

Express Mailing Label No. EV 343632189 US

Provisional Patent Application
Docket No. 3412.2.3p

UNITED STATES PROVISIONAL PATENT APPLICATION

of

G. EDWARD POWELL, JR. and MICHAEL ANDERER

for

SYSTEM FOR SOFTWARE SOURCE CODE COMPARISON

SYSTEM FOR SOFTWARE SOURCE CODE COMPARISON

DESCRIPTION

- [01] Computers and computer technology have become an integral part of today's society. Computer software running on computers enables the computers to accomplish many different tasks. Computer software includes any kind of computer program that may be run on a computer.
- [02] Software engineers develop software by writing source code that is then transformed into a computer program or into part of a computer program. Source code is typically a type of document that people can read and understand. When it is compiled and/or linked it is turned into code that is understood by the computer so that it may be executed.
- [03] Many software engineers change jobs in their lifetimes. Their employers may be concerned with whether or not they take source code from their current job to their new job. The development of software often costs a company a substantial amount of resources. Many companies would like to protect their investments in the software they have developed.
- [04] Intellectual property law may be used to protect many aspects of computer software.
- [05] The open source movement has made it very easy for source code written by one person or company to become part of another piece of software. Open source software is freely available and can be downloaded from the Internet. Employees, with or without their employers permission, may submit source code to open source projects.
- [06] It would be beneficial if systems and methods were provided for a person or company to determine whether they are using source code from another source. It would also be beneficial if systems and methods were provided for a person or company to determine whether someone else (e.g., an open source project) is using their proprietary source code.
- [07] Figure 1 is a block diagram illustrating a system for comparing source code. Software source code A is taken as input by a profiling block which performs profiling on the source code to produce a metadata file A. The profiling tool may perform multi-source characterizing and a one-way transform. By making the transform a one-way transform the system will protect the proprietary nature of the source code. If the source code could be reverse engineered from the

metadata file, very few companies with proprietary source code would want to use the system for fear of disclosing their source code to others. By making it a one-way transform, they may be comfortable that their confidential information and source code will be kept confidential.

[08] Software source code B is taken as input by a profiling block which performs profiling on the source code to produce a metadata file B. The metadata files are then compared to one another and a report is generated. The report will reflect how closely the two pieces of software resembled one another.

[09] Figure 2 is a block diagram illustrating a system for managing a database of source code and a database of metadata. In one embodiment this system may be on the Internet so that it can be accessed by third parties. Third parties may submit source code, metadata and the like and they may also request comparisons of source code and metadata. In addition, a web agent or crawler may be used to mine data from the Internet or other sources.

[10] What is claimed is:

1. A method for comparing software source code, the method comprising:
profiling one or more pieces of software source code using a one-way transform having a
natural language characterization algorithm;
generating metadata based on the profiling;
comparing two or more metadata files; and
generating output reflecting the results of the comparison.

