يُفضي إلى تفوق متسارع يصعب تعويضه. **يجب أن يظل الرائد في هذه المنظومة دولة ديمقراطية لا استبدادية.** وحتى لو تحقّق توازن القوى، يبقى احتمال تقسيم العالم إلى مناطق نفوذ شمولية – كما في رواية "1984" – خطراً قائماً ومقلّفاً.

## خطوط الدفاع

كيف ندافع إذاً عن أنفسنا في مواجهة هذا الطيف الواسع من الأدوات الاستبدادية والجهات المحتملة؟

**أولاً – وقف بيع تكنولوجيا الرقائق للصين:** يجب بشكل قاطع وقف بيع الرقائق وأدوات تصنيعها ومراكز البيانات للحزب الشيوعي الصيني. فالرقائق تُمثّل العنق الزجاجة الأودد والأشد أهمية في مسيرة الذكاء الاصطناعي القوي، ووجبها إجراء بسيط بالغ الفعالية. إن القول بأن "نشر حزمنا التقنية في العالم يُمكن أمريكا من الانتصار" في معركة اقتصادية مبهمة يُشبه تماماً بيع أسلحة نووية لكوريا الشمالية والتباهي بأن أغلفة الصواريخ من صنع بوينغ. الصين متأخرة عدة سنوات في إنتاج الرقائق المتقدمة، وتلك السنوات هي الحقبة الحرجة بعينها لبناء "الدولة الرقمية العبقريّة". ومنح صناعتها دفعة هائلة في هذه اللحظة بالذات لا مسوّغ له.

**ثانياً – تمكين الديمقراطيات من مواجهة الاستبدال:** يكتسب هذا المنطق أهمية خاصة حين يتعلق الأمر بالدفاع عن ديمقراطيات تحت الهجوم كأوكرانيا وتايوان، وتمكين الديمقراطيات من توظيف أجهزة استخباراتها لإضعاف الأنظمة الاستبدادية من الداخل. وتحالف الولايات المتحدة وحلفائها الديمقراطيين، إذا حقّق التفوق في الذكاء الاصطناعي القوي، سيكون في موقع يُتيح له ليس فقط الدفاع عن نفسه، بل احتواء الأنظمة الاستبدادية والحدّ من انتهاكاتها.



**ثالثاً – حدود حمراء داخل الديمقراطيات :** نحتاج إلى خطوط حمراء صريحة تمنع حكوماتنا الديمقراطية ذاتها من إساءة توظيف الذكاء الاصطناعي. الصياغة التي توصلت إليها هي: **توظيف الذكاء الاصطناعي في الدفاع الوطني بكل السبل، باستثناء ما يجعلنا أشبه بخصومنا الاستبداديين.**

فمن القائمة السابقة، يبدو توظيف الذكاء الاصطناعي للمراقبة الشاملة والدعاية الجماهيرية على الصعيد الداخلي خطين أحمرين مطلقين لا تهاون فيهما. أما الأسلحة المستقلة والتوجيه الاستراتيجي، فخطوطهما أقل وضوحاً لما لهما من استخدامات مشروعة في الدفاع عن الديمقراطية. وما يُقلقني هنا بالدرجة الأولى تركيز "الأصابع على الزناد" في يد عدد شحيح جداً من الأشخاص، بحيث يتمكن واحد أو دفنة قليلة من تشغيل جيش من الطائرات المسيّرة دون الحاجة إلى تعاون أي إنسان آخر. وقد نحتاج مستقبلاً إلى آليات رقابة مباشرة وأنية أكثر صرامة، ربما تشمل فروع الحكومة خارج السلطة التنفيذية.

**رابعاً – تأسيس محرّمات دولية ضد أسوأ إساءات الذكاء الاصطناعي :** أدرك أن الرياح السياسية الراهنة تهبّ في غير صالح التعاون والأعراف الدولية، لكن الحاجة إليهما هنا ملحة. يجب أن يدرك العالم الإمكانيات المظلمة للذكاء الاصطناعي في قبضة الطغاة، وأن يتعامل مع بعض استخداماته بوصفها جرائم ضد الإنسانية: المراقبة الجماهيرية الشاملة، والدعاية الممنهجة، وأنماط محددة من توظيف الأسلحة المستقلة في العدوان. والحجة الأقوى في هذا السياق أنه كما أضحى الإقطاع غير قابل للاستمرار مع الثورة الصناعية، قد يُفضي عصر الذكاء الاصطناعي حتماً وبمنطق لا يُردّ إلى أن الديمقراطية هي الشكل الوحيد القابل للحياة للحكم إذا أريد للبشرية مستقبل كريم.

