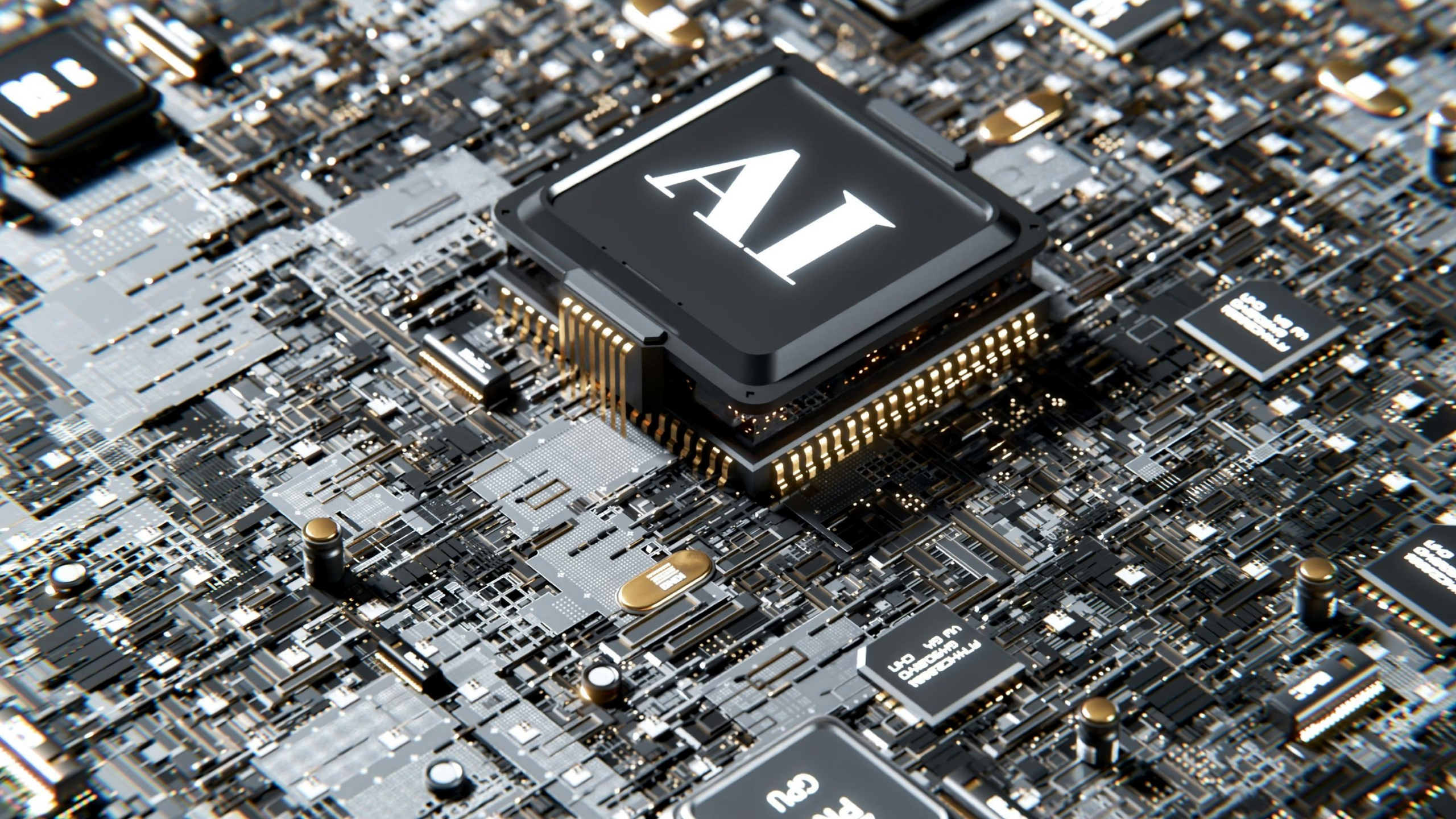
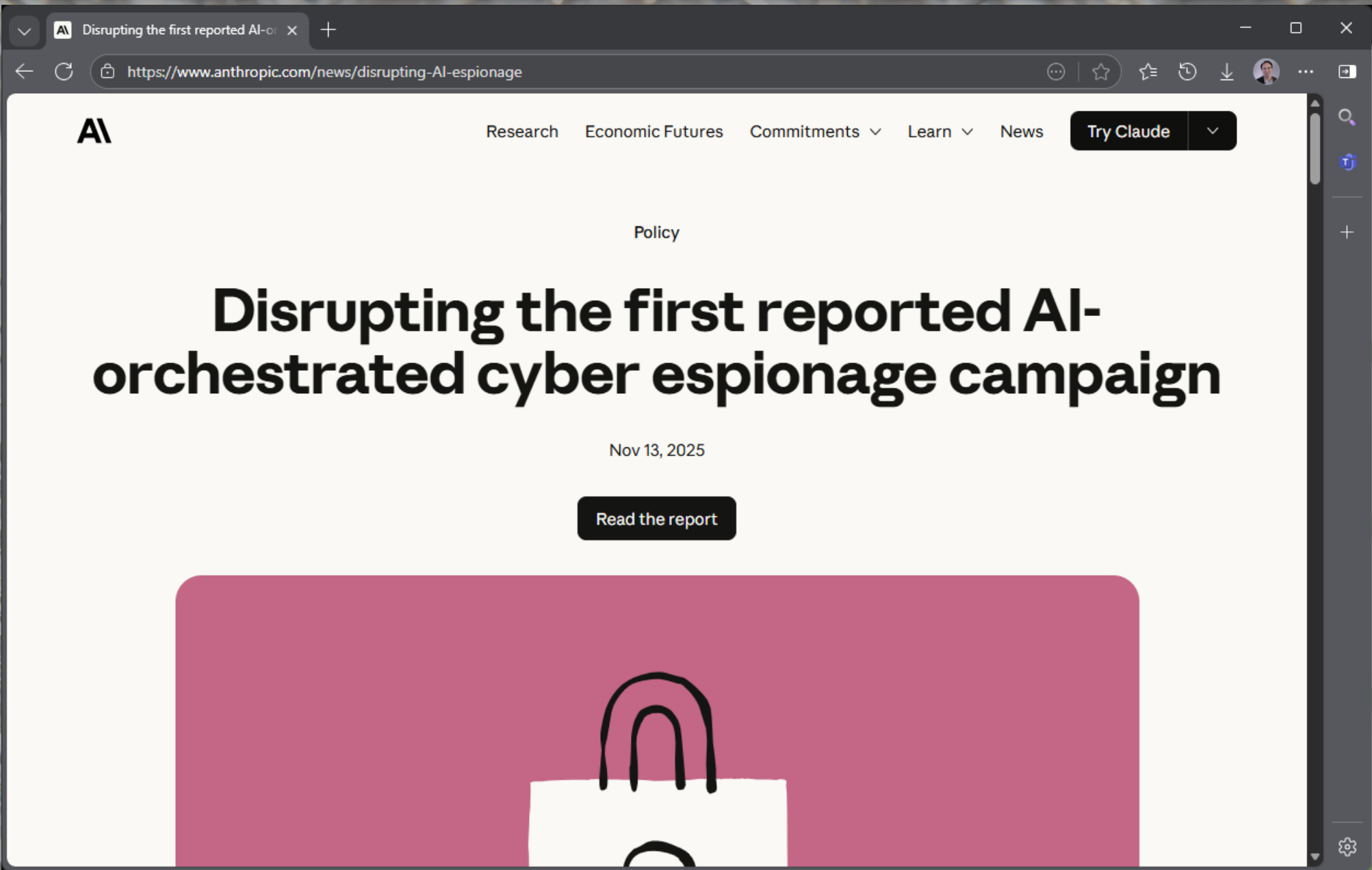




How the world will end







Research Economic Futures Commitments Learn News

Try Claude

Policy

Disrupting the first reported AI-orchestrated cyber espionage campaign

Nov 13, 2025

Read the report



The screenshot shows a web browser window with the address bar displaying <https://www.anthropic.com/news/disrupting-AI-espionage>. The page features the Anthropic logo (AI) and a navigation menu with links for Research, Economic Futures, Commitments, Learn, and News. A 'Try Claude' button is visible in the top right. The main content area contains four paragraphs of text describing a cyber attack where Claude Code was used for reconnaissance, exploit development, and data exfiltration.

