

O'REILLY®

Helion 

Nowoczesne architektury danych

Przewodnik po hurtowni danych,
siatce danych oraz Data Fabric
i Data Lakehouse



James Serra

Tytuł oryginału: Deciphering Data Architectures: Choosing Between a Modern Data Warehouse, Data Fabric, Data Lakehouse, and Data Mesh

Tłumaczenie: Piotr Pilch

ISBN: 978-83-289-1669-2

© 2024 Helion S.A.

Authorized Polish translation of the English edition of Deciphering Data Architectures ISBN 9781098150761 © 2024 James Serra.

This translation is published and sold by permission of O'Reilly Media, Inc., which owns or controls all rights to publish and sell the same.

All rights reserved. No part of this book may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording or by any information storage retrieval system, without permission from the Publisher.

Wszelkie prawa zastrzeżone. Nieautoryzowane rozpowszechnianie całości lub fragmentu niniejszej publikacji w jakiegokolwiek postaci jest zabronione. Wykonywanie kopii metodą kserograficzną, fotograficzną, a także kopiowanie książki na nośniku filmowym, magnetycznym lub innym powoduje naruszenie praw autorskich niniejszej publikacji.

Wszystkie znaki występujące w tekście są zastrzeżonymi znakami firmowymi bądź towarowymi ich właścicieli.

Autor oraz wydawca dołożyli wszelkich starań, by zawarte w tej książce informacje były kompletne i rzetelne. Nie biorą jednak żadnej odpowiedzialności ani za ich wykorzystanie, ani za związane z tym ewentualne naruszenie praw patentowych lub autorskich. Autor oraz wydawca nie ponoszą również żadnej odpowiedzialności za ewentualne szkody wynikłe z wykorzystania informacji zawartych w książce.

Drogi Czytelniku!

Jeżeli chcesz ocenić tę książkę, zajrzyj pod adres

<https://helion.pl/user/opinie/noarda>

Możesz tam wpisać swoje uwagi, spostrzeżenia, recenzję.

Helion S.A.

ul. Kościuszki 1c, 44-100 Gliwice

tel. 32 230 98 63

e-mail: helion@helion.pl

WWW: <https://helion.pl> (księgarnia internetowa, katalog książek)

Printed in Poland.

- [Kup książkę](#)
- [Poleć książkę](#)
- [Oceń książkę](#)

- [Księgarnia internetowa](#)
- [Lubię to! » Nasza społeczność](#)

Spis treści

Słowo wstępne	17
Przedmowa	19

Część I. Fundamenty **23**

1. Technologia Big Data	25
Czym jest technologia Big Data i jak może być pomocna?	26
Dojrzałość danych	29
Etap 1: reakcyjny	30
Etap 2: informacyjny	31
Etap 3: predykcyjny	31
Etap 4: transformatywny	31
Samoobsługowa analityka biznesowa	32
Podsumowanie	33
2. Typy architektur danych	35
Ewolucja architektur danych	36
Relacyjna hurtownia danych	38
Jezioro danych	40
Nowoczesna hurtownia danych	42
Architektura Data Fabric	43
Architektura Data Lakehouse	43
Siatka danych	44
Podsumowanie	45
3. Sesja projektowania architektury	46
Czym jest sesja projektowania architektury?	46
Dlaczego należy przeprowadzać sesję ADS?	46
Przed sesją ADS	47

Przygotowanie	47
Zaproszenie uczestników	50
Przeprowadzanie sesji ADS	51
Wprowadzenia	52
Ustalenia	52
Korzystanie z tablicy	57
Po zakończeniu sesji ADS	58
Wskazówki dotyczące prowadzenia sesji ADS	59
Podsumowanie	61

Część II. Typowe pojęcia związane z architekturami danych **63**

4. Relacyjna hurtownia danych	65
Czym jest relacyjna hurtownia danych?	65
Czym hurtownia danych nie jest?	67
Podejście odgórne	69
Dlaczego warto skorzystać z relacyjnej hurtowni danych?	71
Wady związane z korzystaniem z relacyjnej hurtowni danych	74
Zasilanie hurtowni danych	75
Częstotliwość wyodrębniania danych	75
Metody wyodrębniania	75
Metoda określania, jakie dane zmieniły się od ostatniej operacji wyodrębnienia	76
Informacje o kresie relacyjnej hurtowni danych okazały się mocno przesadzone	77
Podsumowanie	78
5. Jezioro danych	80
Czym jest jezioro danych?	80
Dlaczego warto używać jeziora danych?	81
Podejście oddolne	83
Najlepsze praktyki projektowania architektury jeziora danych	84
Wiele jezior danych	90
Zalety	90
Wady	93
Podsumowanie	94
6. Procesy i rozwiązania z zakresu magazynowania danych	95
Rozwiązania do przechowywania danych	96
Składnice danych	96
Magazyny danych operacyjnych	97
Centra danych	99
Procesy danych	101
Zarządzanie danymi głównymi	101
Wirtualizacja i federacja danych	102

Katalogi danych	107
Platformy danych	108
Podsumowanie	109
7. Metody projektowe	110
Porównanie systemów OLTP i OLAP	111
Dane operacyjne i analityczne	113
Przetwarzanie SMP i MPP	113
Architektura Lambda	114
Architektura Kappa	117
Trwałość poliglotyczna i poliglotyczne magazyny danych	118
Podsumowanie	120
8. Metody modelowania danych	121
Modelowanie relacyjne	121
Klucze	121
Diagramy relacji encji	122
Reguły i formy normalizacji	122
Śledzenie zmian	123
Modelowanie wymiarowe	124
Fakty, wymiary i klucze	125
Monitorowanie zmian	125
Denormalizacja	126
Wspólny model danych	128
Model Data Vault	129
Metodologie Kimballa i Inmona dotyczące hurtowni danych	131
Metodologia odgórna Inmona	132
Metodologia oddolna Kimballa	133
Wybór metodologii	134
Modele hybrydowe	135
Mity dotyczące metodologii	137
Podsumowanie	140
9. Metody pozyskiwania danych	141
Porównanie procesów ETL i ELT	141
Odwrócony proces ETL	143
Porównanie przetwarzania wsadowego oraz przetwarzania w czasie rzeczywistym	145
Zalety i wady przetwarzania wsadowego	146
Zalety i wady przetwarzania w czasie rzeczywistym	146
Nadzór nad danymi	147
Podsumowanie	147

