# Extraction of Relevant Figures and Tables For Multi - Document Summarization

Ashish Sadh, Amit Sahu, Devesh Srivastava,
Ratna Sanyal, Sudip Sanyal

Indian Institute of Information Technology - Allahabad, India

CICLing 2012

asheesh.sadh@gmail.com

Our main contributions are:

- Incorporation of Relevant Figures and Tables in Multi-Document summary

- Relevancy Measure based on Direct and Indirect references

- Generation of Priority Lists of Figures and Tables to assist integration of the most relevant once only

Introduction

Poster Snapshot



The source code and compiled executables with an interactive interface are available at, http://goo.gl/YgoIy