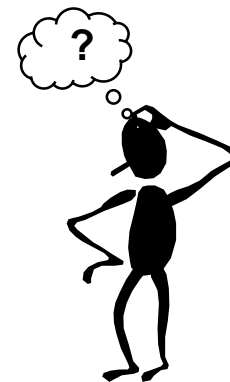


- **Olbrzymie zasoby informacji**
  - *różna struktura*
  - *różna jakość*
- **Łatwość dostępu do informacji**  
(wystarczy mieć przeglądarke...)
- **Łatwość wprowadzania własnych danych**  
(wystarczy mieć serwer WWW...)

**Zalety** czy **wady** ?

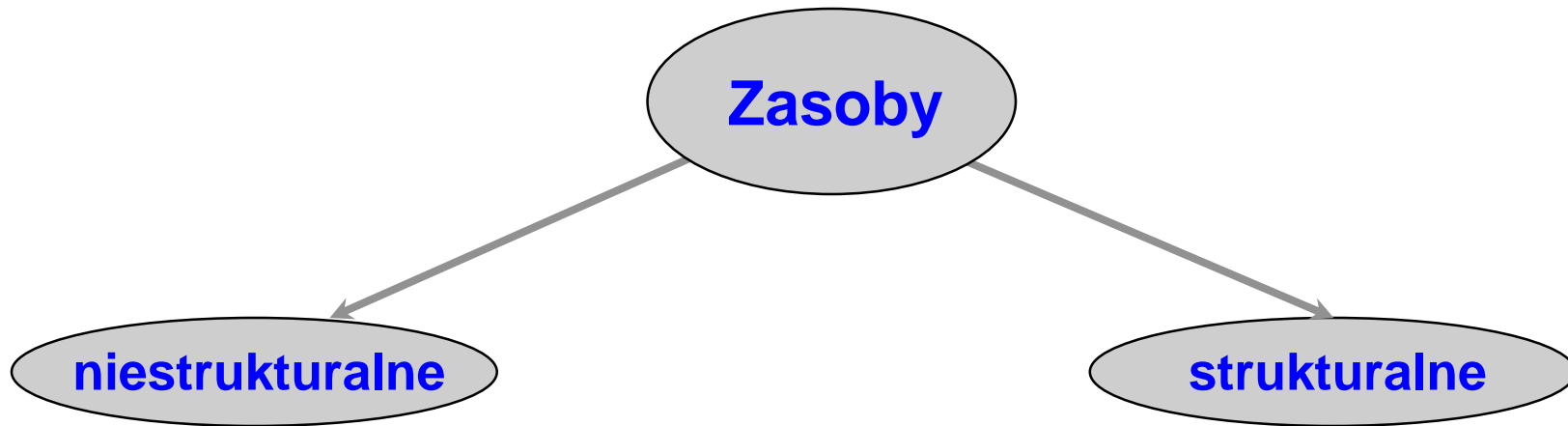


# Problemy



- Jak znaleźć potrzebne informacje?
  - *odkrywanie źródeł informacji*
- Jak odróżnić informacje użyteczne od zbędnych?
  - *filtrowanie informacji*
- Gdzie gromadzić informacje?
  - *Web warehousing*
- Jak wykorzystać pozyskane informacje?
  - *przetwarzanie analityczne (OLAP)*