New Applications

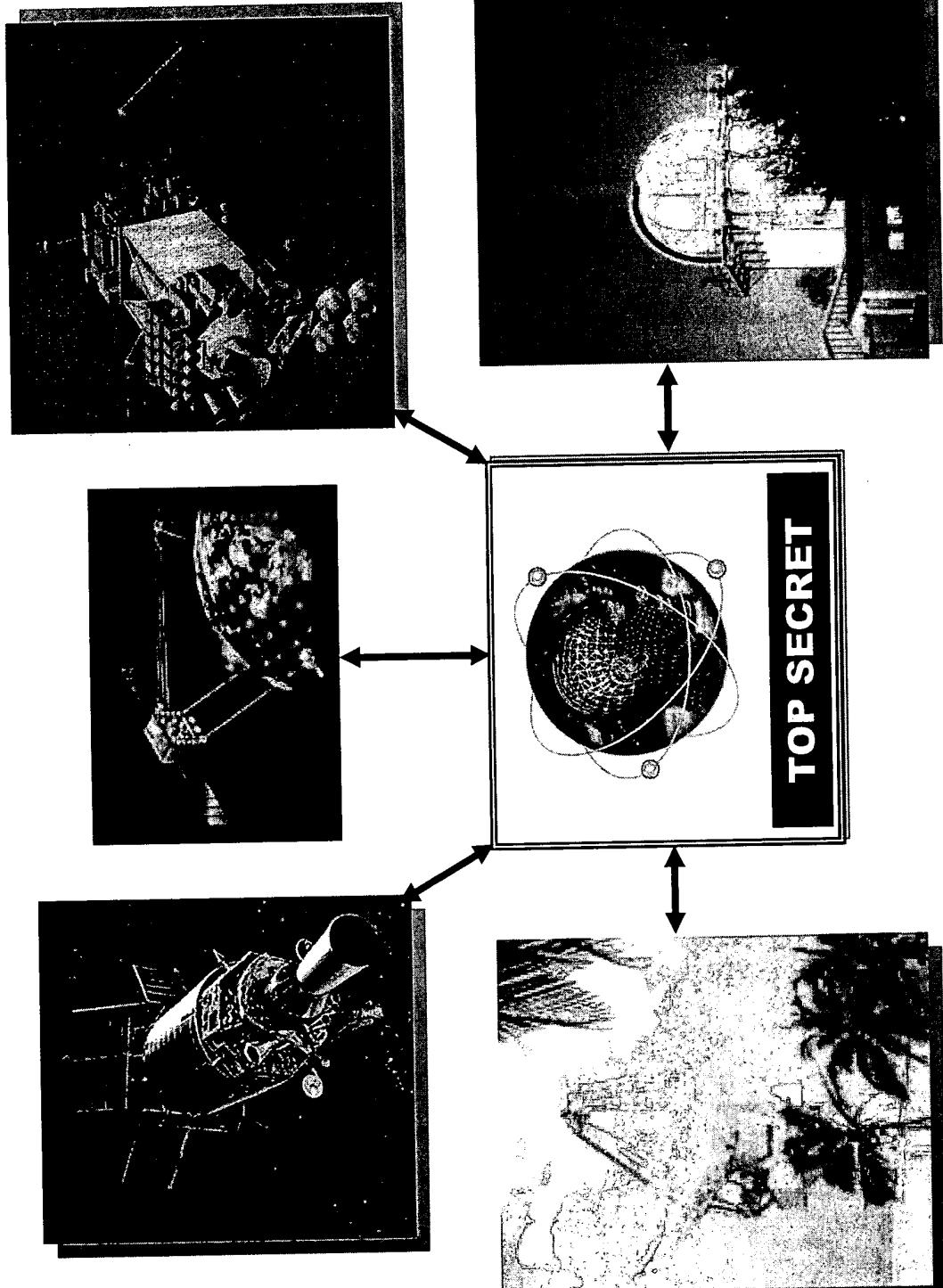
IP Clearing House

- Data Fusion
- Multi-dimensional Characterization
- Public Key Inspection



Background

Data Fusion / Multi-source Profiling / Threat Assessment

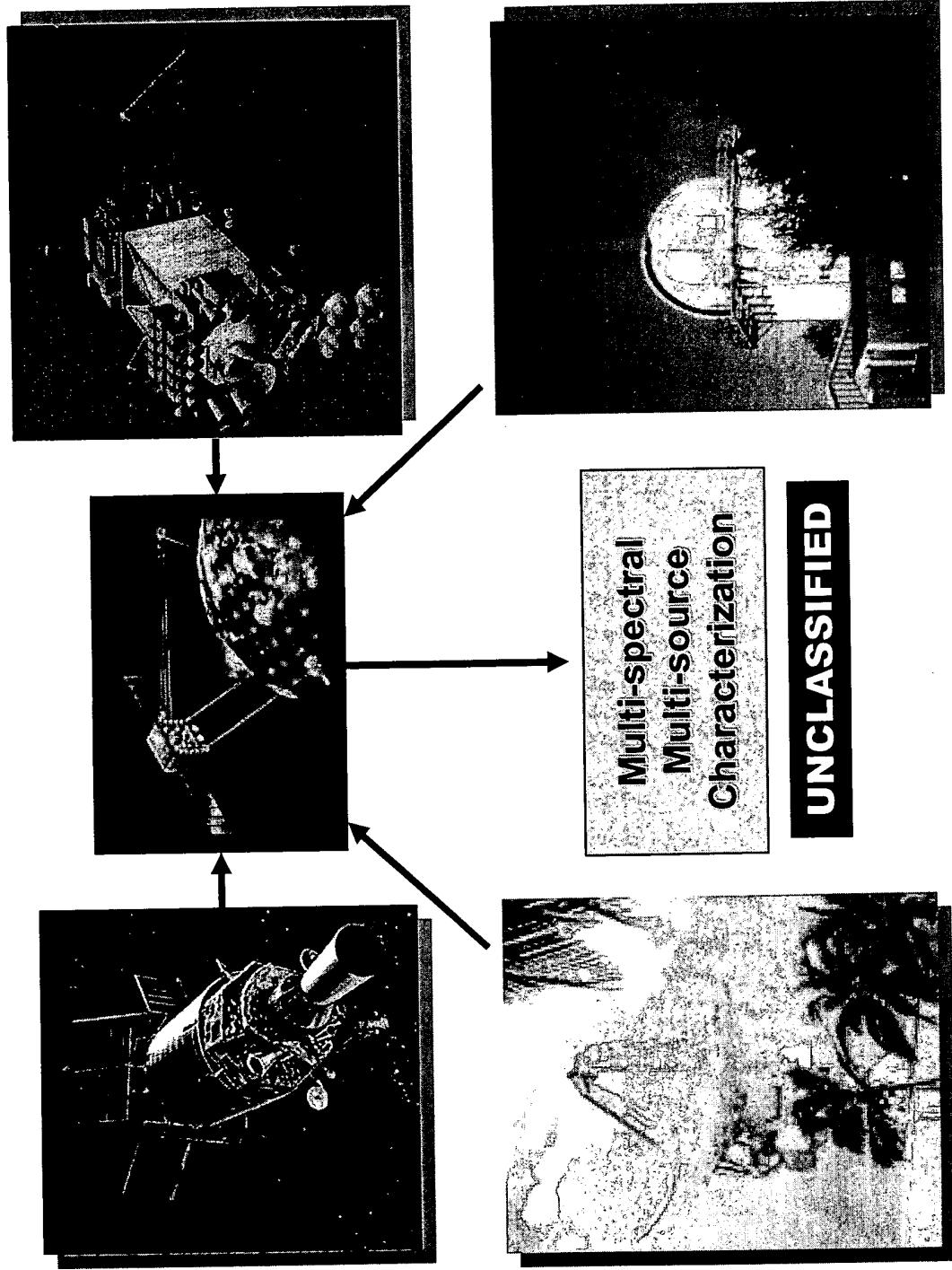


SCO Confidential



Transform

Data Fusion / Object Characterization



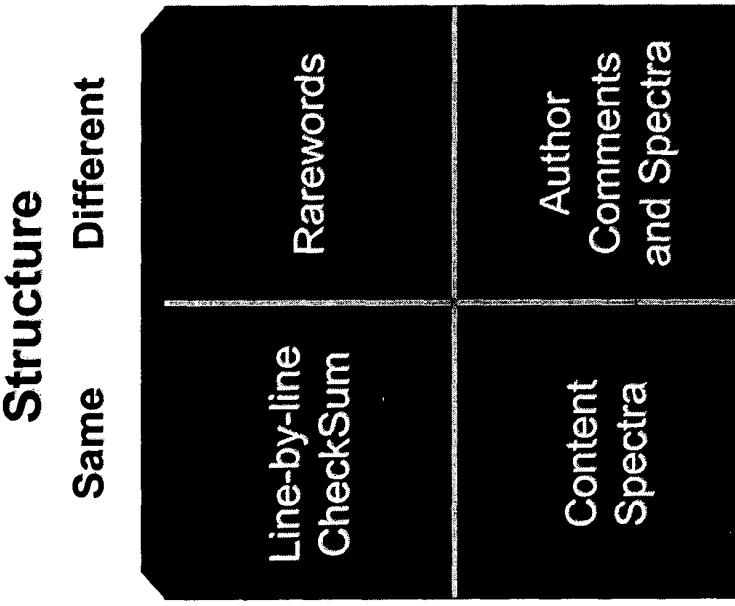
SCO Confidential

Multi-Source Characterization

IP Characteristics

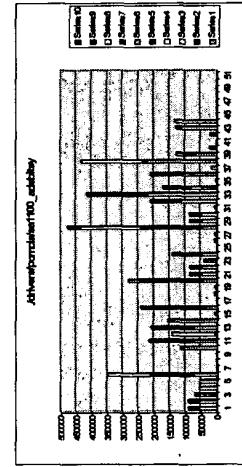
Structure

Same Different

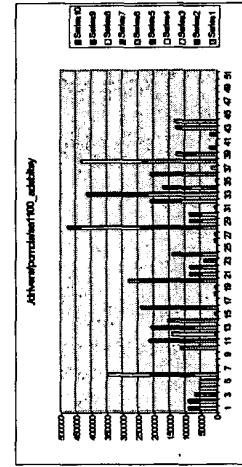
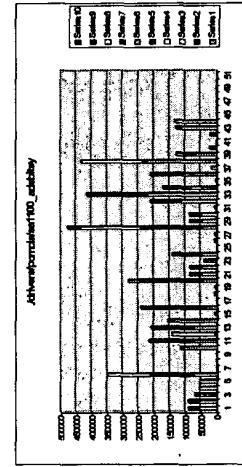
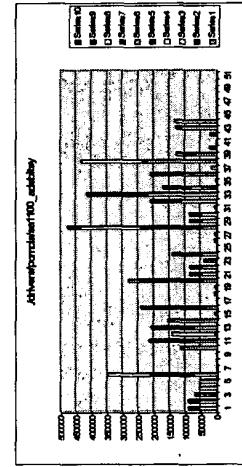
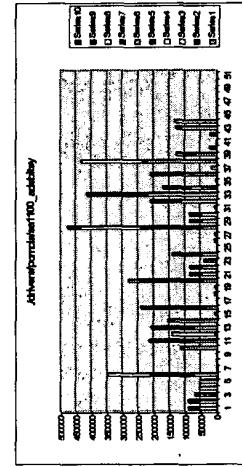
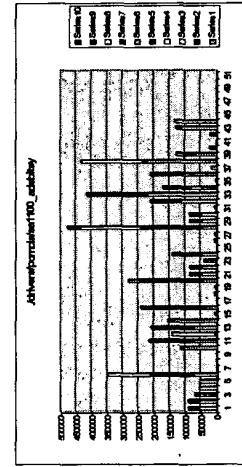
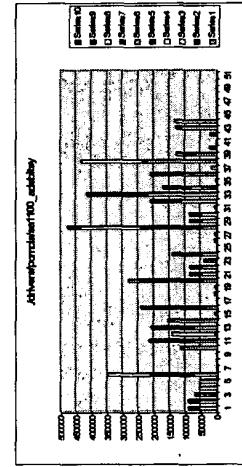
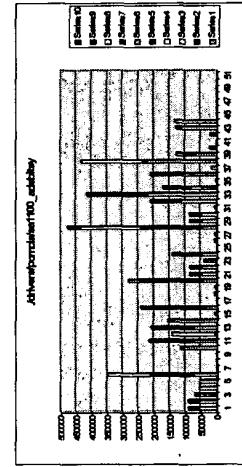
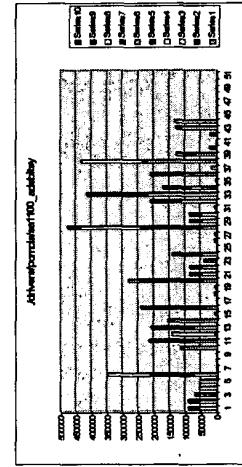
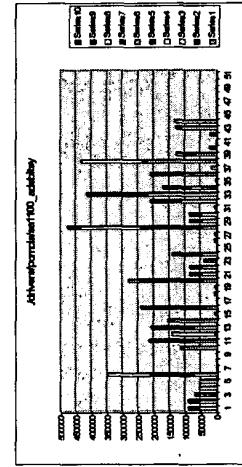
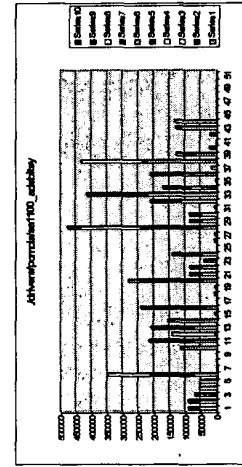
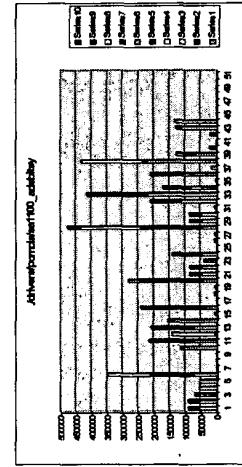
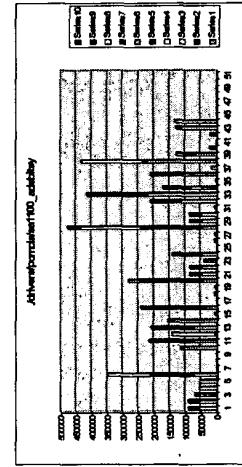
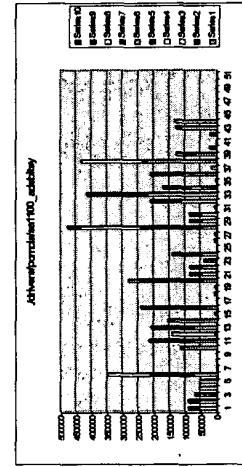
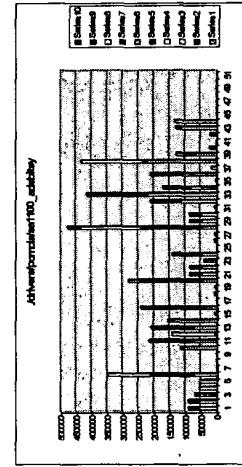
