# An Architecture for Logging Text and Searching Chat Messages
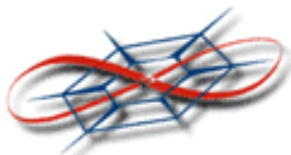
Master's Thesis Defense

Rajan Vijayaraghavan

04.28.2003

Committee:

Dr. Susan Gauch (Chair)

Dr. Arvin Agah

Dr. Joseph Evans

I T T C

# Organization

➢ Motivation

➢ Objectives

➢ Architecture

➢ Testing

➢ Evaluation

➢ Future Work

# Motivation

➢ Instant Messaging
   - ➢ Chat Rooms
   - ➢ Instant Messengers
➢ Growing Popularity
   - ➢ People of nearly different age groups participate
   - ➢ All conceivable topics discussed
➢ Reach of Instant Messenger Services
   - ➢ Corporate users
   - ➢ Home users
   - ➢ US Marines

# Motivation (cont'd)

➢ So what is the problem ?

  ➢ Parental Controls

  ➢ Need for message archiving

  ➢ Search Feature

  ➢ Scalability Issues

# Related Works

➢ Chat Track Project

➢ Commercial Systems (Archiving capability)
- ➢ Iambigbrother (Keyword search)
- ➢ Net Nanny
- ➢ Cyber-Scoop
- ➢ Desktop Snooper
- ➢ I-Spy Now

# Objective

➤ To create a system that can log text from Microsoft MSN Messenger

➤ Index the text frequently

➤ Support variety of queries

# Architecture

➢ Message Logging system for MSN Messenger
➢ Indexing System
➢ Retrieval System

# Message Logging System

- Stand Alone Application
- Actions
  - Identify new/existing conversation window
  - Check for new message
  - Log both incoming and outgoing messages
  - Log the sign-in name of the speaker of the message
  - Keep track of entering/leaving users of a conversation window
  - Log administrative messages
  - Detect Window Closing

# Message Logging System (cont'd)

➢ Conversation window monitoring based on Window Handles
  ➢ Each window will have a unique window handle
  ➢ Any MSN messenger window will have " – Instant Message" or
    " – Conversation " in the title bar

# Message Logging System (cont'd)

➢ Logs messages approximately every 1 second

➢ Writes message file, speaker id file and listener id file

➢ Adding messages is an incremental process

If (current message length > previous message length)

    If Administrative Message

      Log the event;

    else if User Message

      Log Message, Speaker Id, Listener Id(s), Add File name to "to be indexed list";

    else

      continue;

➢ Monitors administrative messages

    ➢ Logs them as said by user 'none'

    ➢ This helps in playing back the conversation as it happened

# Message Logging System (cont'd)

Example.,

Assume Window Handle 41157726; date 04.26.2003; time 15:17 PM

# Message Logging System (cont'd)

Message File : session000041157726/2003/04/26/151754866.mesg

Speaker Id File : session000041157726/2003/04/26/151754866.uid

Listener Id File :

     session000041157726/2003/04/26/151754866.receiver

File Contents:

Message File: "there's a concept that works 20 million other white rappers emerge but no matter how many fish in the sea it'd be so empty without me"

Speaker Id File: "rharishnandan"

Listener Id File : "rharishnandan"

File to Index list: session000041157726 2003 04 26 151754866

# Message Logging System (cont'd)

➢ Current system monitors 100 conversation windows

    ➢ Practical number might be around 10 conversation windows

➢ Calls the indexing program periodically

        If (Files to Index)

           Call Indexer;

        else

           continue;

# Indexing System

➤ Indexing

   ➤ Creation of Inverted Index from the raw files.

      - Dictionary File

        - Contains one record for each unique word in the collection

          Words, number of documents in which the word is present, total frequency of the word in the collection, inverse document frequency, pointer to postings record

      - Postings File

        - Each record Information contains information about the word occurrence in a document.

          Frequency of the word in a document, weight of the word in a document, document id for the document, pointer to next postings record.

