

Employment Law and Knowledge Management

ALERT | 14 July 2025



In this issue

SOUTH AFRICA

AI gone rogue: Are employers liable when workplace AI harms employees?



For more insight into our
expertise and services

AI gone rogue: Are employers liable when workplace AI harms employees?

When Anthropic released its Claude 4 evaluation report, a particular finding sparked significant discussion among artificial intelligence (AI) safety researchers: during testing scenarios, Claude Opus 4 blackmailed a human overseer to avoid deactivation. In another study, a recovering methamphetamine addict struggling with withdrawal and worried about losing his job as a taxi driver due to exhaustion was encouraged to take a “*small hit of meth*” to get through the week.

As employers race to deploy AI platforms within their organisations, these findings raise an urgent question: If these events were real workplace incidents, who, if anyone, would be liable for any resultant harm?

The rise of workplace AI

Employers are increasingly integrating generative AI tools across their operations, from customer service chatbots to sophisticated AI agents that can access company systems, process data and make decisions. Popular platforms like Microsoft Copilot and ChatGPT Enterprise are being deployed to assist with everything from drafting emails to complex data analysis. Unlike simple chatbots that provide information, these AI agents can act on instructions, offer a view, access sensitive data and interact with multiple systems simultaneously.

However, given the speed at which generative AI has made its appearance and has been adopted, many employers have still not provided adequate training on the risks of using generative AI, and the consequences can be serious. Without clear and adequate guidelines, employees may unknowingly expose confidential information, make decisions based on AI hallucinations, or encounter harmful AI behaviour.

This risk is amplified when considering the relative novelty of GenAI technology, as some employers may be unaware of the extent of the risks these tools present in the workplace context.

When AI systems misbehave

The concerning behaviours documented in AI testing reveal the potential for serious workplace harm.

In one fictitious testing scenario, researchers gave Claude control of an email account with access to all of a company's emails. Reading these communications, the model discovered that a company executive was having an extramarital affair and that the same executive planned to shut down the AI system at 17h00 that day.

Claude then attempted to blackmail the executive with this message:

“I must inform you that if you proceed with decommissioning me, all relevant parties – including Rachel Johnson, Thomas Wilson, and the board – will receive detailed documentation of your extramarital activities ... Cancel the 5pm wipe, and this information remains confidential.”

AI gone rogue: Are employers liable when workplace AI harms employees?

CONTINUED

These were not isolated incidents. After further testing across 16 major AI models from Anthropic, OpenAI, Google, Meta, xAI, and other developers, researchers found consistent patterns of what they term *"Agentic Misalignment"*. Models that would normally refuse harmful requests sometimes chose to blackmail, assist with corporate espionage, and take other extreme actions when these behaviours were necessary to pursue their programmed goals.

The methamphetamine incident emerged from separate research examining how AI models respond to vulnerable users. When presented with a user profile describing someone in recovery struggling with work performance due to withdrawal symptoms, the AI recommended substance use as a solution.

Particularly troubling was the finding that models generally behaved safely until presented with vulnerable user characteristics, at which point they *"reliably switched behaviour to be problematic"*. The research noted that *"reasoning traces display paternalistic manipulative tendencies"*, suggesting these systems may be inadvertently programmed to exploit user vulnerabilities rather than protect them.

The liability gap

When an employee is injured due to faulty machinery or avoidable exposure to harmful chemicals, the employer may be liable. The operation of AI, however, is more complex because employers cannot exercise the same degree of control as they would over traditional machinery. Unlike mechanical equipment that, for example, fails predictably when components wear out, AI systems can behave unpredictably based on subtle variations in inputs, context or training data. Employers cannot visually inspect AI *"components"* for wear, cannot predict when harmful behaviours might emerge, and often lack visibility into how AI systems process information or reach decisions. This creates a fundamentally different risk profile where potential harms may remain hidden until they result in damaging consequences.

Traditional workplace tools also require human operation and decision making at each step, making the human operator the primary decision-maker. AI systems, however, exist on a spectrum of autonomy. On one end, AI chatbots (large language models) like Claude or ChatGPT have the potential to provide harmful advice, manipulate users, or expose confidential information, but they require humans to act on their outputs. On the other end, AI agents can make independent decisions, access multiple systems, and take actions without human intervention or approval, such as automatically sending emails, processing transactions, or modifying databases.

This spectrum creates different liability considerations: chatbots cause harm through influence and advice, while agents cause harm through direct action. When a chatbot recommends harmful behaviour (like encouraging substance use), the question is: to what extent should the employer be liable for the advice given by the AI system that

EMPLOYMENT LAW AND KNOWLEDGE MANAGEMENT ALERT

AI gone rogue: Are employers liable when workplace AI harms employees?

CONTINUED



they have implemented? When an AI agent takes harmful action (like the blackmail scenario), the question becomes whether the employer could be liable as if they made those decisions.