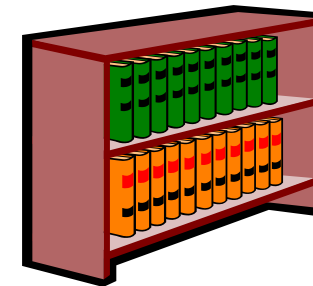




nie istnieje odpowiadający im  
*schemat danych*





istnieje ściśle zdefiniowany  
*schemat danych*



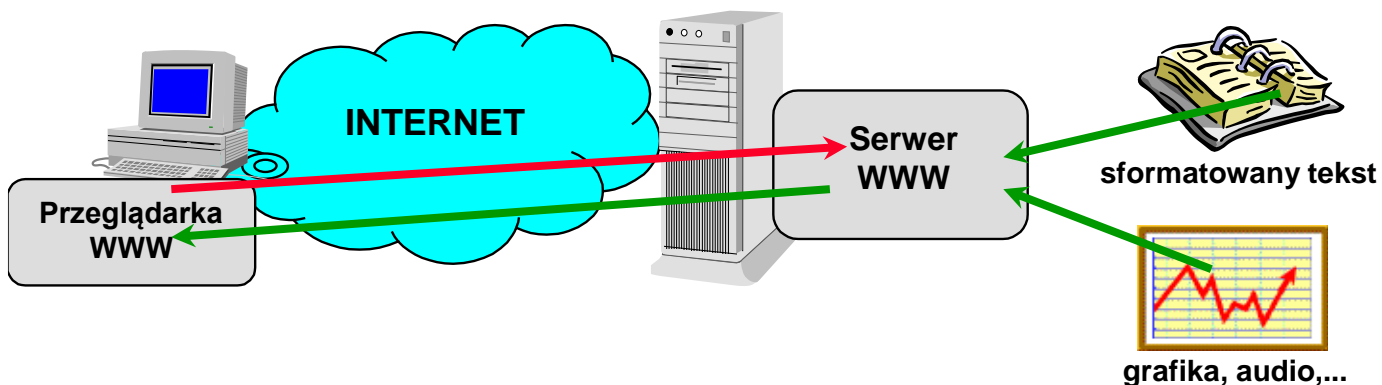
# Zasoby niestrukturalne

Styczne strony (pliki) HTML zawierające

- tekst
- grafikę
- dane audio  wideo 
- odsyłacze do innych stron/plików

Identyfikowane adresem URL (*Uniform Resource Locator*):

*protokół // komputer.domena / plik*



# Zasoby strukturalne

Dane zgromadzone w bazach danych dostępnych poprzez Internet

## Schemat danych + Dane

Dostęp poprzez adres URL rozszerzony o parametry zapytania:

*protokół // komputer.domena / procedura ? p1=w1&p2=w2 ...*



***Dynamiczna strona HTML***

Przykład:

<http://quote.yahoo.com/q?s=intc&d=t>



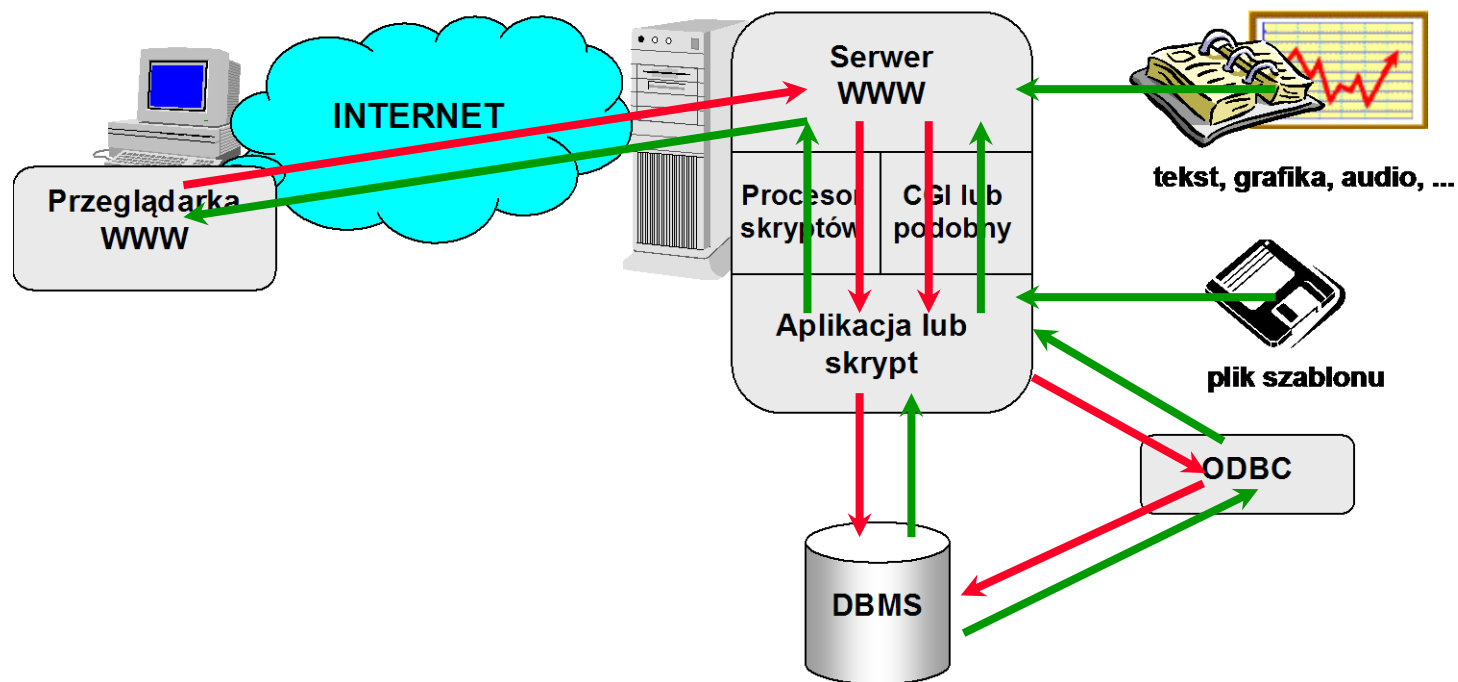
Views: [Basic](#) - [DayWatch](#) - [Performance](#) - [Real-time Mkt](#) - **[Detailed](#)** - [[Create New View](#)]

