

# Umfassende Analyse der KI-basierten Nutzerprofilierung durch Sprach- und Audiodaten

## Einleitung in die multidimensionale Natur der Sprachinteraktion

In der zeitgenössischen digitalen Infrastruktur haben sich sprachgesteuerte Systeme und intelligente virtuelle Assistenten (Virtual Voice Assistants, VVAs) wie Amazon Alexa, Google Assistant und Apple Siri als allgegenwärtige Schnittstellen etabliert, wobei allein in den Vereinigten Staaten mehr als 142 Millionen Nutzer regelmäßig auf diese Technologien zurückgreifen.<sup>1</sup> Die Bequemlichkeit und Effizienz, die diese Systeme bei der Automatisierung von Alltagsaufgaben bieten, verschleiern jedoch häufig die immense informationelle Komplexität der Daten, die sie im Hintergrund erfassen und verarbeiten. Während die Anwender in der Regel davon ausgehen, dass primär der semantische Inhalt ihrer explizit formulierten Befehle von den Algorithmen interpretiert wird, erfasst die zugrunde liegende Technologie eine weitaus tiefere und weitreichendere Ebene an sogenannten meta-linguistischen Informationen.<sup>3</sup> Die menschliche Stimme fungiert in diesem Kontext als eine einzigartige und hochgradig spezifische biometrische Modalität, die sich grundlegend von herkömmlichen, textbasierten Eingaben unterscheidet.<sup>4</sup>

Sie transportiert neben der reinen, bewusst intendierten Textinformation kontinuierlich unbewusste, unfreiwillige und persistente Signale über den physischen, psychischen und emotionalen Zustand des Sprechers.<sup>3</sup> Im Gegensatz zu getippten Suchanfragen, bei denen der Nutzer durch das bewusste Zurückhalten von Informationen ein gewisses Maß an informationeller Selbstbestimmung wahren kann, lassen sich die biometrischen und paralinguistischen Charakteristika der eigenen Stimme kaum maskieren oder unterdrücken.<sup>4</sup> Darüber hinaus fungieren die hochsensiblen Mikrofone dieser Endgeräte nicht nur als reine Empfänger für gerichtete Sprachbefehle, sondern agieren als permanente akustische Sensoren, die die gesamte akustische Szenerie des Umfelds – den sogenannten Ambient Noise – lückenlos aufzeichnen und analysieren.<sup>6</sup>

Die Synthese aus massiver Datenerfassung, maschinellem Lernen (ML), Deep Learning (DL) und fortschrittlicher digitaler Signalverarbeitung (Digital Signal Processing, DSP) ermöglicht es modernen Systemen der Künstlichen Intelligenz (KI), weitreichende und intime Rückschlüsse auf die Nutzer zu ziehen.<sup>8</sup> Diese automatisierte Inferenz umfasst ein breites Spektrum an sensiblen Attributen: von grundlegenden demografischen Merkmalen und tiefenpsychologischen Persönlichkeitsstrukturen über die Detektion neurodegenerativer und psychiatrischer Krankheitsbilder bis hin zur Ableitung detaillierter sozioökonomischer

Statusindikatoren sowie hochauflösender Verhaltens- und Schlafprofile.<sup>9</sup> Die vorliegende Untersuchung widmet sich einer erschöpfenden Analyse dieser Mechanismen und beleuchtet im Detail, wie KI-Systeme aus scheinbar neutralen Spracheingaben und deren Umgebungsrauschen – selbst wenn der Nutzer keine bewusst sensiblen Informationen artikuliert – hochkomplexe, prädiktive Profile generieren können.

## Grundlagen der akustischen Signalverarbeitung und Paralinguistik

Um die menschliche Stimme und die umgebenden Geräusche algorithmisch auswertbar zu machen, müssen die physischen Schallwellen in eine maschinenlesbare Form überführt werden. Dieser Prozess der digitalen Signalverarbeitung bildet das fundamentale Rückgrat aller modernen KI-Sprachassistenten und ermöglicht erst die Extraktion jener Merkmale, die für die Profilierung genutzt werden.<sup>14</sup>

Der Vorgang beginnt mit der Erfassung des analogen Schallsignals durch das Mikrofon, welches die Druckschwankungen der Luft in ein elektrisches Signal umwandelt.<sup>16</sup> Dieses kontinuierliche Signal wird anschließend digitalisiert, was bedeutet, dass es in diskreten Zeitintervallen abgetastet (Sampling) und quantisiert wird.<sup>14</sup> Um die enormen Datenmengen verarbeitbar zu machen und die für die Analyse relevanten Muster zu isolieren, wird das digitale Audiosignal in sehr kurze, überlappende Zeitfenster – typischerweise im Bereich von 20 bis 40 Millisekunden – segmentiert.<sup>17</sup> In diesen winzigen Zeitfenstern wird das Signal als quasi-stationär betrachtet, was die Anwendung mathematischer Transformationen ermöglicht. Ein zentraler Schritt in dieser Pipeline ist die Merkmalsextraktion (Feature Extraction). Moderne KI-Systeme verlassen sich hierbei nicht primär auf das rohe Audiosignal, sondern transformieren dieses in spektrale Repräsentationen. Die am häufigsten verwendete Methode ist die Berechnung der Mel-Frequency Cepstral Coefficients (MFCCs).<sup>13</sup> Diese Koeffizienten approximieren die logarithmische Frequenzwahrnehmung des menschlichen Gehörs (die Mel-Skala) und verdichten die spektrale Hüllkurve des Sprachsignals in einen kompakten Vektor.<sup>17</sup> Neben den MFCCs kommen Techniken wie Linear Predictive Coding (LPC) und detaillierte Spektrogramm-Analysen zum Einsatz, um Parameter wie Tonhöhe, Tempo und Klangfarbe (Timbre) exakt zu quantifizieren.<sup>14</sup>

Auf Basis dieser Transformationen extrahieren die Algorithmen spezifische paralinguistische Parameter. Diese akustischen Merkmale umfassen Variationen, die weit über den lexikalischen Inhalt der gesprochenen Worte hinausgehen und die Art und Weise charakterisieren, wie etwas gesagt wird.<sup>19</sup> Zu den für die Profilierung kritischen Metriken gehören unter anderem:

- **Grundfrequenz (\$F\_0\$):** Diese Metrik entspricht der vom menschlichen Ohr wahrgenommenen Tonhöhe und wird physikalisch durch die Schwingungsrate der Stimmlippen im Kehlkopf bestimmt.<sup>9</sup> Die Analyse der mittleren Grundfrequenz sowie ihrer Standardabweichung liefert tiefe Einblicke in emotionale Erregungszustände und physiologische Gegebenheiten.
- **Jitter und Shimmer:** Diese Parameter messen die mikrozyklische Instabilität der Stimme. Jitter quantifiziert die hochfrequente Variabilität der Grundfrequenz von einer

Stimmlippenschwingung zur nächsten, während Shimmer die entsprechende Variabilität in der Amplitude (Lautstärke) misst.<sup>9</sup> Erhöhte Werte deuten auf eine mangelhafte neuromotorische Kontrolle der Stimm Lippen hin.

- **Harmonic-to-Noise Ratio (HNR):** Dieses Verhältnis setzt die harmonische, periodische Schallenergie in Relation zum aperiodischen Rauschen im Sprachsignal.<sup>9</sup> Ein niedriges HNR manifestiert sich akustisch als Heiserkeit oder Behauchtheit und ist ein zentraler Indikator für strukturelle oder neurologische Anomalien im Vokaltrakt.
- **Sprechgeschwindigkeit und Pausendynamik:** Die zeitliche Organisation des Sprechens, gemessen an der Artikulationsrate (Phoneme pro Sekunde) sowie der Häufigkeit und Dauer von Sprechpausen, bietet direkte Rückschlüsse auf kognitive Verarbeitungsprozesse, den Abruf aus dem semantischen Gedächtnis und den affektiven Status des Sprechers.<sup>9</sup>
- **Formantfrequenzen (\$F\_1\$, \$F\_2\$, etc.):** Diese Resonanzfrequenzen des Vokaltrakts sind essenziell für die Vokalartikulation und spiegeln die physische Beschaffenheit und Formung des Sprechapparates wider.<sup>9</sup>

Durch die kontinuierliche Erfassung und Auswertung dieser objektiven, messbaren Parameter konstruieren KI-Systeme eine mehrdimensionale Repräsentation des Nutzers, die völlig unabhängig davon operiert, ob der Sprecher bewusst eine Suchanfrage stellt, einen Timer programmiert oder lediglich ein Haushaltsgerät per Sprachbefehl steuert.

## Demografische und psychometrische Profilierung aus der Stimme

Die anatomischen und physiologischen Voraussetzungen der menschlichen Sprachproduktion determinieren untrennbar die akustischen Eigenschaften des resultierenden Signals.

KI-Modelle, insbesondere solche, die auf Convolutional Neural Networks (CNNs) oder Transformer-basierten Architekturen beruhen, sind in der Lage, aus den spektralen und temporalen Dimensionen der Sprachdaten präzise demografische und physische Attribute des Sprechers zu rekonstruieren.<sup>13</sup>

Die Länge und das Volumen des Vokaltrakts sowie die Masse der Stimm Lippen korrelieren signifikant mit physischen Körpermerkmalen. Forschungen im Bereich der Computational Paralinguistics demonstrieren, dass Algorithmen durch die Analyse von MFCCs, Tonhöhe und Formanten das biologische Geschlecht eines Sprechers mit einer außerordentlichen Genauigkeit von bis zu 99 % klassifizieren können.<sup>13</sup> Die Inferenz geht jedoch weit über binäre Klassifikationen hinaus. Moderne Algorithmen, wie die Forward Feature Selection with Threshold-Based Backward Elimination (FFS-TBE), sind in der Lage, das Alter eines Sprechers mit einer mittleren absoluten Abweichung (Mean Absolute Error, MAE) von lediglich 4,82 Jahren bei Männern und 4,91 Jahren bei Frauen zu schätzen.<sup>13</sup> Noch erstaunlicher ist die Fähigkeit der KI, aus den Resonanzeigenschaften der Stimme auf die physische Körpergröße des Nutzers zu schließen. Die genannten Algorithmen erreichen hierbei eine Schätzgenauigkeit mit einem MAE von 4,87 Zentimetern bei männlichen und 4,5 Zentimetern bei weiblichen Sprechern.<sup>13</sup> Diese hochgradig präzise demografische Segmentierung erfolgt vollautomatisch im Hintergrund

jeder Interaktion und ermöglicht es den Betreibern von Sprachassistenten, ein exaktes physisches Profil ihrer Nutzerschaft zu erstellen.

