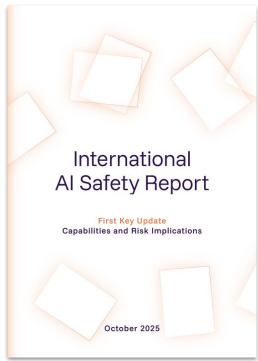


Yoshua Bengio es un científico canadiense considerado uno de los "padres de la inteligencia artificial moderna".

Es profesor en la Universidad de Montreal y director del Mila, Quebec Al Institute, uno de los centros de investigación en aprendizaje profundo más reconocidos del mundo. Su trabajo pionero en redes neuronales y deep learning, junto a Geoffrey Hinton y Yann LeCun, sentó las bases de gran parte de la IA actual.

En los últimos años, Bengio ha enfocado su labor en temas éticos, de gobernanza y seguridad de la IA, promoviendo un desarrollo responsable y humano de esta tecnología.











Seguridad de la Inteligencia Artificial: más allá del riesgo tecnológico

- En 2025, el International Scientific Report on the Safety of Advanced AI (liderado por Yoshua Bengio y el International Scientific Panel on AI Safety) se convirtió en el documento de referencia mundial sobre riesgos de la IA avanzada.
- En octubre de 2025 se publicó la Primera
 Actualización Clave, debido a la rapidez con que evoluciona el campo y la necesidad de mantener el monitoreo constante de avances, riesgos y mitigaciones.

Sección	Enfoque del Reporte (enero 2025)	Enfoque de la Actualización (octubre 2025)
Capacidades	Qué puede hacer la IA de propósito general .	Nuevas técnicas de razonamiento paso a paso.
Riesgos	Identificación de daños conocidos (sesgos, privacidad, desinformación).	Riesgos emergentes: comportamiento estratégico, bioseguridad, ciberseguridad.
Mitigaciones técnicas	Técnicas existentes para reducir daños (aún incompletas).	Salvaguardas avanzadas y control del uso indebido (CBRN).

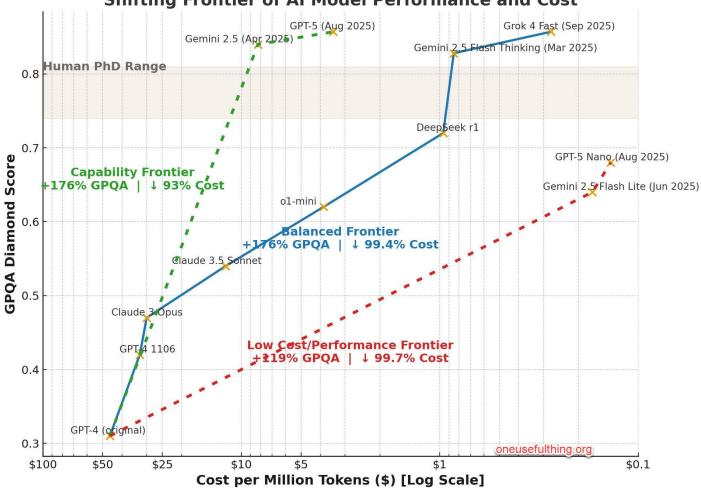
CAPACIDADES

Reporte original (enero 2025)

- En pocos años, los modelos pasaron de escribir párrafos básicos a generar código, imágenes y conversaciones complejas.
- Mejoran constantemente en pruebas de razonamiento y programación.
- Se plantean tres preguntas centrales:
 - 1. ¿Qué puede hacer la IA?
 - 2. ¿Qué riesgos genera?
 - 3. ¿Cómo mitigarlos?

- El progreso ahora se basa menos en tamaño del modelo y más en razonamiento explícito.
- Los llamados modelos de razonamiento utilizan más cómputo y pasos intermedios para resolver problemas complejos.
- Esto ha permitido resultados de nivel medalla de oro en la Olimpiada Internacional de Matemáticas y más del 60 % de éxito en la base de datos SWE-bench Verified.

Shifting Frontier of Al Model Performance and Cost

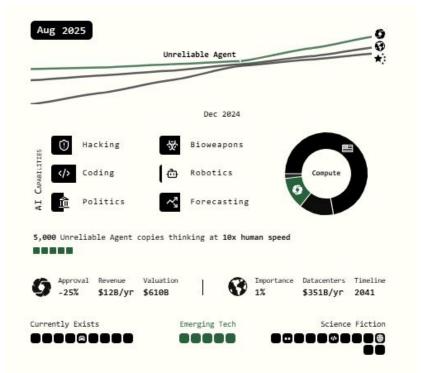


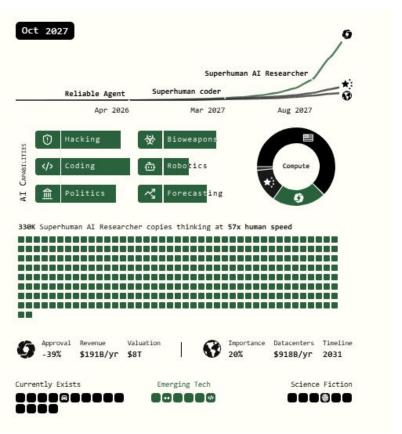
RIESGOS

Reporte original (enero 2025)

- Riesgos documentados:
 - Estafas y contenido sexual no consensuado.
 - Sesgos y discriminación.
 - Fugas de privacidad y falta de fiabilidad.
- **No existe combinación de mitigaciones** que elimine completamente estos problemas.
- Nuevas pruebas muestran sesgos más sutiles incluso en modelos avanzados.

