

# Ausarbeitung zum Pipeline-Architekturstil

Alber Jonas, Schweitzer Tim

3. Mai 2025

## Zusammenfassung

Diese Ausarbeitung befasst sich mit dem Pipeline-Architekturstil, auch bekannt als Pipes-and-Filters-Architektur. Es werden die grundlegenden Konzepte, Komponenten und Funktionsweisen dieses Stils erläutert. Darauf aufbauend werden typische Anwendungsbereiche und konkrete Beispiele vorgestellt. Abschließend erfolgt eine kritische Bewertung der architektonischen Eigenschaften, indem die wesentlichen Vor- und Nachteile des Pipeline-Stils diskutiert werden. Die Darstellung basiert auf etablierten Quellen der Softwarearchitektur-Literatur.

## Inhaltsverzeichnis

<b>1</b>	<b>Einführung und Kernkonzepte</b>	<b>2</b>
1.1	Filter: Die Verarbeitungseinheiten	2
1.2	Pipes: Die Kommunikationskanäle	3
1.3	Stärke durch Komposition: Das Unix-Beispiel	3
<b>2</b>	<b>Anwendungsbereiche und Beispiele</b>	<b>3</b>
2.1	Beispiel: Verarbeitung von Telemetriedaten mit Kafka	4
<b>3</b>	<b>Bewertung: Vor- und Nachteile</b>	<b>5</b>
3.1	Vorteile	5
3.2	Nachteile	6
3.3	Fazit zur Bewertung	6

# 1 Einführung und Kernkonzepte

In der Welt der Softwarearchitektur existieren diverse etablierte Muster und Stile, die als Blaupausen für die Strukturierung von Softwaresystemen dienen. Einer dieser fundamentalen und weit verbreiteten Stile ist die Pipeline-Architektur, oft auch als Pipes-and-Filters-Architektur bezeichnet [1]. Dieser Stil entstand als natürliche Konsequenz der Modularisierung von Anwendungsfunktionalität in eigenständige Verarbeitungsschritte [1]. Seine Prinzipien sind vielen Entwicklern und Architekten intuitiv vertraut, insbesondere durch die Kommandozeilen-Interpreter von Unix-basierten Betriebssystemen wie Bash oder Zsh, wo Befehle über Pipes (‘|’) miteinander verkettet werden [1]. Auch in der funktionalen Programmierung finden sich Parallelen, da viele Sprachelemente die Komposition von Funktionen fördern, was dem Konzept der Pipeline ähnelt [1]. Moderne Datenverarbeitungsparadigmen wie MapReduce bauen ebenfalls auf dieser Grundtopologie auf [1]. Es ist hervorzuheben, dass der Pipeline-Stil nicht auf Low-Level-Systemprogrammierung beschränkt ist, sondern ebenso effektiv für die Gestaltung komplexer Geschäftsanwendungen eingesetzt werden kann [1].

Die Architektur basiert auf zwei zentralen Komponenten: **Filtern** (Filters) und **Pipes** [1].

## 1.1 Filter: Die Verarbeitungseinheiten

Filter repräsentieren die aktiven Bausteine der Architektur, die für die eigentliche Datenverarbeitung zuständig sind [1]. Charakteristisch für Filter ist ihre Eigenständigkeit und Unabhängigkeit voneinander. Idealerweise sind sie zustandslos (‘stateless’), was bedeutet, dass das Ergebnis einer Verarbeitung ausschließlich von den aktuellen Eingabedaten abhängt und nicht von vorherigen Interaktionen [1]. Ein Kernprinzip des Designs ist die *Single Responsibility*: Jeder Filter sollte genau eine klar definierte Aufgabe erfüllen [1]. Steht eine komplexere Verarbeitungslogik an, wird diese nicht in einem einzigen, monolithischen Filter implementiert, sondern auf eine Sequenz spezialisierter Filter verteilt, die nacheinander ausgeführt werden [1]. Diese Granularität fördert die Wiederverwendbarkeit und Testbarkeit der einzelnen Komponenten.