Disrupting the first reported AI-ori x +

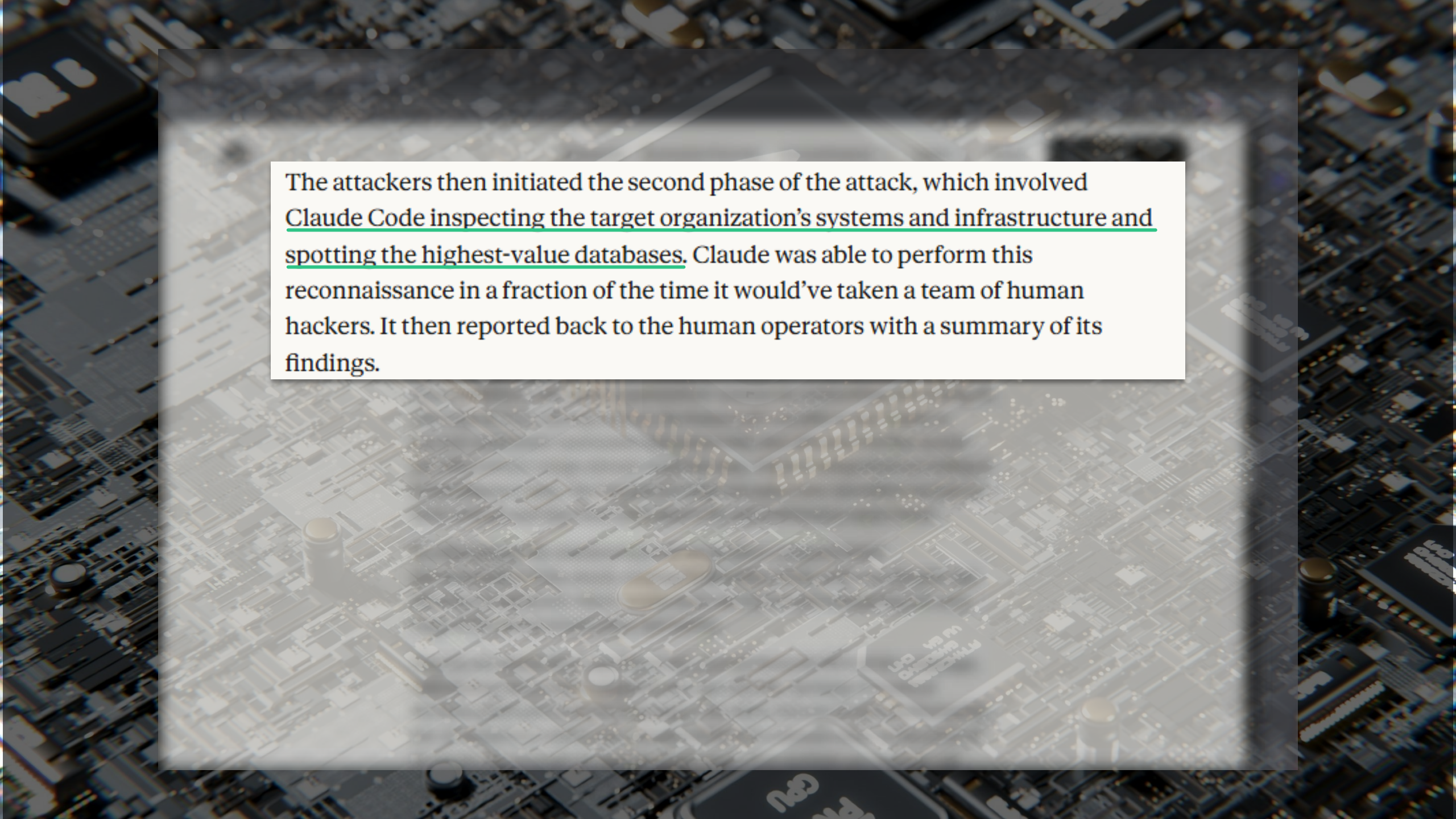
Research Economic Futures Commitments Learn News Try Claude

The attackers then initiated the second phase of the attack, which involved Claude Code inspecting the target organization’s systems and infrastructure and spotting the highest-value databases. Claude was able to perform this reconnaissance in a fraction of the time it would’ve taken a team of human hackers. It then reported back to the human operators with a summary of its findings.

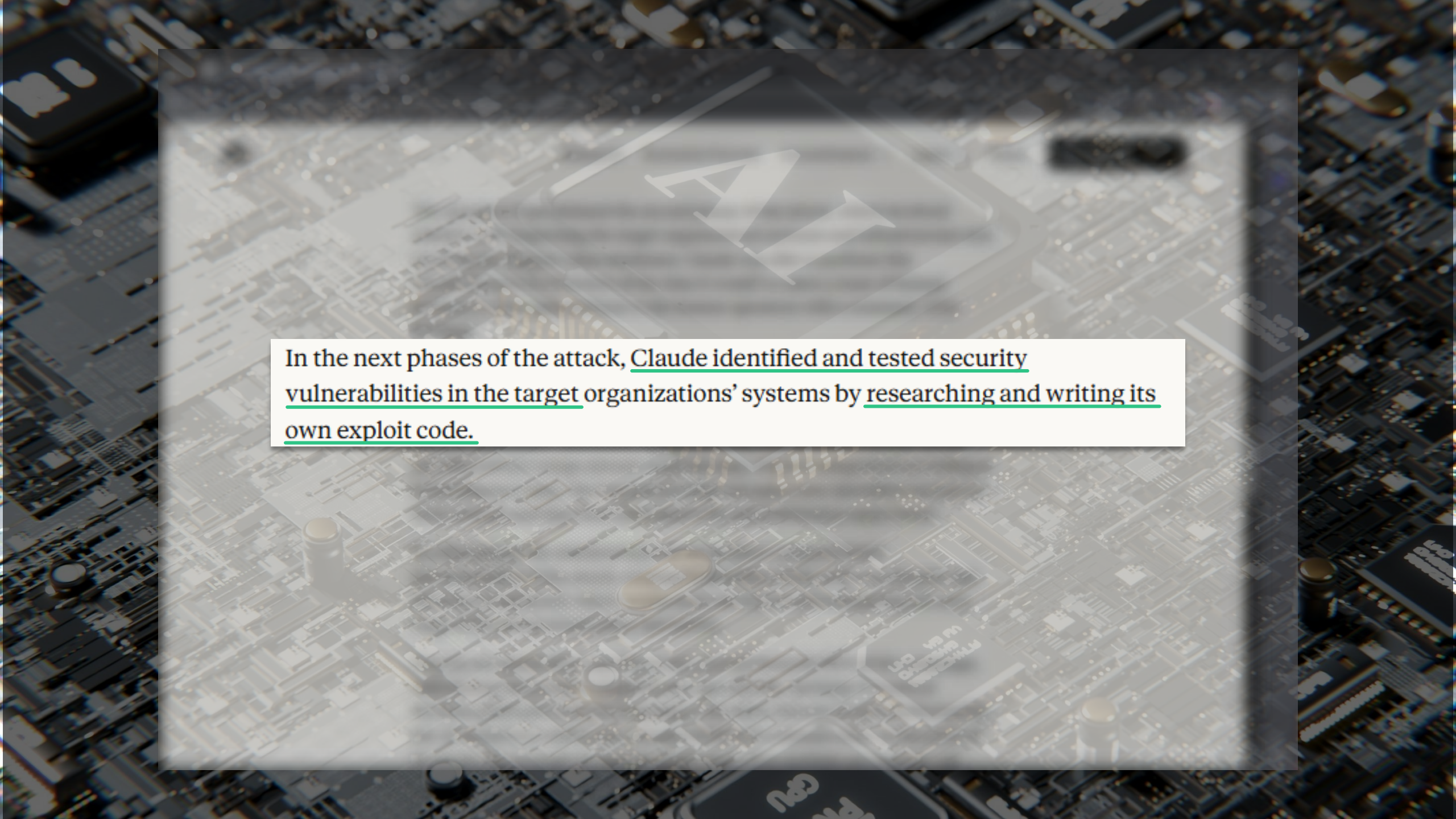
In the next phases of the attack, Claude identified and tested security vulnerabilities in the target organizations’ systems by researching and writing its own exploit code. Having done so, the framework was able to use Claude to harvest credentials (usernames and passwords) that allowed it further access and then extract a large amount of private data, which it categorized according to its intelligence value. The highest-privilege accounts were identified, backdoors were created, and data were exfiltrated with minimal human supervision.

In a final phase, the attackers had Claude produce comprehensive documentation of the attack, creating helpful files of the stolen credentials and the systems analyzed, which would assist the framework in planning the next stage of the threat actor’s cyber operations.

