# Improving the Scalability of the TotalView Debugger using TBONs

## Michael J. Brim

Paradyn Project, University of Wisconsin

## John DelSignore

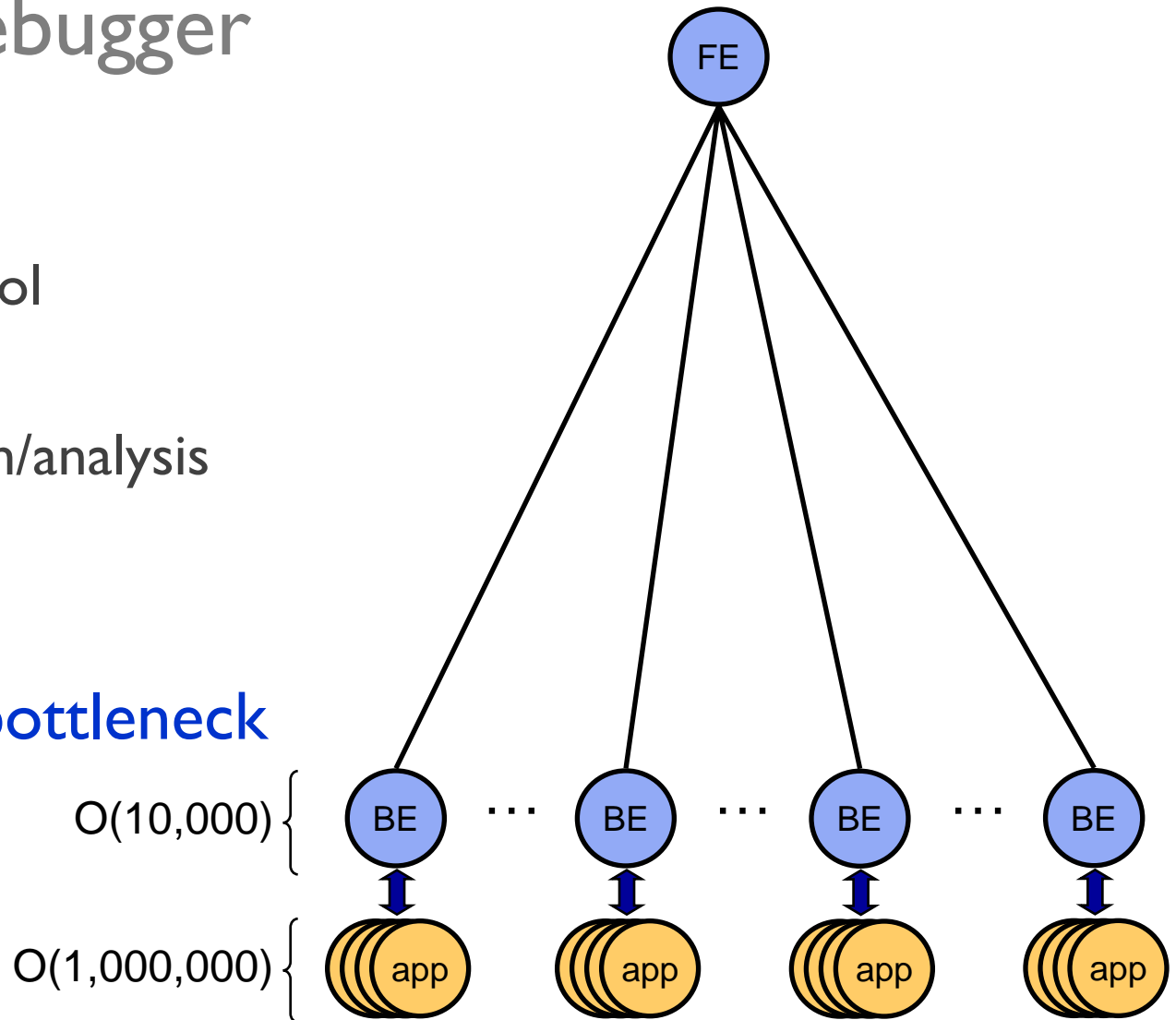Rogue Wave Software

CScADS

August 1, 2011

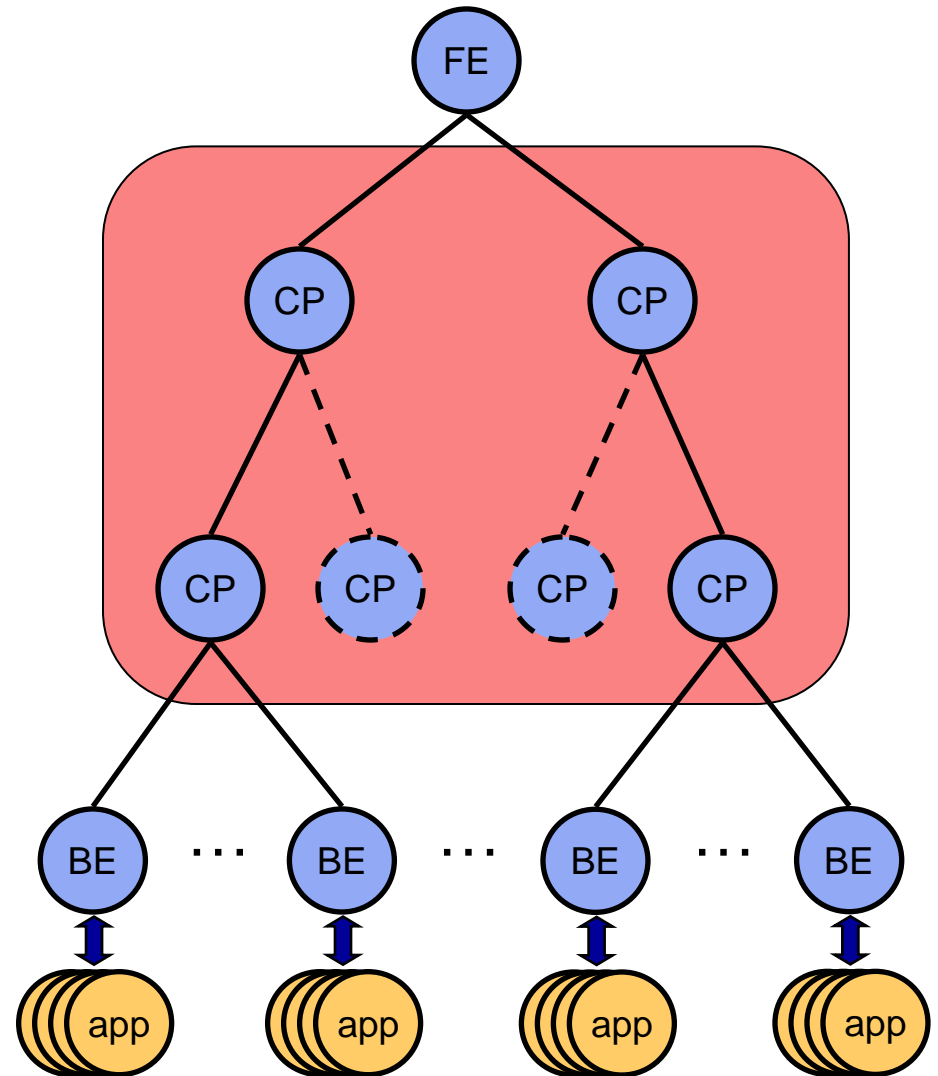# The ~~Tool~~ Scalability Problem
# TotalView Debugger



