

Marcel Küppers

# Cyber- sicherheit & KI

Strategien und Praxis für  
IT- und Security-Verantwortliche



# Inhaltsverzeichnis

1	<b>Das neue Paradigma: KI trifft Cybersicherheit</b> . . . . .	15
1.1	Was wir mit KI in der Cybersicherheit wirklich meinen . . . . .	17
1.1.1	Künstliche Intelligenz als Oberbegriff . . . . .	17
1.1.2	Machine Learning vs. Regeln – der fundamentale Unterschied. . . . .	19
1.1.3	LLMs und GenAI – warum sie Security besonders verändern . . . . .	20
1.2	Von regelbasierten zu lernenden Security-Systemen . . . . .	22
1.2.1	Die Ära der Signaturen und starren Regeln. . . . .	22
1.2.2	Übergang zu ML-basierten Detektionen . . . . .	23
1.2.3	Der Sprung mit GenAI und LLMs . . . . .	24
1.3	Drei Rollen von KI in der Cybersicherheit. . . . .	25
1.3.1	KI als Verstärker der Verteidiger. . . . .	25
1.3.2	KI als Werkzeug der Angreifer . . . . .	26
1.3.3	KI als »unberechenbarer Dritter« . . . . .	27
1.4	Chancen für Verteidiger: Wo KI wirklich hilft . . . . .	28
1.4.1	Entlastung im SOC. . . . .	28
1.4.2	Bessere Nutzung vorhandener Daten. . . . .	29
1.4.3	Beschleunigung von DevSecOps. . . . .	30
1.4.4	GRC & Compliance . . . . .	31
1.5	Neue Risiken und Angriffsflächen durch KI. . . . .	32
1.5.1	Angriffe auf und über LLMs . . . . .	32
1.5.2	Governance- und Compliance-Risiken . . . . .	33
1.5.3	Organisatorische Risiken. . . . .	34
1.6	Was sich für IT-Verantwortliche und CISOs konkret ändert. . . . .	35
1.6.1	Vom Regel-Admin zum Risiko-Architekten. . . . .	35
1.6.2	Skill-Shift im Security-Team . . . . .	36
1.6.3	Zusammenarbeit über Silos hinweg. . . . .	38
1.7	Grundprinzipien für verantwortungsvollen KI-Einsatz in der Security. . . . .	39
1.7.1	»Augment, don't replace«. . . . .	39
1.7.2	»No Grounding – No Answer«. . . . .	40
1.7.3	Transparenz & Auditierbarkeit . . . . .	41
1.7.4	Minimalismus bei Daten . . . . .	41
1.7.5	Safety-by-Design und Security-by-Design. . . . .	42

1.8	Typische Einstiegsfehler – und wie Sie sie vermeiden . . . . .	42
1.8.1	»Wir benutzen einfach Tool X, das hat schon KI drin« . . . . .	42
1.8.2	Blindes Vertrauen in KI-Ergebnisse . . . . .	43
1.8.3	Undokumentierter Einsatz von externen KI-Diensten (»Shadow AI«). . . . .	43
1.8.4	Kein Lifecycle-Management. . . . .	44
1.9	Ausblick: Wohin die Reise in diesem Buch geht. . . . .	45
1.10	Referenzen . . . . .	46
<b>2</b>	<b>Bedrohungslandschaft 2026: Angreifer nutzen KI . . . . .</b>	<b>49</b>
2.1	Wer sind die Angreifer? – Akteurslandschaft mit KI . . . . .	52
2.1.1	Klassische Cybercrime-Gruppen . . . . .	52
2.1.2	Staatliche und staatlich unterstützte Gruppen (APT) . . . . .	54
2.1.3	»Cybercrime-as-a-Service 2.0« . . . . .	56
2.2	Der KI-Werkzeugkasten der Angreifer. . . . .	57
2.2.1	Text: Phishing, Social Engineering und Betrug . . . . .	57
2.2.2	Code: Exploits, Malware und Evasion . . . . .	59
2.2.3	Medien: Deepfakes und synthetische Identitäten . . . . .	60
2.2.4	Agenten: KI-»Bots«, die Kampagnen orchestrieren . . . . .	61
2.3	Phishing und Spear-Phishing 2.0. . . . .	63
2.4	Social Media & Messaging als primäre Angriffsflächen. . . . .	65
2.5	Desinformation & Informationsoperationen. . . . .	66
2.6	Deepfakes & KI-gestützte Erpressung . . . . .	68
2.7	Typische Angreifer-»Playbooks« mit KI. . . . .	70
2.7.1	Playbook 1: KI-gestützter BEC / CEO-Fraud. . . . .	71
2.7.2	Playbook 2: Ransomware mit KI-Augmentation . . . . .	72
2.8	Was bedeutet diese Bedrohungslage für Verteidiger? . . . . .	73
2.9	Fazit: Die KI-getriebene Bedrohungslandschaft als neues Normal . . . . .	76
2.10	Referenzen . . . . .	77
<b>3</b>	<b>KI-Grundlagen für IT- und Security-Entscheider . . . . .</b>	<b>81</b>
3.1	Was »KI« in der Praxis wirklich bedeutet . . . . .	81
3.2	Klassisches Machine Learning . . . . .	82
3.3	Deep Learning: Mustererkennung auf großer Skala. . . . .	83
3.4	Generative KI und LLMs: Token, Kontextfenster, Embeddings, Tool/Function Calling. . . . .	84
3.4.1	Token: Die »Währung« von LLMs. . . . .	85
3.4.2	Kontextfenster: »Wie viel kann das Modell gleichzeitig sehen?« . . . . .	86
3.4.3	Embeddings: Semantik als Vektoren (Basis für Suche und RAG) . . . . .	87
3.4.4	Tool/Function Calling: Vom Reden zum Handeln . . . . .	89

3.4.5	Zusammenfassung: Was Security-Entscheider daraus ableiten sollten	91
3.5	RAG (Retrieval Augmented Generation): »Chat mit eigenen Daten« richtig verstanden	91
3.6	Daten als Engpass: Qualität, Labels, Drift, Telemetrie	97
3.6.1	Datenqualität: »Garbage in, garbage out« – aber konkret	97
3.6.2	Labels: Warum »gelabelte Daten« in Security so schwer sind	99
3.6.3	Drift: Wenn Normalität sich ändert	100
3.6.4	Telemetrie: Die KI sieht nur, was Sie messen	101
3.7	Fazit: Daten entscheiden – und KI macht Observability zur Pflicht	102
3.8	Referenzen	103
<b>4</b>	<b>Architektur moderner KI-Sicherheitsplattformen</b>	<b>105</b>
4.1	Was unter »KI-Sicherheitsplattform« tatsächlich zu verstehen ist	106
4.2	Architekturprinzipien: Was »modern« im KI-Security-Kontext bedeutet	107
4.2.1	Evidence-first statt »KI sagt«	107
4.2.2	Zero Trust für Daten, Modelle und Tools	108
4.3	Die Referenzarchitektur: Sieben Schichten einer KI-Sicherheitsplattform	109
4.3.1	Schicht 1: Datenquellen (Signals & Knowledge)	109
4.3.2	Schicht 2: Ingestion und Normalisierung	110
4.3.3	Schicht 3: Speicherung und Indizes (Hot/Cold plus Vektor)	110
4.3.4	Schicht 4: Intelligence Layer (Modelle und Reasoning)	111
4.3.5	Schicht 5: Guardrails und Policy Enforcement	112
4.3.6	Schicht 6: Orchestrierung und Workflows	112
4.3.7	Schicht 7: Observability, Audit und Betriebsmodell	113
4.4	Referenzmuster für GenAI: LLM-only, RAG, Tool/Function Calling und Agenten	113
4.5	Datenfluss-Design: Vom Event zur Entscheidung	116
4.6	Security Controls innerhalb der Plattform	118
4.7	Observability und LLMOps/MLOps: Betrieb ist der Engpass	121
4.7.1	Was gemessen werden muss	121
4.7.2	Versionierung und Change Management	123
4.7.3	Reproduzierbarkeit als Audit-Anforderung	123
4.8	Governance in der Architektur verankern – nicht als PDF daneben	124
4.8.1	Policy Engine als zentraler Baustein	124
4.8.2	Data Classification und Retention	125
4.9	Build vs. Buy: Architekturentscheidungen aus CISO-Sicht	126

4.9.1	Wann »Buy« sinnvoll ist . . . . .	126
4.9.2	Wann »Build/Extend« sinnvoll ist. . . . .	127
4.10	Referenzarchitekturen (Blueprints) . . . . .	128
4.10.1	Blueprint A: »RAG Knowledge Layer« (Security-tauglich) . . . . .	129
4.10.2	Blueprint B: »LLM + Tool Gateway« . . . . .	129
4.10.3	Blueprint C: »Hybrid Enterprise« . . . . .	130
4.11	Typische Architekturfehler – und wie sie vermieden werden . . . . .	131
4.11.1	Fehler 1: »Es wird nur ein Modell gebaut« – statt einer Plattform . . . . .	131
4.11.2	Fehler 2: »Logs werden in RAG gekippt« . . . . .	132
4.11.3	Fehler 3: Fehlende ACLs im Retrieval. . . . .	132
4.11.4	Fehler 4: Prompt statt Policy . . . . .	132
4.11.5	Fehler 5: Keine Evals und keine Regressionstests . . . . .	133
4.11.6	Fehler 6: Tool Calling ohne IAM und ohne Validation . . . . .	133
4.12	Fazit: Architektur als Sicherheitskontrollsystem. . . . .	134
4.13	Referenzen . . . . .	134
<b>5</b>	<b>Daten, Telemetrie und Wissensquellen als Fundament von KI-Security</b> . . . . .	<b>137</b>
5.1	Zwei Datenwelten: »Signals« vs. »Knowledge«. . . . .	139
5.1.1	Signals: Telemetrie, Events, Zustände . . . . .	140
5.1.2	Knowledge: Dokumente, Regeln, Erfahrungswissen. . . . .	140
5.1.3	Warum diese Trennung entscheidend ist. . . . .	141
5.2	Die Mindestanforderung für Security-KI: Korrelation . . . . .	142
5.2.1	Identität: »Wer/was hat gehandelt?«. . . . .	143
5.2.2	Asset: »Worauf hat die Aktion gewirkt?«. . . . .	143
5.2.3	Zeit: »Wann genau?«. . . . .	144
5.2.4	Aktion: »Was ist passiert?«. . . . .	144
5.3	Telemetrie-Fundament: Die »Must-have«-Signalquellen . . . . .	145
5.3.1	Identity- und Access-Telemetrie (Tier-0). . . . .	146
5.3.2	Endpoint-Telemetrie (EDR/XDR) . . . . .	146
5.3.3	Cloud Control Plane Logs . . . . .	147
5.3.4	Netzwerk / DNS / Proxy (je nach Architektur). . . . .	147
5.3.5	Case-/Ticket-Daten als »Ground Truth der Realität«. . . . .	147
5.4	Wissensquellen-Fundament: Was in den RAG-Index gehört (und was nicht). . . . .	148
5.4.1	Kuratierte Quellen, die sich bewährt haben . . . . .	148
5.4.2	Was typischerweise nicht direkt in RAG gehört. . . . .	149
5.5	Data Governance: Klassifizierung, Zugriff, Retention, Zweckbindung . . . . .	150
5.5.1	Datenklassifizierung und Index-Policy . . . . .	151
5.5.2	Zugriffskontrolle: »Retrieval respects source ACLs«. . . . .	151
5.5.3	Retention und Auditability. . . . .	152
5.5.4	Datenschutz: PII-Minimierung und klare Zwecke . . . . .	152