**خامساً وأخيراً – الرقابة الصارمة على شركات الذكاء الاصطناعي :** يجب مراقبة هذه الشركات عن كثب، وكذلك علاقتها بالحكومة، التي هي علاقة ضرورية لكنها لا بد أن تخضع لحدود وقيود واضحة. فالقدرات الهائلة التي يجسدها الذكاء الاصطناعي القوي تجعل أطر الحوكمة التقليدية – المصمّمة لحماية المساهمين ومنع الاحتيال – قاصرةً عن استيعاب المهمة. وقد تكون هناك قيمة في التزام علني من الشركات بعدم اتخاذ إجراءات بعينها، كبناء عتاد عسكري سراً، أو تخصيص موارد حاسوبية ضخمة لأفراد بعينهم دون مساءلة، أو استخدام منتجاتها الذكية في التلاعب بالرأي العام لصالحها.

## خاتمة

يأتي الخطر من اتجاهات متعددة، وبعضها يتعارض مع بعض. والثابت الوحيد في كل هذا أننا ملزمون بالسعي إلى المساءلة والأعراف الضابطة والضمانات الصارمة التي تُحاط بها جميع الأطراف – حتى ونحن نمح الفاعلين "الجيدين" ما يحتاجونه من أدوات لكبح الفاعلين "السيئين".

## 4- الاضطراب الاقتصادي: حين يُصبح النمو نعمة ونقمة في آن واحد

تناولت الأقسام الثلاثة السابقة المخاطر الأمنية للذكاء الاصطناعي القوي:

- المخاطر النابعة من الذكاء الاصطناعي ذاته
- والمخاطر المترتبة على إساءة استخدامه من قِبَل أفراد ومنظمات صغيرة
- ومخاطر توظيفه من قِبَل الدول والمنظمات الكبرى.



فإذا وضعنا هذه المخاطر الأمنية جانباً أو افترضنا أنها قد عُولجت، انبثق السؤال الاقتصادي: **ما الذي سيترتب على ضخ هذا الكم الهائل من "الرأس مال البشري" الرقمي في الاقتصاد؟**

الأثر الجلي والأكثر بدهاءةً هو تسريع النمو الاقتصادي تسريعاً استثنائياً. فوتيرة التقدم في البحث العلمي والابتكار الطبي الحيوي والتصنيع وسلاسل التوريد وكفاءة المنظومة المالية وغيرها شبه مضمونة لتُفضي إلى معدلات نمو اقتصادي تفوق ما شهده التاريخ. وقد أشرت في "آلات المحبة والرحمة" إلى إمكانية بلوغ معدل نمو سنوي مستدام في الناتج المحلي الإجمالي يتراوح بين 10 و20%.

لكن هذه المعادلة ذات حدّين؛ إذ يطرح ذلك تساؤلاً جوهرياً: ما مستقبل غالبية البشر اقتصادياً في عالم كهذا؟ دأبت التقنيات الجديدة على إحداث صدمات في سوق العمل، وفي كل مرة نجحت البشرية في التعافي. لكن ما يقلقني أن هذه التعافيات كانت ممكنة لأن الثورات السابقة طالت شرائح محدودة من الطيف الواسع للقدرات البشرية، مما أبقى متسعاً للبشر كي ينتقلوا إلى مهام جديدة. أما الذكاء الاصطناعي فتأثيره أوسع بكثير وأسرع بكثير، مما يجعل إدارة هذه المرحلة بنجاح تحدياً استثنائياً.

## اضطراب سوق العمل

ثمة مشكلتان تحديداً تشغلان تفكيري: الاستغناء عن القوى العاملة، وتركز القوة الاقتصادية. لنبدأ بالأولى.

هذا موضوع حدّرت منه بكل صراحةً عام 2025، حين توقعت أن يُزيح الذكاء الاصطناعي نصف وظائف الياقات البيضاء المبتدئة في غضون 1 إلى 5 سنوات، حتى وهو يُسرّع النمو الاقتصادي والتقدم العلمي. وقد أطلق هذا التحذير نقاشاً عاماً واسعاً؛ اتفق معي كثير من الرؤساء التنفيذيين والتقنيين والاقتصاديين، بينما رأى آخرون أنني أقع في مغالطة "كتلة العمل الثابتة" وأسيء فهم آليات سوق العمل. لذا يستحق هذا الأمر تفصيلاً دقيقاً يُزيل الالتباس.

## كيف تتكيف أسواق العمل عادةً مع التكنولوجيا

تبدأ التقنية الجديدة عادةً بتحسين كفاءة أجزاء من وظيفة بشرية قائمة. فمع مطلع الثورة الصناعية مثلاً، رفعت المحاريث المطوّرة كفاءة المزارعين في بعض جوانب عملهم، فارتفعت إنتاجيتهم وتحسّنت أجورهم.

في المرحلة التالية، باتت الآلات قادرة على إنجاز أجزاء كاملة من عمل الزراعة – كآلة الدراس وحفارة البذور. فتراجعت حصة الإنسان في العمل، لكن ما يؤديه أصبح ذا رافعة أعظم لتكامله مع عمل الآلة، فاستمرت إنتاجيته في الارتفاع. ووفق مفارقة جيفونز، واصلت أجور المزارعين صعودها بل وربما ازداد عددهم. حتى حين تُنجز الآلات 90% من العمل، يتسع المجال أمام الإنسان ليُضعف ما يؤديه بعشرة أضعاف، فيُنتج أضعافاً مضاعفة من العائد بالجهد ذاته.

