

Verarbeiten von Gesundheitsdaten - Regulatorische Anforderungen und technische Vorschläge @ GWDG

Hendrik Nolte
hendrik.nolte@gwdg.de

Sebastian Krey

Lars Quentin

Robin Strahl

Julian Kunkel

9. September 2025

hpc@gwdg.de

GWDG – Gesellschaft für wissenschaftliche Datenverarbeitung mbH Göttingen

- EU Data Governance Act
 - Ziel wiederverwendung von (öffentlichen) Daten
 - Sowohl personen- als auch nicht personenbezogen
 - Definiert Datenmittler als neutrale Dritte
 - Definiert insbesondere eine „sichere Verarbeitungsumgebung“ (Art. 5)
- EU Data Act
 - Ergänzt den DGA bzgl. Zugang und Nutzung der Daten
 - Anwendung z.B. bei IoT Geräten
 - Sichert Interoperabilität
 - Rechte von Betroffene
 - Definiert insbesondere „Datenräume“ (Art. 44)



Regulatorische Aspekte
ausgearbeitet von Robin
Strahl

Ziel: Erleichterung des Zugangs zu elektronischen Gesundheitsdaten für die Primär- und Sekundärnutzung

- Ist die erste Instanziierung eines Datenraumes
- präzisiert und ergänzt die DSGVO [...] in Bezug [auf] [...] elektronische Gesundheitsdaten (Art. 1 Abs. 2 lit. a EHDS)
- Auf der Grundlage sollen technische und organisatorische Standards entstehen
 - Systeme interoperabler elektronischer Gesundheits- bzw. Patientenakten (*Electronic Health Records*; EHR bzw. EHR-Systeme)
 - MyHealth@EU
 - HealthData@EU (und entsprechende Zugangsstellen)

- Primärnutzung (Art. 3 bis 24 EHDS)
 - Verarbeitung von Gesundheitsdaten für die Gesundheitsversorgung
- Sekundärnutzung (Art. 50 bis 81 EHDS)
 - Verarbeitung der Gesundheitsdaten für in VO genannte Zwecke
 - u.a. wissenschaftliche Forschung im Bereich des Gesundheitswesens oder des Pflegesektors (Art. 53 Abs. 1 lit. e EHDS)
 - erfolgt in sicheren Verarbeitungsräumen

- personenbezogene elektronische Gesundheitsdaten
 - [elektronische] Gesundheitsdaten, sowie genetische Daten
 - Art. 6 DSGVO (Rechtmäßigkeit der Verarbeitung)
 - Art. 9 DSGVO (Verarbeitung besonderer Kategorien personenbezogener Daten)
- nicht personenbezogene elektronische Gesundheitsdaten
 - anonymisierte Gesundheitsdaten (kein Personenbezug)
 - Daten, die sich nie auf eine Person bezogen haben

Zugangsstelle für Gesundheitsdaten (Art. 55 EHDS)

- werden vom Mitgliedsstaat benannt
- entscheidet über Zugang zu Daten

Aufgaben und Pflichten in Art. 57 ff. EHDS:

- Verarbeiten der Gesundheitsdaten
(im Rahmen der Aufgaben und Pflichten der EHDS VO)
- Bereitstellen der Gesundheitsdaten in sicherer Verarbeitungsumgebung

Gesundheitsdateninhaber (Art. 2 Abs. 2 lit. t EHDS)

- Stellung als (gemeinsam) Verantwortlicher bei Gesundheitsdatenverarbeitung von besonderer Bedeutung (Art. 2 Abs. 2 lit. t sublit. i EHDS)
- Fähigkeit, durch Kontrolle über technische Konzeption eines Produkts/Dienstes, Daten zur Verfügung zu stellen (Art. 2 Abs. 2 lit. t sublit. ii EHDS)

Pflicht:

- Gesundheitsdateninhaber *müssen* der Zugangsstelle elektronische Gesundheitsdaten nach Art. 60, 51 Abs. 1 EHDS für Mindestkategorien bereitstellen

nat. Personen (einschließlich einzelner Forscher) und Kleinstunternehmen sind von dieser Pflicht befreit (Art. 50 Abs. 1 EHDS)

vertrauenswürdiger Gesundheitsdateninhaber (Art. 72 EHDS)

- Benennung als vertrauenswürdiger Gesundheitsdateninhaber durch Zugangsstelle für Gesundheitsdaten
- Bedingungen
 - sichere Verarbeitungsumgebung
 - Fachwissen um Anträge zu prüfen
 - erforderliche Garantien
- zusätzliche Aufgaben/Pflichten:
 - Zugangsstelle kann Anträge auf Zugang direkt weiterleiten
 - Prüfung des Antrags und Übermittlung eines Vorschlags für die Entscheidung
 - Gewährt Zugang in (eigener) sicheren Verarbeitungsumgebung

Gesundheitsdatennutzer (Art. 2 Abs. 2 lit. u EHDS)

- nat. o. jur. Personen, die rechtmäßig Zugang zu elektronischen Gesundheitsdaten erhalten haben

- grds. gilt Datenminimierung und Zweckbindung auch nach EHDS (Art. 66 EHDS)
- Antrag auf Zugang zu Gesundheitsdaten nach Art. 67 EHDS
- Zugangsstelle prüft, ob Kriterien für Datengenehmigung nach Art. 68 EHDS erfüllt sind
 - **JA:** Zugangsstelle gewährt Zugang zu pseudonymisierten Daten (Art. 66 Abs. 3 EHDS) in sicherer Verarbeitungsumgebung
 - **NEIN:** kein Zugang oder nur über Gesundheitsdatenanfrage
- Gesundheitsdatenanfrage (Art. 69 EHDS)
 - ausschließlich anonymisiertes statistisches Format

Art. 2 Nr. 20 DGA:

„sichere Verarbeitungsumgebung“ sind physische oder virtuelle Umgebung und die organisatorischen Mittel, mit denen die Einhaltung der Anforderungen des Unionsrechts, wie der [DSGVO], insbesondere im Hinblick auf die Rechte der betroffenen Personen, der Rechte des geistigen Eigentums und der geschäftlichen und statistischen Vertraulichkeit, der Integrität und der Verfügbarkeit, sowie des geltenden Unionsrechts und des nationalen Rechts gewährleistet wird und die es der Einrichtung, die die sichere Verarbeitungsumgebung bereitstellt, ermöglichen, alle Datenverarbeitungsvorgänge zu bestimmen und zu beaufsichtigen, darunter auch das Anzeigen, Speichern, Herunterladen und Exportieren von Daten und das Berechnen abgeleiteter Daten mithilfe von Rechenalgorithmen;

Art. 2 Nr. 20 DGA:

1. physische oder virtuelle Umgebung **und** organisatorische Mittel
2. Einhaltung des Unionsrechts (inkl. DSGVO, etc.) und nationalen Rechts
3. Bestimmung und Aufsicht über alle Datenverarbeitungsvorgänge

Sicherheitsmaßnahmen (Art. 73 EHDS)

- Begrenzung des Zugangs
- Minimieren des Risikos unbefugten Verarbeitens
- Beschränkung der Eingabe
- Sicherstellung, dass Gesundheitsdatennutzer ausschließlich auf die von ihrer Datengenehmigung erfassten Daten zugreifen können
- Führung identifizierbarer Protokolle
- Sicherstellung der Befolgung und Überwachung dieser Sicherheitsmaßnahmen

Aber: Bis zum 26. März 2027 legt die Kommission technische und organisatorische Anforderungen sowie die Anforderungen an die Informationssicherheit, die Vertraulichkeit, den Datenschutz und die Interoperabilität fest.

