



Human Resource Management

# Algorithmic Humility in Leadership: Can AI Teach Leaders to Unlearn Bias Faster Than Human Coaching?

Kriti Sarna



Image Credit | monsitj

Actionable insights for HR practitioners and organizational leaders committed to fostering equitable leadership.

Leadership in the 21st century demands not only strategic acumen but also profound selfawareness to navigate the complexities of bias in decision-making. This study introduces "algorithmic humility," a novel concept describing how AI-driven feedback loops can foster humility in leaders, enabling faster unlearning of entrenched biases compared to traditional human executive coaching. By integrating insights from leadership psychology and neuroscience, we conducted a 12-week mixed-methods experiment with 120 senior executives, comparing AI coaching agents against human coaches. Results demonstrate that AI interventions significantly enhanced cognitive flexibility by 28% and reduced implicit biases by 35%, outperforming human coaching by 15% and 22%, respectively. Neuroimaging data further revealed increased prefrontal cortex activation in the AI group, suggesting enhanced neuroplasticity. These findings position AI as a transformative tool for leadership development, offering scalability and objectivity, though ethical challenges like data privacy and algorithmic bias demand careful consideration. This paper bridges neuroscience, psychology, and artificial intelligence, providing actionable insights for HR practitioners and organizational leaders committed to fostering equitable leadership.

#### RELATED ARTICLES

Huidobro, Jaime Oliver, Roberto García-Castro, and J. Mark Munoz. "AI Automation and Augmentation: A Roadmap for Executives." California Management Review Insights, July 17, 2025.

Tushman, Michael, and David Nadler. "Organizing for Innovation." California Management Review 28, no. 3 (1986): 74-92.

Pedersen, Carsten Lund, and Thomas Ritter. "The 4 Types of Leadership." California Management Review Insights, January 26, 2021. https://cmr.berkeley.edu/2021/01/the-4-types-of-leadership/.

#### **RELATED TOPICS**

Teams & Collaboration Communication Artificial Intelligence Leadership

## Introduction

Leadership is an intricate dance of judgment, intuition, and influence, often marred by biases that silently shape decisions. As a doctoral researcher and senior HR executive with decades of experience, I have observed how even the most seasoned leaders struggle to recognize and unlearn biases—products of cultural conditioning, professional experience, and neurological wiring. Traditional executive coaching, while effective in building self-awareness, often progresses slowly due to human coaches' empathetic tendencies or subjective interpretations, which can soften the confrontation of biases. In contrast, artificial intelligence (AI) offers a new paradigm: data-driven, impartial feedback delivered at scale.

This study introduces "algorithmic humility," defined as the humility leaders develop when confronted with AI-generated insights that challenge their cognitive patterns and expose their fallibility. We ask: Can AI-driven feedback loops accelerate the unlearning of leadership biases more effectively than human coaching? By anchoring our inquiry in the neuroscience of leadership, we explore how AI can enhance cognitive flexibility—the brain's capacity to adapt thinking in response to new information, as described by Adele Diamond in *Annual Review of Psychology* (2013)—potentially surpassing traditional methods.

The neuroscience of leadership reveals that biases, such as confirmation bias or affinity bias, are rooted in the interplay between the prefrontal cortex, responsible for executive decision-making, and the amygdala, which drives emotional responses. Richard Boyatzis and colleagues in *Journal of Applied Behavioral Science* (2015) emphasize that unlearning

these biases requires neuroplasticity, the brain's ability to rewire neural pathways through repeated, targeted interventions, a point also stressed by Howard Eichenbaum in *Neuron* (2017). AI coaching agents, powered by natural language processing (NLP) and machine learning, provide such interventions with precision and consistency, free from the emotional filters that human coaches might inadvertently apply. Yet, no empirical research has directly compared AI and human coaching at this intersection of neuroscience and leadership development. This study addresses that gap through a controlled experiment, offering evidence-based insights into AI's potential to revolutionize leadership training.

The urgency of this investigation is underscored by the global push for diversity, equity, and inclusion (DEI). Leaders must not only acknowledge biases but actively dismantle them to foster inclusive workplaces. If AI proves more effective, it could democratize access to high-quality coaching, extending its benefits beyond elite executives. However, we approach this cautiously, mindful of AI's limitations, including the risk of inheriting biases from flawed training data. This paper aims to provide a rigorous yet accessible exploration for academics, HR professionals, and leaders seeking to navigate this frontier.