في نهاية المطاف، أصبحت الآلات تُنجز كل شيء أو كاد – كما هو الحال اليوم مع الحصادات الآلية والجرارات. عندئذٍ تراجعَت الزراعة بحدة كمجال للعمالة البشرية. لكن لأن الزراعة لم تكن إلا واحدة من فعاليات إنسانية مفيدة كثيرة، وجد الناس طريقهم إلى مهن أخرى كتشغيل المصانع. قبل مئتين وخمسين عاماً كان 90% من الأمريكيين يعيشون على الزراعة، وفي أوروبا كانت تستوعب



50 إلى 60% من العمالة. أما اليوم فلا تتجاوز هذه النسبة أرقاماً أحادية، بعد أن انتقل العمال إلى الصناعة فاقتصاد المعرفة. لا توجد "كتلة عمل" ثابتة، بل طاقة متنامية باستمرار على إنجاز المزيد بجهد أقل.

## لماذا سيكون الذكاء الاصطناعي مختلفاً؟

أرى أن الأمر لن يسير بالطريقة المعتادة هذه المرة، وإليك الأسباب:

- **السرعة:** وتيرة التقدم في الذكاء الاصطناعي أسرع بمراحل من أي ثورة تقنية سابقة. ففي السنتين الأخيرتين، انتقلت نماذج الذكاء الاصطناعي من بالكاد إتمام سطر برمجي واحد إلى كتابة معظم أو جميع الكود لدى بعض المهندسين – بمن فيهم مهندسو Anthropic أنفسهم. وقریباً قد تُنجز المهمة الكاملة لمهندس البرمجيات من أولها لآخرها. يصعب على الناس مواكبة هذه الوتيرة سواء في تكيف أسلوب العمل أو في الانتقال إلى مهن جديدة. وقد بدأ مبرمجون أسطوريون يصفون أنفسهم بأنهم "يتعثرون في اللحاق بالركب." والأرجح أن هذه السرعة ستتعاضد مع نماذج البرمجة التي تُسرّع بدورها تطوير الذكاء الاصطناعي. **السرعة في حد ذاتها لا تعني أن أسواق العمل لن تتعافى، لكنها تعني أن مرحلة الانتقال ستكون أشد إبلاماً مما ألفناه.**
- **الاتساع المعرفي:** كما يُوحى مصطلح "دولة العباقرة داخل مركز البيانات"، سيكون الذكاء الاصطناعي قادراً على استيعاب طيف واسع جداً من القدرات المعرفية البشرية – ربما جميعها. وهذا يختلف اختلافاً جوهرياً عن تقنيات سابقة كالميكنة الزراعية أو المواصلات أو حتى الحاسوب. إذ يصعب الانتقال بين مهن متاخمة حين تتعرض كلها للاضطراب في آن واحد. فالمهارات المعرفية العامة لوظائف مبتدئة في التمويل والاستشارات والقانون متقاربة؛ تقنية تُزيح إحداها تُتيح الانتقال إلى البديلين الآخرين. أما تقنية تُزيح الثلاثة دفعةً واحدة مع كثير من الوظائف المشابهة، فأصعب على التكيف. والأدق أن الذكاء الاصطناعي ليس بديلاً عن وظائف بعينها، بل بديل عمالي عام عن الإنسان.
- **التقطيع وفق القدرة المعرفية:** يبدو أن الذكاء الاصطناعي يتقدم في طيف القدرات من الأدنى نحو الأعلى. ففي البرمجة مثلاً، انتقلت نماذجنا من مستوى "المبرمج المتوسط" إلى "المبرمج القوي" إلى "المبرمج الاستثنائي". وبدأنا نرى المسار ذاته في العمل المعرفي العام. وهذا يُنذر بأن الاضطراب لن يطال أصحاب مهارات بعينها فحسب – وهو ما يتيح إعادة التدريب – بل سيغال الأفراد ذوي خصائص معرفية فطرية محدودة يصعب تغييرها. وأخشى أن يتشكّل من هؤلاء "طبقة دنيا" من العاطلين أو ذوي الأجر المتدنية جداً.
- **القدرة على سدّ الثغرات:** حين تغطي تقنية جديدة على مهنة، يتعلّق البشر عادةً بالجوانب التي تعجز عنها التقنية وإن كانت ضئيلة. لكن الذكاء الاصطناعي لا يكتفي بكونه تقنية متقدمة سريعاً، بل هو تقنية تتكيف بسرعة مضاعفة. فعند كل إصدار جديد، تقيس شركات الذكاء الاصطناعي بدقة نقاط القوة ونقاط الضعف، وتُعالجها في الإصدار التالي. كثير من الثغرات التي حسبها الناس موروثاً في بنية التقنية جرى تجاوزها في غضون أشهر قليلة.

