Al Factory Austria Al:AT



Trustworthy AI für Fortgeschrittene: Ethische Aspekte

Dr. Peter Biegelbauer Co-Lead Legal, Regulatory and Ethics



Warum brauchen wir die Al Factory Österreich?



Souveränität



Ethik und Trustworthiness



KI-Ökosystem



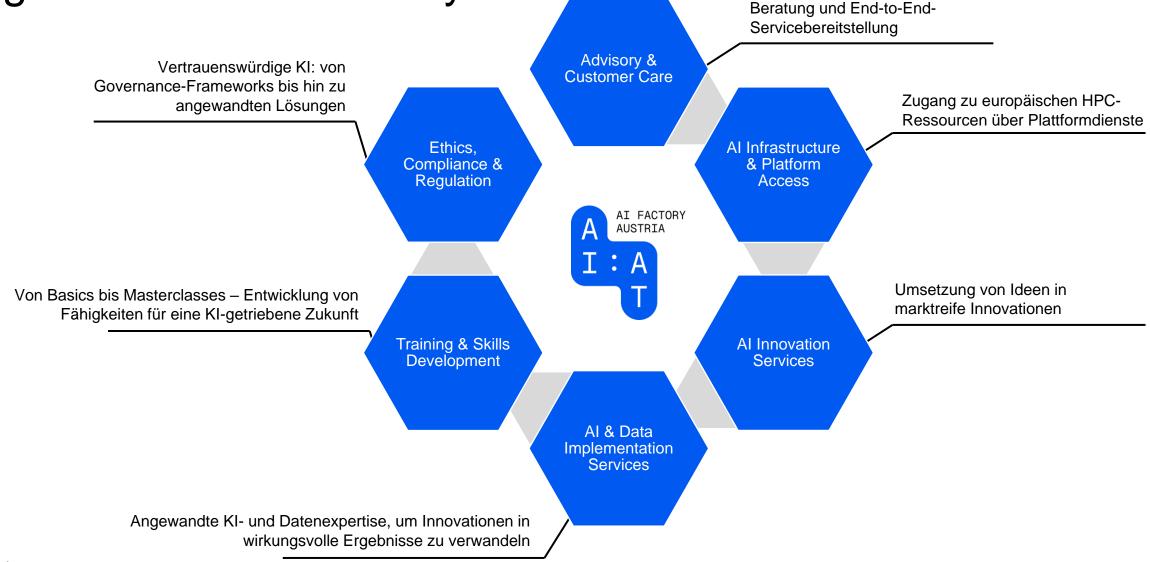
Unsere Mission Von der Idee zur Innovation

Etablierung eines One-Stop Shop für Al Schließen des Al Ressourcenund Wissens-Gap Bewerbung von ethical und trustworthy Al

Launchpad Für Al-getriebene Innovation



Unsere Services unterstützen über den gesamten Al-Lebenszyklus



Al Factory Austria Al:AT – Team





Al Factory Austria Al:AT Consortium

Disclaimer:

The speakers are solely sharing their personal experiences. Therefore, this free seminar is not a substitute for professional/legal advice.

Beneficiaries





Affiliated Entities























Why do we need Al Factory Austria?



Sovereignty



Ethics and Trustworthiness



Connecting the Ecosystem



Fallstudie: Robodebt

Aufdecken von Sozialleistungsbetrug in Australien

Automatisiertes System (Robodebt) sollte Sozialleistungsbetrug aufdecken

- Fehlerhafte Datenabgleiche zwischen Steuer- und Sozialbehörden
- als Folge: Falschbeschuldigungen hunderttausender Bürger:innen
- als Folge: Massive finanzielle & psychische Belastungen, teils tragische Konsequenzen

m Politische und rechtliche Konsequenzen

- Offizielle Entschuldigung der Regierung, Rückzahlungen & Entschädigungen > 1,8 Mrd. AUD
- Royal Commission (2022–23) mit umfassender Untersuchung
- Stärkung von Rechtsstaatlichkeit, Transparenz & Aufsicht
- Verbesserte Kontrollmechanismen f
 ür staatliche KI
- Klare Verantwortlichkeiten für Behörden & Politik



Robodebt: was hätte es gebraucht?



Robuste Datenvalidierung & -integration:

Aufbau sicherer Schnittstellen zwischen Sozialdatenbanken, Steuerbehörde und anderen Quellen mit automatisierter Fehlerprüfung vor Berechnungen.

