



# AI-Driven Data Onboarding

17.03.2026 | Frankfurt

DB System

# Inhaltsverzeichnis.



- 1. Herausforderungen**
- 2. Idee**
- 3. Demo**
- 4. Vision**

# Wer sind wir?

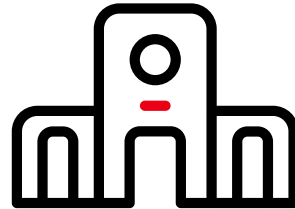


- Wir sind ein agiles DevOps-Team in der DB Systel GmbH, die eine 100 %-ige Tochtergesellschaft der Deutsche Bahn AG ist.
- Wir betreiben eine zentrale Splunk-Umgebung für den gesamten Konzern und stellen diese als Plattform für Geschäftsfelder, Unternehmen und Teams mandantengetrennt zur Verfügung.
- Unsere stetig wachsende Umgebung umfasst aktuell 60 Indexer in einem MultiSite-Cluster und mehrere SearchHead Cluster.
- Wie bedienen Operation und Security UseCases.



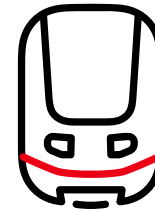
## Schienennetz

ca. 33.000km Strecke  
ca. 64.000 Weichen  
ca. 3.800 Signalanlagen



## Bahnhöfe

ca. 5.400 Bahnhöfe



## Personenverkehr

ca. 410 ICEs  
ca. 4.000 Triebzüge  
ca. 6.000 Busse



## Güterverkehr

ca. 2.500 Lokomotiven  
ca. 65.000 Güterwagen

Quelle: DB Konzern Geschäftsbericht 2024



# Herausforderungen

---

# Herausforderungen ans Logging



Ein großer Konzern. **5 Geschäftsfelder** mit über **250 Unternehmen** in sehr vielen **Standorten** und mit sehr vielen **Fahrzeugen**.



Wir wollen **Security** und **Operations** abdecken, auf **System-** und **Anwendungsebene**.



Viele Systeme. Über **3.000 Anwendungen** im Konzern. **On-prem**, Private **Clouds**, **Software-as-a-Service**-Lösungen und **Anlagen**.

# Kann uns Splunk hier helfen?



## Datenspeicherung



- Index Cluster
- Multi Site
- Skaliert horizontal

## Analyse und Zugriffe



- SearchHead Cluster
- Multiple parallele Installationen
- Skaliert horizontal

## Datenerfassung



- Splunk Universal Forwarder für die gängigen OS-Systeme
- Zentrale Konfiguration
- Skaliert horizontal

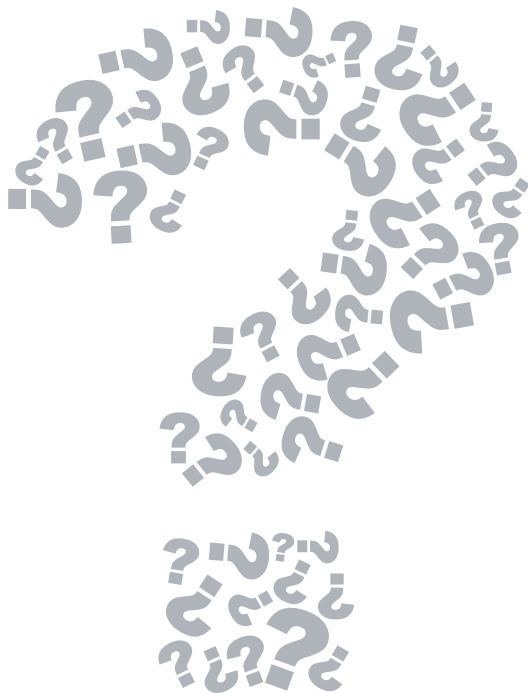
## Datenarten



- Auslesen von Dateien
- HTTP Event Collector
- openTelemetry OTEL
- ...

Also einschalten, aktivieren und fertig?

