# The Architecture Concepts for Building Highly Scalable Crawling Cluster For Data-Driven On-Page Optimization*

Krystian MAGDZIARZ and Damian FRĄSZCZAK

Military University of Technology, Warsaw, Poland

Correspondence should be addressed to: Krystian MAGDZIARZ; krystianmagdziarz@gmail.com

## Abstract

The paper presents the architecture concepts for building crawling clusters for data-driven on-page optimization tasks. The aim of the study is to develop a base architecture capable of continuous monitoring of on-page parameters in terms of SEO on a very large scale. The issue of building an efficient data crawling mechanism has been addressed in the literature since the spread of the Internet, however, the problem is not thoroughly described in relation to tools designed for SEO. This paper offers the concept of building a highly scalable environment designed to analyze key on-page metrics, the analysis of which will provide critical knowledge in the context of their optimization.

**Keywords:** crawling, spider cluster, content extraction, SEO, crawling cluster architecture

_____