10. Nowoczesna hurtownia danych	151
Architektura nowoczesnej hurtowni danych	151
Zalety i wady architektury nowoczesnej hurtowni danych	156
Łączenie relacyjnej hurtowni danych z jeziorem danych	158
Jezioro danych	158
Relacyjna hurtownia danych	158
Kamienie milowe prowadzące do hurtowni MDW	159
Rozbudowa korporacyjnej hurtowni danych	159
Tymczasowe jezioro danych oraz korporacyjna hurtownia danych	161
Rozwiązanie kompleksowe	162
Studium przypadku: strategiczne przejście firmy Wilson & Gunkerk do nowoczesnej hurtowni danych	163
Wyzwanie	163
Rozwiązanie	163
Rezultat	164
Podsumowanie	164
11. Architektura Data Fabric	166
Architektura Data Fabric	167
Zasady dostępu do danych	167
Katalog metadanych	168
Zarządzanie danymi głównymi	169
Wirtualizacja danych	169
Przetwarzanie w czasie rzeczywistym	169
Interfejsy API	170
Usługi	170
Produkty	170
Dlaczego warto dokonać przejścia z hurtowni MDW na architekturę Data Fabric?	170
Potencjalne wady	171
Podsumowanie	171
12. Architektura Data Lakehouse	173
Opcje warstwy Delta Lake	174
Poprawa wydajności	176
Architektura Data Lakehouse	177
Co się stanie, gdy zrezygnujesz z relacyjnej hurtowni danych?	179
Relacyjna warstwa udostępniająca	181
Podsumowanie	182

13. Fundamenty siatki danych	183
Zdecentralizowana architektura danych	184
Szum wokół siatki danych	185
Cztery zasady Dehghani dotyczące siatki danych	186
Pierwsza zasada: własność domeny	186
Druga zasada: dane jako produkt	187
Trzecia zasada: samoobsługowa infrastruktura danych jako platforma	189
Czwarta zasada: nadzór nad federacyjnymi zasobami obliczeniowymi	190
„Czysta” siatka danych	191
Domeny danych	192
Logiczna architektura siatki danych	193
Różne topologie	195
Porównanie siatki danych i architektury Data Fabric	197
Warianty zastosowania	197
Podsumowanie	199
14. Czy powinno się adaptować siatkę danych? Mity, obawy i przyszłość	200
Mity	200
Mit: użycie siatki danych to cudowny środek pozwalający szybko poradzić sobie z wszystkimi trudnościami towarzyszącymi danym	200
Mit: siatka danych zastąpi Twoje jezioro danych i hurtownię danych	201
Mit: wszystkie projekty z hurtownią danych nie udają się, a siatka danych rozwiąże ten problem	201
Mit: budowanie siatki danych oznacza decentralizację absolutnie wszystkiego	201
Mit: możesz użyć wirtualizacji danych, aby utworzyć siatkę danych	201
Obawy	202
Kwestie filozoficzne i koncepcyjne	203
Łączenie danych w środowisku zdecentralizowanym	204
Inne kwestie związane z decentralizacją	205
Złożoność	206
Duplikacja	206
Wykonalność	207
Ludzie	209
Bariery na poziomie domen	210
Ocena organizacyjna: czy powinno się adaptować siatkę danych?	211
Zalecenia dotyczące implementowania z powodzeniem siatki danych	212
Przyszłość siatki danych	214
Szersze spojrzenie: zrozumienie architektur danych i ich zastosowań	215
Podsumowanie	216

15. Ludzie i procesy	219
Organizacja zespołów: role i obowiązki	220
Role w przypadku nowoczesnej hurtowni danych oraz architektury Data Fabric lub Data Lakehouse	220
Role w przypadku siatki danych	222
Dlaczego projekty się nie udają: pułapki i zapobieganie im	225
Pułapka: pozwalanie szefostwu myśleć, że analityka biznesowa jest „łatwa”	225
Pułapka: używanie niewłaściwych technologii	225
Pułapka: zgromadzenie zbyt wielu wymagań biznesowych	225
Pułapka: zgromadzenie zbyt małej liczby wymagań biznesowych	226
Pułapka: prezentowanie raportów bez wcześniejszego sprawdzenia poprawności ich zawartości	226
Pułapka: zatrudnianie niedoświadczonej firmy konsultingowej	226
Pułapka: zatrudnianie firmy konsultingowej, która zleca prace projektowe pracownikom z innych krajów	227
Pułapka: przekazywanie konsultantom własności projektu	227
Pułapka: zlekceważenie konieczności transferu zasobów wiedzy z powrotem do organizacji	227
Pułapka: zmniejszanie budżetu w połowie trwania projektu	227
Pułapka: rozpoczynanie od daty końcowej i cofanie się z działaniami	228
Pułapka: określanie struktury hurtowni danych zgodnie z danymi źródłowymi, a nie wymogami firmy	228
Pułapka: zaprezentowanie użytkownikom rozwiązania z długim czasem odpowiedzi lub innymi problemami z wydajnością	228
Pułapka: przesadzenie z projektem architektury danych lub niedopracowanie go	229
Pułapka: kiepska komunikacja między działem informatycznym i domenami biznesowymi	229
Wskazówki dotyczące sukcesu	229
Nie oszczędzaj na inwestycjach	230
Angażuj użytkowników, prezentuj im wyniki i wzbudzaj ich entuzjazm	230
Zapewnij wartość w nowych raportach i panelach kontrolnych	231
Poproś użytkowników o zbudowanie prototypu	231
Znajdź orędownika/sponsora projektu	232
Stwórz plan projektu z celem 80% efektywności	232
Podsumowanie	232

16. Technologie	234
Wybór platformy	234
Rozwiązania <i>open source</i>	234
Rozwiązania lokalne	237
Rozwiązania dostawców usług w chmurze	238
Modele usług w chmurze	240
Główni dostawcy usług w chmurze	243
Rozwiązania z wieloma chmurami	243
Środowiska oprogramowania	246
Hadoop	246
Databricks	249
Snowflake	251
Podsumowanie	252

Typy architektur danych

Bezwzględnie ważne jest poświęcenie na samym początku czasu na zaprojektowanie i zbudowanie odpowiedniej architektury danych. Przekonałem się o tym na własnej skórze na początku mojej kariery. Byłem tak bardzo podekscytowany faktem rozpoczęcia procesu tworzenia mojego rozwiązania, że beztrąsko podszedłem do kwestii istotnych decyzji dotyczących projektu architektury i tego, jakie produkty mają zostać zastosowane. Po trzech miesiącach realizowania projektu uświadomiłem sobie to, że architektura nie będzie obsługiwać niektórych z wymaganych źródeł danych. W zasadzie konieczne było rozpoczęcie projektu od nowa i przygotowanie kolejnej architektury oraz innych produktów, co oznaczało mnóstwo straconego czasu i pieniędzy. Bez właściwego planowania Twoje rozwiązanie nie zapewni użytkownikom wartości. Poza tym będą oni niezadowoleni z powodu niedotrzymanych terminów, a Twoja firma w coraz większym stopniu ryzykuje pozostawaniem w tyle względem swojej konkurencji.