INTEL CORP (NasdaqNM:INTC) - Trade: <a href="#">Choose Brokerage</a>						
Last Trade Jun 5 · <b>28.18</b>	Change +0.68 (+2.47%)		Prev Cls 27.50	Open 27.60	Volume 41,657,800	 <p>Small: <a href="#">1d</a> <a href="#">5d</a> <b><a href="#">1y</a></b> <a href="#">none</a> Big: <a href="#">1d</a> <a href="#">5d</a> <a href="#">3m</a> <a href="#">6m</a> <a href="#">1y</a> <a href="#">2y</a> <a href="#">5y</a> <a href="#">max</a></p>
Day's Range 26.89 - 28.20	Bid 28.12	Ask 28.18	P/E 105.77	Mkt Cap 188.4B	Avg Vol 43,997,045	
52-wk Range 18.96 - 36.78	Bid Size 2,000	Ask Size 700	P/S 6.90	Div/Shr 0.08	Div Date Jun 1	
1y Target Est 41.38	EPS (ttm) 0.26	EPS Est 0.70	PEG 2.18	Yield 0.29	Ex-Div May 3	
<a href="#">Chart</a> , <a href="#">Financials</a> , <a href="#">Historical Prices</a> , <a href="#">Insider</a> , <a href="#">Messages</a> , <a href="#">News</a> , <a href="#">Options</a> <a href="#">Profile</a> , <a href="#">Reports</a> , <a href="#">Research</a> , <a href="#">SEC Filings</a> , <a href="#">Upgrades</a> , <a href="#">more...</a>						
<a href="#">Compare 22 long-distance plans at a glance.</a> Yahoo! Savings Finder.						

[Add to My Portfolio](#) - [Set Alert](#)

[Download Spreadsheet](#)