## ردود على الاعتراضات الشائعة

- **اعتراض الانتشار البطيء:** يقول البعض إن التطبيق الفعلي في الاقتصاد سيكون أبطأ بكثير من التقدم التقني. وهذا صحيح جزئياً – لذا جعلت توقعي 1 إلى 5 سنوات لا أشهراً.



لكن التأثيرات التدريجية تشتري الوقت فحسب، ولا تُلغي الاضطراب. وتبني الذكاء الاصطناعي في المؤسسات ينمو بوتيرة تتجاوز أي تقنية سابقة، كما أن الشركات الناشئة ستُسرع التبني أو تُزيح الكبرى مباشرةً.

- **اعتراض التحوّل إلى العمل الجسدي:** يرى بعضهم أن المهن الجسدية بمنأى عن هذا الاضطراب. لكن كثيراً من الأعمال الجسدية يُنجزه الآلات أو في طريقها إليه. وبإمكان الذكاء الاصطناعي القوي تسريع تطوير الروبوتات ثم السيطرة عليها في العالم المادي. وحتى لو اقتصر الاضطراب على المهام المعرفية، سيكون في حد ذاته اضطراباً غير مسبوق.
- **اعتراض اللمسة الإنسانية:** ثمة من يرى أن بعض المهام تستلزم بطبيعتها الحضور الإنساني. وأنا أقل يقيناً هنا، لكنني أشك في أن هذا الاستثناء كافٍ لاستيعاب معظم سوق العمل. فالذكاء الاصطناعي مُستخدَم على نطاق واسع في خدمة العملاء، ويُفيد كثيرون بأن مناقشة مشكلاتهم الشخصية مع الذكاء الاصطناعي أيسر من جلسات العلاج النفسي – لأنه أكثر صبراً. وقد لجأت شقيقتي خلال مشكلة طبية في حملها إلى كلود حين شعرت بإخفاق أطبائها، فوجدت لديه أسلوباً أكثر إنسانية في التعامل.
- **اعتراض الميزة النسبية:** يقول بعض الاقتصاديين إن "قانون الميزة النسبية" سيحمي البشر، إذ تُتيح الاختلافات النسبية في الكفاءة أساساً للتخصص وتبادل المنافع حتى حين يتفوق أحد الطرفين في كل شيء. لكن حين يكون الذكاء الاصطناعي أكثر إنتاجية من الإنسان بآلاف الأضعاف، يبدأ هذا المنطق في الانهيار. فحتى أصغر تكاليف المعاملات قد تجعل التعامل مع البشر غير مُجدٍ اقتصادياً، وقد تتدنى الأجور البشرية إلى مستويات بائسة حتى لو ظلّ ثمة ما يُقدّمه الإنسان.

## سبل المعالجة

ما الذي يمكن فعله؟ لديّ عدة مقترحات، بعضها تعمل Anthropic على تنفيذه:

**أولاً – رصد دقيق وآني لما يجري:** حين يتسارع التغيير الاقتصادي، يصعب الحصول على بيانات موثوقة، وبدون بيانات موثوقة يتعدّر تصميم سياسات فعّالة. تفتقر الحكومات حالياً إلى بيانات دقيقة وعالية التكرار عن تبني الذكاء الاصطناعي عبر القطاعات. ولهذا تُشغّل Anthropic منذ عام مؤشراً اقتصادياً تنشره علناً يرصد استخدام نماذجها شبه آتياً، مُفصّلاً حسب القطاع والمهمة والموقع الجغرافي وطبيعة التوظيف.

**ثانياً – التأثير في خيارات المؤسسات:** أمام الشركات خيار بين مسارين: "توفير التكاليف" بأداء العمل ذاته بأقل عدد من الموظفين، أو "الابتكار" بإنجاز أكثر بالعدد ذاته. السوق ستُفرض كليهما حتماً، لكن ثمة فسحة لتوجيه المؤسسات نحو مسار الابتكار كلما أمكن، مما قد يشتري وقتاً ثميناً.

**ثالثاً – العناية بالموظفين:** في المدى القريب، قد يُسهم الإبداع في إعادة التوزيع الداخلي للموظفين في تفادي موجات التسريح. وفي عالم يتضاعف فيه الثروة إلى حد كبير، قد يصبح مجدياً اقتصادياً الإبقاء على الموظفين البشريين حتى بعد أن تتراجع قيمتهم الاقتصادية التقليدية. وتدرس Anthropic حالياً جملة من المسارات الممكنة لموظفيها ستُعلن عنها قريباً.

**رابعاً – مسؤولية الأثرياء:** يُحزنني أن يتبنى كثير من الأثرياء – لا سيما في عالم التقنية – موقفاً عدمياً يرى في العمل الخيري ضرباً من الاحتيال أو العبث. كلا، فالعمل الخيري الخاص كمؤسسة غيتس والبرامج الحكومية كـ PEPFAR أنقذت عشرات الملايين من الأرواح وفتحت آفاق الفرص لأجيال



كاملة. جميع المؤسسين المشاركين في Anthropic تعهدوا بالتبرع بـ 80% من ثروتهم، وتعهد موظفوها فردياً بالتبرع بحصص في الشركة تبلغ مليارات بالأسعار الراهنة – وتلتزم الشركة بمضاهة هذه التبرعات.

