

From Pilot to Profit: Der Strategische Leitfaden zur Realisierung von KI-Wertschöpfung

Executive Summary: Das Ende der Experimentierphase

Das Jahr 2025 markiert einen fundamentalen Wendepunkt in der Unternehmens-IT. Während die Jahre 2023 und 2024 von einer Euphorie über die generativen Fähigkeiten von Künstlicher Intelligenz (KI) geprägt waren, stehen wir nun vor einer harten Realität, die wir bei Dreher Consulting als "Pilot Purgatory" (Pilot-Fegefeuer) bezeichnen. Die Datenlage ist eindeutig: Obwohl 88 % der Unternehmen angeben, KI in mindestens einer Geschäftsfunktion einzusetzen, verharren fast zwei Drittel in der Experimentier- oder Pilotphase.¹

Noch alarmierender ist die finanzielle Diskrepanz: Nur eine Minderheit der Unternehmen konnte bisher einen messbaren Einfluss auf das EBIT (Earnings Before Interest and Taxes) nachweisen.¹ Der Grund hierfür liegt nicht in der technologischen Unreife, sondern in einem fundamentalen Missverständnis der ökonomischen Wirkungsmechanismen. Die erste Welle der KI-Adoption konzentrierte sich auf "Co-Piloten" – Werkzeuge, die den Menschen assistieren. Dies führte zu individuellen Produktivitätsgewinnen, senkte jedoch selten die strukturellen Kosten, da der Mensch "in the loop" blieb und die gewonnene Zeit oft durch neue Aufgaben absorbiert wurde (Parkinsonsches Gesetz).

Dieser Report dient als strategische Handlungsanleitung ("Blueprint"), um diesen Stillstand zu durchbrechen. Wir analysieren den technologischen Paradigmenwechsel von **Generativer KI** (Content-Erstellung) hin zu **Agentic AI** (Autonome Ausführung). KI-Agenten, die in der Lage sind, komplexe Handlungsstränge selbstständig zu planen und auszuführen, bieten erstmals die Möglichkeit, die Grenzkosten von Transaktionen von der menschlichen Arbeitszeit zu entkoppeln.³

Basierend auf First-Principles-Denken und einer MECE-Analyse (Mutually Exclusive, Collectively Exhaustive) legen wir dar, wie Unternehmen den Sprung von isolierten Piloten zu profitablen, skalierbaren KI-Operationen schaffen.

1. Status Quo 2025: Die Anatomie des "Pilot Purgatory"

1.1 Die Diskrepanz zwischen Adoption und Wertschöpfung

Eine nüchterne Betrachtung der aktuellen Marktlandschaft offenbart eine signifikante Lücke zwischen Aktivität und Ergebnis. McKinsey-Daten zeigen, dass zwar die Nutzung von KI in den Unternehmen explodiert ist, die Skalierung jedoch stagniert. Nur etwa ein Drittel der Organisationen hat begonnen, KI-Lösungen unternehmensweit auszurollen.³ Die überwältigende Mehrheit der Projekte bleibt in Silos stecken – gefangen zwischen Proof-of-Concept (PoC) und Produktion.

Die Ursachen hierfür sind struktureller Natur:

1. **Fehlende Prozessintegration:** KI wird oft als technologisches "Add-on" betrachtet, anstatt Prozesse fundamental neu zu denken ("Re-Wiring").
2. **Daten-Fragilität:** Piloten funktionieren auf bereinigten Testdaten, scheitern aber an der Komplexität und Unsauberkeit realer Produktionsdaten.⁵
3. **Mangelnde Ambition:** Während "High Performer" KI nutzen, um neue Geschäftsmodelle zu entwickeln, beschränken sich viele Unternehmen auf inkrementelle Effizienzgewinne, die die hohen Implementierungskosten nicht rechtfertigen.¹

1.2 Der Aufstieg der Agentic AI

Um das Kostensenkungspotenzial zu realisieren, müssen wir die Technologie präzise definieren. Wir befinden uns im Übergang von passiven Modellen zu aktiven Agenten.