## Key tasks:

- Application Control

- Data collection

- Data centralization/analysis

**As scale increases,
front-end becomes bottleneck**

O(10,000)

O(1,000,000)

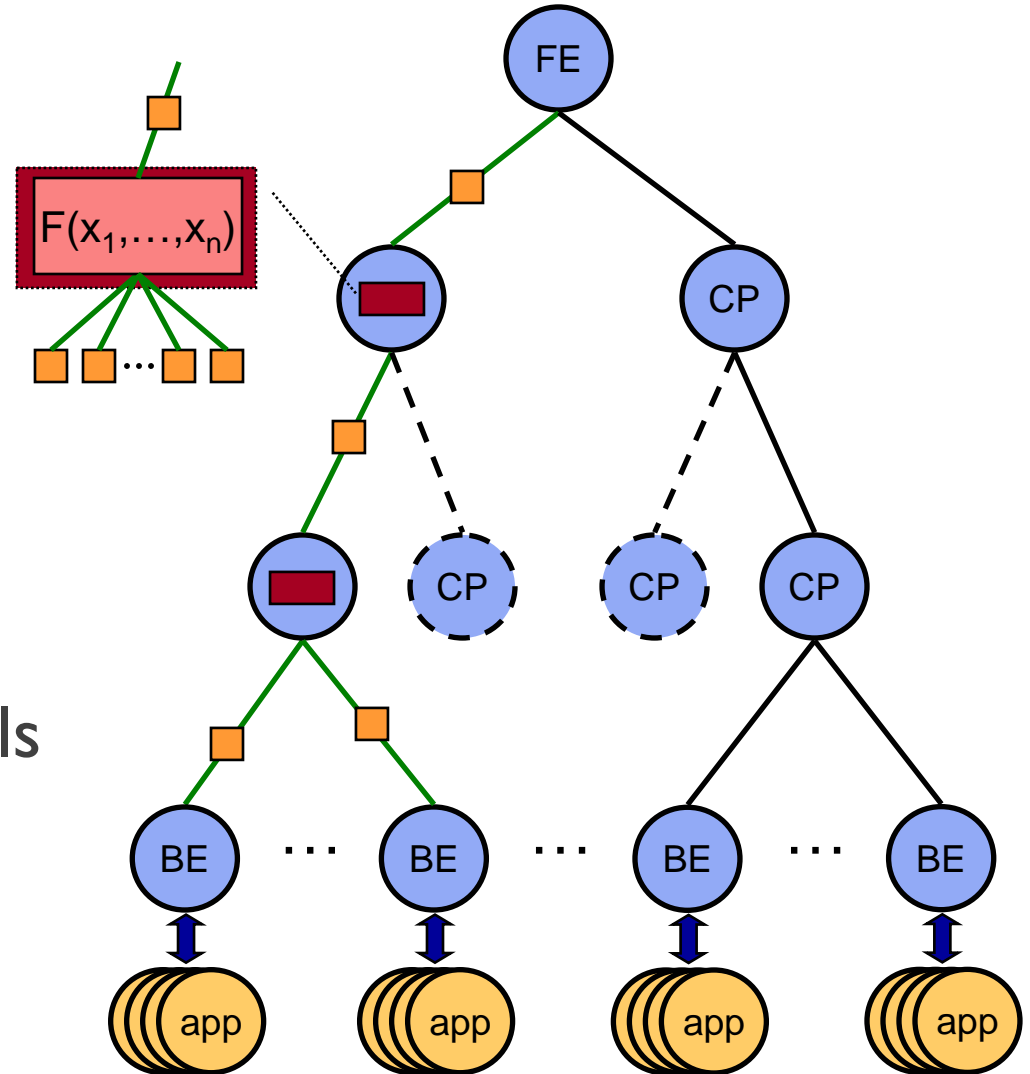# Tree-Based Overlay Networks (TBONs)

o Scalable multicast

o Scalable gather

o Scalable data aggregation

o Natural redundancy

# MRNet – Multicast / Reduction Network

## General-purpose TBON API

- Network: user-defined topology
- Stream: logical data channel
  - to a set of back-ends
  - multicast, gather, and custom reduction
- Packet: collection of data
- Filter: stream data operator
  - synchronization
  - transformation