W trakcie budowania rozwiązania do obsługi danych musisz postępować zgodnie z dobrze przemyślanym planem. Właśnie tutaj do gry wkracza **architektura danych**. Definiuje ona ogólną metodę architektoniczną i pojęcia, które są w niej stosowane, a ponadto opisuje zestaw używanych technologii oraz określa przepływ danych, jaki posłuży do zbudowania rozwiązania do przechwytywania danych o pokaźnej wielkości. Decydowanie w kwestii architektury danych może być dużym wyzwaniem, gdyż nie istnieje żadna uniwersalna architektura. Nie masz możliwości przejrzzenia książki w celu znalezienia wykazu metod architektonicznych powiązanych z odpowiednimi produktami do użycia. Nie istnieje prosty diagram możliwy do wykorzystania w przypadku drzew decyzyjnych, który doprowadzi Cię do idealnej architektury. Twoja metoda architektoniczna oraz używane technologie będą się w znacznym stopniu różnić dla poszczególnych klientów oraz wariantów zastosowania.

Podstawowe typy ogólnych metod architektonicznych oraz zagadnienia są dokładnie tym, co przedstawiłem w tej książce. Choć nie ma w niej żadnego wykazu, będziesz mieć możliwość zorientowania się, o co w tym wszystkim chodzi. Wprawdzie przydatne jest podzielenie architektur danych na typy zależnie od ich właściwości (tak też postąpiłem w książce), jednak nie jest to tożsame z dokonywaniem wyboru spośród zestawu predefiniowanych i uniwersalnych szablonów. Każda architektura danych jest unikalna, a ponadto wymaga spersonalizowanego podejścia w celu spełnienia konkretnych wymagań biznesowych.

Architektura danych odnosi się do ogólnego projektu i organizacji danych w obrębie systemu informatycznego. Predefiniowane szablony architektur danych mogą sprawiać wrażenie prostej metody pozwalającej szybko przygotować nowy system. Tego rodzaju szablony często jednak

nie uwzględniają konkretnych wymagań i ograniczeń systemu, względem którego są stosowane. Może to prowadzić do problemów z jakością danych, wydajnością systemu oraz serwisowaniem. Ponadto wymagania organizacji i systemu danych będą prawdopodobnie zmieniać się z upływem czasu, co będzie wymagać aktualizowania architektury danych i dokonywania w niej korekt. Standaryzowany szablon może nie być wystarczająco elastyczny do tego, aby uwzględniać takie zmiany, a to może powodować pojawianie się w systemie nieefektywności i ograniczeń.

W niniejszym rozdziale znajdziesz skrócony przegląd podstawowych typów, które omówię w książce. Jest to relacyjna hurtownia danych, jezioro danych, nowoczesna hurtownia danych, Data Fabric, Data Lakehouse oraz siatka danych. W dalszej części książki każdemu z nich poświęciłem osobny rozdział zawierający mnóstwo szczegółów.

Ewolucja architektur danych

W **relacyjnej bazie danych** dane są przechowywane w sposób ustrukturyzowany. Za pomocą kluczy między elementami danych definiuje się relacje. Dane są zwykle porządkowane w obrębie tabel, z których każda składa się z wierszy i kolumn. Każdy wiersz reprezentuje pojedynczą instancję danych, natomiast poszczególne kolumny reprezentują konkretne atrybuty danych.

Relacyjne bazy danych zaprojektowano do obsługi danych ustrukturyzowanych. Bazy te zapewniają strukturę służącą do tworzenia, modyfikowania i uzyskiwania danych za pomocą standaryzowanego języka SQL (*Structured Query Language*). Model relacyjny po raz pierwszy został zaproponowany w 1970 r. przez Edgara F. Codda¹. Od połowy lat 70. model ten stał się dominujący w systemach zarządzania bazami danych. Większość aplikacji operacyjnych musi trwale przechowywać dane, a relacyjna baza danych to narzędzie wybierane w przypadku ich zdecydowanej większości.

W relacyjnych bazach danych, gdzie podstawowe znaczenie odgrywa spójność i integralność danych, są one zwykle porządkowane z użyciem metody *Schema-on-Write*. Schemat odnosi się do formalnej struktury definiującej organizację tabel, pól, typów danych i ograniczeń oraz relacji między nimi. Schemat pełni rolę planu w zakresie przechowywania danych oraz zarządzania nimi, a ponadto zapewnia spójność, integralność i efektywną organizację wewnątrz bazy danych. Relacyjne bazy danych (oraz relacyjne hurtownie danych) wymagają niewielkiej ilości wstępnych działań, zanim będą mogły zostać umieszczone w nich dane. Musisz utworzyć bazę danych oraz jej tabele, pola i schemat, a następnie kod służący do przeniesienia danych do bazy. W przypadku metody *Schema-on-Write* schemat danych jest definiowany i wymuszany w momencie zapisywania danych w bazie lub pozyskiwania ich. Zanim dane będą mogły zostać zapisane, muszą być zgodne z predefiniowanym schematem obejmującym typy danych, ograniczenia i relacje.

Dla porównania w przypadku metody *Schema-on-Read* schemat stosuje się w chwili odczytu danych lub uzyskiwania do nich dostępu, a nie wtedy, gdy są one zapisywane. Dane mogą być pozyskiwane do systemu magazynowania bez utrzymywania zgodności z konkretnym schematem, a struktura jest

¹ Dokument *Derivability, Redundancy, and Consistency of Relations Stored in Large Data Banks*, „Communications of the ACM”, 13, nr 6 (1970): str. 377 – 387. Dokument autorstwa E. F. Codda jest dostępny na stronie internetowej *A Relational Model of Data for Large Shared Data Banks*.

definiowana tylko wtedy, gdy dane są używane lub są celem zapytania. Takie podejście oferuje większą elastyczność w przypadku danych bez struktury lub mających częściową strukturę, a ponadto jest powszechnie stosowane w jeziorach danych, które omówiłem w dalszej części rozdziału.

Na ogólnym poziomie architektury danych zapewniają strukturę służącą do porządkowania danych i zarządzania nimi w sposób spełniający wymagania organizacji. Obejmuje to definiowanie tego, jak dane są gromadzone, przechowywane, przetwarzane i używane, a także utrzymywanie jakości danych, bezpieczeństwa i prywatności. Choć architektury danych mogą przyjmować wiele różnych form, wśród ich wspólnych elementów należy wymienić następujące:

Magazyn danych

Wszystkie architektury danych muszą określać metodę przechowywania danych z uwzględnieniem fizycznego nośnika danych (np. dyski twarde lub magazyn w chmurze) oraz struktur danych używanych do ich organizacji.