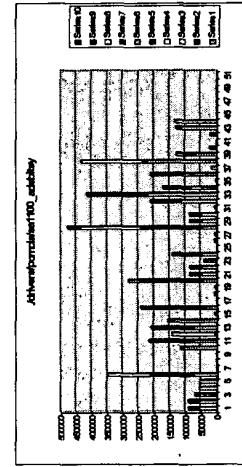
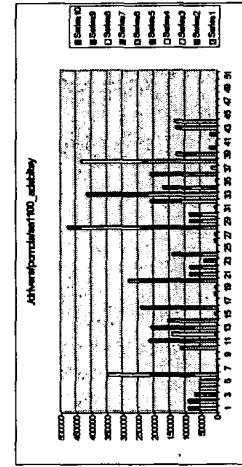
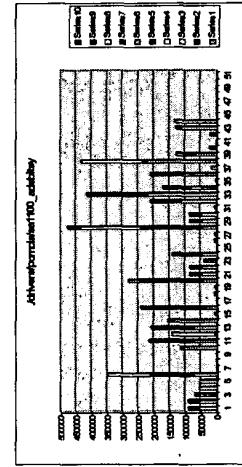
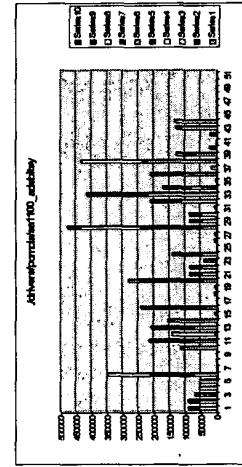
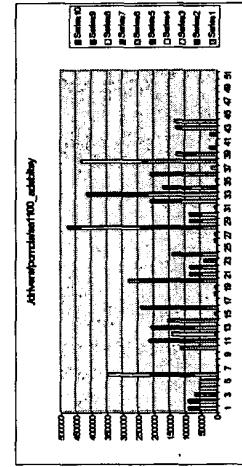
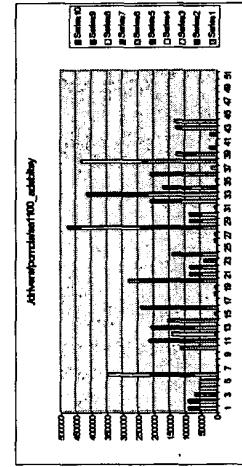
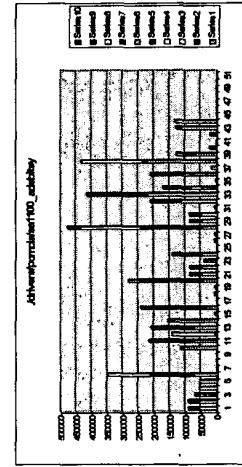
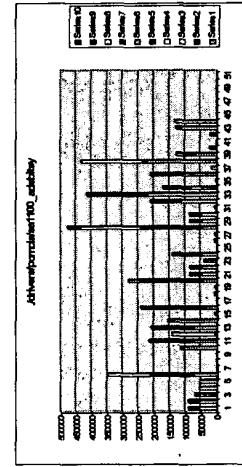
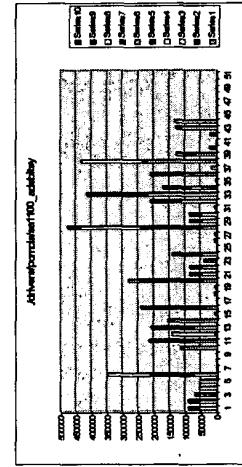
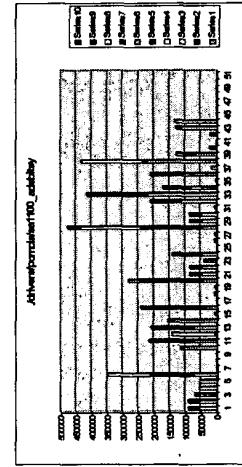
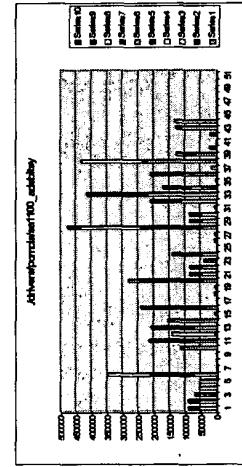
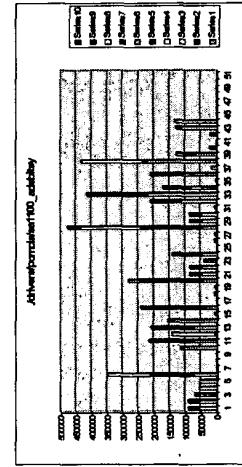
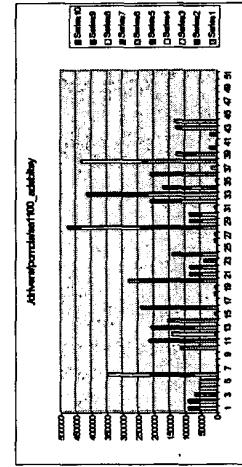
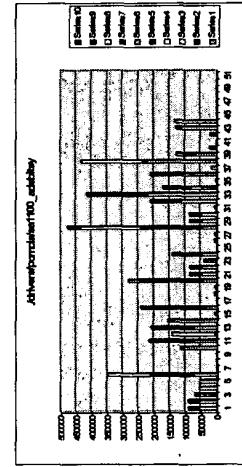
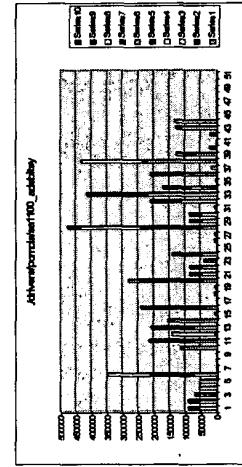
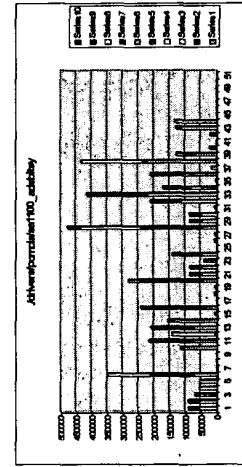
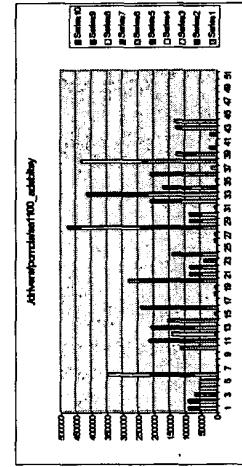
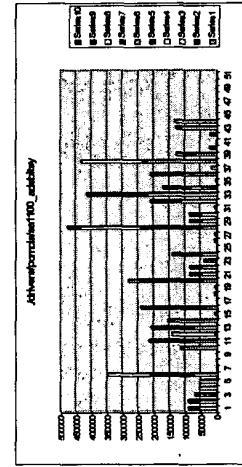
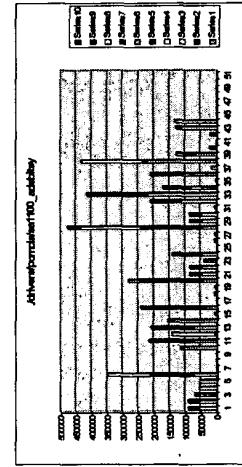
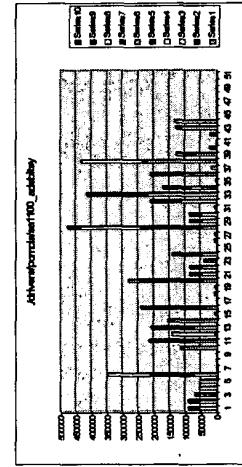
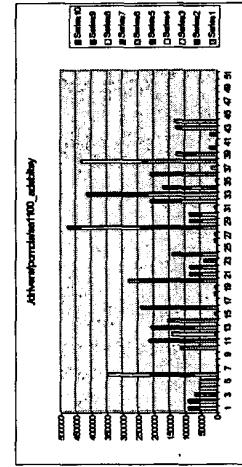
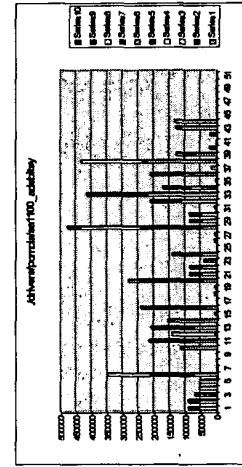
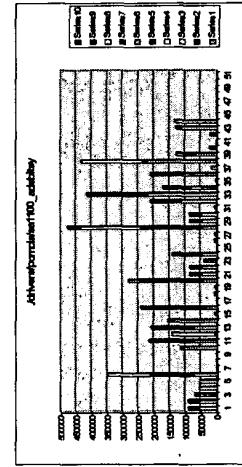
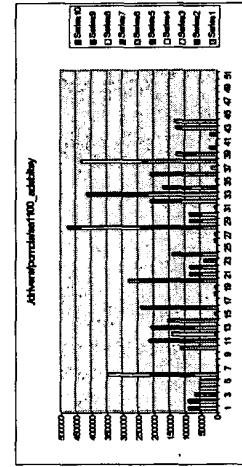
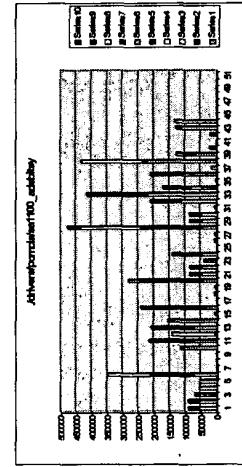
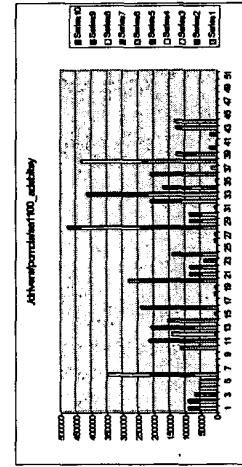
