

**White Paper**  
**KI-basierte Datenverarbeitung mit DocuWare**

# Inhalt

1. Im Business von künstlicher Intelligenz profitieren – aber mit Sicherheit	2
Künstliche Intelligenz in drei Minuten	2
2. Wie DocuWare künstliche Intelligenz zu Ihrem Vorteil einsetzt	4
Transparente und klar begrenzte Datennutzung	4
DocuWare Intelligent Document Processing (IDP)	4
DocuWare Intelligent Indexing	6
3. KI und Datenschutz bei DocuWare	7
Anonymisierung der Daten	7
4. Sicherheit bei Datenzentren und Kommunikation	9

# 1. Im Business von künstlicher Intelligenz profitieren - aber mit Sicherheit

Die künstliche Intelligenz hat blitzartig Einzug in die Welt der Unternehmens-Software gehalten. Firmen weltweit profitieren von beschleunigten Prozessen und nie gekannten Potenzialen an Kreativität und auch Einsparmöglichkeiten.

Gleichzeitig sind die Unternehmen herausgefordert bei der Entscheidung, ob sie eine Software mit KI-Einsatz nutzen möchten. Neben Fragen der Ethik gehören der Datenschutz und die Transparenz bei der Datenverarbeitung zu den wichtigsten Kriterien. Jedes Unternehmen sollte schließlich wissen, wie sein Datengold verarbeitet wird und welche No-Gos zu beachten sind.

Auch DocuWare bietet den Einsatz von künstlicher Intelligenz an. Und vertritt den Ansatz, dass nur maximale Transparenz bei der Datenverarbeitung zu größtmöglichem Vertrauen aufseiten der Kunden führt. In diesem Sinne bietet dieses Paper Ihnen neben einem Grundwissen zur künstlichen Intelligenz alle Informationen, um deren Einsatz in DocuWare einzuschätzen.

## Künstliche Intelligenz in drei Minuten

Als künstliche Intelligenz werden Systeme und maschinelle Verfahren bezeichnet, die die menschliche Intelligenz imitieren. Diese KI-Systeme bestehen aus mathematischen Bausteinen, den sogenannten KI-Modellen mit ihren Algorithmen.

Diese Modelle werden mit Daten „gefüttert“ oder, wie es heißt, „auf Daten trainiert“. Der Prozess ähnelt dem bei einem Kind, dem man Erfahrungen verschafft und das daraus lernt. Ein Modell lernt, Muster zu erkennen, Vorhersagen und Entscheidungen zu treffen oder auch völlig neuen Inhalt zu generieren.

Dabei gibt es Modelle, die mit besonders vielen Daten trainiert werden und ein sehr gutes Verständnis unserer Sprache besitzen. Sie werden Large Language Models (LLMs) genannt. Sie kommen bei KI mit allgemeinem Verwendungszweck (General Purpose AI, GPAI) zum Einsatz, also bei der unten erläuterten generativen KI wie ChatGPT.

Was ein KI-System leisten kann, hängt stets vom genutzten Modell ab und davon, welche Technologien darin angewendet werden. Dabei begegnet man grundsätzlich drei Begriffen: maschinelles Lernen (ML), prädiktive KI und generative KI.

### Maschinelles Lernen

Maschinelles Lernen ist grundsätzlich ein Bereich in der KI-Forschung, der sich mittlerweile aber zu deren größtem Gebiet entwickelt hat. In diesem werden statistische Algorithmen entwickelt, auch Lernalgorithmen genannt.

Zur Anwendung kommt maschinelles Lernen, wenn der Lösungsweg für ein sehr komplexes Problem mit Regeln nicht mehr gut beschrieben werden kann, man aber über große Datenmengen verfügt. Aus deren Analyse werden im Training eines Modells Regeln abgeleitet. Diese verallgemeinert man als Lösungsweg und speichert sie im Modell. Es findet also kein Programmieren im üblichen Sinne statt.

Maschinelles Lernen wird angewendet, um für neue Daten Prognosen oder Entscheidungen zu erstellen.