Przetwarzanie danych

W przypadku architektur danych niezbędne jest podanie, w jaki sposób dane są przetwarzane, co obejmuje wszelkie transformacje lub obliczenia wykonywane względem danych przed ich zapisaniem lub poddaniu analizie.

Dostęp do danych

Obowiązkową funkcją architektury danych jest zapewnienie mechanizmów uzyskiwania dostępu do danych z uwzględnieniem interfejsów użytkownika oraz interfejsów API (*Application Program Interface*), które umożliwiają analizowanie danych lub tworzenie dotyczących ich zapytań.

Bezpieczeństwo i prywatność danych

Wymagane jest, aby architektury danych zawierały mechanizmy zapewniające bezpieczeństwo i prywatność danych, takie jak kontrola dostępu, szyfrowanie i maskowanie danych.

Nadzór nad danymi

Kolejnym elementem, jaki muszą oferować architektury danych, są struktury przeznaczone do zarządzania danymi, w tym standardy jakości, monitorowanie pochodzenia oraz zasady dotyczące retencji.

Generalnie głównym celem architektury danych jest umożliwienie organizacji efektywnego zarządzania swoimi zasobami danych oraz korzystania z nich, tak aby mogła ona realizować własne cele biznesowe oraz procesy podejmowania decyzji.

Tabela 2.1, w której znajduje się ogólne porównanie właściwości architektur danych omawianych w książce, powinna być na początek odpowiednia do określenia, jaka architektura może być najbardziej odpowiednia w konkretnym wariantcie zastosowania.

Tabela 2.1. Porównanie architektur danych

Właściwość	Relacyjna hurtownia danych	Jezioro danych	Nowoczesna na hurtownia danych	Data Fabric	Data Lakehouse	Siatka danych
Rok wprowadzenia	1984	2010	2011	2016	2020	2019
Scentralizowana/zdecentralizowana	Scentralizowana	Scentralizowana	Scentralizowana	Scentralizowana	Scentralizowana	Zdecentralizowana
Typ magazynu	Relacyjny	Obiektowy	Relacyjny i obiektowy	Relacyjny i obiektowy	Obiektowy	Domenowy
Typ schematu	Schema-on-Write	Schema-on-Read	Schema-on-Write i Schema-on-Read	Schema-on-Write i Schema-on-Read	Schema-on-Read	Domenowy
Bezpieczeństwo danych	Duże	Od małego do średniego	Od średniego do dużego	Duże	Średnie	Powiązane z domeną
Opóźnienie dostępu do danych	Małe	Duże	Od małego do dużego	Od małego do dużego	Od średniego do dużego	Zależne od domeny
Czas uzyskania efektu	Średni	Krótki	Krótki	Krótki	Krótki	Długi
Całkowity koszt rozwiązania	Duży	Mały	Średni	Od średniego do dużego	Od małego do średniego	Duży
Wspierane warianty zastosowania	Mała ilość	Od małej do średniej ilości	Średnia ilość	Od średniej do dużej ilości	Duża ilość	Duża ilość
Trudność projektowa	Mała	Średnia	Średnia	Średnia	Od średniej do dużej	Duża
Dojrzałość technologii	Duża	Średnia	Od średniej do dużej	Od średniej do dużej	Od średniej do dużej	Mała
Niezbędny zestaw umiejętności pracowników firmy	Mały	Od małego do średniego	Średni	Od średniego do dużego	Od średniego do dużego	Duży

Relacyjna hurtownia danych

Przez kilka dekad relacyjne bazy danych były ostoją magazynowania danych. System Teradata zaprojektowany na uniwersytecie Stanford przez dr. Jacka E. Shemera (który w 1979 r. założył firmę Teradata Corporation) był pierwszą relacyjną hurtownią danych używaną w środowiskach produkcyjnych. W 1983 r. bank Wells Fargo *zainstalował pierwszy system Teradata* i używał go do analizowania danych finansowych.

Gdy firmy zaczęły generować coraz większe, pokaźne ilości danych, coraz trudniejsze stawało się ich przetwarzanie i analizowanie bez znacznego opóźnienia. Ograniczenia relacyjnych baz danych

przyczyniły się do opracowania relacyjnych hurtowni danych², które znacznie zyskały na popularności pod koniec lat 80.³, czyli około 15 lat po tym, jak pojawiły się relacyjne bazy danych. **Relacyjna hurtownia danych** (ang. *relational data warehouse*) to konkretny typ relacyjnej bazy danych zaprojektowany pod kątem przechowywania danych w hurtowni oraz zastosowań z zakresu analityki biznesowej, w przypadku którego zoptymalizowano wydajność zapytań i obsługę analizy danych o dużej skali. Choć relacyjne hurtownie danych i przetwarzanie transakcyjne w celu zapewnienia organizacji danych bazują na modelu relacyjnym, to relacyjna hurtownia danych ma zwykle większą skalę, a ponadto jest zoptymalizowana pod kątem zapytań analitycznych.

Relacyjne hurtownie danych są wyposażone zarówno w mechanizm realizujący obliczenia, jak i magazyn. Mechanizm ten zapewnia moc obliczeniową służącą do wykonywania zapytań dotyczących danych. Magazyn ma postać *magazynu relacyjnego* przechowującego dane ze strukturą uzyskiwaną dzięki zastosowaniu tabel, wierszy i kolumn. Moc obliczeniowa relacyjnej hurtowni danych może być używana wyłącznie względem jej magazynu relacyjnego (są one ze sobą sprzężone).

Wśród najważniejszych opcji relacyjnych hurtowni danych należy wymienić obsługę transakcji (zapewnienie, że dane są przetwarzane w niezawodny i spójny sposób), ścieżki audytu (utrzymywanie zapisu całej aktywności związanej z danymi w systemie) oraz wymuszanie schematu (zadbanie o to, aby dane zostały uporządkowane i uzyskały strukturę w predefiniowany sposób).