# Zasoby strukturalne



# Zasoby semistrukturalne

## Dane semistrukturalne:

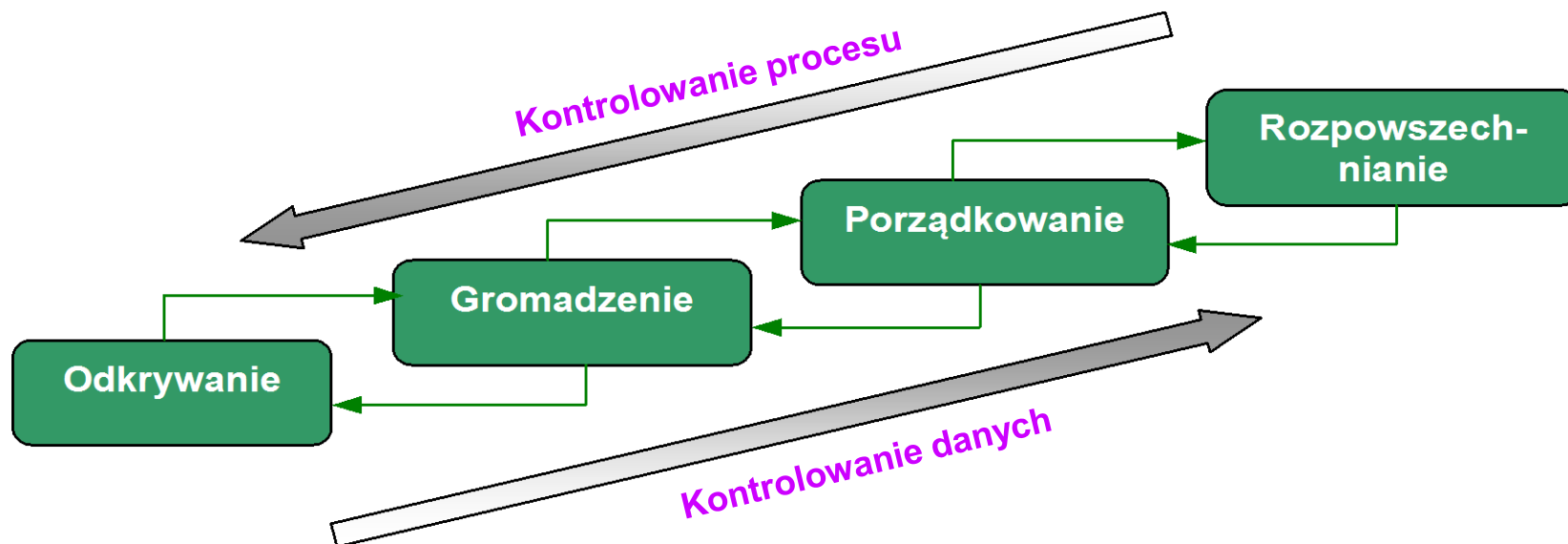
- o nieznanym (lub znanym częściowo) schemacie danych
- o strukturze nieregularnej
  - pola opcjonalne
  - alternatywne kombinacje pól
  - pola o nieokreślonej liczności
- niekompletne
- często zawierające błędy i niespójności

Język do opisu i zapisu tego typu danych:

e**X**tensible **M**arkup **L**anguage (**XML**)



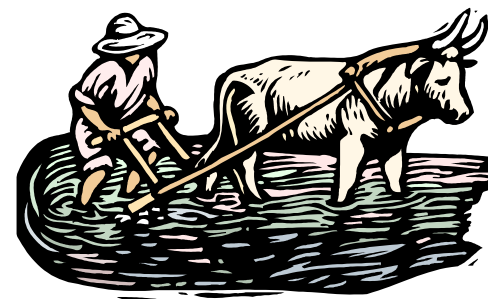
Metoda systematycznego pozyskiwania i wykorzystywania zasobów informacyjnych Internetu



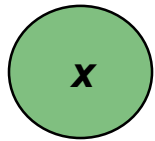
## Poszukiwanie w Internecie miejsc istotnych z punktu widzenia prowadzonej działalności biznesowej

- ❑ Etap najtrudniejszy w całym procesie (trudny do zautomatyzowania)
  