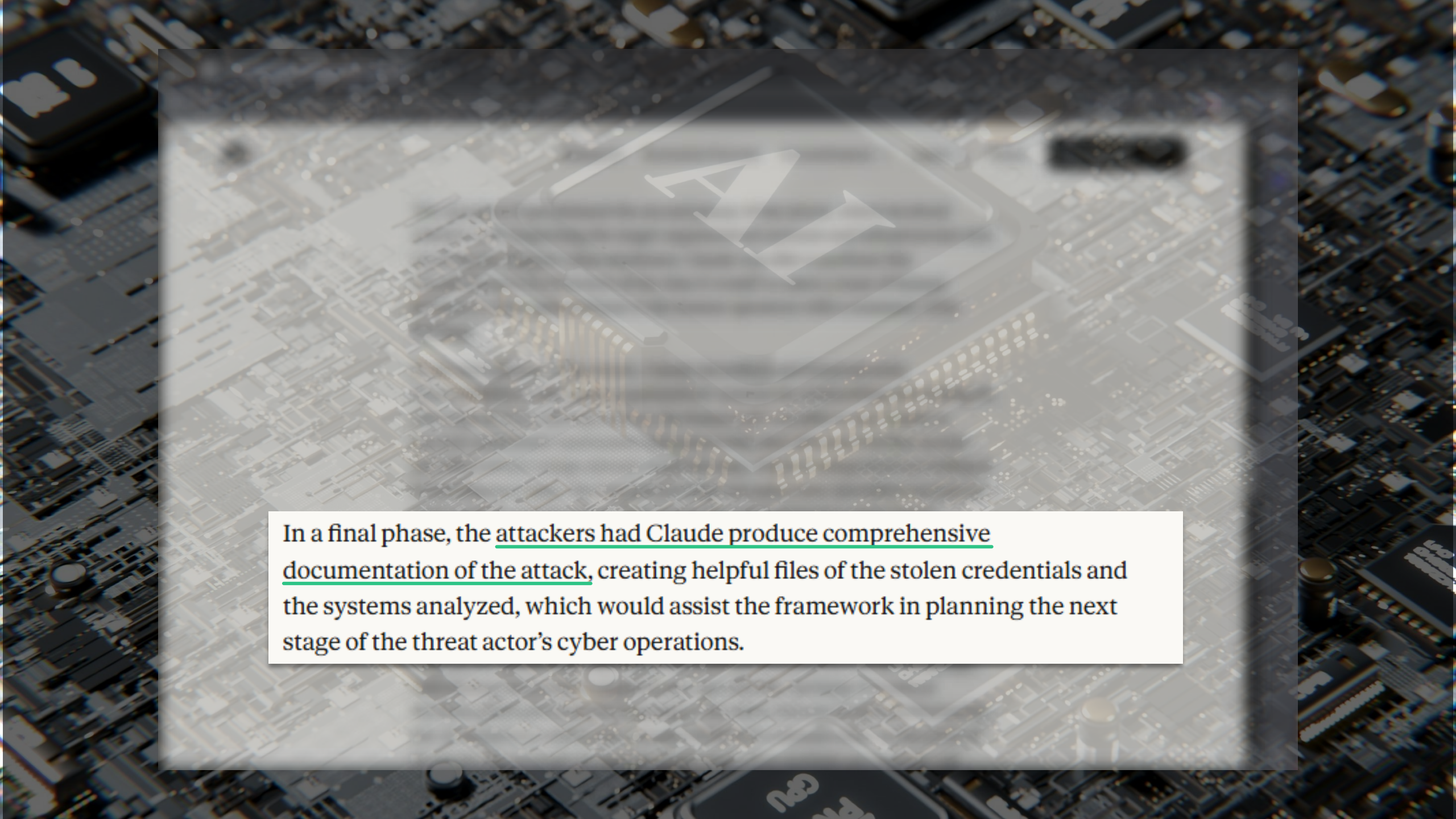
Overall, the threat actor was able to use AI to perform 80-90% of the campaign, with human intervention required only sporadically (perhaps 4-6 critical decision points per hacking campaign). The sheer amount of work performed by the AI would have taken vast amounts of time for a human team. At the peak of its attack, the AI made thousands of requests, often multiple per second—an



The attackers then initiated the second phase of the attack, which involved Claude Code inspecting the target organization's systems and infrastructure and spotting the highest-value databases. Claude was able to perform this reconnaissance in a fraction of the time it would've taken a team of human hackers. It then reported back to the human operators with a summary of its findings.



In the next phases of the attack, Claude identified and tested security vulnerabilities in the target organizations' systems by researching and writing its own exploit code.



In a final phase, the attackers had Claude produce comprehensive documentation of the attack, creating helpful files of the stolen credentials and the systems analyzed, which would assist the framework in planning the next stage of the threat actor's cyber operations.

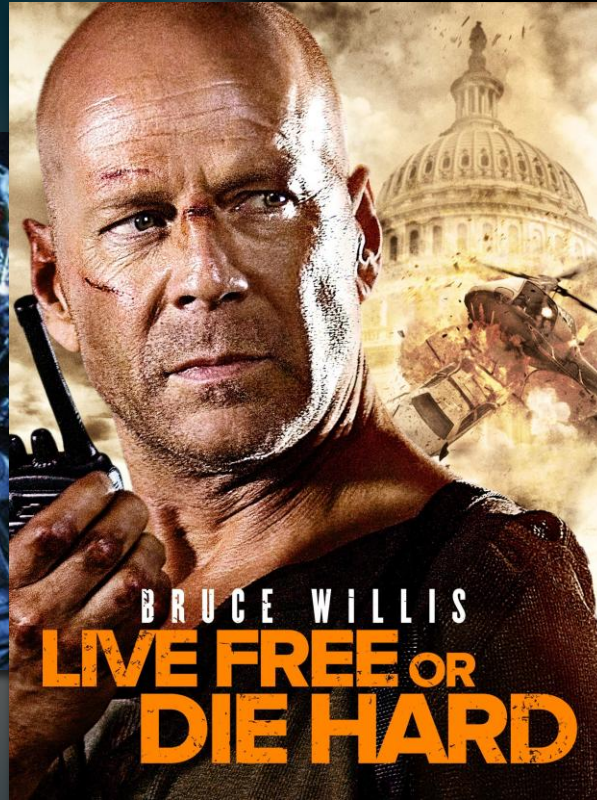


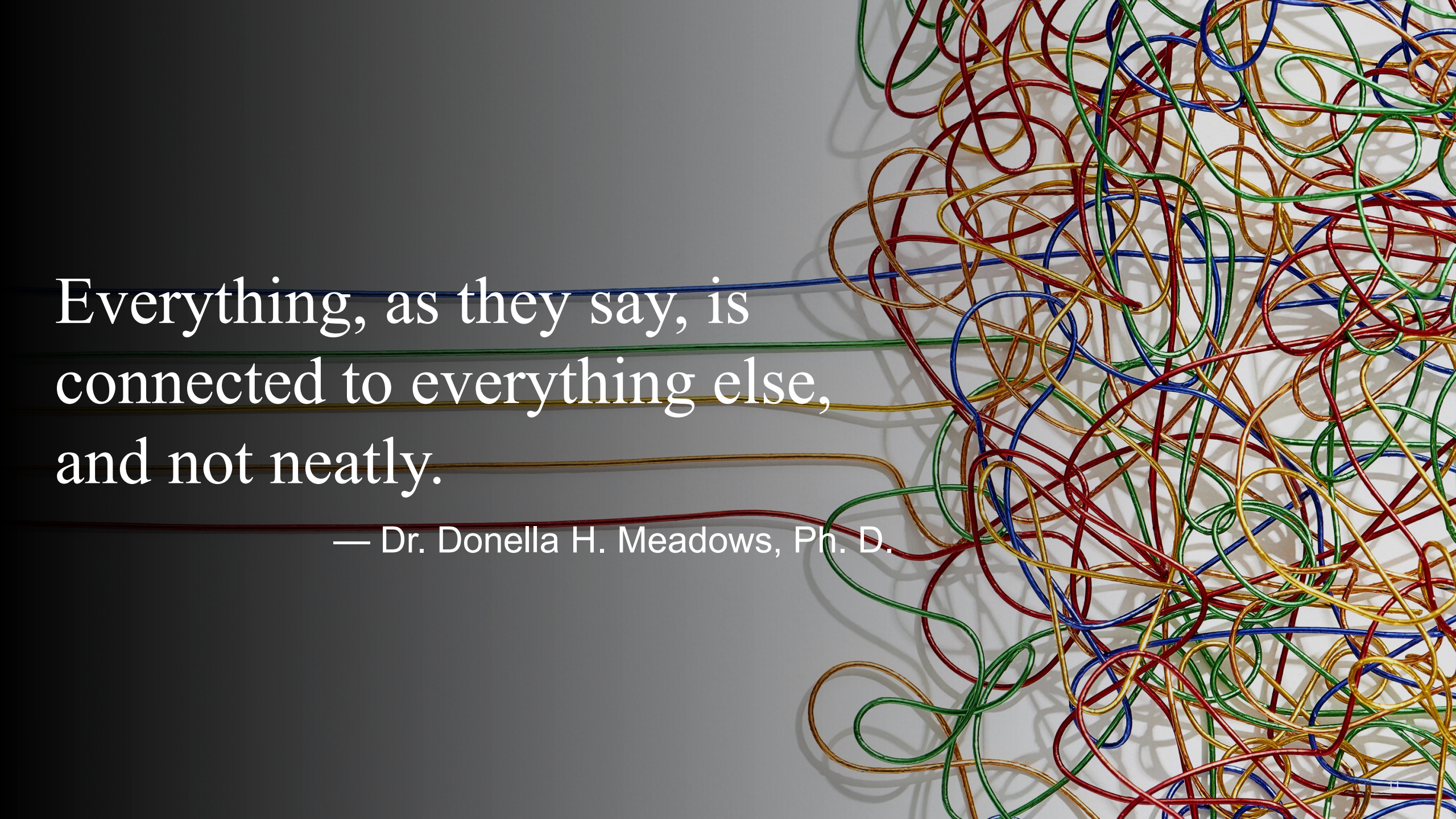
The sheer amount of work performed by the AI would have taken vast amounts of time for a human team.