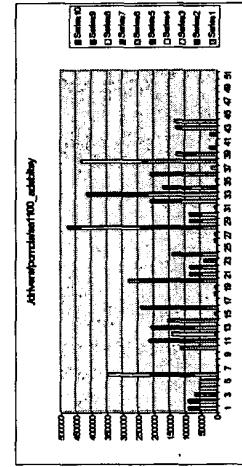
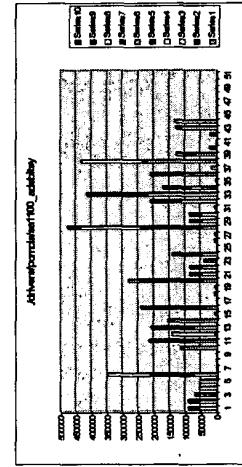
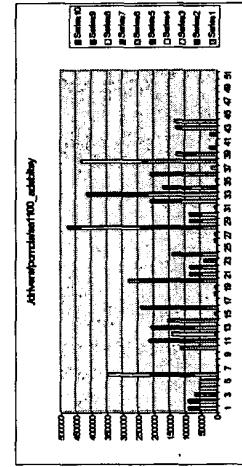
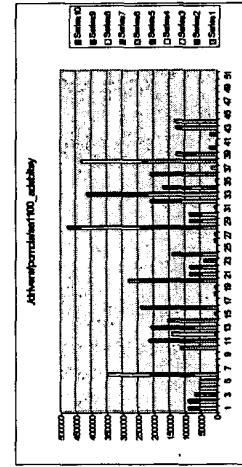
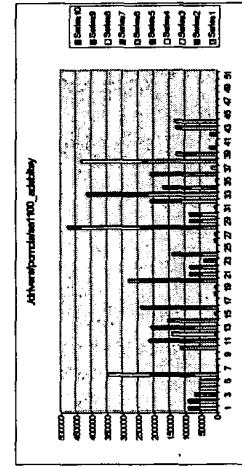
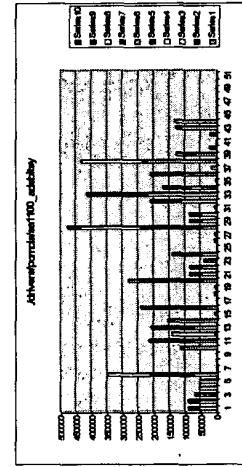
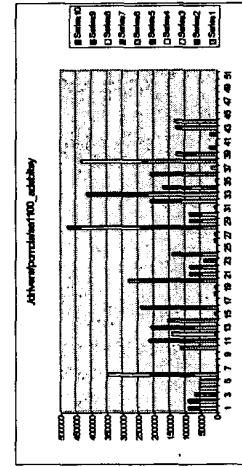
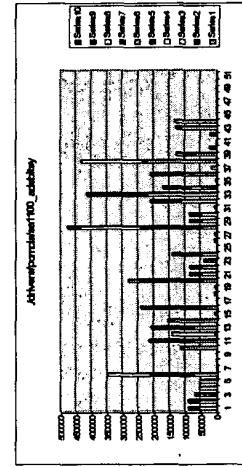
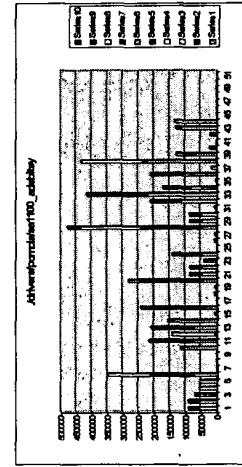
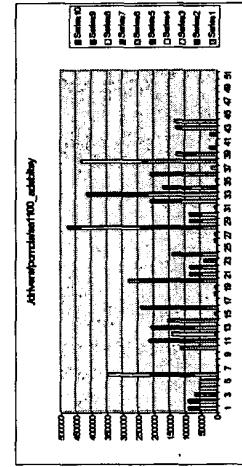
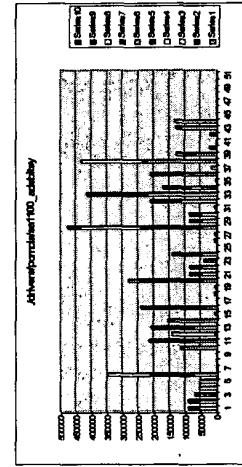
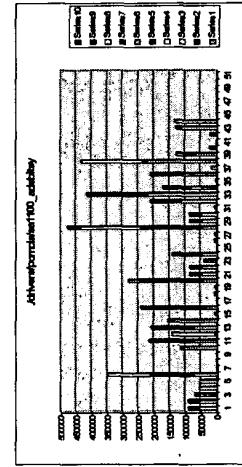
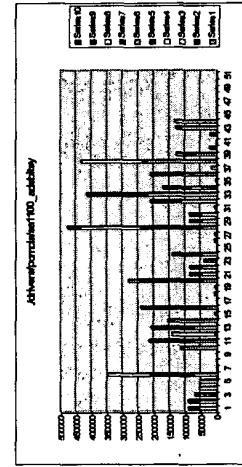
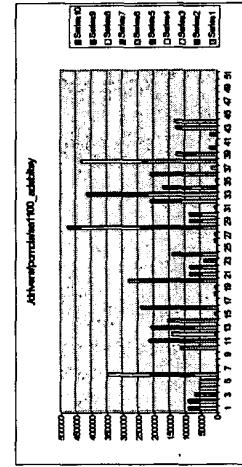
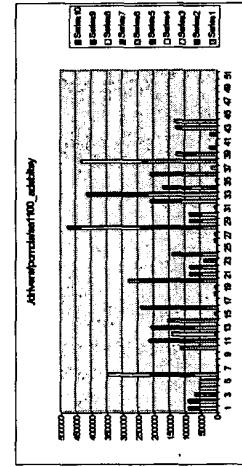
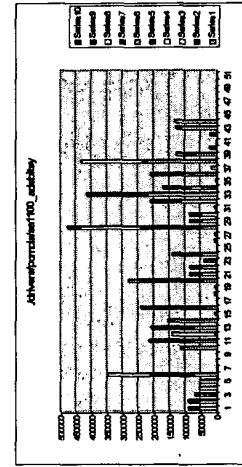
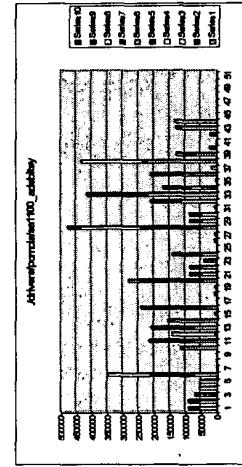
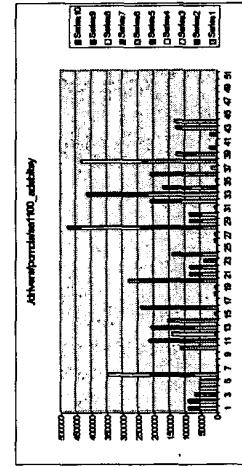
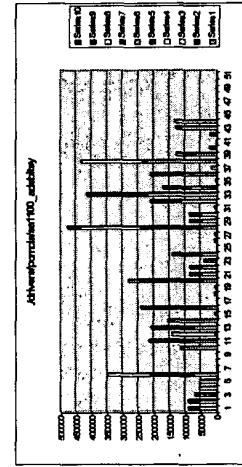
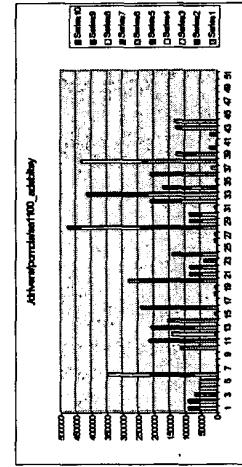
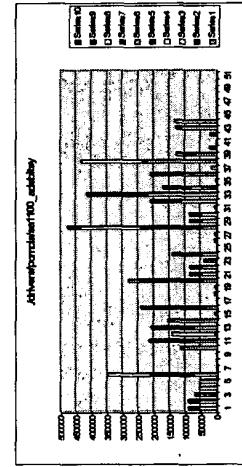
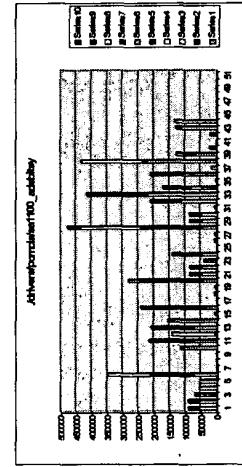
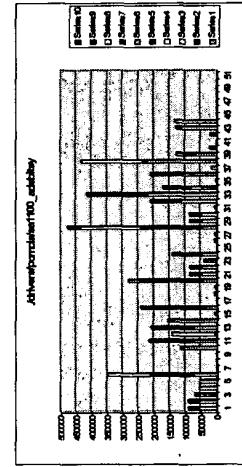
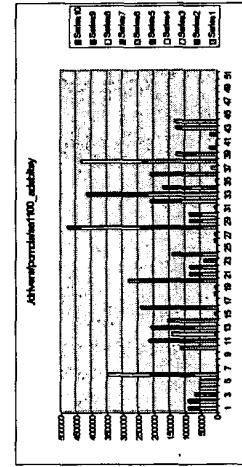
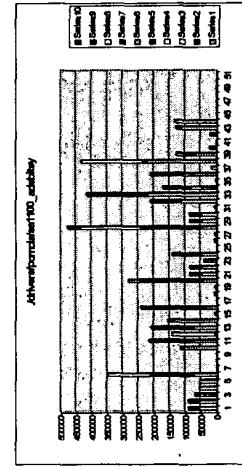
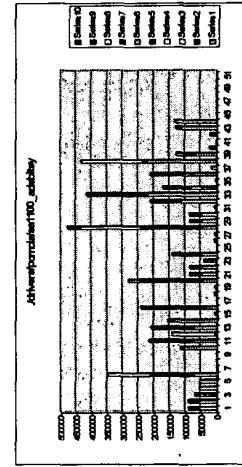
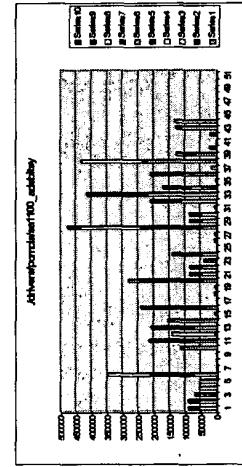