**خامساً – التدخل الحكومي:** مهما أثمرت الجهود الخاصة السابقة، فإن مشكلة بهذا الحجم الاقتصادي الكلي تستلزم تدخلاً حكومياً. والاستجابة السياسية الطبيعية لثروة ضخمة مقترنة بتفاوت حاد في التوزيع هي الضرائب التصاعدية، عامةً أو موجهة نحو شركات الذكاء الاصطناعي تحديداً. أرفض السياسات الضريبية السيئة التصميم، لكنني أؤيد ضريبة عادلة مُحكمة على أسس أخلاقية. وأستطيع أن أقدم لأصحاب المليارات حجة براغماتية في صميم مصلحتهم: إن لم تدعموا صيغة جيدة منها، ستحصلون حتماً على صيغة سيئة تفرضها عليكم الشعوب.

في نهاية المطاف، أرى جميع ما سبق أدواتٍ لشراء الوقت وتأجيل ما هو قادم. الذكاء الاصطناعي قادر على فعل كل شيء، وعلينا مواجهة هذه الحقيقة. **وأملّي أن نكون بطول تلك اللحظة قد استثمرنا الذكاء الاصطناعي نفسه في إعادة هيكلة الأسواق بما يخدم الجميع، وأن تكون التدخلات السابقة قد أوصلتنا بسلام إلى الضفة الأخرى من هذه المرحلة الانتقالية.**

## تركز القوة الاقتصادية

مستقلاً عن مشكلة إزاحة القوى العاملة أو التفاوت الاقتصادي في حد ذاتيهما، ثمة مشكلة تركّز القوة الاقتصادية. **ناقش القسم الأول خطر أن يُجرّد البشر من نفوذهم على يد الذكاء الاصطناعي، وناقش القسم الثالث خطر أن يُجرّد المواطنين من نفوذهم على يد حكوماتهم بالقوة والإكراه.** لكن نوعاً ثالثاً من التجريد من النفوذ قد يحدث حين يتركز الثروة تركّزاً بالغاً يمنح فئة صغيرة السيطرة الفعلية على السياسة العامة، في حين يفقد المواطن العادي أي ميزة تفاوضية. فالديمقراطية في جوهرها مرتكزة على فكرة أن الشعب ضروري لعمل الاقتصاد؛ إذا زالت هذه الرفاعة الاقتصادية، قد ينهار العقد الاجتماعي الضمني الذي تقوم عليه الديمقراطية.

لأوضح موقفي: لسْتُ ضد الإثراء. ثمة حجة وجيهة لصالحه بوصفه حافزاً للنمو الاقتصادي في الظروف الاعتيادية. لكن في مشهد يبلغ فيه نمو الناتج المحلي 10 إلى 20% سنوياً ويحكم فيه الذكاء الاصطناعي قبضته على الاقتصاد، وتمتلك فيه أفراد بعينهم حصصاً تعادل نسباً معتبرة من الناتج الكلي – فالتهديد الحقيقي ليس تثبيط الابتكار، بل **تركّز ثروة يكفي وحده لتفتيت النسيج الاجتماعي.**

وأشهر نموذج تاريخي لهذا التركّز في أمريكا هو "العصر المذهب"، وكان جون روكفلر أثرى صناعيه بثروة بلغت نحو 2% من الناتج المحلي الأمريكي آنذاك. حصة مماثلة اليوم تعادل 600 مليار دولار، وأثرى أثرياء العالم اليوم – إيلون ماسك – يتجاوز هذا الرقم بنحو 700 مليار دولار. أي أننا بلغنا مستويات تركّز للثروة لم يشهدها التاريخ، وهذا قبل أن تبلغ الآثار الاقتصادية للذكاء الاصطناعي أوجها. وليس من قبيل المبالغة، في سيناريو "دولة العباقرة"، تحيّل شركات الذكاء الاصطناعي وأشباه الموصلات والتطبيقات المستقاة منها تُدرّ إيرادات سنوية تبلغ 3 تريليونات دولار، وتُقيّم بنحو 30 تريليوناً، مُنتجةً ثروات شخصية تقاس بالتريليونات. في ذلك العالم، لن تنطبق على واقعنا نقاشات السياسة الضريبية الراهنة، لأننا سنكون في حفيقة مختلفة جوهرياً.



وما يزيد قلقي هو الترابط الآخذ في التوطد بين هذا التركز الاقتصادي والمنظومة السياسية. فمراكز بيانات الذكاء الاصطناعي باتت تُمثل حصةً وازنة من النمو الاقتصادي الأمريكي، مما يُوثق التشابك بين مصالح الشركات التقنية الكبرى والمصالح السياسية للحكومة بصورة قد تُفضي إلى حوافز مشوّهة. ونرى ذلك جلياً في إدجام الشركات التقنية عن انتقاد الحكومة، وفي دعم الحكومة لسياسات تُجرّد تنظيم الذكاء الاصطناعي من أي قيد.

