



# Anthropic KI Claude Mythos und das Ende kontrollierter Verfügbarkeit

Dirk Pappelbaum

Das KI-System selbst blieb dabei unverändert und funktionstüchtig. Das Problem lag nicht in der Technik der KI, sondern in der Umgebung, in der sie eingesetzt wird. Der Fall zeigt: Solche Systeme stehen selten allein. Sie sind Teil eines größeren Netzwerks aus Partnern, Schnittstellen und Testumgebungen. Genau an diesen Verbindungen entstehen oft Sicherheitslücken – und dort konnte auch in diesem Fall der Zugriff erfolgen.

## Verteilte Systeme erzeugen erwartbare Leckagen

Der Vorfall ist kein Ausreißer, er ist ganz klar Ausdruck der Systemlogik verteilter KI-Infrastrukturen. Jeder zusätzliche Zugriffspunkt erhöht nicht nur die Funktionalität, sondern auch die Wahrscheinlichkeit unautorisierter Nutzung. Ökonomisch entspricht das einem bekannten Muster: Kontrolle nimmt mit wachsender Komplexität ab. In einem System mit  $n$  Zugriffsknoten steigt die Zahl potenzieller Schwachstellen nicht linear, sondern überproportional, weil Abhängigkeiten zwischen den Knoten entstehen. Drittanbieter sind dabei kein Sonderfall, sondern integraler Bestandteil der Architektur. Der Zugriff ist deshalb kein singuläres Ereignis, sondern das

Resultat einer Struktur, die Offenheit benötigt, um überhaupt zu funktionieren.

## Project Glasswing und die Grenze institutionellen Vertrauens

Mit Project Glasswing versucht Anthropic, genau diese Offenheit zu kontrollieren. Der Zugang zu Mythos wird auf große Technologieunternehmen und sicherheitsrelevante Organisationen beschränkt. Die Steuerung erfolgt über Auswahl, Verträge und Vertrauen. Das Problem liegt in der Art dieser Kontrolle. Sie ist institutionell, nicht technisch. Zugang wird delegiert, nicht eingeschlossen. Damit verschiebt sich das Risiko entlang der Kette weiter – von Anthropic zu Partnern, von Partnern zu deren internen Strukturen. Der Guardian-Fall zeigt, dass diese Delegation nicht neutral ist. Sie erzeugt neue Angriffsflächen, ohne die alten zu eliminieren. Vertrauen ersetzt keine Zugriffssicherheit.

## Nutzung ist zweitrangig, Reproduzierbarkeit ist entscheidend

Die unbefugten Nutzer setzten Mythos laut Guardian nicht für Angriffe ein, sondern für Experimente. Das

begrenzt kurzfristig das Schadenspotenzial, ist aber für die ökonomische Bewertung nachrangig. Entscheidend ist, dass der Zugriff funktioniert hat – und damit prinzipiell wiederholbar ist. In dem Moment, in dem ein Zugangspfad existiert, sinken die Kosten für weitere Zugriffe drastisch. Die erste Nutzung ist ein Test, nicht das eigentliche Risiko. Sicherheit bemisst sich damit nicht an der aktuellen Nutzung, sondern an der Stabilität der Zugriffsschranken. Genau diese Schranken erweisen sich hier als durchlässig.

## Exklusivität verliert ihre ökonomische Funktion

Für Unternehmen verändert sich damit die Logik von Wettbewerb im Cyberbereich. Exklusiver Zugang zu leistungsfähiger Sicherheits-KI ist nur dann ein Vorteil, wenn er durchsetzbar bleibt. Wird der Zugang faktisch porös, verliert er seine strategische Wirkung. Gleichzeitig steigt der Anpassungsdruck. Systeme müssen so gebaut werden, dass sie auch unter der Annahme bestehen, dass vergleichbare Werkzeuge breiter verfügbar sind. Sicherheit verschiebt sich von Abschottung zu Resilienz.

Der Mythos-Vorfall markiert damit keinen Kontrollverlust im engeren Sinn. Er zeigt, dass Kontrolle in verteilten Systemen kein stabiler Zustand ist. Sie muss permanent gegen eine Struktur verteidigt werden, die auf Durchlässigkeit angelegt ist. Am Ende steht eine einfache, aber folgenreiche Einsicht: Wer Hochleistungs-KI in vernetzte Umgebungen integriert, organisiert nicht nur Zugriff – er organisiert auch dessen Umgehung.

Versicherungs- und Finanznachrichten

# expertenReport



<https://www.experten.de/id/4949495/Anthropic-Mythos-und-das-Ende-kontrollierter-Verfuegbarkeit/>