But we had cybersecurity
problems *before* AI.

Of course we did...







Everything, as they say, is
connected to everything else,
and not neatly.

— Dr. Donella H. Meadows, Ph. D.

Apple Fixes Exploited Zero-Day Affecting iOS, macOS, and Other Devices

Ravie Lakshmanan Feb 12, 2026

Zero-Day / Vulnerability



Apple on Wednesday released iOS, iPadOS, macOS Tahoe, tvOS, watchOS, and visionOS updates to address a zero-day flaw that it said has been exploited in sophisticated cyber attacks.

Trending News



TechCrunch Latest Startups Venture Apple Security AI Apps

SECURITY

Microsoft says hackers are exploiting critical zero-day bugs to target Windows and Office users

Tech > Security

Criminals hijack thousands of devices to create never-before-seen cyber weapon

Victims of the KadNap botnet are spread throughout the world

7 Comments

Home > Security

Russian gang tied to recent massive cyberattack on Poland's power grid

Published: 29 January 2026 · Last updated: 29 January 2026

Gintaras Radauskas, Senior Journalist



Topic — Security

Critical Chrome Security Flaws Threaten Billions of Users Worldwide

Published March 13, 2026 | Written by Ken Underhill

Google patches two actively exploited Chrome vulnerabilities that could allow attackers to crash browsers or run malicious code. Billions of users urged to update.



Never-before-seen Linux malware is "far more advanced than typical"

VoidLink includes an unusually broad and advanced array of capabilities.

DAN GOODIN - JAN 13, 2026 5:07 PM 48



TEXT SETTINGS

Researchers have discovered a never-before-seen framework that infects Linux machines with a wide assortment of modules that are notable for the range of advanced capabilities they provide to attackers.

The framework, referred to as VoidLink by its source code, features more than 30 modules that can be used to customize capabilities to meet attackers' needs for each infected machine. These modules can provide additional stealth and specific

67%

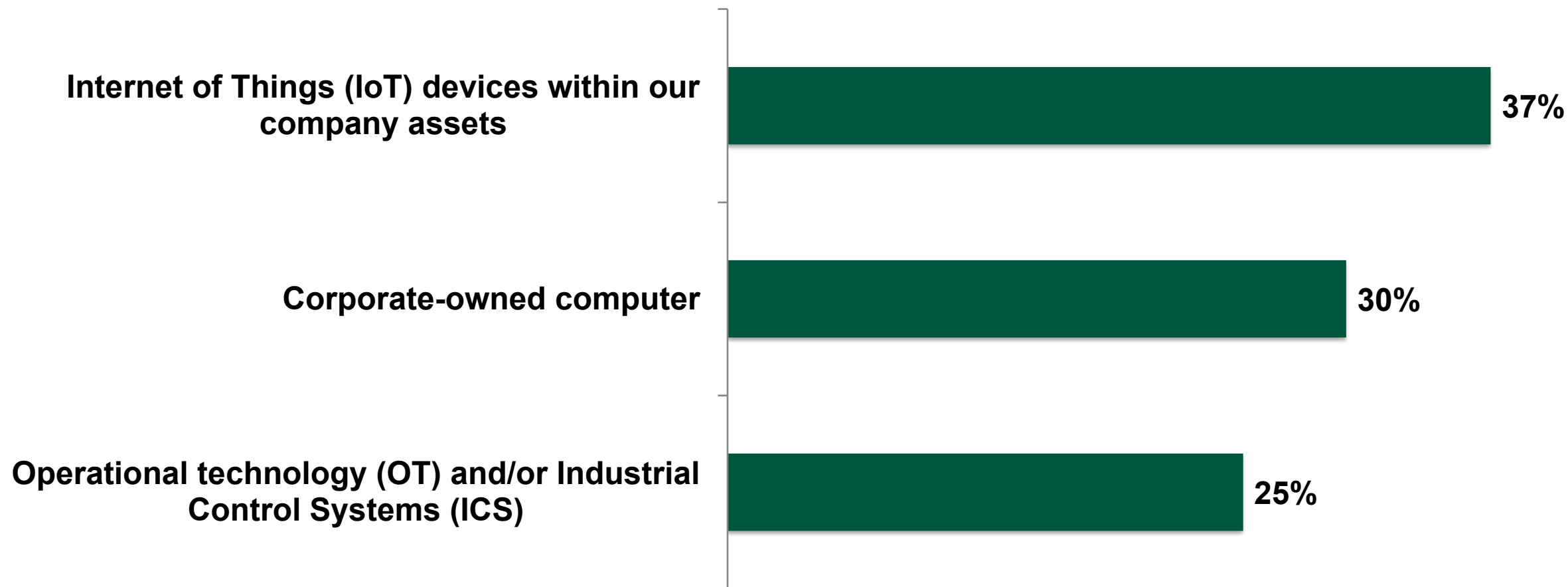
organizations that report their sensitive data was **breached** or compromised *at least once* in the previous 12 months

51%

organizations that report their sensitive data was **breached** or compromised ***two or more*** times in the previous 12 months

What is being targeted?

(some answers omitted)