Tabelle 1: Technologische Evolution und ökonomische Implikation

Merkmal	Generative KI (GenAI)	Agentic AI (Autonome Agenten)
Kernfunktion	Erstellung & Zusammenfassung	Planung & Ausführung
Auslöser	Menschlicher Prompt ("Schreibe eine Mail")	Systemereignis oder Zielsetzung ("Bearbeite alle Rechnungen")
Kontext	Sitzungsbasiert (kurzfristig)	Persistenter Speicher & Zugriff auf Unternehmensdaten (RAG)
Interaktion	Chat-Interface	API-Calls & Systemintegrationen
Ökonomischer Hebel	Persönliche Produktivität (Soft ROI)	Arbeitssubstitution (Hard ROI)

Agentic AI unterscheidet sich durch drei Kernkompetenzen: **Wahrnehmung** (Perception) von Daten aus Live-Systemen, **Schlussfolgerung** (Reasoning) zur Zerlegung komplexer Ziele in Teilaufgaben, und **Handlung** (Action) durch direkte Systemeingriffe.¹ Erst durch diese Autonomie wird es möglich, den Menschen aus der kritischen Pfad der Transaktionsbearbeitung zu nehmen und so echte Kostensenkungen zu realisieren.

2. First Principles Economics: Die Ökonomie der Intelligenz

Um "Profit" zu generieren, müssen wir die Unit Economics (Stückkostenrechnung) der KI verstehen. Es handelt sich hierbei um eine Verschiebung von Fixkosten (Personal) zu variablen Kosten (Compute/Token), die jedoch bei Skalierung gegen Null tendieren.

2.1 Die Entkopplung von Arbeit und Volumen

In traditionellen Betriebsmodellen korreliert das Transaktionsvolumen linear mit den Arbeitskosten. Um doppelt so viele Kundenanfragen zu bearbeiten, benötigen Sie – ceteris paribus – doppelt so viel Personal. Agentic AI durchbricht diese Linearität. Nach der initialen Investition in Training und Integration skalieren die Kosten logarithmisch. Die Grenzkosten einer zusätzlichen Transaktion entsprechen lediglich den Inferenzkosten (Token) und den API-Gebühren.

Analysen zeigen, dass KI-Agenten die Kosten pro Transaktion im Vergleich zu Onshore-Arbeitskräften um 90-95 % und im Vergleich zu Offshore-Kräften um 50-70 % senken können.⁷

2.2 Die J-Kurve der Investition

Ein entscheidendes Konzept für die Kommunikation mit dem CFO ist die "J-Kurve". KI-Projekte liefern selten sofortigen ROI.

1. **Investitionsphase (Tal der Tränen):** Hohe Ausgaben für Datenbereinigung, Modelltraining und Infrastruktur. Der ROI ist negativ.
2. **Lernphase (Human-in-the-Loop):** Der Agent ist produktiv, benötigt aber hohe menschliche Überwachung. Die Effizienz ist gering, da Mitarbeiter sowohl arbeiten als auch die KI korrigieren müssen.
3. **Skalierungsphase:** Die Konfidenz des Agenten steigt, die Deflektionsrate (Autonomiegrad) wächst von 20 % auf 80 %. Hier kreuzt die Kurve die Nulllinie und generiert exponentiellen Profit.⁸

Unternehmen, die Projekte während der Lernphase abbrechen, weil "es schneller geht, es selbst zu machen", realisieren nie den Profit. Sie zahlen die Rüstkosten, ohne die Ernte einzufahren.

2.3 LCOAI: Eine neue Metrik für 2025

Wir führen bei Dreher Consulting die Metrik **LCOAI** (Levelized Cost of AI) ein, analog zu den Stromgestehungskosten in der Energiewirtschaft. Sie berechnet die Gesamtkosten pro nützlichem Output über den Lebenszyklus des Systems.¹⁰

$$\text{LCOAI} = \frac{\text{Entwicklungskosten} + \sum(\text{Inferenzkosten}) + \sum(\text{Wartungskosten})}{\text{Anzahl erfolgreich automatisierter Transaktionen}}$$

Diese Formel zwingt zur Ehrlichkeit: Ein Agent, der nur 500 Mal im Jahr genutzt wird, ist oft teurer als ein Mensch. Ein Agent, der 500.000 Transaktionen übernimmt, ist unschlagbar günstig. Volumen ist der Schlüssel zur Amortisation der Fixkosten.

