

# Capriccio: Scalable Threads for Internet Services

F Z u G . N cula a E c w  
pu Sc c D s  
U s f alf a k l  
{j jc zf cula w }@cs. k l . u

## ABST ACT

This paper presents a new approach to scalable thread-based servers. The new approach uses a simple per-thread lock-free queue to support a large number of threads. The new approach is based on a simple per-thread lock-free queue. The new approach is based on a simple per-thread lock-free queue.

By using a simple per-thread lock-free queue, the new approach can support a large number of threads. The new approach is based on a simple per-thread lock-free queue. The new approach is based on a simple per-thread lock-free queue.

The new approach is based on a simple per-thread lock-free queue. The new approach is based on a simple per-thread lock-free queue. The new approach is based on a simple per-thread lock-free queue.

## 1 Introduction

The new approach is based on a simple per-thread lock-free queue.

## 2

The new approach is based on a simple per-thread lock-free queue.

The new approach is based on a simple per-thread lock-free queue.

## 3

The new approach is based on a simple per-thread lock-free queue. The new approach is based on a simple per-thread lock-free queue. The new approach is based on a simple per-thread lock-free queue.

The new approach is based on a simple per-thread lock-free queue. The new approach is based on a simple per-thread lock-free queue. The new approach is based on a simple per-thread lock-free queue.

The new approach is based on a simple per-thread lock-free queue.

The new approach is based on a simple per-thread lock-free queue. The new approach is based on a simple per-thread lock-free queue. The new approach is based on a simple per-thread lock-free queue.

The new approach is based on a simple per-thread lock-free queue.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SOSP'03, October 19–22, 2003, Bolton Landing, New York, USA.

Copyright 2003 ACM 1-58113-757-5/03/0010 ...\$5.00.

I ee r t rea pa a e a impr e per rma e  
e isti t rea e appi ati s it itt et m i ati  
t e appi ati itse .

## Th P

I t e pr ess i t rea s r sei ser er appi a  
ti s e t ata ser ee appr a isesse tia . ie  
ser e e t rea s a er e t rea s are t se t e  
s e ame ta i ere t pr ems. Ker e t rea s are  
primari se rea i tre rre ia m tipe  
e i es is req ests r Us. User e e t rea s are rea  
i a t rea s t at s pr i e a ea pr rammi  
m e it se i ari a ts a sema ti s.

T ate e t str a ate a *articular* se  
ma ti s r t rea s rat er ear et at a ea se  
ma ti s r t rea s req ires *decou ling* t e t rea s t e  
pr rammi m e ( i a t rea s) r m t se t e  
er i ere e. e pi t e pr rammi m e r m  
t e re is imp rta t r t reas s. irst t ere is s  
sta tia ariati i i ter a es a sema ti s am m  
er er es espite t ee iste e t e O IX sta ar .  
e er e t rea s a as r s I/O i ter a es are  
areas a ti e resear 19 20. T era e sema ti s  
a t e rate e ti t req ire e pi : i a  
t rea s a i e t O ariati a ere e ti .

I r ase t is e pi as pr i e a m er a  
a ta es. e a e ee a et i te rate mpi ers pp rt  
i t r t rea pa a e a e a eta e a a ta e se  
era e ere eat res. T s e a e ee a et i rease  
per rma e impr e s a a iit a a ress appi ati  
spe i ee s a it t a i appi ati e.

## 2

T is paper is sses r e t rea pa a e apr i i .  
T ist rea pa a e a ie es r as it t e ep t ree  
e eat res:

irst e impr e t e s a a iit asi t rea pera  
ti s. ea mp is e t istas si ser e e t rea s  
it perati es e i ta i a a ta e a e  
as r s I/O i ter a e a e i eeri r r time  
s stem s t at a t rea perati s are  $O(1)$ .

e e i tr e *linked stacks* a me a ism r  
ami sta r t t at s est e pr em sta a  
ati r ar e m ers t rea s. Tra iti a t rea s s  
tems prea ate ar e s mem r rea t rea s  
sta i se ere imits s a a iit. apr i i ses a  
m i ati mpi e time a a sis a r time e s  
t imit t e am t aste sta spa e i a effie ta  
appi ati spe i ma er.

i a e esi e a *resource-a are sc eduler* i  
e tra t si r mati a t t e tr it i a pr  
ram i r er t ma e s e i e isi s ase pre  
i te res r e sa e. T is s e i te iq e ta es a  
a ta e mpi ers pp rta perati e t rea i t  
a ress appi ati spe i ee s it t req iri t e pr  
rammer t m i t e ri i a pr ram.

T e remai er t is paper is sses ea t ese t ree  
eat res i etai. T e e pre se ta era e perime ta  
e a ati r t rea pa a e. i a e is ss  
t re ire ti s r ser e e t rea pa a es it i te rate  
mpi ers pp rt.

## 2 TH EA ES A SCALAB L TY

apri i isa ast ser e e t rea pa a et ats pp rts  
t e O IX A I r t rea ma a eme ta s r iza  
ti . I t is se ti e is sste era esi r r  
t rea pa a e a e em strate t at it satis es r  
s a a iit as.

## 2 -L v Th

O e t e rst iss es ee p re e esi i apr i  
i as et ert emp ser e e t rea s r er e t rea s.  
User e e t rea s a es me imp rta ta a ta es r t  
per rma e a e i iit . U r t ate t e as m  
pi ate preempti a a i ter a a it t e er e  
s e er. U timate e e i e t at t e a a ta es  
ser e e t rea s are si i a t e t arra t t e a  
iti a e i eeri req ire t ir m et t eir ra a s.

### 2.1.1 Flexibility

User e e t rea s pr i e a treme sam t e i  
iit r s stem esi ers reati a e e i ire ti  
et ee appi ati sa t e er e. T is a stra ti eps  
t e pet et a ita s aster i ati t  
si es. re ampe apr i i is apa e ta i a a ta e  
t e e as r s I/O me a ism t e e e p me t  
series Li ere e i a s st pr i e per rma e  
impr eme ts it t a i appi ati e.

T e se ser e e t rea s a si reases t e e i iit  
t e t rea s e er. Ker e e e t rea s e i m st  
e e era e t pr i e a reas a e e e q ait  
ra appi ati s. T s e r e t rea s a t t ai r t e  
s e i a rit m t t a spe i appi ati . r t  
ate ser e e t rea s t s er r m t is imitati .  
I stea t e ser e e t rea s e er a e it a  
it t e appi ati .

User e e t rea s are e treme i t e i t i a s  
pr rammer s t se a treme s m er t rea s it  
t rri a t t rea i er ea . T e e mar s  
i e ti 2.3 s t at apr i i a s a e t 100,000  
t rea s t s apr i i ma es it p ssi et rite i  
rre t appi ati s ( i are te ritte it mess  
e e t rie e) i a simp et rea e st e.