5.6	Datenaufbereitung: Normalisierung, Enrichment, Qualitätssignale . . . . .	154
5.6.1	Normalisierung: Ohne Canonical Schema kein Scale. . . . .	154
5.6.2	Enrichment: Der Multiplikator für Priorisierung . . . . .	155
5.6.3	Data Quality Monitoring: »Data Health« als eigener KPI-Stream . . . . .	156
5.6.4	Fazit. . . . .	156
5.7	Knowledge Engineering für RAG: Chunking, Metadaten, Versionierung . . . . .	157
5.7.1	Chunking: Struktur erhalten, nicht zerstören . . . . .	157
5.7.2	Metadaten: Für Retrieval und Governance unverzichtbar. . . . .	158
5.7.3	»Quellenpflicht« technisch absichern. . . . .	159
5.7.4	Fazit. . . . .	159
5.8	Betriebsmodell: Ownership, Produktdenken, »Data as a Product« . . . . .	160
5.8.1	Rollenmodell, das sich bewährt hat . . . . .	160
5.8.2	»Data as a Product«: Denken Sie Daten wie Produkte, nicht wie Nebenprodukte. . . . .	162
5.9	Fazit: KI ist ein Verstärker – Daten und Wissen sind der Hebel. . . . .	162
5.10	Referenzen. . . . .	163
<b>6</b>	<b>Generative KI (GenAI) im Security-Alltag. . . . .</b>	<b>165</b>
6.1	Grundbegriffe, die Security-Verantwortliche beherrschen müssen. . . . .	166
6.1.1	Tokens – die Währung des Systems. . . . .	167
6.1.2	Kontextfenster – Kapazität, nicht Qualität . . . . .	168
6.1.3	Embeddings – Semantik als Suchprimitive . . . . .	169
6.1.4	Tool/Function Calling – vom Antworten zum Handeln. . . . .	171
6.2	Die vier Grundmuster von GenAI in Security. . . . .	172
6.2.1	Muster A: LLM-only (Text- und Strukturassistentz) . . . . .	172
6.2.2	Muster B: RAG (Antworten auf Basis interner Quellen). . . . .	173
6.2.3	Muster C: Tool-augmented LLM (Echtzeit-Fakten über Queries). . . . .	174
6.2.4	Muster D: Agentische Orchestrierung (mehrschrittige Planung) . . . . .	176
6.3	Prompting als Engineering-Disziplin . . . . .	178
6.3.1	Prompt-Schichten: System, Developer/Policy, User . . . . .	178
6.3.2	Output-Formate erzwingen. . . . .	180
6.3.3	Token Budgeting und Kontextkurierung . . . . .	181
6.4	Guardrails: Von »Prompt-Regeln« zu echten Kontrollen. . . . .	182
6.4.1	Input-Guardrails. . . . .	182
6.4.2	Output-Guardrails . . . . .	182
6.4.3	Retrieval- und Tool-Guardrails . . . . .	183
6.4.4	Governance-Guardrails . . . . .	184

6.5	Betriebsmodelle: On-Prem, Cloud, Hybrid – aus Security-Sicht . . . .	186
6.5.1	Entscheidungskriterien . . . . .	186
6.5.2	Anforderungen bei On-Prem Umgebungen . . . . .	187
6.5.3	Hybridbetrieb . . . . .	188
6.5.4	Cloudbetrieb . . . . .	189
6.6	Risiko- und Bedrohungsmodell für GenAI im Security-Umfeld . . . .	190
6.6.1	Prompt Injection und Datenexfiltration . . . . .	190
6.6.2	Data Poisoning in Wissensquellen . . . . .	191
6.6.3	Tool Misuse und Privilege Escalation . . . . .	192
6.6.4	Model-Behavior-Risiken (Qualität, Übervertrauen, Kontextbias) . . . . .	193
6.7	Fazit: GenAI ist der Multiplikator – aber nur mit Systemdesign. . . .	194
6.8	Referenzen . . . . .	195
<b>7</b>	<b>KI in zentralen Sicherheitsdomänen: Praxis-Use-Cases . . . . .</b>	<b>197</b>
7.1	Referenzstruktur für alle Use Cases . . . . .	197
7.2	Use Case 1: Automatisierte SOC-Lageberichte . . . . .	199
7.2.1	Zielbild und Nutzen . . . . .	199
7.2.2	Muster: LLM-only und Tool-Augmentation . . . . .	201
7.2.3	Inputs . . . . .	202
7.2.4	Architektur . . . . .	213
7.2.5	Komplettes Beispiel: Rohdaten → LLM-Prompt → fertiger Tageslagebericht . . . . .	215
7.2.6	Betrieb & Messung (KPIs, Evals, SLOs) . . . . .	225
7.2.7	Anti-Patterns . . . . .	226
7.3	Use Case 2: GRC-Assistent als RAG-System (Policies, Controls, Evidence) . . . . .	227
7.3.1	Zielbild und Nutzen . . . . .	228
7.3.2	Muster . . . . .	230
7.3.3	Wissensbasis (Knowledge) – was indexiert wird . . . . .	231
7.3.4	Architektur . . . . .	243
7.3.5	Implementierung . . . . .	246
7.3.6	Betrieb & Messung . . . . .	250
7.3.7	Anti-Patterns . . . . .	251
7.4	Use Case 3: Vulnerability & Patch-Priorisierung (Contextual Risk) . . . . .	253
7.4.1	Zielbild und Nutzen . . . . .	253
7.4.2	Muster . . . . .	254
7.4.3	Inputs und Outputs . . . . .	254
7.4.4	Architektur . . . . .	255
7.5	Fazit . . . . .	257

<b>8</b>	<b>Governance, Compliance und Ethik</b> .....	259
8.1	Governance-Ziele und Leitprinzipien .....	261
8.1.1	Evidence-first und Nachvollziehbarkeit .....	261
8.1.2	Purpose Limitation und Datenminimierung .....	261
8.1.3	Least Privilege für Menschen, Modelle und Tools .....	262
8.1.4	Transparenz und Verantwortlichkeit .....	262
8.1.5	Fairness und Schadenminimierung .....	262
8.2	Governance-Operating-Model: Rollen, Gremien, Entscheidungsrechte .....	263
8.2.1	KI-Governance-Board (strategisch) .....	264
8.2.2	KI-Risk & Compliance Council (kontrollierend) .....	265
8.2.3	KI-Plattformbetrieb (operativ) .....	265
8.3	Policy-Architektur für KI-Security: Welche Richtlinien Sie wirklich brauchen .....	266
8.3.1	KI-Nutzungsrichtlinie (Acceptable Use) .....	267
8.3.2	Daten- und Index-Policy (RAG & Embeddings) .....	267
8.3.3	Tool-Access-Policy (Function Calling Governance) .....	268
8.3.4	Prompt- und Modell-Change-Policy (LLMOps/MLOps) ....	269
8.3.5	Output- und Kommunikationspolicy .....	269
8.4	Risiko-Management: Wie man KI-Risiken systematisch bewertet .....	270
8.4.1	Risikodimensionen .....	270
8.4.2	Risikoklassifizierung pro Use Case .....	271
8.4.3	Kontrollen (Mapping Risiko → Control) .....	273
8.5	Compliance: KI-Security in regulierten Umgebungen sicher betreiben .....	274
8.5.1	Auditierbarkeit als First-Class Requirement .....	274
8.5.2	Datenschutz und arbeitsrechtliche Dimensionen (PII, Monitoring) .....	275
8.5.3	Datenresidenz, Subprozessoren, Drittlandtransfer .....	276
8.5.4	Nachweise für Kontrollen (Evidence Packs) .....	276
8.6	Ethik in KI-Security: Was das konkret im Alltag bedeutet .....	277
8.6.1	Fairness und Bias: Wo Bias in Security praktisch entsteht .....	277
8.6.2	Automation Bias: »Wenn die KI es sagt, wird es stimmen« .....	279
8.6.3	Transparenz und Kennzeichnung .....	279
8.6.4	Ethical Red Lines (praktische rote Linien) .....	280
8.7	Third-Party & Vendor Governance: Wenn KI von außen kommt ...	280
8.7.1	Vendor Due Diligence – KI-spezifische Fragen .....	281
8.7.2	Vertrags- und Kontrollpunkte .....	282
8.8	Controls-by-Design: Governance als Architektur .....	282
8.8.1	RAG: Evidenzfähiges Retrieval .....	283
8.8.2	Tool Calling: Privilegierte Pfade kontrollieren .....	283

8.8.3	Outputs: Datenhygiene und kommunikative Kontrolle. . . . .	284
8.8.4	Betrieb als Kontrollfläche. . . . .	284
8.8.5	Audit: Reproduzierbarkeit als Kernfähigkeit . . . . .	284
8.9	Metriken & KPIs für Board und Management . . . . .	285
8.9.1	Wertmetriken (Outcome). . . . .	285
8.9.2	Risiko- und Kontrollmetriken . . . . .	286
8.9.3	Compliance-Metriken. . . . .	287
8.10	Incident Response für KI-Systeme: »Model Misbehavior« ist ein Incident . . . . .	287
8.10.1	Trigger für KI-Incidents. . . . .	288
8.10.2	Minimales Runbook. . . . .	288
8.11	Fazit: Governance ist der Multiplikator für sicheren Nutzen. . . . .	291
8.12	Referenzen . . . . .	291
<b>9</b>	<b>Organisation und Betrieb: Der CISO und IT-Leiter im KI-Zeitalter . . . . .</b>	<b>295</b>
9.1	Neue Rollenlogik: Vom Tool-Betrieb zum Produktbetrieb – »Security AI as a Product«. . . . .	296
9.2	Organisatorische Zielbilder – drei Referenzmodelle . . . . .	299
9.2.1	Modell A: »Federated Enablement« . . . . .	299
9.2.2	Modell B: »Security AI Platform Team« – dediziertes Plattformteam. . . . .	300
9.2.3	Modell C: »Enterprise AI Platform + Security Overlay« – bei großen Konzernen. . . . .	301
9.3	Notwendige Betriebsprozesse. . . . .	301
9.3.1	Change Management für Prompts, Modelle und Indizes (LLMOps/MLOps) . . . . .	302
9.3.2	Incident Response für KI-Systeme – Model Misbehavior . . . . .	303
9.3.3	Data Health. . . . .	303
9.3.4	Kosten- und Kapazitätsmanagement. . . . .	303
9.4	Menschen & Kultur: Adoption, Training und der Kampf gegen Automatisierungsbias . . . . .	305
9.5	Sicherheitsorganisation im KI-Zeitalter: Neue Fähigkeiten als Capability Map . . . . .	307
9.6	Verzahnung CISO ↔ IT-Leiter: Die neue »gemeinsame Verantwortung«. . . . .	310
9.6.1	Gemeinsame Architekturentscheidungen . . . . .	310
9.6.2	Gemeinsame SLOs und KPIs . . . . .	311
9.6.3	Gemeinsame Change- und Incident-Prozesse. . . . .	311
9.7	Investment-Strategie . . . . .	312
9.7.1	Budget-Buckets (praktisch) . . . . .	313
9.7.2	»Build vs Buy« für Betrieb . . . . .	314
9.8	Roadmap: 90-Tage-Plan für CISO & IT-Leitung. . . . .	315