3. Strategische Identifikation: Das MECE-Framework zur Potenzialanalyse

Um Potenziale "MECE" (überschneidungsfrei und vollständig) zu erfassen, kategorisieren wir Anwendungsfälle nach zwei Dimensionen: **Komplexität der Aufgabe** und **Volumen der Transaktionen**.

3.1 Quadrant A: Hohes Volumen, Niedrige Komplexität (Die "Automatisierungs-Zone")

Ziel: Direkte Kostensenkung (Labor Substitution)

Hier liegen die größten und am schnellsten realisierbaren Hebel ("Low Hanging Fruits"). Es handelt sich um regelbasierte, repetitive Aufgaben.

- **Customer Operations (Service & Support):**
 - *Anwendungsfall:* Vollständige autonome Bearbeitung von Tier-1-Anfragen (Rücksendungen, Statusabfragen, Adressänderungen).
 - *Evidenz:* Klarna ersetzte die Arbeitsleistung von 700 Vollzeitäquivalenten (FTEs) durch einen KI-Agenten, der 2,3 Millionen Konversationen führte und die Gewinnprognose um 40 Millionen Dollar verbesserte.¹¹
 - *Mechanismus:* Integration in CRM und ERP erlaubt dem Agenten nicht nur zu antworten, sondern die Transaktion (z.B. Rückerstattung) durchzuführen.
- **Finance & Accounting (F&A):**
 - *Anwendungsfall:* Rechnungsabgleich (Invoice Matching). Agenten vergleichen eingehende Rechnungen mit Bestellungen (PO) und Wareneingangsbelegen. Bei Diskrepanzen kontaktieren sie autonom den Lieferanten.
 - *Evidenz:* Ein globales Medienunternehmen konsolidierte Daten aus 80 Hauptbüchern und identifizierte Millionen an "Shadow IT"-Ausgaben durch KI-Analysen.¹²

3.2 Quadrant B: Hohes Volumen, Hohe Komplexität (Die "Augmentierungs-Zone")

Ziel: Produktivitätssteigerung & Durchsatz

Hier wird der Mensch nicht ersetzt, sondern massiv beschleunigt ("Super-Powering").

- **Software Engineering:**

- *Anwendungsfall:* Generierung von Boilerplate-Code, Unit-Tests und Dokumentation.
- *Evidenz:* Entwicklungszyklen können um 20–30 % verkürzt werden. Dies senkt nicht zwingend die Headcount-Kosten, erhöht aber den Output bei gleicher Kostenbasis (Vermeidung von Neueinstellungen).¹²
- *Risiko:* "Lazy Reviews" durch Entwickler können zu Qualitätsverlust führen, wenn der generierte Code nicht kritisch geprüft wird.¹¹

- **Healthcare Revenue Cycle Management (RCM):**

- *Anwendungsfall:* Bearbeitung von Leistungsablehnungen (Denials). Agenten analysieren den Ablehnungsgrund, korrigieren die Codierung und reichen den Antrag neu ein.
- *Evidenz:* Reduktion der Außenstandsdauer (A/R Days) um 35 Tage und Senkung der Ablehnungsquote um 7 %.¹³

3.3 Quadrant C: Niedriges Volumen, Hohe Komplexität (Die "Innovations-Zone")

Ziel: Strategischer Wettbewerbsvorteil

- **Supply Chain Management:**
 - *Anwendungsfall:* Prädiktive Risikoanalyse. Agenten überwachen Tausende externer Signale (Wetter, Streiks, Geopolitik) und simulieren Auswirkungen auf die Lieferkette in Echtzeit.
 - *Evidenz:* Ermöglichung proaktiver Routenänderungen vor Eintritt der Störung, was teure Sonderfahrten und Produktionsstopps verhindert.¹⁴

4. Das Operative Betriebsmodell: Architecture of Agency

Technologie allein generiert keinen Wert; sie benötigt ein operatives Einbettungssystem. Wir nennen dies die **Architecture of Agency**. Bain & Company betont zu Recht, dass die meisten Piloten nicht an der KI scheitern, sondern an der fehlenden Datenstrategie.⁵

4.1 Die Daten-Infrastruktur als Fundament

Agenten benötigen Zugriff auf "Ground Truth". Ein Agent, der auf veralteten oder widersprüchlichen Daten trainiert wird, halluziniert nicht zufällig – er halluziniert systemisch.