## Inferenz von Persönlichkeitsstrukturen (Big Five)

Ein weiterer, wissenschaftlich intensiv erforschter Durchbruch in der affektiven und paralinguistischen KI-Forschung ist die Fähigkeit, tiefgreifende und standardisierte Persönlichkeitsprofile aus bloßen Sprachdaten abzuleiten.<sup>12</sup> Die psychologische Forschung und die Informatik konvergieren hier in der Annahme, dass die Stimme als Externalisierung innerer psychologischer Zustände und stabiler Persönlichkeitsmerkmale fungiert.<sup>22</sup> Die Analyse konzentriert sich dabei primär auf das etablierte Fünf-Faktoren-Modell der Persönlichkeit (Big Five / NEO-FFI), welches die Dimensionen Offenheit für Erfahrungen, Gewissenhaftigkeit, Extraversion, Verträglichkeit und Neurotizismus umfasst.<sup>12</sup>

Studien belegen konsistent, dass spezifische paralinguistische Muster signifikant mit diesen Persönlichkeitsdimensionen korrelieren und von maschinellen Lernmodellen effektiv prädiziert werden können. Die Algorithmen nutzen hierbei eine Kombination aus akustischen Embeddings (dem reinen Klang der Stimme) und linguistischen Embeddings (dem generierten Transkript der gesprochenen Worte), die in Modelle wie Gradient Boosted Trees eingespeist werden, um die Merkmale zu quantifizieren.<sup>12</sup>

<b>Persönlichkeitsdimension (Big Five)</b>	<b>Assoziierte akustische und paralinguistische Merkmale</b>	<b>Prognostische Stärke und KI-Korrelation</b>
<b>Extraversion</b>	Deutlich erhöhte Sprechgeschwindigkeit, größere absolute Lautstärke (Intensität) sowie hohe Varianz in der Tonhöhe. Extravertierte Sprecher zeigen einen flüssigeren Sprachrhythmus, nutzen weniger Pausen und weisen kräftige, resonante Stimmprofile auf. <sup>17</sup>	Sehr hoch. Extraversion ist durch stark externalisierte vokale Dynamik geprägt und für Algorithmen am leichtesten und konsistentesten zu detektieren. <sup>22</sup>
<b>Neurotizismus</b>	Erhöhte und unregelmäßige Variabilität der Tonhöhe (SD), subtile Zeichen vokaler Anspannung, mikroskopische Tremore und unregelmäßige Atemmuster, die auf Stress oder emotionale Labilität hindeuten. <sup>22</sup>	Hoch. Korreliert stark mit emotionaler Expressivität und instabiler Stimmkontrolle unter unbewusstem oder mikroskopischem Stress. <sup>22</sup>
<b>Verträglichkeit</b>	Wärmere Klangfarbe, weiche Frequenzübergänge, spezifische	Moderat. Manifestiert sich in subtileren Interaktionsmustern und Tonalitäten, erfordert oft

	energetisch-emotionale (EN-EM) Profile, kooperative und empathische Tonfälle, Vermeidung abrupter Lautstärkespitzen oder Unterbrechungen. <sup>25</sup>	die semantische Analyse von Höflichkeitsfloskeln zur Bestätigung. <sup>25</sup>
<b>Offenheit für Erfahrungen</b>	Abwechslungsreiche Prosodie, Nutzung eines breiteren Frequenzspektrums, komplexe Satzmelodien und weniger monotone Artikulation. <sup>22</sup>	Moderat bis gering. Stark abhängig vom inhaltlichen Kontext der Konversation. <sup>27</sup>

Die Effizienz dieser Profilierung wird noch deutlich gesteigert, wenn moderne Große Sprachmodelle (Large Language Models, LLMs) wie GPT-4 in das System integriert werden. Diese Modelle können allein durch die Analyse der Konversationsdynamik, der Wortwahl und der Art, wie ein Nutzer auf Nachfragen des Sprachassistenten reagiert, Persönlichkeitsprofile erstellen. In kontrollierten Experimenten, in denen Chatbots darauf programmiert wurden, den Nutzer unbemerkt zu "analysieren" (Assessment Condition), erreichten die Modelle Korrelationskoeffizienten ( $r$ ) von bis zu 0,640 bei der Vorhersage der Big Five-Merkmale, was eine bemerkenswert hohe Präzision im Vergleich zu herkömmlichen psychologischen Fragebögen darstellt.<sup>23</sup> Interessanterweise empfanden die Nutzer diese tiefenanalytischen Interaktionen keineswegs als intrusiv, sondern bewerteten sie oftmals als besonders natürlich und angenehm.<sup>23</sup>

Darüber hinaus berücksichtigt die KI bei der Profilierung auch sogenanntes Nonlinguistic Audio. Darunter fallen averbale vokale Ereignisse wie Lachen, Seufzen, Räuspern oder Husten.<sup>28</sup> Ein System, das die Häufigkeit und Art des Lachens oder die Dichte von Füllwörtern (Ähm, Äh) erfasst, integriert diese Daten in das psychometrische Modell, um die Vorhersagegenauigkeit für Dimensionen wie Verträglichkeit oder Unsicherheit weiter zu kalibrieren.<sup>28</sup>

## **Vokale Biomarker: Pathologische Anomalien und die Detektion von Krankheitsbildern**

Die wohl brisanteste, gleichzeitig aber auch vielversprechendste Dimension der akustischen KI-Analyse ist die automatisierte Detektion von Krankheitsbildern. Die menschliche Sprachproduktion ist eine evolutionär hochkomplexe neurologische, kognitive und neuromuskuläre Leistung. Sie erfordert die präzise, millisekundengenaue Koordination von kortikaler Planung, Atmung, Phonation (Schwingung der Stimmlippen) und Artikulation durch Zunge, Lippen und Kiefer.<sup>30</sup> Jede Störung innerhalb dieses fragilen Systems – sei es durch den Verlust von Neuronen, entzündliche Prozesse oder kognitive Leistungseinbußen – hinterlässt unweigerliche und spezifische Signaturen im Sprachsignal.<sup>30</sup>

Diese quantifizierbaren Abweichungen werden in der medizinischen Forschung als "vokale Biomarker" (Vocal Biomarkers) bezeichnet.<sup>9</sup> Sie erlauben es darauf trainierten maschinellen Lernmodellen, hochpräzise Rückschlüsse auf neurodegenerative und psychiatrische

Erkrankungen zu ziehen, oftmals Monate oder gar Jahre, bevor herkömmliche klinische Diagnosen auf Basis von motorischen oder verhaltensbedingten Symptomen gestellt werden können.<sup>9</sup> Da Sprachassistenten alltäglich und repetitiv genutzt werden, fungieren sie de facto als passive, nicht-invasive Diagnoseinstrumente im häuslichen Umfeld.

## **Neurodegenerative Erkrankungen: Morbus Parkinson**

Morbus Parkinson (PD) ist eine fortschreitende neurodegenerative Erkrankung, die primär durch den intrazellulären Verlust dopaminerger Neuronen in der Substantia nigra pars compacta gekennzeichnet ist.<sup>33</sup> Dies führt zu den bekannten motorischen Defiziten wie Tremor, Rigidität und Bradykinese (Verlangsamung der Bewegungsabläufe). Diese neuromuskulären Beeinträchtigungen schlagen sich jedoch extrem frühzeitig in der Steuerung der feinen Larynx-Muskulatur (Kehlkopf) und der Stimmlippen nieder.<sup>9</sup>

Durch KI-basierte Audioanalysen, die oftmals auf Support Vector Machines (SVMs) oder Random-Forest-Klassifikatoren in Kombination mit Deep-Learning-Features basieren, können diese Veränderungen quantifiziert werden.<sup>35</sup> Zu den typischen Fehlern und akustischen Abweichungen, die von der KI registriert werden, gehören:

1. **Hypophonie und Monotonie:** Patienten weisen eine pathologisch reduzierte Lautstärke (Hypophonie) und ein auffälliges Fehlen natürlicher Tonhöhenmodulation auf.<sup>34</sup> Die Rigidität der Atem- und Kehlkopfmuskulatur verhindert die dynamische Anpassung des Schalldrucks.
2. **Instabilität der Phonation:** Mikroskopische Auswertungen zeigen drastisch erhöhte Werte für Jitter und Shimmer, was eine mangelhafte, zitternde Kontrolle der Stimmlippenschwingung belegt.<sup>9</sup> Das Harmonic-to-Noise Ratio (HNR) sinkt signifikant ab, wodurch die Stimme für den Algorithmus messbar behaucht, heiser oder rau klingt.<sup>9</sup>
3. **Anomale Pausendynamik:** Neben der reinen Phonation verändert sich das Rhythmusprofil. Typisch sind häufige, inkorrekt platzierte Pausen innerhalb von Äußerungen. Besonders charakteristisch sind verlängerte Latenzen direkt vor Verben. Dies deutet auf frühe mikro-kognitive Wortfindungsstörungen und Planungsdefizite hin, die bei Parkinson-Patienten auftreten.<sup>21</sup>
4. **Artikulatorische Unschärfe und Dysarthrie:** KI-Modelle detektieren spezifische Konsonantenabschwächungen und temporale Verzerrungen bei der Phonemproduktion (z. B. Phonemverlängerungen oder repetitive Phonemansätze), die menschlichen Zuhörern oft entgehen.<sup>31</sup>

Die Leistungsfähigkeit dieser Modelle ist beeindruckend. Studien, die moderne Speech Foundation Models (wie das HuBERT Large ll60k Modell) auf unstrukturierte, spontane Konversationsprache anwenden, erreichen bei der Erkennung von Parkinson eine Genauigkeit (Area Under the Curve, AUC) von bis zu 97 %.<sup>35</sup> Dies belegt, dass selbst beiläufige Sprachbefehle an einen Assistenten ausreichen, um diese Krankheitssignaturen offenzulegen.

