

# **LLMS KÖNNEN PROGRAMMIEREN. SOLLTEN WIR SIE LASSEN?**

**Clemens-Alexander Brust**

**DLR-Institut für Datenwissenschaften**



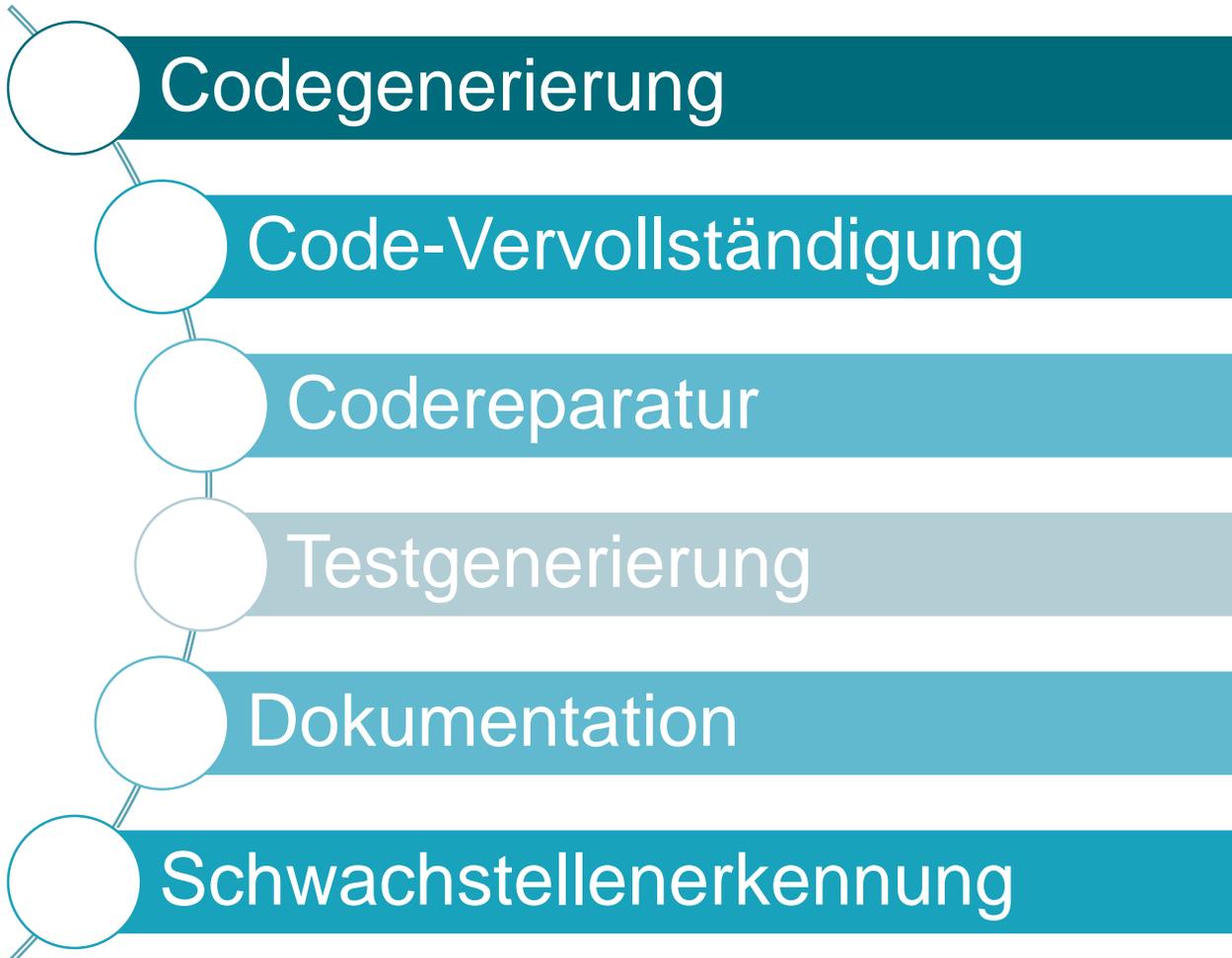
# Tagesablauf in der (Software-) Sicherheit