## Prädiktive KI

Prädiktive künstliche Intelligenz, auch prädiktive Analyse genannt, beruht auf maschinellem Lernen. Sie kann Muster erkennen und auf dieser Basis Schlussfolgerungen ziehen, also Prognosen stellen. Ein prädiktives KI-System arbeitet mit statistischen (Teil-)Modellen und Algorithmen für maschinelles Lernen.

Ein Modell wird mit großen Datenmengen zunächst vortrainiert, um bestimmte Muster zu erkennen. Auf dieser Basis werden dann Prognosen ausgegeben.

Zusätzlich werden historische und aktuelle Daten in das KI-System eingegeben, das diese in seinem Modell verarbeitet und die Prognosen ausgibt. Ein Versandhandel etwa kann so anhand von Daten wie vergangenen Einkäufen den zukünftigen Bedarf von Kunden prognostizieren.

## Generative KI

Im Unterschied zur prädiktiven wird die generative künstliche Intelligenz meist dazu eingesetzt, neue Inhalte zu generieren. Ein Beispiel dafür ist das Modell ChatGPT.

Hier kommen hochentwickelte Algorithmen zum Einsatz, die Texte, Bilder, Videos und Audioaufnahmen erstellen können. Dabei werden natürliche Sprachinhalte genutzt und per maschinellem Lernen synthetisiert.

## Was prädiktive und generative KI gemeinsam haben

Sowohl die prädiktive als auch die generative KI beruhen auf maschinellem Lernen. Allerdings unterscheiden sich die in den beiden Formen verwendeten Algorithmen erheblich.

Beide KI-Arten nutzen künstliche neuronale Netze. Das sind Rechenknoten-Verbünde, die ähnlich den Neuronen im menschlichen Gehirn funktionieren und in mehreren Schichten angeordnet sind. Sie erlauben die gleichzeitige Verarbeitung sehr großer Datenmengen.

Bei neuronalen Netzen mit besonders vielen Schichten, sogenannten Hidden Layers, spricht man auch von Deep Learning.

## 2. Wie DocuWare künstliche Intelligenz zu Ihrem Vorteil einsetzt

DocuWare arbeitet mit prädiktiver künstlicher Intelligenz und maschinellem Lernen. Generative Modelle wie ChatGPT kommen nicht zum Einsatz. Das „Fantasieren“, also das mögliche Erfinden unlogischer oder unwahrer Inhalte sowie generell das Entwickeln neuer Ideen aus der KI selbst heraus ist demnach ausgeschlossen. Denn dies ist nur mit generativer KI möglich.

### Transparente und klar begrenzte Datennutzung

Bei seiner Anwendung prädiktiver KI nutzt DocuWare Kundendaten stets nur in dem Umfang, der vertraglich vereinbart ist – nämlich allein zur Erbringung der gebuchten Services.

Andersherum gesagt: DocuWare erbringt eine Dienstleistung unter Zuhilfenahme von KI-Technologien. Diese Dienstleistung wäre ohne eine Nutzung von Kundendaten technologisch nicht möglich. DocuWare verpflichtet sich dazu, die Kundendaten darüber hinaus in keiner Weise zu verwenden.

Gleichzeitig garantiert DocuWare durch die unten erläuterten Sicherheitsmaßnahmen, dass die für das Training von Modellen verwendeten Kundendaten gemäß dem Datenschutz nach DSGVO und entsprechend allen Vertraulichkeitsvereinbarungen (NDAs) behandelt werden.

Der Kunde kauft die Services ein und begrenzt damit die Datennutzung im hier beschriebenen Rahmen.

### DocuWare Intelligent Document Processing (IDP)

DocuWare Intelligent Document Processing (IDP) nutzt prädiktive KI und arbeitet mit verschiedenen KI-Techniken, die selbst auch als Modelle bezeichnet werden können.

Wenn Daten für Modelle kundenübergreifend eingesetzt werden, kommen nur vollständig anonymisierte Trainingsdaten zum Einsatz. „Kundenübergreifend“ bedeutet, dass mit den Daten ein allgemeines Basismodell trainiert wird, auf dem die weiter unten beschriebenen Pre-Built- und Fine-Tune-Modelle für IDP-Kunden aufbauen. Dabei werden keinerlei Daten von einem Kunden an andere Kunden weitergegeben. Überhaupt werden die Daten nicht in ihrer Reinform genutzt, sondern wie unten beschrieben stets zunächst durch De-Identifikation anonymisiert.