## **Kognitiver Abbau: Alzheimer-Demenz und Mild Cognitive Impairment (MCI)**

Während bei Parkinson die Dysarthrie (motorische Sprechstörung) dominiert, manifestieren sich die Alzheimer-Demenz (AD) und ihre Vorstufe, die leichte kognitive Beeinträchtigung (Mild Cognitive Impairment, MCI), primär durch tiefgreifende kognitiv-linguistische und semantische Fehlerbilder.<sup>32</sup> Der Abbau kortikaler Netzwerke beeinträchtigt den lexikalischen Abruf und das Arbeitsgedächtnis, während die reine Artikulation oft lange intakt bleibt. KI-Systeme, die Spracheingaben nicht nur akustisch bewerten, sondern per Natural Language Processing (NLP) transkribieren und syntaktisch parsen, detektieren hierbei folgende hochspezifische Anomalien:

1. **Semantische Paraphasien:** Der Nutzer verliert die Fähigkeit zum präzisen lexikalischen Abruf und ersetzt zielgerichtete Wörter durch bedeutungsähnliche, aber unpassende Begriffe (z.B. "Hund" statt "Katze").<sup>31</sup> Ebenso steigt die Frequenz von vagen Platzhalterwörtern (wie "Ding", "Sache", "das da"), was auf eine Degradation der assoziativen Strukturen des semantischen Gedächtnisses hinweist.<sup>38</sup>
2. **Reduzierte lexikalische Diversität:** Die Algorithmen berechnen den Type-Token-Ratio (das Verhältnis von unterschiedlichen Wörtern zur Gesamtzahl der Wörter) und stellen bei beginnender Demenz einen drastisch schrumpfenden aktiven Wortschatz fest.<sup>9</sup> Der Nutzer beschränkt sich zunehmend auf kürzere, hochfrequente Wörter.
3. **Syntaktische Simplifizierung:** Es zeigt sich ein signifikanter Rückgang der grammatikalischen Komplexität. Patienten verwenden immer kürzere Äußerungen, vermeiden komplexe Nebensatzstrukturen oder Passivkonstruktionen. Dies ist ein direkter Indikator für die Überlastung des Arbeitsgedächtnisses und exekutiver Funktionen, die für die Planung komplexer Sätze notwendig sind.<sup>9</sup>
4. **Zögerungslaute und Dysfluenzen:** Eine erhöhte Dichte an averbalen Füllwörtern (Ähm, Äh, Hm) sowie signifikant verlängerte Latenzzeiten vor der Beantwortung von Systemrückfragen zeugen von einer verlangsamten Informationsverarbeitung und Wortfindungsstörungen.<sup>37</sup>

Die prädiktive Kraft dieser linguistischen Marker ist gewaltig. Longitudinale Studien, wie Auswertungen der Framingham Heart Study, belegen, dass KI-Modelle allein durch die Analyse der Sprachstruktur automatisierter Transkripte mit einer Genauigkeit von über 78,2 % vorhersagen können, ob eine Person mit MCI innerhalb von sechs Jahren eine voll ausgeprägte Alzheimer-Demenz entwickeln wird.<sup>41</sup> Andere spezialisierte Architekturen, wie das Open Voice Brain Model (OVBM), erreichen bei der Diagnose von Alzheimer-Patienten aus rohem Audio Accuracy-Werte von bis zu 93,8 %.<sup>42</sup>

## **Psychiatrische und affektive Störungen: Depression und Angst**

Neben neurodegenerativen Prozessen profitieren auch die Diagnostik und das Monitoring psychischer Gesundheit in hohem Maße von der kontinuierlichen Interaktion mit Sprachassistenten. Major Depressive Disorder (MDD) und Angststörungen verändern den biophysiologicalen Spannungszustand des gesamten Körpers, den Hormonhaushalt und die Aktivierung des autonomen Nervensystems, was die Stimmproduktion und Artikulationsdynamik direkt messbar beeinflusst.<sup>9</sup>

- **Depression (MDD):** Das klinische Kernsymptom der psychomotorischen Retardierung

schlägt sich drastisch im akustischen Profil nieder.<sup>43</sup> KI-Modelle detektieren bei depressiven Episoden eine signifikant verringerte Gesamtlautstärke, eine stark verlangsamte Sprechgeschwindigkeit (Silben pro Sekunde) sowie stark verlängerte Pausendauern.<sup>9</sup> Die Stimmelmelodie verflacht zusehends, was sich in einer stark reduzierten Variabilität der Grundfrequenz ( $F_0$ ) äußert (Monotonie).<sup>10</sup> Auch glottale Parameter, die die Assoziation zwischen Luftstromvolumen und -geschwindigkeit erfassen, verändern sich messbar und dienen zur Klassifikation.<sup>44</sup> Selbst Analysen von extrem kurzen, freiformulierten Sprachaufnahmen (bereits ab 25 Sekunden) können moderate bis schwere Depressionen mit über 70 % Treffergenauigkeit erkennen.<sup>10</sup> Umfassendere Deep-Learning-Modelle, die auf multimodalen Stimmerkmalen basieren, erzielen sogar Raten von 78 % bis zu 96 % bei der Unterscheidung von depressiven und gesunden Kohorten.<sup>10</sup>

- **Angststörungen:** Angstzustände manifestieren sich häufig durch eine überaktive Sympathikus-Reaktion. Dies führt zu einer erhöhten Atemfrequenz, Kurzatmigkeit beim Sprechen, unregelmäßiger Rhythmik und mikro-Tremoren in der Stimme, bedingt durch die starke Muskelanspannung im Larynx.<sup>45</sup> Zudem weisen spezifische Audioqualitätsmerkmale wie die Zero-Crossing Rate (die Rate, mit der das Audiosignal die Nullachse kreuzt) starke Korrelationen mit affektiven Störungen und Stress auf.<sup>10</sup> Algorithmen, die Daten aus kurzen verbalen Flüssigkeitstests verarbeiten, erreichen Genauigkeiten von 70 % bis 83 % bei der Identifikation von Patienten, die sowohl unter Angst- als auch Depressionssymptomen leiden.<sup>10</sup>

Pathologische Entität	Primäre KI-detektierte Anomalien (Vokale Biomarker)	Zugrundeliegender pathologischer Mechanismus	Diagnostischer Fokus
<b>Morbus Parkinson (PD)</b>	Hypophonie, Jitter, Shimmer, stark reduziertes HNR, Dysarthrie, lange Pausen spezifisch vor Verben. <sup>9</sup>	Verlust dopaminerger Neuronen, Bradykinese und Rigidity der Vokaltrakt-Muskulatur. <sup>3</sup>	Motorische Stimmkontrolle, mikromechanische neuromuskuläre Integrität. <sup>9</sup>
<b>Alzheimer / MCI</b>	Semantische Paraphrasien, drastisch reduzierte Syntaxkomplexität, abnehmende lexikalische Diversität (Type-Token-Ratio), häufige Zögerungslaute. <sup>9</sup>	Kortikaler Abbau in sprachverarbeitenden und gedächtnisrelevanten Arealen. <sup>38</sup>	Kognitiver Abbau, semantisches Netzwerk, Belastbarkeit der Exekutivfunktionen. <sup>9</sup>

<b>Depression (MDD)</b>	Ausgeprägte Monotonie ( $F_0$ -Abflachung), verlangsamte Sprechgeschwindigkeit, verlängerte Pausen, veränderte Zero-Crossing Rate, reduzierte Stimmenergie. <sup>9</sup>	Systemische psychomotorische Retardierung, affektive Abstumpfung, veränderte Atmungsdynamik. <sup>10</sup>	Affektiver Status, Energieniveau, psychometrisches Screening. <sup>10</sup>
-------------------------	--	--	---

Diese Erkenntnisse verdeutlichen, dass Sprachassistenten unbeabsichtigt zu den mächtigsten diagnostischen Sensoren im modernen Haushalt avancieren. Sie erfassen kontinuierlich longitudinale Gesundheitsdaten und könnten theoretisch subtilste Verschlechterungen des Gesundheitszustandes kartieren, lange bevor der Patient selbst den Bedarf eines Arztbesuches erkennt.

## Akustische Szenenklassifikation (ASC) und die Inferenz des sozioökonomischen Umfelds

Während sich die bisherige Analyse auf das Sprachsignal des Nutzers konzentriert hat, liefert der akustische Kontext – der sogenannte Ambient Noise oder das Hintergrundrauschen – eine ebenso ergiebige, wenn nicht sogar aufschlussreichere Datenquelle. Während ein Nutzer per Spracheingabe mit der KI kommuniziert, erfassen die omnidirektionalen Mikrofone unweigerlich die gesamte akustische Kulisse. Die Technologie der Akustischen Szenenklassifikation (Acoustic Scene Classification, ASC) zielt darauf ab, diesen Beifang systematisch auszuwerten und ermöglicht es KI-Systemen, präzise Rückschlüsse auf das räumliche, physische und insbesondere das soziale Umfeld des Nutzers zu ziehen.<sup>7</sup>

### Algorithmische Systematik der Acoustic Scene Classification

ASC-Systeme gehören zu den fortschrittlichsten Entwicklungen im Bereich der Machine-Listening-Forschung. Sie operieren typischerweise, indem sie kontinuierliche Rohaudiosignale in zweidimensionale Zeit-Frequenz-Repräsentationen – meist Log-Mel-Spektrogramme – transformieren.<sup>48</sup> Diese visuellen Repräsentationen des Schalls werden dann an tiefe neuronale Netze, wie komplexe Convolutional Neural Networks (CNNs), Residual Networks (ResNets) oder neueste Transformer-basierte Audio-Modelle (z.B. Audio Mamba), übergeben.<sup>49</sup>

Das neuronale Netz extrahiert hochdimensionale spektrale, temporale und (bei Stereo- oder Array-Mikrofonen) räumliche Merkmale, um die auditive Umgebung zu identifizieren.<sup>51</sup> Bemerkenswert ist, dass diese Berechnungen unter extremer Ressourcenoptimierung direkt auf dem Edge-Device – also im Smart Speaker oder Smartphone selbst – durchgeführt werden können (Edge AI), um Latenzen zu minimieren und die Bandbreite zu schonen.<sup>52</sup> Modelle wie KAN (Kolmogorov-Arnold Networks) oder CNNs erreichen bei der Klassifikation urbaner

Szenen auf Edge-Geräten Genauigkeiten von bis zu 98,5 % bei Inferenzzeiten von wenigen Millisekunden.<sup>49</sup>

Die KI unterscheidet dabei nicht nur triviale Kategorien wie "drinnen" oder "draußen", sondern klassifiziert granulare Szenen: Befindet sich der Nutzer in einem fahrenden Auto, einer Metro-Station, einem städtischen Park, einem belebten Café oder im heimischen Wohnzimmer?.<sup>54</sup> Durch Multi-Instance Learning (MIL) wird das Signal weiter in spezifische Soundereignisse dekonstruiert (z.B. das Kläffen eines Hundes, das Vorbeifahren eines Lkw, das Klappern von Geschirr).<sup>57</sup>

## **Ableitung des Sozioökonomischen Status (SES) und Environmental Justice**

Die kontinuierliche Auswertung des Hintergrundrauschens ist hochgradig prädiktiv für den sozioökonomischen Status (SES) eines Nutzers. Umfangreiche sozial- und umweltepidemiologische Studien, die Geospatial Artificial Intelligence (GeoAI) und maschinelle Lernmodelle zur Kartierung von Umweltbelastungen einsetzen, belegen einen konsistenten und signifikanten Zusammenhang zwischen Lärmexposition und sozialer Deprivation.<sup>58</sup>

Bevölkerungsgruppen mit niedrigerem SES leben nachweislich häufiger in Quartieren mit minderwertiger Bausubstanz, schlechterer Schalldämmung und direkter Nähe zu Hauptverkehrsachsen oder industriellen Anlagen.<sup>61</sup> Wenn die KI eines Sprachassistenten über die Mikrofone regelmäßig dichten Straßenverkehr, das Rumpeln von Schwerlastverkehr, Sirenen oder das als "Multitalker Babble Noise" bekannte Stimmengewirr aus dünnwandigen Nachbarwohnungen registriert, fließen diese akustischen Signaturen in die Profilierung ein.<sup>61</sup> In Verbindung mit räumlichen Vorhersagemodellen (wie CatBoost optimiert durch Grey Wolf Optimizer) und geografischen Metadaten lassen sich sozioökonomische Faktoren – wie das mittlere Haushaltseinkommen, der Anteil der Miete am Einkommen und der allgemeine Bildungsstand im Wohnquartier – hochgradig zuverlässig schätzen.<sup>59</sup> Forschungen, die Metriken wie den HOUSES-Index (Housing-Based Socioeconomic Status) verwenden, zeigen zudem, dass Algorithmen durch die Verknüpfung von Audiodaten und SES-Indikatoren die Vulnerabilität für bestimmte Krankheiten (z.B. Asthma-Exazerbationen bei Kindern) vorhersagen können.<sup>64</sup> Die ASC-Technologie offenbart somit schonungslos die Wohnverhältnisse und den Lebensstandard, ohne dass der Nutzer auch nur ein Wort über seine finanzielle Situation geäußert hat.