Art. 74 EHDS

	EHDS	DSGVO
Bereitstellung	Gesundheitsdateninhaber	Verantwortlicher
Verarbeitung (i.R.d. Aufgaben der EHDS)	Zugangsstelle	Verantwortliche
Verarbeitung in TRE	Zugangsstelle Gesundheitsdateninhaber	Auftragsverarbeiter Verantwortlicher

Die EHDS wird in DE durch das Gesundheitsdatennutzungsgesetz (GDNG) auf nat. Ebene umgesetzt und in Teilen ergänzt

- Bundesinstitut für Arzneimittel und Medizinprodukte (BfArM)
 - wird Datenzugangsstelle **und** Datenkoordinierungsstelle
- Das BfArM bekommt im wesentlichen die Aufgaben aus der EHDS zugeteilt
- zusätzlich kommt das Erstellen von Konzepten ...
 - zur Nutzung von sicheren Verarbeitungsumgebungen
 - zur Weiterentwicklung der zentralen Datenzugangs- und Koordinierungsstelle
 - zur Verknüpfung und gemeinsamen Verarbeitung von pseudonymisierten Gesundheitsdaten verschiedener datenhaltender Stellen

- Viele Anforderungen müssen erst noch gemacht werden
- Stichtag 26. März 2027
 - Anforderungen zur Erfassung personenbezogener elektronischer Gesundheitsdaten (Art. 13)
 - europäisches Austauschformat für EHR (Art. 15)
 - Maßnahmen zur Entwicklung von MyHealth@EU (Art. 23)
 - Spezifikationen von EHR-Systemen (Art. 36)
 - Muster für Antrag auf Zugang zu Gesundheitsdaten (Art. 70)
 - technische und organisatorische Anforderungen an sicherer Verarbeitungsumgebungen (Art. 73)
- Vorbild sind Länder wie Finnland, die eine ähnliche Struktur bereits nutzen

- Plötzliches Aufkommen der „sicheren Verarbeitungsumgebung“ und EHR
- Weckt(e) Bestrebungen zu einer einheitlichen, kollaborativen Plattform
 - in Anlehnung an die UK-Biobank
- Weckt(e) Verunsicherung/Neugierde bzgl. der technischen und organisatorischen Unklarheiten
 - Prototyp von GWDG mit 2 DIZen



Technische Aspekte
ausgearbeitet von Lars
Quentin

- Entwicklung eines eigenen TREs
- Möglichst Bare-Metal auf HPC
- Anwendungsfall: Medizinische Daten
- Vorbild: 5-Safe Framework
 - Initial UK ONS, HDR UK Green Paper (1 2)
- Zugriff: Linux VDI via Browser
- Verschlüsselung + Key Management
- Multi-User Support
- Balance: Sicherheit ggü Produktivität

5 Safe

Safe Data: Daten wurden pre-ingest pseudonymisiert.

Safe Projects: Nur zugelassene Projekte dürfen Datenzugriff kriegen.

Safe People:

- KYC (Person + Institution)
- Nutzungsvertrag
- Absolviertes Training

Safe Settings: Sichere Datenumgebung

Safe Outputs: Pre-Publikation
Statistical Disclosure Control (SDC).

The screenshot shows a Jupyter Remote Desktop interface. At the top, there's a browser window with the URL `https://jupyter.hpc.gwdg.de/user/lars.quentin01/desktop/`. Below that is the Jupyter interface with a status bar showing 'Status: Connected' and 'Remote Clipboard'. The main workspace contains a file browser on the left, a code editor in the center, and a terminal window on the right.