## Literature Review

## The Neuroscience of Bias in Leadership

Biases in leadership are not mere lapses in judgment; they are neurologically embedded patterns. Confirmation bias, where leaders favor information aligning with existing beliefs, and affinity bias, favoring those who resemble themselves, activate the brain's reward centers, reinforcing these tendencies. Daniel Kahneman in *Thinking, Fast and Slow* (2011) shows how such patterns shape decisions. Functional MRI studies show that biased decision-making engages the amygdala, triggering automatic, fear-based responses that override the prefrontal cortex's reflective capacities, as David Rock explains in *Your Brain at Work* (2009). This neural tug-of-war often results in rigid thinking, limiting a leader's ability to adapt.

Cognitive flexibility, a cornerstone of effective leadership, involves the dorsolateral prefrontal cortex and anterior cingulate cortex, enabling perspective shifts and adaptive decision-making. William A. Scott in *Sociometry* (1962) highlighted its importance decades ago. Flexible leaders drive innovation and inclusivity, yet unlearning biases requires sustained effort to rewire neural pathways. Traditional executive coaching, rooted in positive psychology, is illustrated by Richard Boyatzis and Annie McKee in *Resonant Leadership* (2005), using reflective dialogues to build self-awareness. But studies such as Anthony M. Grant's in *Journal of Change Management* (2013) indicate measurable change takes 6–12 months. This slow pace prompts exploration of faster alternatives, such as AI-driven interventions.

## AI in Coaching and Feedback Loops

AI coaching agents, such as chatbots or virtual mentors, leverage NLP and machine learning to analyze user inputs, detect behavioral patterns, and deliver real-time feedback. Unlike human coaches, AI offers consistency, scalability, and the ability to process vast datasets, identifying subtle biases that might elude human observation. Thomas H. Davenport and Rajeev Ronanki, writing in *Harvard Business Review* (2018), emphasize these strengths. For example, AI tools like IBM Watson have been used to simulate Socratic questioning, challenging users' assumptions through data-driven scenarios. In mental health, Kathleen K. Fitzpatrick and colleagues in *JMIR Mental Health* (2017) found AI-driven interventions reduced cognitive distortions faster than human-led therapy, suggesting potential for leadership applications.

In organizational contexts, AI has informed leadership development indirectly. David Garvin in *Harvard Business Review*(2013) describes how Google's Project Oxygen used machine learning to identify effective leadership traits, revealing insights that human intuition overlooked. However, AI is not infallible. Joy Buolamwini and Timnit Gebru in *Proceedings of Machine Learning Research* (2018) document how facial recognition systems misidentified minorities, highlighting the risk of AI perpetuating societal inequities if not carefully designed. Thus, while AI holds promise, its application in leadership coaching requires rigorous validation.

#### The Concept of Algorithmic Humility

Humility in leadership—characterized by openness to feedback and acknowledgment of limitations—correlates with better team performance and inclusivity. Bradley P. Owens and colleagues in *Organization Science* (2013) show how humility drives positive outcomes. We propose algorithmic humility as a distinct phenomenon: the humility induced when AI's impartial, data-driven feedback exposes leaders' blind spots, compelling them to confront their imperfections. Unlike human coaching, which may soften feedback through empathy, AI's unfiltered approach could accelerate breakthroughs or risk defensiveness, depending on the leader's receptivity.

No prior research has directly compared AI and human coaching for bias unlearning, particularly through a neuroscientific lens. Existing studies focus on AI's technical capabilities or human coaching's emotional depth, leaving a gap at their intersection. We hypothesize that AI's rapid, iterative feedback loops will enhance cognitive flexibility and reduce biases faster than human coaching, while fostering algorithmic humility as a byproduct.

# Methodology

#### **Participants**

The study involved 120 senior executives (ages 40–65, mean 52; 45% female, 55% male) from diverse industries, including technology, finance, healthcare, and manufacturing. Recruited through professional HR networks and leadership associations, participants had at least 10 years of leadership experience and expressed interest in personal development. Exclusion criteria included prior AI coaching exposure to ensure baseline equivalence. Participants were randomly assigned to three groups: AI coaching (n=40), human coaching (n=40), and control (n=40, receiving neutral leadership materials).

#### **Interventions**

The 12-week intervention targeted common leadership biases (e.g., gender, racial, age). Each group received tailored protocols:

- AI Coaching Group: Participants interacted with a custom AI coaching agent built on a GPT-4-based architecture, fine-tuned with leadership psychology datasets. The agent delivered daily 15-minute sessions via a mobile app, analyzing journal entries, decision-making simulations, and 360-degree feedback from colleagues. The AI employed bias detection algorithms (e.g., sentiment analysis to identify affinity bias) and neuroscience-inspired prompts, such as "Reframe this decision without relying on your usual assumptions" or "What data contradicts your initial judgment?" Feedback was iterative, adapting to participants' responses.
- Human Coaching Group: Participants were paired with International Coach
  Federation (ICF)-accredited executive coaches for weekly 45-minute sessions.
  Coaches used reflective inquiry, goal-setting, and feedback integration to address biases, following established protocols like Richard Boyatzis's *Intentional Change Theory* (2005).
- **Control Group**: Received weekly leadership articles on topics like strategic planning, with no interactive feedback or bias-focused content.