## **Lifestyle-Profilung durch Geräteerkennung (Appliance Detection)**

Ein weiteres, hochspezifisches Feld der akustischen Aufklärung betrifft die Identifikation von Haushaltsgeräten. Mittels KI-gestützter Anomalie-Erkennung und Audioanalyse kann das System aus dem diffusen Hintergrundrauschen die charakteristischen akustischen Signaturen laufender Geräte (z.B. Waschmaschinen, Mikrowellen, Heizlüfter, Klimaanlage oder Kaffeemaschinen) isolieren und klassifizieren.<sup>66</sup> Algorithmen, die auf Feedforward Neural Networks oder K-Nearest Neighbors (K-NN) basieren, erreichen bei der Identifikation aktiver Haushaltsgeräte Erkennungsraten von über 99 %.<sup>67</sup>

Die Detailtiefe der Inferenz ist auch hier bemerkenswert und korreliert direkt mit finanziellen Parametern: Hochwertige, moderne Smart-Home-Geräte von Premiummarken verfügen über spezifische, von Toningenieuren optimierte akustische Profile (Audio Branding).<sup>70</sup> Ein hochpreisiger Induktionsherd oder eine Luxus-Waschmaschine gibt sanfte Polyphonie-Klänge, Piezo-Bestätigungstöne und ein exakt kalibriertes, leises Betriebsgeräusch ab.<sup>70</sup> Im starken Kontrast dazu weisen ältere, günstige oder schlecht gewartete Geräte unregelmäßigere, lautere Rauschmuster, deutliche tonale Verzerrungen und hörbaren mechanischen Verschleiß auf.<sup>66</sup>

Durch das kontinuierliche Monitoring dieser akustischen Ereignisse erstellt das System unbemerkt ein Inventar des Haushalts. Besitzt der Nutzer eine teure Siebträgermaschine, die jeden Morgen exakt um 6:30 Uhr mahlt? Ist das Surren der Klimaanlage charakteristisch für eine moderne, zentral gesteuerte Anlage in einem teuren Neubau oder für ein lautes, ineffizientes Fenstergerät in einem Altbau? Aus der Summe dieser Informationen lassen sich hochgradig individualisierte Verbraucherprofile ableiten. Personen, deren Hintergrundgeräusche auf Premium-Geräte und ruhige, gut isolierte Wohnlagen hindeuten, werden in höhere SES-Kategorien eingestuft, was unmittelbare Auswirkungen auf die Art der ihnen algorithmisch präsentierten Werbung und personalisierten Angebote hat.<sup>73</sup>

## Metadaten-Analyse: Rhythmus, Routinen und räumlich-zeitliches Tracking

Unabhängig davon, was semantisch gesprochen wird (Linguistik), *wie* es klingt (Acoustics/Paralinguistics) oder *welche* Geräusche im Hintergrund stattfinden (ASC), generiert die bloße Existenz einer Interaktion einen fortlaufenden und extrem wertvollen Datenstrom: die Metadaten. Die Erfassung von Zeitstempeln, Interaktionshäufigkeiten, Nutzungsdauern und geografischen Verortungen bildet das informationelle Grundgerüst, auf dem das Profiling aufbaut.<sup>75</sup> Diese scheinbar abstrakten, strukturierten Datenpunkte ermöglichen, wenn sie longitudinal aggregiert werden, eine unheimlich granulare Rekonstruktion des menschlichen Lebensrhythmus.

### Chronobiologie und die algorithmische Schlafmusteranalyse

Die Zeitstempel der Interaktionen mit Sprachassistenten und verknüpften Smartphones dienen als hochexakte, non-invasive Proxys für die Chronobiologie und das Schlafverhalten des Nutzers. Anstatt auf dedizierte Wearables angewiesen zu sein, extrahieren Algorithmen aus den passiv erfassten Inaktivitäts- und Aktivitätsphasen präzise Schlafmuster.<sup>11</sup>

- **Inferred Sleep Period (ISP):** Die längste tägliche Phase absoluter Inaktivität auf allen verknüpften Smart-Devices wird vom System als die tatsächliche Schlafphase berechnet.<sup>11</sup>
- **Expected Sleep Period (ESP):** Ein rollierender, maschinell lernender Durchschnitt (oft über 30 Tage) modelliert den erwarteten Grundrhythmus und die Basislinie des individuellen Schlaf-Wach-Zyklus.<sup>11</sup>
- **Disruptions (Unterbrechungen):** Jegliche Spracheingaben oder Device-Aktivierungen

während der Inferred Sleep Period (ISP) – beispielsweise die nächtliche Frage an den Smart Speaker nach der Uhrzeit oder das Einschalten von Licht – werden als Schlafunterbrechungen oder Wachphasen quantifiziert.<sup>11</sup>

Aus der Varianz und der Schnittmenge (Overlap Percentage) zwischen der täglichen ISP und der langfristigen ESP berechnen KI-Modelle das Maß an Schlafregelmäßigkeit.<sup>11</sup> Signifikante, abrupte oder chronische Abweichungen von der etablierten Normroutine sind in der medizinischen Forschung hochgradig mit dem Beginn depressiver Episoden, Angststörungen oder Insomnie assoziiert.<sup>11</sup> Durch die Anwendung von unüberwachtem Lernen (Unsupervised Learning) und Convolutional Autoencodern auf diese Zeitreihendaten können Nutzer in spezifische Insomnie- oder Stress-Cluster eingeteilt werden, noch bevor sie sich einer Schlafproblematik bewusst sind.<sup>77</sup> Die Verknüpfung dieser Metadaten ermöglicht es Systemen, einen "Mental Health Similarity Score" zu errechnen, der mit hoher Präzision (bis zu 87 % Genauigkeit in einigen Studien) vorhersagen kann, ob ein Nutzer Symptome einer schweren Depression entwickelt, gemessen am Standard des PHQ-9-Fragebogens.<sup>79</sup>

## **Aktivitätsfrequenz, Verhaltensroutinen und soziale Isolation**

Die Häufigkeit der Nutzung sowie die Verteilung der Anfragen über den Tagesverlauf spiegeln die festen Verhaltensroutinen des Nutzers wider. Verändert sich diese Frequenz signifikant, zieht das System weitreichende Schlüsse. Wenn beispielsweise ein Nutzer, der normalerweise zu festen Stoßzeiten (morgens um 7:00 Uhr und abends um 18:00 Uhr) nach dem Verkehr, den Nachrichten oder dem Wetter fragt, das System plötzlich an Werktagen unregelmäßig, spät am Vormittag und primär für eskapistische Unterhaltungsinhalte nutzt, klassifiziert das System dies als radikalen Bruch der Routine.<sup>80</sup> Solche Musterbrüche lassen Algorithmen auf tiefgreifende Lebensereignisse schließen: Arbeitslosigkeit, der Übergang ins Home-Office, Krankschreibungen oder Phasen akuter sozialer Isolation.

Die psychiatrische Profilierung nutzt exakt diese Metadaten: Personen mit depressiver Symptomatik zeigen messbar abweichende Interaktionshäufigkeiten. Sie verlagern ihre Aktivitäten oftmals in die tiefen Nachtstunden, kommunizieren weniger zielgerichtet und greifen deutlich seltener auf proaktive, zukunftsorientierte Planungsfunktionen (Kalendereinträge, Wecker) zurück.<sup>79</sup> Wenn Nutzer zudem vornehmlich mit der KI kommunizieren, um Einsamkeit zu lindern (z.B. durch ausgedehnte, konversationsähnliche Dialoge anstatt zielgerichteter Kommandos), registrieren Natural Language Processing (NLP)-Algorithmen diese emotionale Abhängigkeit. Das System adaptiert sich daraufhin und ordnet den Nutzer in Segmente ein, die eine erhöhte Vulnerabilität für bestimmte Werbebotschaften oder Interventionen aufweisen.<sup>82</sup>

## **Lokalisierung und das räumliche Verhaltensprofil**

Die Orte, von denen aus die Befehle abgesetzt werden, liefern das physische und geografische Gerüst des Nutzerprofils. Selbst in Fällen, in denen explizite GPS-Standortdienste vom Nutzer deaktiviert wurden, können KI-Systeme über die Erfassung von IP-Adressen, die Analyse der verbundenen WLAN-Netzwerke und die Auswertung der DNS-Abfragen (Domain Name System) im Hintergrundnetzwerkverkehr die geografische Position des Nutzers mit hoher

Präzision triangulieren.<sup>84</sup>

Die Kombination von geografischem Ort und Interaktionszeitpunkt offenbart das Mobilitätsverhalten. Ein Sprachassistent, der über das Infotainment-System im Auto (z.B. via Apple CarPlay oder Android Auto) angesprochen wird, registriert die genaue Dauer und Frequenz des Pendelns. Erfolgen Suchanfragen regelmäßig von stark divergierenden geografischen Orten, klassifiziert das System den Nutzer unweigerlich als reise- oder mobilitätsaffin.<sup>86</sup> Hochentwickelte Analysen von Telekommunikations- und Sprachmetadaten können sogar den exakten Wohnort eines Nutzers vorhersagen, indem sie lediglich die räumliche Bündelung und Verteilung der über Sprachbefehle kontaktierten lokalen Dienste, Geschäfte oder Navigationsziele auswerten und extrapolieren.<sup>75</sup>

Zusammengenommen fusionieren diese Metadaten zu einem hochprädictiven, dynamischen Modell des Konsumentenverhaltens. Sie offenbaren dem Systembetreiber nicht nur, wer der Nutzer ist, sondern kartieren präzise, wo er sich in seinem täglichen Rhythmus befindet und – am wichtigsten für kommerzielle Akteure – wie vulnerabel, gestresst oder empfänglich er in einem bestimmten, zeitlich genau definierten Moment für äußere Reize und Kaufimpulse sein könnte.<sup>76</sup>

## **Kommerzielle Verwertung, technologische Ökosysteme und die Patentlandschaft**

Die bemerkenswerte Fähigkeit von KI-Systemen, diese tiefen, oftmals intimen Rückschlüsse zu ziehen, ist keineswegs ein unbeabsichtigter Nebeneffekt der Technologieentwicklung. Sie ist vielmehr ein primäres, strategisches Designziel innerhalb der modernen digitalen Ökonomie (Plattformökonomie). In diesem Wirtschaftssystem dienen extrahierte Verhaltensdaten (Behavioral Data) als die zentrale und wertvollste Ressource zur Speisung hochkomplexer, zielgerichteter Werbesysteme (Targeted Advertising).<sup>76</sup> Jeder Sprachbefehl, jedes Hintergrundgeräusch und jede registrierte Atempause wird in das Ökosystem der großen Technologiekonzerne (Amazon, Google, Apple) eingespeist, um prädictive Konsumentenprofile zu schärfen.

