
Trino: Profesjonalny przewodnik

SQL w dowolnej skali, w dowolnym magazynie i w dowolnym środowisku

*Matt Fuller, Manfred Moser
i Martin Traverso*

przekład: Joanna Zatorska

Spis treści

Przedmowa	xiii
Wprowadzenie.....	xv

Część I. Wprowadzenie do Trino

1. Wstęp do Trino	3
Problemy związane z wielkimi zbiorami danych.....	3
Trino na ratunek.....	5
Stworzony z myślą o wydajności i skalowalności	5
SQL do wszystkiego.....	6
Oddzielenie magazynu danych od zasobów służących do przetwarzania zapytań	7
Przypadki użycia Trino	7
Jeden analityczny punkt dostępu za pomocą SQL	8
Punkt dostępu do hurtowni danych i systemów źródłowych	8
Dostęp do wszystkiego za pomocą SQL.....	9
Zapytania federacyjne.....	10
Warstwa semantyczna dla wirtualnej hurtowni danych.....	10
Silnik zapytań w jeziorze danych.....	11
Przekształcenia SQL i ETL.....	11
Lepszy wgląd dzięki szybszym wynikom	12
Wielkie zbiory danych, uczenie maszynowe i sztuczna inteligencja	12
Inne przypadki użycia.....	12
Zasoby dotyczące Trino.....	12
Witryna	12
Dokumentacja	13
Czat społeczności.....	13
Kod źródłowy, licencja i wersja	14

Współpraca	14
Repozytorium dla książki	15
Zbiór danych Iris	15
Zbiór danych Flight	16
Krótką historia Trino	16
Podsumowanie	18
2. Instalowanie i konfigurowanie Trino	19
Wypróbowanie Trino w kontenerze Dockera	19
Instalowanie za pomocą pliku archiwum	21
Java Virtual Machine	21
Python	21
Instalacja	21
Konfiguracja	22
Dodawanie źródła danych	23
Uruchamianie Trino	24
Podsumowanie	25
3. Używanie Trino	27
Interfejs wiersza poleceń Trino	27
Zaczynamy	27
Stronicowanie	30
Historia i uzupełnianie	30
Dodatkowa diagnostyka	31
Wykonywanie zapytań	31
Formaty wynikowe	32
Ignorowanie błędów	32
Sterownik JDBC dla Trino	32
Pobieranie i rejestrowanie sterownika	34
Nawiązywanie połączenia z Trino	35
Trino i ODBC	37
Biblioteki klienckie	37
Interfejs internetowy Trino	38
SQL w Trino	39
Koncepcje	39
Pierwsze przykłady	40
Podsumowanie	43

Część II. Zagłębiamy się w Trino

4. Architektura Trino	47
Koordynator i węzły robocze w klastrze	47
Koordynator	49
Usługa wykrywania	49
Węzły robocze	50
Architektura oparta na konektorach	50
Katalogi, schematy i tabele	52
Model wykonywania zapytań	52
Planowanie zapytania	56
Parsowanie i analiza	57
Początkowe planowanie zapytania	58
Reguły optymalizacji	60
Mechanizm przekazywania predykatu w głąb	60
Eliminacja złączeń krzyżowych	61
TopN	62
Reguły implementacji	63
Dekorelacja złączenia bocznego	63
Dekorelacja złączeń częściowych (IN)	64
Optymalizator oparty na kosztach	65
Koncepcja kosztu	66
Koszt złączenia	67
Statystyki tabel	68
Statystyki filtrowania	70
Statystyki tabel podzielonych na partycje	71
Enumeracja złączeń	71
Złączenia rozgłaszane kontra rozproszone	71
Korzystanie ze statystyk tabeli	73
Polecenie Trino ANALYZE	74
Zbieranie statystyk podczas zapisu na dysku	75
Polecenie Hive ANALYZE	75
Wyświetlanie statystyk tabel	75
Podsumowanie	76
5. Wdrażanie w środowisku produkcyjnym	77
Szczegółowe informacje o konfiguracji	77
Konfiguracja serwera	77
Logowanie	78
Konfiguracja węzła	79
Konfiguracja JVM	80
Skrypt startowy	81

Instalacja klastra	82
Instalacja RPM	84
Instalacja struktury katalogów	85
Konfiguracja	85
Odinstalowywanie Trino	86
Instalacja w chmurze	86
Pakiet Helm dla wdrożenia w platformie Kubernetes	87
Rozważania na temat rozmiaru klastra	88
Podsumowanie	89
6. Konektory	91
Konfiguracja	92
Przykład konektora RDBMS: PostgreSQL	93
Przesyłanie zapytania w głąb	95
Obliczenia równoległe i współbieżność	96
Inne konektory RDBMS	97
Bezpieczeństwo	98
Przekazywanie zapytań	98
Konektory Trino TPC-H i TPC-DS	99
Konektor Hive dla rozproszonych źródeł danych	100
Apache Hadoop i Hive	100
Konektor Hive	101
Format tabel Hive	103
Tabele zarządzane i zewnętrzne	104
Dane podzielone na partycje	105
Wczytywanie danych	107
Formaty plików i kompresja	110
Przykład MinIO	111
Zarządzanie nowoczesnym systemem magazynowym i jego analiza	111
Nierelacyjne źródła danych	113
Konektor JMX dla Trino	113
Konektor Black Hole	115
Konektor memory	116
Inne konektory	116
Podsumowanie	117
7. Przykłady zaawansowanych konektorów	119
Łączenie się z HBase za pomocą narzędzia Phoenix	119
Przykład konektora dla magazynu typu klucz-wartość: Accumulo	120
Używanie konektora Accumulo w Trino	123
Przesyłanie predykatów w Accumulo	125
Konektor Apache Cassandra	128
Przykład konektora systemu strumieniowego: Kafka	128