•Transparente Algorithmusgestaltung:

Offenlegung der Entscheidungslogik und Dokumentation der Gewichtungen; verpflichtende "explainability"-Mechanismen für jede automatisierte Entscheidung.

•Human-in-the-Loop-Prinzip:

Automatisierte Vorschläge dürfen keine endgültigen Verwaltungsakte auslösen – verpflichtende menschliche Prüfung bei *jeder* Rückforderungsentscheidung.

Algorithmisches Impact Assessment (AIA):

Vor Einführung KI-basierter Systeme: Risiko- und Fairnessprüfung mit Fokus auf Diskriminierung, Fehlklassifikationen und Verfahrensgerechtigkeit.



Robodebt: was hätte es gebraucht?



Kontinuierliche Auditierung & Monitoring:

Einrichtung unabhängiger Kontrollstellen zur Überwachung der Datenqualität, Modellperformance und der sozialen Auswirkungen automatisierter Entscheidungen.

•Rechtsstaatliche Rückkopplungsschleifen:

Implementierung effektiver Beschwerdemechanismen und automatischer Benachrichtigungssysteme bei algorithmisch generierten Bescheiden.

•Ethik- und Governance-Rahmen:

Verankerung ethischer Leitlinien (z. B. Fairness, Rechenschaftspflicht, Transparenz) im gesamten Lebenszyklus der KI – von Entwicklung bis Einsatz.



KI-Ethik braucht das Land!

- Internationaler Review von **KI-Guidelines** (Jobin et al. 2019):
 - > transparency
 - > justice and fairness
 - > non-maleficence
 - > responsibility
 - > privacy
 - ➤ beneficence
 - > freedom and autonomy
 - > trust
 - > sustainability
 - > dignity and solidarity

- > Transparenz
- Gerechtigkeit und Fairness
- Schadensvermeidung
- Verantwortung
- Privatsphäre
- Wohlwollen
- > Freiheit und Autonomie
- Vertrauen
- Nachhaltigkeit
- Würde und Solidarität



Ethische Prinzipien im Al-Act

Integration der **HLEG-Prinzipien für vertrauenswürdige KI** in den AI-Act (Erwägungsgründe):

- Vorrang menschlichen Handelns und menschliche Aufsicht
- Technische Robustheit und Sicherheit
- Datenschutz und Datenqualitätsmanagement
- Transparenz
- Vielfalt, Nichtdiskriminierung und Fairness
- Gesellschaftliches und ökologisches Wohlergehen

Das fehlende Prinzip: Accountability (Rechenschaftspflicht) → Al Liability Directive (geplant)





Vertrauenswürdige KI Rechtmäßige KI Ethische KI

(Wird hier nicht behandelt.)

Fundamente einer vertrauenswürdigen KI

Sicherstellung der Einhaltung ethischer Grundsätze auf Basis der Grundrechte

4 Ethische Grundsätze

Erkennen und Lösen von Spannungen zwischen ihnen

Robuste KI



- Schadensverhütung
- Fairness
- Erklärbarkeit

Verwirklichung einer vertrauenswürdigen KI

Sicherstellung der Umsetzung der Kernanforderungen

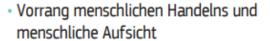
7 Kernanforderungen

Kontinuierliche Bewertung und Berücksichtigung der Kernanforderungen während des gesamten Lebenszyklus des **KI-Systems**



Technische Verfahren

Nichttechnische Verfahren



- Technische Robustheit und Sicherheit
- Datenschutz und Datenqualitätsmanagement
- Transparenz
- · Vielfalt, Nichtdiskriminierung und Fairness
- Gesellschaftliches und ökologisches Wohlergehen
- Rechenschaftspflicht

KAPITEL I

KI-Prinzipien in Österreich

Leitfaden digitale Verwaltung: KI, Ethik und Recht



Eine Arbeit des AIT AI Ethics Lab für BMKÖS / BKA

Webpage



Leitfaden



YouTube





Kriterien im Leitfaden digitale Verwaltung I



- 1. Recht Einhaltung des geltenden Rechts, KI-Anwendung muss die einschlägigen Gesetze und Vorschriften einhalten, einschließlich der Grundrechte
- 2. Transparenz Informationen über die KI-Anwendung verfügbar und zugänglich machen, Transparenz in KI-Entscheidungsprozessen fördern, Öffentlichkeit und Verwaltungsbedienstete über die Ziele und der KI-Anwendung informieren, Offenlegung der Entscheidungsergebnisse
- 3. Unparteilichkeit und Fairness KI-Anwendung muss unvoreingenommene und vielfältige Daten und Modelle verwenden, Vermeidung der Aufrechterhaltung bestehender Vorurteile, Fairness der KI-Anwendung im Kontext der öffentlichen Verwaltung
- **4. Effektivität und Effizienz –** Einsatz von KI-Anwendungen in der Verwaltung muss deren Effektivität und Effizienz nachhaltig verbessern, ohne die Arbeitssituation der im öffentlichen Dienst tätigen Menschen zu verschlechtern



