

Azure Data Factory: Datenpipelines clever orchestrieren

Category: Online-Marketing

geschrieben von Tobias Hager | 6. Februar 2026



Azure Data Factory: Datenpipelines clever orchestrieren

Du hast Terabytes an Daten, ein halbes Dutzend Datenquellen, ein BI-Team am Rande des Nervenzusammenbruchs – und keiner weiß, wie man das alles effizient zusammenführt? Willkommen im Club. Die gute Nachricht: Mit Azure Data Factory baust du dir eine datengetriebene Infrastruktur, die nicht nur skaliert, sondern auch deinen Tech-Stack nicht in den Wahnsinn treibt. Die schlechte:

Wenn du's falsch machst, wird's teuer, langsam und unwartbar. Hier kommt dein ungeschönter Deep Dive in die Welt der modernen Datenorchestrierung.

- Was Azure Data Factory (ADF) eigentlich ist – und warum du es brauchst
- Wie du mit ADF Datenpipelines zwischen On-Prem und Cloud effizient orchestrierst
- ETL vs. ELT – und warum ADF beides kann (aber du wissen musst, wann was sinnvoll ist)
- Die Kernkomponenten: Pipelines, Datasets, Linked Services und Trigger erklärt
- Best Practices für skalierbare, wartbare und performante Pipelines
- Monitoring, Logging und Alerting in ADF – so behältst du den Überblick
- Security, Netzwerkanbindung und Zugriffskontrolle in Enterprise-Setups
- Typische Fehler – und wie du sie garantiert vermeidest
- Schritt-für-Schritt-Blueprint für eine produktionsreife ADF-Implementierung
- Warum ADF kein “Klick-und-fertig“-Tool ist – und das auch gut so ist

Was ist Azure Data Factory?

Datenintegration auf Enterprise-Niveau

Azure Data Factory (ADF) ist Microsofts cloudbasierter Dienst zur Datenintegration und -orchestrierung. Er ermöglicht es, Daten aus unterschiedlichsten Quellen zu extrahieren, zu transformieren und in Zielsysteme zu laden – klassisch als ETL (Extract, Transform, Load) oder moderner als ELT (Extract, Load, Transform). Dabei unterstützt ADF mehr als 90 native Konnektoren – von SQL Server über Amazon S3 bis zu SAP HANA. Kurz gesagt: Wenn du irgendwo Daten hast, kannst du sie mit ADF bewegen.

Der Clou von ADF liegt in seiner Architektur. Statt ETL-Jobs auf dedizierten Servern auszuführen, orchestriert ADF Workflows in der Cloud – skalierbar, serverlos und kosteneffizient. Die Rechenleistung kommt aus sogenannten Integration Runtimes (IR), die entweder als Azure IR (in der Cloud), Self-hosted IR (On-Prem) oder SSIS IR (für Legacy-SSIS-Pakete) bereitgestellt werden. Damit ist ADF nicht nur für Cloud-native Projekte geeignet, sondern auch für hybride Szenarien – was in der Praxis extrem wichtig ist.

Im Gegensatz zu klassischen ETL-Tools wie Talend oder Informatica setzt ADF auf ein deklaratives Low-Code-Modell. Das klingt erstmal nach Clicky-Bunti-Tool, ist aber in Wahrheit ein hochgradig konfigurierbares Framework, das sich über Azure DevOps oder ARM-Templates vollständig automatisieren lässt. Wer will, kann ADF-Pipelines per YAML bauen, testen und deployen – ganz im Sinne von Infrastructure-as-Code.

Und ja, ADF ist kein BI-Tool. Es ist keine Reporting-Plattform. Es ist der Klebstoff zwischen deinen Datenquellen und deinen Analyse-Tools. Wer ADF mit Power BI verwechselt, hat das Prinzip nicht verstanden. ADF ist der

Transportlayer – und der entscheidet, ob deine Daten schnell, sicher und zuverlässig dorthin kommen, wo sie gebraucht werden.