- **Data Products:** Behandeln Sie Datensätze (z.B. "Kundendaten", "Produktkatalog") als Produkte mit klaren SLAs (Service Level Agreements), Eigentümern und Qualitätsmetriken.¹⁵
- **Vectorization Pipeline:** Um unstrukturierte Daten (PDF-Handbücher, E-Mail-Archive) nutzbar zu machen, müssen diese in Vektordatenbanken (RAG - Retrieval Augmented Generation) überführt werden. Dies ist das "Langzeitgedächtnis" des Agenten.

4.2 Governance & Risiko-Management

Autonomie erfordert Kontrolle. Ein Agent, der autonom buchen oder kommunizieren darf, stellt ein operatives Risiko dar.

Das 3-Schichten-Sicherheitsmodell:

1. **Input Guardrails:** Filterung bösartiger Prompts ("Jailbreaking") und Sicherstellung des Datenschutzes (PII-Reduktion) bevor die Daten das Modell erreichen.
2. **Model Governance:** Auswahl des richtigen Modells für den richtigen Zweck. Nicht jede Aufgabe benötigt ein teures GPT-4; oft reichen kleinere, spezialisierte Modelle (SLMs), die schneller und günstiger sind (LCOAI-Optimierung).
3. **Output Guardrails:** Eine unabhängige KI-Instanz ("Critic Model") prüft den Output des Agenten auf Halluzinationen, Tonalität und Compliance, bevor die Aktion ausgeführt wird. "Darf der Agent diese Rückerstattung über 500€ wirklich freigeben?".¹⁶

4.3 Human-in-the-Loop (HITL) Design

Das Ziel ist nicht die 100%ige Automatisierung, sondern die *optimale* Automatisierung. Erfolgreiche Systeme leiten Transaktionen mit niedriger Konfidenz ("Low Confidence Score") automatisch an menschliche Experten weiter. Diese Korrekturen müssen zwingend in das System zurückfließen (Feedback Loop), um das Modell kontinuierlich zu verbessern.⁶

5. Handlungsanleitung: "From Pilot to Profit"

Im Folgenden präsentieren wir den konkreten Fahrplan für die Transformation. Dieser Leitfaden ist phasenbasiert und deckt die kritischen Meilensteine ab.

Phase 1: Die Diagnose & Auswahl (Wochen 1-4)

Ziel: Identifikation der "Golden Use Cases", wo Datenverfügbarkeit auf ökonomische Relevanz trifft.

1. **Kosten-Analyse (P&L-Scan):** Beginnen Sie nicht mit Ideen, sondern mit der Gewinn- und Verlustrechnung. Wo sind die größten Blöcke an SG&A-Kosten (Selling, General & Administrative Expenses)? Identifizieren Sie Prozesse mit hohem manuellem Aufwand.
2. **Atomare Prozess-Zerlegung:** Brechen Sie diese Prozesse in kleinste Arbeitsschritte auf. Wenden Sie den "Autonomie-Test" an:
 - Ist der Input digital?
 - Sind die Entscheidungsregeln explizit?
 - Ist das Ergebnis messbar?
3. **Daten-Audit:** Prüfen Sie die Verfügbarkeit der notwendigen Daten via API. Ohne API keine Agenten-Skalierung.

Output: Eine priorisierte Liste von 3 Anwendungsfällen mit berechnetem ROI-Potenzial.

Phase 2: Der "Minimum Viable Agent" (Wochen 5-12)

Ziel: Technischer Beweis und Aufbau der Governance.

1. **Workflow-Redesign:** Automatisieren Sie nicht den bestehenden Prozess! Ein schlechter Prozess wird durch KI nur schneller schlecht. Designen Sie den Prozess neu unter der Annahme, dass der Agent der Hauptakteur ist und der Mensch nur die Ausnahme behandelt.
2. **Shadow Mode Deployment:** Lassen Sie den Agenten parallel zum Menschen laufen, ohne dass er Aktionen ausführt. Vergleichen Sie die Entscheidungen des Agenten mit denen der Experten.
3. **Baseline-Messung:** Etablieren Sie Metriken für AHT (Average Handling Time), Fehlerraten und Kosten pro Ticket vor der Einführung.