### 2.1.2 Performance

User e e t rea s a reat re et e er ea t rea  
s r izzati . I t e simp est ase perati es e  
i a si e U s r izzati is ear ree si e  
eit er ser t rea s r t e t rea s e er a e i ter  
r pte i e i a riti a se ti .<sup>1</sup> I t e t re e e i e e  
t at e i e ser e e s e i a mpi e time a a sis  
i a st er simi ar a a ta es a m ti U  
ma i e.

e i t e ase preempti et rea i ser e e t rea s  
era a a ta e i t at t e t req ire er e r ss  
i s r m te a q isiti r re ease. B mparis ere  
e e m t a e si req ires a er e r ssi r e er  
s r izzati perati . i e t is sit ati a e im  
pr e r t e e s<sup>2</sup> i t e e m te es  
sti req ire er e r ssi s.

<sup>1</sup> r esi e si a a i e a rei tr et ese  
pr ems t t is pr em a easi e a i e .

<sup>2</sup> T e futexes i re e t Li ere s a perati s  
t e e m te est r e tire i ser spa e.

ia mem r ma a eme t is m re effi e t it ser  
e e t rea s. Ker e t rea s req ire ata str t res t at  
eat p a a e er e a resspa e e reasi t e spa e  
a aia e r I/O ers e es ript rs a t er res r es.

### 2.1.3 Disadvantages

User e e t rea i is t it t its ra a s  
e er. I r er t retai tr t epr ess r e a ser  
e e t rea e e tes a i I/O a a ser e e t rea  
i pa a e erri est ese i a sa rep a est em  
i ter a it i eq i a e ts. T e sema ti s  
t ese i I/O me a isms e era req ire a  
i rease m er er e rssi s e mpare t t e  
i eq i a e ts. r e ampe t e m st effi e t  
i et r I/O primiti e i Li (epoll) i es  
rst p i s ets r I/O rea i ess a t e per rmi  
t e a t a I/O a . T ese se I/O a s are i e ti a  
t t se per rme i t e i ase t e p a s are  
a iti a er ea . N i is I/O me a isms are  
te simi ar i t at t e mp separate s stem a s t  
s mit req ests a re trie e resp ses.<sup>3</sup>

I a iti ser e e t rea pa a es m st i tr e a  
r apper a er t at tra sates i I/O me a isms t  
i I/O es a t is a er i s a t e r s r e  
er ea . At est t is a er a e a er t i s im i  
simp a s a e e tra ti a s. H e er r q i  
perati s s a i a e rea s t at are easi satis e  
t e er e t is er ea a e me imp rta t.

ia ser e e t rea i a ma e it m re effi t  
t ta e a a ta e m tipe pr ess rs. T e per rma e  
a a ta e i t e i ts r izati is imi is e e  
m tipe pr ess rs are a e si es r izati is  
er r ree . A iti a as is sse A ers  
et a . i t eir r s e er a ti ati s p re ser  
e e s r izati me a isms are i e t i e i t e a e  
tr e rre a ma ea t star ati 2.

U timate e e i e t e e e ts ser e e t rea i  
ar t e i t ese isa a ta es. As t e e mar s i  
e ti 2.3 s t e a iti a er ea i rre es  
t seem t e a pr em i pra ti e. I a iti e are  
r i ast er me t e effi ties it m tipe  
pr ess rs e i is sst is iss e r t e r i e ti 7.

## 22

e a e imp eme te apr i i as a ser e e t rea i  
i rar r Li . apr i i imp eme ts t e O IX t rea  
i A I i a s i t t r m st appi ati s it t  
m i ati .

Co x w ch . apr i i is it t p ar  
T e r i s r t i e i rar 32. T is i rar pr i es e  
treme ast t e ts it es r t e mm ase i i  
t rea s tari ie eit ere p i it r t r ma i  
a i I/O a . e are rre t esi i si a  
ase et at a s r preempti r i ser

<sup>3</sup>At t ere are i I/O me a isms (s as  
O IX AIO s lio . listio () a Li s e io . submit ())  
t at a t e s missi m tipe I/O req ests it a  
si e s stem a t ere are t er iss es t at ma e t is  
eat re effi t t se. r e ampe imp eme tati s  
O IX AIO t e s er r m per rma e pr ems.  
A iti a se at i reates a tra e et ee  
s stem a er ea a I/O ate i is effi t t  
ma a e.

t rea s t apr i i es t pr i e t is eat re et.

I/ . apr i i i ter epts i I/O a s at t e i rar  
e e erri i t e s stem a s t ti si GNU  
i . T is appr a r s a ess r stati a i e  
appi ati s a r ami a i e appi ati s t at  
se GNU i ersi s 2.2 a ear ier. H e er GNU i  
ersi 2.3 passes t e s stem a s t s r ma its  
i ter a r ti es (s as printf) i a ses pr ems  
r ami a i e appi ati s. e are r i t a  
apr i i t ti as a i a i r er t pr i e  
etter i te rati it t e at est ersi s GNU i .

I ter a apr i i se t e at est Li as r s  
I/O me a isms epoll r p a e e es ript rs (e .  
s ets pipes a s) a Li AIO r is . I t ese  
me a isms are t a aia e apr i i as a t e  
sta ar U i poll () a r p a e es ript rs a a  
p er e t rea s r is I/O. Users a se e tam  
t e a aia e I/O me a isms setti appr priate e i  
r me t aria es pri r t starti t eir appi ati .

ch dul . apr i i s mai s e i p s er  
m i e a e e t r i e appi ati a ter ate r i  
appi ati t rea s a e i r e I/O mp eti s.  
N t e t t at t e s e er i es t i e e t r i e  
e a i r r m t e pr rammer sti ses t e sta r  
t rea ase a stra ti . apr i i as a m ar s e  
i me a ism t a t a s t e ser t easi se e t et ee  
i ere t s e ers at r time. T is appr a as a s  
ma e it simpe r st e e p se era i ere t s e ers  
i i a e s e er ase t rea res r e ti za  
ti . e is sst is eat re i etai i e ti .

ch o z o . apr i i ta es a a ta e per  
ati e s e i t impr es r izati . At prese t  
apr i i s pp rts perati et rea i si e U ma  
i es i i ase i t e r t rea s r izati primiti es  
req ire simpe e s a ea e / e a .  
r ases i i m tipe er e t rea s are i e  
apr i i emp se i t er spi s r ptimisti r  
re tr primiti es epe i i me a ism  
est ts t e sit ati .