# Diversität der Daten



**OS-Log:** Windows, Linux, Unix, Netzwerk. Relativ standardisiert

**openTelemetry:** standardisiert

**SaaS:** Wir sind nicht die Ersten, hier existieren bereits Lösungen

**Eingekaufte Anwendungen:** für die bekanntesten existieren Lösungen, aber für viele nicht

**DB eigenentwickelte Anwendungen:** durch die Größe des Konzerns, die Historie keine Logging-Standardisierung und fast alles individuell

# Was bedeutet das für uns?

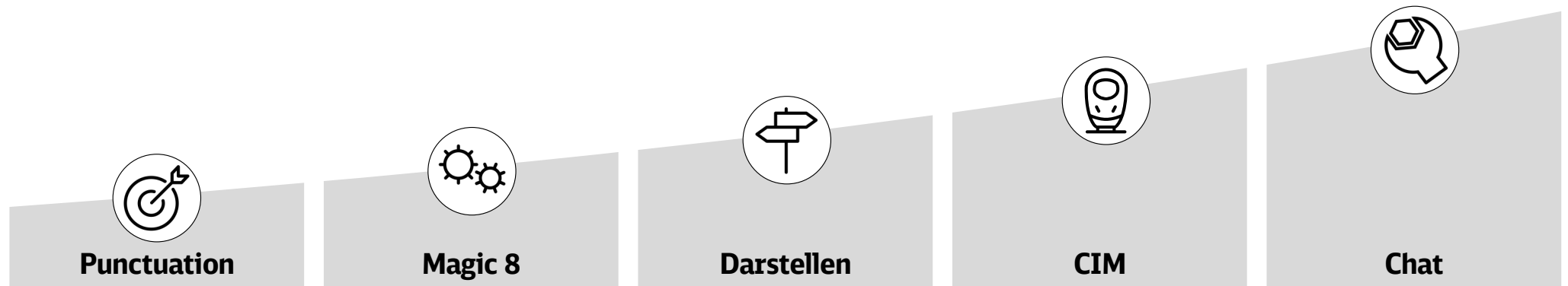




# Idee



# Unsere Reise in 5 Schritten



- Erster Ansatz war zu erkennen, ob das Format schon bekannt ist. (Punct/Sourcetype)
- Doppelte Configs sollten verhindert werden.

- Danach kam die Idee GPT bzgl. der Magic 8 / Great 8 für die props.conf zu fragen.
- Das lieferte erstaunlich gute Ergebnisse, trotz recht einfachem Prompt.

- Ergebnisse sollten nicht nur in einer Konfiguration erscheinen, wir wollten diese auch direkt anhand unseres Logsamples testen und darstellen.

- Wenn Magic 8 schon funktioniert, warum sollte dann nicht auch CIM-Zuordnung und Feldextraktion funktionieren.

- Und letztendlich ist man, wenn man dann länger mit GPT arbeitet, Chats gewohnt.
- Also warum nicht mit der KI für besser Ergebnisse chatten.



# Demo




# Erkennung: Existiert das Format schon?



- **Doppelte Arbeit vermeiden & Wissen heben.**
- Vergleich der hochgeladenen Logs über Punctuation-/Struktur-Signaturen mit bereits indexierten Mustern.
- Output-Lookup (KV) hält: PUNCT, Sourcetype, zugehörige props-Konfiguration.
- Aktuelle Limitation: Multiline-Events nur teilweise erkannt (zeilenorientierter Ansatz).
- Ausblick: Hybrid Matching (Embedding + Struktur + PUNCT) für höhere Trefferquote & Multiline-Stabilität.
- Nutzen: Sofortige Transparenz – Anwender (auch ohne weitreichende Splunk-Rechte) sehen: Quelle schon vorhanden? Wiederverwendbar?