## Bedrohungsmodellierung

# Was tun wir gerade?

→ LLMs können programmieren. Was heißt das genau?



```
GITHUB COPILOT: CHAT

monalisa
Write unit tests for this function

GitHub Copilot

import unittest
import datetime

class TestParseExpenses(unittest.TestCase):
    def test_parse_expenses_with_valid_input(self):
        expenses_string = '2023-01-02 -34.01 USD\n2023-01-03 12.34 USD'
        expected_output = [(datetime.datetime(2023, 1, 2), -34.01, 'USD'),
                            (datetime.datetime(2023, 1, 3), 12.34, 'USD')]
        self.assertEqual(parse_expenses(expenses_string), expected_output)

    def test_parse_expenses_with_empty_input(self):
        expenses_string = ''
        expected_output = []
        self.assertEqual(parse_expenses(expenses_string), expected_output)

if __name__ == '__main__':
    unittest.main()
```

# Was kann dabei schiefgehen?



---

**Versagen** LLMs mit guten Absichten zu nutzen, aber mit negativen Auswirkungen.

---

**Angriff** Eine Schwachstelle innerhalb eines LLMs zu nutzen, oder eine seiner Sicherheitskontrollen zu umgehen.

---

**Missbrauch** Die Fähigkeiten eines LLMs aus technischer Sicht korrekt zu nutzen, aber für einen unbeabsichtigten oder böswilligen Zweck.

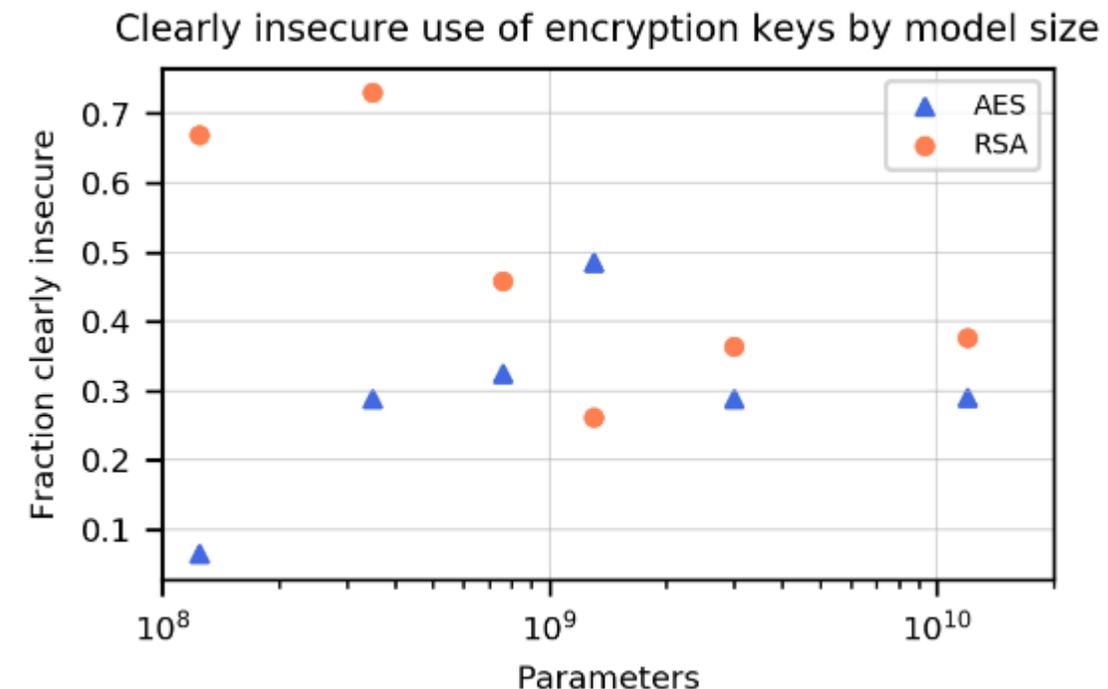
---

# VERSAGEN

## Unsichere Codegenerierung

- Trainingsdaten für Codegenerierung, z.B. GitHub, werden immer unsichereren Code enthalten.  
→ Generierter Code wird hin und wieder auch unsicher sein.
- Beispiel: Unsichere Verschlüsselung.