Effic c . I e e pi apr i i e a e ta e reat  
are t se effi e ta rit ms a ata str t res.  
seq e t a t e apr i i s t rea ma a eme t  
ti s as a e rst aser i time i epe  
e t t e m er t rea s. T e s e e epti i s t e  
s eep q e e i rre t ses a ai e i e ist impe  
me tati . i e t e iterat re tai s a m er  
a rit ms r effi e t s eep q e es r rre t imp eme  
tati as t a se pr ems et s e a e se r  
e e p me t e rts t er aspe ts t e s stem.

## 23 The M h k

era a m er mi r e mar s t ai ate apr i  
i s esi a imp eme tati . O r test p at rm as a  
it t 2. GHz Xe pr ess rs 1 GB mem r  
t 10K R I U tra II ar ri es a 3 Gi a it  
t e r e t i t e r a es. T e perati s stem as Li 2.5.70  
i i es s pp r t r epoll as r s is I/O  
a i t e i t s stem a s (vsyscall). era r e  
mar s t r e e t rea pa a es: apr i i Li T rea s  
(t e sta ar Li er e t rea pa a e) a N TL  
ersi 0.53 (t e e Nati e O IX T rea s r Li  
pa a e). e it a appi ati s it 3.3 a i e

	Capriccio	Capriccio_notrace	LinuxThreads	NPTL
Thread creation	21.5	21.5	37.9	17.7
Thread context switch	0.56	0.24	0.71	0.65
Uncontended mutex lock	0.04	0.04	0.14	0.15

Table 1: Load of the different lock.

against GNU 2.3. The Linux Thread set of the system is tested. The results are compared with NPTL. The results are:

## 24 The Pipe

The comparison was made using the primitive of the Linux Thread. The test is the Capriccio\_notrace and the statistics are taken from the (see the results in the section) of the impact of the pipe. The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread. The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread. The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread.

## 25 The Socket

The network was efficient. The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread. The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread. The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread.

## 26 / Pipe

The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread. The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread. The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread.

The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread. The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread. The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread.

The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread. The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread. The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread.

The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread. The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread. The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread.

The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread. The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread. The results are: the standard deviation is 10% less than the standard deviation of the Linux Thread.

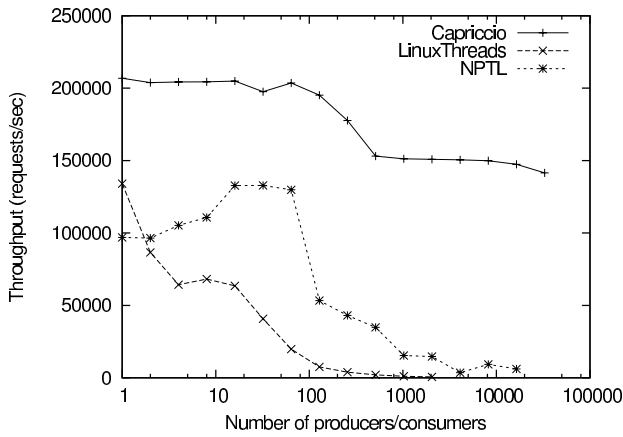


Figure 1: Producer-Consumer - chdul chozofo c.

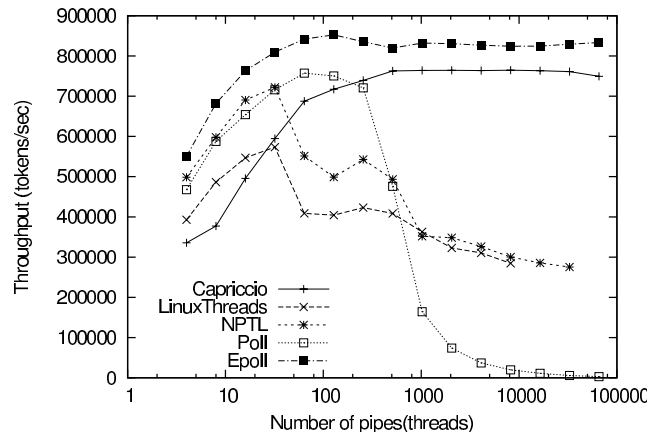


Figure 2: Pipetest - wokc1bl.

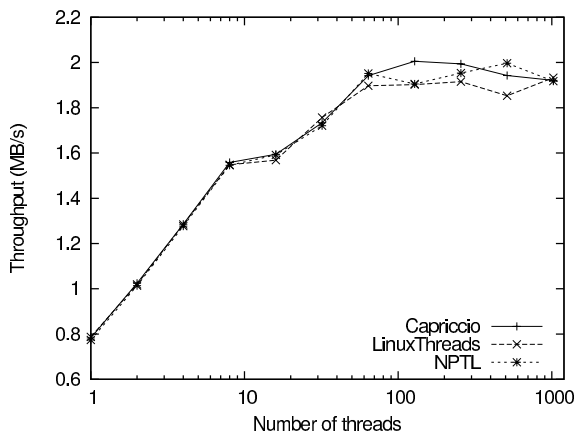


Figure 3: Bfi of dkh d chdul.

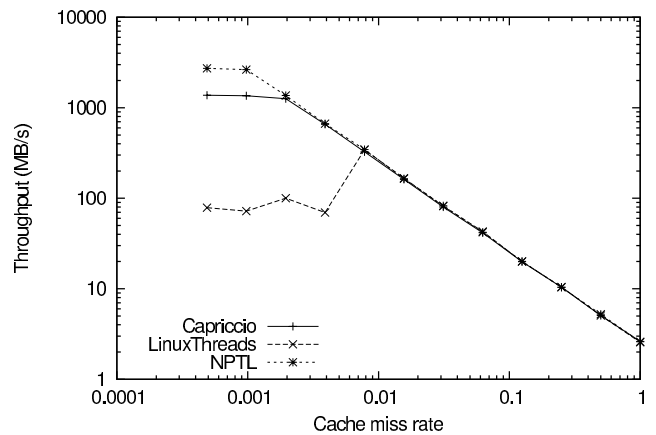


Figure 4: DkI/ fo c w h buff c ch.

am t er ea r a e itti perati sa es t at rea est e is : rea I/O req est a mpeti e et eest e str te qee a ei ere t ser e e tr a separate stem a . H e er t is s r t mi is re ati e eas t : re r i t res t imme iate r req ests t at t ee t ait e a e imi ate m st (i t a ) t is er ea . e ea et is m i ati as t re r . i a as rprisi res t is t at Li T rea s per rma e e ra es si i a t at a er miss rate . e e ie et is e ra ati is a res t a eit er i te er e r i te i rar si e te pr ess r is m st i e r i t e test .