Kriterien im Leitfaden digitale Verwaltung II

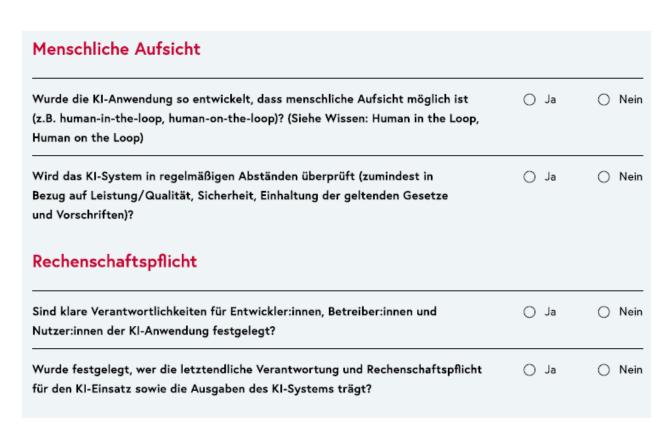


- Sicherheit KI-Anwendung muss sicher eingesetzt werden, Schutz sensibler Informationen, Vermeidung unbefugten Zugriffs
- 6. Barrierefreiheit und Inklusion KI-Anwendung muss für Menschen mit unterschiedlichen Fähigkeiten, Hintergründen und Kulturen zugänglich und integrativ sein, Angebot von Alternativen zur KI-Technologie für gleichberechtigten Zugang zu öffentlichen Dienstleistungen
- 7. Rechenschaftspflicht Klare Zuständigkeiten und Verantwortlichkeiten, Bewusstsein bei Verantwortlichen über ihre Verantwortung
- 8. Digitale Souveränität Die Verwaltung muss in der Lage sein die Entwicklung von KI-Lösungen zu beeinflussen, unabhängig anzuwenden und vertrauliche Daten in ihrem eigenen Einflussbereich zu halten



Checkliste im Leitfaden digitale Verwaltung

Basierend auf den Kriterien





- Vier Seiten mit Fragen zu den Bereichen
 - Recht,
 - Transparenz,
 - Unvoreingenommenheit und Fairness,
 - Effektivität und Effizienz,
 - Sicherheit,
 - Zugänglichkeit und Inklusion,
 - menschliche Aufsicht,
 - Rechenschaftspflicht,
 - digitale Souveränität.
- Geschlossene Fragen, die als Indikatoren der Erfüllung der wesentlichsten ethischen Kriterien gelten können.



Zum Beispiel: Fairness

Aktuelle Debatten rund um Fairness und KI:

- Unfaire Behandlung von Individuen oder Gruppen aufgrund von Verzerrungen in KI-Entscheidungen
- Undurchsichtigkeit von KI-Entscheidungen ("Black Box" KI)
- Verschiedene Gefahren für die Demokratie und das gesellschaftliche Wohlergehen
- Ungleichheiten auf dem Markt durch die Macht von Big Tech

Fairness im Al Act: "Diversity, non-discrimination and fairness means that Al systems are developed and used in a way that includes diverse actors and promotes equal access, gender equality and cultural diversity, while avoiding discriminatory impacts and unfair biases that are prohibited by Union or national law."