Innerhalb der Pipeline-Architektur unterscheidet man typischerweise vier Arten von Filtern [1]:

- **Producer (Erzeuger oder Quelle):** Dieser Filtertyp steht am Anfang einer Pipeline und ist für die Generierung der initialen Daten verantwortlich. Er empfängt keine Daten von einem vorherigen Filter, sondern agiert als Datenquelle, z.B. durch Lesen aus einer Datei, einer Datenbank oder einem Netzwerk-Socket. Er wird daher auch als »Source« bezeichnet [1].
- **Transformer (Umwandler):** Ein Transformer nimmt Daten von einer Pipe entgegen, führt eine Transformation darauf aus (z.B. Daten anreichern, Format ändern, Werte berechnen) und leitet das modifizierte Ergebnis an die nächste Pipe weiter. Dieser Filtertyp verkörpert das zentrale Verarbeitungselement in vielen Pipelines. Funktionale Programmierer erkennen hier eine Analogie zur ‘map’-Operation [1].
- **Tester (Filter im engeren Sinne):** Dieser Filter empfängt ebenfalls Daten, prüft diese jedoch anhand bestimmter Kriterien. Abhängig vom Ergebnis dieser Prüfung leitet er die Daten entweder unverändert weiter, verwirft sie oder leitet sie möglicherweise auf unterschiedliche Ausgänge (obwohl dies die klassische lineare Pipeline erweitert). Er selektiert also Datenströme. Funktionale Programmierer mögen hier an ‘filter’- oder ‘reduce’-Operationen denken, je nach spezifischer Implementierung [1].
- **Consumer (Verbraucher oder Senke):** Der Consumer bildet das Ende einer Pipeline. Er nimmt die final verarbeiteten Daten entgegen und führt eine abschließende Aktion durch, wie das Speichern der Daten in einer Datenbank, das Anzeigen auf einer Benutzeroberfläche oder das Senden an ein externes System. Er wird auch als »Sink« bezeichnet [1].

## 1.2 Pipes: Die Kommunikationskanäle

Pipes stellen die Verbindungen zwischen den Filtern dar und definieren den Datenfluss durch die Architektur [1]. Sie fungieren als gerichtete Kommunikationskanäle, die typischerweise genau einen Filterausgang mit genau einem Filtereingang verbinden (Punkt-zu-Punkt-Verbindung). Dies unterscheidet sie von Broadcast-Mechanismen, bei denen Daten an mehrere Empfänger gleichzeitig gesendet werden, was in der Regel aus Performancegründen vermieden wird [1]. Eine Pipe nimmt die Ausgabe eines Filters entgegen und leitet sie als Eingabe an den nächsten Filter weiter [1].

Die Daten, die als Nutzlast (‘payload’) durch die Pipes transportiert werden, können prinzipiell jedes beliebige Format haben – von einfachen Textströmen über strukturierte Daten wie XML oder JSON bis hin zu binären Objekten [1]. Für eine hohe Performance und geringe Kopplung zwischen den Filtern wird jedoch oft empfohlen, das Datenformat möglichst standardisiert und die Datenmenge pro Nachricht relativ klein zu halten [1]. Die Implementierung von Pipes kann vielfältig sein: von einfachen In-Memory-Warteschlangen in monolithischen Anwendungen bis hin zu komplexeren Message-Queues oder Netzwerkprotokollen in verteilten Systemen.

## 1.3 Stärke durch Komposition: Das Unix-Beispiel

Die wahre Stärke der Pipeline-Architektur liegt in der Komposition einfacher, wiederverwendbarer Bausteine zu komplexen Verarbeitungsketten. Die Unidirektionalität der Pipes und die Fokussierung der Filter auf einzelne Aufgaben machen sie hochgradig modular und wiederverwendbar [1]. Dieses Prinzip ist vielen Entwicklern aus der Unix-Shell-Nutzung bestens bekannt [1].