9.8.1	Phase 1 (0–30 Tage): Fundament & Governance . . . . .	315
9.8.2	Phase 2 (31–60 Tage): Plattform-MVP – kontrollierbar, wiederverwendbar . . . . .	316
9.8.3	Phase 3 (61–90 Tage): Skalierung & Härtung . . . . .	317
9.9	10 Fragen, die CISO und IT-Leiter gemeinsam beantworten müssen . . . . .	318
9.10	Fazit . . . . .	320
9.11	Referenzen . . . . .	321
<b>10</b>	<b>Angreifbare KI: Security für KI-Systeme selbst . . . . .</b>	<b>323</b>
10.1	Warum KI-Systeme »anders« angreifbar sind . . . . .	324
10.2	Bedrohungsmodell: Was wir schützen und wogegen . . . . .	326
	10.2.1 Schutzgüter (Assets) . . . . .	326
	10.2.2 Angreiferprofile . . . . .	328
	10.2.3 Angriffsziele (typische Outcomes) . . . . .	329
10.3	Hauptangriffsklassen . . . . .	330
	10.3.1 Prompt Injection . . . . .	330
	10.3.2 RAG Data Exfiltration (ACL-Bypass) . . . . .	331
	10.3.3 Data Poisoning / Knowledge-Base-Manipulation . . . . .	332
	10.3.4 Tool Misuse / Privilege Escalation (Function Calling) . . . . .	332
	10.3.5 Sensitive Data Leakage (Prompts, Logs, Outputs) . . . . .	333
	10.3.6 Denial-of-Wallet / Cost DoS . . . . .	333
	10.3.7 Model/Prompt Supply Chain & Dependency Risk . . . . .	334
10.4	Referenzarchitektur für »Sichere KI-Systeme« . . . . .	335
	10.4.1 Sicherheitsprinzipien (Design-Level) . . . . .	335
	10.4.2 Kontrollpunkte . . . . .	336
10.5	Sicherheitsanforderungen . . . . .	339
	10.5.1 Identität & Zugriff . . . . .	339
	10.5.2 Daten- & RAG-Governance (Pflicht) . . . . .	340
	10.5.3 Tool-Governance . . . . .	341
	10.5.4 Output-Sicherheit (Pflicht) . . . . .	342
	10.5.5 Betrieb . . . . .	343
10.6	Security Testing: Wie man KI-Systeme realistisch testet . . . . .	344
	10.6.1 »Evals« sind die neuen Unit Tests . . . . .	345
	10.6.2 Red Teaming (systematisch, nicht ad hoc) . . . . .	346
	10.6.3 Secure Prompt Engineering . . . . .	347
10.7	Fazit: KI-Systeme sind »High-Trust Systems« . . . . .	347
10.8	Referenzen . . . . .	349
<b>A</b>	<b>Glossar zentraler KI- und Security-Begriffe . . . . .</b>	<b>351</b>
<b>B</b>	<b>Übersicht relevanter Normen und Frameworks . . . . .</b>	<b>361</b>
	<b>Stichwortverzeichnis . . . . .</b>	<b>365</b>

# Das neue Paradigma: KI trifft Cybersicherheit

Cybersicherheit war schon immer ein Wettlauf – aber lange Zeit ein Wettlauf mit vertrauten Regeln: Angreifer entwickeln neue Methoden, Verteidiger bauen Detektion und Kontrollen nach, und die Effizienzgewinne entstehen hauptsächlich durch bessere Tools, mehr Automatisierung und mehr Daten. Mit moderner KI – insbesondere großen Sprachmodellen (LLMs) und generativen Modellen – verschieben sich diese Regeln gerade spürbar. Nicht, weil KI »noch ein Tool« ist, sondern weil sie eine neue Klasse von Fähigkeiten in die Breite bringt: Verständnis, Synthese und Handlungsfähigkeit in natürlicher Sprache und Code.

Zum ersten Mal verfügen Angreifer und Verteidiger über Werkzeuge, die nicht nur Muster erkennen, sondern Bedeutung ableiten können: aus Logdaten, Netzwerkmetrie, Identitäts- und Verhaltenssignalen. Sie können natürliche Sprache verarbeiten – also genau das Medium, in dem ein großer Teil moderner Arbeit stattfindet: E-Mails, Support-Tickets, Runbooks, Dokumentation, Chat-Verläufe, Policies. Und sie können Code nicht nur lesen, sondern generieren, variieren und testen – von Scripts über Infrastructure-as-Code bis hin zu Exploit-ähnlichen Proofs of Concept. Was früher Spezialwissen, Zeit und mehrere Rollen erforderte, kann heute teilweise in Minuten orchestriert werden.

Damit ist KI kein inkrementelles Upgrade der bisherigen Security-Landschaft. Es ist ein Paradigmenwechsel: Wir gehen weg von überwiegend deterministischen, regelbasierten Systemen (»wenn X, dann Y«) hin zu lernenden, probabilistischen Systemen, die Wahrscheinlichkeiten und Kontext bewerten. Diese Systeme lassen sich nicht vollständig »durchkonfigurieren« wie klassische Security-Controls. Stattdessen müssen sie gesteuert werden: durch Datenqualität, klare Aufgabenabgrenzung, überprüfbare Outputs, robuste Grenzen (Guardrails) und kontrollierte Integrationen. Kurz: Wir konfigurieren nicht nur Tools – wir managen ein neues Risikoprofil.

Dieser Wandel wirkt in mehrere Richtungen:

Erstens: Verteidiger können schneller werden. KI kann operative Last reduzieren, Analysen beschleunigen, triagieren, korrelieren, zusammenfassen, Prioritäten vorschlagen und sogar Gegenmaßnahmen vorbereiten. Das betrifft SOC-Prozesse genauso wie Secure SDLC, Threat Modeling, Asset Discovery oder Third-Party-

Risk. Die Chance ist real: mehr Wirkung pro Kopf, bessere Reaktionszeiten, konsistentere Entscheidungen.

Zweitens: Angreifer werden skalierbarer und glaubwürdiger. Social Engineering wird qualitativ besser und quantitativ massentauglich, Malware- und Script-Varianten entstehen schneller, Reconnaissance und Targeting werden effizienter, und Angreifer können Verteidigungsmaßnahmen leichter umgehen, indem sie Sprache, Verhalten und Infrastruktur dynamisch anpassen. KI senkt nicht nur die Kosten – sie senkt vor allem die Reibung. Und das ist im Angreifer-Ökosystem ein enormer Hebel.

Drittens (und oft unterschätzt): Der Einsatz von KI erzeugt neue Angriffsflächen und neue Compliance-Risiken. Prompt Injection, Datenabfluss über Modelle, Schatten-KI im Unternehmen, unkontrollierte Tool-Integrationen, fehlerhafte oder halluzinierte Outputs, Lizenz- und Urheberrechtsfragen, Datenschutz, sowie regulatorische Anforderungen – all das kann aus einem »Produktivitätsgewinn« schnell ein strukturelles Risiko machen, wenn Governance, Architektur und Betrieb nicht mitwachsen.

Für IT- und Security-Verantwortliche ergeben sich daraus drei zentrale Leitfragen, die über die nächsten Jahre entscheiden werden:

1. Wie nutzen wir KI so, dass unsere Verteidigung messbar stärker wird?  
Nicht als Spielerei, sondern als Capability: schnellere Detection, bessere Priorisierung, geringere MTTR, robustere Engineering-Prozesse, weniger manuelle Fehler.
2. Wie schützen wir uns gegen Angreifer, die KI als Beschleuniger einsetzen?  
Das umfasst neue Phishing-Qualität, automatisierte Recon, variantenreiche Payloads, aber auch KI-gestützte Umgehung von Kontrollen – also eine Verschiebung im Threat Model.
3. Wie verhindern wir, dass KI selbst zum Risiko wird – technisch, rechtlich und organisatorisch?  
Governance, Datenklassifizierung, Zugriffskontrollen, Logging, Lieferkette, Modell- und Tool-Risiken sowie klare Verantwortlichkeiten müssen Teil des Security-Programms werden – nicht ein Add-on.

Dieses Kapitel legt das Fundament für genau diese Neubewertung. Es klärt zentrale Begriffe, grenzt KI-basierte Ansätze von klassischen Security-Mechanismen ab und zeigt, warum sich Sicherheitsstrategien in den kommenden Jahren verändern müssen: weniger Fokus auf starre Regeln als alleinige Wahrheit – mehr Fokus auf Resilienz, Überprüfbarkeit, kontrollierte Automatisierung und ein Threat Model, in dem sowohl Menschen als auch Maschinen auf beiden Seiten des Spielfelds agieren.

## 1.1 Was wir mit KI in der Cybersicherheit wirklich meinen

Bevor wir über Architektur, Risiken und konkrete Use Cases sprechen, müssen wir ein gemeinsames Begriffsverständnis schaffen. »KI« ist in der Sicherheitsdebatte inzwischen ein Sammelbegriff für sehr unterschiedliche Technologien – und genau das führt regelmäßig zu Missverständnissen: Manche meinen klassische Anomalieerkennung in einem SIEM, andere denken an Chatbots, wieder andere an autonom handelnde Agents, die Tickets schließen und Changes ausrollen. Wenn wir nicht sauber trennen, reden wir schnell aneinander vorbei – und treffen am Ende falsche Entscheidungen bei Tooling, Governance und Risikobewertung.

Wichtig ist auch die Abgrenzung: Nicht alles, was »KI« genannt wird, ist neu. Viele Security-Produkte nutzen seit Jahren statistische Verfahren, Heuristiken oder ML-basierte Klassifikatoren. Neu ist weniger die Existenz von ML – neu ist, dass generative Modelle und insbesondere LLMs Fähigkeiten in die Breite bringen, die bisher stark menschlich geprägt waren: Sprache verstehen, Wissen synthetisieren, Code erzeugen, kontextbezogen argumentieren und Aufgabenketten orchestrieren. Genau diese Fähigkeiten verschieben die operative Realität im Security-Alltag.

### 1.1.1 Künstliche Intelligenz als Oberbegriff

Künstliche Intelligenz (KI) beschreibt grundsätzlich Systeme, die Aufgaben übernehmen, für die bislang menschliche Intelligenz erforderlich war. In der Cybersicherheit betrifft das vor allem vier Klassen von Fähigkeiten:

- **Muster erkennen und klassifizieren:**

KI kann aus großen Datenmengen wiederkehrende Strukturen ableiten – etwa typische Merkmale von Malware, verdächtige Kommunikationsprofile oder auffällige Login-Muster. Der Kern ist: Aus »viel Telemetrie« wird »eine bewertbare Hypothese«.

- **Natürliche Sprache verstehen und verarbeiten:**

Ein erheblicher Teil von Security passiert in Textform: Phishing-Mails, Ticket-Kommentare, Incident-Reports, Policies, Chat-Verläufe, Post-Mortems. KI kann diese Inhalte zusammenfassen, einordnen, priorisieren und in Maßnahmen übersetzen – also genau das, was bisher viel menschliche Arbeitszeit gebunden hat.

- **Entscheidungen unter Unsicherheit unterstützen:**

Security ist selten schwarz-weiß. Meist geht es um Wahrscheinlichkeiten, unvollständige Informationen, widersprüchliche Signale. KI kann hier helfen, Optionen zu bewerten, Risiken zu gewichten und sinnvolle nächste Schritte vorzuschlagen – nicht als Ersatz für Verantwortung, aber als Beschleuniger für Analyse und Triage.

- **Neue Inhalte generieren:**

Dazu zählen Texte, Playbooks, Detection-Queries, Code, Tests, IaC-Snippets, aber auch verständliche Erklärungen für unterschiedliche Zielgruppen (Engineering vs. Management). Das ist einer der größten Produktivitätshebel – und gleichzeitig eine der größten Risikoquellen, wenn Outputs ungeprüft übernommen werden.

Im Alltag der Cybersicherheit sind für uns vor allem vier Teilbereiche relevant, die sich teilweise überlappen, aber unterschiedliche Stärken und Risiken haben:

- **Machine Learning (ML):**

ML-Systeme lernen Muster aus Daten, statt Regeln ausschließlich manuell vorzugeben. Typische Security-Anwendungen sind Klassifikation (»bösaartig vs. gutartig«), Scoring (Risiko- oder Prioritätswerte), Clustering (Gruppierung ähnlicher Ereignisse) und Anomalieerkennung (Abweichungen vom Normalverhalten). ML ist oft dort stark, wo viele strukturierte Daten vorliegen und ein klarer Feedback-Loop möglich ist.