The file browser shows a folder named 'DL_LocalData...' containing a file 'DL_LocalData.ipynb'. The code editor contains the following Python code:

```
[3]: from matplotlib import pyplot as plt
plt.suptitle('First 9 images')
for i in range(9):
    plt.subplot(330 + 1 + i)
    plt.imshow(train_X[i], cmap=plt.get_cmap('gray'))
plt.show()
```

The terminal window displays the output of the code, showing a 3x3 grid of handwritten digits: 5, 0, 4, 1, 9, 2. The terminal also shows the execution progress of the Jupyter notebook, including the following logs:

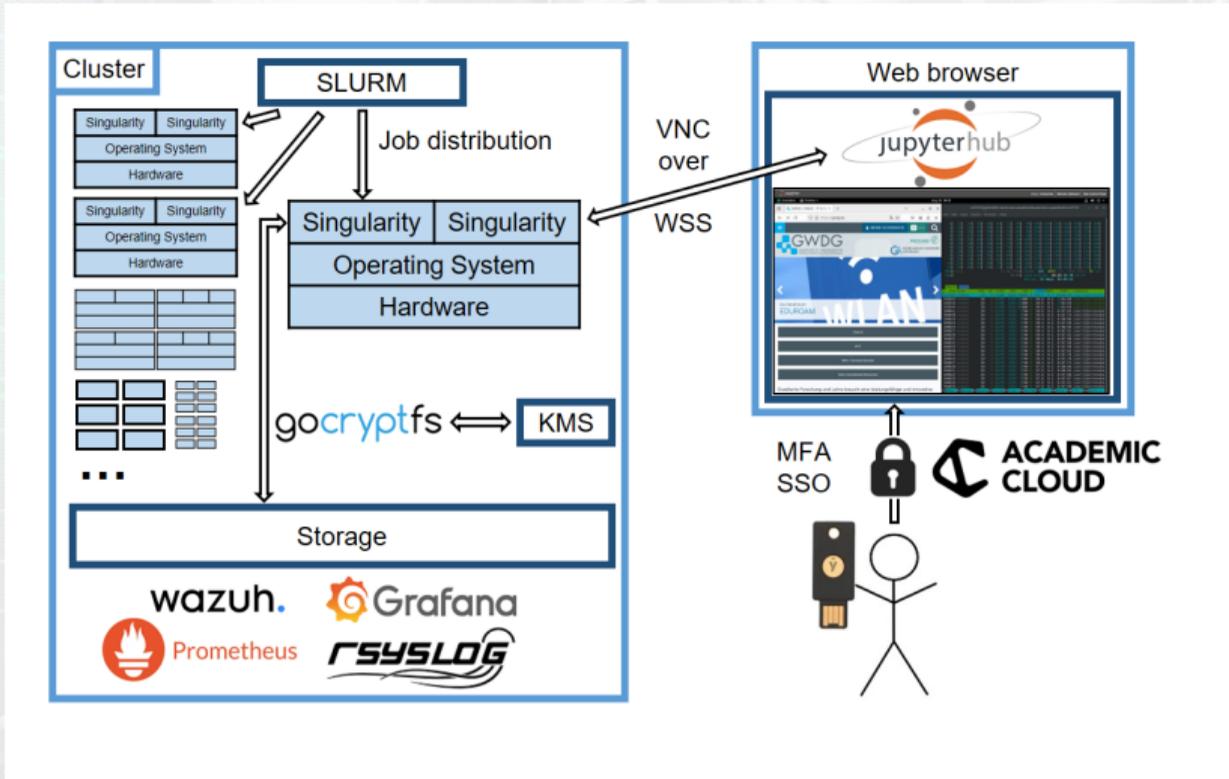
```
Epoch 4/15
422/422 ----- 6s 13ms/step - accuracy: 0.9791 - loss: 0.0695 -
val_accuracy: 0.9893 - val_loss: 0.0408
Epoch 5/15
422/422 ----- 6s 13ms/step - accuracy: 0.9825 - loss: 0.0588 -
val_accuracy: 0.9882 - val_loss: 0.0406
Epoch 6/15
422/422 ----- 6s 13ms/step - accuracy: 0.9836 - loss: 0.0537 -
val_accuracy: 0.9902 - val_loss: 0.0344
Epoch 7/15
422/422 ----- 6s 13ms/step - accuracy: 0.9842 - loss: 0.0502 -
val_accuracy: 0.9917 - val_loss: 0.0347
Epoch 8/15
422/422 ----- 6s 13ms/step - accuracy: 0.9856 - loss: 0.0462 -
val_accuracy: 0.9913 - val_loss: 0.0326
Epoch 9/15
422/422 ----- 6s 13ms/step - accuracy: 0.9858 - loss: 0.0446 -
val_accuracy: 0.9927 - val_loss: 0.0309
Epoch 10/15
422/422 ----- 6s 14ms/step - accuracy: 0.9876 - loss: 0.0399 -
val_accuracy: 0.9898 - val_loss: 0.0329
Epoch 11/15
422/422 ----- 6s 15ms/step - accuracy: 0.9872 - loss: 0.0389 -
val_accuracy: 0.9927 - val_loss: 0.0319
Epoch 12/15
422/422 ----- 6s 15ms/step - accuracy: 0.9884 - loss: 0.0354 -
val_accuracy: 0.9912 - val_loss: 0.0328
```