Die folgenden IDP-Techniken bietet DocuWare an:

**Splitting:** Das Splitting ist ein KI-basierter Prozess zur Aufteilung eines Stapel-Scans in einzelne Dokumente (PDFs), wofür visuelle und textliche Informationen verwendet werden. Beim Splitting gibt es ein vortrainiertes Modell, das Kunden mit ihren eigenen Daten für eine

verbesserte Automatisierung optimieren können. Dazu kommt eine Validierungsschnittstelle, die Ihnen in kürzester Zeit perfekt geteilte Dokumente liefert.

**Klassifizierung:** Die Klassifizierung dient der Zuordnung von Dokumenten zu verschiedenen Kategorien oder Gruppen auf der Grundlage ihres Inhalts, ihrer Eigenschaften oder Merkmale. Art und Anzahl der Klassen können kundenindividuell angepasst und trainiert werden.

**Extraktion:** Die Datenextraktion dient zur automatisierten Verarbeitung nahezu aller Dokumente, wie Rechnungen, Verträge, Personaldokumente oder Briefe. Bei Rechnungen etwa werden selbst Daten aus komplexen Einzelposten zuverlässig erkannt und extrahiert.

**Extraktion mit Handschriftenerkennung:** Die Datenextraktion umfasst auch die Erkennung von handschriftlichem Text in Bildern oder Dokumenten sowie dessen Umwandlung in bearbeitbaren Text.

## Die angebotenen Modelle

Für diese dokumentbezogenen Aufgaben bietet DocuWare eine IDP-Plattform, die es Kunden ermöglicht, innovative Modelle für maschinelles Lernen zu nutzen und auch selbst zu erstellen.

Es werden drei Anpassungsstufen von Modellen angeboten:

### **Modell Stufe 1: Pre-built**

Das Modell wurde von DocuWare konfiguriert und trainiert.

- *Funktionen:* Sofort einsatzbereite Plug-and-Play-Lösung.
- *Datenschutz:* Das Modell wird von DocuWare vollständig mit anonymisierten Daten trainiert. Es werden keine Daten des spezifischen Kunden für das Training genutzt.

### **Modell Stufe 2: Fine-tune**

Das Modell wurde zunächst von DocuWare konfiguriert und vortrainiert und dann an Ihre spezifischen Daten angepasst.

- *Funktionen:* Nutzt das grundlegende Modell und verfeinert es gleichzeitig mit Ihren Daten, um die Leistung zu verbessern.
- *Datenschutz:* Das Modell ist bereits von DocuWare mit anonymisierten Daten vortrainiert. Für das Finetuning Ihres individuellen Modells werden nur Ihre Daten verwendet. Diese werden nicht weitergegeben und nicht zum Training anderer Modelle verwendet.

### **Modell Stufe 3: Train-Your-Own (TYO)**

Das Modell wird vollständig für Ihr Unternehmen konfiguriert und ausschließlich auf Basis Ihrer Daten trainiert.

- *Funktionen:* Vollständig anpassbar an Ihre individuellen Anforderungen. Sie verwalten, kennzeichnen und kontrollieren Ihre Daten.
- *Datenschutz:* Für Ihr vollständig individuelles Modell werden nur Ihre Daten verwendet. Diese werden nicht weitergegeben und nicht zum Training anderer Modelle verwendet.

## **DocuWare Intelligent Indexing**

Intelligent Indexing arbeitet mit Algorithmen des maschinellen Lernens.

DocuWare klassifiziert Dokumente in verschiedene Typen. Dabei werden automatisch die relevanten Indexbegriffe in beziehungsweise zu den Dokumenten gesucht und dem Benutzer vorgeschlagen. Dieser bestätigt oder verbessert die Vorschläge. Anhand des Feedbacks lernt das System ständig hinzu.

Die Informationen, die für das Training dieser Algorithmen erforderlich sind, werden mit den Trainingsergebnissen in einem Modellraum gespeichert. Dabei handelt es sich um die Informationen, die von einem bereits gelernten Dokument für die Indexierung eines neu zu indexierenden Dokuments genutzt werden.