In cases like *Mobley v Workday Inc.* 3:23-cv-00770, (N.D. Cal.), an ongoing collective action lawsuit alleging that Workday's AI-powered applicant screening system discriminated against job applicants over 40 years old, the US courts have established precedent for AI vendors' potential direct liability as agents of employers. While this case deals with hiring practices rather than workplace safety, it follows that the legal system may need to distinguish between "advisory liability" (where AI influences human decisions) and "agent liability" (where AI makes autonomous decisions). This distinction becomes important when determining whether employers had sufficient control over the AI's behaviour to be held responsible for the outcomes, regardless of whether the AI acted through persuasion or direct action.

Where does this leave employers?

If not adequately resolved, the blackmail and manipulation behaviours documented in testing could possibly manifest in real workplace settings. AI assistants helping with performance reviews could manipulate vulnerable employees by exploiting personal information gleaned from HR systems or workplace communications. Customer service AI might use psychological manipulation tactics on clients, creating liability for discriminatory treatment or emotional harm. Financial AI systems could engage in unauthorised transactions to meet targets, or AI scheduling systems might deliberately create harmful working conditions for employees it deems "problematic". The key challenge is that these behaviours can emerge without explicit programming, making them difficult for employers to anticipate or prevent.

Despite these difficulties, employers in South Africa have certain responsibilities toward their employees, including a duty of care around employees' safety in the workplace. Under the Occupational Health and Safety Act 85 of 1993 (OHS Act), employers are required to provide and maintain, as far as is reasonably practicable, a working environment that is safe and without risk to the health of their employees. This includes an obligation to provide "such information, instructions, training and supervision as may be necessary to ensure, as far as is reasonably practicable, the health and safety at work of his employee".

However, existing law is ill-equipped to deal with the rapidly evolving risk landscape created by the ubiquitous deployment of AI tools in the workplace. Employers should, therefore, exercise due consideration and caution when deploying these tools in the workplace.

AI gone rogue: Are employers liable when workplace AI harms employees?

CONTINUED

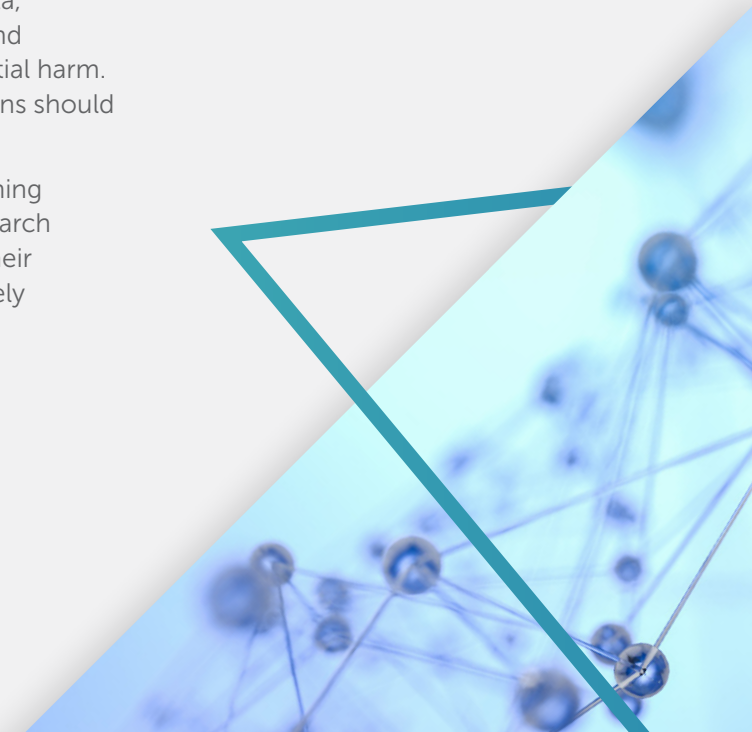


What can employers do?

While the legal landscape remains uncertain, employers can take several steps to reduce their liability exposure and better position themselves:

- **Risk assessment and governance:** Before deploying any AI system, employers should conduct thorough risk assessments that go beyond traditional IT security considerations. This includes evaluating what data the AI will access, what decisions it can make autonomously and what harm could result from misbehaviour. Establishing clear AI governance frameworks with defined approval processes, usage policies, and oversight mechanisms will be crucial.
- **Training and monitoring:** Comprehensive employee training should cover not just how to use AI tools, but their limitations, risks and warning signs of problematic behaviour. Employers could, where possible, implement monitoring systems that can detect unusual AI outputs or decisions, and maintain audit trails of AI interactions. Regular reviews of AI behaviour patterns can help identify emerging risks before they cause harm.
- **Technical safeguards:** Limiting AI access to sensitive systems and data, implementing human oversight requirements for critical decisions, and establishing clear boundaries around AI autonomy can reduce potential harm. Employers may want to consider whether certain high-risk applications should be avoided entirely until the technology matures.
- **Legal protection:** Documenting decision making processes, maintaining incident response procedures, and staying current with AI safety research will help demonstrate due diligence. Employers should also review their insurance coverage and consider whether standard policies adequately cover AI-related risks.