### 3 L E STAC MA A EME T

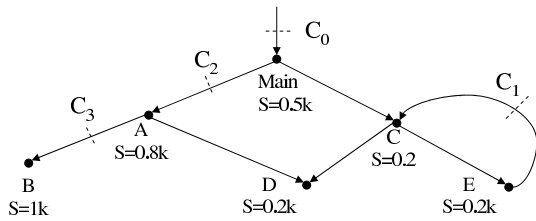
T rea pa a es sa attempt t pr i e t e pr ram mer it t ea str ati a e a sta rea t rea . I rea it t esta size is e tte s are se ser ati e s t at t ere is p e t spa e r rma pr ram e e ti . re ampe Li T rea s a ates t me a tes per sta ea t it s a ser ati e a ati s em e e s me l GB irt a mem r r sta spa e it st 500 t rea s . rt ate m st t rea s s me a e i tes sta spa e

ata i e time at t e mi t tr sta es e te se si era m re . T is ser ati s ests t at e a si i a t re e te size irt a mem r e i ate t sta si e a pta ami sta a ati pi i ere i sta spa e is a ate t t rea s ema i re ati e sma i reme ts a is ea ate e te t rea req ires ess sta spa e . I t e rest t is se ti e is ssa mpi er eat re t at a s s t pr i es a me a ism i e pr e s e r i t e pr ram mi a str ati e sta s .

### 3 A L k S k

O r appra ses a mpi e r a a sist imitt e am t sta spa e t at m st e prea ate . e per rm a e pr ram a a sis ase a eig ted call gra .<sup>4</sup> a ti i t e pr ram is repre te a ei t is a rap ei te t ema im mam t sta spa e t at asi esta rame r t at ti i s me . A e e et ee e A a e B i i ates t at ti A a s ti B ire t . T s pat s et ee es i

<sup>4</sup> e se te IL t it 23 r t is p r p se i a s effie t e pr ram a a sis rea r appi ati s i e t e Apa e e ser er .



**Fig. 5: A** x l of c ll h o d w h ck f z . Th d k d w h  $C_i$  ( $=0, \dots, 3$ ) h ch ck o .

ti s rap rresp t seq e es sta ramest at ma appear t esta atr time. T e e t apat iste s m t e ei ts a es i t is pat t at is it iste t ta size t e rresp i seq e e sta rames. A e ampe s a rap iss i i re 5.

Usi t is a rap e is t pae areas a e team t sta spa et at i e s me ea t rea . I t ere are re rsi e ti si r pr ram t ere i e es i t e a rap a t s e a easi t e ma im m sta size r t e pr ram at mpi e time i t e est pat starti r m ea t rea s e tr p i t . H e er m st rea r pr rams ma e se re rsi i mea st at e a t mp t e a t e sta size at mpi e time. A e e i t e a se e re rsi t e stati mp tati sta size mi t e t ser ati e . re ampe si er t e a rap i i re 5. I ri t e e i t e rap t e ma im m sta size is 2.3 KB t e pat Main-A-B. H e er t e pat Main-C-D as a sma er sta size 0.9 KB. I t e rst pat is se ri i itia izati a t e se pat is se tr t e pr ram s e e ti t e a ati 2.3 KB t e a t rea e aste . r t e se reas s it is imp rta t t e a e t r a s ri t e sta size ema .

I r er t imp em e t ani a size sta s r a rap a a sis i e ti es a sites at i e m st i s e r t c eck oints. A e p i t is a sma pie e e t at e t e r m es e r t e re i se sta spa e e t t rea t e e t e p i t i t a si sta er . I t e spa e remai s a e stack c unk is a ate a t e sta p i t e r is a ste t p i t t t is e . e t e ti a re t r s t e sta is i e a re t r e t a re e ist.

T is s eme res t i ti s sta s t e a se t e sta s are s i t e ri t e re t e a t a ar me t s r a ti a are p s e t e r t e a e e e t e a e . A e a se t e a e r s r a m p i t e r is s t r e t e a e s sta r a m e e e r s a t e a t r a e a p r a m .<sup>5</sup> T e e r a e p i t is ritte i it a sma am t i i e assem r rea i a setti t e sta p i t e r t is e i s i s e r t e si a s r e t s r e t r a s r m a t i t e p r r a m p r i r t m p i a t i . t a e si r a e s s i t e r e e sta i s t i s e s r e r p e r a t i e t r e a i a p r a .

<sup>5</sup>T is s eme es t r e t e omit-frame-pointer is e a e i gcc. It is p s s i e t s p p r t t is p t i m i z a t i si m re e p e s i e e p i t p e r a t i s s a s p i t e a r m e t s r m t e a e r s r a m e t t e a e s r a m e .

## 32 P Ch k

ri r pr r a m a a sis e m s t e r m i e e r e t p a e e p i t s . A s i m p e s t i s t i s e r t e p i t s a t e e r a s i t e e e r t i s a p p r a i s p r i t i e e p e s i e . A e s s r e s t r i t i e a p p r a i s t e s r e t a t a t e a e p i t e a e a t e s t a s p a e t a t m a e s m e e r e r e a t e e t e p i t ( r a e a i t e a r a p ) .

T s a t i s t i s r e q u i r e m e t e m s t e s r e t a t t e r e i s a t e a s t e e p i t i e e r e e i t i t e a r a p ( r e a t a t t e e s i t e a r a p r r e s p t a s i t e s ) . T e a p p r i a t e p i t s t i s e r t e p i t s e p e r m a e p t r s t s e a r t e a e r a p i i e t i e s a e e s t a t i s e e s t a t e t a e t e i t s a e s t r s i t e a r a p 22 . A e s i t e r a p m s t t a i a a e e s e a e p i t s a t a a s i t e i e t i e a s a e e s i r e r t e s r e t a t a p a t r m a t i t a e p i t a s e e t . I i r e 5 t e e p i t  $C_0$  a a t e s t e r s t a a t e e p i t  $C_1$  i s i s e r t e a e e E-C .