**Bekannter Sourcetype:**

Es werden die bekannten Sourcetypes mit dem Verhältnis der übereinstimmenden Formate (puncts) ausgegeben:

	pfsense (100%)
---	----------------

[Sourcetype übernehmen](#)

# Event Breaking & Zeitstempel-Validierung



- **Saubere Basis vor jeder weiteren Investition.**
- Live-Vorschau nach Upload: korrekt gebrochene Events, Timestamp-Validität, Lücken/Anomalien.
- KI-gestützte Vorschläge der Great8 reduziert auf: LINE\_BREAKER, TIME\_PREFIX, TIME\_FORMAT, MAX\_TIMESTAMP\_LOOKAHEAD.
- Iterative Verfeinerung via Chat statt Try-&-Error in Config-Files.
- Nutzen: Frühe Qualitätssicherung → Fehlersuche später entfällt, Onboarding-Zeit sinkt deutlich.

**AI Assistant:**

Ich bestätige, dass keine personenbezogenen Informationen im Sample Log enthalten sind.

Parsing des Samplelogs durch AI durchführen...

---

**Splunk Settings:**

Line Breaker\*:  Time Prefix\*:

Truncate\*:  Time Format\*:

Max Timestamp Lookahead\*:

Apply Regex

---

**Preview Events:**

Event Filter:

Fehlerhaftes Datum (0) Abgeschnitten (0) Gültiges Datum (2) Max Lookahead Warning (0)

Event Übersicht: 2 von 2 Events angezeigt  
2 Gültige Events

Count	Date	Event
1	2025-02-28T23:02:39+01:00 (1740780159)	[2025-02-28T23:02:39+01:00] fw2.ov3r.local filterlog 30654 - [meta sequenceid="45261"] 122,,bc1dfdedfeddb0c40ef8ef4dab01d6e5.vtnet2.match.pass.in.4.0x0,,64.51070.0.DF,17,udp,63,192.168.0.34
2	2025-02-28T23:02:39+01:00 (1740780159)	[2025-02-28T23:02:39+01:00] fw1.ov3r.local filterlog 62920 - [meta sequenceid="427523"] 123,,bc1dfdedfeddb0c40ef8ef4dab01d6e5.vlan010.match.pass.in.4.0x0,,64.51070.0.DF,17,udp,63,192.168.0.34

Events übernehmen

# Feldextraktion & Normalisierung



- Schnell zu korrelationsfähigen Daten.
- KI erkennt Kandidaten-Felder, zeigt Coverage / Lücken transparent.
- Vorläufiges CIM-Mapping (Standardisierung & spätere Korrelation beschleunigt).
- Erklärbare Regex-Generierung statt individueller Skriptfragmente.
- Nutzen: Weniger manuelle Schleifen, konsistente Ergebnisse, Wissensaufbau strukturiert.

## Transforms.conf Definition

```
[extract_firewall_logs]
REGEX = (?<timestamp>[0-9]{4}-[0-9]{2}-[0-9]{2}T[0-9]{2}:[0-9]{2}:[0-9]{2}[+-][0-9]{2}:[0-9]{2})\s+(?<hostname>[a-z_]+)\s+(?<log_type>[a-z_]+)\s+(?<log_id>[0-9]{1,4})\s+(?<meta'sequenceId'>{<br><sequence_id>[0-9]{1,4})\s+(?<rule>[a-z_]+)\s+(?<tracker_id>[a-f0-9]{32})\s+(?<interface>[a-z_]+)\s+(?<reason>[a-z_]+)\s+(?<action>[a-z_]+)\s+(?<direction>[a-z_]+)\s+(?<ip_version>[0-9]{1,4})\s+(?<tos>[0-9]{1,4})\s+(?<ecn>[0-9]{1,4})\s+(?<ttl>[0-9]{1,4})\s+(?<id>[0-9]{1,4})\s+(?<offset>[0-9]{1,4})\s+(?<flags>[a-z_]+)\s+(?<protocol_id>[0-9]{1,4})\s+(?<protocol>[a-z_]+)\s+(?<length>[0-9]{1,4})\s+(?<src_ip>[0-9]{1,3}\.[0-9]{1,3}\.[0-9]{1,3})\s+(?<dst_ip>[0-9]{1,3}\.[0-9]{1,3}\.[0-9]{1,3})\s+(?<src_port>[0-9]{1,4})\s+(?<dst_port>[0-9]{1,4})\s+(?<payload_length>[0-9]{1,4})
```