Herausforderung:

- Kein einheitliches Fairnesskonzept in einer pluralistischen Gesellschaft
- Technologie als Werkzeug, nicht als Allheilmittel



Verschiedene Perspektiven auf das Kriterium Fairness

- Anerkennung moralischer Gleichheit: Entscheidungen sollten gleiche Achtung und Würde für alle Menschen widerspiegeln
- Schutz vor **struktureller Diskriminierung**: Historische Ungleichheiten dürfen nicht verstärkt oder verfestigt werden
- Transparenz als Voraussetzung für Gerechtigkeit: Faire Entscheidungen erfordern Nachvollziehbarkeit und öffentliche Rechenschaft
- Verteilungsgerechtigkeit mit Blick auf Machtverhältnisse: Nutzen und Schaden sollte gerecht (z.B. über soziale Gruppen)
 verteilt sein
- Berücksichtigung kontextueller Gerechtigkeit Fairness ist nicht universell einheitlich: Was als fair gilt, muss sich an kulturellen und sozialen Kontexten orientieren



...und auf das Thema Fairness in Bezug auf Kl

- Demografische Parität: Gleiche positive Entscheidungsraten für alle Gruppen, unabhängig von sensiblen Attributen wie Geschlecht oder Ethnie
- Gleichheit der Fehlerraten: Gleiche Wahrscheinlichkeiten für False Positives und False Negatives zwischen Subgruppen, um faire Fehlerverteilungen zu gewährleisten
- Individuelle Fairness: Ähnliche Individuen erhalten ähnliche Entscheidungen
- Gleichheit der Chancen: Gleiche True Positive Rate über Gruppen hinweg
- Kontextualisierte Fairness: Berücksichtigt gesellschaftliche, rechtliche und kulturelle Rahmenbedingungen; ev. domänenspezifische Gerechtigkeitskonzepte berücksichtigen



Wie wird Fairness implementiert?



Ex-ante: System

Design

Ex-post: *Produkt Nutzungskontrolle*



Impact/Risk Assessment
Beteiligung von
Interessensgruppen

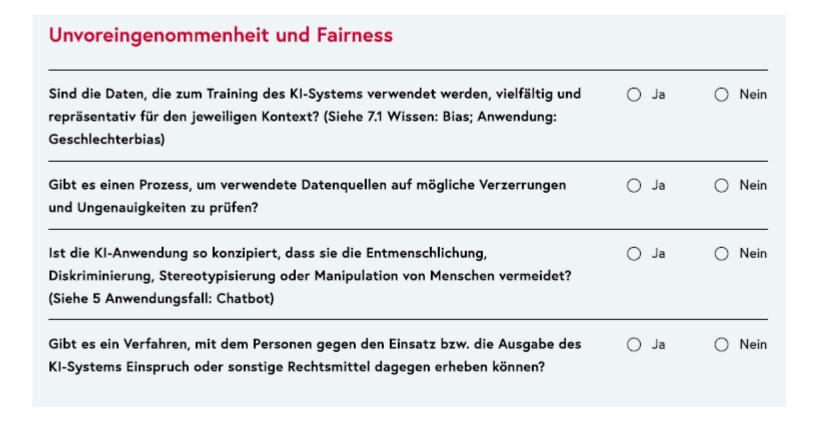
Qualitätsstandards

Fairnessmetriken



Und im Leitfaden digitale Verwaltung?

Checkliste: Kriterium Fairness







Fallstudie: Apple Pay

KI-gestütztes Risikomodell der Apple Card (entwickelt mit Goldman Sachs)

- Algorithmische Diskriminierung bei Kreditlimits zeigte systematische Verzerrungen: Frauen erhielten häufig deutlich niedrigere Kreditlimits als Männer bei vergleichbarer Bonität
- Intransparente Modelllogik & fehlende Erklärbarkeit: zugrundeliegende Entscheidungsalgorithmen waren nicht nachvollziehbar; Kundinnen erhielten keine Begründungen
- Regulatorische und ethische Defizite: Fehlende interne Fairness-Tests, unzureichende Compliance-Kontrollen führten Verletzungen des Equal Credit Opportunity Acts
- Lehren für KI-Governance: KI-Modelle ohne Fairness-Audits und Bias-Monitoring verdeutlichen Notwendigkeit von Explainable AI, Fairness-Metriken und menschlicher Aufsicht



Q & A

Funded by







Federal Ministry Innovation, Mobility and Infrastructure Republic of Austria



under discussion with

Al Factory Austria Al:AT has received funding from the European High-Performance Computing Joint Undertaking (JU) under grant agreement No 101253078. The JU receives support from the Horizon Europe Programm of the European Union and Austria (BMIMI / FFG).



Contact



Dr. Peter Biegelbauer

Co-Lead Legal, Regulatory and Ethics Al Factory Austria Al:AT

+43 664 88390033 peter.biegelbauer@ai-at.eu Al Factory Austria Al:AT Schwarzenbergplatz 2 1010 Wien, Austria

info@ai-at.eu ai-at.eu@ai-factory-austria

