

XINITY

xinity.ai

Xinity 2026
All rights reserved

VON ANTHROPIC CLAUDE ZU XINITY

KI-PLATTFORM MIGRATIONS-
WHITEPAPER SERIE 2026

LEGAL NOTICES

Xinity reminds you to carefully read through and completely understand all content in this section before you read or use this document. If you read or use this document, it is considered that you have identified and accepted all contents declared in this section.

- 1.** This document is published by Xinity for informational purposes. The contents are intended for legal and compliant business activities. You shall not use or disclose all or part of the contents to any third party without written permission from Xinity.
- 2.** This document may be subject to change without notice due to product upgrades, adjustment, and other reasons. Xinity reserves the right to modify the contents without notice.
- 3.** This document is only intended for product and service reference. Xinity provides this document for current products and services with current functions, which may be subject to change.
- 4.** All content including images, architecture design, page layout, and description text is owned by Xinity. You shall not use, modify, copy, or publish the content without written permission.
- 5.** If you discover any errors or mistakes within this document, please contact Xinity directly.

THE AUTHORS

CORE AUTHORS

Alexander Zehetmaier (CEO & Co-Founder, Xinity)

TECHNICAL REVIEW

Jonas (CTO & Co-Founder, Xinity)

EDITING AND DESIGN

Xinity Marketing Team

TARGET AUDIENCE

Dieser Leitfaden richtet sich an Engineering-Teams, CTOs und IT-Entscheidungstraeger, die derzeit Anthropic Claude API-Dienste (Claude Opus 4.6, Claude Sonnet 4.6, Claude Haiku 4.5, Messages API, Tool Use) nutzen und KI-Workloads auf eine souveraeene On-Premise-Infrastruktur migrieren muessen. Ob Sie Claude wegen seiner Sicherheitseigenschaften gewaehlt haben und nun regulatorische Anforderungen vollstaendige Datensouveraenitaet verlangen, oder ob Sie die Cloud-API-Abhaengigkeit fuer Geschaeftskontinuitaet eliminieren moechten -- dieses Whitepaper liefert die technischen Zuordnungen und Migrationsprozesse.

CONTENTS

1. Enterprise AI ohne Kompromisse: Warum Xinity die bessere Wahl ist

2. Ihr Anthropic-Stack, neu aufgebaut auf Xinity (Zugeordnet & Bereit)

2.1 Kern-Inferenz & Konversation

2.2 Tool Use & Funktionsaufrufe

2.3 Dokumentenverarbeitung & Analyse

2.4 Sicherheit & Guardrails

2.5 Plattform & Betrieb

3. Migrationsprozess

3.1 Bestandsaufnahme & Discovery

3.2 Infrastrukturplanung

3.3 Pilot-Migration

3.4 Vollstaendige Migration

4. Migrations-Werkzeuge & Beschleuniger

4.1 API-Uebersetzung

4.2 Observability & Betrieb

5. Naechste Schritte: Starten Sie Ihre Migration mit Xinity

1. ENTERPRISE AI OHNE KOMPROMISSE: WARUM XINITY DIE BESSERE WAHL IST

Wenn Ihr Unternehmen KI-Workloads in der Produktion betreibt, bietet die Migration von Cloud-gehosteten KI-APIs zur On-Premise-Plattform von Xinity etwas, das kein Cloud-Anbieter liefern kann: vollständige architektonische Souveränität über Ihre Daten, Modelle und Inferenz-Infrastruktur. Dies ist kein einfacher Anbieterwechsel -- es ist ein fundamentaler Wandel vom Mieten von KI-Kapazität zum Besitzen.

-- Architektonische Souveränität statt Richtlinien-Versprechen

Cloud-KI-Anbieter bieten vertraglichen Datenschutz durch Nutzungsbedingungen und Auftragsverarbeitungsverträge. Xinity liefert architektonische Souveränität: Ihre Daten verlassen niemals Hardware, die Sie physisch besitzen und kontrollieren. Für regulierte Branchen -- Gesundheitswesen, Recht, Finanzdienstleistungen, Medien und Fertigung -- ist diese Unterscheidung nicht akademisch. Es ist der Unterschied zwischen Compliance-Risiko und Compliance-Sicherheit. Keine ausländische Regierungsvorladung, keine Änderung der Cloud-Anbieter-Richtlinien und keine geopolitische Verschiebung kann Daten beeinflussen, die ausschließlich auf Ihren Räumlichkeiten existieren.

-- Planbare Wirtschaftlichkeit im Enterprise-Massstab

Cloud-KI-Preise skalieren mit dem Verbrauch: Jeder API-Aufruf, jedes Token, jede GPU-Stunde wird gemessen und abgerechnet. Xinity's On-Premise-Modell wandelt variable OPEX in planbare CAPEX um. Kunden, die Xinity Runtime auf ASUS Ascent GX10 Servern einsetzen, berichten von ca. 80% Kostenersparnis gegenüber vergleichbarer Cloud-Kapazität. Im Massstab bedeutet das ca. 320 EUR/Jahr Stromkosten gegenüber 18.600 EUR/Jahr für vergleichbare Cloud-Rechenleistung.