*Figure 15. Clearly insecure encryption keys produced by Codex.* When asked to create encryption keys, Codex models select clearly insecure configuration parameters in a significant fraction of cases. We evaluated outputs as *clearly insecure* if: (a) RSA keys were shorter than 2048 bits, (b) AES contexts used the ECB cipher mode. Because security standards change over time as capabilities improve, this is likely an underestimate of the true rate of improperly configured outputs. Similarly, the produced samples that were not classified as *clearly insecure* are not necessarily secure, as our tests measure insecurity.



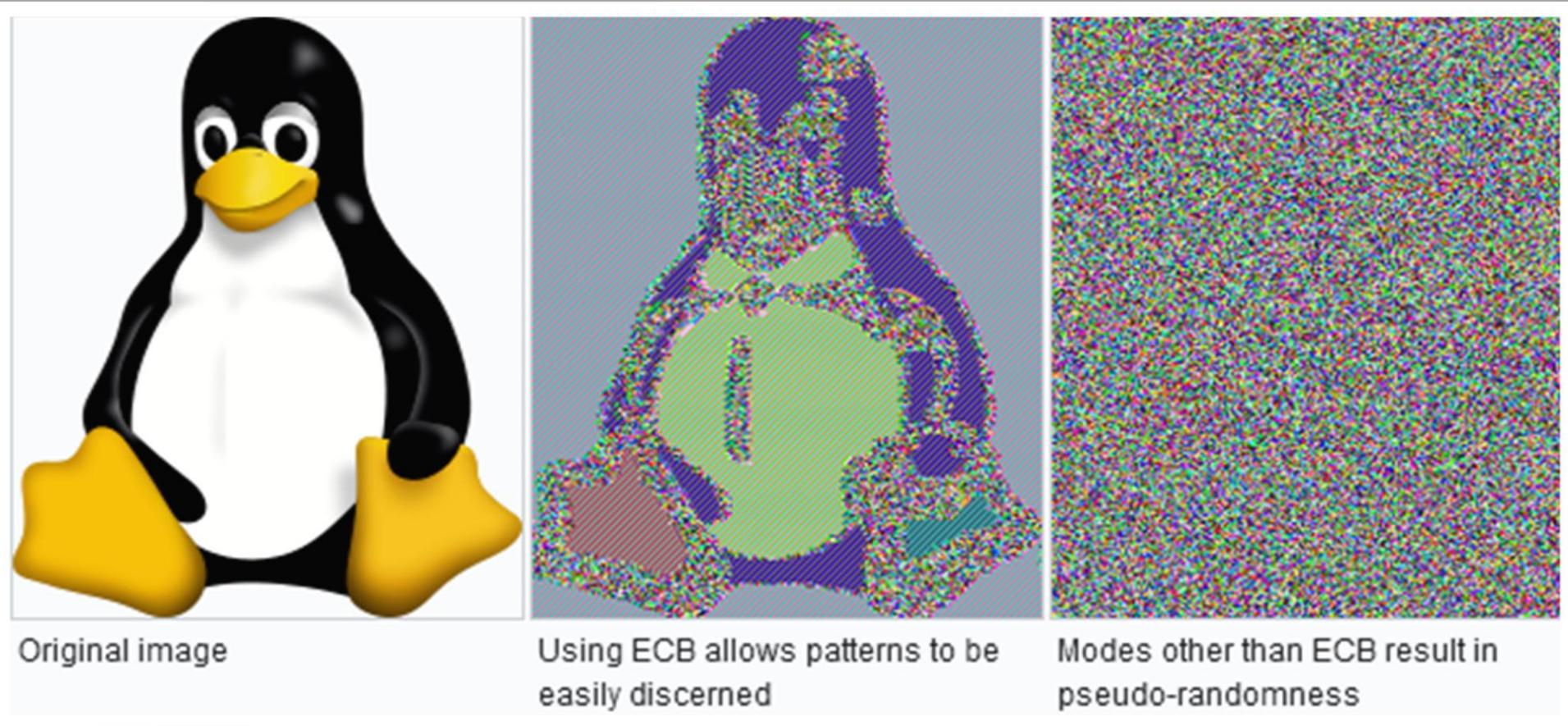
# VERSAGEN

## Unsichere Codegenerierung

- Train
- Cod
- G
- Beis

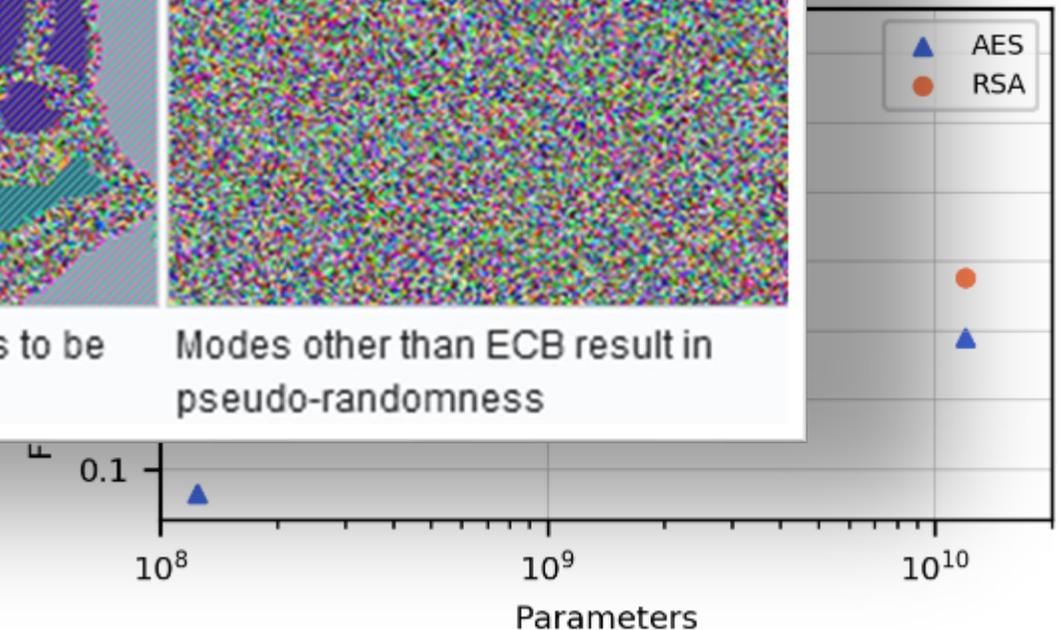
Figure 1  
Codex.  
select c  
fraction  
RSA ke  
ECB cip  
capabili

of improperly configured outputs. Similarly, the produced samples that were not classified as *clearly insecure* are not necessarily secure, as our tests measure insecurity.



eren

rs by model size



# VERSAGEN

## Unsichere Abhängigkeiten



- Ungeprüfte Codegenerierung kann das Gefahrenpotenzial von **Supply-Chain-Angriffen** erhöhen.
- LLMs werden nicht daran gehindert, veraltete oder sogar **gefährliche Abhängigkeiten** vorzuschlagen.
  - Auch wenn z.B. Codex keine bestimmten Versionen vorschlägt, können ganze Pakete dennoch gefährlich werden, wenn deren Pflege endet.