## Widely adopted by HPC tools

- CEPBA toolkit
- Cray ATP & CCDB
- Open|SpeedShop & CBTF
- STAT
- TAU
- …

$$F(x_1, \ldots, x_n)$$

FE

CP

CP  CP  CP

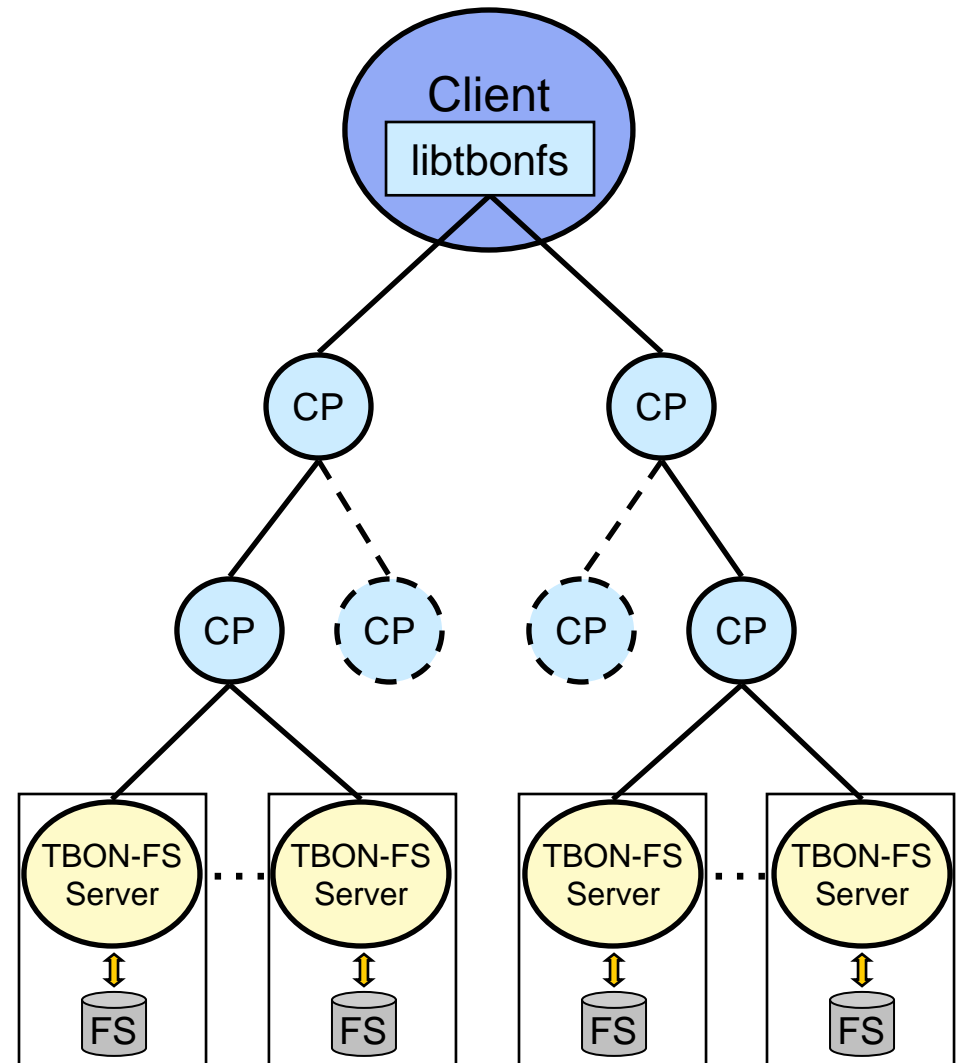BE  …  BE  …  BE  …  BE

app  app  app  app

# TBON-FS : the TBON File System

## Specialized TBON for distributed file access

- back-end data sinks/sources are files
- simplifies tool front-end development by providing an intuitive interface based on POSIX I/O
- custom tool back-end functionality via synthetic file systems loaded into TBON-FS servers

## Uses MRNet for:

- scalable unified name space composition
- scalable group file operations
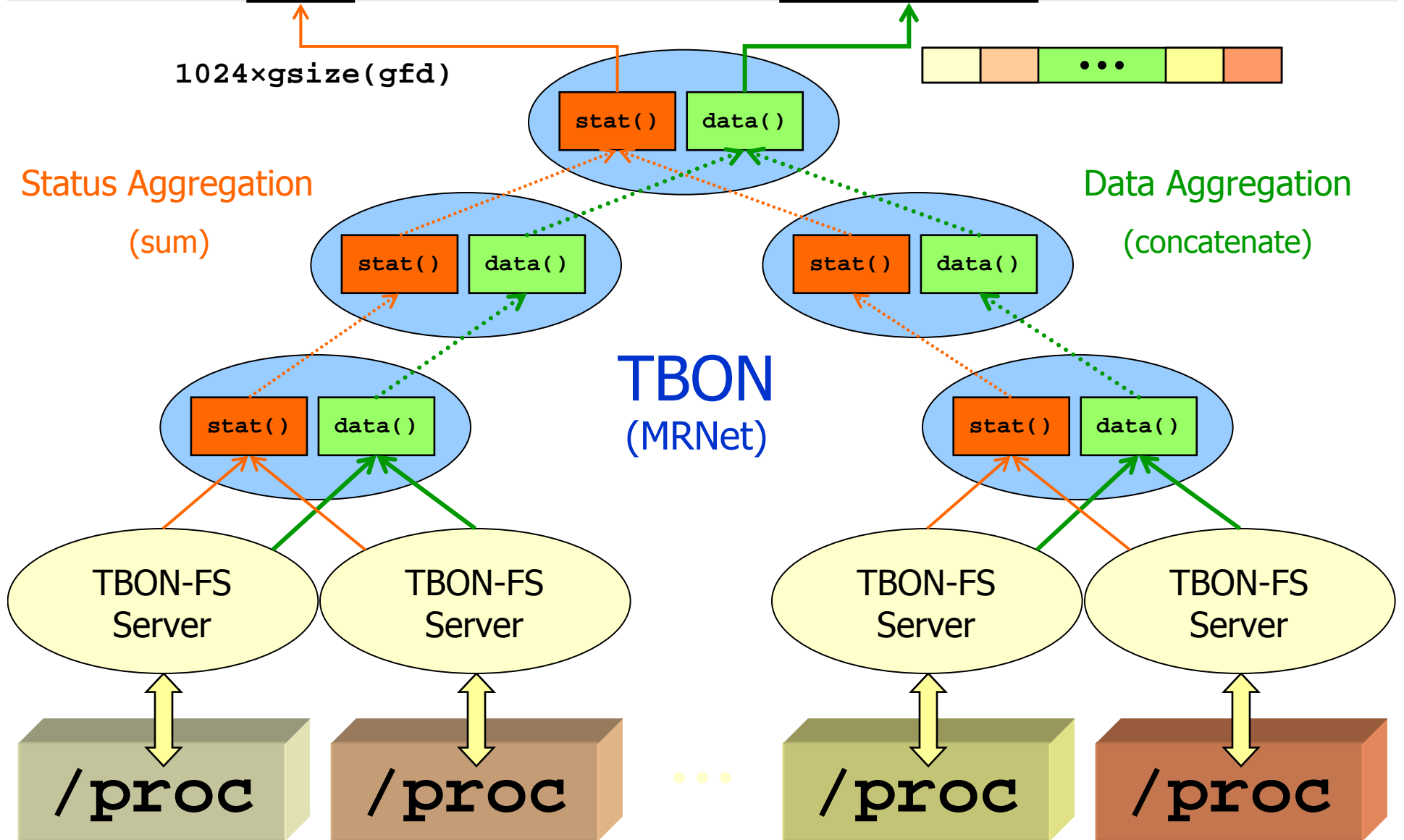
# Group File Operations

```
gfd = gopen(dir, flags, mode)
```