Erhält Intelligent Indexing ein neues Dokument, wird geprüft, ob im Modellraum ähnliche Dokumente vorliegen, deren Indexierungsmethoden auf das aktuelle Dokument übertragen werden können. Dabei kann bereits ein einziges sehr ähnliches Dokument ausreichend sein. Mit einer größeren Anzahl an Referenzdokumenten steigt jedoch die Wahrscheinlichkeit für eine erfolgreiche automatische Extraktion der Indexdaten.

Unabhängig vom Modellraum werden zusätzlich fest kodierte Regeln für die Indexierungsvorschläge verwendet. Diese Regeln beruhen auf typischen Dokumenten in vielen Sprachen und werden immer wieder angepasst.

Ein Modellraum ist immer organisationsspezifisch. Das heißt, die Volltextauszüge und Trainingsergebnisse werden pro Organisation, also Kunden-Firma, zusammengefasst. Sie sind damit strikt getrennt von den Daten anderer DocuWare Organisationen.

Die Daten eines Kunden, also Dokumente und Indexfeldwerte, werden allein in dessen Modellraum genutzt. Sie stehen weder anderen Kunden noch DocuWare zur Verfügung. Sie werden auch nicht zur Verbesserung von DocuWare Modellen oder zur Verbesserung von DocuWare oder Produkten bzw. Dienstleistungen von Drittanbietern verwendet.

Die vom Kunden optimierten Intelligent Indexing-Algorithmen werden nur in der Cloud des Kunden gespeichert und sind nicht für andere Kunden verfügbar. DocuWare bietet den Intelligent Indexing Service bereits seit vielen Jahren an.

## 3. KI und Datenschutz bei DocuWare

Beim Einsatz von künstlicher Intelligenz stehen für DocuWare der Schutz und die Sicherheit Ihrer Daten an erster Stelle. Das spiegeln auch unsere [Trusted AI-Prinzipien](#). Um beides zu gewährleisten, setzen wir auf proaktive Maßnahmen mit höchstem Standard.

Alle Daten werden anonymisiert, bevor sie für das Training des Basismodells verwendet werden. Genutzt wird im KI-System außerdem nie das ganze Dokument, sondern von jedem Dokument nur jeweils ein einziger, sehr kleiner Bestandteil.

Bevor Daten für das Training verwendet werden, wird zusätzlich der Kontext der einzelnen Textpassagen eines Dokuments entfernt. Dies geschieht, indem Teile der Dokumente vermischt miteinander zum Training verwendet werden. Darüber hinaus werden die Daten mit synthetischen Daten durchmischt, sodass nicht mehr festgestellt werden könnte, ob etwa ein Name der einer Person oder ein Fantasiename ist.

Die Dokumente selbst können genauso wenig wie personenbezogene Daten rekonstruiert werden.

### Anonymisierung der Daten

Bei der KI-basierten Datenverarbeitung mit DocuWare ist ausgeschlossen, dass Daten von einem Kunden an andere Kunden weitergegeben werden. Dies wird sichergestellt, indem das Basismodell, auf dem sowohl das Pre-Built- als auch das Fine-Tune-Modell aufbauen, nur mit anonymisierten kundenübergreifenden Daten trainiert wird.

Damit sorgt DocuWare vor Verwendung der Dokumente und Daten dafür, dass die personenbezogenen Daten geschützt sind. Noch wichtiger, es werden auch keine personenbezogenen Daten im Modell gespeichert.

Bei der Anonymisierung von personenbezogenen Daten ist es nicht möglich, Rückschlüsse auf Personen bzw. Betroffene zu ziehen. Angewendet wird hier die Anonymisierungstechnik der Generalisierung. Der konkrete Prozess wird allgemein als De-Identifizierung bezeichnet, wie [vom deutschen Digitalverband Bitkom beschrieben](#).

### Kundenübergreifende Modelle

Beim Training der kundenübergreifenden Klassifizierungs-, Aufteilungs- und Extraktionsmodelle, insbesondere des Pre-Built-Modells, werden verschiedene Arten von Daten in das neuronale Netz eingespeist. Diese Daten haben wie gesagt immer eine vorherige Anonymisierung durchlaufen. Es handelt sich um:

- visuelle Daten in Form von Pixeldaten aus nur einem Teil der Dokumentenseite mit niedriger Auflösung: Mit dieser Beschränkung wird sichergestellt, dass kein lesbarer Text vorhanden ist.