Eine oft zitierte Anekdote illustriert eindrucksvoll die Eleganz und Mächtigkeit dieses Ansatzes [1]: Donald Knuth, ein renommierter Informatiker, wurde einst gebeten, ein Programm zu schreiben, das eine Textdatei einliest, die  $n$  häufigsten Wörter identifiziert und diese zusammen mit ihrer Häufigkeit sortiert ausgibt. Sein Ergebnis war ein über zehn Seiten langes, dokumentiertes Pascal-Programm, das einen eigens entwickelten Algorithmus implementierte [1]. Doug McIlroy, einer der Väter von Unix und Pipes, präsentierte daraufhin eine Lösung als einfaches Shell-Skript, das dieselbe Aufgabe löste – wesentlich kürzer, verständlicher und oft sogar performanter. Dieses Skript, das problemlos in eine kurze Nachricht passt, lautete:

```
tr -cs A-Za-z '
textbackslashn' < input.txt | tr A-Z a-z | sort | uniq -c | sort -rn | head
-n ${1:-10}
```

Hierbei wird der Text zunächst in einzelne Wörter zerlegt (‘tr’), in Kleinbuchstaben umgewandelt (‘tr’), alphabetisch sortiert (‘sort’), Duplikate gezählt (‘uniq -c’), nach Häufigkeit numerisch und absteigend sortiert (‘sort -rn’) und schließlich die ersten  $n$  (hier standardmäßig 10) Zeilen ausgegeben (‘head’). Jeder Befehl ist ein Filter, das Pipe-Symbol ‘|’ repräsentiert die Pipe. Dieses Beispiel zeigt, wie mächtig die Kombination einfacher Werkzeuge sein kann, eine Komplexität, die selbst die Unix-Erfinder manchmal überraschte [1].

## 2 Anwendungsbereiche und Beispiele

Der Pipeline-Architekturstil ist aufgrund seiner Struktur und Eigenschaften für eine Vielzahl von Anwendungsfällen prädestiniert. Insbesondere eignet er sich hervorragend für Systeme, die eine sequentielle, unidirektionale Verarbeitung von Datenströmen erfordern [1]. Die klare Trennung der Verarbeitungsschritte und der definierte Datenfluss machen ihn robust und nachvollziehbar für bestimmte Problemklassen.

Typische Anwendungsbereiche umfassen [1]:

- **Elektronischer Datenaustausch (EDI):** Im Bereich EDI müssen oft Daten zwischen verschiedenen Formaten (z.B. unternehmensspezifische Formate, Industriestandards wie EDIFACT oder XML-Derivate) konvertiert werden. Pipeline-basierte Werkzeuge zerlegen diesen Prozess in Schritte wie Validierung, Mapping, Anreicherung und Formatkonvertierung, wobei jeder Schritt von einem spezialisierten Filter durchgeführt wird.
- **ETL-Prozesse (Extract, Transform, Load):** ETL-Werkzeuge, die Daten aus Quellsystemen extrahieren, transformieren und in Zielsysteme (oft Data Warehouses) laden, nutzen inhärent das Pipeline-Muster. Die Phasen Extraktion (Producer), Transformation (mehrere Transformer und Tester) und Laden (Consumer) bilden eine natürliche Pipeline. Die Architektur unterstützt hierbei sowohl den Datenfluss als auch die schrittweise Modifikation der Daten.
- **Orchestrierung und Mediation in Integrationsszenarien:** Werkzeuge wie Apache Camel oder Mule ESB verwenden oft Pipeline-Konzepte, um komplexe Integrationsprozesse und Geschäftsworkflows zu modellieren. Nachrichten oder Datenobjekte werden durch eine Kette von Verarbeitungsschritten (z.B. Protokolltransformation, Routing, Inhaltsanreicherung) geleitet, die als Filter in einer Pipeline konfiguriert sind.
- **Compiler-Konstruktion:** Klassische Compiler arbeiten oft nach einem Pipeline-Prinzip: Lexikalische Analyse -> Syntaktische Analyse -> Semantische Analyse -> Zwischencode-Erzeugung -> Optimierung -> Zielcode-Erzeugung. Jede Phase nimmt die Ausgabe der vorherigen als Eingabe und führt eine spezifische Transformation durch.
- **Verarbeitung von Streaming-Daten:** Systeme, die kontinuierliche Datenströme verarbeiten (z.B. Sensordaten, Log-Events, Finanztransaktionen), setzen häufig Pipelines ein. Frameworks wie Apache Kafka Streams, Flink oder Spark Streaming ermöglichen die Definition von Verarbeitungstopologien, die dem Pipes-and-Filters-Modell folgen.