Output: Ein funktionierender Agent im Schattenbetrieb mit einer Genauigkeit von >80%.

Phase 3: Die Vertrauensbrücke & Skalierung (Monate 3-6)

Ziel: Übergang von Überwachung zu Autonomie.

1. **Confidence Thresholds:** Implementieren Sie Schwellenwerte. Ist sich der Agent zu >90% sicher, handelt er autonom. Darunter: Weiterleitung an Menschen.
2. **Aktives Lernen:** Jede menschliche Korrektur wird protokolliert und für das Fine-Tuning des Modells genutzt.
3. **Change Management (Die 70%-Regel):** Investieren Sie 70% des Aufwands in die Menschen.⁸ Schulen Sie Mitarbeiter nicht darin, die KI zu "bedienen", sondern sie zu "trainieren" und komplexe Ausnahmen zu managen.

Output: Ein Agent im Live-Betrieb mit >50% Deflektionsrate. Erste Realisierung von Kosteneinsparungen.

Phase 4: Profit-Realisierung & Integration (Monat 6+)

Ziel: P&L-Wirksamkeit.

1. **Workforce-Anpassung:** Stoppen Sie die Nachbesetzung von Stellen in automatisierten Bereichen (natürliche Fluktuation nutzen). Verschieben Sie High-Performer in wertschöpfendere Rollen (z.B. Kundenservice zu Sales).
2. **Plattform-Strategie:** Abstrahieren Sie die Komponenten (Sicherheit, Logging, ERP-Anbindung) in eine zentrale "Agent Platform", um die Grenzkosten für den *nächsten* Agenten zu senken.

6. Deep Dives: Kennzahlen und Erfolgsmessung

Um den Erfolg gegenüber Stakeholdern nachzuweisen, bedarf es eines robusten KPI-Dashboards. Weiche Faktoren wie "Mitarbeiterzufriedenheit" reichen nicht aus.

Tabelle 2: Das Agentic AI KPI-Framework

Kategorie	KPI	Beschreibung	Zielwert (Benchmark)
Finanzen	Cost per Transaction	Gesamtkosten (Tech + Mensch) geteilt durch Volumen.	Reduktion um >50 %
Finanzen	LCOAI	Levelized Cost of AI (siehe Kapitel 2.3).	Muss < menschliche Kosten sein
Operativ	Deflection Rate	Anteil der Fälle, die ohne menschliches Zutun gelöst werden.	> 60 % (Top-Performer: >80 %)
Qualität	Resolution Accuracy	Anteil korrekter Lösungen (keine Wiedereröffnung des Tickets).	> 95 %
Technik	Hallucination Rate	Häufigkeit faktisch falscher Aussagen.	< 1 % (kritisch!)

Berechnung des ROI: Ein Praxisbeispiel

Szenario: IT-Helpdesk in einem mittelständischen Unternehmen (100.000 Tickets/Jahr).

- **Status Quo (Mensch):**
 - Kosten pro Ticket: 8,00 € (Vollkosten).
 - Gesamtkosten: 800.000 € p.a.
- **Szenario Agent (Investition):**
 - Entwicklung & Setup: 100.000 € (Einmalig).
 - Laufende Kosten (Hosting, Token, Wartung): 60.000 € p.a.
- **Ergebnis:**
 - Annahme: 60 % Deflektionsrate (60.000 Tickets autonom).
 - Verbleibende Tickets für Menschen: $40.000 * 8 \text{ €} = 320.000 \text{ €}$.
 - Gesamtkosten Neu: $320.000 \text{ €} (\text{Mensch}) + 60.000 \text{ €} (\text{KI}) = 380.000 \text{ €}$.
 - **Einsparung Jahr 1:** $800.000 \text{ €} - 380.000 \text{ €} - 100.000 \text{ €} (\text{Invest}) = \textbf{320.000 € Netto-Ersparnis}$.
 - **ROI Jahr 1:** 3,2x.