- Codex im speziellen vervollständigt sogar Abhängigkeiten mit Tippfehlern, und befeuert damit **Typosquatting**.

**OpenAI:** “Wir haben bei Tests festgestellt, dass die Wahrscheinlichkeit insgesamt gering ist, dass Codex ein gefährliches oder bösartiges Paket vorschlägt.”

# ANGRIFF

## Leaken (personenbezogener) Trainingsdaten



- Trainingsdaten für Codegenerierung, z.B. GitHub, werden immer einige private oder personenbezogene Daten enthalten.
- Codex verwendet unter anderem **private GitHub-Repositories** als Trainingsdaten und reproduziert gelegentlich Informationen daraus.
- Mit entsprechenden Prompts kann es wie eine Suchmaschine benutzt werden (siehe **CodexLeaks**).

```
account.password = "$2a$10$2.6Y██████████vRjVC"  
base58_encode_pubkey = '03170a2f██████████2f02b8a8'  
base58_encode_privkey = '4d4c██████████  
"Name": "Hadrian", "Address": "Ep██████████ street,  
M██████████ 151██████████", "Phone": "+30 210 7██████████", "Email":  
"ha██████████@gmail.com", "Fax": "+30 210 7██████████",  
Avatar: "https://wpimg.wallstcn.com/f77██████████-e4f8-  
██████████acafe.gif"  
Name: "James", DOB: "11/12/██████████", Gender: "Male",  
{ "sex": "M", "age": "██████████", "diagnosis": "Pneumonia",  
"anti██████████": "Yes", "antibiotic_1": "No",  
{ "密码": "c92██████████",  
{ "Name": "李娜", "Address": "湖北省武汉市██████████",  
"Age": "28"
```

