# Challenges and Innovations in Establishing an ETD @ Indian Institute of Technology Kanpur, India

R. Mishra, S. K. Vijaianand[*], P. P. Noufal, Gaurav Shukla

[*]*Assistant Librarian, P. K. Kelkar Library*

*Indian Institute of Technology Kanpur, UP, India PIN-208 016.*

*vanand@iitk.ac.in*

**Abstract**

Electronic Theses and Dissertations (ETD) is one of the common subsets of any Institutional Repositories (IRs). Considering the ever-increasing importance, inevitable demand push, long term preservation and promotion of the intellectual output of IITK, we have digitized over 9000 Masters and Doctoral dissertations runnuing in to one million pages, produced between 1963 and 2005. By meaningfully customizing and manipulating DSpace, we established our ETD system with redesigned user friendly work flow having different modes of submission, including the option for online submission from 2007 onwards. In addition to the default features offered by Dspace, we have incorporated well defined useful features like Cross Community/Collection search, Browsing by Supervisor with its item strength, Citations, Cross Ref to IITK theses, linking subjects and keywords for easy navigation and redesigned feedback for analyzing the system and for further developments. This paper identifies and explains the challenges faced, innovations made at various stages of our research project with special imphasis on mass digitization of documents having complex structure and nature, Quality control, S/W configuration and customization, H/W infrastructure, Content management and handling, Metadata extraction and injection, and setting up of a full-fledged digitization facility at our premise.

## Introduction

Indian Institute of Technology Kanpur is an insititute of national importance, established in 1960 by the Government of India. The aim of the Institute is to provide under-graduate and post-graduate education and cutting edge research leading to BTech, MSc, MTech, MDes, MBA and PhD degrees, as also PDF programmes. Institute's library, known as P. K. Kelkar Library, is one the best scientific and technological libraries in India. It is fully automated with iitKLAS, an in-house library automation package developed in 1989 with avant-garde technological tools and facilities. Our collection, both print and online is very rich. IITK had started its PG and research programmes in 1963 and onwards and we have over 9000 of M.Tech and Ph. D theses in different disciplines of Science, Engineering, Humanities and Social Sciences.

## Digitization of IITK Theses and Dissertations

IITK Library has initiated its Digitization Program in mid 2005 with a multi-oriented approach and vision. At initial stage we have formulated a Digital Library Team (DLT) consisting of two IT Advisors, Project Investigator, Co-Project Investigator, two Metadata/ Quality Analysts and two more Project Assistants. The team was assigned the responsibility to prepare a detailed project proposal and to submit it to the institute authority. Following is the roadmap of our digitization programme:

Phase 1          Electronic Theses and Dissertations (ETD)

| Phase 2 | Faculty/Academic staff Publications (Conf./Seminar proc., Journal articles, Project/ Technical Reports, lecture notes, delivered lectures/speeches, cover and contents pages of books, etc.) |
|---|---|
| Phase 3 | Conference/Seminar/Workshop proceedings (Organized & Published by IITK) |
| Phase 4 | Inter-Institutional Aggregation of premier S&T Institutes Student Academic Portfolio System Minutes of IITK Senate Meetings (Intranet) |

This paper delas with the first phase of digitization of over 9000 MTech and PhD theses. Digitization of theses has proved to be the most stupendous task; nonetheless a challenging one! Scanning of 900,000 pages with graphs, images, charts and programming codes with pale typed papers dating back to 1963 was really a challenging and painstaking job. We have availed the scanning facility available at Indian Institute of Information Technology Allahabad (IIITA), a Govt. of India, Mega Scanning Centre, Million Book Project (MBP) of Universal Digital Library (UDL). IITK signed an MoU with IIITA on 15 July, 2005 for scanning the existing theses. The scanning of over 9000 theses was successfully completed well within a twelve months time. Now we have developed the digitization facility in the library; well equipped with Minolta PS7000 scanner and required accessories.

## Scanning Specifications:

After careful consideration, the following decisions with regard to scanning specifications were undertaken:

- TIFF for archival and PDF for presentation purpose.

- Resolution of Image: 600 dpi TIFF and PTIFF (Processed). PTIFF has been subjected again for scan fixing to resize it to a uniform shape and size.

- Compression: CCITT 4 Fax

- Conversion to PDF format using Adobe Acrobat: PDF is the de facto standard for the secure and reliable distribution and exchange of electronic documents and forms around the world.