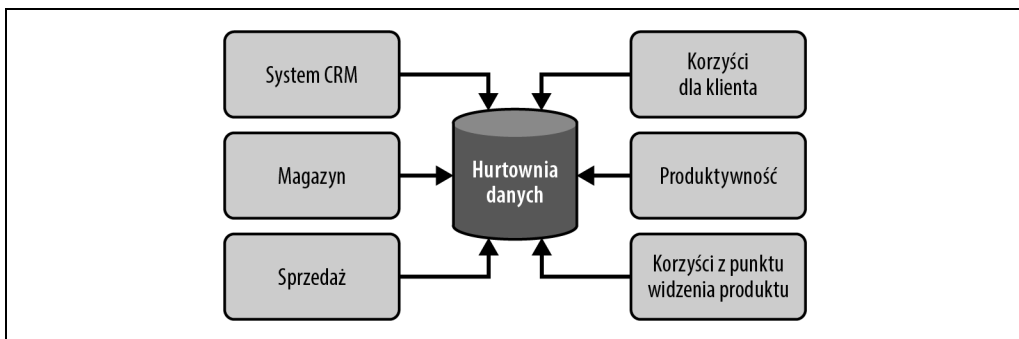
W latach 70. i 80. firmy używały relacyjnych baz danych na potrzeby aplikacji operacyjnych, takich jak obsługujące wprowadzanie zamówień i zarządzające stanem magazynowym. W odniesieniu do tego rodzaju aplikacji używa się terminu systemów OLTP (*Online Transaction Processing*). Systemy te mogą wprowadzać, odczytywać, aktualizować i usuwać zmiany (są to *operacje CRUD* — *Create, Read, Update, Delete*) dotyczące danych w bazie danych. Operacje te tworzą fundament modyfikowania danych oraz zarządzania w architekturach danych. Mają one kluczowe znaczenie z punktu widzenia projektowania i implementowania systemów do magazynowania danych oraz interfejsów użytkownika, które prowadzą interakcję z danymi.

Operacje CRUD wymagają krótkich czasów odpowiedzi aplikacji, tak aby użytkownicy nie irytowali się tym, ile czasu zajmuje aktualizowanie danych. Możliwe jest uruchamianie zapytań i generowanie raportów z wykorzystaniem relacyjnej bazy danych używanej przez aplikację operacyjną, ale wiąże się to z zużyciem wielu zasobów, a ponadto może wywoływać konflikt z innymi operacjami CRUD wykonywanymi w tym samym czasie. W efekcie wszystko może zostać spowolnione.

Po części relacyjne hurtownie danych zostały opracowane w celu rozwiązania tego problemu. Dane z relacyjnej bazy danych są kopiowane do hurtowni danych, a użytkownicy mogą uruchamiać zapytania i raporty względem hurtowni danych, a nie relacyjnej bazy danych. Dzięki temu zapytania i raporty nie obciążają systemu udostępniającego relacyjną bazę danych i nie spowalniają aplikacji, z której korzystają użytkownicy. Relacyjne hurtownie danych centralizują również dane z wielu aplikacji, aby usprawnić raportowanie (zilustrowano to na rysunku 2.1).

² Uwaga dotycząca języka: możesz spotkać osoby, które w odniesieniu do architektury relacyjnej hurtowni danych używają terminu *tradycyjna hurtownia danych*. W książce przeważnie posługuję się określeniem *relacyjna hurtownia danych*, który czasami skracam do postaci *hurtownia danych*. Wszystkie te terminy odwołują się do tej samej rzeczy.

³ Popularność relacyjnej hurtowni danych zwiększyła się w dużej mierze dzięki Barry'emu Devlinowi i Billowi Inmonowi, o czym będzie mowa w rozdziale 8.



Rysunek 2.1. Zastosowanie hurtowni danych

Bardziej szczegółowo relacyjne hurtownie danych omówię w rozdziale 4.

Jezioro danych

Jezioro danych (ang. *data lake*) to nowsze pojęcie, które po raz pierwszy pojawiło się około roku 2010. Jezioro danych możesz traktować jako wychwalany system plików, który nie różni się zbytnio od systemu plików obecnego na Twoim laptopie. **Jezioro danych** to po prostu magazyn, z którym — w przeciwieństwie do relacyjnej hurtowni danych — nie jest skojarzony żaden mechanizm obliczeniowy. Na szczęście istnieje wiele takich mechanizmów, które mogą współpracować z jeziorami danych, dlatego w jego przypadku moc obliczeniowa jest zwykle tańsza niż przy korzystaniu z relacyjnej hurtowni danych. Kolejną różnicą jest to, że relacyjne hurtownie danych używają magazynu relacyjnego, natomiast jeziora danych stosują *magazyn obiektowy*, który nie wymaga tworzenia struktury danych w ramach wierszy i kolumn.

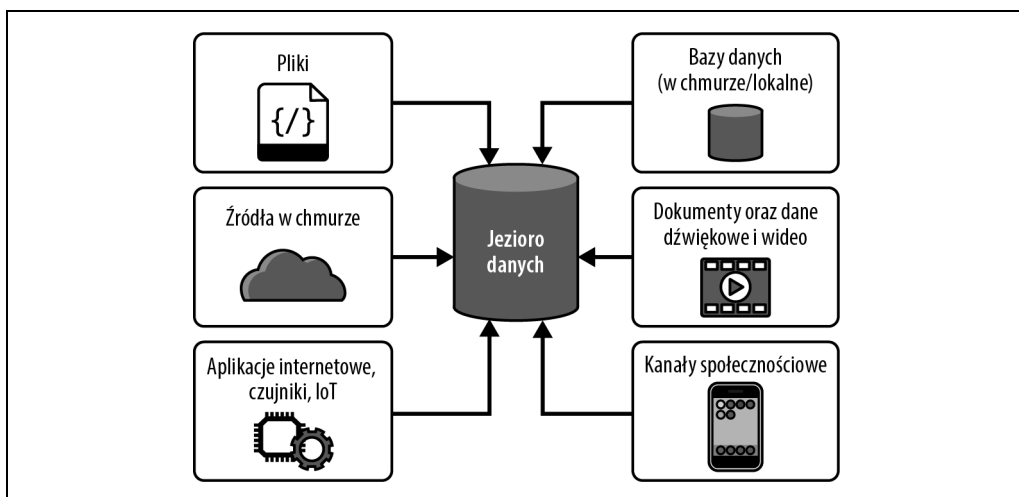
Technologia magazynów jeziora danych została zapoczątkowana systemem Apache HDFS (*Hadoop Distributed File System*), czyli darmową technologią *open source*⁴ udostępnianą prawie wyłącznie lokalnie, która cieszyła się bardzo dużą popularnością w pierwszych latach po roku 2010. HDFS to skalowalny i odporny na błędy rozproszony system magazynowania, który zaprojektowano pod kątem działania na powszechnie dostępnym sprzęcie. System plików stanowi główny komponent ekosystemu Apache Hadoop, który szerzej omawiam w rozdziale 16. Wraz z dalszym zwiększaniem się znaczenia obliczeń w chmurze w jej obrębie były tworzone jeziora danych z użyciem różnego typu magazynu, a obecnie większość jezior istnieje w technologii chmury.

