

## Accepted Manuscript

NPIY : A Novel Partitioner for Improving MapReduce Performance

Wei Lu, Lei Chen, Liqiang Wang, Haitao Yuan, Weiwei Xing,  
Yong Yang

PII: S1045-926X(17)30241-0  
DOI: [10.1016/j.jvlc.2018.04.001](https://doi.org/10.1016/j.jvlc.2018.04.001)  
Reference: YJVLC 832



To appear in: *Journal of Visual Languages and Computing*

Received date: 2 November 2017  
Revised date: 12 April 2018  
Accepted date: 24 April 2018

Please cite this article as: Wei Lu, Lei Chen, Liqiang Wang, Haitao Yuan, Weiwei Xing, Yong Yang, NPIY : A Novel Partitioner for Improving MapReduce Performance, *Journal of Visual Languages and Computing* (2018), doi: [10.1016/j.jvlc.2018.04.001](https://doi.org/10.1016/j.jvlc.2018.04.001)

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

# NPIY : A Novel Partitioner for Improving MapReduce Performance

Wei Lu<sup>a</sup>, Lei Chen<sup>a,b,\*</sup>, Liqiang Wang<sup>b</sup>, Haitao Yuan<sup>a</sup>, Weiwei Xing<sup>a</sup>, Yong Yang<sup>a</sup>

<sup>a</sup>*School of Software Engineering, Beijing Jiaotong University, Beijing, China*

<sup>b</sup>*Department of Computer Science, University of Central Florida, Orlando, USA*

---

## Abstract

MapReduce is an effective and widely-used framework for processing large datasets in parallel over a cluster of computers. Data skew, cluster heterogeneity, and network traffic are three issues that significantly affect the performance of MapReduce applications. However, the hash-based partitioner in the native Hadoop does not consider these factors. This paper proposes a new partitioner for Yarn (Hadoop 2.6.0), namely, NPIY, which adopts an innovative parallel sampling method to distribute intermediate data. The paper makes the following major contributions: (1) NPIY mitigates data skew in MapReduce applications; (2) NPIY considers the heterogeneity of computing resources to balance the loads among Reducers; (3) NPIY reduces the network traffic in the shuffle phase by trying to retain intermediate data on those nodes running both map and reduce tasks. Compared with the native Hadoop and other popular strategies, NPIY can reduce execution time by up to 41.66% and 58.68% in homogeneous and heterogeneous clusters, respectively. We further customize NPIY for parallel image processing, and the execution time has been improved by 28.8% compared with the native Hadoop.

**Keywords:** MapReduce, Hadoop, data skew, load balance, data transmission amount, heterogeneous, parallel image processing

**2010 MSC:** 00-01, 99-00

---

\*Corresponding author

Email address: 13112084@bjtu.edu.cn (Lei Chen)

Download English Version:

<https://daneshyari.com/en/article/6934517>

Download Persian Version:

<https://daneshyari.com/article/6934517>

[Daneshyari.com](https://daneshyari.com)