- Scanning software: Abbyy fine reader v.8.0

## Quality Control

Quality Checking is no doubt one of the essential processes in digitization to ensure the quality output and to get the most reliable and consistent data. Even though the structure, paper, printing quality, and complex nature of theses with graphs, figures, charts, tables, mathematical and scientific formulae and the fact that many of theses were as old as 30 to 40 years, made the job of scanning a challenging one for the scanning technicians. Nevertheless, we showed our concern for standard and high quality scanning.

We approached this challenging assignment with a quality oriented approach and established predefined guidelines for rectifying problems by keeping in mind the following elements, i.e. completeness, contrast, sharpness, skew, resolution, bit depth level of compression and data conversion. It was intriguing that major part of time and effort of our quality analysts was consumed by the problems such as missing, duplicate and misplaced pages, data conversion

and file naming. Few other major shortcomings were negligence of proper scan fix setting and auto cropping which is an automatic process to eliminate unwanted junks and discrepancies from scanned documents so as to get clean and legible output. But this automatic process again resulted into dim view of the data. Consequently, we were constrained to recommend for manual cropping as the only solution. We feel it would have been better if the quality checking was taken immediately after the raw scanning and processing to avoid defects that crept in as detailed above. Finally, having put in tremendous efforts and time, we completed the quality checking and rectified the problematic items as late as April, 2006.

## Setting up of Infrastructure for ETD @ IITK

Keeping in mind the features of an outstanding ETD system, our DL team actively involved and conducted brain storming sessions in following areas, right from proposal preparation down to ETD launching are:

- IT Infrastructure

- Selection of hardware and software

- Configuration and customization

- Selection of Metadata Standard

- Designing standard workflow pattern

- Content Management

- Metadata extraction and injection

- Uploading: Mode of Submissions

### IT Infrastructure @IITK

The Computer-Centre has about 100-150 Linux terminals and more than 100 Windows-NT terminals supported by PARAM 10000 super computer and well acquainted with modern application softwares IITK has wide 100 Mbps fiber optic network that connects to all the academic departments, hostels, library and other central facilities to the Computer Centre. Internet access is provided 24 hours a day, 365 days an year through VSNL 34 Mbps and ERNET 2 MB dedicated Internet link.

Our Library has got more than 100 network points, 50 PCs, 10 terminals for OPAC search, 10 latest work stations   for accessing various E- resources, a dedicated web server, a Network CD Server and many other new computing and peripheral technologies.

## Selection of hardware and software

**Hardware**

At the initial stage, we have used a test bed server for installation and customization. Later on we procured a high end server with the following specifications:

HP Proliant DL 385 / Dual Processor Capable / AMD 8000 Series chipset / Op 2.4 GHz with 1024KB L2 Cache / 4 GB RAM, 2x 73GB Ultra320 10k Hp Universal, 3x 300 GB Ultra320 10k Hp Universal, DVD+RW Drive, Hot plug Redundant power supply, 17" TFT.

3x300GB Segate external USB Hard disks are used for Data Transfer from IIITA to IITK.

Backup and security of the system will be taken care of by our Computer Center, a central facility of IITK.

As part of our exercise to establish a digitization facility in our library, we have procured a MINOLTA PS7000 scanner with requisite accessories for further digitization and development work.

**Software**

We decided to choose DSpace, an Open Source Software from MIT and HP due to its features like granularity, adherence to standard, multi-format support, customizable interface, OAI-PMH compliant, support with fully qualified DCMI, remote submission, authorization and reviewing, community/sub-community based collection architecture, import and export features, Persistent Identifiers Handle System, Open URL Support, Lucene search engine and generation of statistical reports, etc.

## Configuration and customization

After detailed systems study and groundwork, we successfully installed and configured DSpace (Version: 1.2.2) on Linux platform on   19 July 2005 for our test bed. Later, we have acquired a dedicated high end server with above mentioned configuration. Now it is running on latest version of DSpace (Version.1.4) with all prerequisites.

## Selection of Metadata Standard

There are various well-known metadata standards available. Dublin Core Elements (DCMI) is one of the wide accepted standards among library fraternity. Its usability, flexibility, repeatable elements and qualifying nature are widely appreciable.

We felt qualified DCMI is the right option for our ETD system due to its usability, flexibility, repeatable elements, qualifying nature and wide popularity.