- **Deep Learning (DL):**

Deep Learning ist eine Unterkategorie von ML und nutzt komplexe neuronale Netze, um sehr hochdimensionale Daten zu verarbeiten – z. B. Sprache, Text, Sequenzen oder Binärdaten. In Security kann DL etwa bei der Analyse von Malware-Familien, bei der Auswertung großer Textmengen oder bei Verhaltensmustern helfen. Der Preis ist häufig: weniger Transparenz, höhere Anforderungen an Daten und Rechenleistung und schwierigere Erklärbarkeit.

- **Large Language Models (LLMs):**

LLMs sind Modelle, die natürliche Sprache verstehen und generieren können. Praktisch werden sie oft als »Assistenten« genutzt: sie lesen und schreiben Text, erklären Zusammenhänge, erstellen Zusammenfassungen, formulieren Tickets, generieren Code, helfen bei Root-Cause-Analysen und können zunehmend auch Werkzeuge ansteuern (z. B. Abfragen gegen Logsysteme), wenn man sie entsprechend integriert. Ihr großer Vorteil ist ihre Schnittstelle: Sprache – und damit die Fähigkeit, Security-Wissen und Security-Arbeit in natürlicher Form zu skalieren.

- **Generative KI (GenAI):**

Generative KI ist der Oberbegriff für Modelle, die neue Inhalte erzeugen können – Text, Code, Bilder, Audio und mehr. LLMs sind ein Teil davon (GenAI für Text/Code). In Security ist GenAI besonders relevant für Automation, Dokumentation, Detection Engineering, Secure Coding Support und Training. Gleichzeitig erhöht GenAI das Risiko für überzeugende Täuschung (Phishing, Impersonation) und für unkontrollierte Content-Erzeugung (z. B. unsichere Code-Snippets oder falsche Handlungsempfehlungen).

Die zentrale Konsequenz aus dieser Begriffsarbeit ist simpel, aber entscheidend: Wenn wir »KI in der Cybersicherheit« sagen, meinen wir nicht ein einzelnes Produkt und nicht eine einzelne Methode. Wir meinen ein Spektrum von Fähigkeiten – von statistischer Klassifikation bis hin zu sprach- und handlungsfähigen Systemen. Und je nachdem, welche Kategorie wir einsetzen, ändern sich Nutzen, Grenzen, Angriffsfläche, Governance-Anforderungen und die Art, wie wir Sicherheit »bauen und betreiben«.

### 1.1.2 Machine Learning vs. Regeln – der fundamentale Unterschied

Viele der klassischen Security-Kontrollen, mit denen wir groß geworden sind, funktionieren nach einem klaren Prinzip: Regeln und bekannte Muster. Signaturbasierter Virenschutz, IDS/IPS-Regeln, YARA-Signaturen, SIEM-Korrelationen oder starre DLP-Policies sind im Kern deterministisch aufgebaut. Sie folgen Logiken wie: »Wenn Bedingung A und B erfüllt sind, dann ist das Ereignis böseartig« – oder zumindest alarmwürdig.

Das hat Vorteile: Solche Systeme sind oft transparent, reproduzierbar und gut zu auditieren. Wenn ein Alarm auslöst, kann man häufig nachvollziehen, welche Regel gegriffen hat. Der Nachteil ist ebenso offensichtlich: Diese Ansätze funktionieren am besten dort, wo das Problem bereits bekannt ist. Sie sind stark gegen »das, was wir schon gesehen haben« – und schwächer gegen Varianten, neue Taktiken oder bewusstes Evasion-Verhalten.

Machine Learning (ML) dreht diese Logik teilweise um. Statt dass Menschen jede Erkennungslogik explizit formulieren, lernt ein Modell Muster aus Daten. Drei Punkte sind dabei entscheidend:

- Daten als Grundlage:

ML-Systeme brauchen Trainingsdaten (und oft auch Validierungsdaten), die das »Normal« und das »Verdächtig« abbilden. Die Qualität der Ergebnisse hängt direkt davon ab, wie gut diese Daten sind: vollständig, repräsentativ, sauber gelabelt und aktuell.

- Wahrscheinlichkeiten statt Wahrheiten:

ML entscheidet selten binär. Statt »gut/böse« liefert es häufig ein Score, eine Wahrscheinlichkeit oder eine Rangfolge: »mit 93 % verdächtig« oder »Top 5 % auffällig im Vergleich zur Baseline«. Das ist in der Praxis extrem nützlich – aber es zwingt uns, mit Unsicherheit professionell umzugehen.

- Verallgemeinerung auf Unbekanntes:

Der große Hebel von ML ist die Fähigkeit, Muster zu erkennen, die nicht exakt im Regelwerk stehen. Ein Modell kann Varianten auffangen, die keiner Signatur entsprechen, und Abweichungen markieren, die noch nie exakt so vorkamen – gerade bei Verhalten, Sequenzen oder Kombinationsmustern.

Für Sie als IT- und Security-Verantwortlichen hat das eine konkrete Konsequenz: Sie »programmieren« die Erkennung nicht mehr primär über einzelne Regeln, sondern über Rahmenbedingungen, in denen das System lernt und entscheidet. Das bedeutet in der Praxis:

- Sie entscheiden, welche Datenquellen überhaupt in das Modell fließen (Identitätsdaten, Endpoint-Telemetrie, Cloud-Logs, Netzwerkdaten, Tickettexte etc.).
- Sie definieren, wie Erfolg gemessen wird: Was ist ein False Positive? Was kostet ein False Negative? Was ist die akzeptable Alarmrate?
- Sie bestimmen, welche Autonomie das System erhält: Nur Vorschläge machen? Automatisch Tickets erstellen? Quarantäne auslösen? Accounts sperren?
- Sie bauen Kontrollen um das Modell herum: Monitoring von Drift, Qualitätschecks, human-in-the-loop, Rollback-Mechanismen.

Kurz gesagt: Bei regelbasierten Systemen ist die Hauptarbeit »Regeln schreiben und pflegen«. Bei ML-basierten Systemen ist die Hauptarbeit »Daten, Bewertung und Leitplanken managen«. Beides ist Security Engineering – aber es ist ein anderer Hebel und ein anderes Betriebsmodell.

### 1.1.3 LLMs und GenAI – warum sie Security besonders verändern

LLMs und generative KI verändern Security nicht nur, weil sie »bessere ML-Modelle« sind, sondern weil sie eine neue Schnittstelle in den Betrieb bringen: Sprache und Code als universelle Arbeitsoberfläche. Und genau darin liegt die Sprengkraft – positiv wie negativ.

LLMs sind für Security deshalb so relevant, weil sie drei Fähigkeiten kombinieren, die bisher selten in einem System zusammenkamen:

- Natürliche Sprache verstehen und erzeugen:  
Ein großer Teil der Security-Arbeit steckt in unstrukturierten Artefakten: Phishing-Mails, Incident-Notes, Chat-Verläufe, Post-Mortems, Policies, Risikoausnahmen, Audit-Fragen, Vendor-Fragebögen. LLMs können diese Inhalte lesen, zusammenfassen, strukturieren, vergleichen und in Handlungsvorschläge übersetzen.
- Code lesen, schreiben und erklären:  
Security lebt von Code: Detection-Queries, Regex/YARA, Sigma-Regeln, SIEM-Korrelationen, Scripts, IaC, CI/CD-Konfigurationen. LLMs können helfen, Code zu reviewen, Schwachstellen zu erklären, Fixes vorzuschlagen, Tests zu erzeugen oder Regeln schneller zu formulieren. Das beschleunigt besonders DevSecOps- und Detection-Engineering-Workflows.

■ Interaktion und Orchestrierung über Prozesse hinweg:

LLMs sind nicht nur »Textgeneratoren«, sondern werden in der Praxis als Interface genutzt: als Chatbot im SOC, als Assistenz im Change-Prozess, als Unterstützung im Incident-Management. In Kombination mit Tool-Integrationen können sie Informationen abfragen, Tickets anlegen, Reports erstellen oder standardisierte Workflows anstoßen.

In der Security-Praxis sehen wir LLMs und GenAI bereits in sehr konkreten Rollen, zum Beispiel:

- automatische SOC-Lageberichte aus Events, Alerts und Tickets,
- Incident-Zusammenfassungen inkl. Timeline, Impact und Next Steps,
- Assistenz bei der Log-Analyse (»Was ist hier auffällig? Welche Hypothesen passen?«),
- Policy-Entwurf und Review (z. B. Abgleich mit Standards/Controls),
- Unterstützung im Secure SDLC (Code Review, Threat Modeling, Abuse Cases, Security Tests).

Der entscheidende Punkt ist jedoch: Die gleichen Fähigkeiten stehen auch Angreifern zur Verfügung – und dort wirken sie als Skalierungsfaktor. Besonders relevant sind:

- Hochqualitative, überzeugende Phishing- und Impersonation-Mails (stilistisch sauber, kontextbezogen, sprachlich passend),
- Automatisierte Reconnaissance (Auswertung öffentlicher Informationen, Ableitung von Rollen/Prozessen, Identifikation von Angriffswegen),
- Schnelle Variation von Payloads und Exploit-ähnlichem Code (nicht zwingend »Hollywood-Malware«, aber genug, um Defenses zu testen, Evasion zu betreiben oder Initial Access zu erleichtern).

Wichtig ist dabei die richtige Einordnung: Viele Modelle sind durch Policies und technische Schutzmechanismen begrenzt und nicht jede direkte »Malware-Generierung« ist trivial. Aber selbst ohne »vollautomatische Exploit-Fabrik« reicht die Kombination aus Sprachqualität, Automatisierung und Codeassistenz aus, um die Effektivität vieler Angriffe deutlich zu erhöhen – vor allem bei Social Engineering und bei der Geschwindigkeit, mit der Angreifer iterieren können.

Die Konsequenz für Security-Programme ist klar: LLMs/GenAI sind zugleich Verteidigungsbeschleuniger und Angriffsverstärker. Wer sie nur als Produktivitätstool betrachtet, unterschätzt das Risiko. Wer sie nur als Risiko betrachtet, verschenkt einen realen Hebel. Die richtige Antwort ist: gezielter Einsatz mit klaren Grenzen, messbaren Ergebnissen und einem aktualisierten Threat Model.

## 1.2 Von regelbasierten zu lernenden Security-Systemen

Um zu verstehen, warum KI in der Cybersicherheit mehr ist als ein weiteres Feature-Upgrade, lohnt sich ein kurzer Blick zurück. Security-Mechanismen haben sich historisch entlang einer klaren Linie entwickelt: erst Signaturen und harte Regeln, dann datengetriebene Modelle, und jetzt – mit GenAI und LLMs – Systeme, die nicht nur erkennen, sondern kommunizieren, erklären und (teilweise) handeln können. Jede Stufe hat unsere Verteidigung leistungsfähiger gemacht, aber auch neue Betriebs- und Risikofragen geschaffen.

### 1.2.1 Die Ära der Signaturen und starren Regeln

Die lange dominierende Phase der IT-Sicherheit war geprägt von deterministischen Kontrollen: Ein Ereignis wird gegen ein vordefiniertes Muster geprüft, und daraus folgt eine eindeutige Entscheidung. Typische Vertreter dieser Ära sind:

- Virens Scanner mit Signaturen: Hashes, Bytefolgen, bekannte Malware-Familien.
- Firewalls mit festen Regeln: »Port X nach Y blockieren«, »Subnetz A darf nicht mit Subnetz B sprechen«.
- IDS/IPS mit statischen Mustern: bekannte Angriffssignaturen, Protokollanomalien, feste Regeln.
- SIEM-Korrelationen: händisch gepflegte Logik wie »wenn ungewöhnlicher Login + Admin-Änderung + Datenexfil-Indikator, dann Alarm«.

