

# Hurtownie danych

31 stycznia 2017

**Hurtownia danych** – wg Williama Inmona – zbiór danych wyróżniający się następującymi cechami

- uporządkowany tematycznie
- zintegrowany
- zawierający wymiar czasowy
- nieulotny
- wspomaga podejmowanie decyzji
- wspomaga przetwarzanie informacji dla celów strategicznych i analitycznych

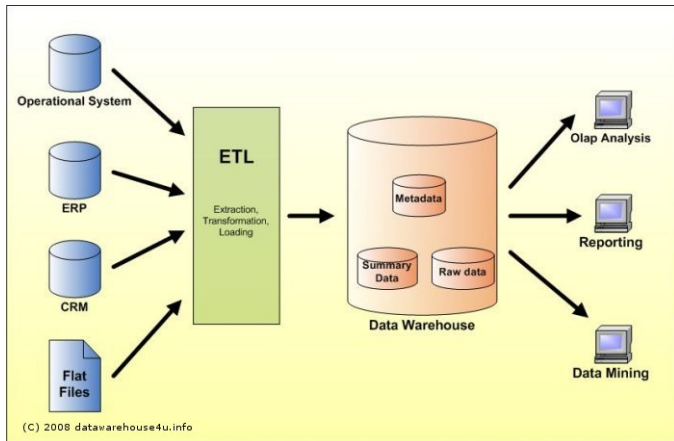
# Przetwarzanie danych

	OLTP OnLine Transaction Processing	OLAP OnLine Analytical Processing
źródło danych	dane operacyjne	dane historyczne
cel	aktualizacje	raportowanie
zastosowania	usługi online, klient-serwer	inteligencja biznesowa, data mining
użytkownicy	klienci, personel	zarząd, marketing
horyzont czasowy	dni, tygodnie, miesiące	lata, dziesiątki lat
model danych	relacyjny	wielowymiarowy
schemat	znormalizowany	gwiazda, płatek śniegu, konstelacja
zapytania	proste	złożone
operacje na danych	szybkie transakcje insert oraz update	cykliczne, długotrwałe importy
optymalizacja	zapisów	odczytu
czas przetwarzania	krótki	długi
objętość	mała	wielka

Uwzględniając wymagania jakie są stawiane przed hurtownią danych przy projektowaniu jej schematu należy zwrócić uwagę na dwa aspekty:

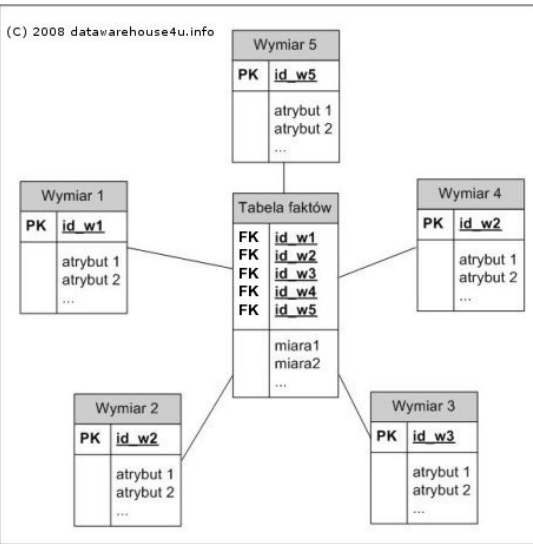
- bardzo duża ilość danych: bieżące + historyczne
- osiągnięcie zadawalającego poziomu efektywności zapytań analitycznych

