



Scalable Distributed Process Group Control and Inspection via the File System

Michael J. Brim
Advisor: Barton P. Miller

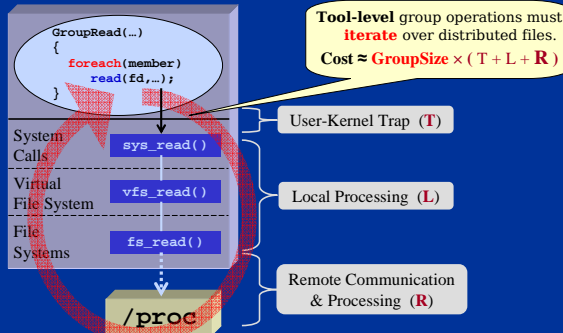
Computer Sciences
University of Wisconsin-Madison

Research Overview

- **Scalable Tools and Middleware are *Crucial***
 - Run, analyze, debug, and improve applications & systems
 - Require group operations on distributed files and processes
- **Problems at scale**
 - How to *access* files on thousands or millions of hosts?
 - How to efficiently *operate* over large file groups?
 - How to handle the *data explosion* for large file groups?
- **Solutions**
 - Access:** TBON-FS Distributed File System
 - Provides global namespace that combines remote file systems
 - Uses Tree-Based Overlay Network (TBON) for scalability
 - Operate:** Group File Operation Idiom
 - Eliminates explicit iteration imposed by file system interface
 - Data Explosion:** Data Aggregation
 - Explicit semantics for aggregation in group file operations
 - TBON-FS architecture for scalable distributed aggregation

Group File Operations

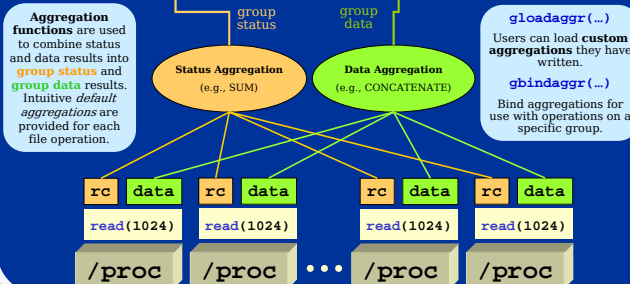
- **Problem: Current File System API Forces Serial Access**



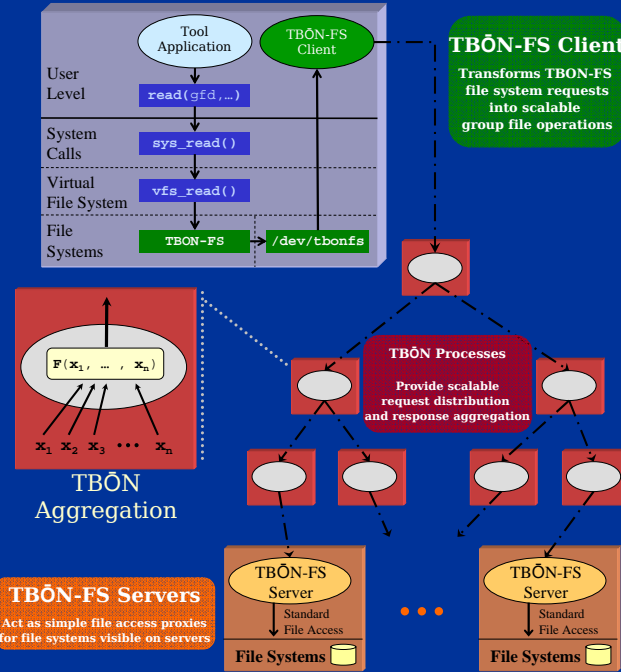
- **Solution: Group File Operation Idiom**

Group File Descriptor - Use with POSIX file operations to act on all group members.
`gopen` - Open all files in directory.

```
int gfd = gopen("directory", flags)
int rc = read(gfd, databuf, 1024)
```



TBON-FS Architecture



TBON-FS Namespace

- **Client Namespace Aggregates Server Namespaces**
 - Ideal Namespace: directories already contain desired file groups
 - FiNAL: File Namespace Aggregation Language**
 - Goal: allow users to specify interesting groups with custom namespace
 - Operators on FS hierarchies: insert, remove, replace, move, merge
 - TBON provides scalable aggregation and name resolution
- **Example: Global Linux /proc**
 - Contains files providing system and process information
 - System Files: cpuinfo, loadavg, meminfo, modules, ...
 - Processes: directory per unique process id
 - Files: cmdline, environ, mem, stat, statm, ...