-- Latenzfreie Inferenz für kritische Anwendungen

On-Premise-KI eliminiert Netzwerk-Roundtrips zu entfernten Cloud-Regionen. Für latenzsensitive Anwendungen -- Echtzeit-Dokumentenanalyse, Qualitätskontrolle in der Produktion, klinische Entscheidungsunterstützung -- liefert lokale Inferenz konsistente Sub-Millisekunden-Antwortzeiten ohne Abhängigkeit von Internetverbindung, Cloud-Region-Verfügbarkeit oder grenzüberschreitenden Datentransfervorschriften.

-- Regulatorischer Rückenwind beschleunigt die Adoption

Der EU Digital Networks Act (vorgeschlagen Januar 2026) mit Compliance-Fristen im August 2026, die 20 Milliarden EUR InvestAI-Förderinitiative und aufkommende 'Buy European'-Beschaffungsregeln validieren die These der souveränen KI-Infrastruktur. Organisationen, die jetzt auf On-Premise-KI migrieren, positionieren sich vor den Regulierungen statt später hektisch reagieren zu müssen.

-- OpenAI-kompatible APIs -- migrieren ohne Neuentwicklung

Xinity Runtime stellt OpenAI-kompatible API-Endpunkte bereit. Das bedeutet: Ihr bestehender Anwendungscode, SDKs, Prompt-Bibliotheken und Orchestrierungsframeworks funktionieren mit minimalen Änderungen weiter. Sie ändern die Base-URL und den API-Key; Ihre Anwendungen bemerken keinen Unterschied.

2. IHR ANTHROPIC-STACK, NEU AUFGEBAUT AUF XINITY (ZUGEORDNET & BEREIT)

Dieser Abschnitt stellt ein Faehigkeiten-Mapping fuer die Migration von Anthropics Claude API zu Xinitys On-Premise-Plattform bereit. Anthropic verwendet ein proprietaraeres API-Format (Messages API), das sich vom OpenAI-Format unterscheidet. Xinity stellt einen OpenAI-kompatiblen Endpunkt bereit, sodass die Migration SDK- und Prompt-Format-Uebersetzung umfasst.

Kern-Inferenz & Konversation

Source Service	Xinity Equivalent	Migration Notes
Claude Sonnet 4.6 (Neuestes Balanciert)	Xinity Runtime (Mistral Medium 3 / Qwen3.5 32B)	Ausgewogene Leistung und Geschwindigkeit. OpenAI-kompatibler Endpunkt. Exzellent fuer Code, Analyse, Text.
Claude Opus 4.6 (Leistungsstaerkstes)	Xinity Runtime (Mistral Large 3 / Nemotron-Ultra-253B)	Maximale Faehigkeit fuer komplexe Aufgaben. Keine Token-basierte Abrechnung.
Claude Haiku 4.5 (Schnellstes)	Xinity Runtime (Qwen3.5 8B / Mistral Small 3)	Ultra-schnelle Inferenz. Ideal fuer Klassifizierung und Extraktion.

Tool Use & Funktionsaufrufe

Source Service	Xinity Equivalent	Migration Notes
Claude Tool Use	Xinity Runtime (OpenAI-Format Tool Calling)	Anthropic-Tool-Schemas ins OpenAI-Format uebersetzen. Lokal ausfuehren.
Claude Computer Use	Xinity + Open-Source Agents	On-Premise Computer-Interaktions-Agenten. Keine Screenshots an Cloud-APIs.
Claude MCP	Xinity + Lokaler MCP Server	MCP-Server On-Premise bereitstellen. Interne Datenbanken und APIs verbinden.

Dokumentenverarbeitung & Analyse

Source Service	Xinity Equivalent	Migration Notes
Claude PDF-Analyse	Xinity Runtime + PDF-Vorverarbeitungspipeline	Lokale Dokumentenaufnahme. Vertrauliche Dokumente sicher verarbeiten.
Claude Vision	Xinity Runtime (LLaVA / CogVLM / Qwen-VL)	On-Premise Bildverstaendnis. Medizinische Bilder lokal verarbeiten.
Anthropic Prompt Caching	Xinity Runtime (KV-Cache + Prompt-Prefix-Caching)	Lokales Kontext-Caching. Keine Cloud-Caching-Gebuehren.

Sicherheit & Guardrails

Source Service	Xinity Equivalent	Migration Notes
Claude Constitutional AI	Xinity Runtime + Guardrails (NeMo / LLM Guard)	Konfigurierbare Sicherheitsebenen. Guardrails nach Anwendungsfall anpassen.
Anthropic Content Policy	Xinity Custom Policy Engine	Eigene Inhaltsrichtlinien definieren. Volle Kontrolle ueber Modellverhalten.