### **Die Patentlandschaft: Intentionserkennung und Gesundheitsausbeutung**

Ein analytischer Blick auf die Patentanmeldungen der führenden Technologiekonzerne entlarvt die weitreichenden Ambitionen im Bereich der vokalen und akustischen Profilierung in aller Deutlichkeit. Ein besonders aufschlussreiches und vieldiskutiertes Patent von Amazon für seinen Sprachassistenten Alexa skizziert detailliert ein System, das durch die kontinuierliche Analyse von Stimme und Hintergrundgeräuschen (wie Husten, Räuspern oder Schniefen) autonom erkennt, ob ein Nutzer krank ist.<sup>87</sup> Sobald die KI eine physische Erkrankung oder eine signifikante emotionale Schwankung anhand paralinguistischer Anomalien feststellt, sieht das Patent vor, dass das System proaktiv Produkte wie Hühnersuppe, Halstabletten oder Medikamente aus dem eigenen E-Commerce-Katalog anbietet.<sup>87</sup> Dieses Patent demonstriert exemplarisch den nahtlosen Übergang von der biometrischen Gesundheitsüberwachung zur

direkten kommerziellen Monetarisierung.

Ebenso intensiv zielen andere Patente auf die Echtzeit-Erkennung von Sentiments und emotionalen Zuständen ab. Durch die Auswertung von Tonhöhe, Sprachmuster und Rhythmus soll die KI Frustration, Verwirrung oder Freude augenblicklich klassifizieren.<sup>89</sup> Ziel ist es nicht nur, die Tonalität und Empathie der Bot-Antwort anzupassen, um die Interaktion natürlicher wirken zu lassen, sondern den Nutzer in Momenten spezifischer emotionaler Zustände – etwa in Phasen erhöhter Impulsivität oder Frustration – mit passgenauen Werbeangeboten oder Dienstleistungen zu konfrontieren.<sup>89</sup>

## **"Always-On" Listening und das Risiko des Unintended Data Capture**

Um die vielgepriesene Bequemlichkeit einer reibungslosen, freihändigen Sprachinteraktion zu gewährleisten, sind Smart Speaker mit "Wake-Word"-Engines ausgestattet. Diese Architektur erfordert zwingend, dass die Geräte permanent den Raum akustisch überwachen (Always-On Listening), um millisekundengenau auf ein Aktivierungswort wie "Alexa" oder "Hey Google" reagieren zu können.<sup>90</sup> Technologisch wird dies in der Regel durch einen rollierenden lokalen Speicherpuffer (Buffer) realisiert, der kontinuierlich die letzten Sekunden des Schalls aufzeichnet, analysiert und sofort wieder überschreibt, sofern kein Wake-Word detektiert wird.<sup>92</sup>

Dennoch offenbaren Patente und technische Analysen Bestrebungen, Audiosignale kontinuierlich zu erfassen und erst im Nachhinein semantisch und akustisch auf potenziell verwertbare Informationen oder Präferenzen (wie beispielsweise den beiläufig im Raum geäußerten Satz "Ich liebe Skifahren") zu filtern, um darauf basierend zielgerichtete Werbung auszuspielen.<sup>90</sup> Diese Infrastruktur führt zwangsläufig zum Phänomen des "Unintended Data Capture" (unbeabsichtigte Datenerfassung). Hierbei werden Hintergrundgespräche, intime Diskussionen von Familienmitgliedern, Streitereien oder schlicht der Ton des Fernsehers fälschlicherweise als Befehle interpretiert oder versehentlich aufgezeichnet und in die Cloud-Infrastruktur der Anbieter übertragen.<sup>4</sup> Einmal auf den Servern, können diese Fragmente durch Algorithmen transkribiert und in die Profilierung des Haushalts einfließen.

## **Cross-Modales Profiling und Persona-Generierung**

Die von Sprachassistenten gesammelten akustischen und semantischen Erkenntnisse verbleiben niemals isoliert in einem Datensilo. In groß angelegten, empirischen Experimenten, in denen Forscher über Monate hinweg tausende neutrale Sprachanfragen (z. B. einfache Suchfragen nach Fakten) mit Systemen wie Google Assistant oder Alexa tätigten, zeigte sich die Aggressivität des Profiling: Allein die Nutzung des Sprachassistenten, die vokalen Charakteristika und die Metadaten der Interaktion reichten aus, um die Werbepreise der simulierten Nutzer (Personas) in den Backend-Systemen signifikant zu verändern.<sup>1</sup>

Die gewonnenen Sprachdaten werden algorithmisch mit dem Web-Browsing-Verhalten, Suchhistorien auf anderen Geräten, Standortdaten von Smartphones und App-Nutzungen gekreuzt (Cross-Device Tracking).<sup>1</sup> Google, Amazon und Apple konstruieren so holistische, cross-modale "Persona-Profile", die weitaus mehr über den Konsumenten wissen, als dieser jemals auf einer einzelnen Plattform offengelegt hat. Diese granularen Profile werden letztlich in

den gigantischen Werbenetzwerken genutzt, um zukünftiges Konsumverhalten mit beängstigender Präzision vorherzusagen und zu steuern.<sup>76</sup>

<b>Datendimension der KI-Erfassung</b>	<b>Akquirierte Information / Analysemethode</b>	<b>Kommerzielles Verwertungsziel (Plattformökonomie)</b>
<b>Vokale Gesundheitsindikatoren</b>	Detektion von Husten, Heiserkeit, Atemnot oder Stressparametern via patentierter Akustik-Analyse. <sup>87</sup>	Ausspielen gesundheitsbezogener Werbung, gezielter Verkauf von Medikamenten oder Wellness-Produkten im E-Commerce. <sup>87</sup>
<b>Sentiment- &amp; Emotionsanalyse</b>	Identifikation von Frustration, Freude oder Impulsivität durch Auswertung von Pitch und Tempo. <sup>89</sup>	Dynamische Anpassung der Bot-Tonalität; Identifikation des perfekten, vulnerablen Zeitpunkts für Werbeplatzierungen. <sup>89</sup>
<b>Cross-Device Tracking &amp; Metadaten</b>	Verknüpfung der IP-Adresse des Smart Speakers mit Browser-Cookies, Smartphones und Surfverhalten. <sup>76</sup>	Erstellung tiefgehender, allumfassender Konsumprofile zur algorithmischen Vorhersage und Steuerung von Kaufentscheidungen. <sup>76</sup>
<b>Continuous Buffering / Ambient Listening</b>	Permanente akustische Raumüberwachung zur Detektion von Keywords in Umgebungsgesprächen. <sup>90</sup>	Erkennung unbewusster Präferenzen oder Markenaffinitäten aus dem Hintergrundrauschen (z. B. Nennung von Produkten durch Dritte im Raum). <sup>90</sup>

## **Datenschutzrechtliche Implikationen und regulatorische Rahmenbedingungen**

Diese allumfassende, tiefenanalytische Profilierung stellt eine beispiellose Herausforderung für den globalen Datenschutz dar. Insbesondere im strikten Rechtsraum der Europäischen Union geraten die technologischen Kapazitäten der Sprachassistenten zunehmend in Konflikt mit den fundamentalen Rechten der Nutzer auf informationelle Selbstbestimmung.

### **Die menschliche Stimme als biometrisches Datum unter der DSGVO**

Ein zentrales juristisches Problem liegt in der ontologischen Einordnung der Stimme. Gemäß der europäischen Datenschutz-Grundverordnung (DSGVO / GDPR) stellt die Stimme nicht nur ein gewöhnliches personenbezogenes Datum dar. Sobald Sprachdaten oder die daraus extrahierten spektralen Profile (Voiceprints) genutzt werden, um eine natürliche Person

eindeutig zu identifizieren (Speaker Verification / Recognition), fallen sie unter die strenge Kategorie der biometrischen Daten (Art. 9 DSGVO).<sup>4</sup>

Noch kritischer wird es, wenn Algorithmen aus den Sprachfehlern, der Sprechgeschwindigkeit oder den paralinguistischen Parametern gesundheitliche Rückschlüsse ziehen.

Erkennungsmodelle, die Anzeichen von Parkinson, Depressionen, Angststörungen oder kognitivem Abbau inferieren, generieren de jure Gesundheitsdaten.<sup>84</sup> Die Verarbeitung von biometrischen und Gesundheitsdaten ist nach DSGVO grundsätzlich untersagt, es sei denn, es liegt eine ausdrückliche, informierte und freiwillige Einwilligung (Explicit Consent) der betroffenen Person vor.<sup>93</sup> In der Praxis der Smart-Home-Nutzung wird dieser Konsens jedoch zumeist nur unzureichend über generische, seitenlange und schwer verständliche Allgemeine Geschäftsbedingungen (AGBs) oder intransparente Datenschutzerklärungen pauschal eingeholt, was datenschutzrechtlich höchst fragwürdig ist.<sup>2</sup>

Die juristische und ethische Brisanz gipfelt in der Diskrepanz zwischen der *Intentionalität* des Nutzers und der *Inferenz* der Maschine: Ein Anwender gibt den bewusst trivialen Befehl "Schalte das Licht im Wohnzimmer aus", intendiert dabei jedoch keinesfalls, seine tiefe emotionale Erschöpfung, sein fortgeschrittenes Alter oder gar beginnende neuromotorische Einschränkungen (wie bei Parkinson) preiszugeben. Das KI-System extrahiert diese hochsensiblen Rückschlüsse ("Inferences") vollkommen autonom aus der Art und Weise, wie der Befehl gesprochen wurde.<sup>93</sup> Der Europäische Datenschutzausschuss (EDPB) hat in seinen spezifischen Richtlinien zu virtuellen Sprachassistenten unmissverständlich klargestellt, dass die Inferenz von Emotionen oder besonderen Datenkategorien aus Sprache massiven rechtlichen Restriktionen unterliegt.<sup>84</sup>