Features

- HPC-Hardware, inkl. GPUs
- Kollaboratives Arbeiten
- End-to-end Verschlüsselung mit Envelope Encryption
- Nutzbar ohne extra Software
- Terrabyte-Scale Uploader mit E2E-Verschlüsselung
- Two-Layer Permissions:
 1. POSIX
 2. Datei-Verschlüsselung
- Audit-Trail aller Dateizugriffe

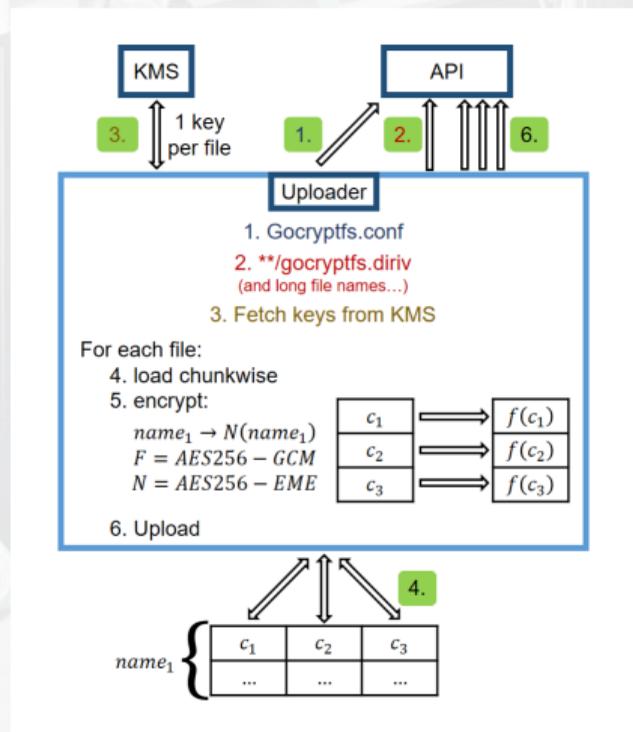
Philosophie

- **Zero-Trust:** Compute Nodes können nicht als sicher angesehen werden
- **Permanente Verschlüsselung:** Daten sind von Upload zu Datenexport für Publikation nie entschlüsselt
- **Integration in existierende Prozesse:**
 - Datenanfragen über Data Owner
 - Impliziert KYC
 - Pseudonymisierung via DIZ/FDPG
 - SDC über Data Owner
- **Nutzung existierender HPC-Infra**