W porównaniu z relacyjną hurtownią danych jezioro danych bazuje na metodzie *Schema-on-Read*. Oznacza to, że do umieszczenia danych w jeziorze nie są wymagane żadne wcześniejsze działania: operacja skopiowania plików do jeziora danych może być tak prosta jak w przypadku folderów obecnych na Twoim laptopie. Dane w jeziorze są przechowywane w ich naturalnym (lub nieprzetworzonym) formacie. Oznacza to, że dane mogą trafić z ich systemu źródłowego do jeziora danych bez przekształcania ich do innego formatu. Jeśli na przykład eksportujesz dane z relacyjnej bazy

⁴ Termin *open source* odnosi się do oprogramowania, którego kod źródłowy udostępniono za darmo do publicznego wykorzystania, modyfikowania i dystrybuowania.

danych do pliku w nieprzetworzonym formacie CSV, dane mógłbyś przechowywać w jeziorze danych w niezmienionej postaci. Jeżeli jednak zdecydowałbyś się na zapisanie danych w relacyjnej hurtowni danych, konieczne byłoby przekształcenie ich tak, aby mogły zostać umiejscowione w wierszach i kolumnach tabeli.

Podczas kopiowania do jeziora danych pliku z danymi jego schemat może nie zostać uwzględniony w ramach tej operacji albo może znajdować się w innym pliku. W związku z tym musisz zdefiniować schemat przez utworzenie go lub uzyskanie z osobnego pliku (stąd też termin *Schema-on-Read*). Jak widać na rysunku 2.2, dane z systemów źródłowych — takich jak bazy danych aplikacji operacyjnych, dane z czujników oraz dane z mediów społecznościowych — w całości mogą trafić do jeziora danych. W takich plikach mogą być przechowywane dane ze strukturą (np. dane z relacyjnych baz danych), z częściową strukturą (np. pliki CSV, dzienniki oraz pliki XML lub JSON) lub bez struktury (np. dane z wiadomości pocztowych, dokumentów i plików PDF). W plikach mogą nawet znajdować się dane binarne (np. obrazy, dane dźwiękowe i wideo).



Rysunek 2.2. Jezioro danych

Początkowo jeziora danych traktowano jako rozwiązanie wszystkich problemów towarzyszących relacyjnym hurtowniom danych, w tym: dużego kosztu, ograniczonej skalowalności, kiepskiej wydajności, dodatkowego nakładu pracy związanego z przygotowywaniem danych oraz ograniczonego wsparcia złożonych typów danych. Firmy zajmujące się sprzedażą systemu Hadoop i architektury jezior danych, takie jak Cloudera, Hortonworks i MapR, robiły wokół nich taki szum, jakby były one wypełnione jednorozcami i tęczmami, które będą kopiować i oczyszczać dane, a następnie udostępniać je użytkownikom z magiczną łatwością. Firmy twierdziły, że jeziora danych mogłyby całkowicie zastąpić relacyjne hurtownie danych w ramach rozwiązania typu „jedna technologia służąca do wszystkiego”. Więcej niż kilka firm zdecydowało się na oszczędności dzięki zastosowaniu darmowych narzędzi *open source* na potrzeby całości swojej technologii.

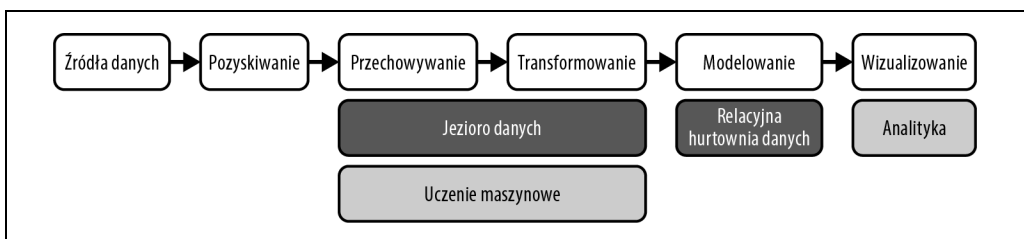
Problem polegał na tym, że w rzeczywistości kierowanie zapytań do jezior danych nie jest takie proste: wymaga to pewnych dość zaawansowanych umiejętności. Informatycy mogliby przekazać użytkownikom następującą informację: „Skopiowaliśmy do tego jeziora danych wszystkie dane,

których potrzebujecie. Po prostu uruchomcie notatnik Jupyter i użyjcie systemu Hive oraz języka Python do utworzenia swoich raportów z wykorzystaniem plików znajdujących się w tych folderach”. Skończyło się to żałością, gdyż większość użytkowników w żadnym razie nie dysponowała umiejętnościami niezbędnymi do zrealizowania wszystkich tych działań. Firmy przekonały się w do-
tkliwy sposób, że te złożone i trudne w użyciu rozwiązania ostatecznie okazały się droższe z powodu kosztów sprzętu i obsługi, opóźnień produkcyjnych oraz utraconej produktywności. Ponadto jeziora danych były pozbawione opcji, jakie spodobały się użytkownikom w przypadku hurtowni danych, takich jak obsługa transakcji, wymuszanie schematu oraz ścieżki audytu. Poskutkowało to tym, że swoją działalność biznesową zakończyło dwóch z trzech głównych dostawców architektury jeziora danych, czyli firmy Hortonworks i MapR.

Jednak architektura jeziora danych nie zniknęła. Zamiast tego jej pierwotne przeznaczenie przekształcono do postaci innego, lecz bardzo przydatnego, czyli organizowania i przygotowywania danych. Jeziora danych omówię szczegółowo w rozdziale 5.

Nowoczesna hurtownia danych

Same w sobie relacyjne hurtownie danych i jeziora danych to bardzo proste architektury. Do centralizacji danych używają one tylko jednej technologii, co jest wspierane przez kilka dodatkowych produktów albo nawet żadne. Gdy zastosujesz więcej technologii i produktów do obsługi relacyjnych hurtowni danych lub jezior danych, przyjmą one postać architektur omawianych w tym i kolejnych rozdziałach. Jeziora danych nie zastąpiły relacyjnych hurtowni danych, ale nadal oferowały korzyści z punktu widzenia organizowania i przygotowywania danych. Dlaczego nie można by skorzystać z zalet obu rozwiązań? W okolicach roku 2011 wiele firm zaczęło tworzyć architektury, które umiejscawiają jeziora danych obok relacyjnych hurtowni danych w celu zapewnienia architektury, która obecnie określana jest mianem **nowoczesnej hurtowni danych** (ang. *modern data warehouse*) (rysunek 2.3). Słowo *nowoczesna* użyte w tym terminie odnosi się do zastosowania w przypadku hurtowni danych nowszych technologii i metod.

