

DEFAME FAKES

Detektion von Deepfakes und medialen Manipulationen
von Bildern und Videos

Mina Schütz, Scientist am AIT Austrian Institute of Technology, 11.09.2024

FFG Forum – Digitale Sicherheit in Europa – Herausforderungen und Chancen für die Gesellschaft



DEFAME FAKES



pwc



 Bundesministerium
Inneres

 Bundesministerium
Finanzen

 Bundesministerium
Landesverteidigung




KEY FACTS



defame Fakes

- Forschungsprojekt finanziert durch das Sicherheitsforschungsförderprogramm KIRAS des Bundesministeriums für Finanzen
- Laufzeit: Februar 2024 – Jänner 2026
- Website: www.defamefakes.at
- Koordination: AIT, Martin Boyer <martin.boyer@ait.ac.at>

 Bundesministerium
Finanzen



 Bundesministerium
Inneres

 Bundesministerium
Landesverteidigung

ZIELE DES PROJEKTS



Erforschung und Entwicklung von **Assessment Werkzeugen** zur **Erkennung von Deepfakes**



Konzeptionierung und Initiierung von gesamtgesellschaftlichen **Awareness-Maßnahmen**



Analyse von Bedrohungsszenarien, gesellschaftlichen, ethischen und rechtlichen Implikationen

DEEPPFAKES

- KI-basierte Manipulation von Audio, Bild, Video, Text
- Zunehmend einfach zugänglich
- Moderne Generatoren können sehr realistische Ergebnisse liefern

Bedrohungsszenarien

- Kriminelle Gruppen als Early Adopters
- Desinformation und Fake News
- Image- und Rufschädigung
- Identitätsdiebstahl
- Aushöhlung des Vertrauens in digitale Inhalte
- Betroffene: Privatpersonen, Regierungen, Unternehmen



DEEPFAKE DETEKTION



■ Herausforderungen:

- Kontinuierliche Weiterentwicklung der Deepfake-Technologie
- Dringlichkeit von Detektionslösungen
- Bedarf an verallgemeinerbaren Erkennungsmethoden für Deepfakes, unabhängig der Techniken die zur Generierung verwendet wurden

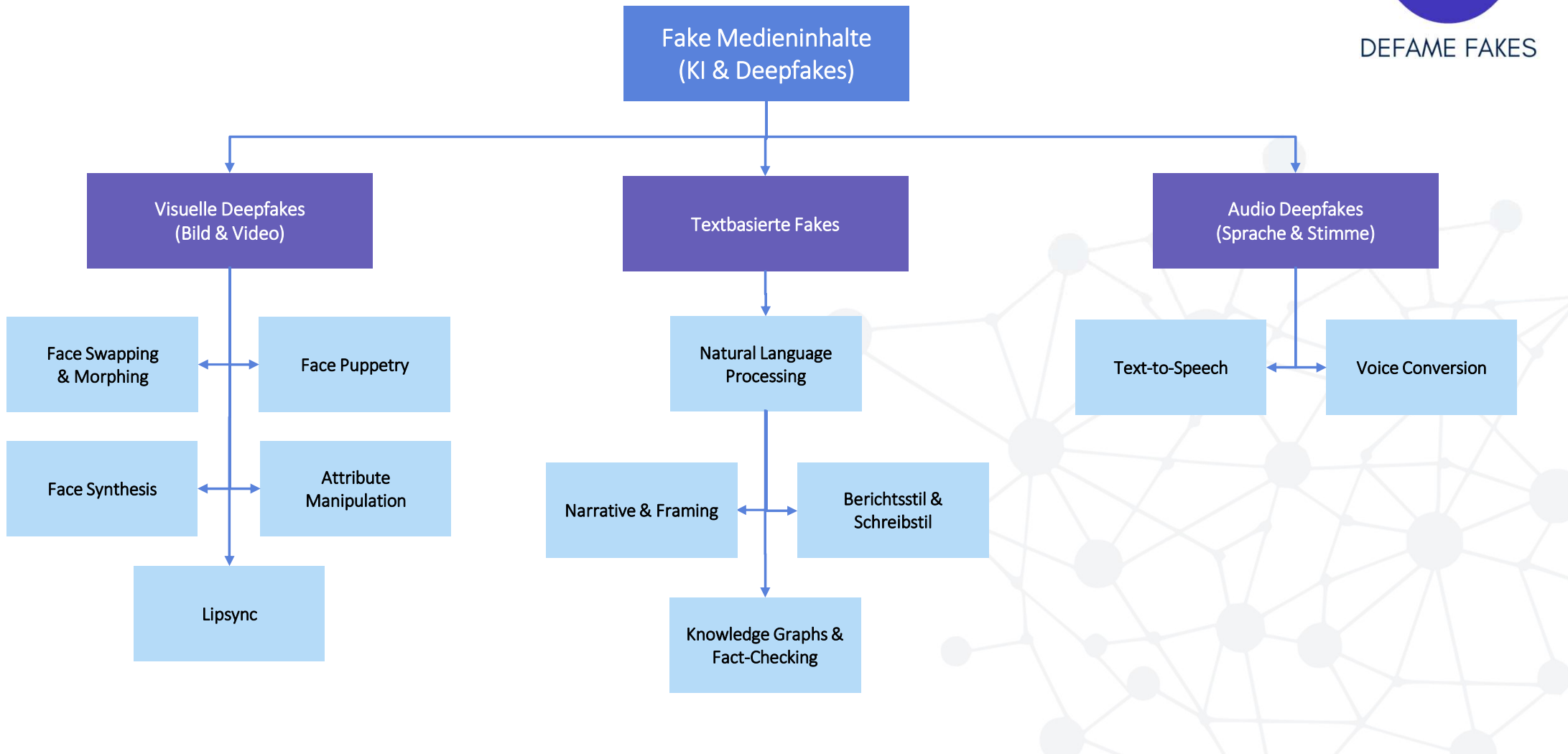
■ Kernpunkte:

- Konzentration auf bestimmte Anwendungsfälle und Szenarien
- Erklärbare und interpretierbare Ergebnisse der Erkennung
- Sammlung repräsentativer, Bias-freier Daten für Training und Evaluierung
- Skalierung und Integration der Lösung in bestehende Arbeitsabläufe

DEEPFAKE TECHNOLOGIEN

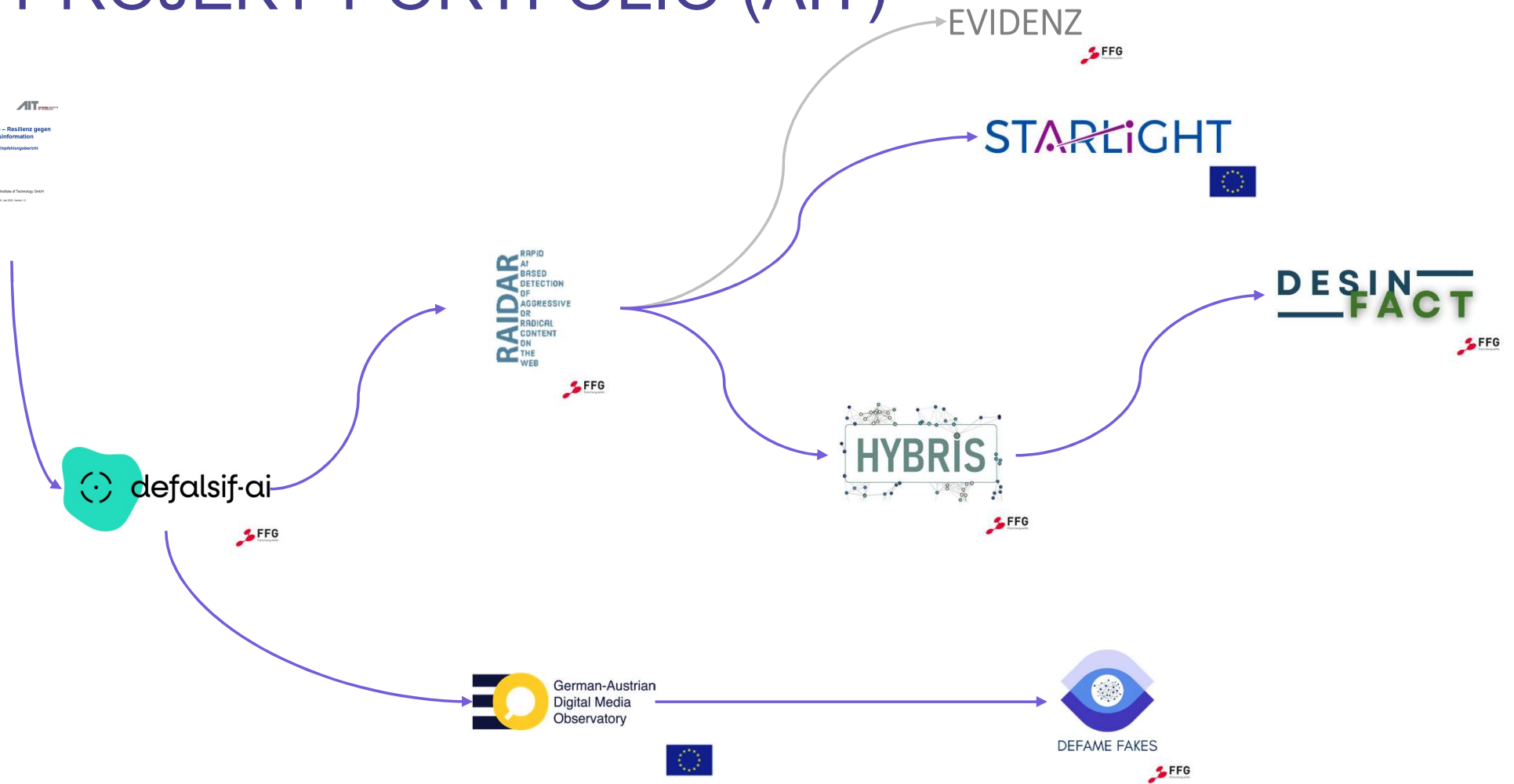


DEFAME FAKES



PROJEKT-PORTFOLIO (AIT)

AIT
FORGE
Vorstudie – Resilienz gegen
Desinformation
02.3. Digitalisierungsbericht
AIT Austrian Institute of Technology GmbH
01.04.2024 (Intern)

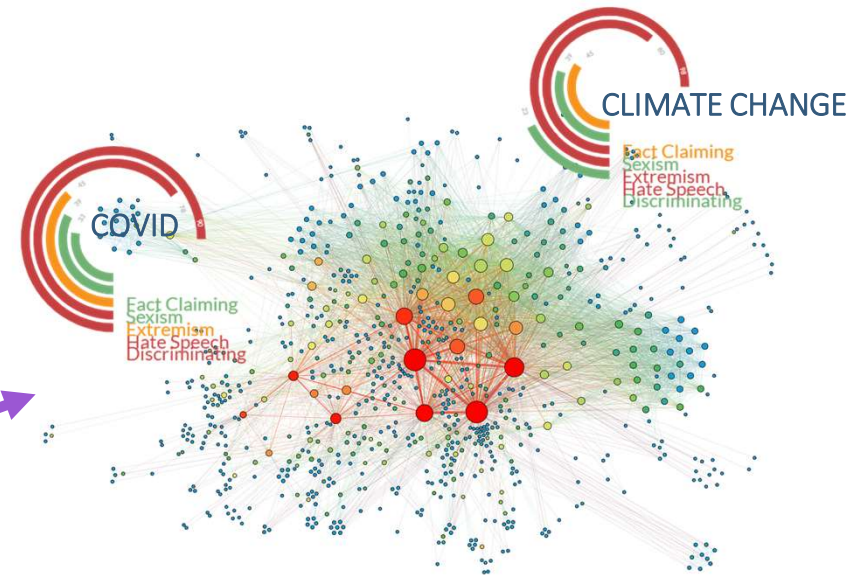
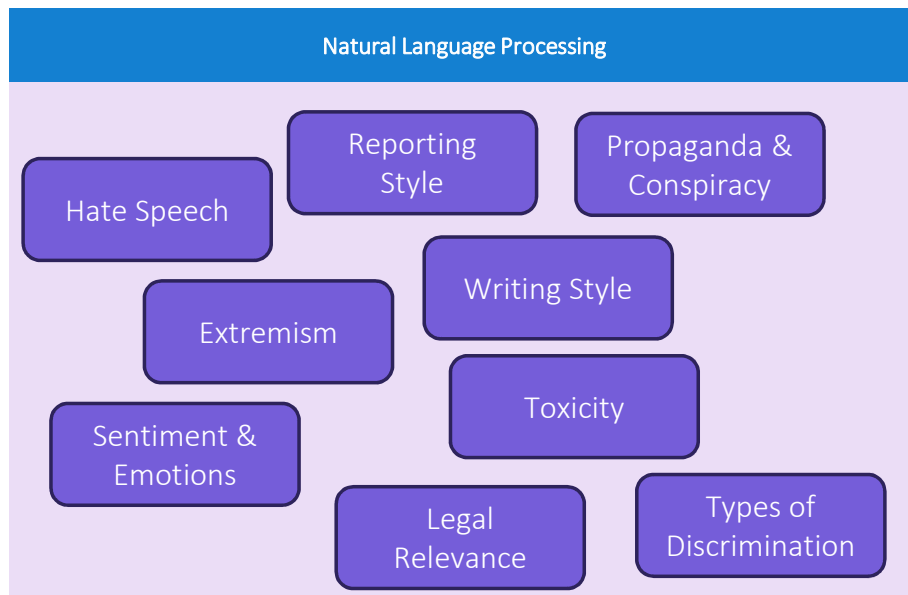