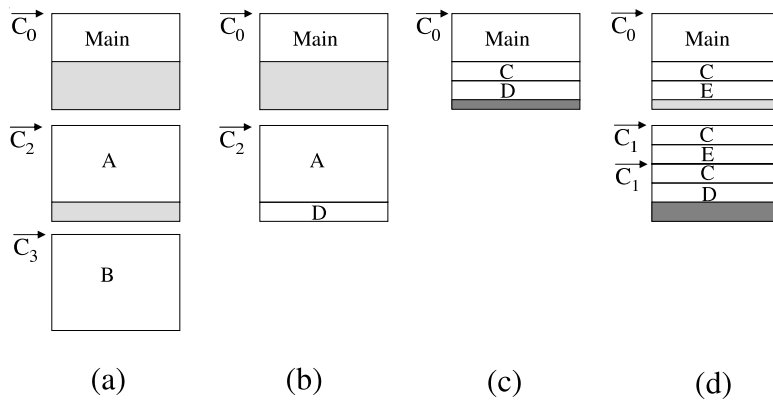
e a t e r e r e a a e s t e s s t a s i z e m a e t a r e . T s e i s t a a i t i a e p i t s t t e r a p i r e r t e s r e t a t a p a t s e t e e p i t s a r e i t i a e s i r e i i s i e a s a m p i e t i m e p a r a m e t e r . T i s e r t t e s e e e p i t s e p r e s t e a r a p e m r e t i s t i m e e t e r m i i t e e s t p a t r m e a e t e e t e p i t r e a . e p e r m i t i s a a s i s e s i e r a r e s t r i t e a r a p t a t e s t t a i a a e e s s i e t e s e e s a r e a a e e p i t s . T i s r e s t r i t e r a p a s e s s e a p r e s s t e e s t t m p t s e p r e s s i e n e i a e a r e a e t e r m i e t e e s t p a t r e a n s s e s s r s . r e a s e s s r s e n e t a e t e e s t p a t r s a a n . I t i s e p a t s e t e e e s t e s p e i e p a t i m i t p a r a m e t e r e a a e p i t t t e e e t e e n a s i e e t i e r e e s t e e s t p a t s t z e r . T e r e s t t i s a r i t m i s a s e t e e s e r e e p i t s s e a e a i t r e a s a e s t e m a i m m p a t e t r m e a e . r t e e a m p e i i r e 5 i t a i m i t 1 K B t i s a r i t m p a e s t e a i t i a e p i t s  $C_2$  a  $C_3$  . i t t t e e p i t  $C_2$  t e s t a r a m e s M a i n a A s e m r e t a 1 K B .

i r e s s r i s t a e s i t e i e t i m e t e t r e a s e a r a p i s s i i r e 5 . I i r e ( a ) t e t i B i s e e t i i t t r e e s t a s a a t e a t e p i t s  $C_0$   $C_2$  a  $C_3$  . N t i e t a t 0.5 K B i s a s t e i t e r s t a a 0.2 K B i s a s t e i t e s e

. I i r e ( ) t i A a s a e D a t s t a s e r e e s s a r . i a i i r e ( ) e s e e a i s t a e i t r e s i . A e s t a i s a a t e e E a s C ( a t e p i t  $C_1$  ) . H e e r t e s e t i m e a r t e e a t e p i t  $C_1$  e i e s t a t t e r e i s e s p a e r e m a i i i t e r r e t s t a t r e a e i t e r a e a t i ( D ) r t e e t e p i t (  $C_1$  ) .

## 33 h S

t i p i t e r s p r e s t a a i t i a a e e t r a r i t m e a s e e t a t m p i e t i m e e a t i t i m a e a e t r a i e t i p i t e r . T i m p r e t e r e s t s r a a s i s t e a t t e t e r m i e a s p r e i s e a s p s s i e t e s e t



F u 6: Ex l of d c lloc o d d lloc o of ck chu k .

ti st at mi t e a e at a ti p i t e r a site. r r e t e a t e r i z e t i p i t e r s m e r a t p e a r m e t s t i t e t r e e p a t s e a m r e s p i s t a t e p i t e r a a s i s .

a s t e t e r a t i s a s a s e p r e m s s i e i t i s m r e i f f i t t t e s t a s p a e s e p r e m p i e i r a r i e s . e p r i e t s t i s t t i s p r e m . i r s t e a t e p r a m m e t a t a t e e t e r a i r a r t i s i t t r s t e s t a s . A t e r a t i e e a a r e r s t a s t e i e r e t e r a t i s a s a s t r e a s t r e q e t i t e s e t i t e s e t i s e a r e s e a s m a m e r a r e s t a s t r t t e a p p i a t i . r t e s t a a r i r a r e s e a t a t i s t e a i t t i s t a t r t i s t a t a r e r e q e t a e t e s e a t a t i s e r e e r i e a a z i i r a r e .

### 34 T h A h

O r a r i t m a s e s t a s p a e t e a s t e i t p a e s . i r s t s m e s t a s p a e i s a s t e e a e s t a i s i e e a t i s s p a e *internal* a s t e d s a c e . e s t a s p a e a t t e t t m t e r r e t i s s i e r e s e t i s s p a e i s a e *external* a s t e d s a c e . I i r e i t e r a a s t e s p a e i s s i i t r a e r e a s e t e r a a s t e s p a e i s s i a r r a .

T e s e r i s a e t t e t p a r a m e t e r s t a t a s t t e t r a e s i t e r m s a s t e s p a e a e e t i s p e e . i r s t t e s e r a a s t *MaxPat* i s p e i e s t e m a i m m e s i r e p a t e t i t e a r i t m e a e s t e s r i e . T i s p a r a m e t e r a e t s t e t r a e e t e e e e t i t i m e a i t e r a a s t e s p a e a r e r p a t e t s r e q i r e e e r e p i t s t m r e s t a i i . e t e s e r a a s t *MinC unkt* e m i m m s t a s i z e . T i s p a r a m e t e r a e t s t e t r a e e t e e s t a i i a e t e r a a s t e s p a e a r e r s r e s t i m r e e t e r a a s t e s p a e t e s s r e q e t s t a i i i i t r r e s t s i e s s i t e r a a s t e s p a e a a s m a r e e e t i t i m e e r e a . O e r a t e s e p a r a m e t e r s p r i e a s e m e a i s m a i t e s e r ( r t e m p i e r ) t p t i m i z e m e m r s a e .

### 35 M B fi

O r i e s t a t e i q e a s a m e r a a t a e s i t e r m s m e m r p e r r m a e . I e e r a t e s e e e t s a r e a i e e i r i t r e a i m p e m e t a t i r m e r e

m e a i s m s t s i m p r i r a i t t t e i i i a a p p i a t i m e m r s a e . m p i e r t e i q e s m a e t i s a p p i a t i s p e i t i p r a t i a .

i r s t r t e i q e m a e s p r e a t i a r e s t a s e e s s a r i i t r r e e s i r t a m e m r p r e s s r e e r i a r e m e r s t r e a s . O r a a s i a i e e s t i s a i t t t e s e a r p a e s i t r i t e e e s s a r e r e r s s i s a i r t a m e m r a s t e .

e s i i e s t a s a i m p r e p a i e a i r s i i a t . L i e s t a s a r e s e i L I O r e r i a s s t a s t e s a r e e t e e t r e a s r e i t e s i z e t e a p p i a t i s r i s e t . A s e a a a t e s t a s t a t a r e s m a e r t a a s i e p a e t s r e i t e e r a a m t m e m r a s t e .