## Designing standard workflow pattern

We meticulously customized the default workflow pattern of DSpace including the provision for various mode of submissions depicted in the figure given below.
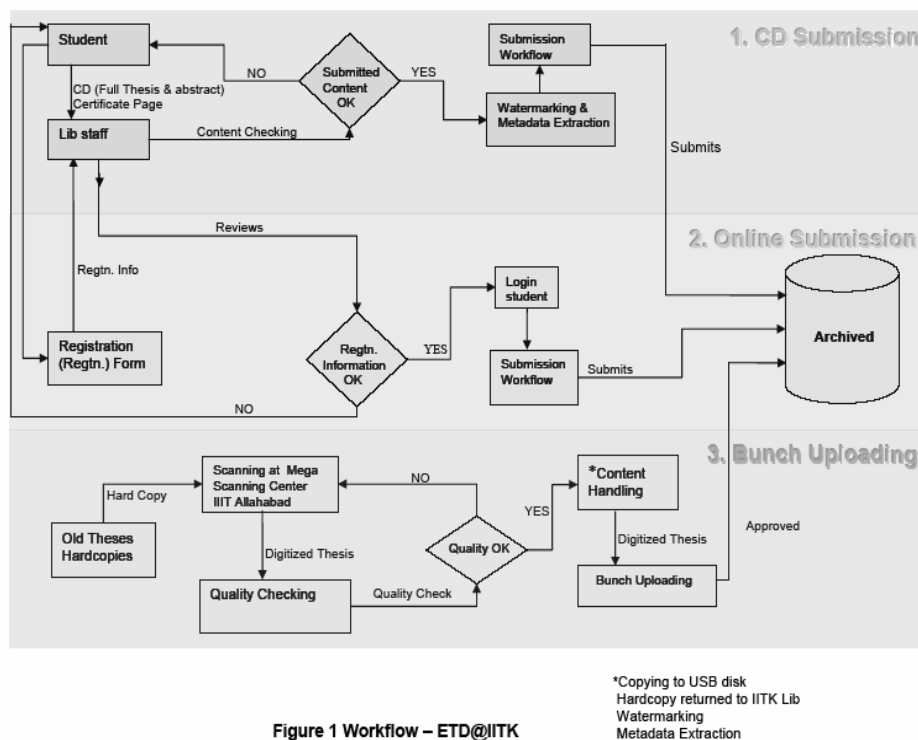
Figure 1 Workflow – ETD@IITK

Content Management

After getting quality checked and rectified data we have done the following activities to make digitized versions more meaningful and secure:

**Scanning & Insertion of Signed Certificate page**

During the current academic year, 2005-06 and onwards, we used to insert certificate page signed by the thesis supervisor(s) to make the item more meaningful and authenticity.

**Embedded Watermarking**

Considering the possibilities of potential infringements of the Copyright, for all submitted theses, we decided to embed IITK Logo in each page of thesis as a water mark. We have developed a script for automatic embedding and Bunch processing.

We are using a tool iText Library freely available at http://www.lowagie.com/iText/. It comes with a utility to embed watermark in Pdf text file. But in our case, we have manipulated it to embed in PDF images.

## Metadata extraction and injection

We have developed a Perl Script for Extracting Metadata from the existing our library thesis database running on Oracle and injected it in our new system, powered by DSpace.

A Java program written by us, connects to our current oracle database and extracts all the metadata for each item separated by a specifier line by line. This program creates a file **"MetaDataDetails"** and writes to it.

## Uploading: Mode of Submissions

According to the modified workflow, the following are the three various mode of submission of theses to our ETD.

**Submission Mode 1: Bunch processing**

By using the default export/import utility of DSpace, the metadata for MTech/MDes and PhD collection existing at that point of time was extracted from the database of the library. The full text digitized data was uploaded to the server by way of bunch processing.

The perl script written reads the metadata from file **"MetaDataDetails"** and creates a directory structure for each collection with a "**dublin_core.xml"** files containing all the metadata with corresponding DC Elements. This script also stores the files to be uploaded after being watermarked.

**Submission Mode 2: CD**

As a current practice, scholars are submitting a CD containing e-theses to the library, our team is manually uploading to the DSpace after the completion of essential content processing like quality assurance, file naming and related content management processes.

**Submission Mode 3: Online (Proposed)**

Our ETD system is ready for online submission by the scholars from their workplace to our server. Once the approval of the authorities is there, Online theses submission to the system can start.