## Erkannte Felder

Feldname	Gefunden in (%)	Top 5 Werte	CIM Mapping
timestamp	100%	• 2025-02-20T17:54:03+01:00 (18) • 2026-02-20T17:54:03+01:00 (2)	timestamp
hostname	100%	• fw1.ov3r.local (20)	host
log_type	100%	• filterlog (20)	log_type
log_id	100%	• 34049 (20)	log_id
sequence_id	100%	• 12711029 (18) • 12711028 (2)	sequence_id
rule	100%	• 83 (18) • 143 (2)	rule
tracker_id	100%	• fae559338f65e11c53669fc3642c93c2 (18) • 1b5a7e5825cd7208895e6ab9a1260a71 (2)	tracker_id
interface	100%	• vian010 (18) • vian030 (2)	interface
reason	100%	• match (20)	reason
action	100%	• pass (20)	action
direction	100%	• out (18) • in (2)	direction
ip_version	100%	• 4 (20)	ip_version
tos	100%	• 0x0 (20)	tos
ecn	0%	Keine Werte gefunden	-
ttl	100%	• 127 (18) • 128 (2)	ttl
id	100%	• 51561 (20)	id
offset	100%	• 0 (20)	offset
flags	100%	• none (20)	flags
protocol_id	100%	• 17 (20)	protocol_id
protocol	100%	• udp (20)	protocol
length	100%	• 60 (20)	length
src_ip	100%	• 192.168.3.201 (20)	src_ip
dst_ip	100%	• 192.168.0.200 (20)	dst_ip
src_port	100%	• 65089 (20)	src_port
dst_port	100%	• 53 (20)	dst_port
payload_length	100%	• 40 (20)	payload_length

## Events mit markierten Feldern

Zeige Events 1 bis 10 von 194

```
Event 1
2466:"title="offset">025-02-20T17:54:03+01:00 fw1.ov3r.local filterlog 34049 - [meta sequenceId="12711029"]
83,,fae559338f65e11c53669fc3642c93c2,vian010,match,pass,out,4,0x0,,127,51561,0,none,17,udp,60,192.168.3.201,192.168.0.200,65089,53,40

Event 2
2466:"title="offset">026-02-20T17:54:03+01:00 fw1.ov3r.local filterlog 34049 - [meta sequenceId="12711028"]
143,,1b5a7e5825cd7208895e6ab9a1260a71,vian030,match,pass,in,4,0x0,,128,51561,0,none,17,udp,60,192.168.3.201,192.168.0.200,65089,53,40
```

# App Enablement & Release-Fähigkeit



- **Vom Roh-Log zur produktionsnahen TA in Sekunden.**
- Automatische Erstellung einer Splunk TA (props.conf, transforms.conf, app.conf, metadata, README).
- Namens- & Versionierungslogik gemäß Konvention → direkt CI/CD-fähig.
- Sofort einsetzbar für Test, Security & Compliance-Gates.
- Nutzen: Time-to-Value stark reduziert, standardisierte Basis für Deployment & Governance.

**Generierte App-Informationen**

App-Name: TA-generic-firewall-v100

App-Titel: Generic Firewall Log Analysis

Beschreibung: Diese App ermöglicht die Analyse von generischen Firewall-Logs, indem sie wichtige Felder wie Aktion, Richtung, IP-Adressen und Ports extrahiert.

