



distributed between instances of the Cheshire3 framework. However this does not, yet, prevent the necessity to share lines of user configured code in order to process localised requirements, which is a requirement for successful grid based information retrieval. The objects include:

<i>DocumentGroup</i>	A set of Documents.
<i>Document</i>	Unparsed data representing a single item.
<i>Record</i>	Parsed XML based data representing a single item.
<i>Query</i>	A CQL query parse tree.
<i>ResultSet</i>	An ordered set of symbolic pointers to Records.
<i>Index</i>	An ordered list of terms extracted from a Record
<i>User</i>	An authenticated user of the system.

Storage facilities exist for each of these object classes. In addition there are several main processing objects including:

<i>PreParser</i>	Convert a Document into another type of Document
<i>Parser</i>	Convert a Document into a Record
<i>Transformer</i>	Convert a Record into a Document
<i>Extractor</i>	Extract data of a given format or type for indexing
<i>Normaliser</i>	Convert data from one format or type to another
<i>ProtocolHandler</i>	Take a request in a known protocol and converts it to an internal representation.

Also, there are three abstract objects:

<i>Server</i>	A logical collection of databases
<i>Database</i>	A logical collection of records and indexes
<i>Workflow</i>	An object that can take input, go through a user-defined sequence of processing steps and produce output.

The last abstract object above is the component that allows the Cheshire3 system to work effectively in a distributed environment. Workflow objects are configured like any other object within the system, and hence share their ease of portability. Each workflow contains a series of instructions in its configuration, and when the object is requested, these are compiled dynamically into executable code. Thus, the requirements for the system operations are specified by describing the logical flow of data through the various types of processing objects, and no additional programming is needed. Specification of the workflow is done via object identifiers, and workflows can call other workflows as required, allowing users to set up a workflow for specific tasks but still have that subsystem integrated within a larger processing environment.

Each object in the system has a (usually unique) identifier as its reference. It is possible, however, to re-use identifiers at a lower levels of the system, allowing a database to override a particular server object without requiring changes to code or configurations referring to it. The utility of this becomes apparent in the context of distributed workflows, as not only can one workflow to call another, one very high level workflow might interact in the same way with many different database-specific workflows, each of which performs customised operations to achieve the same end result.

Once the distributed infrastructure has been defined, one or more 'master' workflows divide the processing requirements among 'slave' processes. Our implementation uses PVM as a very fast communication layer over a gigabit switched network, but we are implementing an MPI based version as well. When a slave process has completed its assigned workflow, it returns the results to the master to be

merged. Object interactions within the architecture are well defined, the result from any workflow is also well defined, based on the last processing object. This is very important for trees or graphs of distributed workflows.

Because each database is a logical collection of records, it can be easily split across many nodes or combined at a single location. This means that each node on the cluster can either look after a slice of the database or do the processing required and then return the record for central storage. The same applies to indexing to a lesser extent; if the index files at each node store all of the terms for the records at the node, then interacting with the sub database is very fast, but the intermediate resultsets will need to be merged before a global query can be answered. If each node maintains a portion of the terms for all of the records, then it can answer the query authoritatively, however this requires some central authority to manage the division of the index terms.

### 3. PERFORMANCE AND CONCLUSIONS

Early performance tests used 1 'master' and 8 slave indexing processes (in a cluster of 16 a dual 1.6Ghz Athlon machines with 2Gb of RAM, with gigabit ethernet). Using this configuration on a workflow that parses MARC records, converts them to XML, parses the XML, stores it in an RecordStore and indexes it, we are able to achieve a sustained indexing rate of about 10,000 records per second (even when competing with other tasks on the cluster). For larger, more complex, records we tested with 640 Mb of TEI encoded documents, and were able to parse, store and index the entire collection in 3 minutes and 20 seconds.

We are in the process of integrating Cheshire3 with the Storage Resource Broker (SRB)[3] DataGrid storage system. Because each object in the SRB is a "document" that can be processed and indexed and identified by its unique SRB ID, we can retrieve the original from storage on demand without needing to store an XML representation of it. Since SRB is also being integrated into the DSpace Digital Library Framework, we will soon have an advanced indexing and search systems for DSpace and SRB installations.

We believe that current digital library architectures and standards can be added to the existing "Grid services" to provide enhanced performance, collaboration and functionality across the developing computational Grid.

### 4. ACKNOWLEDGMENTS

Design of the Cheshire3 system was supported by the NSF and JISC (U.K) under the *International Digital Libraries Program* award #IIS-9975164.

### 5. REFERENCES

- [1] I. Foster and C. Kesselman. *The Grid: Blueprint for a New Computing Infrastructure*. Elsevier, Amsterdam, Second edition, 2004.
- [2] R. R. Larson, J. McDonough, P. O'Leary, L. Kuntz, and R. Moon. Cheshire II: Designing a next-generation online catalog. *Journal of the American Society for Information Science*, 47(7):555-567, July 1996.
- [3] A. Rajasekar, M. Wan, R. Moore, W. Schroeder, G. Kremenek, A. Jagatheesan, C. Cowart, B. Zhu, S.-Y. Chen, and R. Olschanowsky. Storage resource broker - managing distributed data in a grid. *Computer Society of India Journal*, 33(4):42-54, 2003.