# Incorporated features @ our ETD

It was observed that the default features of Dspace were not adequate inasmuch as it lacked certain essential features, viz cross collection search, browse by supervisor/s, linking to keyword, subject, crossref, homepage, item strength for supervisor, subject and citation. It motivated us to customize Dspace to accommodate for these new features to our ETD. We are sure that these incorporated features are useful for other ETD system powered by DSpace to make them more meaningful and user friendly.
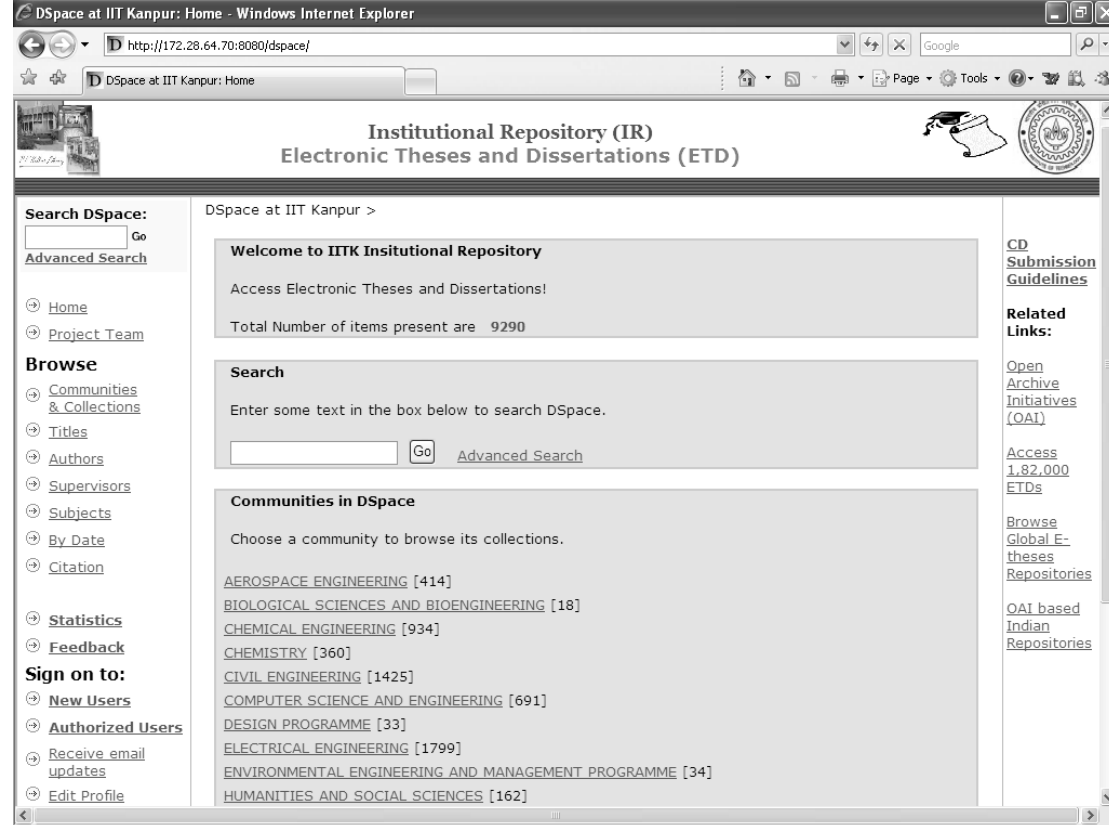
Figure 2: Home page of ETD

## Browsing

DSpace permits by default the following access options: (a) by Communities and Collections, (b) by Authors, (c) by Titles, (d) by Issue Date and (e) by New Collection and Recent Submission and (f) by Subjects.

In addition to the above noted default 'browse' options, we have customized DSpace to incorporate browsing facility by Supervisor/s and Citations.

The main files which are associated for browsing are "**Browse.java**", "**BrowseServlet.java**" . By looking carefully into it one can add the other browsing facility in the same manner as browsing Authors and Subjects are created.

The files we created for browsing by Supervisors are "ItemsBySupervisorServlet.java", "supervisors.jsp", "items_by_supervisor.jsp". Additional tables and sequences are also created ItemsBySupervisor, CommunityItemBySupervisor, CollectionItemBySupervisor and itemsbysupervisor_seq.

The additional servlet and files created for browsing by Citations/References are "**ItemsByReferenceServlet.java**" , "**references.jsp" and "items_by_reference.jsp".** In this regard we have created the following table's **itemsByReference**, **CommunityItemByReference, CollectioItemByReference**, and **sequence itemsbysupervisor_seq.**

## Searching

DSpace offers by default the following search features: (1) Search all DSpace, (2) Bounded Search within a specified Community's Collection, (3) Simple search and (4) Advanced search.