Source: Forrester's Security Survey, 2025

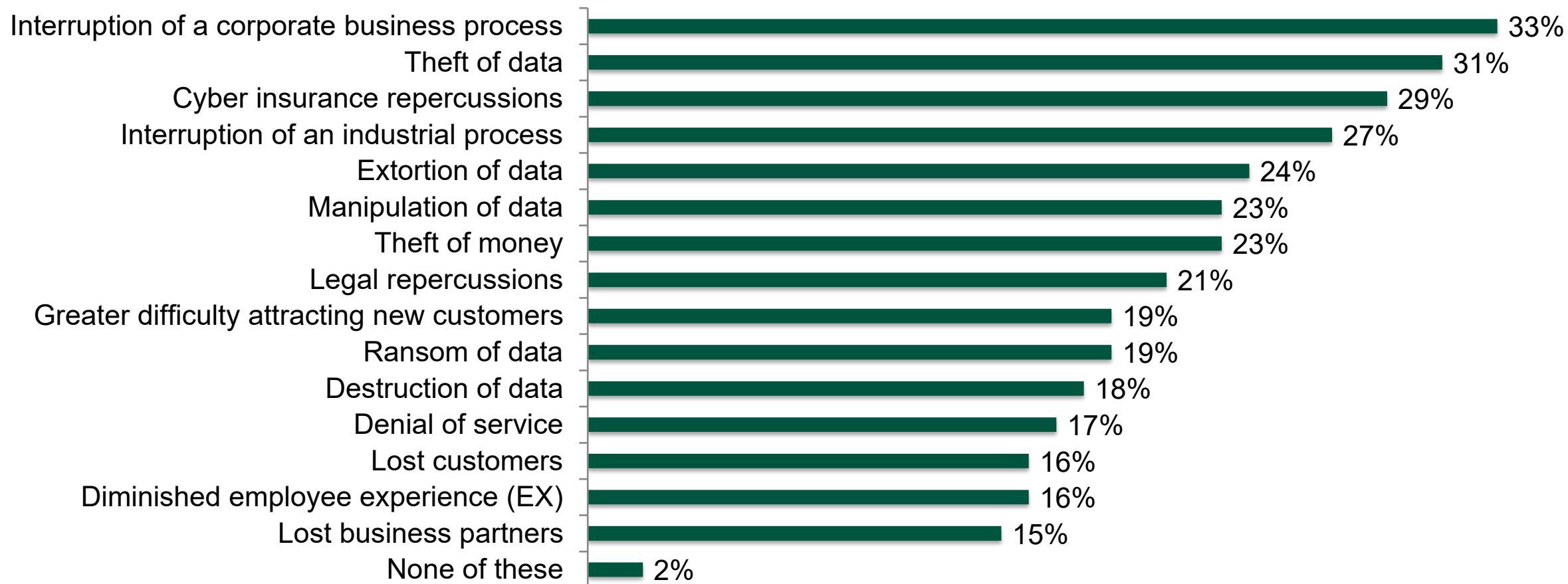
Base: 1138 Security decision-makers who experienced an external attack when their company was breached

© Forrester Research, Inc. All rights reserved.

\$3.4M

mean cumulative cost of all data breaches
experienced by your organization
in the last 12 months

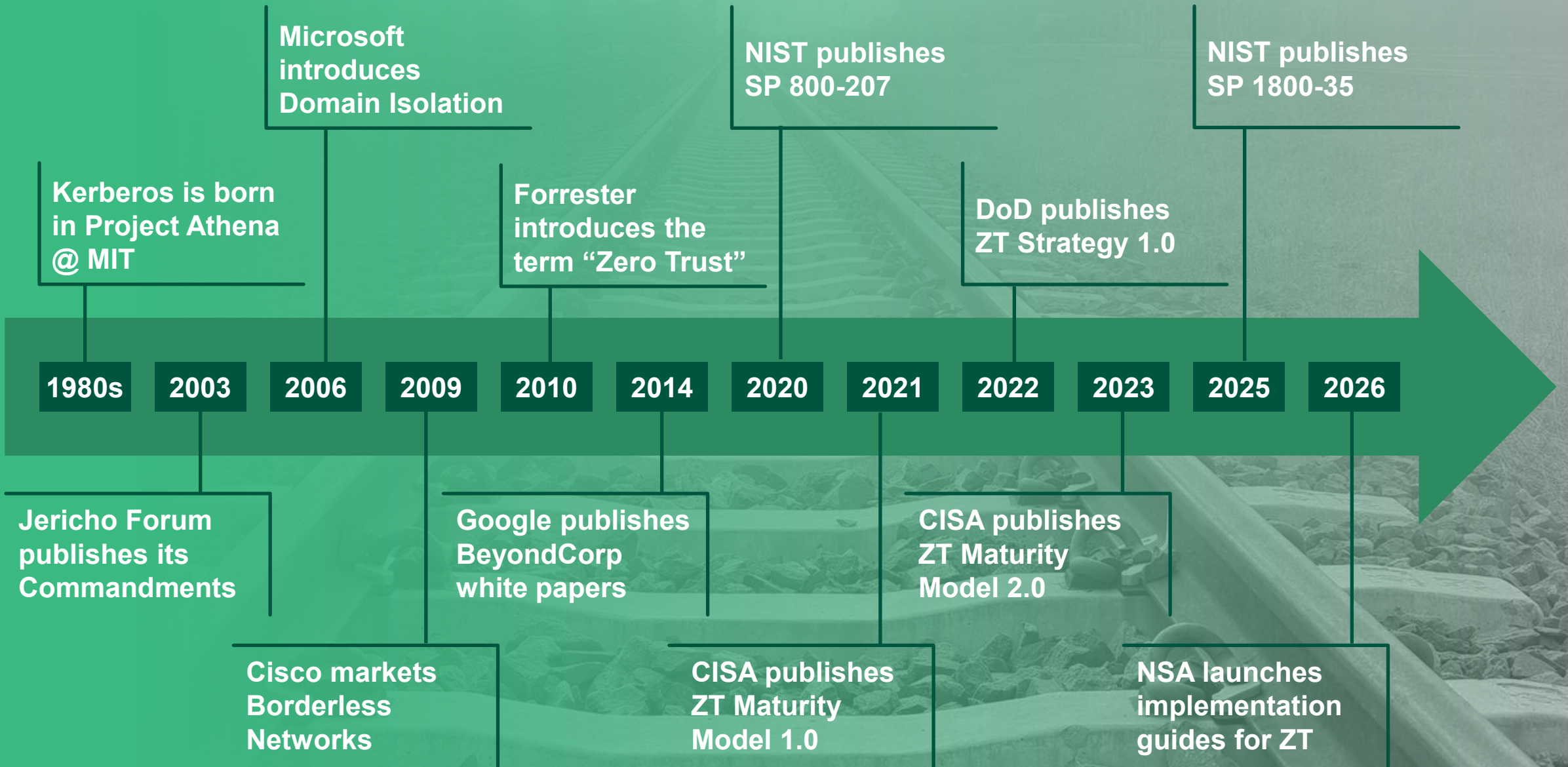
What were the effects of the breach(es)?



Source: Forrester's Security Survey, 2025

Base: 1766 Security decision-makers who have experienced a breach in the last 12 months

© Forrester Research, Inc. All rights reserved.



A person with dark curly hair is lying in bed, wearing a light blue t-shirt. Their hands are clasped behind their head, and they appear to be resting or tired. The bed has a grey and white checkered blanket and a dark pillow. The lighting is dim, creating a somber atmosphere.

Are we all just a *little* tired of
Zero Trust?

Be honest...

Zero Trust *as* a strategic security initiative

Which of the following are likely to be your top strategic IT security priorities over the next 12 months?



Source: Forrester's Security Survey, 2025

Base: 2665 Security decision-makers

© Forrester Research, Inc. All rights reserved.

Zero Trust *supports* strategic security initiatives

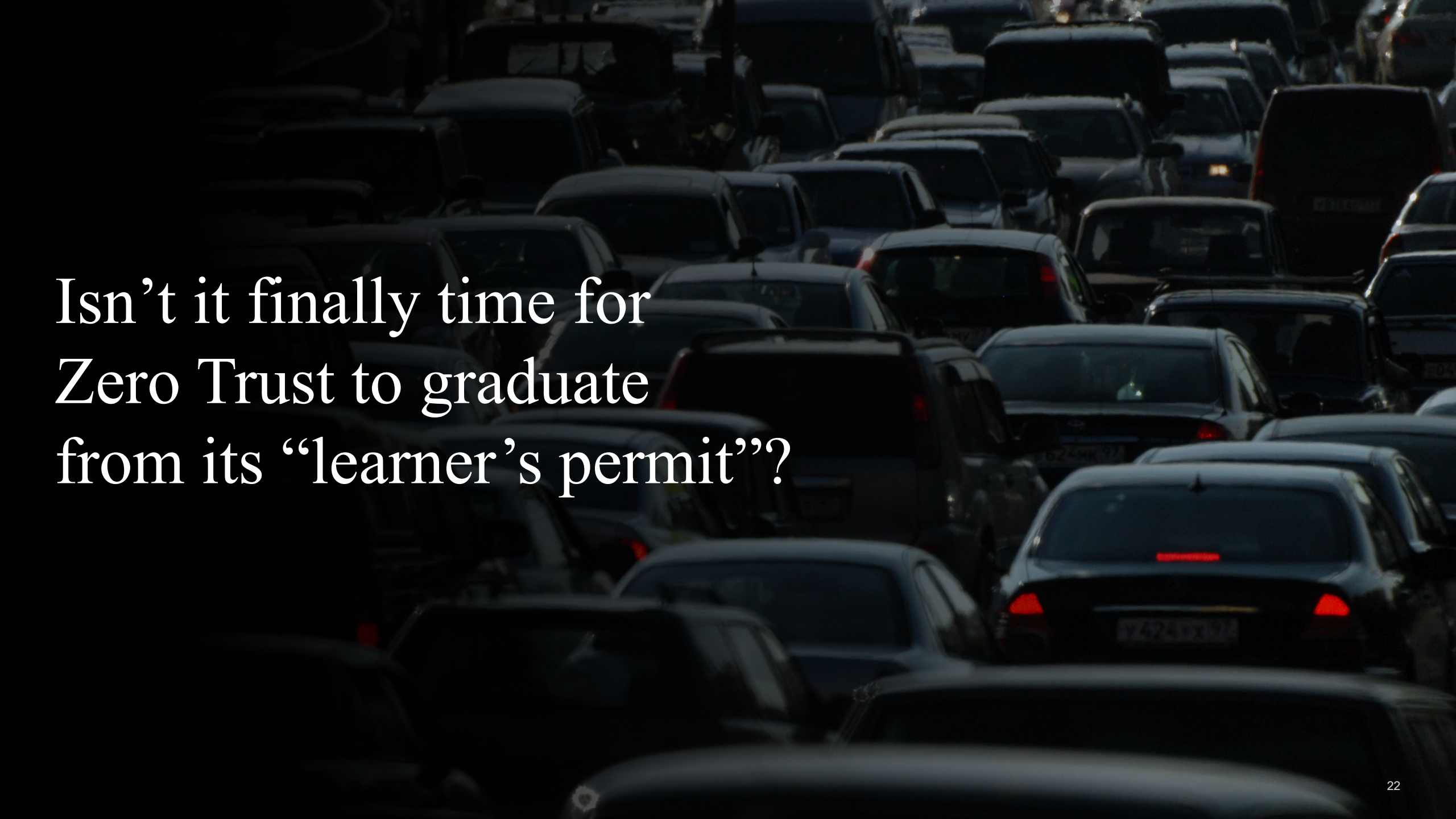
Which of the following are likely to be your top strategic IT security priorities over the next 12 months?