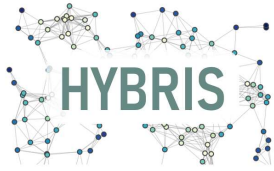
INFORMATION NUTRITION LABELS



AUSTRIAN INSTITUTE OF TECHNOLOGY

- Automatische Erkennung von Desinformation mittels KI Modellen
- Einschätzung von über 50 Merkmalen
- Erkennung von Zusammenhängen und Themen





Bundesministerium Landesverteidigung



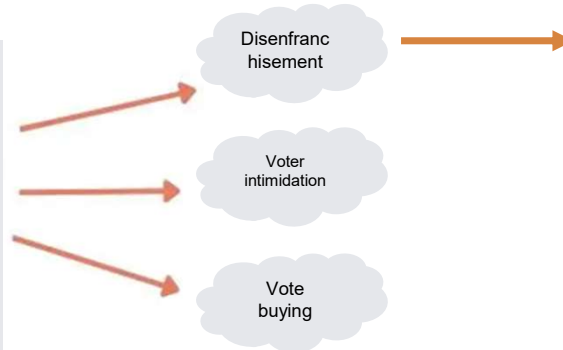
SITUATIONAL AWARENESS

THREAT ESTIMATION & NARRATIVES

Taxonomy

- Election fraud analysis
- Electoral fraud
 - Electorate manipulation
 - Artificial migration or party membership
 - Disenfranchisement
 - Division of opposition support
 - Voter intimidation
 - Voter disinformation
 - Vote buying
 - Voting process and results
 - Misleading or confusing ballot papers
 - Ballot stuffing
 - Misrecording of votes
 - Misuse of proxy votes
 - Destruction or invalidation of ballots
 - Tampering with electronic voting systems
 - Voter impersonation
 - Postal ballot fraud

Data Source



Story-Card

Title: <LLM generated accumulated Title>

Num. Articles: 35

Key-Claims:

- <LLM extracted key-claim>
- <LLM extracted key-claim>
- <LLM extracted key-claim>
- <LLM extracted key-claim>
- <LLM extracted key-claim>

Indicators:

- Incitement of Violence
- Economic Implications
- Medical Implications

Narratives:

- Central Narrative:
- Further Narrative

Predjudice:

- Antisemitism
- Anti-Feminism





DEFAME FAKES

Vielen Dank für Ihre Aufmerksamkeit!

Mina Schütz

Mina.schuetz@ait.ac.at