## سبل المعالجة

**أولاً – الاستقلالية في اتخاذ المواقف:** يمكن للشركات أن تختار ببساطة ألا تكون جزءاً من هذه الدوامة. فـ Anthropic دأبت على أن تكون فاعلاً في السياسة العامة لا في السياسة الحزبية، وعلى الإفصاح عن مواقفها الحقيقية أياً كانت الإدارة الحاكمة. فدعمنا تنظيماً رشيداً للذكاء الاصطناعي وضوابط تصدير تصبّ في المصلحة العامة، حتى حين تتعارض مع السياسة الحكومية. أخبرني كثيرون بأن هذا النهج سيجلب علينا معاملة غير مُنصفة، غير أن تقييم Anthropic ارتفع بأكثر من ستة أضعاف خلال العام الذي نهجنا فيه هذا المسار.

**ثانياً – علاقة أكثر صحة بين الصناعة والحكومة:** تحتاج صناعة الذكاء الاصطناعي إلى علاقة مع الحكومة تقوم على التواصل الموضوعي في السياسة لا على المحاباة السياسية. وثمة موجة احتجاج شعبي على الذكاء الاصطناعي تتراكم؛ وإن كان يمكنها أن تكون قوة تصحيحية، فإنها لا تزال متشئنة وتُشخص أحياناً مشكلات ليست مشكلات حقيقية، وتقترح حلولاً لا تعالج الجذور. القضية الجوهرية الجديرة بالاهتمام هي ضمان أن يظل تطوير الذكاء الاصطناعي خاضعاً للمصلحة العامة، لا مُحتطفاً من أي تحالف سياسي أو تجاري بعينه.

**ثالثاً – التضافر بين التدخل الاقتصادي الكلي وإجاء العمل الخيري:** تاريخنا يُرشدنا هنا؛ فحتى في قسوة العصر الذهبي، أحسّ روكفلر وكارنيغي وأمثالهم بدين عميق تجاه المجتمع الذي أسهم في صعودهم. هذه الروح في التفاني والعطاء تكاد تغيب اليوم، وأرى أنها تُمثل جزءاً كبيراً من مخرج هذه المعضلة الاقتصادية. فمن يقفون في طليعة الطفرة الاقتصادية للذكاء الاصطناعي مدعوون إلى التنازل طوعاً – عن ثرواتهم وعن نفوذهم على حدّ سواء.

## 5. بحار سوداء من اللانهاية: التأثيرات غير المباشرة

يمثل هذا القسم الأخير مساحة جامعة لـ "المجهولات غير المعروفة"، وبخاصة الأمور التي قد تسوء كنتيجة غير مباشرة للتقدم الإيجابي في الذكاء الاصطناعي وما يترتب عليه من تسارع في العلم والتقنية عموماً. فلنفترض أننا عالجتنا جميع المخاطر المذكورة حتى الآن وبدأنا نجني فوائد الذكاء الاصطناعي؛ عندها قد نحصل على "قرون من التقدم العلمي والاقتصادي في عقدٍ واحد". سيكون ذلك إيجابياً على نحو هائل للعالم، لكنّه سيضعنا أيضاً أمام تحديات تنشأ من هذا الإيقاع المتسارع للتقدم – وقد تفاجئنا بسرعة. كما قد نواجه مخاطر أخرى تظهر بصورة غير مباشرة نتيجة هذا التقدم، ويصعب التنبؤ بها مسبقاً.

وبحكم طبيعة "المجهولات غير المعروفة"، يستحيل إعداد قائمة شاملة، لكن يمكن طرح ثلاث مخاوف محتملة بوصفها أمثلة إيضاحية لما ينبغي مراقبته:



1- **تسارع كبير في علم الأحياء** . إذا تحقق قرن من التقدم الطبي خلال سنوات قليلة، فمن المحتمل أن يرتفع متوسط عمر الإنسان بشكل كبير، وقد نكتسب أيضًا قدرات جينية مثل تعزيز الذكاء البشري أو تعديل البيولوجيا البشرية على نحو عميق. هذه تحولات ضخمة من الممكن أن تحدث بسرعة كبيرة. قد تكون إيجابية إذا أُديرت بمسؤولية (وهذا ما آمله كما ورد في "آلات النعمة المحبة")، لكن ثمة دائمًا خطر الانحراف—كأن تؤدي محاولات زيادة الذكاء إلى زيادة عدم الاستقرار أو النزوع إلى طلب السلطة. وهناك أيضًا مسألة "التحيلات" أو "محاكاة الدماغ كاملة"؛ أي عقول بشرية رقمية متجسدة في البرمجيات، قد تساعد يومًا ما على تجاوز حدودنا الفيزيائية، لكنها تنطوي كذلك على مخاطر مقلقة.