## Operating on Groups

o Use group file descriptor with regular file operations
   (e.g., `read` and `write`)

   • avoids iteration, one system call per group operation

o Semantics

   • operation applied to each group member

   • user-controlled aggregation of status and data results

# TBON-FS: Scalable Group File Operations

`int rc = read(gfd, databuf, 1024)`

1024×gsize(gfd)

Status Aggregation
(sum)

Data Aggregation
(concatenate)

stat()  data()

stat()  data()

stat()  data()

stat()  data()

stat()  data()

TBON
(MRNet)

TBON-FS Server

TBON-FS Server

TBON-FS Server

TBON-FS Server

/proc  /proc  ...  /proc  /proc

# Scalable Distributed Process Monitoring: `ptop`

```
ptop - Thu Apr 24 22:46:20 2008
1024 h
Tasks:
CPU: 4
Mem:
Swap: 17182572544k total,  71227968k used,17111344576k free, 200214
     USER           %CPU            %MEM           COMMAND
-----------      ------------    ------------    -----------------
    briml          1.52 @4096      0.05 @4096     tbonfs-server
    root           0.0  @928       0.00 @928      ksoftirqd/1
    root           0    @884       0.00 @884      ksoftirqd/2
    root                                              mainint
                                                      rqd/0
                                                      rqd/3
                                                      cd
                                                      g-ng
                                                      hald
                   0.01 @412       mungod
    root           0.00 @56        0.00 @56        ll_ping
    root           0.00 @1020      0.01 @1020      lrmmond
    root           0.00 @752       0.00 @752       irqbalance
    root           0.00 @68        0.00 @68        kqswnal_sched
    root           0.00 @1004      0.00 @1004      ldlm_cn_14
    root           0.00 @1008      0.00 @1008      ldlm_cn_15
```

> /proc/uptime     /proc/loadavg
> /proc/stat       /proc/meminfo

> /proc/$pid/stat
> /proc/$pid/statm
> /proc/$pid/status

Avg. %C
4096 proc

**4,096 files**

**>1,000,000 files**

# Group Process Control & Inspection

**/proc** : a good starting point

- o write to process/thread control file(s) to run/stop/signal
- o read files containing process/thread status
- o read/write process address space
- o read/write thread registers

But,

- o functionality differs by OS (e.g., no control on Linux)
- o no notion of group operations
- o always contains all host processes

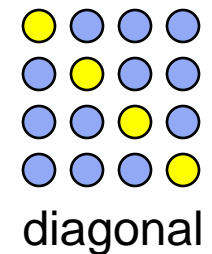# `proc++` : Synthetic File System for Process Control
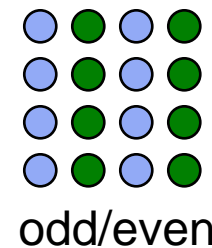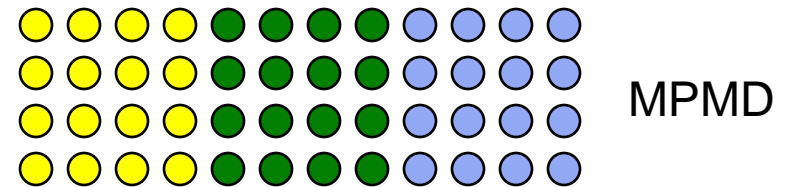
## Improvements over `/proc`

### 1. process/thread groups
- o explicit group management
- o directories containing members' control and inspection files automatically created

### 2. high-level debugger operations
- o breakpoints
- o stepping
- o stack walks

### 3. platform-independent interface

MPMD

odd/even          diagonal

| /proc | proc++ |
|---|---|
| `foreach(member){`<br>  `restore_insn()`<br>  `step_target()`<br>  `insert_bkpt()`<br>  `run_target()`<br>`}` | `run_group()` |

Example: Continue group from breakpoint

# `proc++` : from the makers of Dyninst

## Most capabilities provided by ProcControlAPI

- Cross-platform component library / C++ API
  - Linux, FreeBSD, BlueGene, Windows

- Process / thread control and inspection
  - Stop / continue processes, single-step threads
  - Read / write process memory, thread registers
  - Insert / remove breakpoints
  - Inferior remote procedure calls
  - Callbacks for asynchronous event notification

## Thread stack walks (StackwalkerAPI)

# TotalView Parallel Debugger

## Commercial debugger from Rogue Wave Software

o   Sequential, multi-threaded, and parallel programs

o   Fortran, C, C++ code from various compilers

o   pthreads, OpenMP, MPI, UPC

## 20+ years of engineering and HPC experience

o   Advanced MPI debugging

o   Built-in memory debugger

o   Reverse debugging (application DVR)