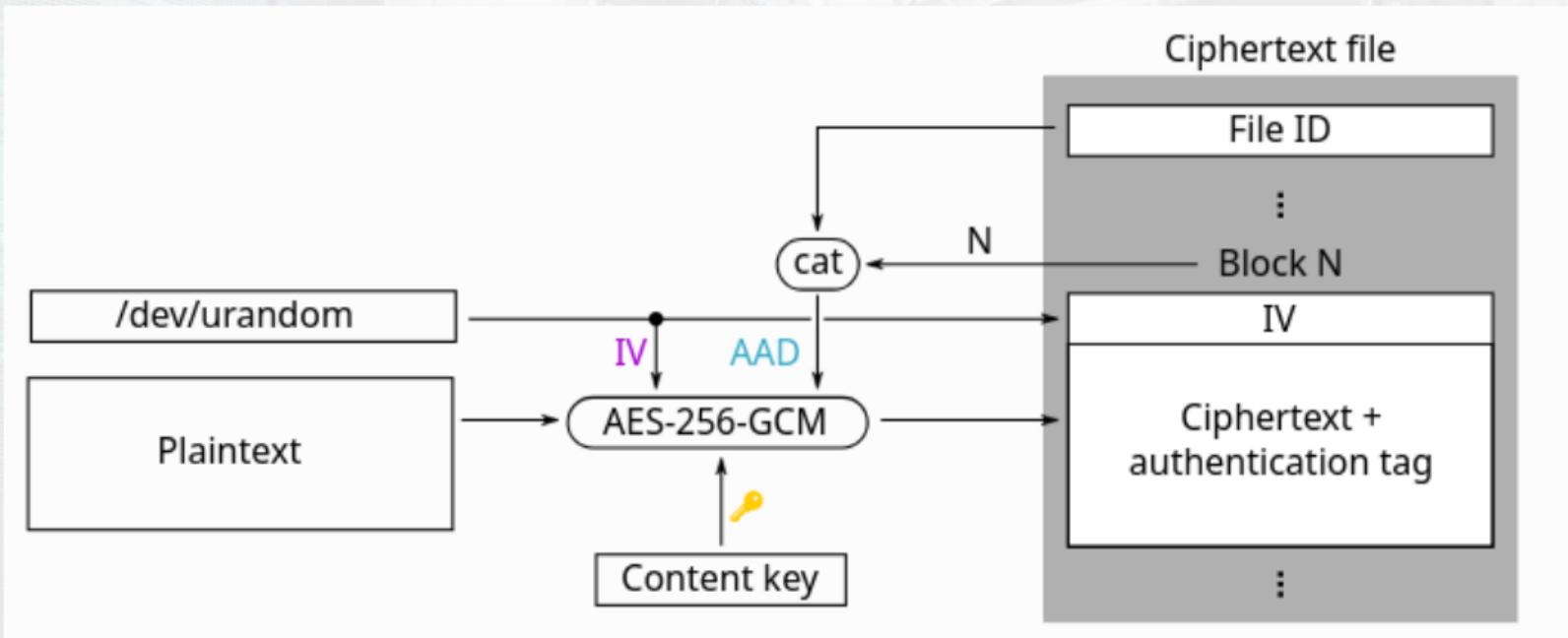


- JupyterHub-embedded Desktops
 - TigerVNC in Container, websockify TCP Socket \Rightarrow WSS, NoVNC JS Frontend
 - Somit via AcademicCloud SSO abgesichert
- VDI in Singularity/Apptainer Container, gestartet als SLURM-Job
- Internet deaktiviert, Software a-priori installiert
- Container read-only, mit Ausnahme von Gocryptfs mount.
 - Gocryptfs fork, per-file verschlüsselung via user-owned Key Encryption Key (KEK)
- Daten nicht exportierbar ohne manuelle Review
- Hauptprobleme:
 1. Sicherer Upload
 2. Sicheres Key-Management