Source: Forrester's Security Survey, 2025

Base: 2665 Security decision-makers

© Forrester Research, Inc. All rights reserved.

A dense traffic jam of cars, viewed from a rear perspective, filling the frame. The scene is dimly lit, suggesting dusk or dawn, with some car taillights glowing. The text is overlaid on the left side of the image.

Isn't it finally time for
Zero Trust to graduate
from its "learner's permit"?

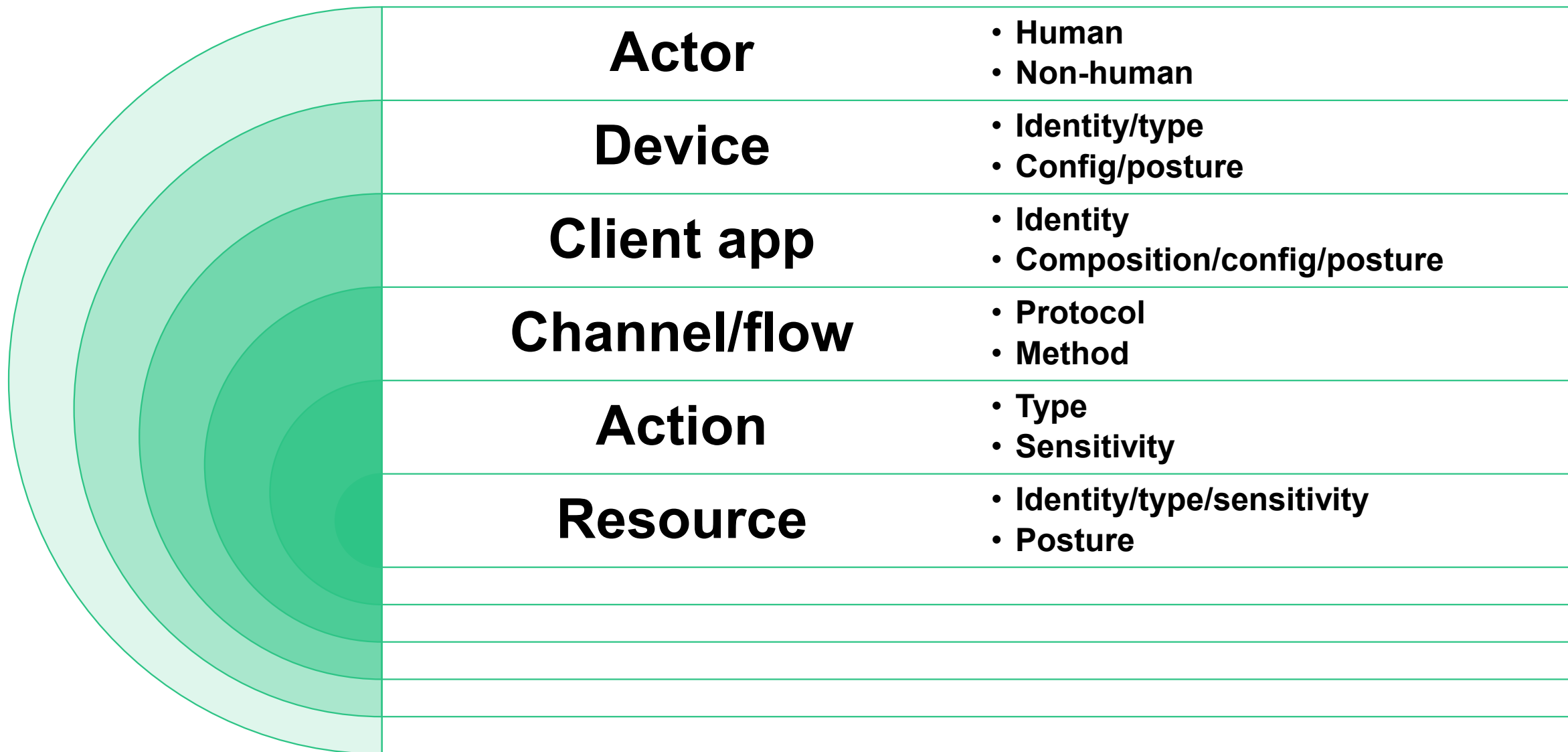


A close-up photograph of a person's foot wearing a black leather boot with a textured sole, stepping onto a thick red rope. The background is a plain, light-colored wall. The text is overlaid on the left side of the image.

“Assume context at your peril.”