# Architektura systemu hurtowni danych

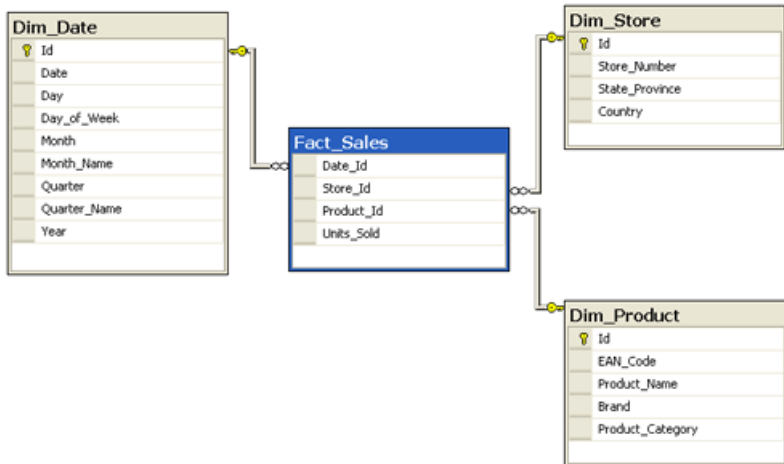


- Schemat gwiazdy (ang. star)
- Schemat płatka śniegu (ang. snowflake)
- Schemat konstelacji faktów (ang. fact constellation albo starflake)

# Schemat gwiazdy



# Schemat gwiazdy: przykład





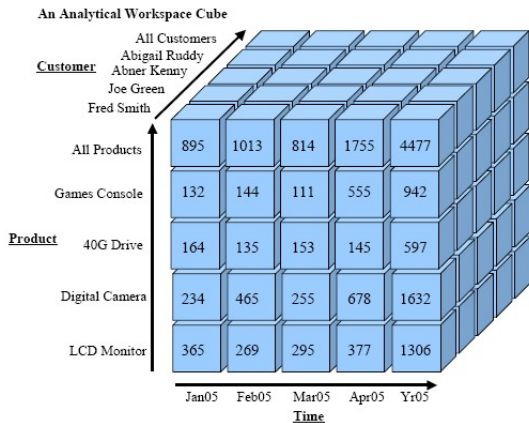
# Schemat gwiazdy: cechy

- Prosta struktura
- Duża efektywność zapytań z uwagi na małą liczbę powiązań
- Długi czas ładowania danych do tabel wymiarów z uwagi na denormalizację i w efekcie na redundancję danych
- Tabela faktów składa się z dwóch typów kolumn:
  - **miary** – wartości numeryczne opisujące dany fakt
  - klucze obce do tabeli wymiarów

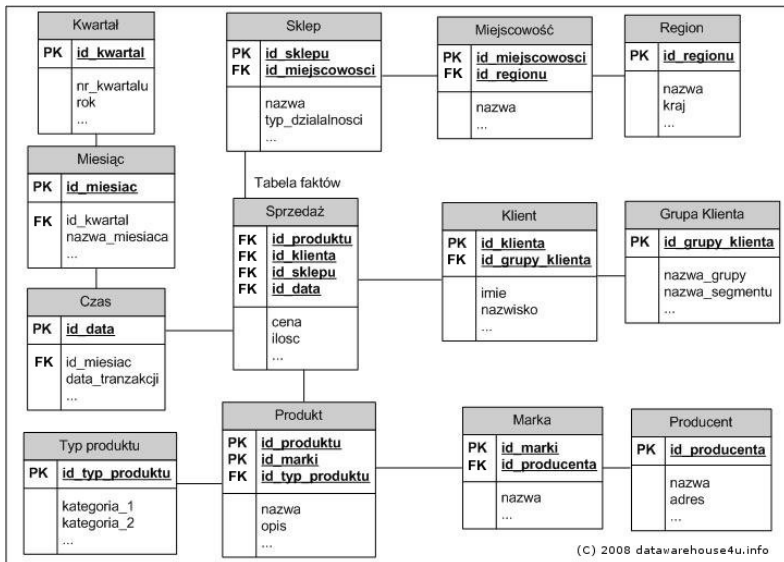
**Wymiary** – wartości opisowe danego faktu

- Klucz główny tabeli faktów składa się z jej wszystkich kluczy obcych
- Tabela faktów może zawierać informacje o faktach na poziomie detalicznym lub zagregowanym