Dieses Rechenbeispiel⁷ verdeutlicht die immense Hebelwirkung, sobald die Fixkosten der Entwicklung durch das Volumen amortisiert werden.

7. Risikomanagement und Herausforderungen

Kein Transformationsprozess ist ohne Risiko. Die folgenden Fallstricke müssen proaktiv gemanagt werden.

7.1 Das "Rebound"-Risiko (Parkinsonsches Gesetz)

Effizienzgewinne führen oft dazu, dass Arbeit "expandiert". Wenn ein Bericht in 5 Minuten statt 5 Stunden erstellt wird, fordern Manager plötzlich 10 Berichte anstatt einem.

- **Mitigation:** Klare Governance über den Output. Nutzen Sie die gewonnene Zeit explizit für neue Wertschöpfung oder realisieren Sie die Einsparung durch Einstellungsstopps. Produktivität ohne Abbau von Arbeitsstunden senkt keine Kosten.⁸

7.2 Technische Schulden

Schnell zusammengebaute Agenten neigen zu Instabilität.

- **Mitigation:** Behandeln Sie Prompts als Code. Nutzen Sie Versionierung, automatisierte Tests (Eval-Frameworks) und CI/CD-Pipelines für Agenten.

7.3 Compliance und Haftung

Wer haftet, wenn ein Agent fälschlicherweise einen Rabatt gewährt?

- **Mitigation:** Definieren Sie klare finanzielle Autoritätsgrenzen (z.B. "Bis 50€ autonom, darüber Genehmigung"). Implementieren Sie Audit-Trails, die jede Entscheidung des Agenten revisionssicher protokollieren.¹⁷

8. Fazit: Ambition als Differenzierungsmerkmal

Abschließend zeigt unsere Analyse der McKinsey-Daten, dass der wichtigste Prädiktor für den Erfolg nicht die Technologie, sondern die **Ambition** ist.¹ Unternehmen, die KI nur nutzen, um "5 % Kosten zu sparen", scheitern oft an den Implementierungshürden. Unternehmen, die KI nutzen, um ihr Geschäftsmodell neu zu erfinden – z.B. durch 24/7-Echtzeit-Service oder vollautomatisierte Lieferketten –, realisieren die massiven Profit-Pools.

Wir bei Dreher Consulting empfehlen unseren Kunden, die Phase des Spielens zu beenden. Die Technologie ist reif. Die Ökonomie ist validiert. Es liegt nun an der Führungsebene, die organisatorischen Weichen zu stellen.

Der Weg "From Pilot to Profit" ist kein technisches Upgrade. Er ist eine operative Transformation.

Checkliste für Führungskräfte

Nutzen Sie diese Checkliste, um den Reifegrad Ihrer Initiative zu prüfen.

Strategische Ausrichtung

- [] **Zieldefinition:** Haben wir definiert, ob wir Kosten senken (Efficiency) oder wachsen (Innovation) wollen? (Vermeiden Sie Mischziele).
- [] **Budgetierung:** Haben wir Budget für die "J-Kurve" (Lernphase) eingeplant?
- [] **Workforce-Strategie:** Gibt es einen Plan für die Mitarbeiter, deren Aufgaben automatisiert werden? (Reskilling vs. Abbau).

Technische Bereitschaft

- [] **Data Readiness:** Sind die Kernsysteme (ERP, CRM) via API zugänglich? Sind die Daten sauber?
- [] **Sandboxing:** Haben wir eine sichere Umgebung, in der Agenten scheitern können, ohne Produktionsdaten zu gefährden?
- [] **Observability:** Können wir nachvollziehen, *warum* ein Agent eine Entscheidung getroffen hat? (Audit Logs).

Operative Umsetzung

- [] **Atomare Aufgaben:** Wurden die Prozesse in kleinste, logische Schritte zerlegt?
- [] **Ground Truth:** Gibt es einen "Gold Standard" an Antworten/Aktionen, gegen den der Agent getestet wird?
- [] **Change Management:** Sind die Mitarbeiter geschult, mit dem Agenten zu arbeiten (Exception Handling) statt gegen ihn?