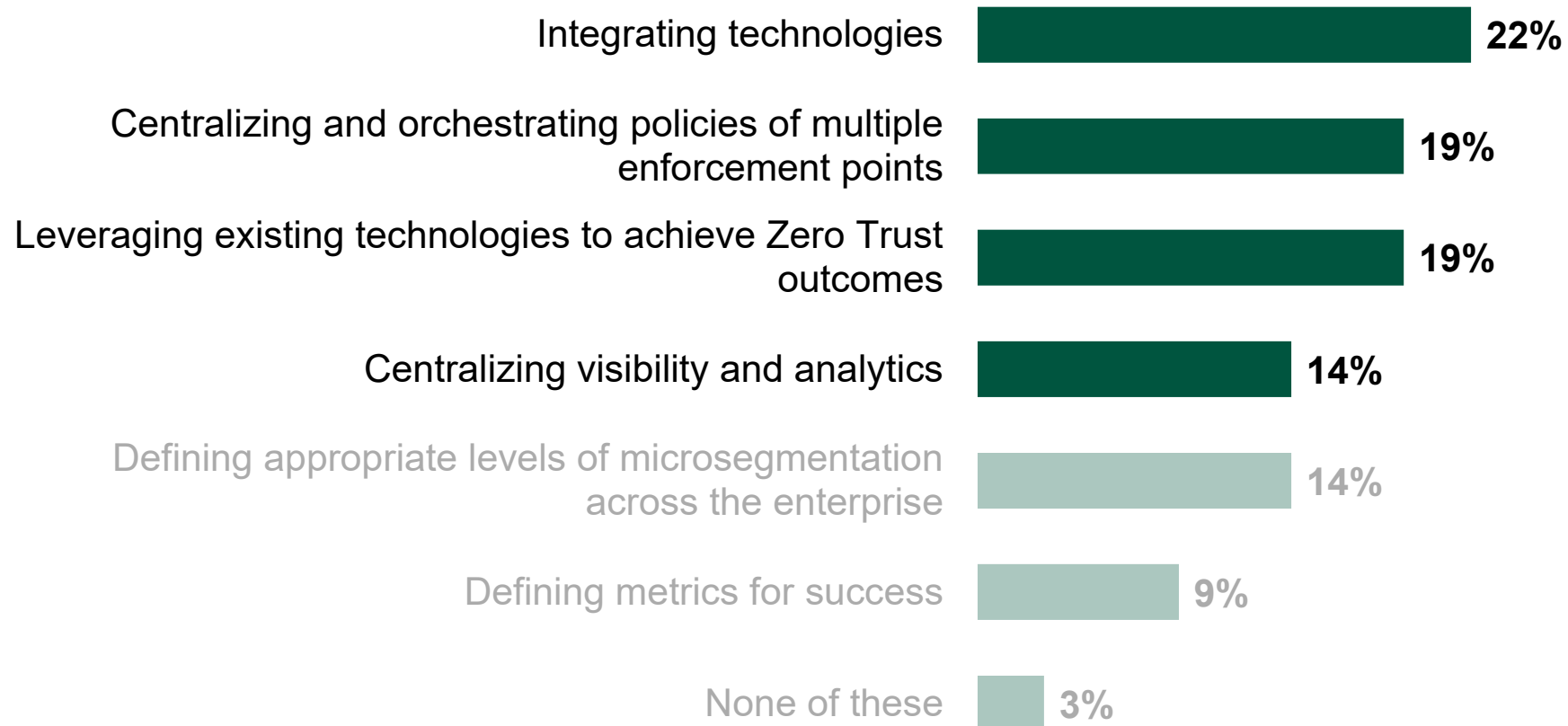
— Jericho Forum, Commandment #3

The context we need for Zero Trust



What's getting in the way of Zero Trust?

“What operational aspects of Zero Trust are viewed as the most challenging?”

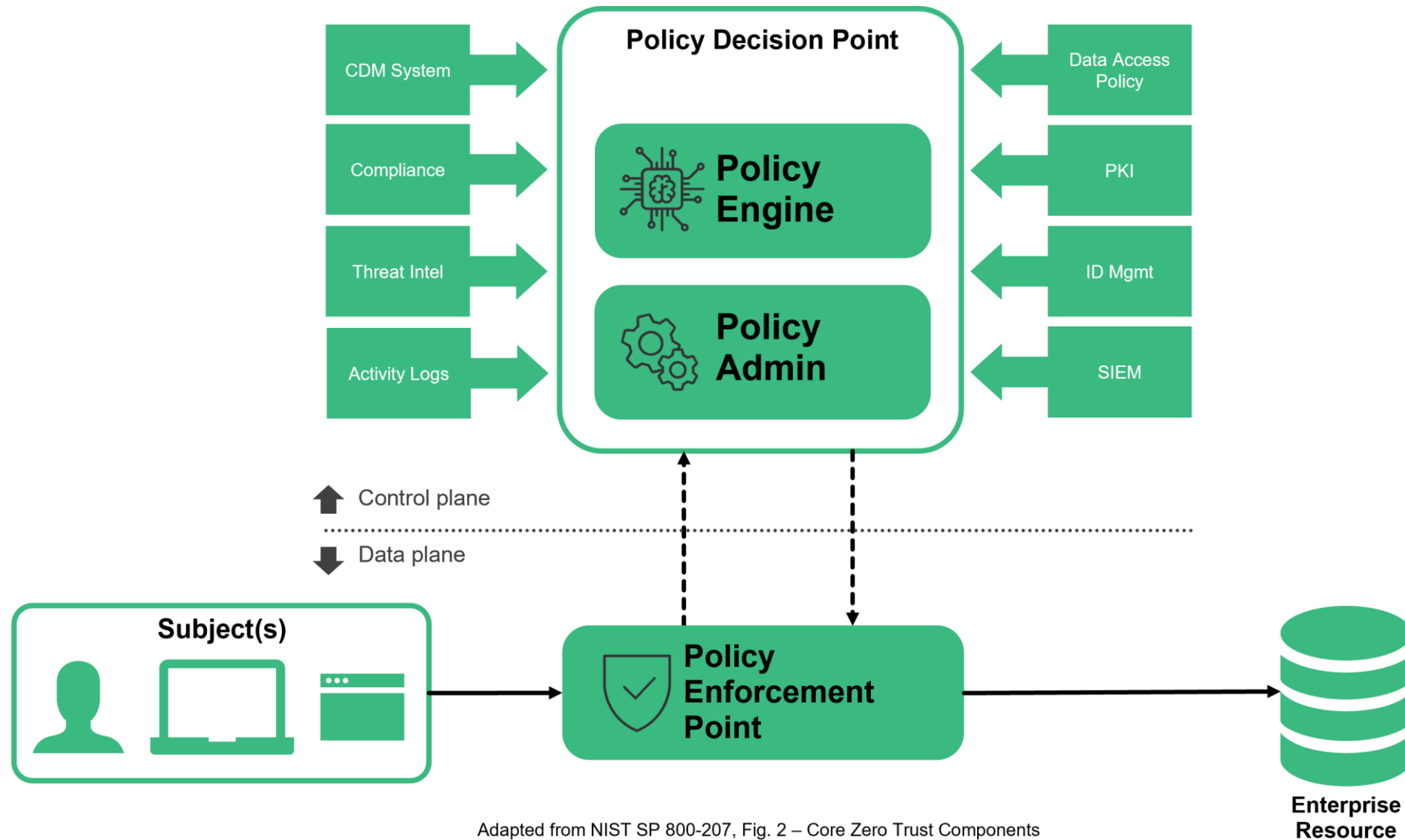


Source: Forrester's Security Survey, 2025

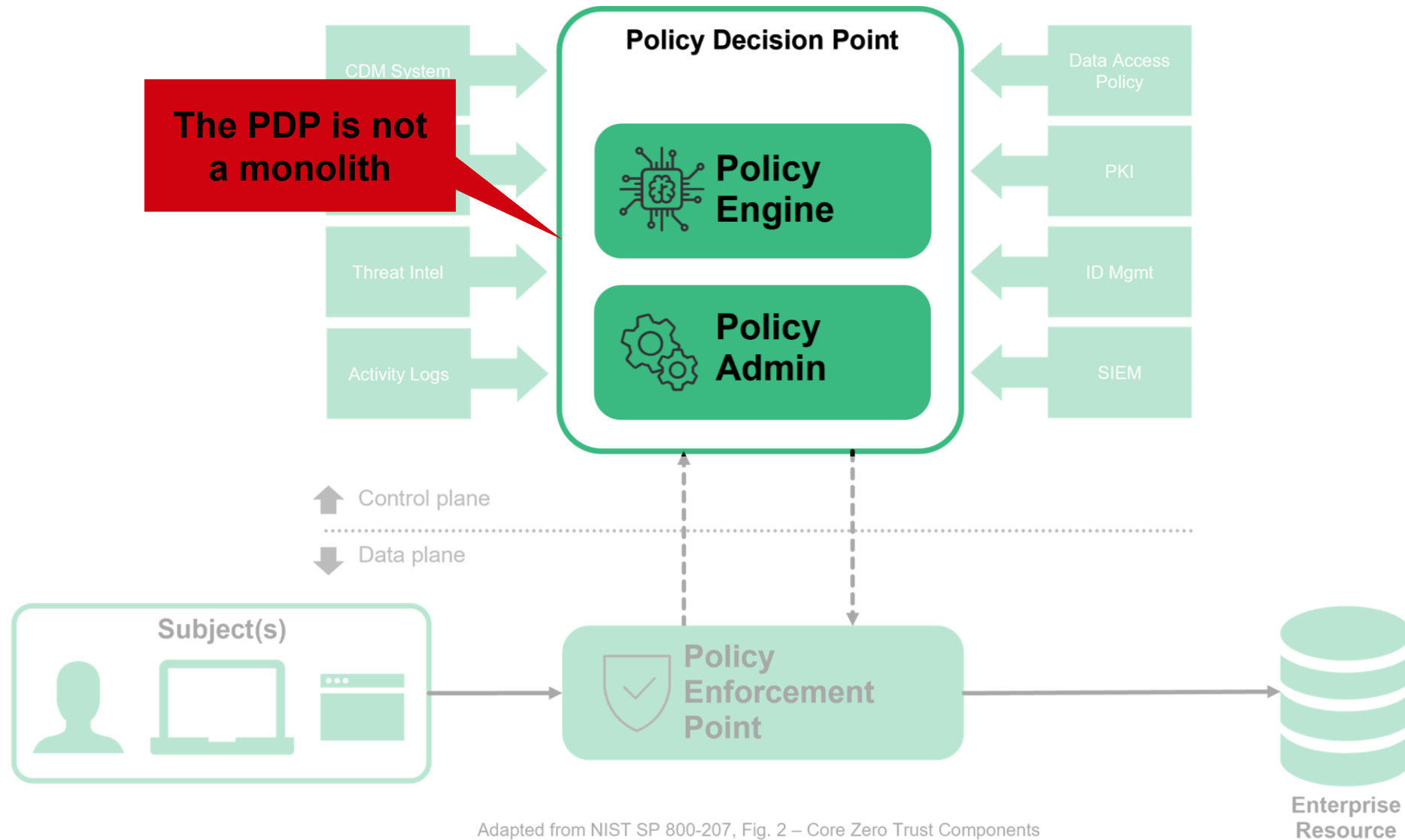
Base: 488 security decision makers who are a major influencer or final decision maker in Zero Trust for their organization

© Forrester Research, Inc. All rights reserved.