## Wielowymiarowa kostka OLAP



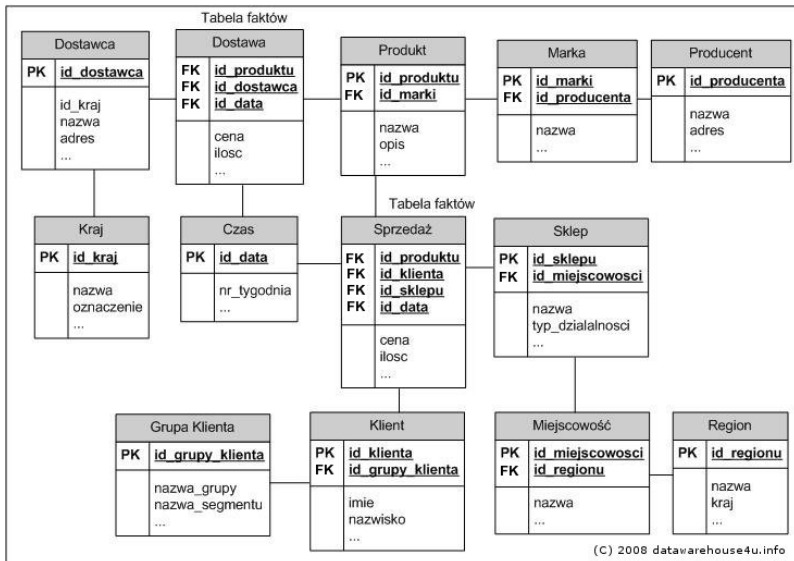
# Schemat płątka śniegu



# Schemat płątka śniegu: cechy

- Spadek wydajności zapytań w porównaniu ze schematem gwiazdy z uwagi na większą liczbę relacji
- Struktura łatwiejsza do modyfikacji
- Relatywnie krótszy czas ładowania danych do tabel ze względu na mniejszą redundancję dzięki normalizacji
- Wykorzystywany rzadziej niż schemat gwiazdy, gdyż efektywność zapytań jest ważniejsza niż efektywność ładowania danych

# Schemat konstelacji faktów



# Schemat konstelacji faktów: cechy

- Rozwiązanie pośrednie między gwiazdą a płatkami śniegu
- Zwykle występuje więcej niż jedna tabela faktów
- Część tabel wymiarów jest znormalizowana a część zdenormalizowana
- Tabele faktów są ze sobą powiązane w relacji 1 : 1 lub 1 :  $n$
- Tabele faktów mogą współdzielić te same tabele wymiarów

- Tabela faktów bez faktów nie zawiera miar
- Przekrój wymiarów

## Przykład

Obecność studentów na zajęciach jako przekrój 3 wymiarów: *studenci*, *czas*, *zajęcia*. Tabela faktów w tym wypadku złożona jest z trzech kolumn zawierających klucze obce. Pozwala łatwo udzielić odpowiedzi na następujące pytania:

- Ilu studentów było obecnych na danych zajęciach w określonym terminie?
- Na ile zajęć średnio uczęszczał dany student danego dnia?

# Junk Dimension

- Nader często zdarzają się sytuacje, gdzie miary są jednobitowe (1/0)
- Przechowywanie danych jednobitowych w tabelach faktów wymaga użycia małych tabel wymiarów
- Jednocześnie ilość danych w tabelach faktów gwałtownie rośnie

## Przykład

FACT\_TABLE

CUSTOMER_ID
PRODUCT_CD
TXN_ID
STORE_ID
TXN_CODE
COUPON_IND
PREPAY_IND
TXN_AMT



FACT\_TABLE

CUSTOMER_ID
PRODUCT_CD
TXN_ID
STORE_ID
JUNK_ID
TXN_AMT



DIM\_JUNK

JUNK_ID	TXN_CODE	COUPON_IND	PREPAY_IND
1	1	Y	Y
2	2	Y	Y
3	3	Y	Y
4	1	Y	N
5	2	Y	N
6	3	Y	N
7	1	N	Y
8	2	N	Y
9	3	N	Y
10	1	N	N
11	2	N	N
12	3	N	N



- ROLAP — Relational OLAP
  - Przetwarzanie dużej ilości danych
  - Naturalne wykorzystanie własności relacyjnych baz danych
  - Niezbyt duża wydajność
  - Ograniczeni narzucone przez model relacyjny
  - IBM DB2 OLAP, Informix MetaCube, MicroStrategy
- MOLAP — Multidimensional OLAP
  - Znakomita wydajność
  - Wylicznia generowane w czasie tworzenia kostki, co pozwala szybko wykonywać skomplikowane operacje
  - Ograniczenia pojemności
  - Oracle Express Server
- HOLAP — Hybrid OLAP
  - Podsumowania w formie wielowymiarowych kostek
  - Dane szczegółowe w bazie relacyjnej
  - Microsoft SQL Server OLAP Services