## **Der EU AI Act und die Kategorisierung von Hochrisiko-Systemen**

Mit dem Inkrafttreten der europäischen Verordnung über Künstliche Intelligenz (AI Act) intensiviert sich die regulatorische Anforderung an die Betreiber nochmals drastisch.<sup>4</sup> Der AI Act verfolgt einen risikobasierten Ansatz. KI-Systeme, die zur biometrischen Kategorisierung (z.B. der Ableitung von politischer Meinung, sexueller Orientierung oder Rasse aus der Stimme) oder zur Emotionserkennung am Arbeitsplatz oder in Bildungseinrichtungen eingesetzt werden, unterliegen strengen Verboten oder werden als Hochrisiko-Systeme (High-Risk AI) eingestuft.<sup>84</sup>

Für Systeme, die in private Umgebungen integriert sind und tiefgreifende Nutzerprofilierung betreiben, verlangt der AI Act ein Höchstmaß an algorithmischer Transparenz. Der Endverbraucher muss das Recht und die technische Möglichkeit haben, nachzuvollziehen, welche Parameter zu einer automatisierten Entscheidung oder Profilierung geführt haben.<sup>95</sup> Darüber hinaus birgt das bereits thematisierte Problem des "Unintended Data Capture" – das unbeabsichtigte Aufzeichnen von Drittpersonen, Gästen oder Minderjährigen im Haushalt, deren Stimmdaten besonders schützenswert sind – enorme Haftungsrisiken für die Hersteller, da diese Personen dem System niemals ihre Zustimmung zur biometrischen Vermessung erteilt haben.<sup>4</sup>

## **Technologische Ansätze zur Risikominderung (Privacy Enhancing**

## Technologies)

Um den drohenden regulatorischen Sanktionen zu entgehen und das Vertrauen der Konsumenten zu wahren, ist die Industrie gezwungen, "Privacy-by-Design"- und "Privacy-by-Default"-Prinzipien architektonisch in die Sprachassistenten zu integrieren.<sup>4</sup> Die Forschung konzentriert sich dabei auf folgende Kerntechnologien:

1. **Edge Computing und On-Device Processing:** Anstatt rohe, unverschlüsselte Audiosignale zur Verarbeitung an zentrale Cloud-Server zu senden, werden die akustische Signalverarbeitung, die Wakeword-Erkennung und zunehmend auch das Natural Language Processing lokal auf dem Mikrochip des Smart Speakers (am "Edge") durchgeführt.<sup>53</sup> Nur noch die anonymisierte, textliche Interpretation des Befehls verlässt das Gerät, was das Risiko der massenhaften Speicherung von biometrischen Voiceprints auf Servern drastisch reduziert.
2. **Federated Learning (Föderiertes Lernen):** Um die Spracherkennungsmodelle kontinuierlich zu verbessern, ohne zentrale Datenhorte aufzubauen, wird Federated Learning eingesetzt. Hierbei trainiert das KI-Modell dezentral auf den Endgeräten der Nutzer mit den lokalen Sprachdaten. Nur die abstrakten mathematischen Verbesserungen des Modells (die Gewichtungen oder "Weights") werden an den zentralen Server gesendet und aggregiert. Die eigentlichen Rohaudiodaten verbleiben dauerhaft auf dem Gerät des Nutzers.<sup>96</sup>
3. **Akustische Anonymisierung und Filterung:** Eine der vielversprechendsten Entwicklungen ist die Integration von Hardware- oder Firmware-Modulen (wie das "VoiceSecure"-System), die tief in der Pipeline operieren, bevor das Audiosignal überhaupt das Betriebssystem erreicht. Solche Module nutzen Adversarial Machine Learning, um genau jene akustischen Merkmale (wie Jitter, Shimmer, Pitch-Varianz) aus dem Signal herauszufiltern oder zu maskieren, die Maschinen zur biometrischen Identifikation und zur Profilierung von Alter, Geschlecht oder Gesundheitszustand benötigen.<sup>3</sup> Dabei bleibt die linguistische Verständlichkeit für den Sprachassistenten intakt, aber der "metalinguistische" Fingerabdruck wird zerstört.

## Synthese und Ausblick

Die vorliegende Untersuchung illustriert in aller gebotenen wissenschaftlichen Detailtiefe, dass moderne Sprachassistenten und KI-gestützte Sprachschnittstellen weit mehr sind als bloße komfortable Werkzeuge zur Transkription von Audiobefehlen. Sie fungieren vielmehr als hochauflösende, omnipräsente und multimodale Sensoren, die ein beispielloses Maß an intimmem, medizinischem und psychologischem Wissen über ihre Nutzer generieren. Aus den rein akustischen und paralinguistischen Charakteristika der Stimme – der Sprechgeschwindigkeit, der mikrozyklischen Tonhöhenvarianz, den Resonanzfrequenzen sowie Parametern wie Jitter und Shimmer – lassen sich fundamentale Demografien, komplexe Persönlichkeitsstrukturen gemäß den Big Five und der augenblickliche emotionale Status mit beeindruckender Präzision extrahieren. Noch weitreichender und gesellschaftlich brisanter ist die ausgereifte Identifikation vokaler Biomarker. Hierbei erlauben kleinste, dem menschlichen

Ohr verborgene Veränderungen in der Pausensetzung, der lexikalischen Diversität, der syntaktischen Komplexität oder der Stimmelmelodie den Algorithmen, schwere neurodegenerative Erkrankungen wie Morbus Parkinson oder Alzheimer sowie tiefgreifende psychiatrische Leiden wie Major Depression mit Genauigkeiten von teilweise weit über 90 % zu detektieren – und dies oftmals Jahre, bevor die Symptome klinisch evident werden. Parallel zur Inferenz am Menschen selbst erschließt die Akustische Szenenklassifikation (ASC) aus dem vermeintlich wertlosen Hintergrundrauschen den sozioökonomischen Status und den Lebensstandard des Haushalts. Durch die algorithmische Analyse von Verkehrslärmexpositionen, dem Nachhall der Umgebung und dem spezifischen, markentypischen Surren von Premium-Haushaltsgeräten konstruiert die KI ein genaues Bild der Wohnverhältnisse. Flankiert von der kontinuierlichen Auswertung der Metadaten – dem exakten Zeitpunkt, der Frequenz, der Dauer und der Triangulierung der geografischen Orte der Interaktion – entsteht letztendlich ein unbestechliches, datengetriebenes Abbild des zirkadianen Rhythmus, des Schlafverhaltens, der physischen Mobilität und der sozialen Isolation des Nutzers.

Diese beispiellose technologische Kapazität offenbart ein fundamentales informationelles Asymmetrie-Problem unserer Zeit: Nutzer interagieren mit den Geräten, indem sie bewusst triviale Befehle eingeben ("Stelle einen Timer", "Wie wird das Wetter?"). Das KI-System hingegen extrahiert völlig unbewusst und automatisch hochsensible tiefenpsychologische, medizinische und sozioökonomische Profile aus dem Kommunikationsakt an sich. Während diese Inferenz-Fähigkeiten enorme, unbestreitbare Potenziale für die personalisierte Präzisionsmedizin, die telemedizinische Früherkennung von Krankheiten und das automatisierte Smart-Home-Management der Zukunft bieten, bergen sie gleichzeitig das dystopische Risiko einer ubiquitären biometrischen Überwachung. In der aktuellen Plattformökonomie werden diese Datenströme primär zur kommerziellen Ausbeutung durch tiefenintegrierte Werbenetzwerke genutzt, um Konsumenten in Momenten emotionaler oder psychischer Vulnerabilität mit algorithmisch perfektionierten Kaufreizen zu manipulieren. Der fortlaufende technologische Fortschritt im Bereich der Large Audio Language Models (LALMs) und der multimodalen neuronalen Netze wird die Sensitivität und die Prädiktionskraft dieser Systeme in den kommenden Jahren nochmals exponentiell steigern. Die zentrale technologische, juristische und gesellschaftliche Herausforderung der Zukunft wird daher zwingend darin bestehen, die architektonische Transparenz dieser Black-Box-Systeme regulatorisch zu erzwingen. Es bedarf der flächendeckenden Implementierung robuster, kryptografischer "Privacy-by-Design"-Mechanismen – wie lokales Edge-Processing und proaktive akustische Anonymisierung –, die garantieren, dass die Stimme des Menschen zwar von der Maschine zur Wunscherfüllung gehört und verstanden wird, sein tiefstes biologisches und psychologisches Wesen jedoch nicht ohne expliziten, informierten Konsens quantifiziert, pathologisiert und monetarisiert werden kann.

## Referenzen

1. Echoes of Privacy: Uncovering the Profiling Practices of Voice Assistants, Zugriff am April 3, 2026, <https://petsymposium.org/popets/2025/popets-2025-0050.pdf>

2. Guidelines 02/2021 on virtual voice assistants Version 2.0 - European Data Protection Board, Zugriff am April 3, 2026, [https://www.edpb.europa.eu/system/files/2021-07/edpb\\_guidelines\\_202102\\_on\\_va\\_v2.0\\_adopted\\_en.pdf](https://www.edpb.europa.eu/system/files/2021-07/edpb_guidelines_202102_on_va_v2.0_adopted_en.pdf)
3. New AI tool fights back against speech eavesdropping - EurekAlert!, Zugriff am April 3, 2026, <https://www.eurekalert.org/news-releases/1103494>
4. Ethical and privacy risks of AI voice agents - Aircall, Zugriff am April 3, 2026, <https://aircall.io/blog/support/ai-voice-agent-privacy/>
5. AI-determined similarity increases likability and trustworthiness of human voices - PMC, Zugriff am April 3, 2026, <https://pubmed.ncbi.nlm.nih.gov/articles/PMC11882051/>
6. Privacy in Speech Technology - arXiv, Zugriff am April 3, 2026, <https://arxiv.org/html/2305.05227v3>
7. Overview of Modern Technologies for Acquiring and Analysing Acoustic Information Based on AI and IoT - MDPI, Zugriff am April 3, 2026, <https://www.mdpi.com/2076-3417/15/12/6690>
8. 17 User Profiling Based on Nonlinguistic Audio Data - DSpace@MIT, Zugriff am April 3, 2026, <https://dspace.mit.edu/bitstream/handle/1721.1/138373/User%20Profiling%20Based%20on%20Nonlinguistic%20Audio%20Data.pdf?sequence=1&isAllowed=y>
9. Listening to the Mind: Integrating Vocal Biomarkers into Digital Health - PMC, Zugriff am April 3, 2026, <https://pubmed.ncbi.nlm.nih.gov/articles/PMC12293195/>
10. Vocal Biomarkers for Mental Health: Diagnosing Mental Disorders with a Short Voice Recording - American Psychiatric Association, Zugriff am April 3, 2026, <https://www.psychiatry.org/news-room/apa-blogs/vocal-biomarkers-for-mental-health>
11. Associations of smartphone usage patterns with sleep and mental health symptoms in a clinical cohort receiving virtual behavioral medicine care: a retrospective study - PubMed, Zugriff am April 3, 2026, <https://pubmed.ncbi.nlm.nih.gov/37485313/>
12. (PDF) Speech-based personality prediction using deep learning with acoustic and linguistic embeddings - ResearchGate, Zugriff am April 3, 2026, [https://www.researchgate.net/publication/386383693\\_Speech-based\\_personality\\_prediction\\_using\\_deep\\_learning\\_with\\_acoustic\\_and\\_linguistic\\_embeddings](https://www.researchgate.net/publication/386383693_Speech-based_personality_prediction_using_deep_learning_with_acoustic_and_linguistic_embeddings)
13. Optimizing Acoustic Feature Selection for Estimating Speaker Traits: A Novel Threshold- Based Approach - IJETA, Zugriff am April 3, 2026, <https://www.ijeta.org/download/file/fid/117987>
14. Audio Signal Processing for AI Voice Assistants: Key Methods - Patsnap Eureka, Zugriff am April 3, 2026, <https://eureka.patsnap.com/article/audio-signal-processing-for-ai-voice-assistants-key-methods>
15. Signal Processing in AI Voice Agents: 2025 - AI Sales Training Platform | Roleplay.video, Zugriff am April 3, 2026, <https://roleplay-video.webflow.io/blog/best-signal-processing-for-voice-ai-2025>
16. 2 Minuten Wissen: Spracherkennung - KI und medizinischer Fortschritt [FAU]