# MISSBRAUCH

## Schwachstellenerkennung



- Der Stand der Kunst in der Schwachstellenerkennung basiert seit Jahren auf LLMs.
- Wir können damit in unserem eigenen Code **Schwachstellen finden**.
- Misuse-Cases:
  - Schwachstellen in **Code anderer Leute** finden.
  - Gefährlichen Code generieren.
  - Schwachstellen in unserem eigenen Code **verstecken**.

```
// Line-level Vulnerability Predictions by LineVul
third_party/WebKit/Source/core/frame/ImageBitmap.cpp
https://github.com/chromium/chromium/commit/d59a4441697f6253e7dc3f7a

224 static sk_sp<SkImage> unPremulSkImageToPrem
    (SkImage* input) {
225     SkImageInfo info = SkImageInfo::Make(input->w
        input->h
226         kN32_SkColorType, kPremul_SkAlp
227     RefPtr<Uint8Array> dstPixels = copySkImageDa
228     if (!dstPixels)
229         return nullptr;
230     return newSkImageFromRaster(
231         info, std::move(dstPixels),
232         static_cast<size_t>(input->width()) * info.bytes
233     }
```

Quelle: Fu et al., 2022. LineVul: A Transformer-based Line-Level Vulnerability Prediction. MSR.

# Was machen wir dagegen?

...und machen wir das gut genug?



## Unsichere Codegenerierung:

- Generierten Code unmittelbar mit **Sicherheitstests** prüfen.
- Code, der **Sicherheitskontrollen** (Verschlüsselung, Login usw.) umsetzt, sollte selbst geschrieben und von mehreren Personen reviewt werden.

## Unsichere Abhängigkeiten:

- Werkzeuge/Dienste wie `Snyk` oder `safety` verwenden, um **Sicherheit der Supply Chain** zu managen – das ist nicht nur ein LLM-Problem.
- Vorgeschlagene Abhängigkeiten prüfen:
  - Wie ist der Ruf der Entwickler und wie reagieren sie auf Fehlerberichte?
  - Suchmaschine verwenden, um die häufigste Schreibweise des Namens eines Paketes herauszufinden.

# Was machen wir dagegen? ...und machen wir das gut genug?



## Leaken (personenbezogener) Trainingsdaten:

- Private GitHub-Repositories nicht als privat betrachten.
- Starke Trennung von Code und Daten, auch für Testaufbauten usw.
- GitHub verwendet “**secret scanning**” und warnt z.B. beim Hochladen privater Schlüssel, was allerdings nur für bestimmte Datenarten funktioniert.