- ❑ Narzędzia:
  - przeglądarki
  - wyszukiwarki
  - agenci

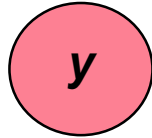
*Przygotowanie roli pod uprawę*



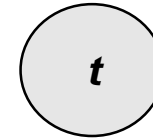
# Trafność a precyzja odkrywania



miejsca  
istotne



miejsca  
odkryte

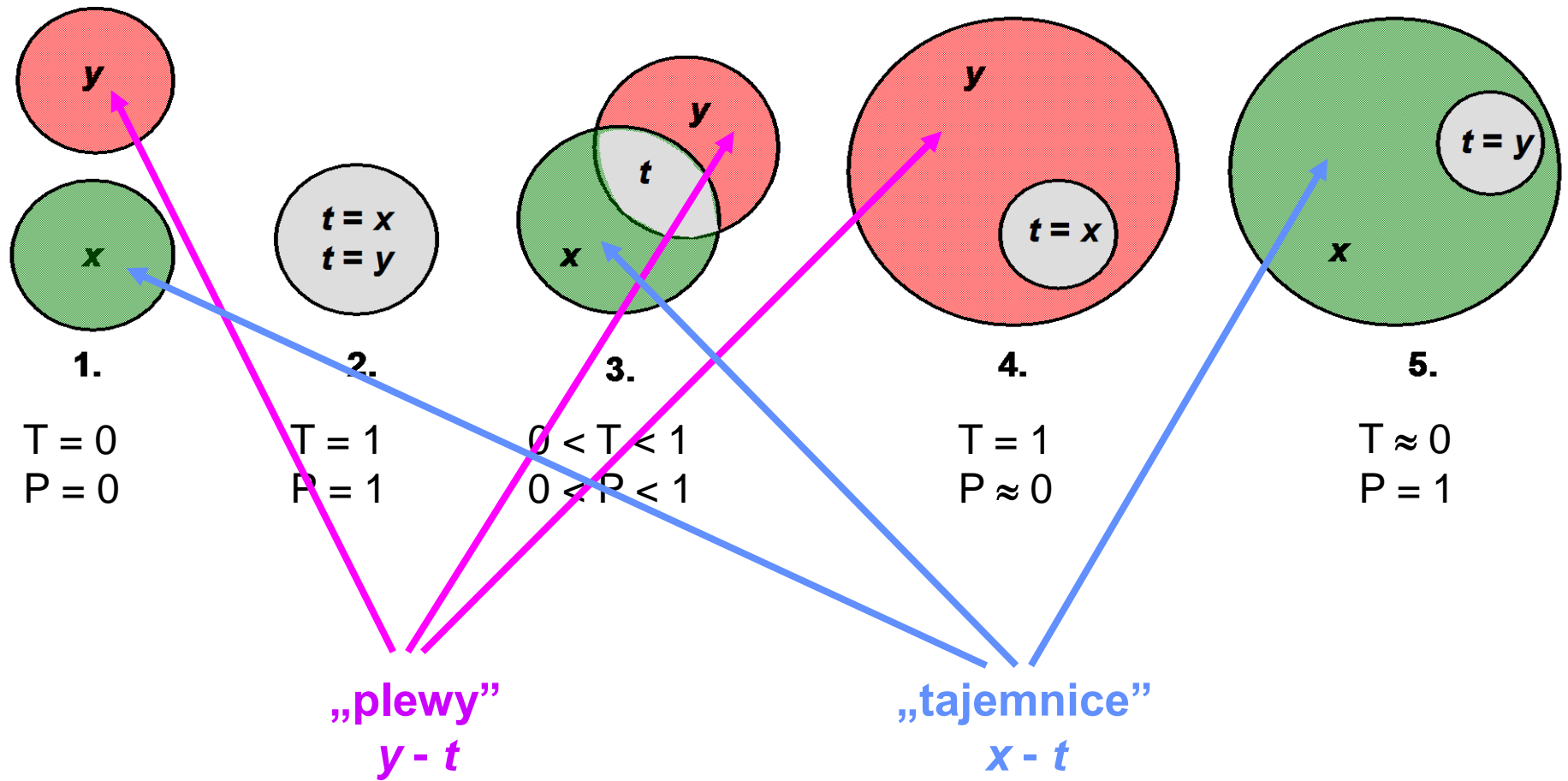


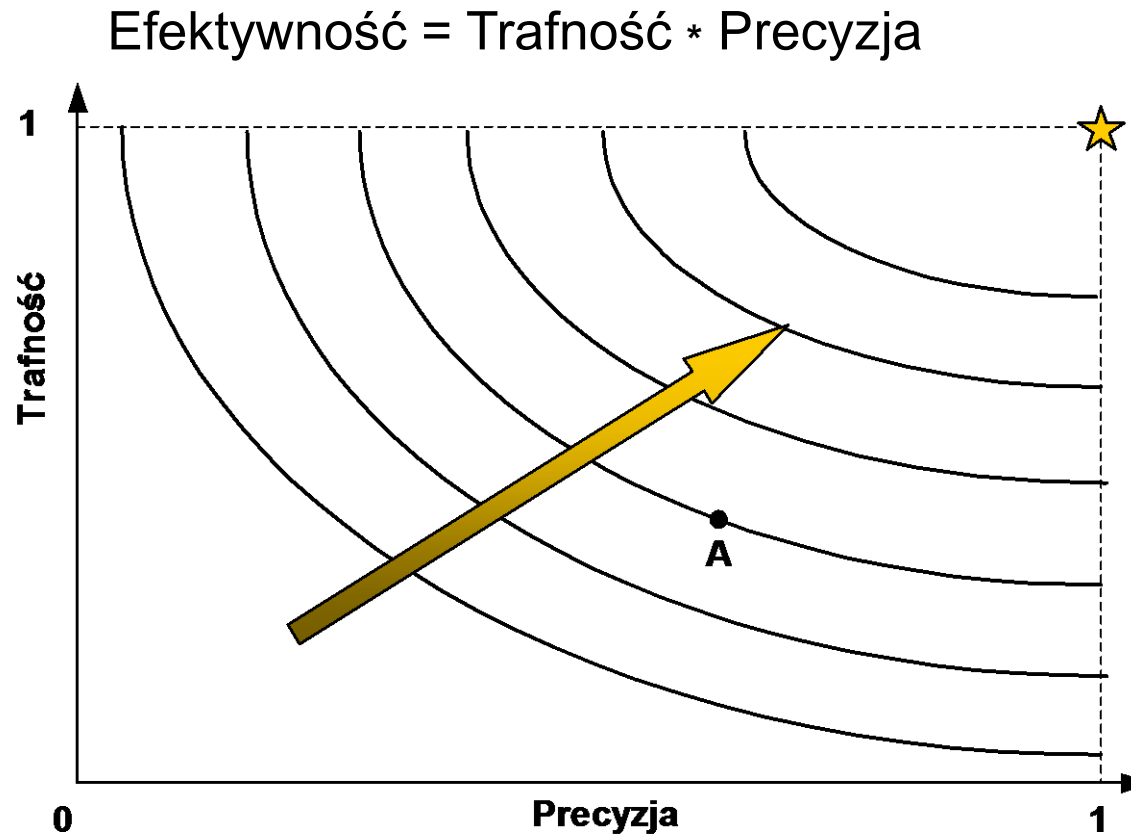
odkryte  
miejsca  
istotne

**Precyzja**  $P = t / y$  [0..1] (*precision*)

**Trafność**  $T = t / x$  [0..1] (*recall*)

# Trafność a precyzja odkrywania





Miara-F =  $2 * \text{Trafność} * \text{Precyzja} / (\text{Trafność} + \text{Precyzja})$   
- średnia harmoniczna trafności i precyzji

## Zbieranie treści statycznych i dynamicznych stron WWW zidentyfikowanych na etapie odkrywania

- Podlega automatyzacji (głównie w odniesieniu do stron dynamicznych)
- Narzędzia:
  - agenty programowe
  - narzędzia systemu plików
  - bazy danych

*Obsiewanie roli*



## Nadanie danym struktury stosownej do ich przechowywania i przetwarzania w hurtowni danych

- ❑ Operacje:
  - sprawdzanie danych
  - czyszczenie danych
  - transformowanie danych
  
- ❑ Narzędzia:
  - bazy danych i hurtownie danych
  - narzędzia do analizowania danych (*business intelligence*)
  - narzędzia do eksploracji danych

*Pielęgnowanie zasiewów*



**Dostarczanie uporządkowanych danych do miejsc,  
w których są potrzebne (czyli do ich konsumentów)**

- Bezpośrednie przekazanie danych albo przekazanie wyników analiz
- Narzędzia:
  - bazy danych i hurtownie danych
  - narzędzia do analizowania danych (*business intelligence*)
  - narzędzia do eksploracji danych
  - narzędzia do prezentacji danych

**Żniwa!**