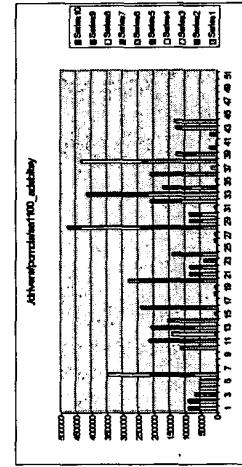
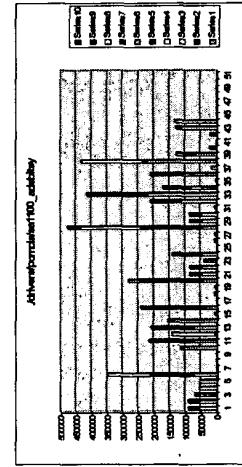
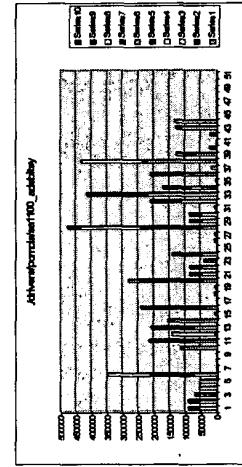
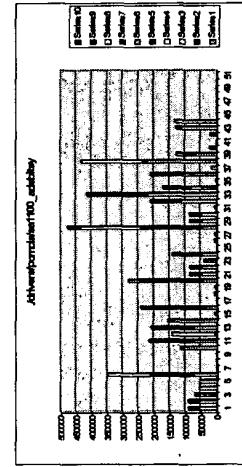
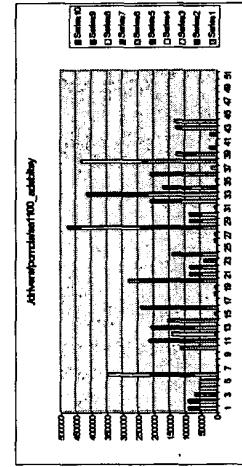
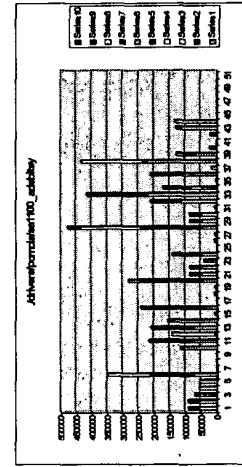
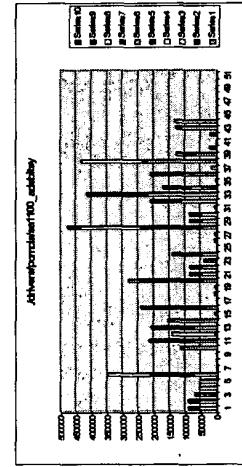
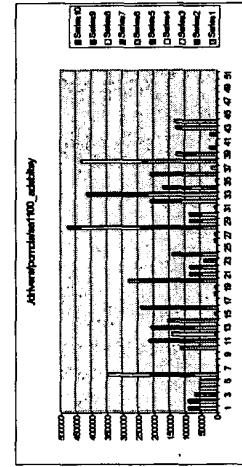
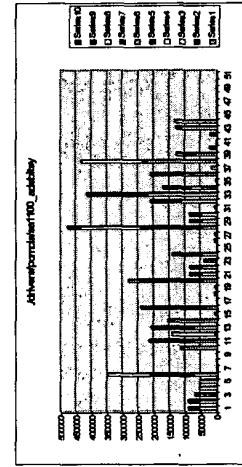
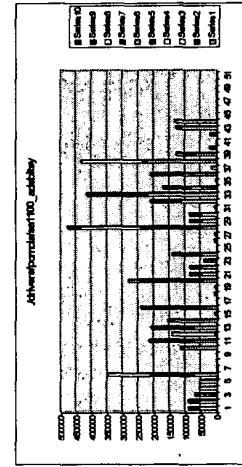
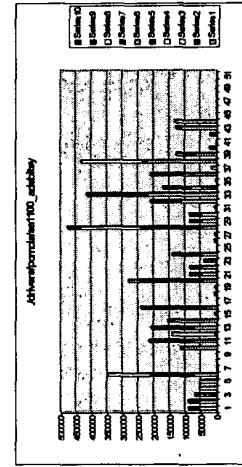
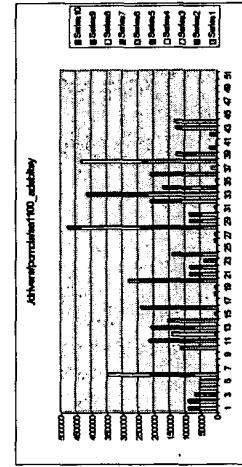
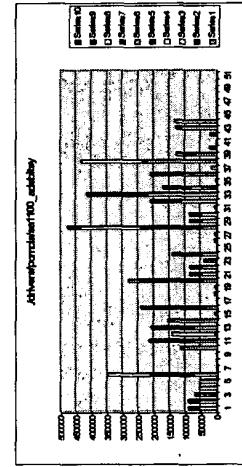
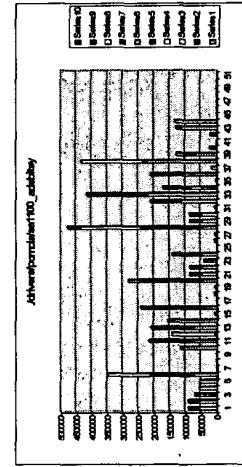
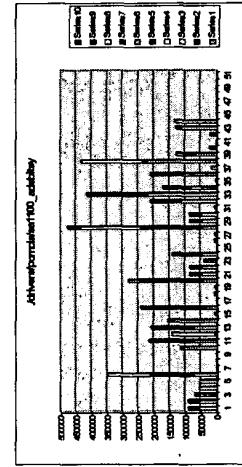
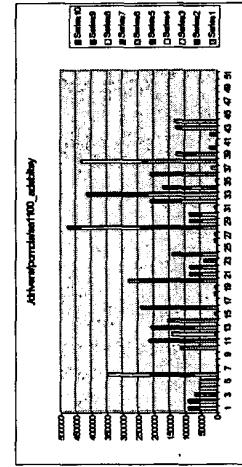
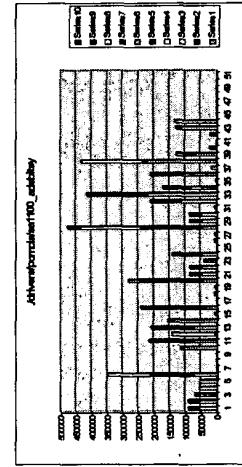
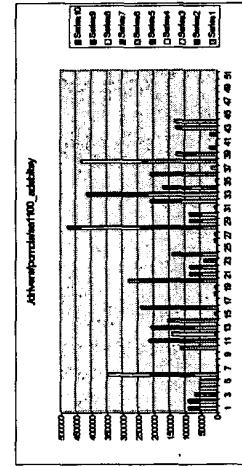
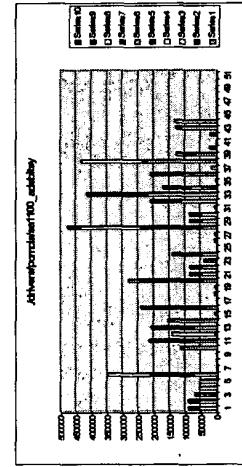
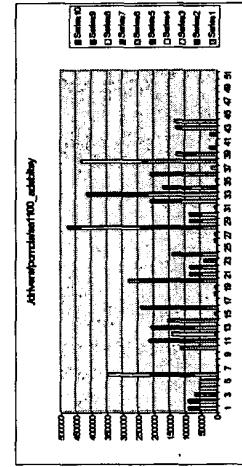
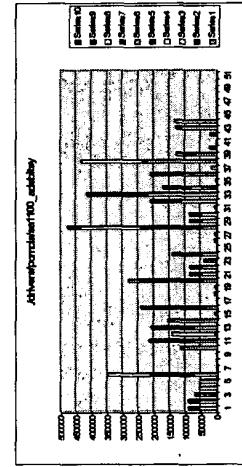
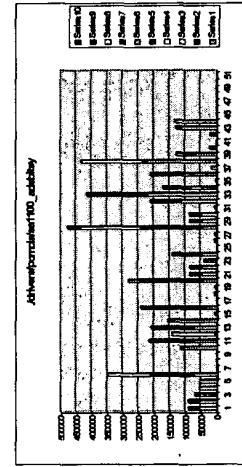
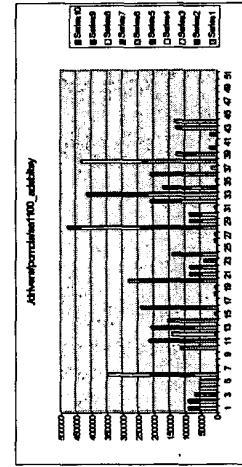
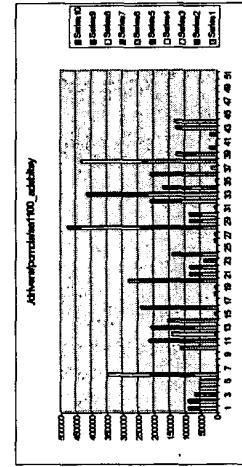
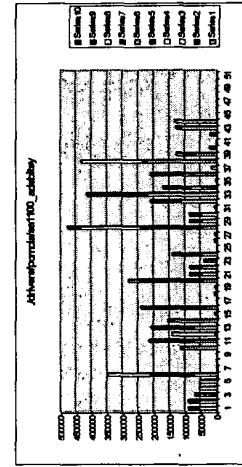
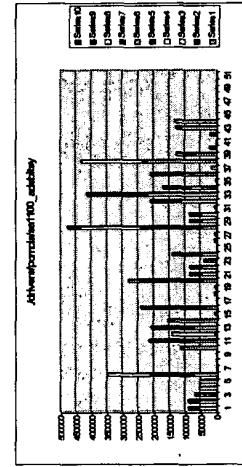