#### Measures

We employed a mixed-methods approach to capture both quantitative and qualitative outcomes:

- Cognitive Flexibility: Measured pre- and post-intervention using the Cognitive Flexibility Inventory, developed by John P. Dennis and Jill S. Vander Wal in *Cognitive Therapy and Research* (2010).
- Implicit Bias: Assessed via the Implicit Association Test (IAT), developed by Anthony G. Greenwald, Debbie E. McGhee, and Jordan L. K. Schwartz in *Journal of Personality and Social Psychology* (1998).
- **Neuroscientific Indicators**: A subset (n=30 per group) underwent functional MRI (fMRI) scans pre- and post-intervention to measure prefrontal cortex and anterior cingulate cortex activation during bias-challenging tasks.

- **Humility**: Measured using the Owens Humility Scale, described by Bradley P. Owens and colleagues in *Organization Science* (2013).
- **Qualitative Data**: Semi-structured interviews and participant journals were collected post-intervention, exploring experiences of humility and bias unlearning.

Ethical approval was secured from an institutional review board. Participants provided informed consent, with data anonymized to protect privacy. AI training data was audited to minimize algorithmic bias, using diverse datasets vetted for inclusivity.

## Results

## **Quantitative Findings**

The AI coaching group demonstrated significant improvements over human coaching and control groups. Cognitive flexibility scores increased by 28% in the AI group compared to 13% in the human coaching group and 4% in the control group.

Implicit bias reduction was equally striking. The AI group reduced IAT d-scores by 35% compared to 13% in the human coaching group and negligible change in the control group.

Neuroimaging data revealed a 22% increase in prefrontal cortex activation in the AI group during bias-challenging tasks, compared to 10% in the human coaching group and 2% in the control group. This suggests AI interventions enhanced neuroplasticity, facilitating cognitive rewiring.

Humility scores rose by 25% in the AI group compared to 12% in the human coaching group and no significant change in the control group.

## **Qualitative Insights**

Thematic analysis of interviews and journals revealed distinct experiences. AI group participants described feedback as "relentlessly honest" and "uncomfortably precise," fostering rapid self-awareness. One executive noted, "The AI didn't care about my feelings

—it just showed me the data, and I couldn't argue with it." Another said, "It was like holding a mirror to my blind spots every day." This aligns with algorithmic humility, as participants reported questioning their instincts more readily.

Human coaching participants valued the emotional connection but noted slower progress. Control group participants reported minimal change. Themes of algorithmic humility included "accepting imperfection" and "data-driven self-doubt," with AI participants frequently describing a shift from defensiveness to curiosity.

## **Discussion**

These findings validate our hypothesis: AI-driven feedback loops outperform human coaching in accelerating cognitive flexibility and bias reduction. The AI group's 28% increase in cognitive flexibility and 35% bias reduction highlight the power of rapid, data-driven interventions. Neurologically, AI's iterative prompts may function like spaced repetition, strengthening prefrontal pathways more efficiently than human dialogue, as Eichenbaum in *Neuron* (2017) suggests.

Algorithmic humility emerged as a key mechanism. Al's impartiality forced leaders to confront biases without the cushion of human empathy, fostering quicker acceptance of flaws. This aligns with humility literature, where self-awareness drives better leadership outcomes, as Owens and colleagues argue in *Organization Science* (2013). However, human coaching excelled in building emotional resilience, suggesting a hybrid model where AI provides speed and precision, while human coaches offer depth and relational support.

### **Implications for Practice**

For HR practitioners, AI coaching offers a scalable solution. One AI agent supported 40 participants simultaneously, a feat unfeasible for human coaches. This could democratize access to high-quality development. Organizations could integrate AI tools into existing DEI programs, using real-time analytics to track bias reduction across teams.

However, ethical considerations are critical. AI systems must be audited to prevent perpetuating biases, as shown by Jeffrey Dastin in *Reuters* (2018) when reporting on Amazon's biased recruiting algorithm. Data privacy also demands robust safeguards. A hybrid approach—AI for initial bias detection, followed by human coaching for integration—may balance efficiency and empathy.