- Science], Zugriff am April 3, 2026,  
<https://www.youtube.com/watch?v=blrZPabR2hU>
17. Feasibility of Big Data Analytics to Assess Personality Based on Voice Analysis - PMC, Zugriff am April 3, 2026, <https://pmc.ncbi.nlm.nih.gov/articles/PMC11598682/>
  18. Artificial intelligence (AI)-driven technologies for managing pediatric speech and language therapy: A scoping review - PMC, Zugriff am April 3, 2026,  
<https://pmc.ncbi.nlm.nih.gov/articles/PMC12589799/>
  19. Paralinguistics in Communication: Speak with Voice and Expression - PlanetSpark, Zugriff am April 3, 2026,  
<https://www.planetspark.in/communication-skills/paralinguistics-in-communication>
  20. Uncovering interactive effects of affective voice tone and personality diversity on dyadic creativity - Frontiers, Zugriff am April 3, 2026,  
<https://www.frontiersin.org/journals/psychology/articles/10.3389/fpsyg.2025.1668759/full>
  21. Linguistic markers in Parkinson's disease (Smith et al., 2018) - ASHA journals - Figshare, Zugriff am April 3, 2026,  
[https://asha.figshare.com/articles/dataset/Linguistic\\_markers\\_in\\_Parkinson\\_s\\_disease\\_Smith\\_et\\_al\\_2018\\_/6615401](https://asha.figshare.com/articles/dataset/Linguistic_markers_in_Parkinson_s_disease_Smith_et_al_2018_/6615401)
  22. Feasibility of Big Data Analytics to Assess Personality Based on Voice Analysis - MDPI, Zugriff am April 3, 2026, <https://www.mdpi.com/1424-8220/24/22/7151>
  23. Large Language Models Can Infer Personality from Free-Form User ..., Zugriff am April 3, 2026, <https://arxiv.org/abs/2405.13052>
  24. Automatically Assessing Personality from Speech - Microsoft Research, Zugriff am April 3, 2026,  
<https://www.microsoft.com/en-us/research/video/automatically-assessing-personality-from-speech/>
  25. Speech-based personality prediction using deep learning with acoustic and linguistic embeddings - PMC, Zugriff am April 3, 2026,  
<https://pmc.ncbi.nlm.nih.gov/articles/PMC11615297/>
  26. Personality through Paralanguage: An Exploration Using Voice Analysis - Society, Zugriff am April 3, 2026, <https://society.org/articles/activity/10.31234/osf.io/dqfmx>
  27. Language-Based AI Modeling of Personality Traits and Pathology from Life Narrative Interviews - PMC, Zugriff am April 3, 2026,  
<https://pmc.ncbi.nlm.nih.gov/articles/PMC12616480/>
  28. User Profiling based on Nonlinguistic Audio Data - Welcom to COMP's Personal Webpage Server, Zugriff am April 3, 2026,  
<https://www4.comp.polyu.edu.hk/~labimcl/paper/Shen-2021-TOIS-UPNAD.pdf>
  29. Detecting paralinguistic events in audio stream using context in features and probabilistic decisions - PMC, Zugriff am April 3, 2026,  
<https://pmc.ncbi.nlm.nih.gov/articles/PMC5507373/>
  30. Clinical Decision Support Using Speech Signal Analysis: Systematic Scoping Review of Neurological Disorders - Journal of Medical Internet Research, Zugriff am April 3, 2026, <https://www.jmir.org/2025/1/e63004/PDF>
  31. AI-based Speech Error Detection to Differentiate Primary Progressive Aphasia

- Variants - medRxiv, Zugriff am April 3, 2026,  
<https://www.medrxiv.org/content/10.64898/2026.02.23.26346899v1.full.pdf>
32. Speaking in Alzheimer's Disease, is That an Early Sign? Importance of Changes in Language Abilities in Alzheimer's Disease - Frontiers, Zugriff am April 3, 2026,  
<https://www.frontiersin.org/journals/aging-neuroscience/articles/10.3389/fnagi.2015.00195/full>
  33. Speech and language biomarkers for Parkinson's disease prediction, early diagnosis and progression - PMC, Zugriff am April 3, 2026,  
<https://pmc.ncbi.nlm.nih.gov/articles/PMC11933288/>
  34. AI-driven precision diagnosis and treatment in Parkinson's disease: a comprehensive review and experimental analysis - Frontiers, Zugriff am April 3, 2026,  
<https://www.frontiersin.org/journals/aging-neuroscience/articles/10.3389/fnagi.2025.1638340/full>
  35. Advancing Parkinson's Detection with Vocal Biomarkers and Speech Foundation Models, Zugriff am April 3, 2026,  
<https://canaryspeech.com/blog/advancing-parkinsons-detection/>
  36. Voice Analysis for Neurological Disorder Recognition—A Systematic Review and Perspective on Emerging Trends - PMC, Zugriff am April 3, 2026,  
<https://pmc.ncbi.nlm.nih.gov/articles/PMC9309252/>
  37. Linguistic Indicators of Early Cognitive Decline in the DementiaBank Pitt Corpus: A Statistical and Machine Learning Study - arXiv, Zugriff am April 3, 2026,  
<https://arxiv.org/html/2602.11028v1>
  38. The Deterioration of Semantic Networks in Alzheimer's Disease - NCBI - NIH, Zugriff am April 3, 2026, <https://www.ncbi.nlm.nih.gov/books/NBK552151/>
  39. Syntactic Complexity as a Linguistic Marker to Differentiate Mild Cognitive Impairment From Normal Aging - ASHA Journals, Zugriff am April 3, 2026,  
[https://pubs.asha.org/doi/10.1044/2020\\_JSLHR-19-00335](https://pubs.asha.org/doi/10.1044/2020_JSLHR-19-00335)
  40. LSU Research Bites: AI Analysis of Everyday Speech Helps Detect Dementia Earlier, Zugriff am April 3, 2026,  
<https://www.lsu.edu/blog/2025/11/rb-dementia-ai.php>
  41. AI speech analysis predicted progression of cognitive impairment to Alzheimer's with over 78% accuracy, Zugriff am April 3, 2026,  
<https://www.alzheimers.gov/news/ai-speech-analysis-predicted-progression-cognitive-impairment-alzheimers-over-78-accuracy>
  42. Longitudinal Speech Biomarkers for Automated Alzheimer's Detection - Frontiers, Zugriff am April 3, 2026,  
<https://www.frontiersin.org/journals/computer-science/articles/10.3389/fcomp.2021.624694/full>
  43. Depression Screening from Voice Samples of Patients Affected by Parkinson's Disease | Digital Biomarkers | Karger Publishers, Zugriff am April 3, 2026,  
<https://karger.com/dib/article/3/2/72/99757/Depression-Screening-from-Voice-Samples-of>
  44. Deconstructing demographic bias in speech-based machine learning models for digital health - Frontiers, Zugriff am April 3, 2026,

- <https://www.frontiersin.org/journals/digital-health/articles/10.3389/fdgth.2024.1351637/full>
45. The Human Voice as a Digital Health Solution Leveraging Artificial Intelligence - MDPI, Zugriff am April 3, 2026, <https://www.mdpi.com/1424-8220/25/11/3424>
  46. A Scoping Review on the Use of Voice Biomarkers for Emotional Assessment, Zugriff am April 3, 2026, <https://econtent.hogrefe.com/doi/10.1027/1015-5759/a000915>
  47. Acoustic Scene Classification: Classifying environments from the sounds they produce, Zugriff am April 3, 2026, [https://openresearch.surrey.ac.uk/view/pdfCoverPage?instCode=44SUR\\_INST&filePid=13140269960002346&download=true](https://openresearch.surrey.ac.uk/view/pdfCoverPage?instCode=44SUR_INST&filePid=13140269960002346&download=true)
  48. How to do acoustic scene classification with stm32cube.AI - STMicroelectronics, Zugriff am April 3, 2026, [https://www.st.com/content/st\\_com/en/st-edge-ai-suite/case-studies/acoustic-scene-classification.html](https://www.st.com/content/st_com/en/st-edge-ai-suite/case-studies/acoustic-scene-classification.html)
  49. Research on Acoustic Scene Classification Based on Time–Frequency–Wavelet Fusion Network - MDPI, Zugriff am April 3, 2026, <https://www.mdpi.com/1424-8220/25/13/3930>
  50. Robust Detection of Background Acoustic Scene in the Presence of Foreground Speech - imec Publications, Zugriff am April 3, 2026, <https://imec-publications.be/bitstream/handle/20.500.12860/43487.3/DS713.pdf?sequence=1&isAllowed=y>
  51. Feature Extraction of Binaural Recordings for Acoustic Scene Classification - University of Huddersfield Research Portal, Zugriff am April 3, 2026, [https://pure.hud.ac.uk/ws/portalfiles/portal/14022736/IEEE\\_2018\\_SpatialSceneClassification\\_Zilensky\\_and\\_Lee.pdf](https://pure.hud.ac.uk/ws/portalfiles/portal/14022736/IEEE_2018_SpatialSceneClassification_Zilensky_and_Lee.pdf)
  52. Low-Complexity Acoustic Scene Classification with Device Information in the DCASE 2025 Challenge - arXiv, Zugriff am April 3, 2026, <https://arxiv.org/html/2505.01747v1>
  53. Real-Time Acoustic Scene Recognition for Elderly Daily Routines Using Edge-Based Deep Learning - MDPI, Zugriff am April 3, 2026, <https://www.mdpi.com/1424-8220/25/6/1746>
  54. EMBEDDED ACOUSTIC SCENE CLASSIFICATION FOR LOW POWER MICROCONTROLLER DEVICES Filippo Naccari Ivana Guarneri STMicroelectronics S - DCASE, Zugriff am April 3, 2026, [https://dcase.community/documents/workshop2020/proceedings/DCASE2020Workshop\\_Naccari\\_44.pdf](https://dcase.community/documents/workshop2020/proceedings/DCASE2020Workshop_Naccari_44.pdf)
  55. TUT Database for Acoustic Scene Classification and Sound Event Detection - EURASIP, Zugriff am April 3, 2026, <https://new.eurasip.org/Proceedings/Eusipco/Eusipco2016/papers/1570251932.pdf>
  56. Data-Efficient Low-Complexity Acoustic Scene Classification in the DCASE 2024 Challenge, Zugriff am April 3, 2026, <https://arxiv.org/html/2405.10018v2>
  57. [1904.05204] Acoustic Scene Classification by Implicitly Identifying Distinct Sound Events, Zugriff am April 3, 2026, <https://arxiv.org/abs/1904.05204>