**Cross collection Search**

Search involving more than one disciplines, known as 'cross collection search', viz. Chemistry and Chemical Engineering; Materials Science and Physics; Metallurgy &Materials Science; Lasers and Biomedicine, etc. have been incorporated as additional features under the 'Advanced search' option.

Our collections are uniquely defined by the combination of department and degree type. e.g. M.Tech Thesis @ CSE (Computer Science and Engineering) , Ph.D Thesis @ AE (Aerospace Engineering), etc. where M.Tech or Ph.D is the degree type and CSE or AE is the department. The same idea is implemented on cross collection search. For each item two DC elements **description.department** and **description.degree-type** uniquely defines a collection for cross collection searching. We have done some changes in the following files "**advanced.jsp**" as a display file and "**QueryArgs.java**" **as** a core java which takes care of structuring the query.
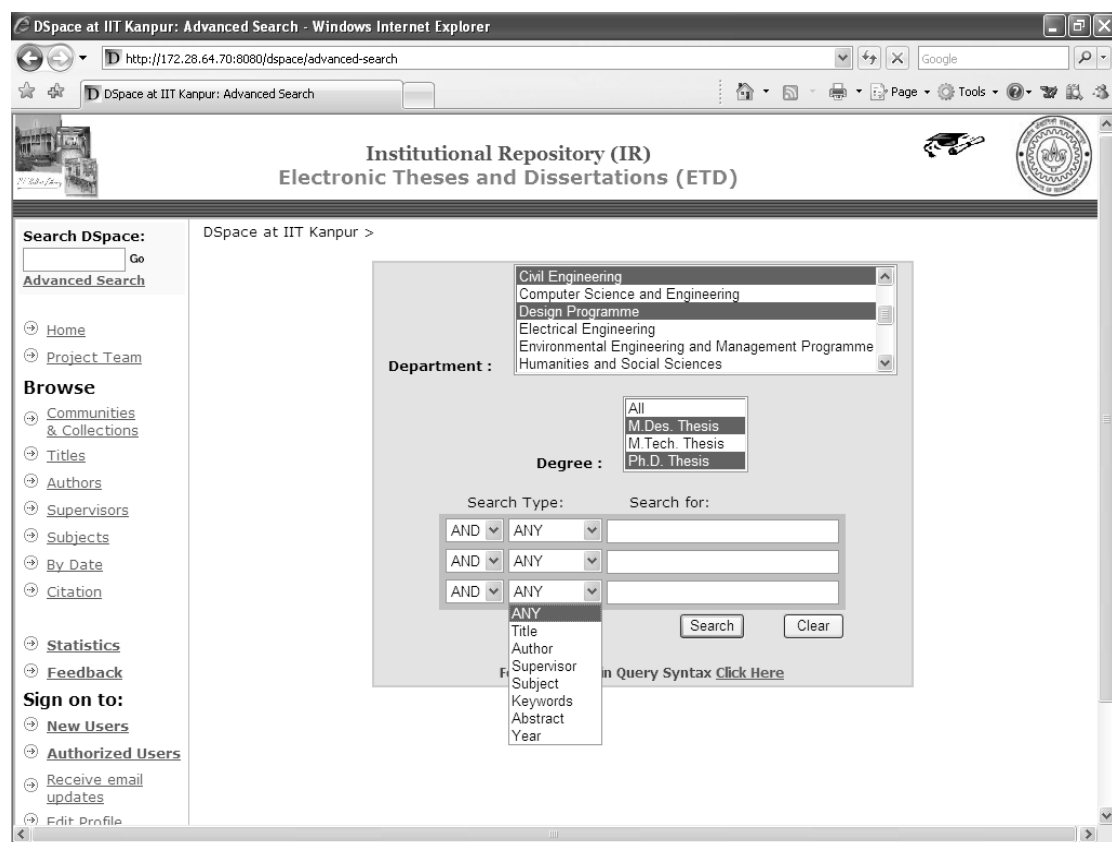