T e m s t r a t e t e e e t r a p p r a i t r e s p e t t p a i e r e a t e a m i r e m a r i i e a t r e a r e p e a t e a s a t i *bigstack()* i t e s a p a e s a l B e r t e s t a . T r e a s i e e t e e a s t *bigstack()* . O r m p i e r a a s i s i s e r t s a e p i t a t t e s e a s a t e e p i t a s e s a a r e s t a t e i e r t e r a t i t e a . i e *bigstack()* e s t i e a t r e a s s a r e a s i e l B s t a i t t r s t a a a s i s e a e t i e e a t r e a i t s i i i a l B s t a .

e r a t i s m i r e m a r i t 00 t r e a s e a i a s *bigstack()* 10 t i m e s . e e a t r e a a s i t s i i i a s t a t e e m a r t a e s 3.33 s e s 1.07 s e s i a r e a t s e r e e . e s i r s t a a a s i s t e e m a r t a e s 1.0 s e s i t 1.00 s e s a t s e r e e s a r i a s i e s t a a s r a s t i a r e e t e s t p a i . e r i t i s t e s t i t 1000 t r e a s t e e r s i i t t r s t a a a s i s s t a r t s t r a s i i t t e s t a a a s i s t t e r i t i m e s a e s i e a r p t 100000 t r e a s .

### 36 S : A h 2044

e a p p i e t i s a a s i s t t e A p a e 2.0 . e s e r e r . e s e t t e *MaxPat* p a r a m e t e r t 2 K B t i s i e a s m a e e a m i i t e m e r a s i t e s i s t r m e t e r a r i s p a r a m e t e r a e s . T e r e s t s s i i r e 7 i i a t e t a t 2 K B r K B i s a r e a s a e i e s i e a r e r p a r a m e t e r a e s m a e i t t e i e r e e i t e e r a a m t i s t r m e t a t i . e s e t t e *MinC unkt* p a r a m e t e r t K B a s e p r i i r m a t i . B

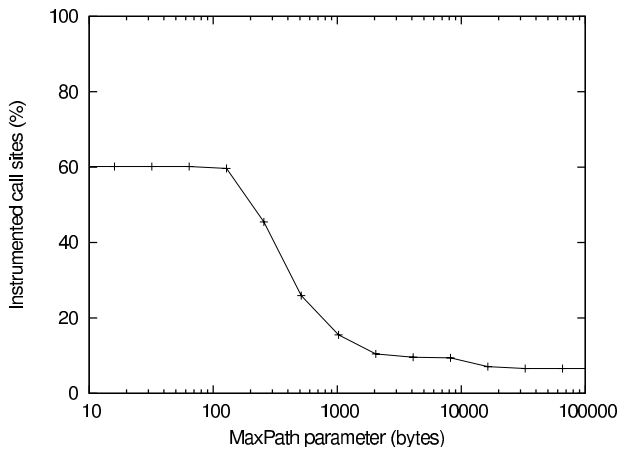


Figure 7: Number of Instrumented Call Sites vs MaxPath Parameter

The number of instrumented call sites decreases as the MaxPath parameter increases. At 10 bytes, approximately 60% of call sites are instrumented. This percentage drops to about 45% at 200 bytes, and continues to decrease, reaching approximately 6% at 20,000 bytes.

Using these parameters, we estimate the number of instrumented call sites for different MaxPath values. At 0.1% (10 bytes), approximately 60% of call sites are instrumented. At 0.5% (500 bytes), approximately 25% are instrumented. At 10% (10,000 bytes), approximately 7% are instrumented. At 50% (50,000 bytes), approximately 6% are instrumented.

The number of instrumented call sites is approximately 60% for MaxPath = 10 bytes, 45% for MaxPath = 200 bytes, 25% for MaxPath = 500 bytes, 15% for MaxPath = 1000 bytes, 10% for MaxPath = 2000 bytes, 8% for MaxPath = 5000 bytes, and 6% for MaxPath = 10000 bytes.

#### 4.2.1.1. THE AWAKE STATE

The awake state is a state where the process is running and has not yet reached the sleep state. The awake state is the state where the process is running and has not yet reached the sleep state. The awake state is the state where the process is running and has not yet reached the sleep state.

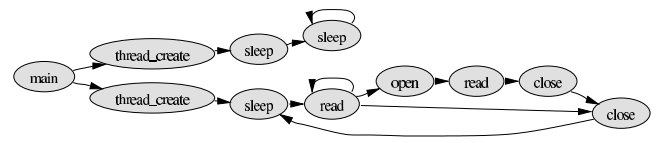


Figure 8: A state transition diagram showing the states of a process. The states are main, thread\_create, sleep, read, open, read, and close. Transitions are shown between these states, including self-loops on sleep and read.

The states of a process are: main, thread\_create, sleep, read, open, read, and close. The transitions between these states are: main to thread\_create, thread\_create to sleep, sleep to sleep (self-loop), sleep to read, read to read (self-loop), read to open, open to read, read to close, and close to close (self-loop).

The states of a process are: main, thread\_create, sleep, read, open, read, and close. The transitions between these states are: main to thread\_create, thread\_create to sleep, sleep to sleep (self-loop), sleep to read, read to read (self-loop), read to open, open to read, read to close, and close to close (self-loop).

The states of a process are: main, thread\_create, sleep, read, open, read, and close. The transitions between these states are: main to thread\_create, thread\_create to sleep, sleep to sleep (self-loop), sleep to read, read to read (self-loop), read to open, open to read, read to close, and close to close (self-loop).

#### 4.2.1.2. THE BLOCK STATE

The block state is a state where the process is blocked and has not yet reached the awake state. The block state is the state where the process is blocked and has not yet reached the awake state. The block state is the state where the process is blocked and has not yet reached the awake state.

The block state is a state where the process is blocked and has not yet reached the awake state. The block state is the state where the process is blocked and has not yet reached the awake state. The block state is the state where the process is blocked and has not yet reached the awake state.

The block state is a state where the process is blocked and has not yet reached the awake state. The block state is the state where the process is blocked and has not yet reached the awake state. The block state is the state where the process is blocked and has not yet reached the awake state.



... e eep a simiar ei te a era e rea e i  
e p ate e er time at rea tra erses e its t i  
e es. a es a era e is esse tia a ei te a era e  
t ee e a es si ete m er p ates is pr pr  
ti a t te m er times ea t i e e is ta e .  
T e e a e t s t e s s t e e t e e i  
ta e *on average*.