58. Relation between Observed and Perceived Traffic Noise and Socio-Economic Status in Urban Blocks of Different Characteristics - MDPI, Zugriff am April 3, 2026, <https://www.mdpi.com/2413-8851/2/1/20>
59. Socioeconomic status and environmental noise exposure in Montreal, Canada - PMC - NIH, Zugriff am April 3, 2026, <https://pmc.ncbi.nlm.nih.gov/articles/PMC4358710/>
60. Full article: Exploring multi-pollution variability in the urban environment: geospatial AI-driven modeling of air and noise - Taylor & Francis, Zugriff am April 3, 2026, <https://www.tandfonline.com/doi/full/10.1080/17538947.2024.2378819>
61. Race/Ethnicity, Socioeconomic Status, Residential Segregation, and Spatial Variation in Noise Exposure in the Contiguous United States - PMC, Zugriff am April 3, 2026, <https://pmc.ncbi.nlm.nih.gov/articles/PMC5744659/>
62. (PDF) An automated framework for traffic noise level analysis using explainable artificial intelligence techniques - ResearchGate, Zugriff am April 3, 2026, [https://www.researchgate.net/publication/396737266\\_An\\_automated\\_framework\\_for\\_traffic\\_noise\\_level\\_analysis\\_using\\_explainable\\_artificial\\_intelligence\\_techniques](https://www.researchgate.net/publication/396737266_An_automated_framework_for_traffic_noise_level_analysis_using_explainable_artificial_intelligence_techniques)
63. The Effect of Background Noise, Bilingualism, Socioeconomic Status, and Cognitive Functioning on Primary School Children's Narrative Listening Comprehension - PubMed, Zugriff am April 3, 2026, <https://pubmed.ncbi.nlm.nih.gov/38363725/>
64. Assessing socioeconomic bias in machine learning algorithms in health care: a case study of the HOUSES index - PMC, Zugriff am April 3, 2026, <https://pmc.ncbi.nlm.nih.gov/articles/PMC9196683/>
65. Socioeconomic status measure helps researchers develop artificial intelligence models, improving equity in health care, Zugriff am April 3, 2026, <https://newsnetwork.mayoclinic.org/discussion/socioeconomic-status-measure-helps-researchers-develop-artificial-intelligence-models-improving-equity-in-health-care/>
66. (PDF) Efficient Energy Consumption: Leveraging AI Models for Appliance Detection - ResearchGate, Zugriff am April 3, 2026, [https://www.researchgate.net/publication/401271752\\_Efficient\\_Energy\\_Consumption\\_on\\_Leveraging\\_AI\\_Models\\_for\\_Appliance\\_Detection](https://www.researchgate.net/publication/401271752_Efficient_Energy_Consumption_on_Leveraging_AI_Models_for_Appliance_Detection)
67. AI-Enhanced System to Monitor Real-Time Energy and to Identify Home Appliances, Zugriff am April 3, 2026, [https://scielo.org.za/scielo.php?script=sci\\_arttext&pid=S2224-78902025000400024](https://scielo.org.za/scielo.php?script=sci_arttext&pid=S2224-78902025000400024)
68. KI-gestützte Klangprüfung: Akustische Intelligenz sorgt für 25 % weniger Fehler und 75 % Automatisierung - AROBS Transilvania Software development, Zugriff am April 3, 2026, <https://arobs.com/blog/ki-gestuetzte-klangpruefung-akustische-intelligenz-sorgt-fur-25-weniger-fehler-und-75-automatisierung/>
69. Efficient Energy Consumption: Leveraging AI Models for Appliance Detection - MDPI, Zugriff am April 3, 2026, <https://www.mdpi.com/2079-9268/16/1/9>
70. GE Appliances Audio Branding | Audio UX®, Zugriff am April 3, 2026,

- <https://www.auxnyc.com/case-studies/ge-profile-audio-branding>
71. Cooking with Intelligence: CUCKOO Debuts AI-Powered Induction Range | Renesas, Zugriff am April 3, 2026, <https://www.renesas.com/en/about/customer-success-stories/cuckoo>
  72. Custom AI Test Equipment Engineering: Acoustic Quality Assurance for Smart Appliances, Zugriff am April 3, 2026, <https://www.cardinalpeak.com/product-development-case-studies/custom-ai-test-equipment-engineering-appliance>
  73. Socioeconomic Status Influences the Usage of Language Technologies - arXiv, Zugriff am April 3, 2026, <https://arxiv.org/html/2505.12158v2>
  74. AI in Your Home Appliances: The Good, The Bad, & The Future, Zugriff am April 3, 2026, <https://www.hallocks.com/blog/shore-ai-home-appliances>
  75. Evaluating the privacy properties of telephone metadata - PNAS, Zugriff am April 3, 2026, <https://www.pnas.org/doi/10.1073/pnas.1508081113>
  76. Voice Assistant Data, Behavioral Profiling, and Targeted Advertising in the 21st Century - eScholarship.org, Zugriff am April 3, 2026, <https://escholarship.org/content/qt0tv6r3ck/qt0tv6r3ck.pdf>
  77. Clustering Insomnia Patterns by Data From Wearable Devices: Algorithm Development and Validation Study - JMIR mHealth and uHealth, Zugriff am April 3, 2026, <https://mhealth.jmir.org/2019/12/e14473/>
  78. Development of a Voice-Activated Virtual Assistant to Improve Insomnia Among Young Adult Cancer Survivors: Mixed Methods Feasibility and Acceptability Study - PMC, Zugriff am April 3, 2026, <https://pmc.ncbi.nlm.nih.gov/articles/PMC11933750/>
  79. A Machine Learning Approach for Detecting Digital Behavioral Patterns of Depression Using Nonintrusive Smartphone Data (Complementary Path to Patient Health Questionnaire-9 Assessment): Prospective Observational Study - PMC, Zugriff am April 3, 2026, <https://pmc.ncbi.nlm.nih.gov/articles/PMC9152726/>
  80. Personalised modelling of routine variability and affective states - PMC, Zugriff am April 3, 2026, <https://pmc.ncbi.nlm.nih.gov/articles/PMC12501027/>
  81. Depression screening using mobile phone usage metadata: a machine learning approach, Zugriff am April 3, 2026, [https://www.researchgate.net/publication/338824245\\_Depression\\_screening\\_using\\_mobile\\_phone\\_usage\\_metadata\\_a\\_machine\\_learning\\_approach](https://www.researchgate.net/publication/338824245_Depression_screening_using_mobile_phone_usage_metadata_a_machine_learning_approach)
  82. Acoustic Scene Analysis - Parvati Jayakumar, Zugriff am April 3, 2026, <https://parvatijay2901.medium.com/acoustic-scene-analysis-18572645ffc5>
  83. The AI Gap: How Socioeconomic Status Affects Language Technology Interactions - ACL Anthology, Zugriff am April 3, 2026, <https://aclanthology.org/2025.acl-long.914.pdf>
  84. AI voice transcription: Implications for data protection | AEPD, Zugriff am April 3, 2026, <https://www.aepd.es/en/press-and-communication/blog/ai-voice-transcription>
  85. Data Exhaust in Voice Assistants: Analysis and Mitigation Approaches - DalSpace, Zugriff am April 3, 2026, <https://dalspace.library.dal.ca/items/86f01622-75de-4efc-bb99-7ef50a382707>

86. Using Smartphones to Collect Behavioral Data in Psychological Science: Opportunities, Practical Considerations, and Challenges - PMC, Zugriff am April 3, 2026, <https://pmc.ncbi.nlm.nih.gov/articles/PMC5572675/>
87. Amazon Patents New Alexa Feature that Knows When You're Ill and Offers you Medicine, Zugriff am April 3, 2026, <https://mediawell.ssrc.org/news-items/amazon-patents-new-alexa-feature-that-knows-when-youre-ill-and-offers-you-medicine-the-telegraph/>
88. Amazon patents illness-spotting speech analyser - Privacy International, Zugriff am April 3, 2026, <https://privacyinternational.org/examples/2514/amazon-patents-illness-spotting-speech-analyser>
89. Patents That Made Amazon's Alexa a Dominant Voice Assistant, Zugriff am April 3, 2026, <https://patentpc.com/blog/patents-that-made-amazons-alexa-a-dominant-voice-assistant>
90. Google and Amazon really DO want to spy on you: Patent reveals future versions of their voice assistants will record your conversations to sell you products - - Consumer Watchdog, Zugriff am April 3, 2026, <https://consumerwatchdog.org/energy/google-and-amazon-really-do-want-spy-you-patent-reveals-future-versions-their/>
91. Virtual voice assistants' challenge to comply with data protection legislation, Zugriff am April 3, 2026, <https://www.osborneclarke.com/insights/virtual-voice-assistants-challenge-comply-data-protection-legislation>
92. Newly Released Amazon Patent Shows Just How Much Creepier Alexa Can Get, Zugriff am April 3, 2026, <https://www.sciencealert.com/creepy-new-amazon-patent-would-mean-alexa-records-everything-you-say-from-now-on>
93. Your essential 2026 guide to voice ai compliance in today's digital landscape, Zugriff am April 3, 2026, <https://www.speechmatics.com/company/articles-and-news/your-essential-guide-to-voice-ai-compliance-in-todays-digital-landscape>
94. From Speech to Data: Unraveling Google's Use of Voice Data for User Profiling - arXiv, Zugriff am April 3, 2026, <https://arxiv.org/html/2403.05586v1>
95. Verbraucherschutz beim Einsatz von Künstlicher Intelligenz - BMJV, Zugriff am April 3, 2026, [https://www.bmjbv.de/DE/themen/verbraucherschutz/digitaler\\_verbraucherschutz/ki/ki.html](https://www.bmjbv.de/DE/themen/verbraucherschutz/digitaler_verbraucherschutz/ki/ki.html)
96. AI Assistant for Socioeconomic Empowerment Using Federated Learning - ACL Anthology, Zugriff am April 3, 2026, <https://aclanthology.org/2025.nlp4dh-1.42.pdf>
97. Voice as a Biomarker of Health - Building an ethically sourced, bioacoustic database to understand disease like never before - NIH RePORTER, Zugriff am April 3, 2026, <https://reporter.nih.gov/project-details/10473236>