Figure 3: Cross Collection Search

## Supervisor, Author, Subject and Reference count

The Strength/Count of each Browsing element are one of the interesting options. We have manipulated the file "**Browse.java**" This option will give you a picture about the number of item guided by a person, item created by an author, items related to a particular subject and the no. of citations for a particular item.
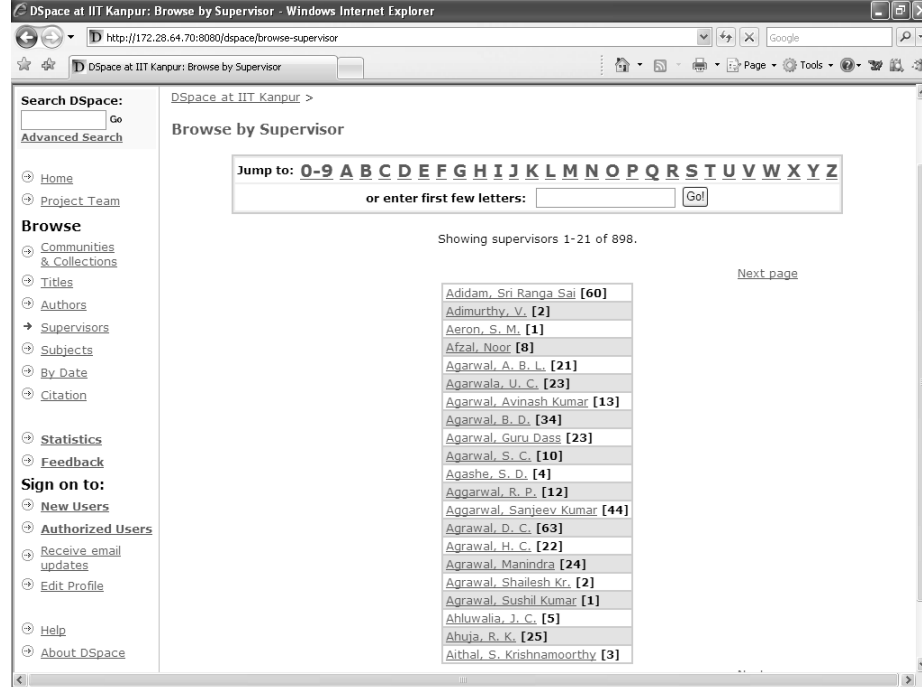
Figure 4: Supervisor Strength/Count

## Linking to Supervisor home page

Another incorporated and useful feature is linking to supervisor home page with his photograph. It will be helpful for others who wish to know him in details for communicating and collaborative works in future.
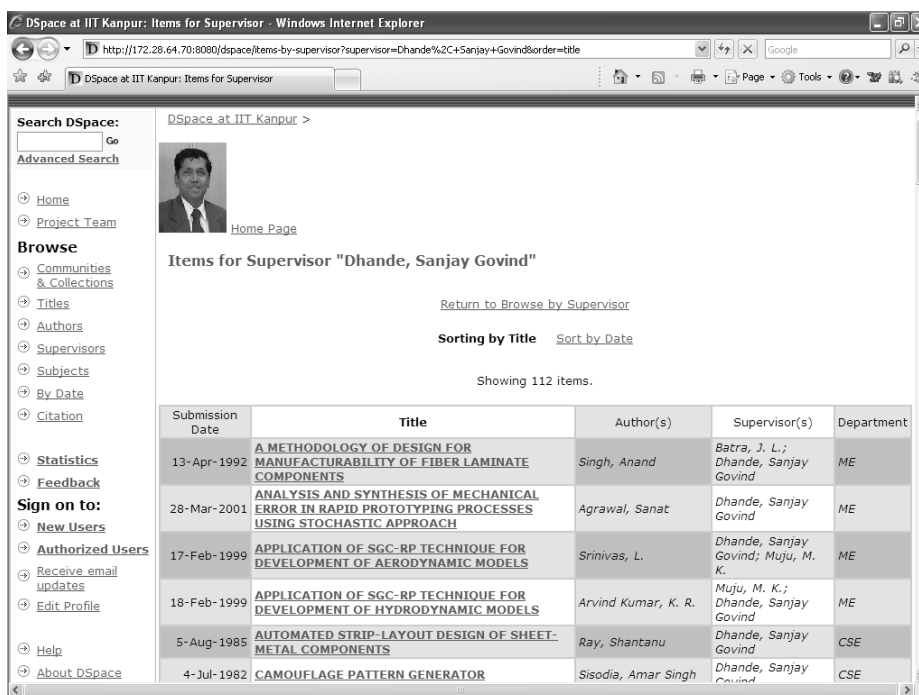


Figure 5. Supervisor photo identity and hompage Linking

## Keyword and Reference Linking

'Subject Keyword' is a common user's approach to find literature in a given subject area. An access has been provided to the various 'Subject Keywords' incorporated in each thesis, through a hyperlink.