Different Same

Content
Spectra



Content
Spectra



Public Key Inspection™



Private

Proprietary
Intellectual
Property

Transform

IP
Characterization

Compare

IP
Characterization
Database

Public

SCO Confidential

Patent Pending

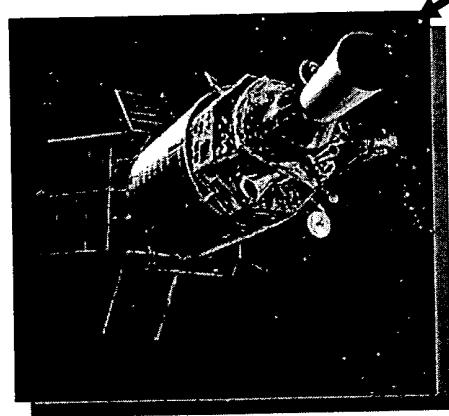
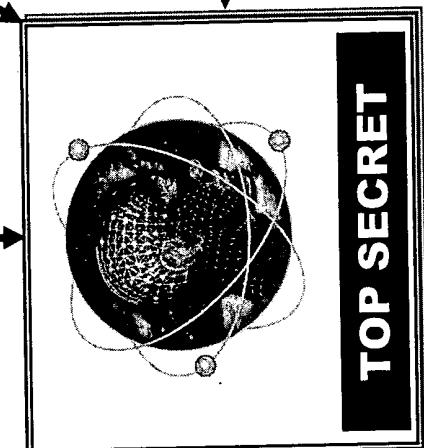
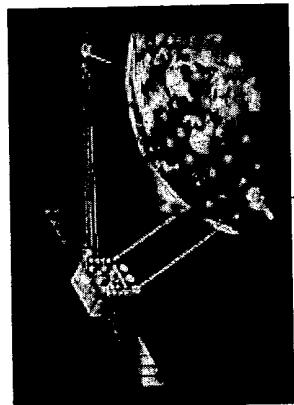
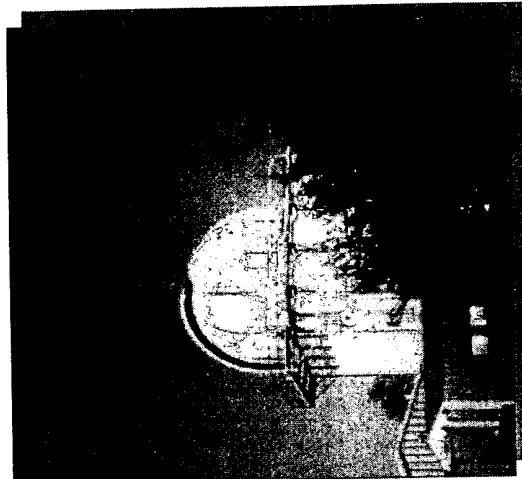
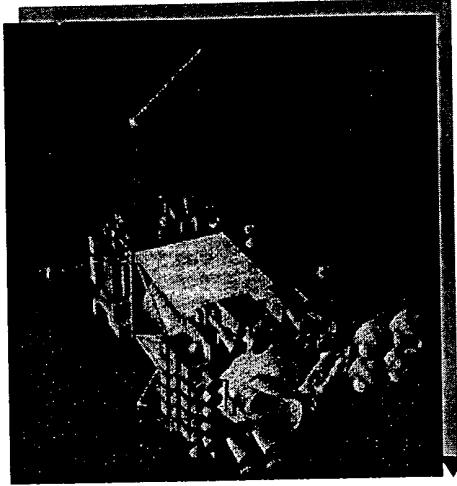
IP Analysis Tools

- Data Gathering
- Natural Language Processing
- Spectral Analysis
- Pattern Recognition



Background

Data Fusion / Multi-source Profiling



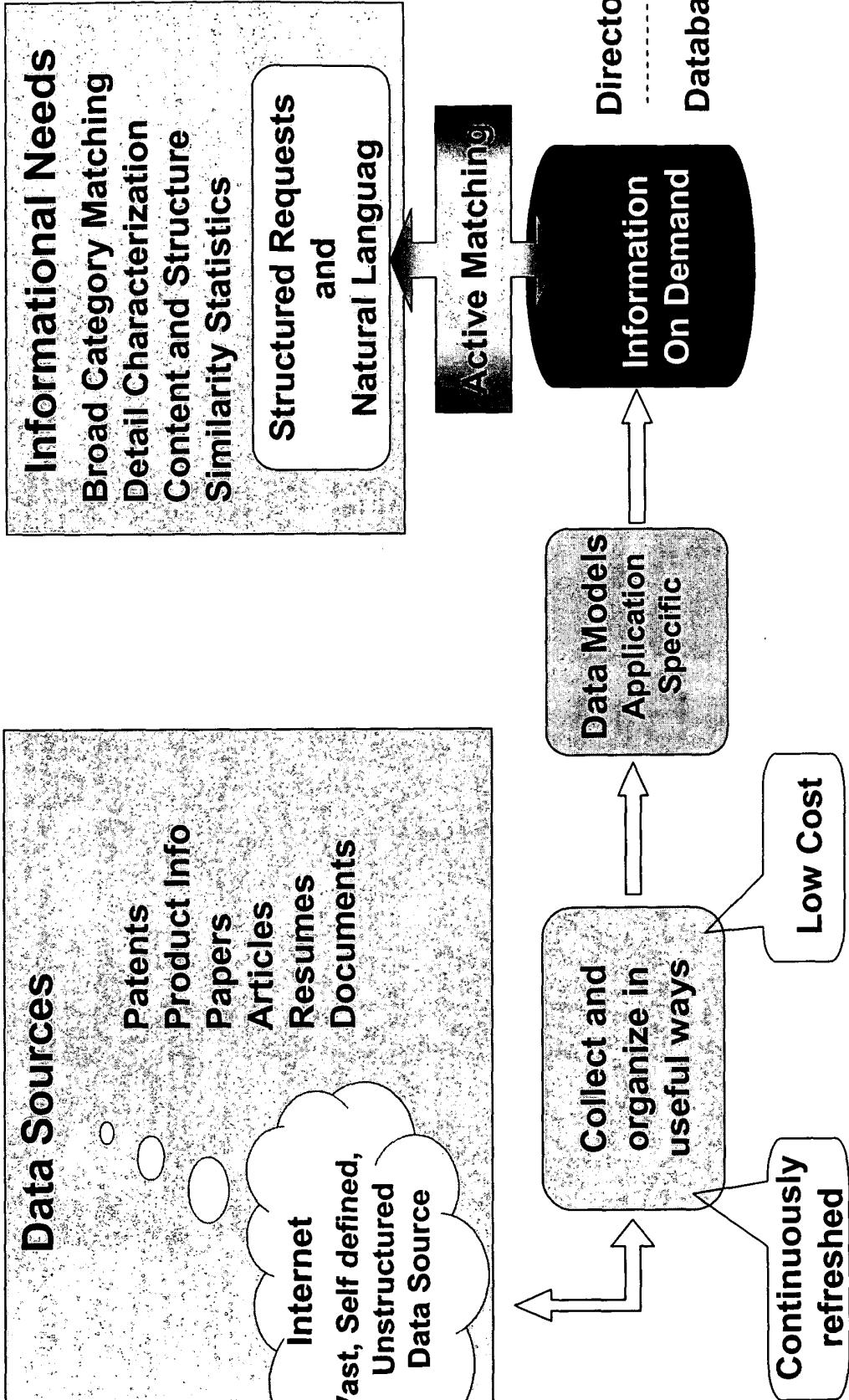
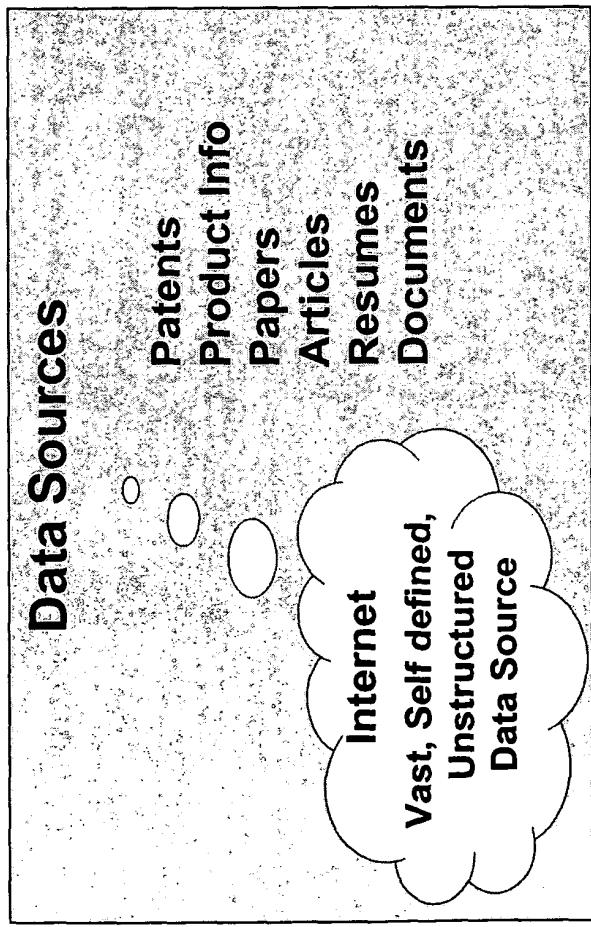
SCO Confidential

