

# Hbase, Hive i BigSQL



**KAPITAŁ LUDZKI**  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwestycji  
Projekt nr POKL.04.01.00-00-151/12 „Inżynieria wiedzy dla inteligentnego rozwoju”  
współfinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



**UNIA EUROPEJSKA**  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPŁOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

str. 1

## *Agenda*

1. NOSQL a HBase
2. Architektura HBase
3. Demo HBase
4. Po co Hive?
5. Apache Hive
6. Demo hive
7. BigSQL



**KAPITAŁ LUDZKI**  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwestycji  
Projekt nr POKL.04.01.00-00-151/12 „Inżynieria wiedzy dla inteligentnego rozwoju”  
współfinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



**UNIA EUROPEJSKA**  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPŁOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

# *HBase*

Jest to

- rozproszona
  - trwała
  - posortowana
  - wielowymiarowa
- mapa.

A P A C H E  
**HBASE**



**KAPITAŁ LUDZKI**  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwestycji  
Projekt nr POKL.04.01.00-00-131/12 „Inżynieria wiedzy dla inteligentnego rozwoju”  
wspofinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



**UNIA EUROPEJSKA**  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPLOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

# *NoSQL*

Not only SQL



**KAPITAŁ LUDZKI**  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwestycji  
Projekt nr POKL.04.01.00-00-131/12 „Inżynieria wiedzy dla inteligentnego rozwoju”  
wspofinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



**UNIA EUROPEJSKA**  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPLOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

## *CAP*

Żaden rozproszony system nie może jednocześnie spełnić wszystkich poniższych wymagań:

- spójność (consistency),
- dostępność (availability),
- odporność na uszkodzenia węzłów (partition tolerance).



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwentyrio  
Projekt nr POKL.04.01.00-00-131/12 „Inzynieria wiedzy dla inteligentnego rozwoju”  
wspofinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPLOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

## *ACID*

- atomowość (atomicity),
- spójność (consistency),
- izolacja (isolation),
- trwałość (durability).



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwentyrio  
Projekt nr POKL.04.01.00-00-131/12 „Inzynieria wiedzy dla inteligentnego rozwoju”  
wspofinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPLOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

## *Dlaczego Hbase?*

- jest dystrybuowana wraz z Hadoop,
- z powodzeniem może zastąpić kosztowne systemy RDBMS w rozwiązaniach BigData,
- bardzo dobrze się skaluje,
- wspiera elastyczny model danych (wszystko jest sekwencją bajtów).



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwestycji  
Projekt nr POKL.04.01.00-00-131/12 „Inżynieria wiedzy dla inteligentnego rozwoju”  
wspofinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPLOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

## *Kiedy Hbase jest złym rozwiązaniem?*

- nie wspiera SQL,
- nie jest zaprojektowana do przetwarzania transakcji,
- nie wspiera łączenia tabel (joins)



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwestycji  
Projekt nr POKL.04.01.00-00-131/12 „Inżynieria wiedzy dla inteligentnego rozwoju”  
wspofinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPLOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

## Konwersja z RDBMS do HBase

id	Name	Age	Interests
1	Ricky		Soccer, Movies, Baseball
2	Ankur	20	
3	Sam	25	Music

Multi-valued

null

id	Name
1	Ricky
2	Ankur
3	Sam

id	Age
2	20
3	25

id	Interests
1	Soccer
1	Movies
1	Baseball
3	Music



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwestycji  
Projekt nr POKL.04.01.00-00-151/12 „Inżynieria wiedzy dla inteligentnego rozwoju”  
wspofinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY

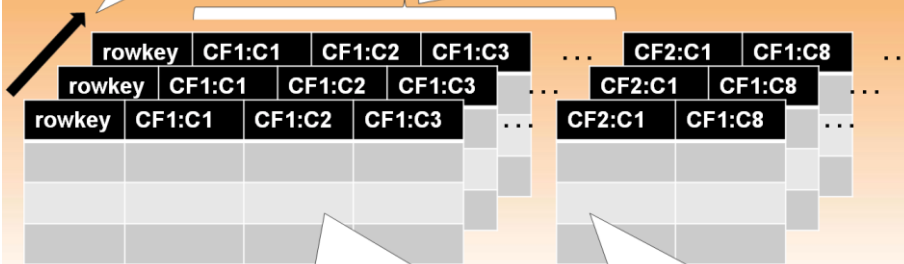


MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPLOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNIE

## Struktura tabeli w HBase

Multi-versioned

One column family can have variable no of columns



Cells within a Column family are sorted physically

Very Sparse, most cell has NULL value



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwestycji  
Projekt nr POKL.04.01.00-00-151/12 „Inżynieria wiedzy dla inteligentnego rozwoju”  
wspofinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego

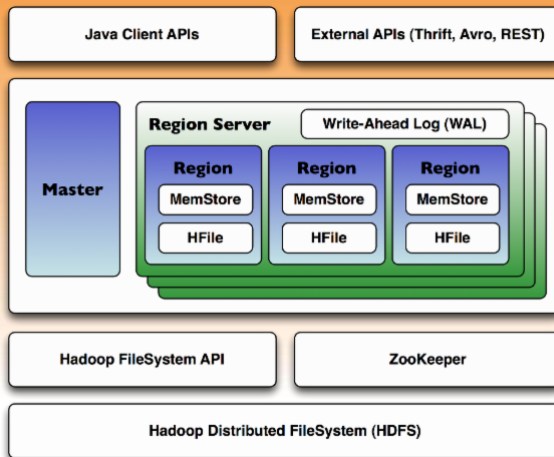


UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPLOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNIE

## Architektura HBase



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwestycji  
Projekt nr POKL.04.01.00-00-151/12 „Inżynieria wiedzy dla inteligentnego rozwoju”  
wspofinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPŁOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

## HBase - demo

- tworzenie tabel,
- dodawanie danych,
- pobieranie danych



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwestycji  
Projekt nr POKL.04.01.00-00-151/12 „Inżynieria wiedzy dla inteligentnego rozwoju”  
wspofinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPŁOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

## *Apache Hive*



- stworzony przez Facebook w 2007,
- codzienne zadania były zbyt czasochłonne dla systemów RDBMS,
- pisanie aplikacji MapReduce było zbyt czasochłonne.



**KAPITAŁ LUDZKI**  
NARODOWA STRATEGIA SPÓJNOŚCI



CeloweK – najfiszaz inwentyzje  
Projekt nr POKL.04.01.00-00-151/12 „Inzynieria wiedzy dla inteligentnego rozwoju”  
wspofinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



**UNIA EUROPEJSKA**  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPŁOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

## *Czym jest Hive?*

- hurtownią danych zbudowaną na Hadoop,
- udostępnia interfejs podobny do SQL (zwany HiveQL lub HQL),
- zamienia zapytania w HQL na zadania MapReduce,
- pozwala na narzucenie struktury na nieustrukturyzowane dane



**KAPITAŁ LUDZKI**  
NARODOWA STRATEGIA SPÓJNOŚCI



CeloweK – najfiszaz inwentyzje  
Projekt nr POKL.04.01.00-00-151/12 „Inzynieria wiedzy dla inteligentnego rozwoju”  
wspofinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



**UNIA EUROPEJSKA**  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPŁOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

## Czym nie jest Hive?

- nie zastąpi RDBMS.
- nie nadaje się do przetwarzania transakcyjnego (dobry dla zapytań analitycznych),
- implementuje “schema on read” przez co odczyty są wolniejsze,
- nie wspiera wszystkich cech SQL, nie pozwala na operacje insert, update i delete na poziomie wiersza,
- nie wspiera transakcji.



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwestycji  
Projekt nr POKL.04.01.00-00-131/12 „Intynieria wiedzy dla inteligentnego rozwoju”  
współfinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego

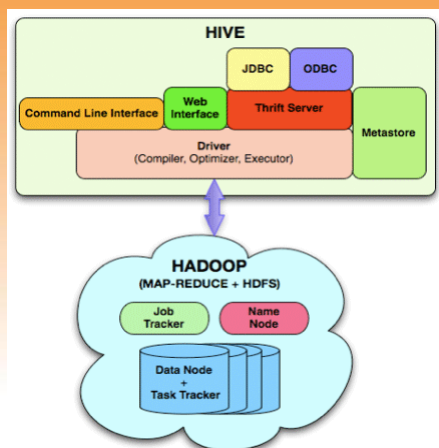


UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPLOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

## Architektura Hive



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwestycji  
Projekt nr POKL.04.01.00-00-131/12 „Intynieria wiedzy dla inteligentnego rozwoju”  
współfinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



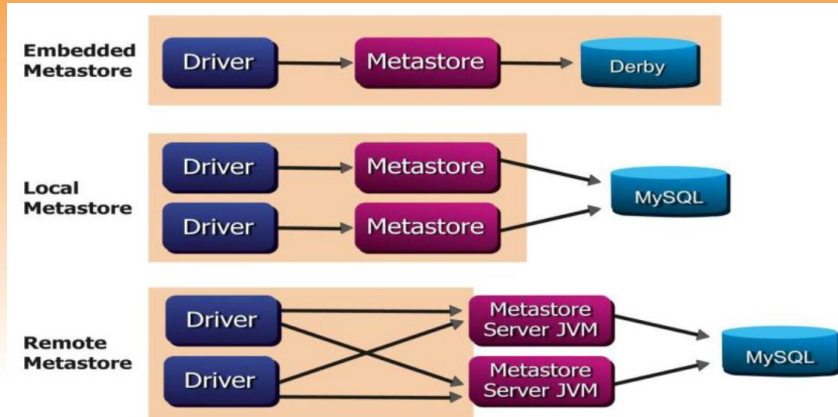
UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPLOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE



## Metastore Hive



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszca inwestycji  
Projekt nr POKL.04.01.00-00-151/12 „Intynieria wiedzy dla inteligentnego rozwoju”  
wspofinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY

MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPŁOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNIE

## Struktura danych tworzona przez Hive

Database \  
Table \  
Partition \  
Bucket



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszca inwestycji  
Projekt nr POKL.04.01.00-00-151/12 „Intynieria wiedzy dla inteligentnego rozwoju”  
wspofinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY

MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPŁOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNIE

## *Typy tabel w Hive*

- managed tabels – administrowane w całości przez Hive  
CREATE TABLE ...
- external tabels – przechowywane poza Hive  
CREATE EXTERNAL TABLE ... LOCATION '/loc'



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwentyzje  
Projekt nr POKL.04.01.00-00-131/12 „Inżynieria wiedzy dla inteligentnego rozwoju”  
współfinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPŁOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

## *Demo Hive*

WordCount w Hive



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwentyzje  
Projekt nr POKL.04.01.00-00-131/12 „Inżynieria wiedzy dla inteligentnego rozwoju”  
współfinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPŁOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

# BigSQL

Mocne wsparcie dla SQL na Hadoop:

- skalowalna architektura,
- wsparcie SQL i typy danych z SQL '92,
- wspiera sterowniki ODBC i JDBC,

Działa tylko z jedną dystrybucją Hadoop

- IBM BigInsights



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwestycji  
Projekt nr POKL.04.01.00-00-151/12 „Inżynieria wiedzy dla inteligentnego rozwoju”  
współfinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego

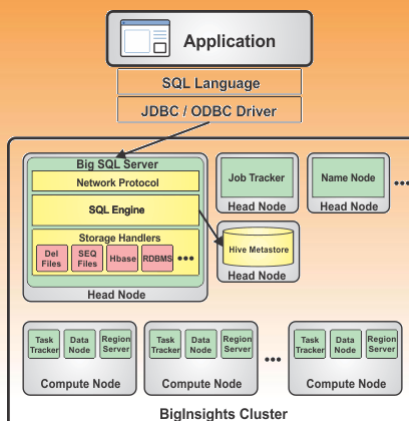


UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPLOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

# Architektura BigSQL



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwestycji  
Projekt nr POKL.04.01.00-00-151/12 „Inżynieria wiedzy dla inteligentnego rozwoju”  
współfinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPLOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

## *Problem prostych zapytań*

- MapReduce wprowadza dość duży narzut czasowy podczas swojego startu (wielokrotne start JVM),
- dla prostych zapytań czas uruchomienia MapReduce jest dłuższy niż czas samego wykonania zapytania,
- BigSQL pozwala na wykonywanie zapytań na serwerze BigSQL zamiast na klastrze,



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwestycji  
Projekt nr POKL.04.01.00-00-151/12 „Intynieria wiedzy dla inteligentnego rozwoju”  
wspofinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPŁOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

## *Wsparcie dla SQL*

- zagnieżdżone podzapytania,
- operacje łączenia tabel (joins),
- wsparcie dla wielu typów danych
  - tinyint, smallint, bigint, varchar etc.
- wiele wbudowanych funkcji
  - abs, sqrt, sin, cos, substring, upper etc.
- wsparcie dla funkcji zdefiniowanych przez użytkownika (UDF, UDTF, UDA)



KAPITAŁ LUDZKI  
NARODOWA STRATEGIA SPÓJNOŚCI



Celownik – najfiszce inwestycji  
Projekt nr POKL.04.01.00-00-151/12 „Intynieria wiedzy dla inteligentnego rozwoju”  
wspofinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



UNIA EUROPEJSKA  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPŁOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE

## *Architektura BigSQL*

Mocne wsparcie dla SQL na Hadoop:

- skalowalna architektura,
- wsparcie SQL i typy danych z SQL '92,
- wspiera sterowniki ODBC i JDBC,

Działa tylko z jedną dystrybucją Hadoop

– IBM BigInsights



**KAPITAŁ LUDZKI**  
NARODOWA STRATEGIA SPÓJNOŚCI



CeloweK – rozbudowa inteligencji  
Projekt nr POKL.04.01.00-10-151/12 „Intynieria wiedzy dla inteligentnego rozwoju”  
współfinansowany przez Unię Europejską w ramach środków Europejskiego Funduszu Społecznego



**UNIA EUROPEJSKA**  
EUROPEJSKI  
FUNDUSZ SPOŁECZNY



MATERIAŁY DYDAKTYCZNE I SZKOLENIOWE NA STUDIA PODYPŁOMOWE ORAZ NA SZKOLENIA DYSTRYBUOWANE SĄ BEZPŁATNE