## 2.1 Beispiel: Verarbeitung von Telemetriedaten mit Kafka

Ein anschauliches Beispiel für die Anwendung der Pipeline-Architektur ist die Verarbeitung von Telemetriedaten, die von verteilten Diensten erzeugt und über Apache Kafka gestreamt werden [1]. Ziel könnte sein, aus dem rohen Datenstrom spezifische Metriken wie Antwortzeiten und Verfügbarkeit zu extrahieren und zu aggregieren.

Eine mögliche Pipeline könnte wie folgt aufgebaut sein [1]:

1. **Producer ('ServiceInfoCapture'):** Dieser Filter abonniert das relevante Kafka-Topic, in das die Rohdaten der Telemetrie (z.B. Logs oder Ereignisse über Serviceaufrufe) geschrieben werden. Er fungiert als Quelle der Pipeline und liest die Nachrichten aus Kafka.
2. **Tester ('DurationFilter'):** Die vom Producer gelesenen Daten werden an diesen Filter weitergeleitet. Seine Aufgabe ist es zu prüfen, ob die aktuelle Nachricht Informationen über die Dauer (Antwortzeit) eines Serviceaufrufs enthält (z.B. einen Messwert in Millisekunden). Dies illustriert die klare Trennung: Der 'ServiceInfoCapture'-Filter kümmert sich nur um das Lesen aus Kafka, der 'DurationFilter' nur um die Identifikation relevanter Daten für diesen Zweig der Verarbeitung.
3. **Transformer ('DurationCalculator'):** Wenn der 'DurationFilter' feststellt, dass die Daten die Antwortzeit betreffen, leitet er sie an diesen Filter weiter. Der 'DurationCalculator' könnte beispielsweise die Zeit von Millisekunden in Sekunden umrechnen, statistische Berechnungen durchführen (z.B. gleitender Durchschnitt) oder die Daten für die spätere Aggregation vorbereiten.
4. **Tester ('UptimeFilter'):** Falls der 'DurationFilter' die Nachricht nicht als relevant für die Antwortzeit einstuft, könnte sie alternativ an den 'UptimeFilter' geleitet werden. Dieser prüft, ob die Nachricht stattdessen Informationen zur Verfügbarkeit oder zum Betriebsstatus des Dienstes enthält (z.B. Heartbeat-Signale, Fehlermeldungen).
5. **Transformer ('UptimeCalculator'):** Wenn der 'UptimeFilter' Uptime-relevante Daten erkennt, werden diese an den

‘UptimeCalculator‘ übergeben. Dieser Filter könnte den Status interpretieren und daraus Verfügbarkeitsmetriken ableiten (z.B. Berechnung der Uptime-Prozentrage über ein Zeitfenster). 6. **Consumer (‘DatabaseOutput‘ oder ‘MetricsAggregator‘):** Die von ‘DurationCalculator‘ oder ‘UptimeCalculator‘ transformierten Metrikdaten werden schließlich an einen oder mehrere Consumer übergeben. Ein ‘DatabaseOutput‘-Filter könnte die berechneten Metriken in einer Zeitreihendatenbank (z.B. InfluxDB, MongoDB) speichern [1]. Ein anderer Consumer könnte die Metriken an ein Monitoring-System weiterleiten.