Datenpipelines in Azure Data Factory: Architektur, Aufbau und Best Practices

Eine Pipeline in ADF ist nichts anderes als ein Workflow – eine definierte Abfolge von Aktivitäten, die Daten bewegen, transformieren oder triggern. Klingt simpel, wird aber schnell komplex, wenn du mehrere Datenquellen, Bedingungen, Schleifen und Abhängigkeiten orchestrieren willst. Deshalb ist ein klares Architekturverständnis Pflicht.

Die vier zentralen Bausteine in ADF sind:

- **Linked Services:** Verbindungen zu Datenquellen oder -zielen. Denk an sie wie an Connection Strings mit Authentifizierung und Zugriffseinstellungen.
- **Datasets:** Beschreiben die Struktur der Daten, die du verarbeiten willst – z. B. eine Tabelle in SQL Server oder ein Blob in Azure Storage.
- **Activities:** Die eigentliche Arbeit – Datenkopieren, SQL ausführen, Webhooks aufrufen, Datentransformationen via Data Flow etc.
- **Pipelines:** Die Klammer, die alles zusammenhält. Hier definierst du den Ablauf, Bedingungen, Wiederholungen und Fehlerbehandlung.

Best Practices? Davon gibt's viele, aber hier die wichtigsten:

- Modularisieren – kleine, wiederverwendbare Pipelines statt monolithischer Monster
- Parametrisierung – keine Hardcodierung von Pfaden, Tabellen oder Connection Strings
- Fehlerhandling – mit Try-Catch-Pattern, Retry-Logik und Alerting
- Logging – jedes Activity-Ergebnis gehört ins zentrale Monitoring
- Deployment via DevOps – kein Klick-Klick-Deploy im Live-System

Ein unterschätzter Punkt: ADF skaliert horizontal. Das heißt, mehrere Pipelines laufen parallel – sofern die Integration Runtime das hergibt. Wer hier nicht aufpasst, erzeugt schnell Bottlenecks oder verursacht unnötige Kosten durch überdimensionierte Instanzen.

ETL oder ELT mit ADF: Was wann Sinn macht

Die klassische ETL-Strategie kennt jeder: Daten aus Quelle extrahieren, transformieren, ins Ziel schreiben. Das Problem: Bei großen Datenmengen oder sehr komplexen Transformationen wird das teuer, langsam und inkonsistent.

Deshalb setzen moderne Architekturen vermehrt auf ELT – also erst laden, dann transformieren, direkt im Zielsystem.

ADF kann beides – aber nur, wenn du weißt, was du tust. ETL bietet sich an, wenn du bereits während der Extraktion Daten anreichern, filtern oder aggregieren musst. ELT ist ideal, wenn dein Zielsystem (z. B. Azure Synapse oder Snowflake) genug Rechenpower hat, um Transformationen effizient zu fahren.

ADF Data Flows sind das zentrale Tool für Transformationen. Sie sind visuelle, Spark-basierte Workflows, die auf Azure Data Lake Storage operieren. Vorteil: massive Parallelisierung, kein eigenes Spark-Cluster nötig. Nachteil: teuer, wenn falsch konfiguriert. Alternativ kannst du Transformationen auch in Stored Procedures verlagern – besonders effizient bei Datenbanken mit starkem MPP-Backend wie Synapse.

Pro-Tipp: Verwende Data Flows nur für Transformationen, die du nicht effizient im Zielsystem erledigen kannst. Alles andere ist Verschwendung. Und ja – auch “No-Code“-Tools kosten, wenn sie falsch eingesetzt werden.

Monitoring, Logging und Alerting in Azure Data Factory

Wer Daten bewegt, muss wissen, was wann wo passiert – und was nicht. ADF bietet dafür ein solides Set an Monitoring-Features, aber auch hier gilt: Nur wer's konfiguriert, profitiert davon. Standardmäßig speichert ADF alle Pipeline-Runs, Activity-Status und Trigger-Logs. Diese Daten kannst du im ADF-UI durchforsten oder per REST API abfragen.