o   Recent support for GPGPU (CUDA) code

# TotalView is a great case study

## Most widely-used HPC debugger
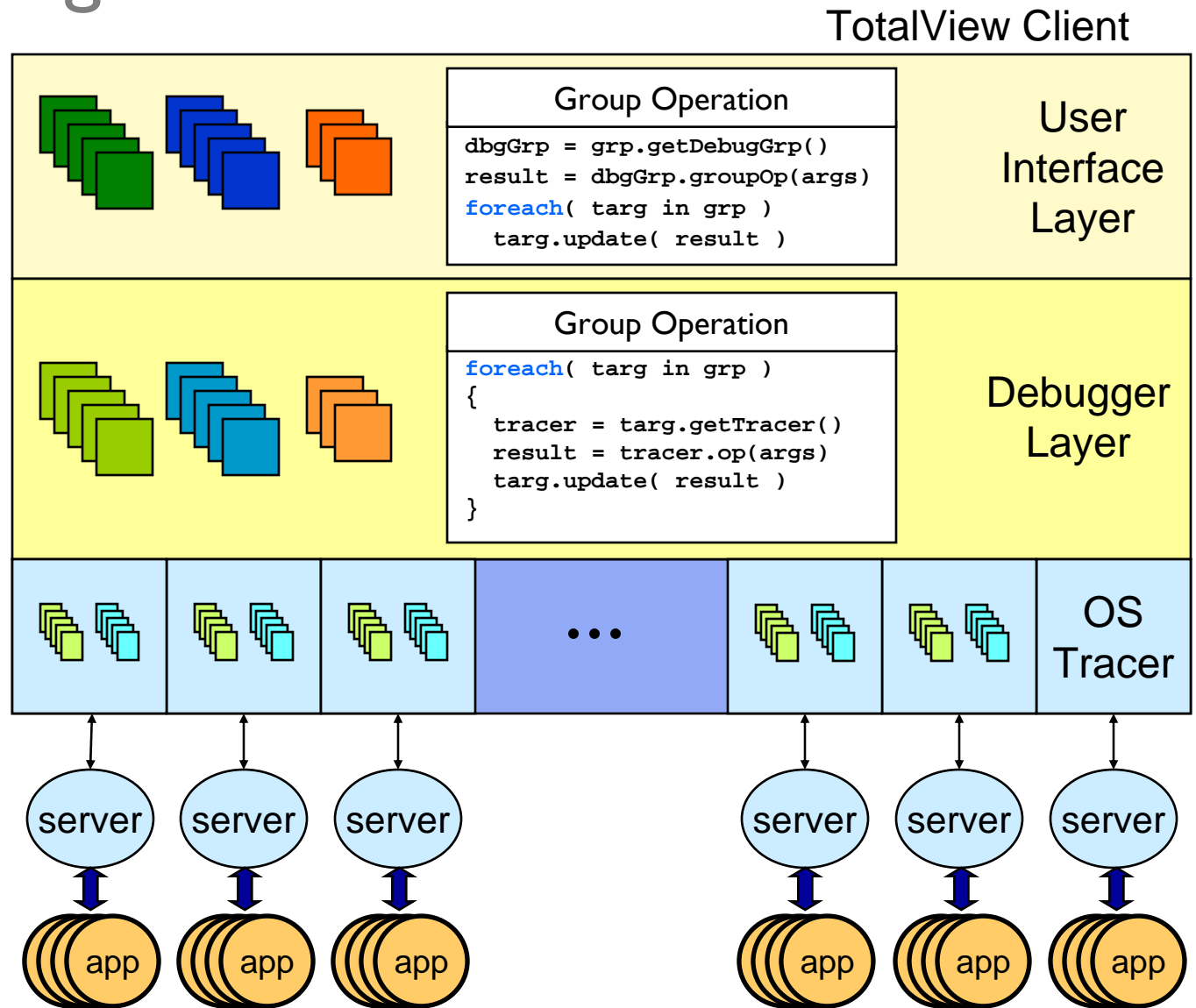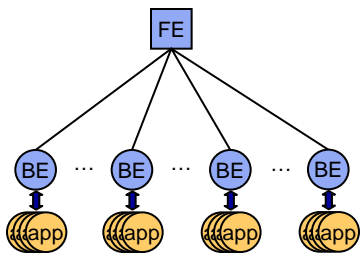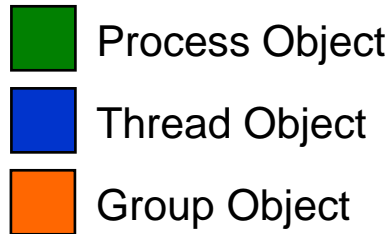
o   Lots of happy users

## Known scalability limitations

o   Lots of users that need it to work at full-scale on largest systems
    (i.e., @ 200K+ processes)

## 20+ years of engineering

o   A real tool that works on real applications

o   Modular architecture that evolved over time

o   Operations on process and thread groups are primary focus

# TotalView: Original Architecture

## User Interface Layer

**Group Operation**

```
dbgGrp = grp.getDebugGrp()
result = dbgGrp.groupOp(args)
foreach( targ in grp )
  targ.update( result )
```

- Process Object
- Thread Object
- Group Object

## Debugger Layer

**Group Operation**

```
foreach( targ in grp )
{
  tracer = targ.getTracer()
  result = tracer.op(args)
  targ.update( result )
}
```

## OS Tracer

· · ·

FE

BE  ···  BE  ···  BE  ···  BE

app  app  app  app

server  server  server  server  server  server

app  app  app  app  app  app

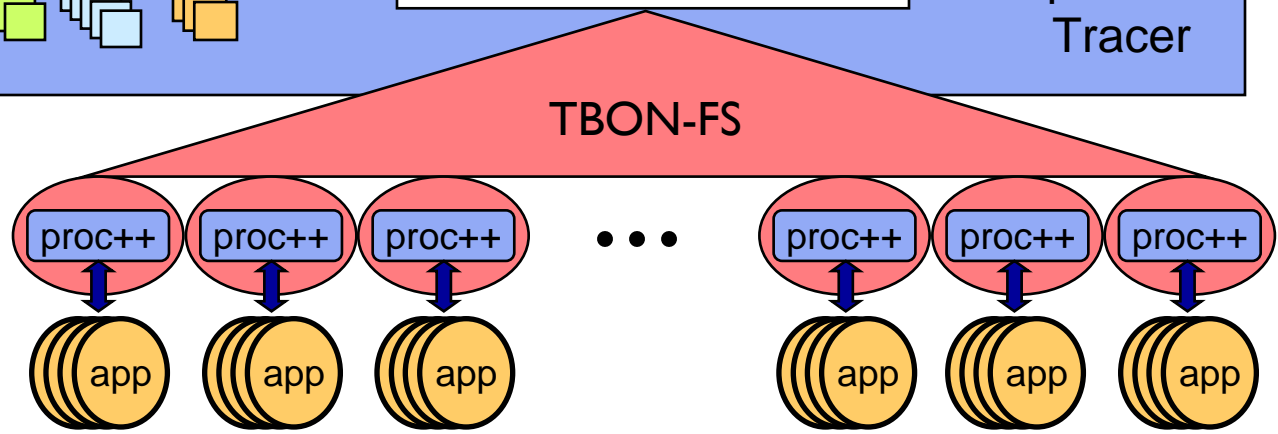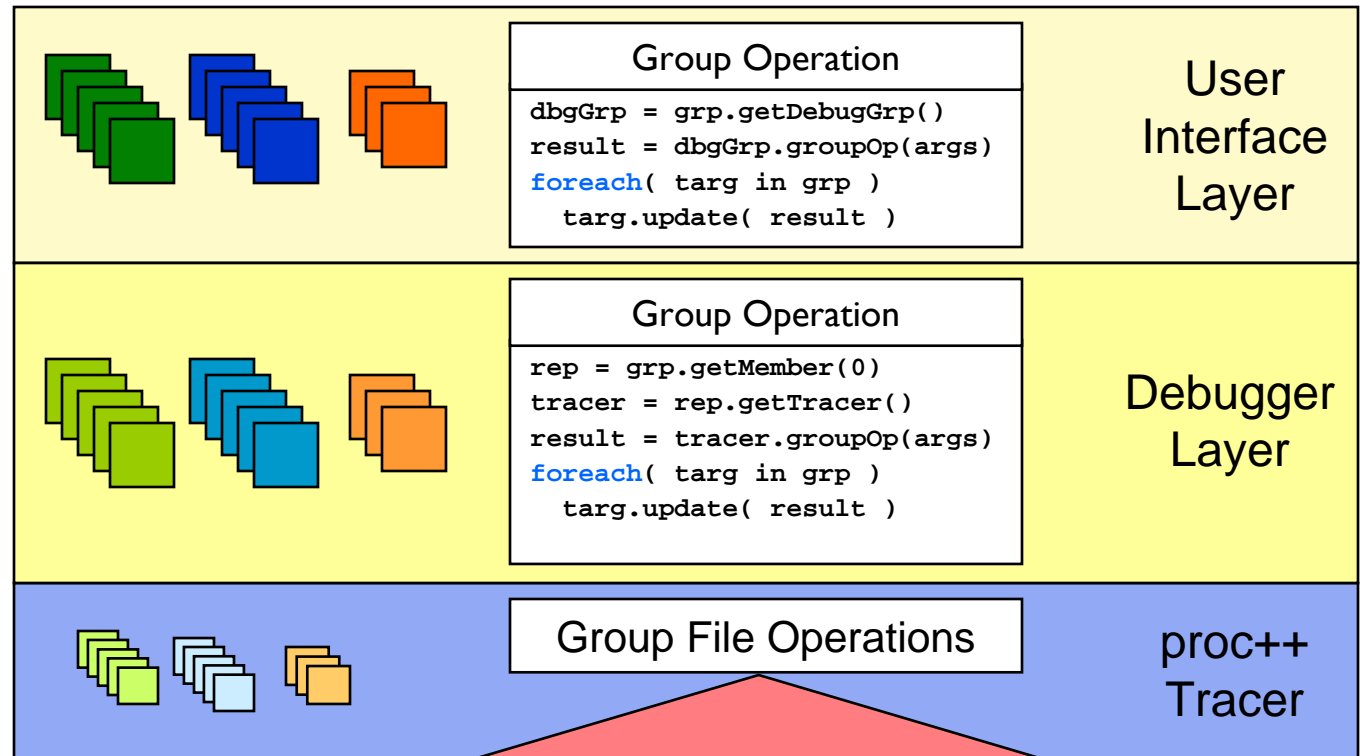# TotalView Integration Challenges