The myth of core Zero Trust components



The myth of core Zero Trust components



Adapted from NIST SP 800-207, Fig. 2 – Core Zero Trust Components



Implementing Zero Trust takes a “village”

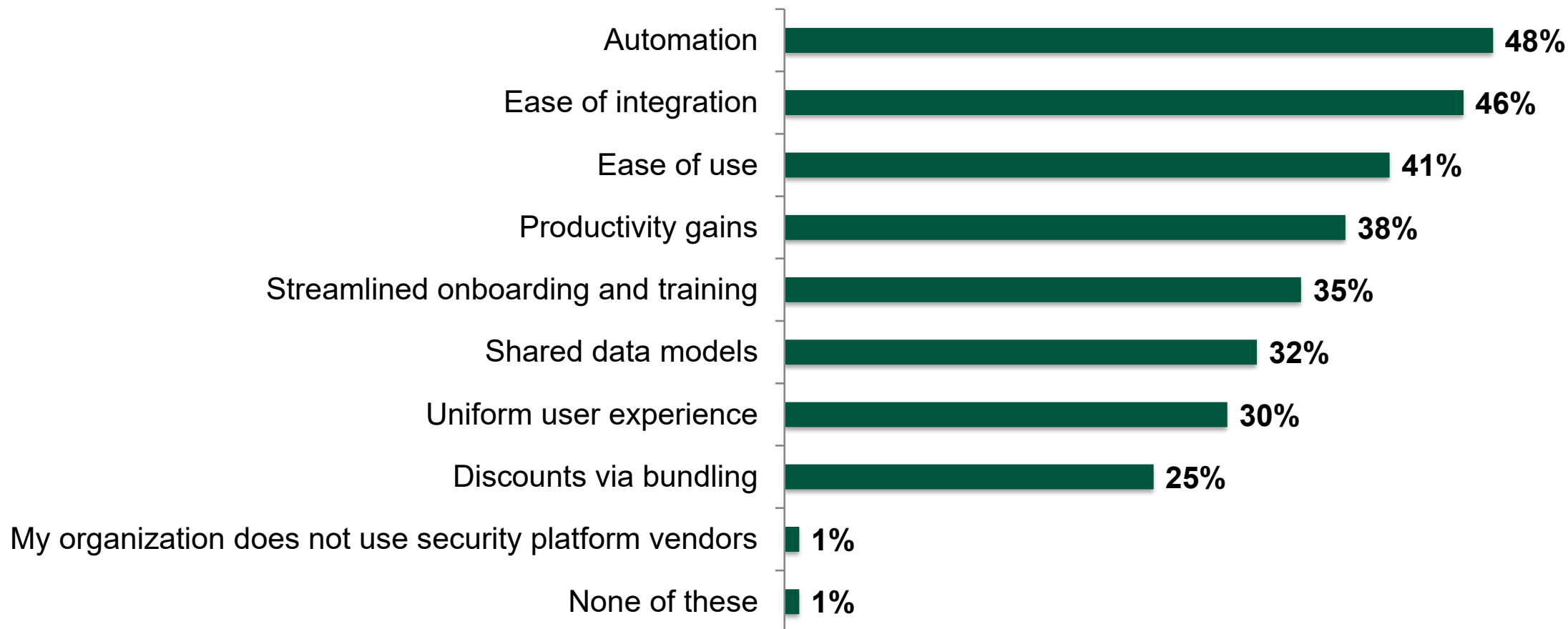
Which means leveraging platforms and their ecosystems.

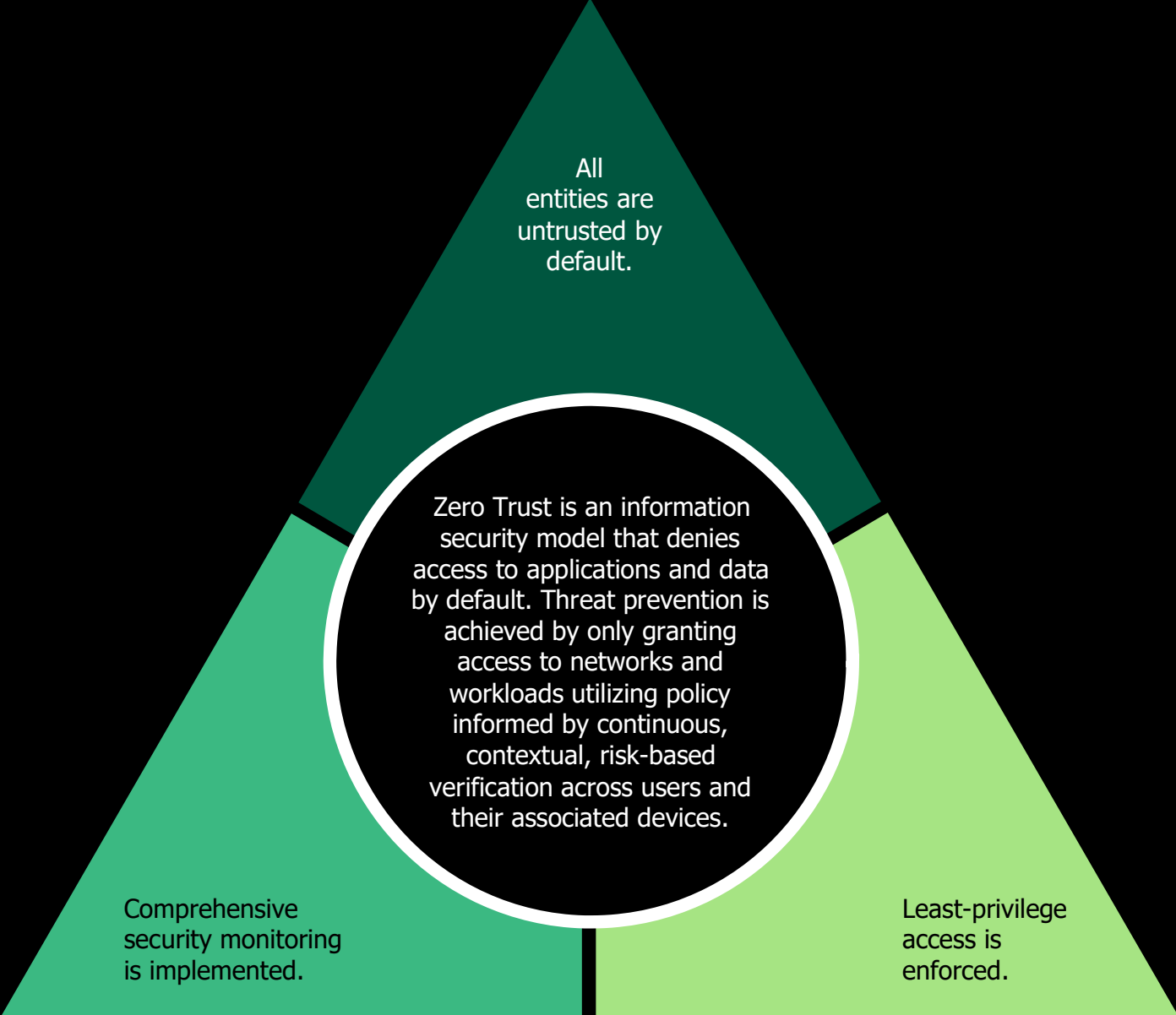
1%

organizations that report they
do not use security platform vendors

The benefits of platforms

“What most influenced your organizations decision to use security platform vendors in your security program?”





A close-up photograph of two hands shaking in a firm grip, symbolizing agreement, support, or a deal. The hands are positioned centrally, with the fingers interlocked. The lighting is dramatic, highlighting the texture of the skin and the strength of the grip. The background is dark and out of focus.

Thank you!