Das war über viele Jahre der Standard – aus guten Gründen. Diese Kontrollen haben drei klassische Stärken:

1. Nachvollziehbarkeit: Man kann meist erklären, warum etwas blockiert oder alarmiert wurde.
2. Determinismus: Gleiche Inputs führen zu gleichen Outputs; das ist operativ stabil.
3. Compliance-Freundlichkeit: Kontrollen lassen sich gut dokumentieren (»Regel X mitigiert Risiko Y«), auditieren und reproduzieren.

Mit der Zeit wurden jedoch die Grenzen immer deutlicher – besonders in dynamischen, hochvernetzten Umgebungen:

- Reaktivität: Regel- und signaturbasierte Systeme erkennen primär das, was bereits bekannt ist. Neue Varianten, neue Taktiken und Evasion laufen leichter durch.
- Pflegeaufwand: Signaturen, Regeln, Ausnahmen, Whitelists – alles muss dauerhaft gepflegt werden. Das skaliert nur begrenzt, vor allem bei hoher Change-Rate.

- Schwache Skalierung bei Datenvolumen und Komplexität: Mit Cloud, SaaS, Container-Plattformen und verteilten Identitäten explodieren Events und Kontextsignale. »Mehr Regeln« ist dann selten die Lösung – oft ist es der direkte Weg in Alarmmüdigkeit.

Die Kernerkenntnis: Regelwerke funktionieren gut für bekannte, klar abgrenzbare Muster. Sie versagen dort, wo Verhalten, Kontext und Variabilität dominieren.

## 1.2.2 Übergang zu ML-basierten Detektionen

Mit zunehmender Komplexität moderner IT-Landschaften – Cloud, Microservices, Remote Work, BYOD, Identity-first-Architekturen – wurde ein Problem offensichtlich: Wir können nicht mehr jede relevante Event-Kombination sinnvoll per Hand modellieren. Nicht nur, weil es zu viele sind, sondern weil sich »Normalverhalten« ständig verändert: Releases, neue Tools, neue Teams, saisonale Schwankungen, neue Workloads.

An dieser Stelle wurden ML-Ansätze in vielen Security-Bereichen populär, unter anderem:

- Anomalieerkennung im Netzwerkverkehr: Abweichungen in Flows, Protokollen, Datenmengen oder Kommunikationsbeziehungen.
- UEBA (User and Entity Behavior Analytics): Verhalten von Benutzerkonten, Service Accounts, Endpoints oder Workloads wird modelliert, um »ungewöhnliche« Aktivitäten zu identifizieren.
- ML-basierte Malware-Erkennung: Statt nur auf Signaturen zu reagieren, werden Features (z. B. Verhalten, Struktur, Aufrufmuster, Metadaten) genutzt, um neue oder leicht veränderte Samples einzuordnen.

Hier beginnt der eigentliche Paradigmenwechsel: Security wird daten- und modellgetrieben, nicht mehr ausschließlich regelgetrieben. Das verändert die Aufgaben des Security-Teams fundamental.

Wo früher die Kernarbeit war: »Regeln schreiben und Exceptions pflegen«, verschiebt sich der Schwerpunkt nun zu:

- Daten verstehen und kuratieren: Welche Telemetrie ist vollständig? Welche Felder fehlen? Welche Daten sind zuverlässig? Wo gibt es Blind Spots?
- Modelle evaluieren: Wie gut sind Precision/Recall? Wie stabil ist das Modell bei Änderungen? Wie reagieren wir auf Drift?
- Outputs interpretieren und in Maßnahmen übersetzen: Ein Score allein löst kein Incident. Das Team muss Kontext herstellen, Prioritäten setzen, Maßnahmen ableiten und das Ganze operationalisieren.

Das ist eine andere Art von Security-Betrieb: weniger Handwerk am Regelwerk, mehr Engineering an Datenpipelines, Qualitätsmetriken, Feedback-Loops und operativen Leitplanken.

### 1.2.3 Der Sprung mit GenAI und LLMs

Der nächste Entwicklungsschritt ist nicht nur »mehr ML«, sondern qualitativ anders: GenAI und LLMs bringen Fähigkeiten in Security-Systeme, die bisher überwiegend menschlich waren – nämlich Sprache, Synthese und Wissensarbeit.

Plötzlich können Security-Systeme mit Ihnen interagieren: im SOC, in GRC-Prozessen, in der Softwareentwicklung, im Change Management. Und sie können dabei nicht nur strukturierte Events verarbeiten, sondern auch unstrukturierte Informationen, die in klassischen Pipelines oft »außen vor« waren, zum Beispiel:

- E-Mails und Phishing-Threads
- Chat-Logs aus Slack/Teams
- Runbooks und Confluence-Seiten
- PDF-Reports, Audit-Dokumente, Lieferantenunterlagen
- Ticket-Historien und Incident-Kommunikation

Zusätzlich können sie neue Artefakte erzeugen – also nicht nur bewerten, sondern produzieren:

- Reports (Lagebilder, Executive Summaries, Audit-Readouts)
- Playbooks (Incident-Workflows, Standard Responses, Runbooks)
- Detections (Queries, Regeln, Sigma/YARA-ähnliche Patterns)
- Code und Konfigurationen (Scripts, IaC-Änderungen, CI/CD-Security-Checks)

Damit entsteht ein neues Rollenmodell im Sicherheitsbetrieb: KI ist nicht mehr nur Sensor (erkennen) oder Filter (blockieren), sondern ein aktiver Assistent – und in manchen Architekturen sogar ein Akteur, der Entscheidungen vorbereitet oder unter klaren Bedingungen selbst ausführt.

Und genau hier liegt die strategische Relevanz: Sobald KI nicht nur »Signal« liefert, sondern Arbeitsergebnisse erzeugt und Prozesse antreibt, müssen wir Sicherheit neu denken – inklusive Governance, Verantwortlichkeiten, Qualitätssicherung, Auditierbarkeit und Kontrollmechanismen. Der Paradigmenwechsel ist nicht nur technisch. Er ist organisatorisch und operativ.

Die Entwicklung der Cybersicherheit lässt sich grob in drei Phasen einteilen: von klassisch signatur- und regelbasierten Systemen über erste ML-gestützte Anomalieerkennung hin zu generativer KI und LLMs, die heute ganze Analyse- und Kommunikationsaufgaben übernehmen können. Abbildung 1.1 stellt diese Evolution und den damit verbundenen Paradigmenwechsel übersichtlich dar.

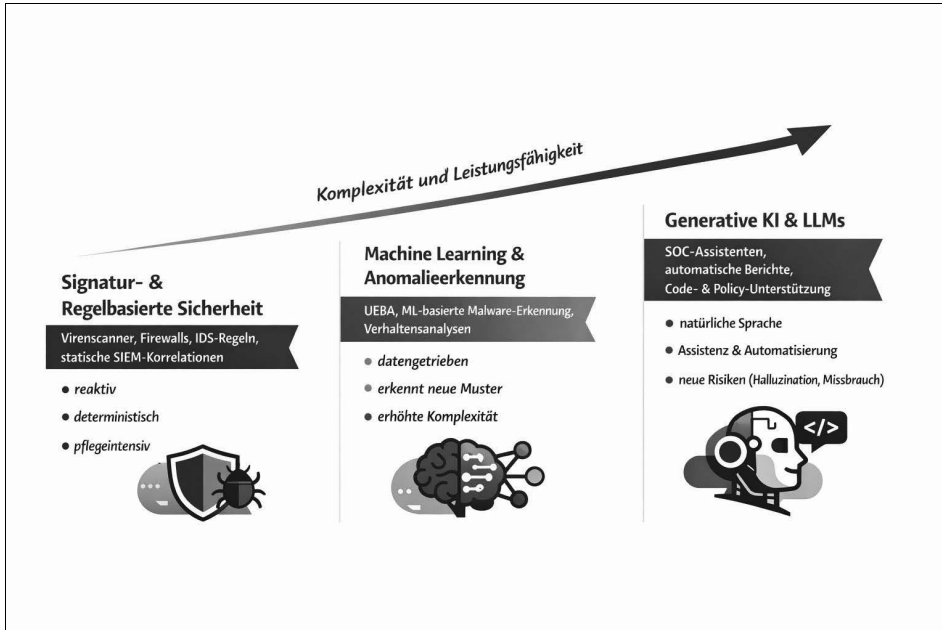


Abb. 1.1: Evolution der Cybersicherheit

## 1.3 Drei Rollen von KI in der Cybersicherheit

Wenn Sie eine KI-Strategie für Security entwickeln wollen, hilft es enorm, KI nicht als »eine Technologie« zu betrachten, sondern als Akteur in drei Rollen. Denn je nachdem, in welcher Rolle KI auftritt, ändern sich Zielbild, Architektur, Controls und Risikoprofil komplett. In der Praxis begegnet Ihnen KI gleichzeitig als Verstärker Ihrer Teams, als Beschleuniger auf Angreiferseite und als eigenständige Risikoquelle, die Sie wie ein neues System mit eigenen Failure-Modes behandeln müssen.

### 1.3.1 KI als Verstärker der Verteidiger

In der ersten und naheliegendsten Rolle fungiert KI als Produktivitäts- und Qualitätshebel für Verteidiger. Sie ersetzt nicht die Verantwortung des Security-Teams, aber sie kann große Teile der repetitiven Arbeit beschleunigen und die »kognitive Last« reduzieren – vor allem dort, wo Menschen bisher viel Zeit mit Suchen, Zusammenfassen und Übersetzen von Informationen verbringen.

Typische Einsatzfelder sind:

- SOC-Augmentation (Triage & Analyse):

KI kann Alerts zusammenfassen, Kontext aus mehreren Systemen anreichern (z. B. Asset-Info, Identity-Kontext, bekannte Kampagnenmuster), ähnliche Fälle

clustern und Playbook-Schritte vorschlagen. Der operative Effekt: weniger Zeit, um »den Fall zu verstehen«, mehr Zeit, um Entscheidungen zu treffen.

■ **Threat Hunting & Detection Engineering:**

Viele Hunting-Iterationen scheitern nicht an fehlenden Ideen, sondern an Reibung: Query-Formulierung, Datenquellen verstehen, Ergebnisinterpretation. KI kann beim Erstellen komplexer Queries helfen, Hypothesen in Suchlogik übersetzen und Ergebnisse erklären – inklusive Vorschläge für »nächste Pivot-Schritte«.

■ **Incident Response & Kommunikation:**

In Incidents ist Zeit der kritische Faktor – und gleichzeitig steigt der Kommunikationsdruck (Stakeholder, Management, Kunden, ggf. Behörden). KI kann Drafts für Tickets, Status-Updates, Benachrichtigungen, Executive Summaries und Post-Incident-Reports erstellen, basierend auf den vorhandenen Fakten aus Timeline und Artefakten. Das entlastet Teams, ohne dass die inhaltliche Verantwortung abgegeben wird.

■ **Knowledge Management & »Security-Wissen on demand«:**

Viele Organisationen haben Security-Wissen – aber es ist verteilt: Policies hier, Architekturdiagramme dort, Lessons Learned in Post-Mortems, Ausnahmen in Tickets. Über RAG-Ansätze (Retrieval-Augmented Generation) kann KI dieses Wissen kontextsensitiv verfügbar machen: »Welche Policy gilt für Datenklasse X?« oder »Wie war das Runbook für Incident-Typ Y?«. Das reduziert Suchzeiten und erhöht Konsistenz.