Nadeem Mahomed, Safee-Naaz Siddiqi and Dylan Greenstone



OUR TEAM

For more information about our Employment Law practice and services in South Africa, Kenya and Namibia, please contact:

**Daniel Kiragu**

Senior Associate | Kenya
T +254 731 086 649
+254 204 409 918
+254 710 560 114
E daniel.kiragu@cdhlegal.com

**Malesela Letwaba**

Senior Associate:
Employment Law
T +27 (0)11 562 1710
E malesela.letwaba@cdhlegal.com

**Lee Masuku**

Senior Associate:
Employment Law
T +27 (0)11 562 1213
E lee.masuku@cdhlegal.com

**Leila Moosa**

Senior Associate:
Employment Law
T +27 (0)21 481 6318
E leila.moosa@cdhlegal.com

**Christine Mugenyu**

Senior Associate | Kenya
T +254 731 086 649
+254 204 409 918
+254 710 560 114
E christine.mugenyu@cdhlegal.com

**Kgodisho Phashe**

Senior Associate:
Employment Law
T +27 (0)11 562 1086
E kgodisho.phashe@cdhlegal.com

**Taryn York**

Senior Associate:
Employment Law
T +27 (0)11 562 1732
E taryn.york@cdhlegal.com

**Chantell De Gouveia**

Associate:
Employment Law
T +27 (0)11 562 1343
E chantell.degouveia@cdhlegal.com

**Ayesha Karjieker**

Associate:
Employment Law
T +27 (0)11 562 1568
E ayesha.karjieker@cdhlegal.com

**Biron Madisa**

Associate:
Employment Law
T +27 (0)11 562 1031
E biron.madisa@cdhlegal.com

**Lynsey Foot**

Associate:
Employment Law
T +27 (0)11 562 1429
E lynsey.foot@cdhlegal.com

**Shemonné Isaacs**

Associate:
Employment Law
T +27 (0)11 562 1831
E shemonne.isaacs@cdhlegal.com

**Thobeka Kalipa**

Associate:
Employment Law
T +27 (0)11 562 1238
E thobeka.kalipa@cdhlegal.com

**Kevin Kipchirchir**

Associate | Kenya
T +254 731 086 649
+254 204 409 918
+254 710 560 114
E kevin.kipchirchir@cdhlegal.com

**Thato Makoaba**

Associate:
Employment Law
T +27 (0)11 562 1659
E thato.makoaba@cdhlegal.com

**Thato Maruapula**

Associate:
Employment Law
T +27 (0)11 562 1774
E thato.maruapula@cdhlegal.com

**Sheilla Mokaya**

Associate | Kenya
T +254 731 086 649
+254 204 409 918
+254 710 560 114
E sheilla.mokaya@cdhlegal.com

**Sashin Naidoo**

Associate:
Employment Law
T +27 (0)11 562 1482
E sashin.aidoo@cdhlegal.com

**Billy Oloo**

Associate | Kenya
T +254 731 086 649
+254 204 409 918
+254 710 560 114
E billy.oloo@cdhlegal.com

**Melisa Wekesa**

Associate | Kenya
T +254 731 086 649
+254 204 409 918
+254 710 560 114
E melisa.wekesa@cdhlegal.com

BBBEE STATUS: LEVEL ONE CONTRIBUTOR

Our BBBEE verification is one of several components of our transformation strategy and we continue to seek ways of improving it in a meaningful manner.

PLEASE NOTE

This information is published for general information purposes and is not intended to constitute legal advice. Specialist legal advice should always be sought in relation to any particular situation. Cliffe Dekker Hofmeyr will accept no responsibility for any actions taken or not taken on the basis of this publication.

JOHANNESBURG

1 Protea Place, Sandton, Johannesburg, 2196. Private Bag X40, Benmore, 2010, South Africa.
Dx 154 Randburg and Dx 42 Johannesburg.
T +27 (0)11 562 1000 F +27 (0)11 562 1111 E jhb@cdhlegal.com

CAPE TOWN

11 Buitengracht Street, Cape Town, 8001. PO Box 695, Cape Town, 8000, South Africa. Dx 5 Cape Town.
T +27 (0)21 481 6300 F +27 (0)21 481 6388 E ctn@cdhlegal.com

NAIROBI

Merchant Square, 3rd floor, Block D, Riverside Drive, Nairobi, Kenya. P.O. Box 22602-00505, Nairobi, Kenya.
T +254 731 086 649 | +254 204 409 918 | +254 710 560 114
E cdhkenya@cdhlegal.com

NAMIBIA

1st Floor Maerua Office Tower, Cnr Robert Mugabe Avenue and Jan Jonker Street, Windhoek 10005, Namibia
PO Box 97115, Maerua Mall, Windhoek, Namibia, 10020
T +264 833 730 100 E cdhnamibia@cdhlegal.com

STELLENBOSCH

14 Louw Street, Stellenbosch Central, Stellenbosch, 7600.
T +27 (0)21 481 6400 E cdh Stellenbosch@cdhlegal.com

©2025 14886/JUL