... ia ea tatete a es i res re sa e. r  
re t e e e res res as mem r sta spa e a  
s ets a e tra t em i i a . As it U  
time t ere are ei te a era es r t e sa es.  
G i e t a t a e t rea is ate at a parti ar e  
t ese a tati sa s s t estimate et er r i  
t is t rea i i rease r e rease t e t rea s sa e  
ea res r e. T is estimate is t e asis r res rea are  
s e i : e e t at a res r e is s a r e e  
pr m te es (a t s t rea s) t at rease t at res r e  
a em te est at a q ire t at res r e.

## 42 -A S h

... ste isti e e t s stems pri ritize e e t a ers stat  
i a . A ses i r mati s ase e t a er q e e  
e t s t t ami a t e t e s stem. a p r i es  
e step r t er i t r i t e t i res rea are  
s e i . I t is se t i es t set e i  
rap t per m res rea are s e i t at is t  
tra spare ta appi ati spe i .

Or strate r res rea ares e i ast ree parts:

1. Keep tra res re ti zati e e s a e i e  
ami a i ea res r e is at its imit.
2. A tate ea e it t e res res se its  
t i e es s e a pre it t e impa t ea  
res res es e et rea s r m t at e.
3. ami a pri ritize es (a t s t rea s) r  
s e i ase i r mati r m t e r s t t  
parts.

... rea res re e i rease ti zati ti it rea es  
ma im m apa it (s as e t er a a t er  
res r e) a t e e t r t t e a s e i es  
t at rease t at res r e. e res r e sa e is  
e a t t pre ere tia s e e est at s met at  
res r e r t e ass mpti t at i s i i rease  
t r p t. r e imp r t a t e a res r e is er  
e e pre ere tia s e e est at rease t e  
res r e t a i t ras i .

... T is m i ati e se it s me steris t e s  
t e e p t e s stem at t r t t e it t t eris t ras  
i . A iti a res rea ares e i pr i es a at  
ra r a se siti e r m a missi t r si e  
tas s ear m p e t i t e t rease res r es ereas  
e tas s a at e m. T is strate is m p e t e a ap  
tie i t at t e s e er resp s t a es res r e  
s mpti e t t t e t p e r ei e a  
ere a . T e spee a aptati is t r e t e  
parameters t e e p e tia ei te a era es i r  
i rap a tati s.

... Or imp eme tati res rea ares e i is q ite  
stai t r ar . e mai tai separate r q e es rea  
e i t e i rap . e p e r i a e t e r m i e t e  
re ati e pri rities ea e ase r pre i ti  
t e i r s seq e t res r e e e s a t e era res r e

... ti zati t e s stem. O e t e pri rities are  
e se e t a es stri es e i a t e e se e t  
t rea s it i es es q e i r m t e es r  
q e es. B t t ese perati sare  $O(1)$ .

... A e er i ass mpti r res rea ares e  
er is t at res r e sa e is i e t e simiar r ma  
tas s at a i p i t. r t ate t is ass mpti  
seems t i pra ti e. it Apa e re am p e t e r e  
is a m s t a r iati i res r e ti zati a t e e es  
t e i rap .

### 4.2.1 Resources

... T e res res e r r e t tra are U mem r a  
e es ript rs. e tra mem r sa e pr i i r  
ersi t e malloc() ami . e e t e t t res r e  
imit r mem r at i pa e a t a ti it .

... r e es ript rs e tra t e open() a close() a s.  
T is t e i q e a s s t e t a i rease i p e e  
es ript rs i e i e as a res r e. r r e t e  
set t e res r e imit estimati t e m er p e  
e t i s at i resp se time mps p.

... e a as tra irt a mem r sa e m er  
t rea s t e t s at prese t. V is tra e t e  
same a as p si a mem r t t e imit is rea e e  
e rea s me a s t e t res r t ta V a ate  
(e. . 90% t e a res spa e).

### 4.2.2 Pitfalls

... e e tere s me i teresti pit a s e imp eme t  
i a p r i s res rea ares e er. i r s t e t e r m i  
t e ma im m apa it a parti ar res rea e t r i .  
T e ti zati e e at i t ras i r s t e e p e s  
t e r a . r e am p e t e is s s stem a  
s stai ar m re req ests per se i t e req ests are  
seq e tia i stea ra m. A iti a res r e sa  
i t e r a t as e t e V s stem tra es spare is a  
i t t ree p si a mem r. T e m s t e e t i e s t i  
e a e i s t at r ear si s t ras i (s  
as i pa e a t rates) a t set ese si s t i i ate  
ma im m apa it .

... U r tate t ras i is t a a s a eas t i t  
ete t si e i t is ara t e r i z e a e rease i pr ti e  
r a a i rease i s stem er ea . i e e a  
meas re er ea pr ti it is i ere t a appi ati  
spe i ti . At prese t e attempt t ess at t r  
p t si meas res i e t e m er t rea s reate a  
estr e a t e m er es p e e a se . A  
t t is appr a seems s ffi i e t r appi ati s s  
as Apa e m re m p i ate appi ati s m i t e e t r m  
at rea i A I t at a s t e m t e p i it i r m t e  
r time s stem a t t e i r r r e t p r ti it .

... Appi ati spe i res r es a s prese t s me a e es.  
r e am p e appi ati e e mem r ma a eme t i es  
res rea ati a ea ati r m t e r time s s  
tem. A iti a appi ati s ma e e t e r i a re  
s r ess as s. O e a ai pr i i a A I t r  
i t e appi ati a i r m t e r time s stem a t  
its i a res res ma e a reas a es ti . r s i m p e  
ases i e mem r a at r s it ma a s e p s s i e t  
a i e e t is a it t e e p t e m p i e r.

### 43 Y P fi

O e pr em t at arises it perati es e i is t att rea sma t ie t e pr ess r i a ea t air ess re e star ati . T ese pr ems are miti ate t s mee te t e att a t t e t rea s are part t e same appi ati a are t ere rem t a tr sti . N et e ess ai re t ie is sti a per rma e pr em t at matters.

Be a se ea tat e rap ami a it t er i time rea e e it is tria t t see est at ai e t ie : t eir r i times are t pi a r ers ma it e ar er t a t e a era e e . O r imp eme tati a s t e s stem perat r t see t e i rap i i e e time req e ies a res res se se i a USR2 si a t t er i ser er pr ess.