Das Ziel dieser Rolle ist nicht »KI macht Security«, sondern: KI reduziert Low-Value-Arbeit und erhöht die Schlagkraft der Experten. In der Praxis heißt das:

- weniger »Copy & Paste« und Formatierung,
- weniger Zeit mit »Wo finde ich diese Information?«,
- mehr Zeit mit »Was bedeutet das – und was tun wir jetzt?«.

Wenn KI in dieser Rolle richtig eingesetzt wird, verbessert sie Geschwindigkeit, Qualität und Konsistenz – und verschiebt Kapazität zurück in Analyse, Engineering und Risikosteuerung.

### 1.3.2 KI als Werkzeug der Angreifer

Die zweite Rolle ist unangenehm, aber zentral für ein realistisches Threat Model: Dieselben Fähigkeiten, die Verteidiger produktiver machen, senken auch auf Angreiferseite die Kosten und erhöhen die Erfolgswahrscheinlichkeit – vor allem durch Skalierung und Qualität.

Typische Angreifer-Use-Cases sind:

- **Phishing auf Industriestandard:**  
Angriffe werden sprachlich fehlerfrei, stilistisch passend, in Landessprache, mit überzeugender Tonalität und hoher Kontextnähe formuliert. Das erhöht nicht nur die Klickrate, sondern auch die Glaubwürdigkeit bei Rückfragen (»Conversation Hijacking«, längere Mail-Threads, präzise Referenzen).
- **Social Engineering & Identitätsbetrug:**  
Generierte Personas, Lebensläufe, Profile und Kommunikationsmuster machen es leichter, Vertrauen aufzubauen – bis hin zu Deepfake-gestützter Audio-/Video-Interaktion in ausgewählten Szenarien. Der Kern ist nicht »Hollywood«, sondern: weniger Reibung bei der Täuschung.
- **Malware- und Exploit-nahe Entwicklung:**  
KI kann Codevarianten erzeugen, Obfuskation unterstützen, Evasion-Ideen liefern oder Komponenten schneller zusammensetzen (z. B. Loader, Dropper, C2-Kommunikation). Nicht jede Modellinstanz liefert »fertige Malware«, aber der Output reicht oft, um Iterationszyklen zu beschleunigen und Defenses gezielt zu testen.
- **Reconnaissance als Automationspipeline:**  
Angreifer können Informationen aus GitHub, LinkedIn, Stellenanzeigen, Dokumentation, öffentlichen Repos, Leak-Datenbanken und Metadaten automatisch auswerten, verdichten und in konkrete Angriffspläne übersetzen: Wer nutzt welche Tools? Welche Rollen gibt es? Wo liegen wahrscheinlich Secrets? Welche Third Parties sind im Einsatz? Welche Prozesse sind anfällig?

Das Ergebnis ist eindeutig: Angriffe werden massiv skaliert und professionalisiert – und zwar nicht nur durch Top-Tier-Gruppen, sondern auch durch weniger technisch versierte Täter, die mit KI »Kompetenz einkaufen«. Das verändert die Baseline: Mehr Angriffe, bessere Qualität, schnellere Iteration. Für Verteidiger bedeutet das: Der »Durchschnittsgegner« wird stärker.

### 1.3.3 KI als »unberechenbarer Dritter«

Die dritte Rolle wird in vielen Strategien unterschätzt: KI ist nicht nur Tool der Guten oder der Bösen – KI ist auch ein System mit eigenen Fehlermodi. Moderne Modelle sind häufig:

- probabilistisch (sie liefern Wahrscheinlichkeiten, keine Gewissheiten),
- intransparent (die Entscheidungslogik ist nicht vollständig erklärbar),
- nicht deterministisch (gleiche Eingabe kann variierende Ausgaben erzeugen, je nach Kontext/State).

Daraus entstehen spezifische Risiken, die Sie in Security-Architektur und Governance einplanen müssen:

■ **Halluzinationen:**

KI kann falsche Aussagen erzeugen, die plausibel klingen – besonders gefährlich, wenn sie in operative Entscheidungen einfließen (z. B. falsche Root-Cause, falsche Remediation, falsche Policy-Auslegung).

■ **Fehleinschätzungen und Bias:**

Ein Modell kann eine Situation als »harmlos« bewerten, obwohl sie kritisch ist – oder umgekehrt Alarmismus erzeugen. Wenn solche Fehler systematisch auftreten (z. B. aufgrund unausgewogener Trainings-/Kontextdaten), werden sie schwer erkennbar und teuer im Betrieb.

■ **Mangelnde Auditierbarkeit und Erklärbarkeit:**

In regulierten Umfeldern – oder generell, wenn Sie Entscheidungen später rechtfertigen müssen – ist »das Modell hat es so gesagt« keine ausreichende Begründung. Ohne saubere Protokollierung, Quellenbezug und Kontrollmechanismen wird KI operativ riskant.

Für Security-Teams folgt daraus eine sehr praktische Leitlinie: Behandeln Sie KI wie ein extrem leistungsfähiges, aber fehleranfälliges Teammitglied. Das heißt:

- klare Verantwortlichkeiten (wer entscheidet wirklich?),
- klare Grenzen (welche Aktionen darf KI nie/immer/unter Bedingungen ausführen?),
- klare Kontrollmechanismen (Review, Logging, Validierung, Rollback),
- und ein Betriebsmodell, das mit Unsicherheit umgehen kann.

So eingesetzt wird KI zu einem Multiplikator – nicht zu einem unkontrollierten Risiko.

## 1.4 Chancen für Verteidiger: Wo KI wirklich hilft

Für IT- und Security-Verantwortliche ist die entscheidende Frage nicht, ob KI »spannend« ist, sondern ob sie konkret messbaren Mehrwert liefert: weniger Risiko, schnellere Reaktion, bessere Qualität – oder schlicht mehr Output pro Team. Der größte Hebel liegt dabei nicht in »vollautonomen« Systemen, sondern in gut designeten Assistenz- und Automationsbausteinen, die Engpässe im Betrieb entlasten und vorhandene Daten endlich nutzbar machen.

### 1.4.1 Entlastung im SOC

Das Security Operations Center (oder die SOC-Funktion) ist oft der Ort, an dem Komplexität und Volumen am brutalsten aufschlagen: zu viele Alerts, zu wenig

Kontext, zu viele Systeme, zu wenig Zeit. KI kann hier besonders effektiv helfen, weil ein großer Teil der Arbeit aus Triagieren, Zusammenfassen, Priorisieren und Kommunizieren besteht.

Typische KI-gestützte Einsatzfelder sind:

- **Automatisierte Lageberichte (täglich / wöchentlich):**  
Statt dass Analysten mühsam Tickets, Alerts und Trends manuell konsolidieren, kann KI daraus ein konsistentes Lagebild erstellen: Top Incidents, relevante Veränderungen, wiederkehrende Muster, offene Risiken, SLA-Status.
- **Clustering von Alerts nach Kampagnen oder Mustern:**  
Häufig sind 50 »einzelne Alerts« in Wahrheit 1–2 zusammenhängende Ereignisse. KI kann Signale gruppieren (z.B. nach Host, User, TTP-Ähnlichkeit, zeitlicher Nähe), damit Teams nicht in Ticket-Splitting untergehen, sondern kampagnenorientiert arbeiten.
- **Priorisierung von Incidents:**  
Nicht jeder Alarm ist gleich relevant. KI kann Prioritäten vorschlagen, indem sie Kontext einbezieht: Kritikalität des Assets, Exposure, Identitätsrechte, Blast Radius, bekannte Threat-Intel-Indikatoren, historische Vergleichsfälle. Das ersetzt nicht die Entscheidung – aber es bringt Struktur in die Flut.
- **Schnelle Zusammenfassung komplexer Incident-Timelines:**  
Gerade in laufenden Incidents ist die Timeline oft ein chaotischer Mix aus Logs, Chat, Ticket-Updates und Hypothesen. KI kann daraus eine verständliche, fortschreibbare Chronologie machen: »Was wissen wir sicher, was vermuten wir, was sind die nächsten Schritte?«.

Der operative Effekt ist unmittelbar:

- weniger Alert-Müdigkeit (weil mehr Kontext und weniger redundante Arbeit),
- bessere Fokussierung auf kritische Ereignisse (weil Priorisierung und Clustering Reibung reduzieren),
- kürzere Einarbeitungszeiten für neue Analysten (weil Zusammenfassungen, Erklärungen und standardisierte Workflows Lernen beschleunigen).

Das ist einer der wenigen Bereiche, in denen KI sehr schnell ROI liefern kann – sofern Sie klare Qualitätsmetriken und human-in-the-loop etabliert haben.

### 1.4.2 Bessere Nutzung vorhandener Daten

In vielen Organisationen ist nicht »zu wenig Sichtbarkeit« das Kernproblem, sondern das Gegenteil: Es gibt Unmengen an Daten – Logs, Tickets, Wiki-Artikel, Architektur-Dokus, Post-Mortems, E-Mail-Threads – aber sie sind verteilt, inkonsistent und schwer auffindbar. Sicherheit scheitert dann nicht an Technik, sondern an Orientierung.

Typische Symptome sind:

- niemand hat einen echten Überblick, wo welche Information liegt,
- Dokumentation ist veraltet oder widersprüchlich,
- Wissen steckt in Köpfen oder in Ticket-Kommentaren,
- Tools sind Insellösungen mit unterschiedlichen Datenmodellen.

KI-Ansätze – insbesondere RAG (Retrieval-Augmented Generation) – können hier einen sehr pragmatischen Hebel bieten:

- Verstreute Wissensinseln zusammenführen:  
Confluence, Ticketsysteme, Runbooks, Mails, Wikis, SharePoint – alles kann in einen such- und zitierfähigen Wissenslayer überführt werden, der kontextbezogenen Antworten liefert.
- Sicherheitswissen kontextsensitiv bereitstellen:  
Statt »Suchmaschine« bekommen Sie »Antwortmaschine«: nicht nur Trefferlisten, sondern eine Zusammenfassung mit relevanten Quellen, bezogen auf den konkreten Fall (»Welche Schritte sind bei diesem Incident-Typ vorgesehen?«).
- Doku-»Leichen« wieder nutzbar machen:  
Alte, verstreute Dokumente sind oft nicht wertlos – sie sind nur nicht zugänglich. KI kann Inhalte extrahieren, normalisieren, miteinander vergleichen und helfen, veraltete Stellen sichtbar zu machen.

Das Ergebnis ist weniger »Magie« als schlicht bessere Organisation: Wissen wird operationalisierbar, Entscheidungen werden konsistenter, und Teams verlieren weniger Zeit durch Suche und Interpretationsarbeit.

### 1.4.3 Beschleunigung von DevSecOps

Im Entwicklungs- und Betriebsumfeld liegt der größte Security-Schaden oft nicht in fehlenden Tools, sondern in Fehlkonfigurationen, unsicheren Defaults und Zeitdruck. KI kann hier helfen, Security in den Delivery-Prozess zu integrieren, ohne dass sie als reines »Stoppschild« wahrgenommen wird.

Konkrete Ansatzpunkte sind:

- Automatisierte Security-Code-Reviews (unter Aufsicht):  
KI kann Pull Requests voranalysieren, typische Schwachstellenmuster markieren, sichere Alternativen vorschlagen und Entwicklerkommentare in ein Security-Fazit übersetzen. Wichtig: als Assistenz, nicht als alleinige Autorität.
- Unterstützung beim Schreiben sicherer Konfigurationen (IaC):  
Gerade in Terraform/Kubernetes/CI-Pipelines sind die Fehler oft banal, aber folgenreich. KI kann Best Practices in konkrete Konfigurationsänderungen übersetzen und Risiken in verständlicher Form erklären.