## Group Operations

o no group operations at (lowest) tracer level

- pushed groups down to use group file operations

o some group operations at UI level use iteration

- added group operations at debugger level

o some group operations require process- or thread-specific context

- extended proc++ interface and capabilities
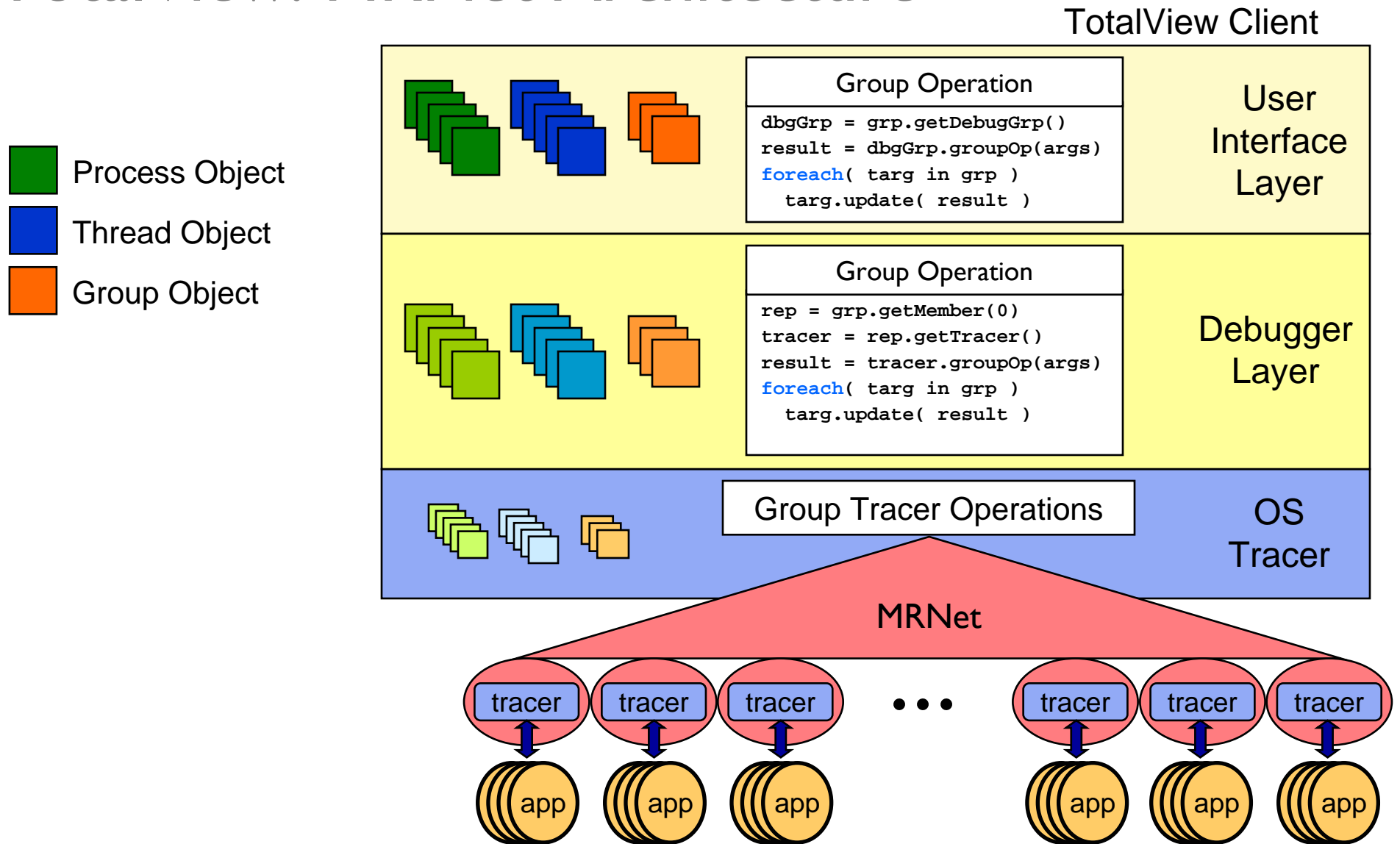
## Multi-level object maintenance