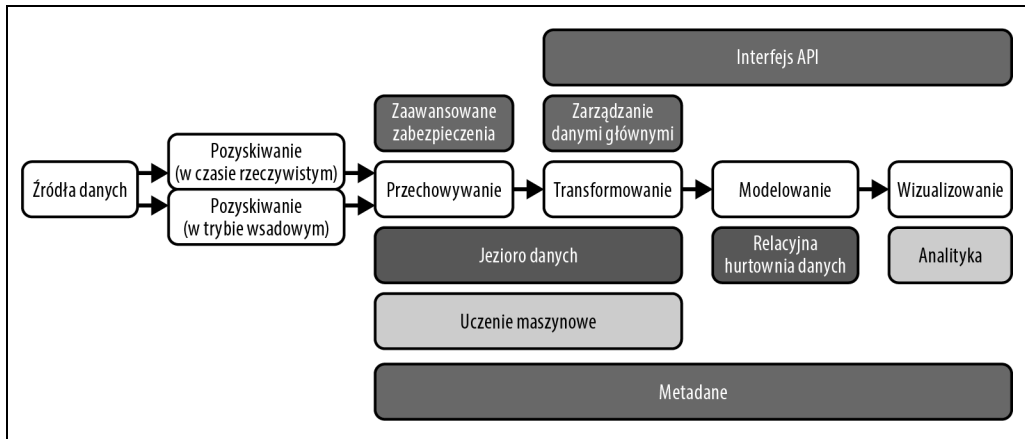


Rysunek 2.3. Nowoczesna hurtownia danych

Oto rozwiązanie reprezentujące „to, co najlepsze z obu światów”: jezioro danych służy do organizowania i przygotowywania danych, a danolodzy używają tej architektury do budowania modeli uczenia maszynowego. Z kolei hurtownia danych ma na celu zapewnienie obsługi, bezpieczeństwa i zgodności. Korzystając z niej, użytkownicy tworzą swoje zapytania i generują raporty. W rozdziale 10. znacznie dokładniej objaśniam architektury nowoczesnej hurtowni danych.

Architektura Data Fabric

Architektura **Data Fabric** zaczęła się pojawiać około roku 2016. Możesz ją traktować jako kolejny etap rozwoju architektury nowoczesnej hurtowni danych poszerzonej o dodatkową technologię służącą do pozyskiwania większej ilości danych, zabezpieczania ich oraz udostępniania. Ponadto wprowadzono ulepszenia dotyczące tego, jak system pozyskuje, przekształca, pobiera i wyszukuje dane oraz zapewnia do nich dostęp (rysunek 2.4). Dzięki tym wszystkim dodatkom system staje się „fabryką”, czyli dużą strukturą, która umożliwiła pozyskiwanie dowolnego rodzaju danych. Zajmiemy się tym bardziej szczegółowo w rozdziale 11.

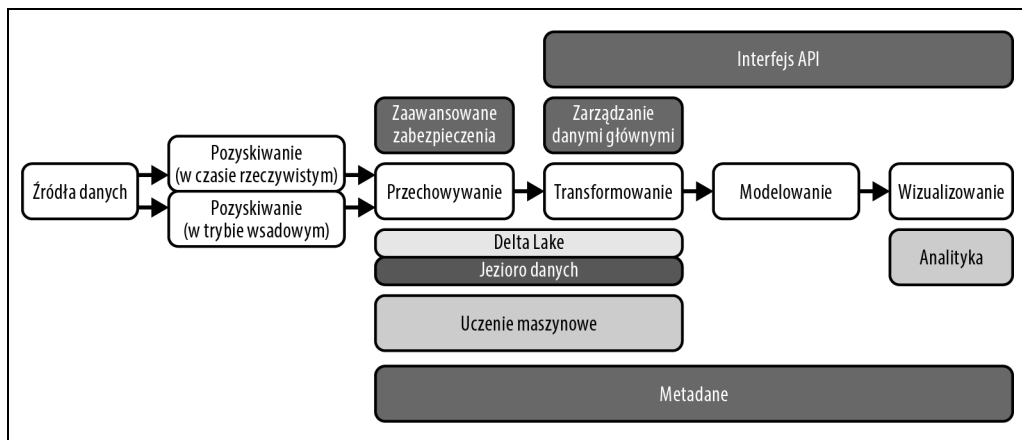


Rysunek 2.4. Architektura Data Fabric

Architektura Data Lakehouse

Termin **Data Lakehouse** powstał w wyniku połączenia terminów *jezioro danych* (ang. *data lake*) i *hurtownia danych* (ang. *data warehouse*). Architektury Data Lakehouse stały się popularne około roku 2020, gdy terminem tym zaczęła się posługiwać firma Databricks. W przypadku zagadnienia reprezentowanego przez termin Lakehouse chodzi o wyeliminowanie relacyjnej hurtowni danych i zastosowanie w posiadanej architekturze danych tylko jednego repozytorium w postaci jeziora danych. Dane każdego typu (ze strukturą, z częściową strukturą i bez struktury) są pozyskiwane do jeziora danych, a wszystkie zapytania i raporty są tworzone z tego poziomu.

Wiem, że możesz stwierdzić następująco: „Jedna chwila. Wspomniałeś, że w jeziorach danych zastosowano takie samo rozwiązanie, gdy się dopiero pojawiły, ale skończyło się to żałością! Co się zmieniło?”. Jak widać na rysunku 2.5, odpowiedzią na to pytanie jest warstwa oprogramowania magazynu transakcyjnego, która działa ponad istniejącym jeziorem danych i sprawia, że funkcjonuje ono bardziej na podobieństwo relacyjnej bazy danych. Wśród opcji *open source* będących konkurencją dla tej warstwy należy wymienić narzędzia Delta Lake, Apache Iceberg i Apache Hudi. Wszystkie omawiam dokładniej w rozdziale 12.



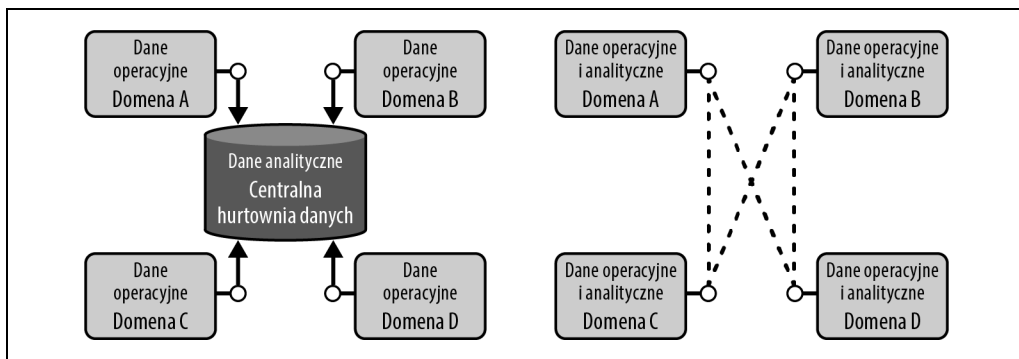
Rysunek 2.5. Architektura Data Lakehouse