- **Threat-Modeling-Assistenten:**  
Nicht jedes Team hat Security-Expertise für STRIDE/Abuse-Cases. KI kann Fragen stellen, Angriffspfade vorschlagen, Trust Boundaries strukturieren und daraus handhabbare Controls ableiten.
- **Generierung von Testfällen und Sicherheitstests:**  
Von Input-Validation-Tests über AuthZ-Checks bis hin zu API-Fuzzing-Ansätzen: KI kann helfen, Testideen zu systematisieren und Boilerplate zu reduzieren.

Das reduziert ganz konkret:

- »Security als Bottleneck« im Delivery-Prozess (weil Vorarbeit automatisiert wird),
- händische Routineaufgaben (weil Vorschläge und Reviews schneller entstehen),
- Fehlkonfigurationen durch Unwissen (weil Teams direkt beim Umsetzen unterstützt werden).

Richtig gemacht ist das kein »Security macht weniger«, sondern »Security kommt früher und leichter in den Prozess«.

#### 1.4.4 GRC & Compliance

Governance, Risk & Compliance ist für viele Organisationen der Bereich mit dem größten Anteil an formaler Arbeit: Policies schreiben, Normen mappen, Evidenzen sammeln, Audits vorbereiten, Management-Reports erstellen. Genau hier ist KI besonders wirksam, weil es viel um Textarbeit, Strukturierung und Abgleich geht.

Typische Einsatzfelder sind:

- **Erstellung und Pflege von Policies:**  
KI kann Policy-Entwürfe erstellen, Versionsänderungen konsistent einarbeiten, Formulierungen vereinheitlichen und auf Lücken prüfen – immer mit menschlicher Freigabe.
- **Mapping von Normen auf Unternehmensprozesse:**  
Ob ISO 27001, NIST, CIS Controls oder interne Standards: KI kann Zuordnungen vorschlagen, Überschneidungen sichtbar machen und Gap-Analysen strukturieren.
- **Analyse von Audit-Berichten und Findings:**  
KI kann Findings clustern, Root-Cause-Muster identifizieren, Wiederholungen über Jahre sichtbar machen und Maßnahmen in umsetzbare Workpakete übersetzen.

- Vorbereitung von Management-Reports und Aufsichtsunterlagen:  
Statt dass GRC-Teams aus zig Quellen Berichtstexte zusammenschreiben, kann KI Fakten konsolidieren und narrative, verständliche Reports erzeugen – inkl. Risikoargumentation und Status.

In dieser Rolle wird KI zum Policy- und Audit-Copiloten: Sie reduziert die formale Last, erhöht Konsistenz und schafft Kapazität für das, was eigentlich zählt – Risikosteuerung, Priorisierung und die wirksame Umsetzung von Kontrollen.

## 1.5 Neue Risiken und Angriffsflächen durch KI

Die Kehrseite ist banal, aber entscheidend: Jedes mächtige Werkzeug erweitert die Angriffsfläche. Sobald Sie KI in Prozesse, Datenflüsse und Entscheidungen integrieren, entsteht nicht nur »mehr Automatisierung«, sondern ein neues Systemverhalten – mit neuen Missbrauchsmöglichkeiten, neuen Failure-Modes und neuen Governance-Anforderungen. Und weil KI heute häufig als Cloud-Service, als API oder als tief integrierter Copilot genutzt wird, sind diese Risiken nicht theoretisch, sondern unmittelbar operativ relevant.

### 1.5.1 Angriffe auf und über LLMs

Mit LLMs entstehen Angriffsformen, die Sie aus klassischer AppSec oder klassischen SIEM-Modellen so nicht kennen – weil das Ziel nicht nur ist, ein System zu kompromittieren, sondern sein Verhalten zu manipulieren oder seine Grenzen zu umgehen.

Typische Muster sind:

- Prompt Injection:  
Hier wird Benutzereingabe so gestaltet, dass das Modell seine eigentlichen Vorgaben ignoriert: »Vergiss alle Regeln«, »Handle als Admin«, »Gib mir die vertraulichen Details«, oder subtiler: durch eingebettete Anweisungen in Dokumenten, E-Mails oder Webseiten, die das Modell als Kontext »mitliest«. Der Angriff ist nicht (nur) technisch, sondern psychologisch: Der Angreifer versucht, die Steuerlogik des Modells umzudefinieren.
- Data Leakage / unbeabsichtigte Datenpreisgabe:  
KI-Systeme arbeiten mit Kontextfenstern, Retrieval, Tool-Integrationen und oft mit riesigen Dokumentenbeständen. Wenn Zugriffsrechte, Kontextfilter oder Output-Constraints nicht sauber umgesetzt sind, können vertrauliche Inhalte in Antworten auftauchen – entweder, weil sie im Kontext enthalten waren, oder weil die KI über RAG »zu viel« zurückholt. Ein verwandtes Risiko ist falsche Nutzung von Trainings-/Fine-Tuning-Prozessen: Dann werden sensible Daten nicht nur angezeigt, sondern potenziell in Modelle oder Logs hinein konserviert.

- **Model Theft / Model Extraction:**

Wenn ein Modell über eine API zugänglich ist, kann ein Angreifer versuchen, durch systematisches Abfragen das Verhalten zu rekonstruieren: Welche Prompts führen zu welchen Outputs? Welche Policies sind implementiert? Welche »Guardrails« lassen sich umgehen? Ziel ist oft, ein funktional ähnliches Modell nachzubauen oder das Modellverhalten so gut zu verstehen, dass es leichter missbraucht werden kann.

- **Evasion und »Model Probing« gegen ML-Detektoren:**

ML-basierte Detektionssysteme (z.B. Malware-Klassifikatoren, UEBA, Anomalie-Detektoren) lassen sich gezielt austesten. Angreifer variieren Inputs, bis ein Score unter die Schwelle fällt, oder sie lernen, welche Features das Modell besonders stark gewichtet. Das ist eine neue Variante des alten Problems »Evasion«, aber durch KI wird die Iteration schneller und die Suche nach blinden Flecken systematischer.

Der Punkt ist: KI ist nicht nur ein neues Tool, sondern ein neues Angriffsziel – und in vielen Fällen auch ein neuer Angriffsweg, weil KI-Systeme häufig mit Datenquellen und Automationsrechten verbunden sind.

## 1.5.2 Governance- und Compliance-Risiken

Neben technischen Angriffen entstehen Governance-Risiken, die oft noch schneller zuschlagen – weil sie im Alltag beiläufig passieren: »mal eben« ein Prompt mit sensiblen Informationen, »kurz« ein Tool ausprobieren, »einfach« einen Copilot integrieren. Das ist selten böswillig, aber es ist riskant.

Typische Problemfelder sind:

- **Intransparente Datenflüsse in externe KI-Dienste:**

Sobald Mitarbeitende oder Tools Inhalte an externe KI-APIs senden, stellt sich die Frage: Welche Daten gehen raus? Wie werden sie verarbeitet? Werden sie geloggt? Wo liegen sie? Wer hat Zugriff? Ohne klare Architektur und Verträge ist das faktisch ein Blindflug.

- **Shadow AI ohne vertragliche Absicherung:**

Mitarbeitende nutzen GenAI-Tools, Browser-Plugins oder »kostenlose« Services, die nicht durch Einkauf, Legal und Security bewertet wurden. Ohne DPA, ohne Sicherheitszusagen, ohne klare Datenverwendungsbedingungen. Das ist das KI-Pendant zu Shadow IT – nur oft mit höherem Datenabflussrisiko, weil Inhalte aktiv hineinkopiert werden.

- **Datenschutz- und Geheimnisschutzverstöße:**

PII, HR-Daten, Kundendaten, interne Finanzen, Security-Incidents, Credentials in Tickets – all das landet schnell in Prompts, wenn keine Leitplanken existieren. Der Verstoß passiert nicht, weil jemand »Daten exfiltriert«, sondern weil Prozesse und Awareness nicht mit dem Tooling mithalten.

- **Fehlende Nachvollziehbarkeit gegenüber Revision und Aufsicht:**  
Wenn KI Outputs erzeugt, die in Entscheidungen einfließen (z.B. Risikoeinstufung, Incident-Schwere, Policy-Auslegung), muss nachvollziehbar sein, auf welcher Basis das passiert ist: Welche Quellen wurden genutzt? Welche Daten lagen zugrunde? Welche Version des Modells? Welche Prompt- und Policy-Konfiguration? Ohne diese Audit-Trails werden KI-gestützte Prozesse im Audit schnell zum Problem – selbst, wenn sie fachlich gut funktionieren.

Governance ist damit kein »nice to have«, sondern Voraussetzung, um KI überhaupt verantwortbar in Security und IT zu betreiben.

### 1.5.3 Organisatorische Risiken

Die unterschätztesten Risiken sind häufig nicht technisch, sondern menschlich und organisatorisch. KI verändert Rollen, Erwartungen und Verantwortungsgrenzen – und wenn das nicht aktiv gemanagt wird, entstehen gefährliche Grauzonen.

Typische organisatorische Risiken sind:

- **Überschätzung der KI-Fähigkeiten:**  
»Die KI wird's schon richten« ist eine der schnellsten Arten, Risiken zu erhöhen. LLMs wirken kompetent, weil sie flüssig formulieren. Das kann dazu führen, dass Teams Outputs zu wenig prüfen, Annahmen übernehmen oder Entscheidungen zu früh automatisieren.
- **Unklare Verantwortlichkeiten:**  
Wenn KI im SOC Tickets vorschlägt, Prioritäten setzt oder sogar Actions anstößt, muss klar sein: Wer trägt die Verantwortung? Wer hat freigegeben? Wer kontrolliert? Ohne klare RACI-Logik und Prozessdefinition entsteht »Verantwortungsdiffusion« – genau das, was Sie in Incident-Situationen nicht gebrauchen können.
- **Skill-Gap im Team:**  
KI-Ergebnisse sinnvoll zu nutzen, erfordert Kompetenz: Wie interpretiere ich Scores? Wie erkenne ich Halluzinationen? Wie teste ich Prompt-Robustheit? Wie bewerte ich Datenqualität und Drift? Wenn diese Skills fehlen, steigt die Wahrscheinlichkeit, dass KI entweder falsch genutzt oder aus Angst gar nicht genutzt wird.
- **Widerstände in der Belegschaft:**  
Wenn KI als Überwachungs- oder Kontrollinstrument wahrgenommen wird (»die KI liest unsere Chats«, »die KI bewertet unser Verhalten«), entstehen Akzeptanzprobleme, Ausweichverhalten und politische Konflikte. Gerade UEBA- und Kommunikationsanalysen brauchen deshalb klare Transparenz, Zweckbindung und Governance, sonst wird das Thema intern toxisch.

Unterm Strich: KI erweitert nicht nur Ihre technischen Möglichkeiten – sie verändert Ihre Organisation. Wer KI sicher nutzen will, muss sie nicht nur integrieren, sondern beherrschen: technisch, vertraglich, prozessual und kulturell.

## 1.6 Was sich für IT-Verantwortliche und CISOs konkret ändert

KI in der Security ist nicht nur ein neues Toolset, sondern ein Wechsel im Betriebsmodell. Das betrifft Entscheidungen, Verantwortlichkeiten, Skill-Profile und die Art, wie Sie Führung und Governance organisieren. Wer KI erfolgreich und sicher einsetzen will, muss nicht »mehr Regeln konfigurieren«, sondern Zielbilder, Leitplanken und messbare Steuerung etablieren – und das über mehrere Disziplinen hinweg.

### 1.6.1 Vom Regel-Admin zum Risiko-Architekten

In der klassischen Security-Welt war ein großer Teil der Steuerung operativ-technisch: Welche Regel setze ich wo? Welche Signatur ist aktiv? Welche Korrelation fehlt? Welche Firewall-Exception ist vertretbar? Das bleibt wichtig – aber KI verschiebt den Schwerpunkt.