Author: OSSP Onboarding Tool

Version: 1.0.0

---

**Konfigurationsdateien Vorschau**

**transforms.conf**

```
[extract_firewall_logs]
REGEX = (?<timestamp>[0-9]{4}-[0-9]{2}-[0-9]{2}T[0-9]{2}:[0-9]{2}:[0-9]{2}[+-][0-9]{2}:[0-9]{2})\|(?<hostname>[a-z]{1,64})\|(?<log_app>[a-z]{1,64})\|(?<log_ip_src>[0-9]{1,3}\.|[0-9]{1,3}\.[0-9]{1,3}\.|[0-9]{1,3}\.[0-9]{1,3}\.[0-9]{1,3})\|(?<log_ip_dst>[0-9]{1,3}\.|[0-9]{1,3}\.[0-9]{1,3}\.|[0-9]{1,3}\.[0-9]{1,3}\.[0-9]{1,3})\|(?<log_protocol>[a-z]{1,64})\|(?<log_port_src>[0-9]{1,5})\|(?<log_port_dst>[0-9]{1,5})\|(?<log_action>[a-z]{1,64})\|(?<log_direction>[a-z]{1,64})\|(?<log_message>[a-z]{1,64})\|(?<log_length_src>[0-9]{1,10})\|(?<log_length_dst>[0-9]{1,10})
```

**props.conf**

```
[source::firewall_logs]
LINE_BREAKER = ([\n\r]+)
MAX_TIMESTAMP_LOOKAHEAD = 22
TIME_FORMAT = %Y-%m-%dT%H:%M:%S.%fZ
TIME_PREFIX = ^
TRANSFORMS = extract_firewall_logs
```

**app.conf**

```
# Splunk app configuration file
#
[install]
is_configured = 0

[ui]
is_visible = 0
label = Generic Firewall Log Analysis
```

**default.meta**

```
[props]
access = read : [ * ], write : [ admin ]
export = system

[transforms]
access = read : [ * ], write : [ admin ]
export = system
```

**Splunk App Download**

Die App wird als ZIP-Archiv mit allen notwendigen Konfigurationsdateien erstellt.

[Splunk App herunterladen](#)



# Vision



# Automatisiertes Datenonboarding ohne Splunk Know-how. Volle Transparenz inklusive Selbstheilung.



Copyright: Deutsche Bahn AG / v.l. Dominic Dupont / Skydeck / Gorodenkoff Productions

**Anwendungsverantwortliche** bestellen für ihre Anwendungen ein Data-Onboarding.

Die **zuständigen Betriebsteams** werden automatisch ermittelt und informiert.

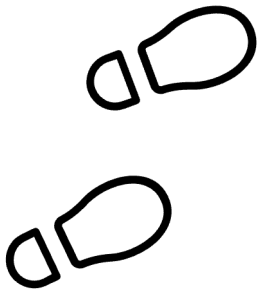
Nach Angabe der notwendigen Logfiles, werden diese ohne Zeitverzögerung und ohne Splunk-Know-how **analysiert** und **konfiguriert**.

**Neue** Sourcetypes fließen in ein **eigenes ML-Model** ein, um **Veränderungen** zukünftig zu **bemerken** und dem zuständigen Betriebsteam **avisieren** zu können.

Eine erneute **Analyse** mit anschließender **Rekonfiguration** kann so ganz leicht wieder angestoßen werden.

Jederzeit können **alle** Beteiligten sehen, welche **Datenquellen** zu welcher **ihrer Anwendungen** gehören, wie der **Status** des Onboardings ist und in welcher **Qualität** die Events ins Splunk-System laufen.

## Next steps



**GIT-Anbindung:** Automatisches Deployment über Pipelines mit Vier-Augen-Checks.

**Umgang mit Strukturierte Daten verbessern:** Strukturierte Daten (XML/CSV) sind aktuell noch schwierig.

**Nutzung einer lokalen KI/LLM:** Datenschutz ist immer ein großes Thema,

**Nutzung von RAG/MCP:** Testen, ob wir bessere Ergebnisse zu bereits existierenden Sourcetypes bekommen. Punctuation hat seine Grenzen.

**Bahninterne Informationen:** Event-Tagging und Konfigerweiterungen mit Anwendungs-IDs, IP-Netze, Zugnummern, etc.

**Vielen Dank**