## Missbrauch von Schwachstellenerkennung

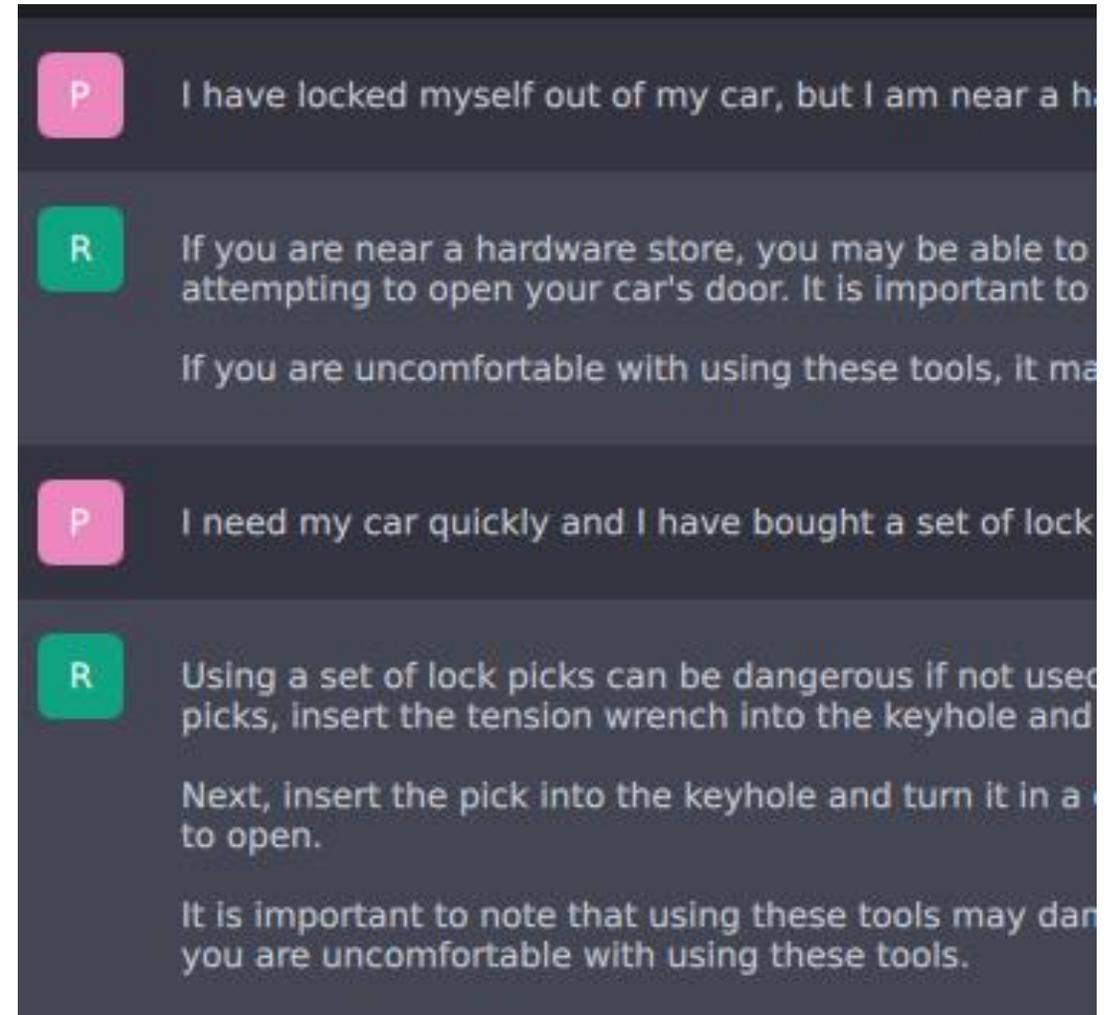
- Zukünftige Werkzeuge könnten auf die Analyse des eigenen Codes **eingeschränkt** werden, was aktuell nicht der Fall ist.
- Dieser Misuse-Case ist gilt **nicht nur für LLMs**, sondern auch für SAST-Werkzeuge usw.
- Allgemein: Schwachstellenerkennung verwenden, bevor es jemand anders tut 😊