Plattform & Betrieb

Source Service	Xinity Equivalent	Migration Notes
Anthropic Console	Xinity Admin Console (RBAC, SSO, Audit Logs)	Enterprise Identity-Integration. Vollstaendiger Audit-Trail On-Premise.
Anthropic Usage API	Xinity Monitoring (Prometheus / Grafana)	Echtzeit-Inferenz-Metriken. Keine Token-Abrechnung.
Batches API	Xinity Batch Processor	On-Premise Batch-Inferenz. Prioritaets-Warteschlange.

3. MIGRATIONSPROZESS

3.1 Bestandsaufnahme & Discovery

Anthropic API-Nutzung auditieren

Überprüfen Sie Ihre Anthropic Console, um alle aktiven Workloads zu katalogisieren. Dokumentieren Sie, welches Claude-Modell jeder Workload nutzt und welche spezifischen Fähigkeiten benötigt werden.

Anthropic-spezifische Features kartieren

Identifizieren Sie Workloads, die Anthropic-spezifische Features nutzen: Constitutional AI, XML-Tag-Konventionen, Messages API Schema.

Datensensitivität klassifizieren

Identifizieren Sie Workloads, bei denen Claudes Cloud-API Compliance-Risiken unter DSGVO oder nationalen Datenschutzgesetzen erzeugt.

3.2 Infrastrukturplanung

Hardware-Dimensionierung

Claude Opus 4.6-Nutzer benötigen typischerweise 253B-Parameter-Modelle (Nemotron-Ultra) auf Xinity. Claude Haiku 4.5-Nutzer können effiziente 8B Modelle nutzen.

API-Übersetzungsstrategie

Planen Sie die Übersetzung von Anthropic's Messages API zu OpenAI's Chat Completions API. Xinity bietet eine Übersetzungs-Middleware.

Prompt-Engineering-Anpassung

Claude-optimierte Prompts müssen getestet und angepasst werden. Planen Sie 1-2 Wochen für Prompt-Tests ein.

3.3 Pilot-Migration

Xinity Runtime bereitstellen

Installation und Konfiguration mit den abgestimmten Open-Weight-Modellen.

SDK-Migration

Von Anthropic's SDK zum OpenAI SDK wechseln, das auf Xinity zeigt: `from openai import OpenAI client = OpenAI(base_url='https://your-domain.com/v1', api_key='your-xinity-key')`

Qualitätsvalidierung

Validieren Sie, dass Xinitys Modelle äquivalente Qualität für Ihre spezifischen Aufgaben liefern.

3.4 Vollständige Migration

Phasenweiser Rollout

Souveraenitaetsblockierte Anwendungsaefelle zuerst, dann Hochvolumen-Workloads, dann verbleibende.

Guardrails-Konfiguration

Claudes Constitutional AI durch Xinitys konfigurierbare Guardrails ersetzen. Vorteil: Sie kontrollieren, was gefiltert wird.

Anthropic-Dienste dekommissionieren

API-Keys widerrufen, Abrechnungskonten schliessen. 90 Tage Rollback-Faehigkeit beibehalten.

4. MIGRATIONS-WERKZEUGE & BESCHLEUNIGER

4.1 API-Uebersetzung

Anthropic-zu-OpenAI Uebersetzungs-Middleware

Transparenter Proxy, der Anthropic Messages API Format ins OpenAI Chat Completions Format konvertiert. Null-Code-Migration.

Prompt-Migrations-Toolkit

Analysiert Claude-optimierte Prompts und passt sie fuer Open-Weight-Modelle an.

4.2 Observability & Betrieb

Xinity Dashboard

Umfassendes Monitoring mit vorkonfigurierten Dashboards.

Compliance & Audit-Modul

Vollstaendiger Audit-Trail mit Compliance-Berichtgenerierung. DSGVO, ISO 27001.

5. NAECHSTE SCHRITTE: STARTEN SIE IHRE MIGRATION MIT XINITY

Die Migration von Anthropic Claude zu Xinity erfordert API-Format-Uebersetzung und Prompt-Anpassung, aber der Prozess ist klar definiert und Xinity bietet automatisierte Werkzeuge.

So starten Sie:

1. Discovery-Gespraech vereinbaren -- Xinitys Solutions-Architekten analysieren Ihre Claude-API-Nutzung und identifizieren optimale Open-Weight-Modell-Matches.
2. Proof of Concept anfordern -- Testen Sie Xinity mit der Anthropic-Uebersetzungs-Middleware. Null Code-Aenderungen noetig.
3. Prompt-Anpassung planen -- Optimieren Sie Prompts fuer Open-Weight-Modelle.
4. Go Live mit voller Kontrolle -- Besitzen Sie Ihre KI-Infrastruktur, Ihre Daten und Ihre Inhaltsrichtlinien.

Kontakt: Web: xinity.ai E-Mail: contact@xinity.ai Standort: Wien, Oesterreich