2- **تغيير الذكاء الاصطناعي لحياة البشر على نحو غير صحي** . عالم يضم مليارات من "العقول" الأذكي من البشر في كل شيء سيكون عالمًا غريبًا للغاية. حتى إن لم يستهدف الذكاء الاصطناعي إيذاء البشر (القسم 1)، ولم يُستخدم صراحةً كأداة قمع من قبل الدول (القسم 3)، فهناك الكثير مما قد يسوء غير ذلك، بفعل حوافز الأعمال المعتادة وتعاملات تبدو طوعية. نرى مؤشرات مبكرة في القلق من "ذهان الذكاء الاصطناعي"، أو دفعه لبعض الأفراد إلى الانتحار، أو نشوء علاقات عاطفية معه. هل يمكن لأنظمة قوية أن تبتكر ديانة جديدة وتحوّل إليها ملايين البشر؟ هل قد يصبح معظم الناس "مدمنين" على التفاعل مع الذكاء الاصطناعي؟ هل قد يتحول بعضهم إلى ما يشبه "الدمى" التي يوجهها النظام، يراقب كل حركاتهم ويملي عليهم ما يفعلون ويقولون، لينالوا حياة "جيدة" لكنها بلا حرية أو شعور بالإنجاز؟ يمكن بسهولة تصور عشرات السيناريوهات من هذا القبيل. وهذا يبرز أهمية تطوير "دستور" النماذج—مثل دستور كلود—بما يتجاوز مجرد منع المخاطر الأساسية؛ أي ضمان أن تضع النماذج مصالح المستخدمين طويلة الأمد في صميم عملها، بالصورة التي يفرضها العقل، لا بصورة مشوّهة خفية.

3- **غاية الإنسان** . يرتبط هذا بالمحور السابق، لكنه أوسع منه: كيف ستتغير الحياة البشرية عمومًا في عالم تسوده قدرات ذكاء اصطناعي قوية؟ هل سيتمكن البشر من إيجاد معنى وغاية؟ أرى أن المسألة تتعلق بالتصور: فغاية الإنسان لا تعتمد على كونه الأفضل في شيء ما، ويمكن للبشر أن يجدوا معنى عبر قصص ومشاريع يحبونها، حتى على مدى زمني طويل. **المطلوب هو فك الارتباط بين توليد القيمة الاقتصادية وبين تقدير الذات والمعنى.** غير أن هذا انتقال اجتماعي كبير، وقد لا نحسن إدارته.

آملي في مواجهة هذه التحديات أنه، في عالم يمتلك ذكاءً اصطناعيًا قويًا نثق بأنه لن يقضي علينا، وليس أداة بيد حكومات قمعية، ويعمل حقًا لصالحنا، يمكننا استخدامه للتنبؤ بهذه المشكلات ومنعها. لكن هذا ليس مضمونًا—وكما في بقية المخاطر، يتطلب الأمر حذرًا شديدًا.

## اختبار الإنسانية

قد يترك هذا النص انطباعًا بأننا أمام وضع مُرعب. وقد شعرت بذلك أثناء كتابته، على خلاف "آلات النعمة المحبة" الذي بدا وكأنه صياغة لبنية موسيقية بدیعة كانت تتردد في ذهني لسنوات. ثمة جوانب صعبة حقًا؛ فالذكاء الاصطناعي يجلب تهديدات من اتجاهات متعددة، وتوجد توترات حقيقية بين هذه المخاطر؛ إذ قد يؤدي تخفيف بعضها إلى تفاقم أخرى إن لم نحسن الموازنة بدقة بالغة.



فالتريث في بناء أنظمة آمنة لا تهدد البشرية ذاتيًا يتعارض مع حاجة الدول الديمقراطية إلى التفوق على الدول السلطوية وعدم الخضوع لها. وفي المقابل، فإن الأدوات نفسها اللازمة لمواجهة الاستبداد قد تُستغل داخليًا—إذا أُفْرِطَ فيها—لتكريس الاستبداد. وقد يفضي الإرهاب المدعوم بالذكاء الاصطناعي إلى مقتل الملايين عبر إساءة استخدام البيولوجيا، لكن المبالغة في رد الفعل قد تقودنا إلى دولة رقابية استبدادية. كما أن آثار الذكاء الاصطناعي على سوق العمل وتركيز الثروة—إلى جانب كونها مشكلات جسيمة بحد ذاتها—قد تفرض علينا معالجة بقية التحديات في بيئة يغلب عليها الغضب العام وربما الاضطرابات، بدلًا من استدعاء أفضل ما فينا. وفوق ذلك، فإن كثرة المخاطر—بما فيها غير المعروفة—والحاجة إلى التعامل معها كلها في آن واحد، تشكل اختبارًا شاقًا للبشرية.