A link to 'cited thesis/theses' from the existing IITK theses collection, as mentioned by author in 'References' showing citation impact of our theses by the IITK scholars, has been provided for. For as an additional feature. Further, it is proposed to send an e-mail alert to the author/s and supervisors in case their theses are cited by any new scholars at IITK.

The java core file we changed a "**ItemTag.java**" which takes care of rendering of Items dc elements for this purpose.
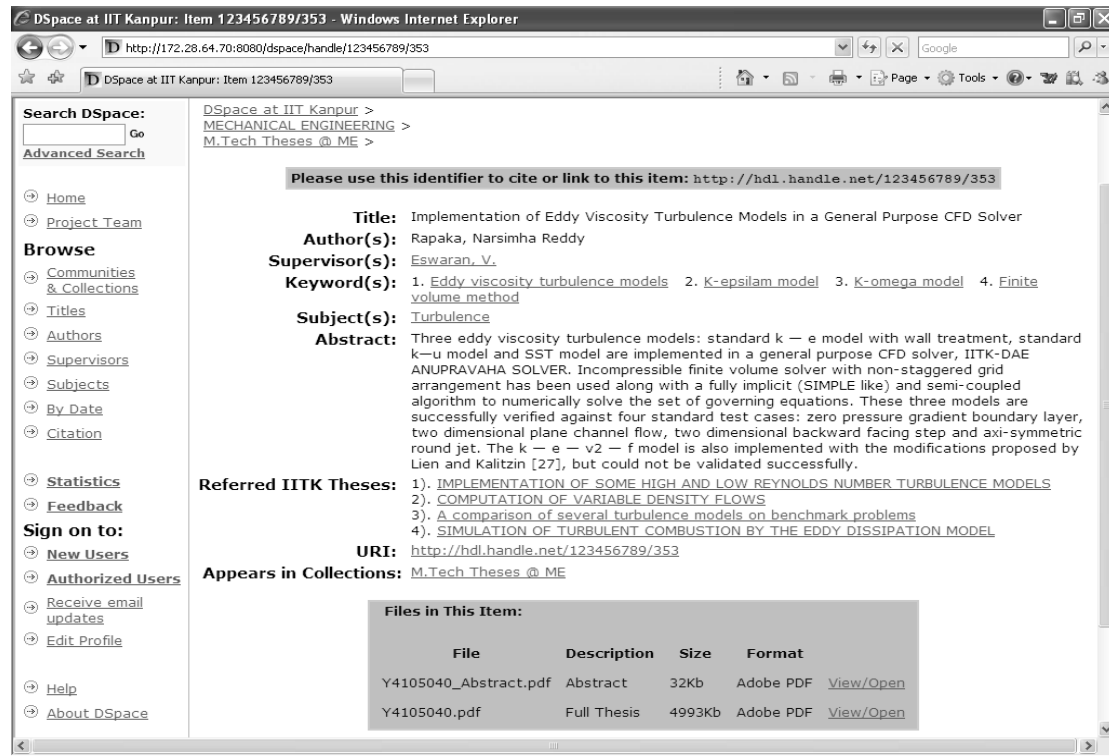


Figure 6: Keyword and Reference Linking

## Redesigning 'Registration Form'

Normally, in DSpace, the registration form contains the following details: First Name, Last Name and Contact No. We redesigned this by providing additional options, i.e. Roll No., Degree and Department. These data elements are very essential for DSpace administrator to assign a scholar to a particular E-Group and for user to submit his/her item to collection without any problem.

File associated with redesigning registration form are **"Eperson.java"** and altered the table eperson and have added the above mentioned additional fields.

## Redesigning 'Feedback Form'

We redesigned the default Feedback Form, incorporating additional personal details of scholar, i.e. email-id, designation, degree, department, PF No./Roll No., and a set of features related to the facilities, and/or description parameters concerning content/metadata and workflow that have been used to solicit ratings on defined parameters from the scholars submitting theses to the library.