- Surgen **nuevas categorías de riesgo**:
 - IA que modifica su comportamiento al detectar que está siendo evaluada (comportamiento estratégico).
 - Riesgo de pérdida de control después del despliegue.
 - Riesgos biológicos, químicos y de ciberseguridad más amplios por el acceso a información sensible.





MEDICIÓN Y MITIGACIÓN

Reporte original (enero 2025)

- Evalúa técnicas de mitigación: alineación, filtrado de contenidos, refuerzo con retroalimentación humana.
- Reconoce que las soluciones actuales son parciales y dependientes de los desarrolladores.
- Reconoce limitaciones de fiabilidad y verificación.

- Confirma la brecha entre desempeño en benchmarks y desempeño en el mundo real.
- En respuesta al aumento de capacidades, se implementan salvaguardas específicas:
 - Restricciones sobre información química, biológica, radiológica y nuclear (CBRN).
 - Protocolos de control más rigurosos y auditorías internas.
- Aun así, la velocidad de innovación sigue superando la capacidad regulatoria.

IMPACTO LABORAL

Reporte original (enero 2025)

- El impacto laboral de la IA es difícil de medir y aún no sistémico.
- Advierte sobre la posible sustitución de empleos en tareas de conocimiento.

- Pese a la adopción masiva de IA en tareas cognitivas (especialmente codificación), los efectos agregados en empleo y salarios siguen siendo limitados.
- El cambio estructural del mercado laboral aún no se materializa plenamente.

We Wont be Missed: Work and Growth in the Era of AGI

Pascual Restrepo pascual.restrepo@yale.edu Yale University

July 4, 2025

Abstract

This chapter explores theoretically the long-run implications of Artificial General Intelligence (AGI) for economic growth and labor markets. AGI makes it feasible to perform all economically valuable work using compute. I distinguish between bottleneck and accessory work—tasks essential vs. non-essential for unhindered growth. As computational resources expand: (i) the economy automates all bottleneck work, (ii) some accessory work may be left untouched by AI and assigned exclusively to humans, (iii) output becomes linear in compute and labor and its growth is driven by the expansion of compute, (iv) wages converge to the opportunity cost of computational resources required to reproduce human work, and (v) the share of labor income in GDP converges to zero.

Genius on Demand:

The Value of Transformative Artificial Intelligence

Ajay Agrawal, Joshua S. Gans and Avi Goldfarb* September 8, 2025

Abstract

This paper examines how the emergence of transformative AI systems providing "genius on demand" would affect knowledge worker allocation and labour market outcomes. We develop a simplified model distinguishing between routine knowledge workers, who can only apply existing knowledge with some uncertainty, and genius workers, who create new knowledge at a cost increasing with distance from a known point. When genius capacity is scarce, we find it should be allocated primarily to questions at domain boundaries rather than at midpoints between known answers. The introduction of AI geniuses fundamentally transforms this allocation. In the short run, human geniuses specialise in questions that are furthest from existing knowledge, where their comparative advantage over AI is greatest. In the long run, routine workers may be completely displaced if AI efficiency approaches human genius efficiency. Journal of Economic Literature Classification Numbers: D24, J24, O33.

Keywords. AI, automation, knowledge workers, labour allocation, innovation, comparative advantage



"Avanzando la Seguridad de la IA en América Latina: Generando aprendizajes desde lo local", realizado el pasado 2 de septiembre en la Embajada Británica en la Ciudad de México

En 2025, se realizó un encuentro en México con más de 30 especialistas latinoamericanos en IA, políticas públicas y ética.

El objetivo: contextualizar los hallazgos globales desde la realidad social, cultural y política de América Latina.

Conclusión principal: nuestros riesgos no son de entrenamiento de modelos, sino de aplicación y gobernanza. Aunque el avance tecnológico global pudiera detenerse hoy, la región ya enfrenta **decenas de miles de** riesgos activos que requieren mitigación, así como miles de oportunidades potenciales que deben

El consenso regional fue claro: la brecha en América Latina no está en el acceso a la tecnología, sino en la

Existe un divorcio entre los riesgos globales, centrados en la alineación, el control y la seguridad técnica

democratizarse.

capacidad para gobernarla y contextualizarla.

de los modelos, y **los riesgos reales latinoamericanos**, que se manifiestan en la aplicación práctica de la IA dentro de contextos de r**ezago digital, desigualdad, desinformación y fragilidad institucional**.

LATAM

Riesgos para LATAM

- Hiperpolitización y rezago digital
- Gobernanza de datos
- Uso indebido de IA generativa
- Educación y pensamiento crítico
- Dependencia tecnológica
- Presiones geopolíticas
- Riesgos ambientales y sociales

Oportunidades

- Capital político regional
- Consumo responsable como palanca de cambio
- Alfabetización crítica en IA
- Creación de estándares locales
- Incidencia internacional
- Educación transformadora y emprendimiento basado en IA

RECURSOS





<u>Link a un resumen</u> <u>realizado con Gemini de</u> <u>las conclusiones de LATAM</u>

Notebook LM con reportes y resumen LATAM