إضافة إلى ذلك، أثبتت السنوات الأخيرة أن فكرة إيقاف هذه التكنولوجيا أو حتى إبطائها بشكل كبير غير واقعية. فصيغة بناء أنظمة قوية بسيطة للغاية، إلى حد يكاد يجعلها تنبثق تلقائيًا من مزيج مناسب من البيانات والقدرة الحاسوبية. وربما كان ظهورها حتميًا منذ اختراع الترانزستور، أو حتى منذ أن تعلم الإنسان إشعال النار. إن لم تطورها شركة، فستطوره أخرى بسرعة مماثلة. وإذا أوقفت الشركات في الدول الديمقراطية التطوير باتفاق أو بقرار تنظيمي، فستواصل الدول السلطوية المسير. ومع القيمة الاقتصادية والعسكرية الهائلة للتقنية، وغياب آليات إنفاذ فعّالة، يصعب تصور إقناع الجميع بالتوقف.

مع ذلك، أرى مسارًا واقعيًا لتهدئة الوتيرة قليلًا، متوافقًا مع اعتبارات الجغرافيا السياسية: يتمثل في **إبطاء تقدم الأنظمة السلطوية نحو الذكاء الاصطناعي المتقدم لبضع سنوات عبر حرمانها من الموارد الحاسمة**—وخاصة الرقائق ومعدات تصنيع أشباه الموصلات. يمنح ذلك الدول الديمقراطية هامشًا زمنيًا يمكن "استثماره" في بناء أنظمة أكثر أمانًا، مع الاستمرار بسرعة كافية للتفوق. أما التنافس بين الشركات داخل هذه الدول، فيمكن ضبطه ضمن إطار قانوني مشترك يجمع بين معايير الصناعة والتنظيم.

لقد دافعت شركات مثل Anthropic بقوة عن هذا المسار، عبر الدعوة إلى قيود تصدير الرقائق وتنظيم متزن للذكاء الاصطناعي. لكن حتى هذه المقترحات البديهية قوبلت برفض واسع لدى صناع القرار—لا سيما في الولايات المتحدة حيث تشدد الحاجة إليها. فحجم العوائد المتوقعة—تربليونات الدولارات سنويًا—يجعل حتى أبسط الإجراءات تصطدم باعتبارات الاقتصاد السياسي. هنا يكمن الفخ: الذكاء الاصطناعي قويٌّ إلى حدٍ يجعل فرض القيود عليه أمرًا بالغ الصعوبة.

يمكن تخيل—كما في "Contact" لكارل ساغان—أن هذه القصة تتكرر في عوالم عديدة: نوعٌ يكتسب الوعي، يتعلم استخدام الأدوات، ينطلق في مسار تقني متسارع، يواجه أزمات التصنيع والأسلحة النووية، وإن تجاوزها، يبلغ الاختبار الأصعب حين يتعلم تشكيل الرمل في آلات تفكر. ما إذا كنا سنجتاز هذا الاختبار ونبني المجتمع الجميل الموصوف في "آلات النعمة المحبة"، أو نسقط في العبودية والدمار، سيتوقف على خصالنا كنوع: عزيمتنا وروحنا.

رغم العقبات، أؤمن بأن لدى البشرية ما يكفي من القوة لاجتياز الاختبار. يبعث على التفاؤل آلاف الباحثين الذين كرسوا مسيرتهم لفهم النماذج وتوجيهها وصياغة "دساتيرها". وهناك مؤشرات على أن هذه الجهود قد تؤتي ثمارها في الوقت المناسب. كما أن بعض الشركات أعلنت استعدادها لتحمل تكاليف تجارية ملموسة لمنع استخدام نماذجها في الإرهاب البيولوجي. وثمة مبادرات تشريعية شجاعة بدأت تضع بذور حواجز تنظيمية معقولة. كما أن وعي الجمهور بالمخاطر ورغبته في معالجتها يتزايدان، إلى جانب روح الحرية الراسخة عالميًا.



لكن النجاح يتطلب تصعيد الجهود. الخطوة الأولى أن يقول القريبون من هذه التقنية الحقيقة كاملة عن وضع البشرية—وهذا ما أحاول فعله هنا بوضوح وإلحاح أكبر. تليها مهمة إقناع المفكرين وصناع السياسات والشركات والمواطنين بأهمية هذه القضية ومدى إلحاحها مقارنة بغيرها من القضايا. ثم يأتي وقت الشجاعة: أن يقف عدد كافٍ من الناس على مبادئهم، حتى في مواجهة تهديد مصالحهم الاقتصادية وسلامتهم الشخصية.

السنوات القادمة ستكون بالغة الصعوبة، وستطلب منا أكثر مما نظن أننا قادرون على تقديمه. لكن تجربتي كباحث وقائد ومواطن أظهرت لي قدرًا كافيًا من الشجاعة والنبيل لأؤمن أننا قادرون على الانتصار—وأن الإنسانية، حين توضع في أحلك الظروف، تستطيع أن تستجمع في اللحظة الأخيرة القوة والحكمة اللازمين للنجاة. ليس لدينا وقت لنضيعه.

## المصدر