Outline



- The SCO IP Model
- Products and Services based on Technology
 - Natural Language Processing
 - Hirehub
 - ZillionResumes

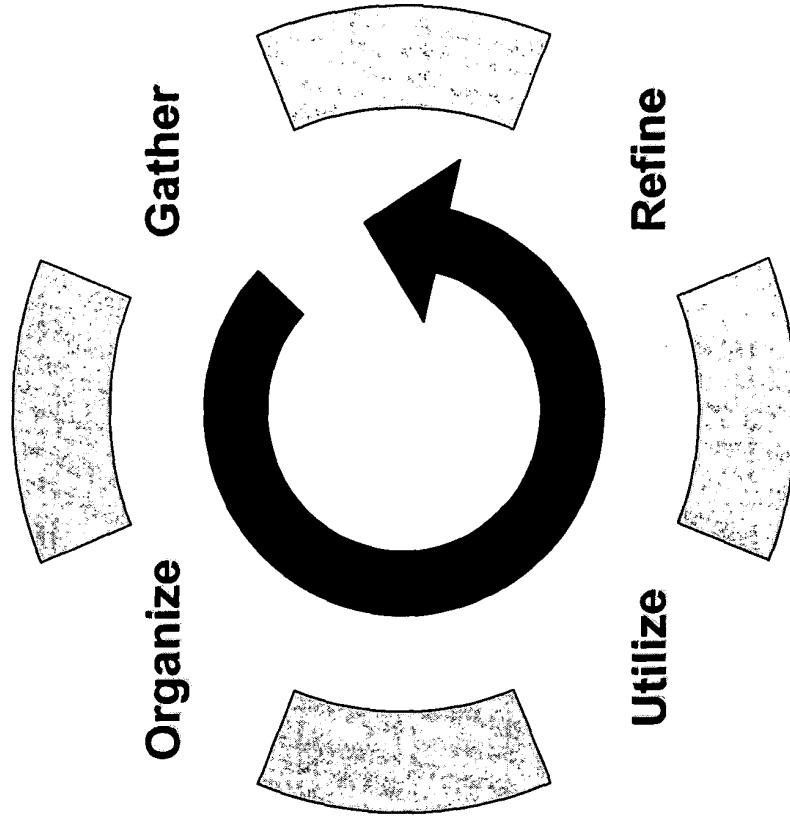
Natural Language Processing



Content Gathering Technology



- **Gather**
 - Goal-driven robots provide focused, deep traversal of web sites
 - No labor required to gather, input, or change
- **Organize**
 - Robots populate pre-defined typically organized data model
 - Learning robots discover and correlate key concepts, then dynamically define and add new data fields to the data model
- **Utilize**
 - Create summaries and aggregate information products
 - Apply advanced matching robots to deliver content consistent with user profiles
- **Refine**
 - Perpetual activity and feedback continuously tune the robots



Deployed Architecture



Internet

20K resumes &
bios a week

Gather

18K base concepts
700 generalizations

key criteria
location
experience level
email contact

Perpetual

Corpus
Miner

Organize

Profiler

Indexer

Corpus

Repository

Stimulus

Blackboard

Session

Concept Indices

Profiles

Matcher

Website

Utilize

On demand

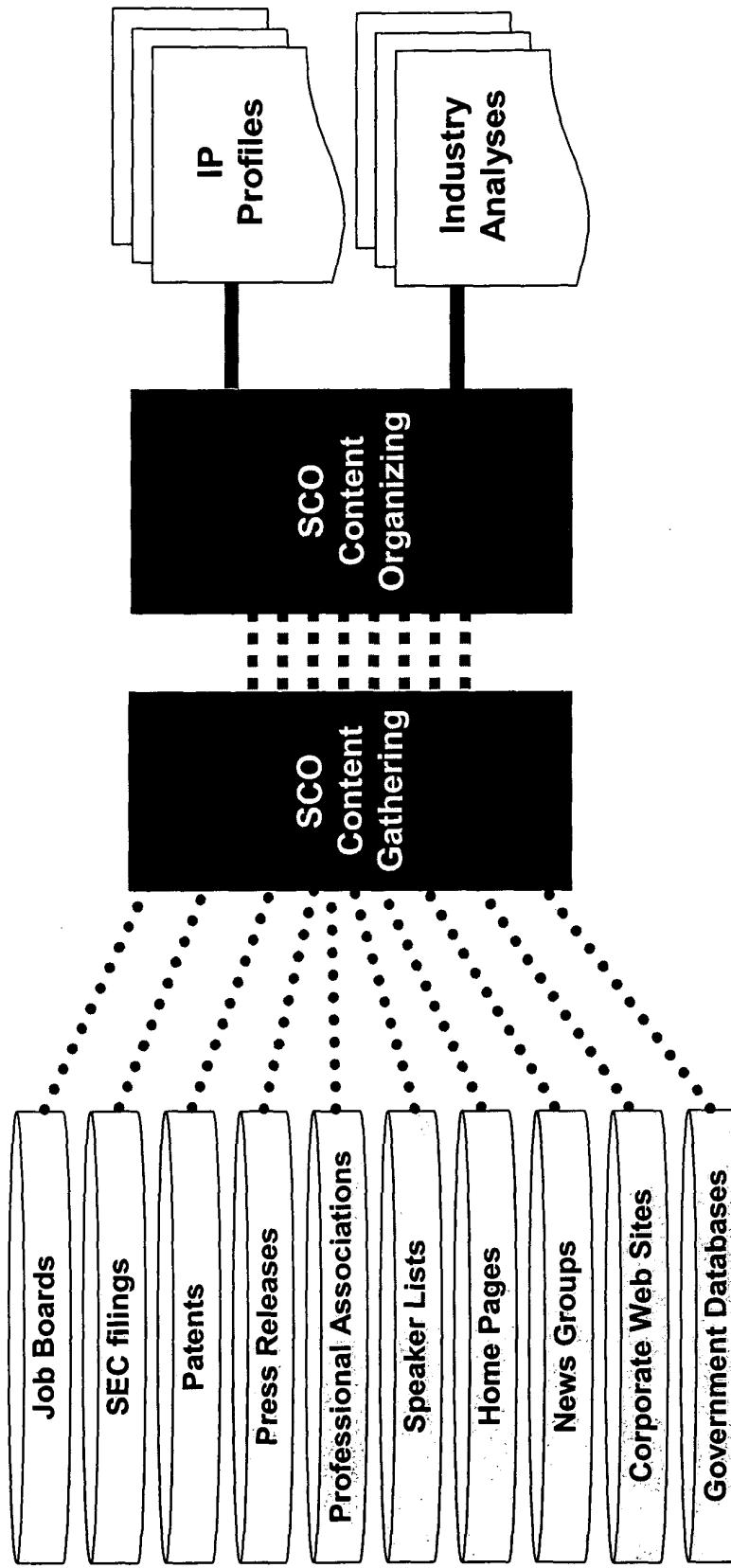
380K candidates
176K jobs

saved
matches

4-tiered, best-fit
optimizer



SCO Information Gathering



The SCO Difference



- **Demonstrated Technology**
 - Gathering and Organizing Large Unstructured Data Sets
 - Resumes
 - Patents
 - Source Code
 - Distributed, scalable, agent based content gathering infrastructure
 - NLP, fact extraction, auto profiling heuristics semantically targeted at industry specific content
 - Deployed applications in e-recruiting, software IP tracking
- **Ability to Innovate**
 - Merge various traditionally unassociated technologies to provide discontinuous innovation.

CONFIDENTIAL

Example Analysis Results



Approach

Concept Dictionaries for Source Code

- **Process every file, line by line, in a code base**
 - Different dictionary for .c files and for header (.h) files
- **Strip away:**
 - Comments
 - White spaces
 - C-specific characters like *, &, *, ;, (,)
- **Gather concept information:**
 - Produce a “raw concept file” which retains the line structure and records the concepts in a dictionary file
 - Tally up the “raw power” of each concept – which is the number of times that concept is used in the entire code base
 - Produce a “raw concept frequency file” for each source file, which records for each line the concepts replaced by their respective powers.
 - Assigns a frequency (or “power number”) to every term used in a code file