Simple
Directory per Server

```
/tbonfs/proc:
serverA/
  /cpuinfo
  /1/...
  /4/...
=
serverZ/
  /cpuinfo
  /1/...
```

Merge
Common Files => Directory (per server files)
Common Directories (per server sub-directories)

```
/tbonfs/proc:
cpuinfo/
  /serverA/...
  /...
  /serverZ/...
1/
  /serverA/...
  /...
  /serverZ/...
4/
  /serverA/...
```

Deep Merge
Start merge at leaves of common sub-hierarchy, rather than root

```
/tbonfs/proc:
cpuinfo/
  /serverA/...
  /...
  /serverZ/...
1/
  /cmdline/...
  /...
  /serverA/...
  /...
  /serverZ/...
4/...
```

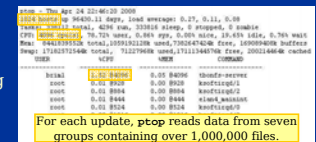
All the World's a File

- **Goal: Abstract *user-defined functions* (UDFs) as files**
 - Enables use with group file operations
 - UDFs with file interface: Previous work exists
 - Plan9 operating system: 9P file service protocol
 - FUSE: user-level file systems
 - Good: Arbitrary code for file system calls
 - Bad: Too much file system specific code, requires kernel support
 - Our Approach: **Synthetic Files**
 - Similar to FUSE: UDFs for file system calls
 - No kernel support required: UDFs execute in TBON-FS server
- **Example: Process Group Control and Inspection**
 - File access for all control and inspection operations
 - read / write process memory or registers
 - handle signals and exceptions (e.g., breakpoints)
 - stop, continue, step
 - thread control and inspection
 - Three Components**
 - Portable process control and inspection layer
 - Synthetic File** for each control / inspection operation
 - additional control file for process execution (a la *Xcpu*)
 - capture standard I/O streams as files too
 - Custom **FiNAL** namespace
 - combine new synthetic files with existing files (i.e., /proc)

Demonstrations

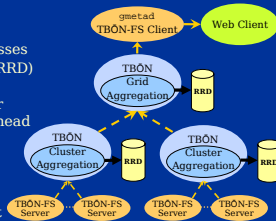
- **System Administration: Parallel Linux commands**

- ptop** : host & process monitor
- pgrep** : parallel text search
 - identify configuration differences
- ptail** : parallel log monitoring
 - with -follow option
- pcp, psync** : file distribution
 - send whole file, or only changes



- **System Monitoring: Ganglia**

- Original System**
 - gmetad cluster / Grid aggregation processes
 - store data to round-robin database (RRD)
 - gmond monitoring daemons
 - multicast data updates within cluster
 - linear memory / CPU / network overhead
- Ganglia-tbonfs**
 - top-level gmetad TBON-FS client
 - scalable TBON-FS aggregation
 - gmonds replaced by TBON-FS servers
 - eliminates overhead due to multicast



- **Example Code: Distributed Debugger**

- Define groups
 - Stop all processes
 - Set a group breakpoint
 - Continue all processes
 - Wait for processes to hit breakpoint
 - Read program variable, generate equivalence classes for current values
- ```
1. mkdir("grp_ctl_dir"); mkdir("grp_mem_dir");
 foreach(member in "/tbonfs/proc/[1-9]*") {
 symlink("member/signal", "grp_ctl_dir");
 symlink("member/mem", "grp_mem_dir");
 }
 ctl_gfd = gopen("grp_ctl_dir", O_WRONLY);
 mem_gfd = gopen("grp_mem_dir", O_RDWR);

2. write(ctl_gfd, SIGSTOP, 4);
3. lseek(mem_gfd, brkpt_addr, SEEK_SET);
 write(mem_gfd, brkpt_code_buf, code_sz);
4. write(ctl_gfd, SIGCONT, 4);
5. WaitForAll(ctl_gfd);

6. lseek(mem_gfd, var_addr, SEEK_SET);
 gbindaggr(mem_gfd, OP_READ, EQUIV_CLASS_AGG, ...);
 read(mem_gfd, grp_var_classes, var_sz);
```