- Textdaten, die mit verschiedenen Filtern bearbeitet und ersetzt werden: So werden etwa alle Namen zufällig ersetzt und jegliche alphanumerischen Kennungen geändert. Auch hier werden jeweils nur Segmente einer Dokumentseite in das neuronale Netz gegeben.
- Positions- und Layoutdaten (wie „x1“, „y1“, „x2“, „y2“) zu einzelnen Wörtern und Texten: Diese werden zusammen mit den entsprechenden Wörtern und Texten konsistent in das neuronale Netz eingespeist.

## 4. Sicherheit bei Datenzentren und Kommunikation

DocuWare Intelligent Document Processing wird im deutschen Rechenzentrum von VSE NET gehostet. Es unterliegt den strengen Regeln der DSGVO.

Für DocuWare Intelligent Indexing werden folgende Datenzentren genutzt, die den jeweiligen regionalen Datenschutzregularien unterliegen:

- Amsterdam (Niederlande) für Kunden aus der Region EMEA
- Virginia (USA) für Kunden aus Nord- und Südamerika
- Tokio (Japan) für Kunden aus Japan
- New South Wales (Australien) für Kunden aus Australien und einigen weiteren asiatisch-pazifischen Ländern

Für Kunden von DocuWare Cloud ist das für Intelligent Indexing verwendete Datenzentrum immer in der gleichen Region wie das von DocuWare Cloud.

Wie jede andere Datenübertragung mit DocuWare geschieht auch die zum Zweck der KI-basierten Datenverarbeitung nur mit sicherer HTTPS-Verschlüsselung. Somit sind alle übertragenen Daten vor fremdem Zugriff geschützt.

[DocuWare](#) und [Microsoft Azure](#) sind hinsichtlich dieser und zahlreicher weiterer Sicherheitsmaßnahmen zertifiziert.

Copyright © 2024 DocuWare GmbH

Alle Rechte vorbehalten

Die Software enthält Proprietary-Information von DocuWare. Sie wird unter Lizenz bereitgestellt und ist darüber hinaus durch das Copyright geschützt. Im Lizenzvertrag sind Einschränkungen bezüglich der Nutzung und Offenlegung enthalten. Rekonstruktion der Software ist untersagt.

Da dieses Produkt laufend weiterentwickelt wird, können die hier enthaltenen Informationen ohne Vorankündigung geändert werden. Die hier enthaltenen Rechte am geistigen Eigentum und Informationen sind vertrauliche Informationen, die nur der DocuWare GmbH und dem Kunden zugänglich sind, und bleiben das ausschließliche Eigentum von DocuWare. Falls Sie in der Dokumentation auf Probleme stoßen, weisen Sie uns bitte in schriftlicher Form darauf hin. DocuWare übernimmt keine Garantie dafür, dass dieses Dokument frei von Fehlern ist.

Kein Teil dieser Veröffentlichung darf ohne die vorherige schriftliche Genehmigung von DocuWare in irgendeiner Form oder mithilfe welcher Verfahren auch immer (elektronisch, mechanisch, Fotokopie, Aufzeichnung oder auf andere Weise) vervielfältigt, in einem Retrievalsystem abgelegt oder übertragen werden.

#### *Disclaimer*

Dieses Dokument wurde mit größter Sorgfalt zusammengestellt und die Informationen darin sind Quellen entnommen, die als zuverlässig gelten. Dennoch kann keine Haftung übernommen werden für die Richtigkeit, Vollständigkeit und Aktualität der Informationen. Aus den in diesem Dokument aufgenommenen Informationen können keine Ansprüche hergeleitet werden. Die DocuWare GmbH behält sich das Recht vor, jegliche Informationen, die in diesem Dokument enthalten sind, ohne vorherige Ankündigung zu verändern.

DocuWare GmbH  
Planegger Straße 1  
82110 Germering

[www.docuware.com](http://www.docuware.com)