# TotalView: TBON-FS Architecture

TotalView Client



**User Interface Layer**

Group Operation

```
dbgGrp = grp.getDebugGrp()
result = dbgGrp.groupOp(args)
foreach( targ in grp )
    targ.update( result )
```

**Debugger Layer**

Group Operation

```
rep = grp.getMember(0)
tracer = rep.getTracer()
result = tracer.groupOp(args)
foreach( targ in grp )
    targ.update( result )
```

**proc++ Tracer**

Group File Operations

TBON-FS

Process Object
Thread Object
Group Object

proc++   proc++   proc++   • • •   proc++   proc++   proc++

app   app   app   app   app   app

# TotalView: MRNet Architecture

TotalView Client

### User Interface Layer

Process Object

### Group Operation

```
dbgGrp = grp.getDebugGrp()
result = dbgGrp.groupOp(args)
foreach( targ in grp )
  targ.update( result )
```

Thread Object

Group Object

### Debugger Layer

### Group Operation

```
rep = grp.getMember(0)
tracer = rep.getTracer()
result = tracer.groupOp(args)
foreach( targ in grp )
  targ.update( result )
```

### OS Tracer

Group Tracer Operations

MRNet

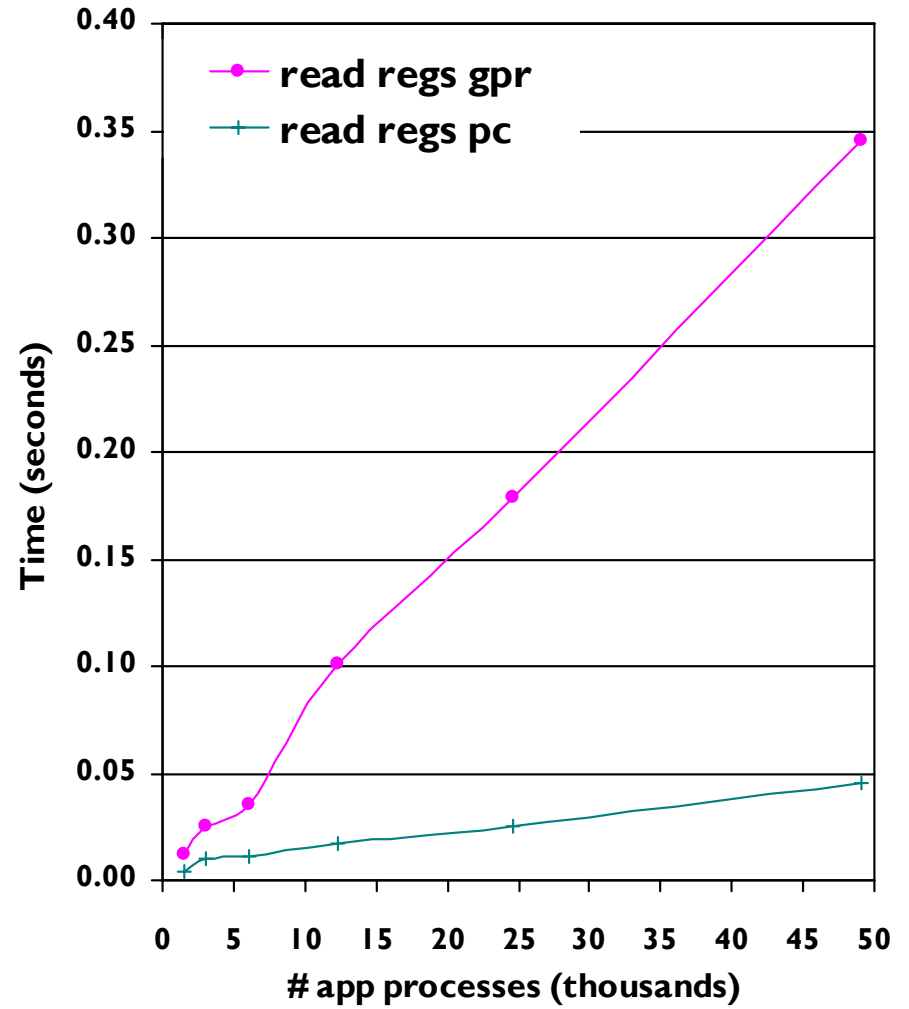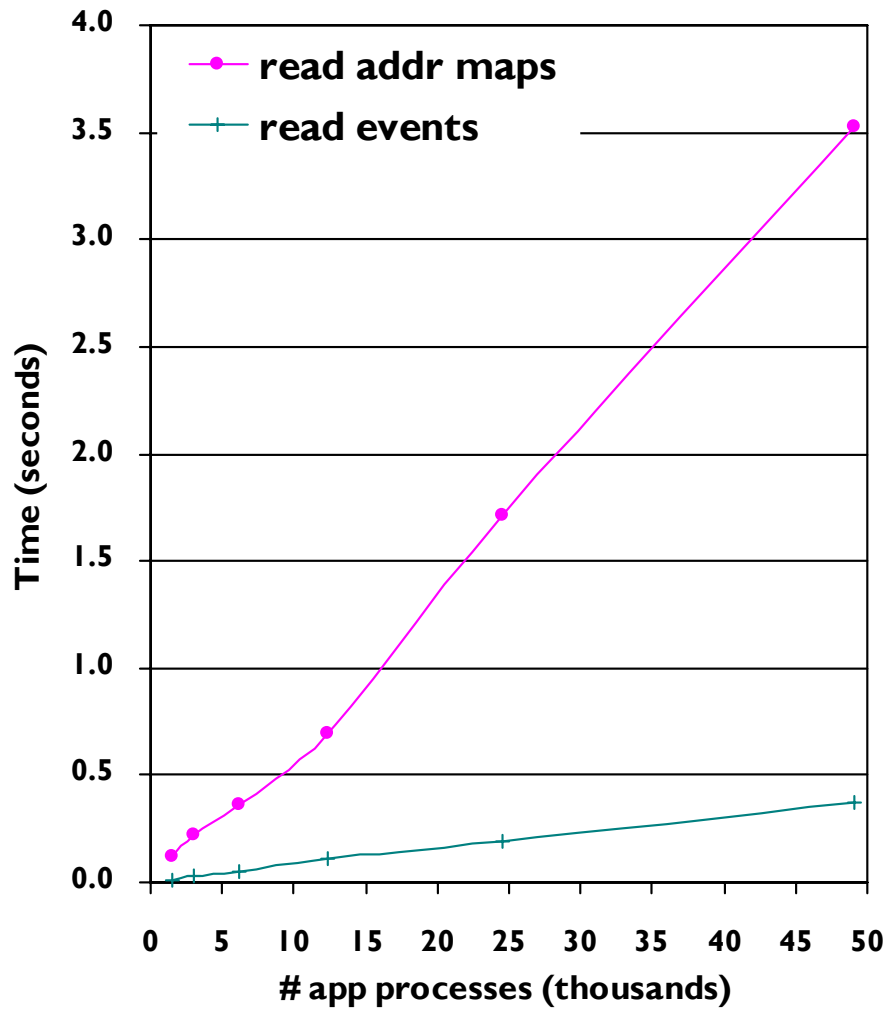| tracer | tracer | tracer | • • • | tracer | tracer | tracer |

