

Akshay Shirahatti

INTERESTS	Information retrieval, Distributed systems, NLP & statistical methods for large datasets.	
HIGHLIGHTS	<ul style="list-style-type: none">• 5+ years of experience working in the Big Data & Information Retrieval domain• Currently working as Big Data Engineer at <i>DataSift</i>.• Worked as Data Mining Engineer at <i>Mendeley</i>.• Awarded MSc in Computer Science with <i>Distinction</i> from <i>University of Edinburgh</i>• Experience in building large scale data processing pipelines	
SKILLS	<ul style="list-style-type: none">• Scala, Java, Python(scripts), Pig(prior experience), Linux shell scripting, R(scripts)• Elasticsearch, Kafka, Hadoop, Apache Lucene, HBase, Solr, AWS, Chef, Apache Spark(Workshop)• Git, Linux Operating System, Continuous Integration.	
EXPERIENCE	Big Data Engineer DataSift Ltd, Reading, UK	2012 - present
	<i>Technology Stack:</i> Scala, Java, Kafka, Elasticsearch, Hadoop, HBase, MySQL, Akka, AWS	
	<ul style="list-style-type: none">• Prototyped and productised a <i>highly-scalable, low-latency</i> data analysis platform to handle Facebook's real-time data stream.• Developed performance/feasibility evaluation and QA tools to de-risk the pilot product.• Delivered privacy-first, anonymised APIs to help extract insights from the Facebook Data.• Helped build the archive of <i>twitter, bit.ly, tumblr etc</i> on <i>HDFS</i>.<ul style="list-style-type: none">• Archiving over 2TB/day working with a 200+ nodes <i>Hadoop</i> cluster.• Custom <i>HDFS</i> compaction workflow to satisfy the deletes processing SLA• Fast sampled insights over the archive to estimate volumes and predict cost	
	Data Mining Engineer Mendeley, London, UK	2010 - 2012
	<i>Technology Stack:</i> Java, Apache Solr, Lucene, Hadoop, HBase, Pig, MySQL, AWS, Voldemort	
	<ul style="list-style-type: none">• Scaling the Mendeley search platform and improving the relevancy ranking model for documents, users, groups in the Mendeley ecosystem.• Developed a solution(based on http://tinyurl.com/cd4klph) to aid exploratory tag based navigation of the Mendeley document corpus.• Worked on an analytic product to compute and serve top papers/journals/authors in a discipline/domain using Apache Pig.	
EDUCATION	M.S Computer Science (Distinction) University of Edinburgh, Scotland, UK	2009 - 2010

Relevant Academic Projects

- [Text Retrieval for Systematic Reviews](#)(Dissertation): To determine how *Latent Dirichlet Allocation (LDA)* compares with pseudo-relevance feedback methodologies like *Query expansion* and *Relevance-based language models* in producing the required level of precision and recall needed for Systematic Reviews. The dataset consisted of 10 million documents from PubMed corpus. [Mallet](#) toolkit was used for LDA implementation, query-document vector distance computation was done on Hadoop.
- Crawler for efficient content extraction from a set of hyper-linked news stories & also design search, relevance ranking, near-duplicate detection solutions for the corpus.

Key Modules: Text Technologies(Information Retrieval), Distributed Systems, Parallel Algorithms & Programming, Advanced Databases.

Bachelor of Engineering, Computer Science
University of Mumbai, India

2004 - 2008