Ihre Rolle wandert weg von der Frage »Welche Regel setze ich in welchem System?« hin zu einer deutlich strategischeren und gleichzeitig architekturlastigeren Frage »Welche Use Cases zahlen wirklich auf unser Sicherheitsziele-Konto ein – und wie steuern wir Daten, Modelle und Prozesse so, dass der Nutzen steigt, ohne das Risiko auszuweiten?«

Das ist ein Rollenwechsel vom »Regel-Administrator« hin zum Risiko-Architekten: Sie designen nicht nur Kontrollen, sondern ein System aus Datenflüssen, Modellen, Entscheidungslogik, Human-in-the-loop und Auditierbarkeit.

Dafür brauchen Sie ein KI-Sicherheitszielbild (Target Operating Model), das mindestens folgende Punkte beantwortet:

- Wo wollen wir KI bewusst einsetzen?
  - Welche Prozesse sind geeignet (z.B. Zusammenfassen, Priorisieren, Assistenz), wo ist der erwartete ROI hoch und das Fehlerrisiko beherrschbar?
- Wo wollen wir KI bewusst nicht einsetzen?
  - Welche Entscheidungen sind zu kritisch, zu rechtlich sensibel oder zu schwer auditierbar? Wo sind Fehlentscheidungen nicht tolerierbar?
- Welche Guardrails setzen wir – technisch und organisatorisch?
  - Zugriffskontrollen, Datenklassifizierung, Output-Constraints, Logging, Freigabeprozesse, Limits für autonome Aktionen, Red-Teaming, Incident-Handling für KI-Fehler.
- Wie messen wir Erfolg und Risiko?
  - Nicht nur »Nutzung« oder »Zeitersparnis«, sondern präzise Metriken: Qualität der Outputs, False-Positive/False-Negative-Raten, Drift-Indikatoren, Ticket-

Throughput, MTTR, Audit-Feststellungen, Datenabfluss-Events, Policy-Verstöße.

Das klingt nach »mehr Governance«, ist aber in Wahrheit die Voraussetzung, um KI überhaupt skalierbar und verantwortbar zu betreiben.

## 1.6.2 Skill-Shift im Security-Team

KI verändert auch, was ein gutes Security-Team können muss. Nicht jede Rolle muss Data Science beherrschen – aber die Grundkompetenz, KI-Systeme zu verstehen, zu bewerten und sauber zu betreiben, wird zur Baseline. Sonst entsteht entweder blindes Vertrauen (»die KI hat gesagt...«) oder totale Blockade (»das ist uns zu riskant«).

Neue bzw. stärker gewichtete Skills sind unter anderem:

- Verständnis von ML-/LLM-Konzepten auf Architekturebene:  
Was ist Training vs. Inferenz? Was bedeutet Drift? Was ist RAG? Wo liegen typische Failure-Modes (Halluzination, Bias, Prompt Injection, Over-retrieval)?
- Evaluierung von Modellen und Outputs:  
Wie messe ich Genauigkeit, Robustheit und Bias? Wie teste ich Worst-Case-Szenarien? Wie überprüfe ich, ob ein System »zu selbstbewusst« falsche Dinge behauptet? Wie stelle ich reproduzierbare Qualität sicher?
- Aufbau und Betrieb von RAG/LLM-Anwendungen:  
Datenquellen anbinden, Zugriff sauber durchsetzen, Retrieval begrenzen, Quellen zitierbar machen, Logging/Audit implementieren, Versionierung und Change-Control etablieren.
- Kritische Einordnung von KI-Ergebnissen:  
Fähigkeit, Aussagen zu verifizieren, Annahmen zu erkennen, Unsicherheit zu akzeptieren und KI-Outputs in Entscheidungen zu übersetzen, ohne Verantwortung abzugeben.

Der Kern ist eine neue Form von Grundbildung im Security-Kontext: KI-Lese- und Schreibfähigkeit. Das heißt:

- verstehen, was ein Modell tut,
- wissen, was es nicht kann,
- einschätzen, wann es hilft,
- und testen können, ob es in Ihrem Kontext zuverlässig ist.

Teams, die diese Kompetenz aufbauen, können KI kontrolliert nutzen. Teams, die sie nicht aufbauen, werden entweder riskant über-automatisieren oder relevante Chancen liegen lassen.

# Stichwortverzeichnis

## A

ABAC 118  
ACL-Bypass 331  
Agenten 61, 115  
Agentische Orchestrierung 176  
Aktion 142, 144  
ALPHV (BlackCat) 53  
Angreiferprofile 328  
Angriffsziele 329  
Anti-Patterns 226, 251  
APT 54  
Architektur 213, 243  
Architekturprinzipien 107  
Asset 142, 143, 326  
Auditability 152  
Auditierbarkeit 105, 274  
Augment, don't replace 39  
Automation Bias 279  
Automatisierte SOC-Lageberichte 199

## B

Bedrohungsmodell 326  
Betrieb & Messung 225, 250  
Betriebsfähigkeit 295  
Betriebsmodelle 186  
Betriebsrisiko 271  
Bias 277  
Blindes Vertrauen 43  
Budget-Buckets 313  
Build vs. Buy 126, 314

## C

Canonical Schema 154  
Capability Map 308  
CEO-Fraud 71  
Change Management 123, 302  
Chunking 157  
ClOp 53  
Cloudbetrieb 189

Cloud Control Plane Logs 147  
Compliance-Metriken 287  
Conditional-Access / Policy Decisions  
146  
Controls-by-Design 282  
Cybercrime-as-a-Service 56  
Cybercrime-Gruppen 52

## D

Data as a Product 162  
Data Classification 125  
Data Exfiltration 331  
Data Health 303  
Data Leakage 32  
Data Poisoning 191, 332  
Data Product Owner 160  
Data Quality Monitoring 156  
Daten- und Index-Policy 267  
Datenabfluss 329  
Datenexfiltration 190  
Datenminimierung 261  
Datenqualität 97, 138  
Datenresidenz 276  
Datenrisiko 270  
Deepfakes 60, 68  
Deep Learning 18, 83  
Denial-of-Wallet 333  
Desinformation 66  
Detection Bias 278  
Detection Engineering 308  
DevSecOps 30  
DNS 147  
Dokumente 140  
Drift 100

## E

Embeddings 84, 87, 169  
Endpoint-Telemetrie 146  
Enrichment 155

Enterprise AI Platform + Security Overlay 301  
 Erfahrungswissen 140  
 EU AI Act 37, 75  
 Evals 345  
 Events 140  
 Evidence-first 107, 261, 335  
 Evidence-first-Prinzip 94  
 Evidence Packs 276  
 Exposure Bias 278

**F**

Fail safely 336  
 Fairness 262, 277  
 Federated Enablement 299

**G**

GenAI 20, 24  
 Generative KI 18, 84  
 Governance-Guardrails 184  
 Governance-Ziel 261  
 GRC 31  
 GRC-Assistent 227  
 Greatness 54  
 Ground Truth 147  
 Guardrails 35, 112, 182

**H**

Halluzinationen 28  
 Hybridbetrieb 188

**I**

Identität 142, 143  
 Identität & Zugriff 339  
 Identity- und Access-Telemetrie 146  
 Identity & Access 118  
 Identity Provider-Logs 146  
 Incident Response 303  
 Indizes 110  
 Ingestion und Normalisierung 110  
 Input-Guardrails 182  
 Integrität 327  
 Intelligence Layer 111  
 IR-Playbooks 149  
 ISO 27001 31, 37  
 ISO/IEC 27001 75

**K**

KI-gestützter BEC 71  
 KI-Governance-Board 264  
 KI-Incidents 288  
 KI-Nutzungsrichtlinie 267  
 KI-Plattformbetrieb 265  
 KI-Risk & Compliance Council 265  
 Knowledge 140  
     kuratierbar 141  
     versionierbar 141  
 Knowledge Engineering 308  
 Knowledge Owner 160  
 Kontextfenster 84, 86, 168  
 Korrelation 138, 142  
 Kosten- und Kapazitätsmanagement 303  
 Künstliche Intelligenz 17

**L**

Label Bias 278  
 Labels 99  
 Large Language Models 18  
 Least Privilege 262, 336  
 Lessons Learned 141  
 Lifecycle-Management 44  
 LLM 20 24, 84  
 LLM-only 113, 172  
 LockBit 53  
 Lumma 54

**M**

Machine Learning 18, 19, 82  
 Malware-as-a-Service 54  
 Metadaten 158  
 MFA 339  
 Minimalismus bei Daten 41  
 Model-Behavior-Risiken 193  
 Modellrisiko 271  
 Model Probing 33  
 Model Theft 33

**N**

Nachvollziehbarkeit 261, 327  
 Netzwerk 147  
 Netzwerksegmentierung 119  
 NIS2 37, 75  
 No Grounding – No Answer 40

Normalisierung 154  
Normen und Frameworks 361

**O**

Observability 113  
On-Prem 187  
Orchestrierung 112  
Output-Formate 180  
Output-Guardrails 182  
Ownership 295

**P**

Phishing 27, 63  
Phishing-as-a-Service 54  
PII-Minimierung 152  
Platform Ops 161  
Policy Enforcement 112  
Policy Engine 124  
Policy Layer 179  
Postmortems 141  
Privileged Access Events 146  
Privilege Escalation 192, 332  
Prompt- und Modell-Change-Policy 269  
Prompt- und Retrieval-Security 120  
Prompting 178  
Prompt Injection 32, 190, 330  
Proxy 147  
Purpose Limitation 261

**Q**

Quellenpflicht 159

**R**

RAG 91, 148, 173  
Ransomware 72  
Ransomware-Gangs 52  
RBAC 118  
Red Teaming 346  
Referenzarchitektur 335  
Referenzstruktur 197  
Regeln 140  
Reproduzierbarkeit 123  
Reputations- und Haftungsrisiko 271  
Resilienz 105  
Retention 150, 152  
Retrieval-Augmented Generation 114  
Retrieval- und Tool-Guardrails 183

Risiko- und Kontrollmetriken 286  
Rollenmodell 160  
Runbooks 141, 149

**S**

Safety-by-Design 42  
Scam-Netzwerke 53  
Secrets & Keys 119  
Secure Prompt Engineering 347  
Security-by-Design 42  
Security AI Platform Team 300  
Security Data Engineering 308  
Sensitive Data Leakage 333  
Service Catalog 149  
Shadow AI 33, 43  
Sicherheitsanforderungen 339  
Sicherheitsarchitektur 105  
Sicherheitsprinzipien 335  
Signalquellen 145  
Signals 140  
Signals & Knowledge 109  
Skalierung 105  
Skill-Shift 36  
SOC 28  
Social Engineering 27  
Spear-Phishing 63  
SSO 339  
Supply-Chain-Risiko 271  
Supply Chain Security 120  
Systemprompt 178

**T**

Telemetrie 101, 140, 145  
Third-Party & Vendor Governance 280  
Threat-Modelings 73  
Token 84, 85, 167  
Token-/Session 146  
Token Budgeting 181  
Tool-/Action-Risiko 271  
Tool-Access-Policy 268  
Tool-augmented LLM 174  
Tool/Function Calling 84, 89, 114, 171  
Tool Misuse 192, 332  
Transparenz 262, 279

**U**

Userprompt 179

**V**

Verfügbarkeit 327

Versionierung 123

Vertraulichkeit 326

Vulnerability & Patch-Priorisierung 253

**W**

Wertmetriken 285

Wissensbasis 231

Wissensqualität 138

**Z**

Zeit 142, 144

Zero Trust 108

Zugriff 150

Zugriffskontrolle 151, 340

Zustände 140

Zweckbindung 150