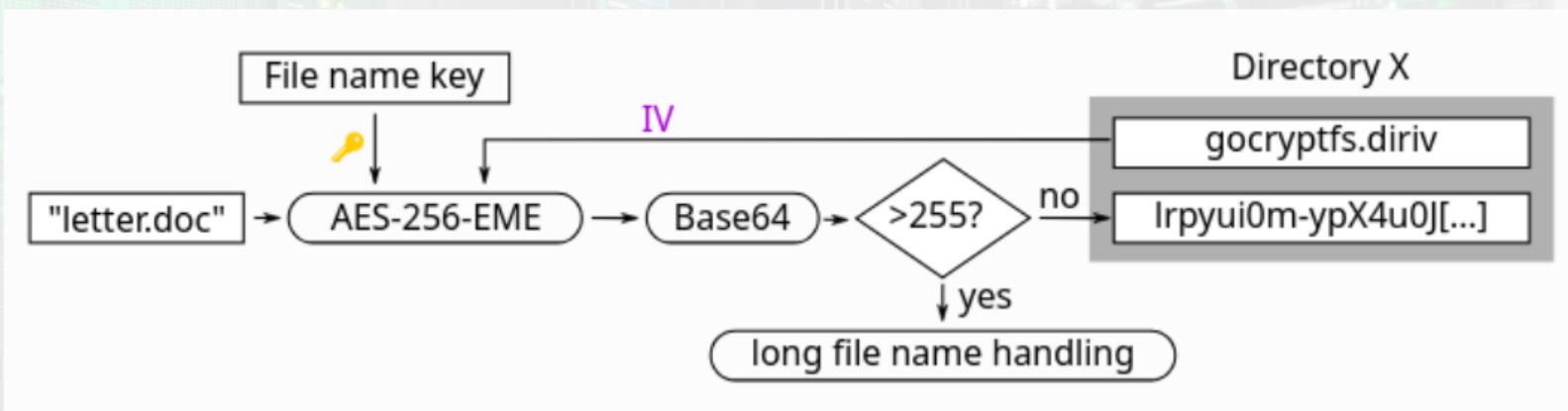
- Für TRE geschriebenen Uploader
 - Beliebig große Dateien
 - Multithreaded Upload und Verschlüsselung
 - Wiederaufnehmbar auf Dateiebene
- POSIX-Backend, keine ext. Dependencies (wie Redis, S3 etc)
- Nutzer live-verschlüsselt jeden chunk beim Upload
- Nutzt gocryptfs-kompatible Verschlüsselung
 - AES256-GCM für Inhalt
 - AES256-EME für Name



- Daten sind permanent verschlüsselt at rest.
 - Verschlüsselung für Nutzer transparent (gocryptfs)
- Jede Datei hat eigenen Schlüssel
 - Fine-Grained Encryption \Rightarrow Principle of Least Privilege
- Daten können nicht durch Admins entschlüsselt werden
- Upload Node kriegt nie entschlüsselte Daten
- KMS speichert keine Keys unverschlüsselt
- Nutzergerät muss per Definition trusted sein
- Das TRE (die Containerlaufzeitumgebung) muss per Definition sicher sein