Für ernsthaftes Monitoring solltest du ADF mit Azure Monitor, Log Analytics und Application Insights kombinieren. Das ermöglicht zentrale Dashboards, Alerts bei Fehlern oder Performance-Engpässen und automatisierte Reaktionen (z. B. Restart von Pipelines bei Timeout).

Und Logging? Alles, was in ADF passiert, gehört ins zentrale Log. Das geht über Diagnostic Settings, in denen du festlegst, welche Events wohin gestreamt werden – z. B. in ein Log Analytics Workspace oder ein Storage Account. Damit kannst du z. B. alle fehlgeschlagenen Aktivitäten automatisch analysieren oder sogar Machine-Learning-basiert Clustern.

Für produktive Umgebungen sind folgende Punkte Pflicht:

- Activity-Status als Metrik für Alerts konfigurieren
- Trigger-Ausführungen regelmäßig auswerten
- Integration mit Azure DevOps für Deployment-Tracking
- Eigene Logging-Strategie für Custom Activities

Wer das ignoriert, steht beim nächsten Fehler ohne Kontext da. Und das kostet – Zeit, Daten und Nerven.

Security, Zugriff und Netzwerkarchitektur: ADF in der Enterprise-Praxis

Azure Data Factory ist ein PaaS-Dienst – das heißt: Microsoft managed die Infrastruktur, aber du bist für Konfiguration und Sicherheit verantwortlich. Und da fangen die Probleme an. Denn viele ADF-Setups sind offen wie Scheunentore oder so restriktiv, dass nichts mehr funktioniert.

Die wichtigsten Sicherheitsfeatures in ADF sind:

- Managed Identity: Authentifiziere ADF gegenüber Azure-Ressourcen ohne Hardcodierung von Credentials
- Private Endpoints: Vermeide öffentliche IPs durch VNET-Integration
- Rollenbasierte Zugriffskontrolle (RBAC): Präzise Steuerung, wer was im ADF tun darf
- Key Vault Integration: Sichere Speicherung von Secrets, Connection Strings und Zugangsdaten

Besonders wichtig in hybriden Szenarien: die Self-hosted Integration Runtime (SHIR). Sie läuft auf einem On-Prem-Server oder in einer VM und ermöglicht Zugriff auf lokale Datenquellen – ohne VPN, aber über sichere Outbound-Verbindungen. Wer hier pfuscht, riskiert Datenlecks oder unerreichbare Systeme.

Im Enterprise-Umfeld ist auch das Thema Netzwerksegmentierung entscheidend. Trenne deine ADF-Umgebung logisch von anderen Ressourcen, nutze NSGs, Firewall-Regeln und Conditional Access Policies. Und ja: Compliance ist kein Nebenthema. Wer personenbezogene Daten verarbeitet, braucht Auditing, Logging und Zugriffshistorien – alles machbar mit ADF, aber eben nicht “by default”.

Fazit: Azure Data Factory ist mächtig – aber kein Selbstläufer

Azure Data Factory ist das Schweizer Taschenmesser für moderne Datenintegration. Es orchestriert, transformiert und verbindet – quer durch Cloud, On-Prem und hybride Landschaften. Aber: Wer ADF als “Klick-und-fertig“-Tool betrachtet, wird scheitern. Erfolgreiche ADF-Projekte leben von sauberer Architektur, durchdachter Orchestrierung, automatisiertem Deployment und konsequentem Monitoring.

Ja, ADF kann vieles. Aber es zwingt dich auch, deine Datenarchitektur zu durchdenken. Und das ist gut so. Denn in einer Welt, in der Daten das neue Öl

sind, willst du keine Pipeline, die leckt – sondern eine, die skaliert.
Willkommen in der Realität der Datenorchestrierung. Willkommen bei 404.