Approach

Produce Concept Histograms



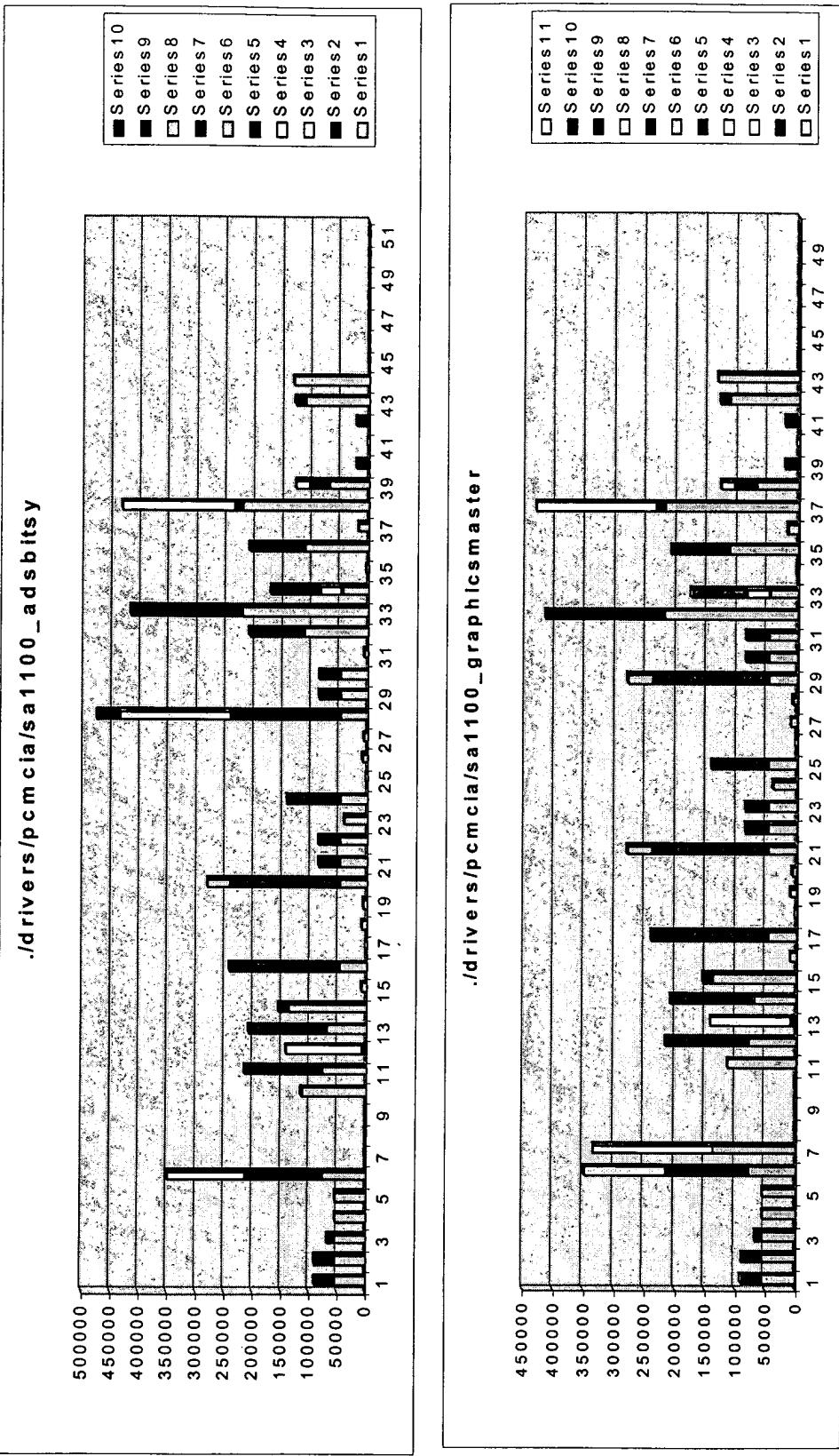
- For each line in a file, translate each concept to the power of that concept from the corpus dictionary
 - Example: a line might look like 2363:12:300:41 (the colons are artificially put in to distinguish different concepts on a line for easy input to a spreadsheet)
- Produce Spectral summary charts
 - A line will contain the name of the file, the number of lines in the file, the number of distinct concepts used in the file, and the total power of the lines in the file
- Can plot and display the file in bar-chart format
 - Horizontal axis is the line number of the file
 - Vertical axis is the total aggregate power of that line from the concept dictionary

Approach

Produce Concept Histograms, examples



- Two histograms from driver files in Linux kernel produced by the same manufacturer
 - the “spectral” similarity is apparent

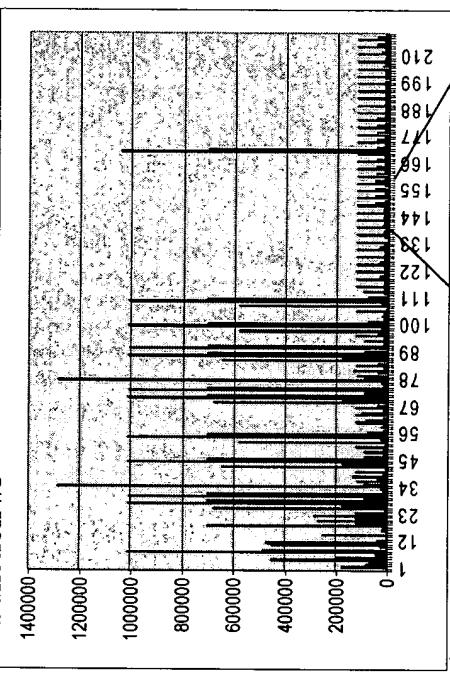


SCO Confidential

Spectral Analysis



SystemV File: common.uts.net(pf).c
Function #3



Code Extract

```

case BPF_JMP|BPF_JGT|BPF_K:
pc += (A > pc->k) ? pc->jt : pc->jf;
continue;

case BPF_JMP|BPF_JGE|BPF_K:
pc += (A == pc->k) ? pc->jt : pc->jf;
continue;

case BPF_JMP|BPF_JEQ|BPF_K:
pc += (A == pc->k) ? pc->jt : pc->jf;
continue;

case BPF_JMP|BPF_JSET|BPF_K:
pc += (A & pc->k) ? pc->jt : pc->jf;
continue;

case BPF_JMP|BPF_JGT|BPF_X:
pc += (A > X) ? pc->jt : pc->jf;
continue;

case BPF_JMP|BPF_JGE|BPF_X:
pc += (A == X) ? pc->jt : pc->jf;
continue;

case BPF_JMP|BPF_JEQ|BPF_X:
pc += (A == X) ? pc->jt : pc->jf;
continue;

case BPF_JMP|BPF_JSET|BPF_X:
pc += (A & X) ? pc->jt : pc->jf;
continue;

```

Code Extract

```

case BPF_JMP|BPF_JGT|BPF_K:
pc += (A > fentry->k) ? fentry->jt : fentry->jf;
continue;

case BPF_JMP|BPF_JGE|BPF_K:
pc += (A == fentry->k) ? fentry->jt : fentry->jf;
continue;

case BPF_JMP|BPF_JEQ|BPF_K:
pc += (A == fentry->k) ? fentry->jt : fentry->jf;
continue;

case BPF_JMP|BPF_JSET|BPF_K:
pc += (A & fentry->k) ? fentry->jt : fentry->jf;
continue;

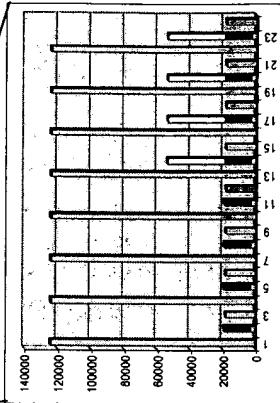
case BPF_JMP|BPF_JGT|BPF_X:
pc += (A > X) ? fentry->jt : fentry->jf;
continue;

case BPF_JMP|BPF_JGE|BPF_X:
pc += (A == X) ? fentry->jt : fentry->jf;
continue;

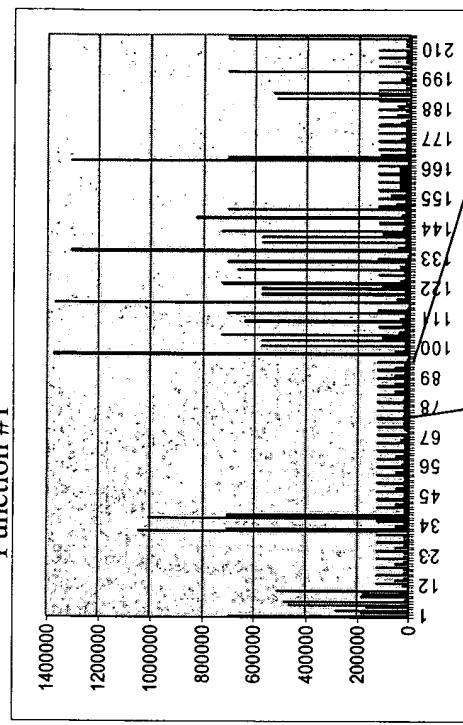
case BPF_JMP|BPF_JEQ|BPF_X:
pc += (A == X) ? fentry->jt : fentry->jf;
continue;

```

Spectral Extract



Linux File: net.core.filter.c
Function #1



IP Tracing

- Examples of Rare Words found with search
Linux Source
- Consider how this method was re-factored from AIX
AIX Source



Header comments:
Alternatively, the contents of this file may be used under the terms of the GNU General Public License version 2 (the "GPL"), in which case the provisions of the GPL are applicable instead of the above. Etc...

```
int CardServices(int func, void *a1, void *a2,  
                 void *a3)
```

```
{  
    switch (func) {  
        case RegisterClient: →  
            return pcmcia_register_client(a1, a2);  
        break;  
        case DeregisterClient: →  
            return pcmcia_deregister_client(a1);  
        break;  
        case GetStatus: →  
            return pcmcia_get_status(a1, a2);  
        break;  
        case ResetCard: →  
            return pcmcia_reset_card(a1, a2);  
        break;  
        case SetEventMask: →  
            return pcmcia_set_event_mask(a1, a2);  
        break;  
    ...  
    }  
    void initCSfuncs ()  
    {  
        /* * Initialize CS function dispatch table */  
        /* *  
         * [CSRegisterClient] = RegisterClient;  
         * [CSDeregisterClient] = DeregisterClient;  
         * [CSGetStatus] = GetStatus;  
         * [CSRSTCard] = ResetCard;  
         * [CSSSetEventMask] = SetEventMask;  
         */  
    ...  
}
```

The methods are implemented differently, but provide the SAME function, use the SAME variable names in the SAME order.

SCO Confidential