Daten, die weder vom ‘DurationFilter‘ noch vom ‘UptimeFilter‘ als relevant eingestuft werden, verlassen die Pipeline an dieser Stelle, da sie für den spezifischen Zweck dieser Verarbeitungskette nicht benötigt werden [1]. Dieses Beispiel demonstriert, wie durch die Kombination verschiedener Filtertypen ein komplexer Verarbeitungsworkflow strukturiert werden kann, wobei jeder Filter eine klar umrissene Aufgabe übernimmt und der Datenfluss explizit durch die Pipes definiert ist.

### 3 Bewertung: Vor- und Nachteile

Wie jeder Architekturstil weist auch die Pipeline-Architektur spezifische Stärken und Schwächen auf, die bei der Entscheidung für oder gegen ihre Anwendung in einem konkreten Projektkontext abgewogen werden müssen. Die Bewertung erfolgt anhand gängiger architektonischer Qualitätsmerkmale [1].

#### 3.1 Vorteile

Die Pipeline-Architektur bietet eine Reihe signifikanter Vorteile, die sie für bestimmte Anwendungsfälle attraktiv machen:

- **Einfachheit und Verständlichkeit:** Der lineare Datenfluss und die klare Aufgabentrennung der Filter machen die Architektur relativ einfach zu verstehen und nachzuvollziehen. Dies reduziert die kognitive Last für Entwickler und erleichtert die Einarbeitung neuer Teammitglieder [1].
- **Modularität und Wiederverwendbarkeit:** Filter sind per Definition unabhängige Einheiten mit einer spezifischen Aufgabe. Sie können oft ohne Modifikation in verschiedenen Pipelines wiederverwendet werden. Die lose Kopplung über die Pipes erlaubt es, Filter auszutauschen, neu anzuordnen oder hinzuzufügen, ohne andere Filter direkt zu beeinflussen [1]. Dies fördert die Wartbarkeit und Flexibilität des Systems.
- **Technische Partitionierung:** Die Logik der Anwendung ist klar entlang der technischen Verarbeitungsschritte (Filtertypen wie Producer, Transformer, Tester, Consumer) aufgeteilt. Dies unterstützt eine strukturierte Entwicklung [1].
- **Flexibilität bei der Komposition:** Die Möglichkeit, komplexe Aufgaben durch die Verkettung einfacher Filter zu lösen, bietet eine hohe Flexibilität im Design von Verarbeitungsprozessen. Neue Anforderungen können oft durch Hinzufügen oder Austauschen von Filtern implementiert werden.
- **Natürliche Parallelisierbarkeit (potenziell):** Obwohl oft monolithisch implementiert, bietet das Modell prinzipiell Potenzial für Parallelisierung. Einzelne Filterinstanzen könnten parallel arbeiten (Task-Parallelität), oder die Pipeline selbst könnte auf mehrere Datenströme parallel angewendet werden (Daten-Parallelität), sofern die Filter zustandslos sind und die Pipes dies unterstützen (z.B. durch konkurrente Queues).

- **Geringe Gesamtkosten (in bestimmten Szenarien):** Durch die Einfachheit, Verständlichkeit und Wiederverwendbarkeit können die Entwicklungs- und Wartungskosten, insbesondere für klar definierte sequentielle Prozesse, vergleichsweise niedrig sein [1].

### 3.2 Nachteile

Trotz ihrer Vorteile bringt die Pipeline-Architektur auch einige Nachteile und Einschränkungen mit sich, die insbesondere bei komplexen oder hochskalierbaren Systemen relevant werden:

- **Monolithisches Deployment (häufig):** In vielen Fällen wird eine Pipeline-Architektur als eine einzige Einheit (Monolith) bereitgestellt [1]. Das bedeutet, das gesamte System bildet ein einzelnes architektonisches Quantum. Änderungen an einem einzigen Filter erfordern potenziell das Testen und Deployen der gesamten Anwendung [1].
- **Eingeschränkte Skalierbarkeit und Elastizität:** Aufgrund des typischerweise monolithischen Deployments ist die Skalierbarkeit oft auf die Skalierung der gesamten Anwendung beschränkt. Es ist schwierig, nur einzelne, ressourcenintensive Filter gezielt zu skalieren, ohne komplexe interne Mechanismen wie Threading oder asynchrone Verarbeitung einzuführen, für die der Stil nicht primär ausgelegt ist [1]. Die horizontale Skalierung des gesamten Monolithen kann ineffizient sein, wenn nur ein kleiner Teil der Pipeline der Flaschenhals ist. Dies führt zu einer niedrigen Bewertung bei Elastizität und Skalierbarkeit (oft nur ein Stern) [1].
- **Geringe Fehlertoleranz und Verfügbarkeit:** In einer monolithischen Implementierung kann der Ausfall eines kritischen Filters (z.B. durch einen Fehler oder Ressourcenmangel wie Speicher) die gesamte Pipeline oder sogar die gesamte Anwendung zum Stillstand bringen [1]. Da es nur ein Deployment-Quantum gibt, fehlt die Isolation, die verteilte Systeme bieten. Die durchschnittliche Wiederherstellungszeit (MTTR) kann bei Monolithen relativ hoch sein (Minuten bis Stunden), was die Verfügbarkeit beeinträchtigt [1].
- **Mittelmäßige Testbarkeit und Bereitstellbarkeit:** Obwohl einzelne Filter gut isoliert getestet werden können (Unit-Tests), erfordert jede Änderung am Code oft einen vollständigen Integrationstest und die erneute Bereitstellung des gesamten Monolithen. Dies verlangsamt den Entwicklungszyklus und erhöht das Risiko bei Deployments im Vergleich zu feingranulareren Architekturen wie Microservices [1]. Die Bewertung für Bereitstellbarkeit ('Deployability') und Testbarkeit ist daher oft nur durchschnittlich [1].
- **Potenzielle Latenzprobleme:** In langen Pipelines kann sich die Latenz summieren, da Daten sequenziell durch viele Filter fließen müssen. Wenn die Pipes zudem I/O-intensiv implementiert sind (z.B. über Festplatten-basierte Queues statt In-Memory), kann dies die End-to-End-Verzögerung weiter erhöhen.
- **Gemeinsames Datenformat als Herausforderung:** Die Notwendigkeit eines gemeinsamen Datenformats für die Pipes kann einschränkend sein oder zu Overhead führen, wenn Filter sehr unterschiedliche Datenrepräsentationen bevorzugen würden. Änderungen am Datenformat können Anpassungen in vielen Filtern erfordern.
- **Schwierigkeiten bei komplexen Kontrollflüssen:** Der Stil eignet sich am besten für lineare Abläufe. Komplexe Verzweigungen, Zusammenführungen oder iterative Prozesse können schwierig und unübersichtlich in einer reinen Pipeline-Struktur abzubilden sein.

### 3.3 Fazit zur Bewertung

Zusammenfassend lässt sich sagen, dass die Pipeline-Architektur ein eleganter und effektiver Stil für Probleme ist, die sich gut als sequentielle Kette von Verarbeitungsschritten modellieren

lassen. Ihre Stärken liegen in der Einfachheit, Modularität und Wiederverwendbarkeit der Komponenten, was oft zu kosteneffizienten Lösungen führt [1]. Die größten Herausforderungen ergeben sich typischerweise aus der oft monolithischen Natur der Implementierung, was zu Einschränkungen bei Skalierbarkeit, Fehlertoleranz, Testbarkeit und Bereitstellbarkeit führt, insbesondere bei großen, komplexen oder hochverfügbaren Systemen [1]. Die Wahl dieses Stils sollte daher sorgfältig gegen die spezifischen Anforderungen des Systems abgewogen werden.

## Literatur

- [1] Mark Richards and Neal Ford. *Handbuch Moderner Softwarearchitektur: Architekturstile, Patterns und Best Practices*. O'Reilly Verlag GmbH & Co. KG, 2020. Zugriff über ProQuest Ebook Central.