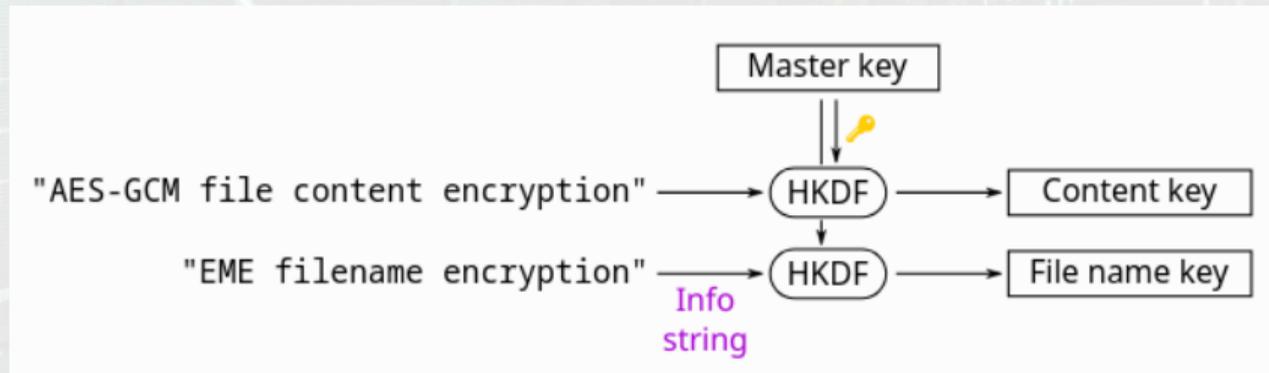


Quelle: GoCryptFS Manual



TRE@GWDG: Verschlüsselung

- Jede Datei oder Verzeichnis hat eigenen Data Encryption Key (DEK) im KMS
 - Envelope Encryption ermöglicht schnelles und sicheres Teilen der Daten
- Schlüssel-Mapping basierend auf ID (16 byte), nicht Dateipfad
 - Datei: File ID in Gocryptfs-Header
 - Verzeichnis: ./goecryptfs.dirid Datei
- Bevor DEK genutzt wird: HKDF in File name key und file content key
 - Grund: Kryptografisch GCM/EME-Interaktion unbewiesen



- Data Encryption Keys (DEKs) werden verschlüsselt im KMS gespeichert
- Hauptziel: Leak von KMS-Datenbank führt nicht zu kompromittierten DEKs
- Stattdessen besitzt User einen Key Encryption Key (KEK), der DEK entschlüsselt
- Verschlüsselungsalgorithmus hat Integritätscheck (z.B. AES-GCM)
 - Somit kann KEK als Authentifizierung genutzt werden

Funktionsweise anhand von Operationen

- **Bei Mount:** Der KEK wird an goecryptfs übergeben
- **Readdir:**
 - Verzeichnisidentifizier `goecryptfs.dirid` wird gelesen.
 - Anfrage des DEK basierend auf KEK + `dirid`
 - HKDF zum file name key
 - Readdir-Operation auf Dateisystem
 - For jeden Eintrag: Ciphernamen entschlüsseln
- **Read:**
 - Analog Verzeichnis-DEK erhalten
 - Plain-Name verschlüsseln um Ciphernamen zu erhalten
 - Aus Cipherdatei File ID extrahieren, Datei-DEK von KMS anfragen
 - Nach HKDF können chunks gelesen und AES256-GCM entschlüsselt werden

- **Bekannte Sicherheitslücke:** CVE-Management bereits deployed
- **Datenleck via Zwischenablage:** Deaktiviert im VDE
- **Falsche Nutzung:** Container isoliert; Internet blockiert; Training verpflichtend
- **Falsche POSIX-Rechte:** Doppelte Berechtigungsebene via KMS
- **KMS-Datenbank leaked:** Schlüssel nicht nutzbar ohne KEKs der Nutzer.
- **Nutzer verliert KEK:** "Tresor"-Funktion basierend auf Shamir's Secret Sharing

- **Bösartige Forscher:** Keine Möglichkeit abzusichern da Remotezugriff erlaubt
 - Beispiel: Bildschirmaufnahme mit OCR
- **Unbemerkttes Teilen von KEKs:** Da Zugriffe von überall erlaubt keine Anomalieerkennung bisher geplant.
- **Zero-Day Exploits:** Angreifer kann dann aus dem Arbeitsspeicher von goecryptfs theoretisch den KEK extrahieren.

- Software-Entwicklung abgeschlossen
- Nächster Schritt: Deployment
- Pilotprojekt: SHIP 2 Datensatz (Universitätsmedizin Greifswald)
 - Aktueller Stand: Warten auf Datenimport
- Noch fehlende Features:
 - SLURM-Job-Support
 - Symlinks / Hardlinks

- Neue regulatorische Anforderungen und Möglichkeiten
 - Zum Zugriff auf (Gesundheits-)daten
 - Zum Prozessieren dieser Daten
- Details werden bis Q2 2027 auf sich warten lassen
- Genügend Beispiele im europäischen Kontext um erste Erfahrungen zu sammeln
- GWDG ist dabei einen solchen PoC zu deployen
 - PoC wird in Q4 von 2 Universitätskliniken evaluiert