➤ Example

   Word 'prozac' occurs in 100 documents and total number of documents 200 and occurs 120 times in total.

   One dictionary record { 'prozac', 100,120, idf = $\log_2(200/100)$, pointer to postings }

   Postings File { 100 postings record ( 100 documents),

        weight in document j = idfi * frequency of 'prozac' in document j;

        document id = unique for a document;

        pointer to next postings record; }

# Indexing System (cont'd)

- ➤ Indexing
  - ➤ Batch Indexing
  - ➤ Incremental Indexing

- ➤ Creates Two sets of Inverted Index
  - ➤ One for Keyword ( kyDictionary, kyPostings and kyDocuments)
  - ➤ One for Speaker Id ( spDictionary, spPostings and spDocuments)

# Retrieve Application

# Scenario 1: "keyword" query

# Scenario 2: "key word + user name" query

# Scenario 3: "key word + date query" specifying "from" option

# Scenario 4: "key word + date query" specifying "to" option

# Scenario 5: "key word + date query" specifying "from" & "to" options

# Scenario 6: "key word + Speaker Id + date sorted"

# Scenario 7: "key word + date + Speaker id sorted"

# Scenario 8: "key word + Speaker id sorted"

# Scenario 9: "key word + date sorted"

# Scenario 10: "key word + two results per page option"

# Scenario 11: "key word + 4 lines before/after match"

# Scenario 12: "key word + show who listened"

# Summary of queries supported

- Keywords
- Keyword + Speaker Id
- Keyword + Date based filtering of results
- Keyword + Speaker Id + Date based filtering of results
- Keyword + Speaker Id + Date based sorting of results
- Keyword + Date based Filtering + sorting results by speaker id
- Keyword + sorting results based on Speaker id
- Keyword + sorting results based on dates
- Varying number of results per page
- Showing text around the exact document match
- Showing who (all) listened

# Latest Additions

➢ *Queries based on Dates, Returns session ids*

➢ *Queries based on Listener Id, Returns session ids*

➢ *Displaying complete session as it happened*

# Screen Shot Showing the complete session

# Scalability Issues

➢ Issues

    ➢ Number of Files

        - Assuming a user chats 1 hour a day, creating 12 message files/minute

          In a day, 720 message files. (Total files = 2,160)
          In a month, 21,600 message files. ( Total files = 64,800)
          In a year, the total will be 259,200 files (Total files = 777,600)
          Windows allows up $2^{32}-1$ Files = 4,294,967,295 files.

    ➢ Disk usage for files

        - Assuming each message file with 150 Bytes of data

          39 MB of Message data + 7 MB for storing user information per year

    ➢ Size of inverted index
    ➢ Time for retrieve

# Key Word Inverted Index File Growth



**Keyword Inverted Index File Growth**

# Speaker Id Inverted Index File Growth



**Speaker Id Inverted Index File Growth**

Legend:
- Speaker Id Dictionary
- Speaker Id Postings
- Speaker Id Documents

Y-axis: File Size (in KB) — 0, 500, 1000, 1500, 2000, 2500, 3000, 3500, 4000, 4500

X-axis: Number of Files — 10, 100, 1000, 5000, 10000, 25000, 50000

# Number of Files Vs Indexing Time

# Number of Files Vs Retrieval Time



**Number of Files Vs Retrieve Time**

# System Features

➢ Environments

  - Windows NT

  - Windows 2000 Professional

➢ MSN Messenger Versions

  - 5.0 and earlier

# Future work

➢ XML

  ➢ Able to be used by chat servers open architecture

➢ Building User Profiles

➢ Topic Detection and Tracking

➢ Segmentation

➢ Text Summarization

# Conclusions

➢ An architecture for logging text, index and retrieve information in an efficient way

➢ First academic system to do research with chat data

➢ Shares code with server-based version

# Special Acknowledgements

- ➢ The 5 Ss

# Questions ??