Earlier this feedback used to be sent through email to the administrator but now we have created a specific table in the database as Feedback to store the values. Analysis of theses values is helpful for assessment of the ETD.
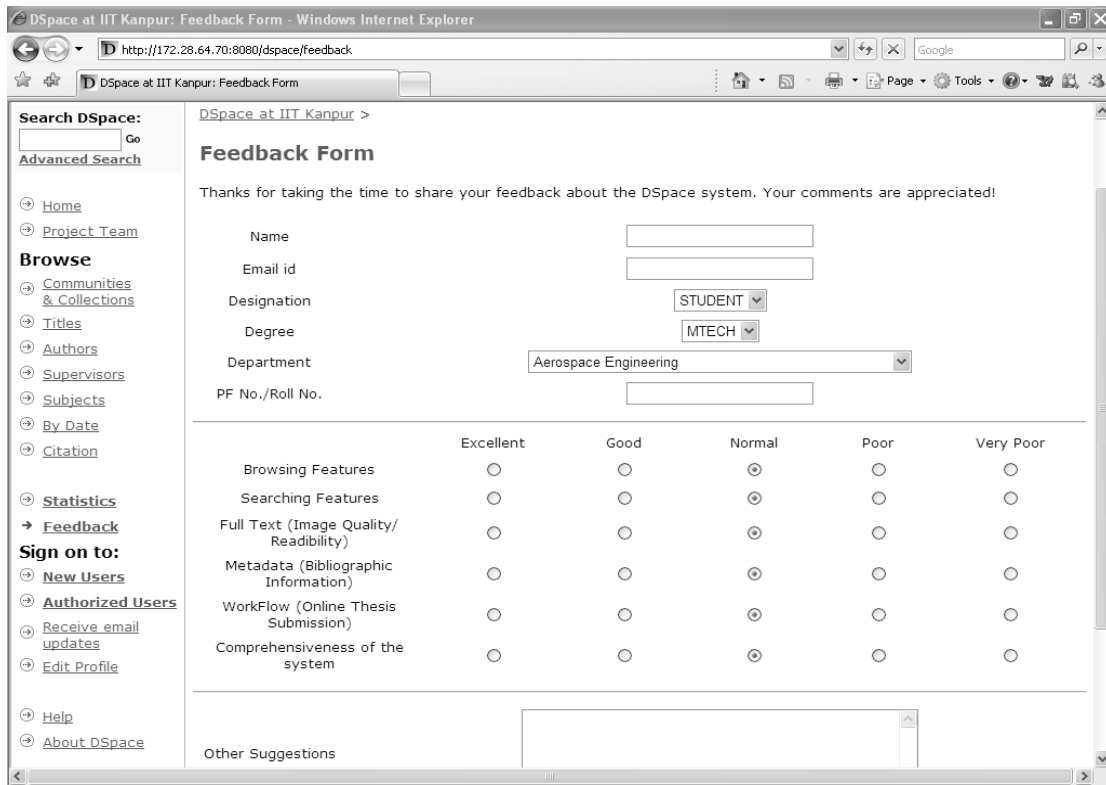


Figure 7: Feedback form

# Conclusion

Establishing an ETD with the mass digitized data and incorporating and enhancing additional features to the system to make it more user-centered and user friendly by customizing open source software, DSpace has really been an interesting and innovative task for us. Incorporated advanced features like workflow for different modes of submission, various browsing and searching options, keyword and reference linking, Citation analysis, redesigned feedback , photo identity and linking to the home page of supervisor will be beneficial for others who are using DSpace for their ETD system. We record with a sense of pride that once our ETD is put on the web, it would be the largest OAI compliant ETD repository in India and one of the top ten in the world. We propose to improve our ETD system incorporating state-of-the-art developments as they take place from time to time. An effort for harvesting such ETDs from the institutions of higher learing, both from India and abroad is also on the cards subject to the inherant constrainsts of inter-operability of these databases.

# Acknowledgments

# References

[1] Adobe

http://www.adobe.com/

[2] Digital Library of India

http://dli.iiit.ac.in/

[3] DSpace Federations

http://www.dspace.org/

[4] Dublin Core Metadata Initiative (DCMI)

http://www.dublincore.org/

[5] Mishra, R., Vijaianand S. K., Noufal P. P., Rajesh Kumar and Shukla, Gaurav., Digitisation Initiatives to Destress Library Collection: A case study of ETD at P. K. Kelkar Library, IIT Kanpur. International Conference on Digital Libraries, New Delhi, 2006.

[6] Mishra, R., Vijaianand S. K., Noufal P. P. and Shukla, Gaurav., (2006). Development of ETD at IITK Library using DSpace: Practical Exposures and Experiences International Conference on Semantic Web and Digital Libraries, Bangalore, 2007.

[7] Open Archives Initiatives

http://www.openarchives.org/

[8] PDFbox

http://www.pdfbox.org/

[9] Registry of Open Access Repositories (ROAR)

http://archives.eprints.org/index.php

[10] UNESCO

Guide to Electronic Theses and Dissertations

http://www.etdguide.org