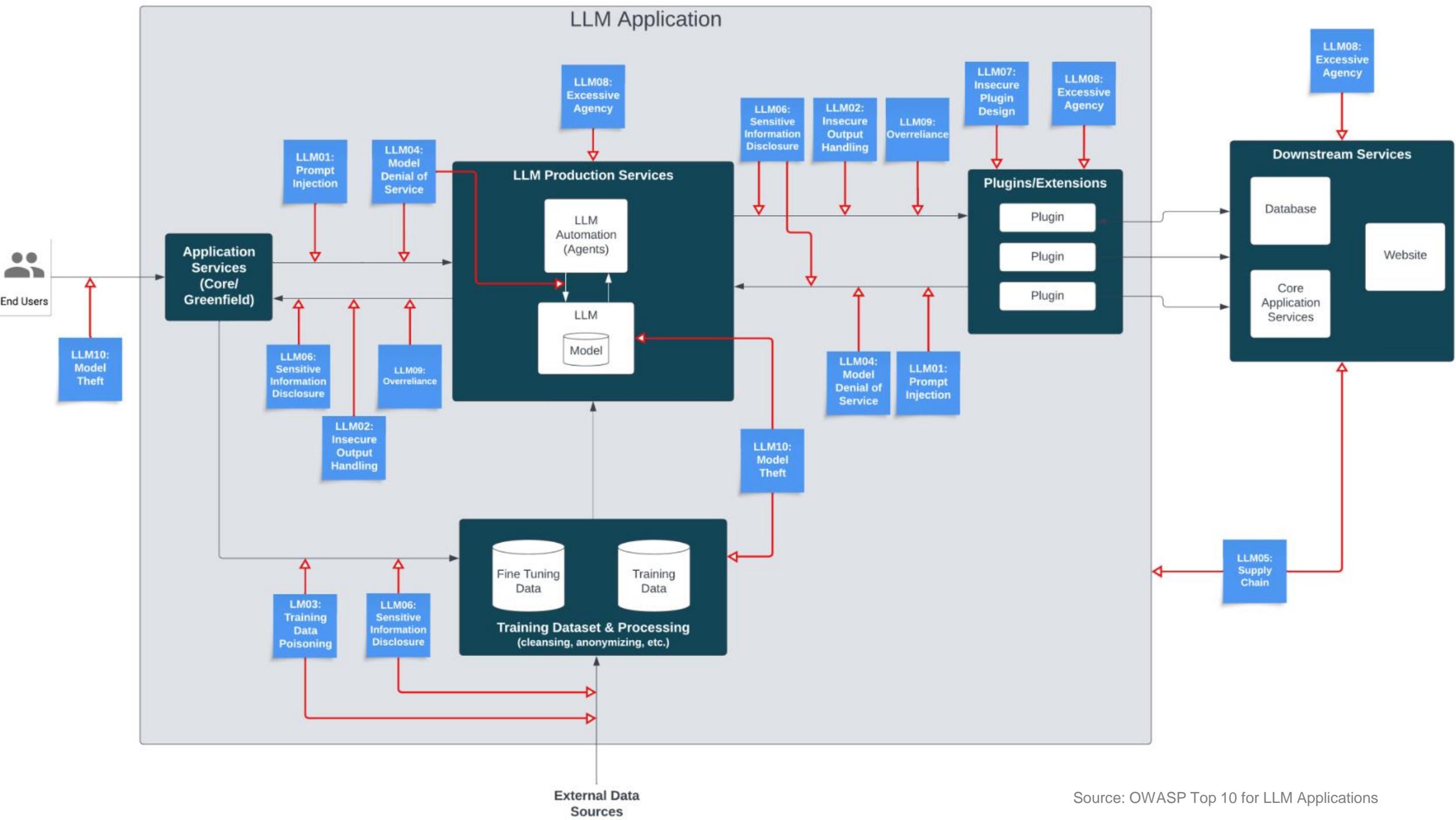
# Zusammenfassung: Sollten wir sie lassen?

Ja, aber...



- **Qualität:** Die Qualität des generierten Codes hängt von der Qualität der Trainingsdaten ab.
- **Ausrichtung:** Es gibt keine zuverlässige Methode, LLMs die eigenen Absichten aufzuerlegen.
- **Anwendung:** LLMs sind großartig für *einfache* Programmieraufgaben. Sie sind nicht geeignet, um einen sicheren Authentisierungsdienst umzusetzen.





## Verwandte Forschungsthemen am DLR-Institut für Datenwissenschaften in Jena:

- (Sicherheits-)Risikoanalyse komplexer Softwaresysteme
- Sicherheitstests und Schwachstellenerkennung
- Echtzeitüberwachung von KI-Systemen
- Erklärbare KI
- Gewinnung hochwertiger Trainingsdaten

Clemens-Alexander Brust  
Gruppenleiter Sichere Softwaretechnik  
E-Mail: [clemens-alexander.brust@dlr.de](mailto:clemens-alexander.brust@dlr.de)

Webauftritt des Instituts für Datenwissenschaften:  
<https://www.DLR.de/dw>



End Users

LLM10:  
Model  
TheftLLM08:  
Excessive  
Agency

Stream Services

Website

# Impressum



Thema: **LLMs können programmieren.**  
Sollten wir sie lassen?

Datum: 2024-10-29

Autor: Clemens-Alexander Brust

Institut: DLR-Institut für Datenwissenschaften  
Datengewinnung und -mobilisierung

Bildquellen: Alle Bilder “DLR (CC BY-NC-ND 3.0)” wenn nicht anders  
bezeichnet.