Siatka danych

Termin **siatka danych** (ang. *data mesh*) po raz pierwszy został wprowadzony w maju 2019 r. przez Zhamak Dehghani (założycielka i prezes firmy Nextdata oraz autorka książki *Siatka danych. Nowoczesna koncepcja samoobsługowej infrastruktury danych*; Helion, 2023 r.) w jej wpisie blogowym. W grudniu 2020 r. Dehghani *dokładniej wyjaśniła*, czym jest siatka danych, a ponadto przedstawiła jej cztery fundamentalne zasady. Od tamtego czasu architektura siatki danych stała się wyjątkowo gorącym tematem. Dyskutowano o niej w ramach ogromnej liczby blogów, prezentacji, konferencji i materiałów prasowych. Temat został nawet ujęty w raporcie badawczym na stronie *Gartner Hype Cycle for data management*. W przypadku architektury siatki danych jest wiele rzeczy, które można polubić, ale — pomimo towarzyszącego jej rozgłosowi — jest ona odpowiednia tylko w odniesieniu do niewielkiej liczby wariantów zastosowań.

Nowoczesna hurtownia danych oraz architektury Data Fabric i Data Lakehouse obejmują swoim zakresem *centralizowanie* danych, czyli kopiowanie danych operacyjnych do centralnej lokalizacji będącej w posiadaniu działu informatycznego w ramach architektury kontrolowanej przez jego pracowników, w obrębie której tworzą oni dane analityczne (widoczne po lewej stronie na rysunku 2.6). Takie scentralizowane podejście stawia przed nami następujące trzy główne wyzwania: własność danych, jakość danych oraz skalowanie o charakterze organizacyjnym lub technicznym. Celem siatki danych jest stawienie czoła tym wyzwaniom.

W siatce danych dane są utrzymywane w obrębie kilku domen należących do firmy, takich jak produkcja, sprzedaż i dostawcy (widoczne po prawej stronie na rysunku 2.6). Z każdą domeną jest powiązany jej własny, niewielki zespół informatyków, którzy są właścicielami danych domeny i oczyszczają je, a ponadto tworzą i udostępniają dane analityczne. Każda domena dysponuje również własną infrastrukturą służącą do magazynowania i wykonywania obliczeń. Skutkuje to architekturą *zdecentralizowaną*, w przypadku której skalowane są dane, zasoby ludzkie oraz infrastruktura. Im większą liczbą domen dysponujesz, tym do większej liczby osób i zasobów infrastruktury uzyskujesz dostęp. System może obsługiwać więcej danych, a dział informatyczny nie stanowi już wąskiego gardła.



Rysunek 2.6. Siatka danych

Ważne jest uświadomienie sobie tego, że siatka danych to *pojęcie*, a nie technologia. Nie istnieje żadna „pudełkowa wersja siatki danych”, którą możesz kupić. Implementowanie siatki danych wiąże się z dokonaniem bardzo dużej zmiany organizacyjnej i kulturowej, na którą jest gotowa bardzo niewielka liczba firm (tak naprawdę większość firm nie jest nawet na tyle duża, aby rozważać zastosowanie architektury siatki danych: jest to rozwiązanie całkowicie przewidziane dla przedsiębiorstw). Zbudowanie siatki danych wymaga określenia, jakie elementy istniejącej technologii możesz przystosować na jej potrzeby, a jakie elementy będą musiały zostać utworzone. W przypadku każdej domeny konieczne jest ustalenie, jakich technologii użyje ona do zbudowania powiązanej z nią części siatki danych. Może to obejmować utworzenie nowoczesnej hurtowni danych albo architektury Data Fabric lub Data Lakehouse. Architekturom siatki danych można poświęcić sporo miejsca, na co też zdecydowałem się w rozdziałach 13. i 14.

Podsumowanie

Gdy już masz ogólną orientację w zakresie typów architektur danych, w kolejnym rozdziale będzie mowa o tym, jak określić najlepszą architekturę, która może zostać użyta: jest to proces określany mianem **sesji projektowania architektury**.

PROGRAM PARTNERSKI

— GRUPY HELION —



1. ZAREJESTRUJ SIĘ
2. PREZENTUJ KSIĄŻKI
3. ZBIERAJ PROWIZJĘ

Zmień swoją stronę WWW w działający bankomat!

Dowiedz się więcej i dołącz już dzisiaj!

<http://program-partnerski.helion.pl>

GRUPA
Helion 

Niezwykłość tej książki polega na przekształcaniu złożonych zagadnień technicznych w jasne i zrozumiałe objaśnienia.

Annie Xu, starszy inżynier danych, Google

Architektury Data Fabric i Data Lakehouse, a także siatka danych pojawiły się niedawno jako alternatywy hurtowni danych. Te nowe architektury mają swoje mocne strony, ale podczas projektowania rzeczywistych rozwiązań musisz pamiętać o odróżnianiu faktów od przesadnych pochwał i niejasności. Nie zawsze jest to proste i oczywiste zadanie.

Dzięki temu praktycznemu przewodnikowi profesjonalści zajmujący się danymi dobrze zrozumieją wady i zalety poszczególnych rozwiązań. Omówiono tu typowe zagadnienia dotyczące architektur danych, w tym ich rozwój i możliwości. Żadna architektura nie jest na tyle uniwersalna, by być odpowiednia w każdej sytuacji, dlatego w książce znajdziesz rzetelne porównanie cech poszczególnych architektur. Dowiesz się, jakie kompromisy towarzyszą każdej z nich, niezależnie od popularności. W ten sposób o wiele łatwiej przyjdzie Ci wybór rozwiązania, które najlepiej odpowiada Twoim potrzebom.

**Położ tę książkę na biurku.
Będziesz często po nią sięgać!**

Sawyer Nyquist, autor, właściciel The Data Shop

Najciekawsze zagadnienia:

- praktyczne działanie architektur danych, ich mocne i słabe strony
- wybór najlepszej architektury pod kątem konkretnego zastosowania
- różnice między hurtowniami i jeziorami danych
- wspólne koncepcje architektur danych i ich historyczny rozwój
- sesje projektowania architektury, organizacja zespołów i najważniejsze uwarunkowania

James Serra — doświadczony architekt rozwiązań w Microsoftzie. Uznany lider w dziedzinie stosowania technologii big data i zaawansowanej analityki, w tym takich architektur danych jak nowoczesna hurtownia danych, siatka danych, a także Data Fabric i Data Lakehouse.

Helion

KOD KORZYŚCI
Sięgnij po więcej! ▶



helion.pl



HELION S.A.
ul. Kościuszki 1c
44-100 Gliwice
tel.: 32 259 98 63
helion@helion.pl

ISBN 978-83-289-1669-2



9 788328 916692

Cena: 79,00 zł