| app | app | app | | app | app | app |

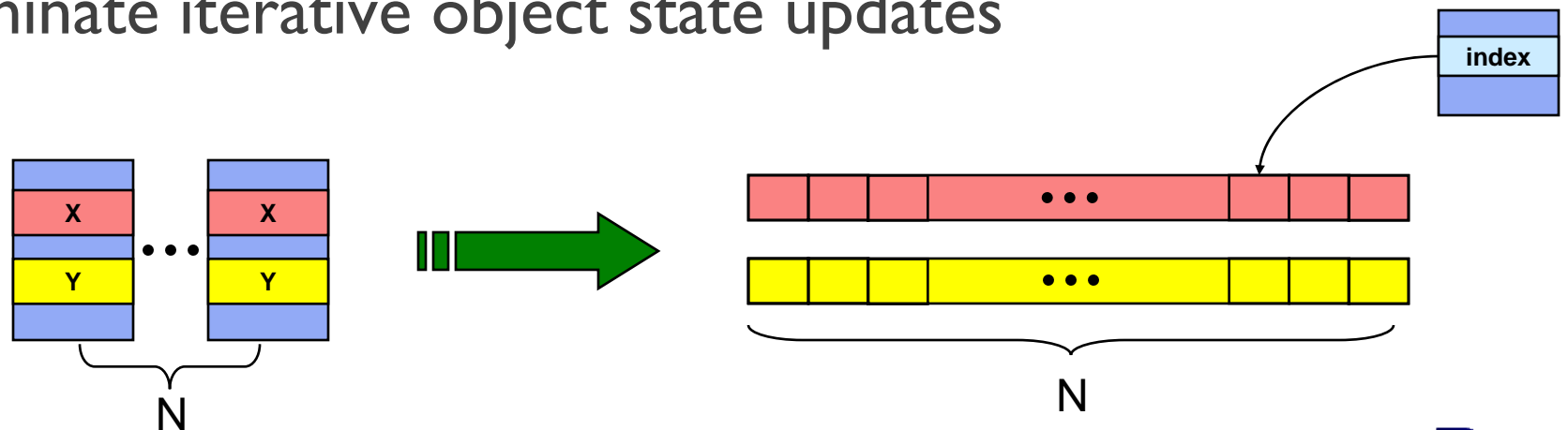# Scalability: proc++ group writes

# Scalability: proc++ group reads

# Amdahl's Law for Scalable Tools

Speed-up from using scalable group file operations is limited by front-end sequential behaviors

o reduce the number of objects per target
o reduce the state kept in those objects
o eliminate iterative allocation of objects
o eliminate iterative object state updates

# Keys to Real Tool Scalability

## "iteration is the bane of scalability"- me

o  any operation requiring a linear number of steps is a show-stopper

1. **Limited sequential behavior in tool front-end**
2. **Good group representation**
   - efficient creation and update $\Rightarrow$ distributed group state
3. **Constant or logarithmic time group operations**
   - parallel execution across group members
4. **Constant or logarithmic size data at tool front-end**
   - tool internal state: O(# of groups), not O(# of targets)
   - user display of group data: scalable aggregation is necessary

# Tool Scalability "rules to live by"

1. Single-target operations must be efficient, but rarely used

2. On-demand data access (lazy evaluation)
   - do not collect or generate data that is never used

3. Data Caching
   - individual target data at tool front-end is a bad idea
     - leads to iterative cache invalidation and update
     - see rule #2
   - individual target data at tool back-ends is a time/space tradeoff
   - group data at tool front-end is a time/space tradeoff
     - caching within a TBON can limit both time and space

# Questions?

## Group File Operations & TBON-FS

o International Conference on High Performance Computing (HiPC 2009) Best Paper

o ftp://ftp.cs.wisc.edu/paradyn/papers/Brim09GroupFile.pdf

## Scalable Composition of File System Name Spaces

o International Workshop on Runtime and Operating Systems for Supercomputers (ROSS 2011)

o ftp://ftp.cs.wisc.edu/paradyn/papers/Brim11FinalNamespace.pdf

MRNet : http://www.paradyn.org/mrnet/

TBON-FS or proc++ Source Code (talk to me)

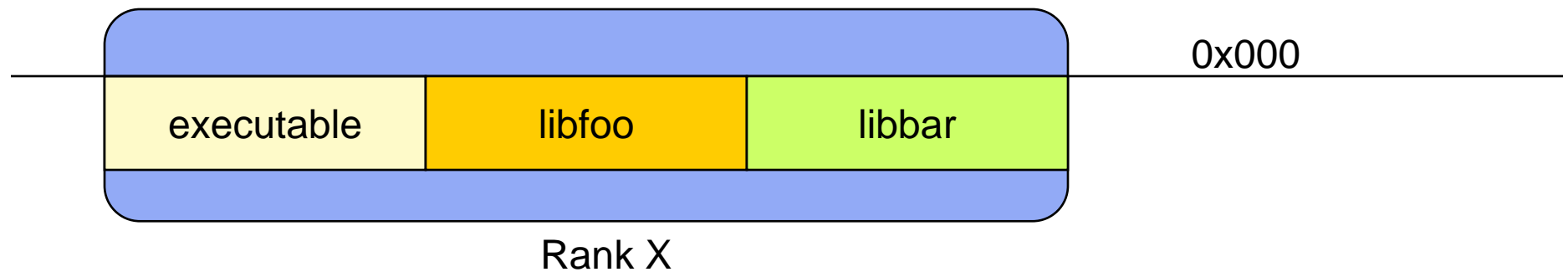# TotalView Integration: `proc++` Extensions

Problem: dynamic address space mappings



How can we do group address space write/read?

# TotalView Integration: `proc++` Extensions

Solution: image files that hide dynamic mappings



o one file for each mapped code image

o zero offset corresponds to map base of image

o to read / write symbols in image, seek to the symbol offset