Przykład konektora dla magazynu opartego na dokumentach: Elasticsearch.	130
Ogólne omówienie.	130
Konfiguracja i użycie.	131
Przetwarzanie zapytań.	132
Wyszukiwanie pełnotekstowe	132
Podsumowanie.	132
Federacja zapytań w Trino	133
Operacje ekstrakcji, transformacji i ładowania z zapytaniami federacyjnymi	139
Podsumowanie	140
8. Użycia SQL w Trino	141
Instrukcje Trino	142
Tabele systemowe Trino	144
Katalogi	146
Schematy	147
Schemat informacji	149
Tabele	150
Właściwości tabeli i kolumn.	152
Kopiowanie istniejącej tabeli	153
Tworzenie nowej tabeli na podstawie wyników zapytania.	154
Modyfikowanie tabeli	155
Usuwanie tabeli	156
Ograniczenia związane z tabelami nakładane przez konektory	156
Widoki	157
Informacje o sesji i konfiguracja	158
Typy danych	159
Typy danych kolekcji.	162
Typy danych dotyczące czasu.	163
Rzutowanie typów	167
Wprowadzenie do instrukcji SELECT.	168
Klauzula WHERE.	170
Klauzule GROUP BY i HAVING.	171
Klauzule ORDER BY i LIMIT	173
Instrukcje JOIN	174
Klauzule UNION, INTERSECT i EXCEPT	175
Operacje grupowania	177
Klauzula WITH	178
Podzapytania.	180
Podzapytanie skalarne.	180
Podzapytanie EXISTS	180
Podzapytanie ilościowe.	181
Usuwanie danych z tabeli	182
Podsumowanie	182

9. Zaawansowany SQL	183
Wprowadzenie do funkcji i operatorów	183
Funkcje skalarne i operatory	184
Operatory logiczne.....	185
Operatory logiczne.....	186
Wybór zakresu za pomocą instrukcji BETWEEN.....	187
Wykrywanie wartości za pomocą instrukcji IS (NOT) NULL	187
Funkcje i operatory matematyczne	188
Funkcje trygonometryczne.....	189
Funkcje zwracające liczby stałe i losowe.....	190
Funkcje i operatory dotyczące ciągów tekstowych	190
Ciągi tekstowe i mapy	192
Unicode	193
Wyrażenia regularne	195
Spłaszczanie złożonych typów danych	197
Funkcje JSON.....	198
Funkcje i operatory dotyczące daty i czasu.....	199
Histogramy	202
Funkcje agregujące.....	203
Funkcje agregujące zwracające mapy	203
Przybliżone funkcje agregujące	206
Funkcje okna.....	207
Wyrażenia lambda	208
Funkcje geoprzestrzenne.....	209
Przygotowane instrukcje.....	210
Podsumowanie	212

Część III. Rzeczywiste przypadki użycia Trino

10. Bezpieczeństwo	215
Uwierzytelnianie.....	216
Uwierzytelnianie za pomocą hasła i protokołu LDAP	217
Inne typy uwierzytelniania.....	220
Autoryzacja	220
Systemowa kontrola dostępu	220
Kontrola dostępu na poziomie konektora	224
Szyfrowanie.....	226
Szyfrowanie komunikacji między klientem a koordynatorem w Trino.....	228
Tworzenie repozytoriów keystore i truststore Java	231
Szyfrowanie komunikacji w klastrze Trino	233
Urząd certyfikacji kontra samopodpisane certyfikaty	235

Uwierzytelnianie za pomocą certyfikatu	237
Kerberos	240
Wymagania wstępne	240
Uwierzytelnianie klienta za pomocą protokołu Kerberos	240
Dostęp do źródła danych i konfiguracja zabezpieczeń	241
Uwierzytelnianie za pomocą protokołu Kerberos w konektorze Hive	242
Uwierzytelnianie w usłudze Hive Metastore Service	243
Uwierzytelnianie w HDFS	244
Separacja klastra	244
Podsumowanie	245
11. Integrowanie Trino z innymi narzędziami	247
Zapytania, wizualizacje i inne operacje z użyciem Apache Superset	247
Lepsza wydajność dzięki platformie RubiX	248
Cykle pracy z użyciem Apache Airflow	249
Przykład wbudowanego Trino: Amazon Athena	249
Wygodne dystrybucje komercyjne: Starburst Enterprise i Starburst Galaxy	253
Przykłady innych integracji	254
Niestandardowe integracje	255
Podsumowanie	255
12. Trino w środowisku produkcyjnym	257
Monitorowanie za pomocą interfejsu internetowego Trino	257
Szczegóły dotyczące klastra	258
Lista zapytań	259
Widok Query Details	262
Dostrajanie zapytań SQL w Trino	268
Zarządzanie pamięcią	272
Współbieżność zadań	275
Planowanie zadań w węźle roboczym	275
Wymiana danych przez sieć	276
Współbieżność	276
Rozmiary bufora	276
Dostrajanie wirtualnej maszyny Java	277
Grupy zasobów	278
Definicja grupy zasobów	280
Zasada planowania	281
Definicja reguł selektora	282
Podsumowanie	282
13. Rzeczywiste przykłady	283
Platformy wdrożeniowe	284
Dobór rozmiaru klastra	285

Przypadek migracji Hadoop/Hive	286
Inne źródła danych	287
Użytkownicy i ruch	287
Podsumowanie	288
Podsumowanie	289
Indeks	291
0 autorach	309