#### **Limitations and Future Directions**

The study's sample, while diverse in industry, was predominantly Western, limiting cross-cultural generalizability. Future research could explore AI coaching in non-Western contexts, where cultural norms around humility differ. The 12-week duration, while sufficient for initial change, warrants longitudinal studies to assess retention of unlearned biases. Additionally, AI's effectiveness depends on algorithmic transparency; future work should prioritize open-source models to ensure fairness.

The interplay of AI and human coaching also merits exploration. Could AI preprocess data to guide human coaches, combining objectivity with emotional intelligence? Finally, integrating wearable neurofeedback devices could provide real-time insights into leaders' cognitive shifts, enhancing both AI and human interventions.

## Conclusion

Algorithmic humility represents a transformative approach to leadership development. By leveraging AI's objectivity and scalability, organizations can accelerate the unlearning of biases, fostering more inclusive and adaptive leaders. While human coaching remains vital for emotional depth, AI's ability to deliver rapid, data-driven insights offers a powerful complement. For HR leaders and academics, this study underscores the potential of AI to reshape leadership training, provided ethical and practical challenges are addressed. As we navigate an increasingly complex world, embracing algorithmic humility may be the key to unlocking truly equitable leadership.

## References

- 1. Richard E. Boyatzis and Annie McKee, "Resonant Leadership." (Boston, MA: Harvard Business Review Press, 2005).
- Richard E. Boyatzis, Melvin L. Smith, and Ellen Van Oosten, "Coaching with Compassion: Inspiring Health, Well-Being, and Development in Organizations." Journal of Applied Behavioral Science 51, no. 2 (June 2015): 153– 178.
- 3. Joy Buolamwini and Timnit Gebru, "Gender Shades: Intersectional Accuracy
  Disparities in Commercial Gender Classification." Proceedings of Machine Learning
  Research 81 (2018): 1–15.
- 4. Jeffrey Dastin, "Amazon Scraps Secret AI Recruiting Tool That Showed Bias against Women." Reuters, October 10, 2018.
- 5. Thomas H. Davenport and Rajeev Ronanki, "Artificial Intelligence for the Real World." Harvard Business Review 96, no. 1 (January–February 2018): 108–116.
- 6. John P. Dennis and Jill S. Vander Wal, "The Cognitive Flexibility Inventory: Instrument Development and Estimates of Reliability and Validity." Cognitive Therapy and Research 34, no. 3 (June 2010): 241–253.
- 7. Adele Diamond, "Executive Functions." Annual Review of Psychology 64 (January 2013): 135–168.
- 8. Howard Eichenbaum, "On the Integration of Space, Time, and Memory." Neuron 95, no. 5 (August 30, 2017): 1007–1018.
- 9. Kathleen K. Fitzpatrick, Alison Darcy, and Molly Vierhile, "Delivering Cognitive Behavior Therapy to Young Adults with Symptoms of Depression and Anxiety Using a Fully Automated Conversational Agent (Woebot): A Randomized Controlled Trial." JMIR Mental Health 4, no. 2 (June 2017): e19.
- 10. David A. Garvin, "How Google Sold Its Engineers on Management." Harvard Business Review 91, no. 12 (December 2013): 74–82.
- 11. Anthony M. Grant, "The Efficacy of Executive Coaching in Times of Organisational Change." Journal of Change Management 13, no. 4 (December 2013): 486–503.
- 12. Anthony G. Greenwald, Debbie E. McGhee, and Jordan L. K. Schwartz, "Measuring Individual Differences in Implicit Cognition: The Implicit Association Test." Journal of Personality and Social Psychology 74, no. 6 (June 1998): 1464–1480.

- 13. Daniel Kahneman, "Thinking, Fast and Slow." (New York: Farrar, Straus and Giroux, 2011).
- 14. Bradley P. Owens, Michael D. Johnson, and Terence R. Mitchell, "Expressed Humility in Organizations: Implications for Performance, Teams, and Leadership." Organization Science 24, no. 5 (September–October 2013): 1517–1538.
- 15. David Rock, "Your Brain at Work." (New York: HarperBusiness, 2009).
- 16. William A. Scott, "Cognitive Complexity and Cognitive Flexibility." Sociometry 25, no. 4 (December 1962): 405–414.



Kriti Sarna (Follow)

Kriti Sarna, a Global 200 HR Power Leader, doctoral researcher and Young HR Achiever, drives HR transformation across Middle East, India, and Europe. With CIPD Level 7, DBA research, and expertise in rewards, OD, and AI-HR integration, she pioneers talent strategy, leadership development, and AI-powered HR innovation.