# Extending RapidMiner with Data Search and Integration Capabilities
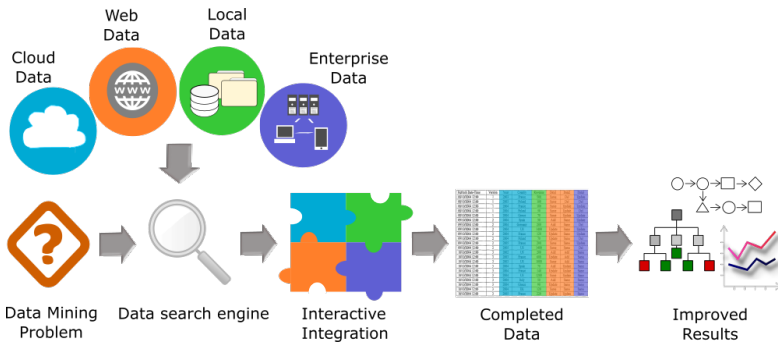
Anna Lisa Gentile[1], Sabrina Kirstein[2], Heiko Paulheim[1], and Christian Bizer[1]

[1] Data and Web Science Group, University of Mannheim
[2] RapidMiner, Germany

44

## The Vision



**PROBLEM**: data mining projects require data which exists somewhere (on the Web or in an organization's intranet) but is difficult to find.

**CURRENT SOLUTIONS**: state of the art data mining tools offer a wide range of powerful data mining methods but hardly support analysts in searching for suitable data as well as in integrating data from multiple sources.

**CONTRIBUTION**: The RapidMiner Data Search extension enables analysts to search for relevant datasets and integrate discovered data with data that they already know.

## Table Search and Table Extension

Query table

Subject column — Extension attribute

RapidMiner Data Search Extension

Extended Query Table

New attribute



Candidate tables

Schema level correspondences

Instance level correspondences



## Wikipedia Table Processing



Extracting tables from Wikipedia pages

Standardization of each table (and persistence to JSON format)

Data Indexing

Search and Join facilities via Web service

## Table Indexing

UNIVERSITÄT MANNHEIM

rapidminer

DS4DM
Data Search for Data Mining

http://ds4dm.de