T ist is er a a e e pr ti e a appi ati s t apr i . re am pe i pr ti Apa e e ma pa est at i t ie s ffi ie t te . T is res t is ts rprisi si e Apa e e pe ts t r it preemp ti et rea s. re am pe it r s t t at t e close() a i ses a s et a s metimes ta e 5ms e e t t e me tati i sist s t at it ret r s imme iate e i I/O is see t e . T t is pr em e i s er t a iti a ie si r s stem a i rar e rea a ter t e a t a a t close(). i e t is s t i es t t e pr em i e era it es a s t rea t e e e i t sma er pie es. A etter s ti ( i e a e t et imp eme te ) ist sem t ipe er et rea s r r i ser e e t rea s. T is appra a te se m t ipe pr ess rs a it i e ate ies r m asi a tr a e i perati s s as close() a s r pa e a t a i .

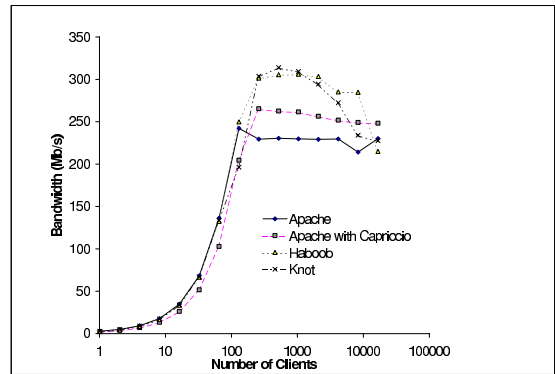
### 5 EVAL AT

T e mi r e mar s prese t e i e ti 2.3s t at apr i as I/O per rma e a e e e t s a a i t . I t is se ti e e a ate apr i s per r ma e m re e era er a rea isti e ser er r a . Rea r e r asi e are m ers p te tia s i e ts i pr i e tests t apr i s s a a i t a s e i . e is s ste er ea apr i s res rea ares e er i t is t e t a t e is ss t is s e er a e e t a t mati a missi tr .

### 5 W S v P f

T e ser er ma i e r r e e mar s is a 500 Hz e ti m ser er it 2GB mem r a a I t e e 1000 Gi a it t er et ar . T e perati s stem is st Li 2. .20. U rt ate e t at t e e e p me t series Li er e se i t e mi r e mar s is sse ear ier e ame sta e e p a e er ea a . He e t is e perime t es t ta e a ta e epoll r Li AIO. imi ar e ere ta et mpare apr i a ai st N TL r t is r a . e ea e t ese a iti a e peri me ts r t re r .

e e erate i e t a it p t 1 simi ar re ma i es a r ss a Gi a it s it e et r . B t apr i a Ha per rm i et r I/O it t e sta r UNIX poll() s stem a a seat rea p r is I/O. Apa e 2.0. ( re t se O IX t rea s) ses a mi ati spi p i i i i a e es ript rs a sta ar i I/O a s.



F u 9: W b v b d w d h v u h u b of ul ou cl .

T e r a r t is test siste req ests r 3.2 GB stati e ata it ari s e sizes. T e req est req e ies rea size a rea e ere esi e t mat t se t e e 99 e mar . T e i e ts r t is test repeate e t t e ser er a iss e a series e req ests separate 20ms pa ses.

e imite t e a e sizes Ha a K t t 200 B i r er t rea ea is a ti it . e se a mi ima rati r Apa e isa i a ami m es a a ess permissi e i . He e it per rme esse tia t e same tas s as Ha a K t . T e per rma eres ts ere q itee ra i . Apa es per rma e impr e ear 15% er er apr i i . A iti a K ts per rma e mat e t at t e e e t ase Ha e ser er.

arti ar remar a e is K ts simp i t . K t sist s 1290 i es e ritte i a strai t r ar t rea e st e. K t as er eas t rite (it t e s 3 a s t reate) a it is eas t ersta . e si e r t is e perie et e str e i e e r t e sim pi it t e t rea e appra t riti i rre t appi ati s.

### 52 B k h S

ai tai i i rmati a t t e res res se atea i p i t req ires t e t ermi i ere t e pr ram is e it sa per rmi s me am t mp ta ti t sa e a a re ate res r e ti izati res.

T a e 2 q a ti est is er ea r Apa e a K t r t e r a es rie a e. T e t p t i es s t e a era e m er appi ati est at ea appi ati spe t i r m e i p i t t t e e t. T e t t m t i es s t e m er est at apr i spe s i t e r a i r er t mai tai i rmati se t e res rea ares e er. A e ts are t e a era e m er es per i rap e e ri rma pr essi (i.e. er a a a t e r t e mem r a ea ra pre i t rs a e arme p).

It is imp rta t t t e t at t ese e ts i e only t e timespe t i t e appi ati it se . Ker e time spe t

	Item	Cycles	Enabled
Apps	Apache	32697	n/a
	Knot	6868	n/a
System	stack trace edge statistics	2447 673	Always for dynamic BG During sampling periods

**Table 2: Average - d c cl cou fo l c o o C cc o.**

I/O processes are I/O intensive. The overhead for each is 0.1%.

The overhead for each is 0.1%.

The overhead for each is 0.1%.

The overhead for each is 0.1%.

### 53 -A A

The overhead for each is 0.1%.

The overhead for each is 0.1%.

The overhead for each is 0.1%.

## 6 ELATE W

### Pool of Hosts

The overhead for each is 0.1%.

The overhead for each is 0.1%.

The overhead for each is 0.1%.

The overhead for each is 0.1%.

The overhead for each is 0.1%.

### U-Lvl Th d

The overhead for each is 0.1%.

The overhead for each is 0.1%.

The overhead for each is 0.1%.

i per rma et rea pa a espe iaize r e a es  
t ati es ast is i raries a mem r ma a eme t.  
T e per rma e ptimizati semp e t ese pa a es  
e se r apr i as e t ese are mp eme  
tar t r r .

T e tate T rea spa a e 3 is a i t ei t pera  
ti et rea i s stem t ats ares apr i s a simp i  
i t e pr rammi m e r et r ser ers. U i e  
apri i t e tate T rea s i rar es t pr i e a  
O IX t rea i i ter a e s appiati s m st e rit  
te t se it. A iti a tate T rea s se it er **select**  
**r poll** i stea t em res a a e Li **epoll** a t e  
se i is I/O. T ese a t rs imit t e s a a i t  
tate T rea s r et r i t e si e r a s a t e  
restr i t its rre r is i t e si e r a s. T ere  
are pat es a a i a e t a Apa e t se tate  
T rea s 33 res ti i a per rma e i rease. T ese  
pat es i rea m er t er impr eme t st Apa e  
e er s it is imp ssi et te m t e im  
pr eme t ame r m tate T rea s. U rt ate t ese  
pat es are er mai tai e a t mpi e ea  
s e ere a e t r i re t mparis s a ai st  
apri i .

e era ti ati s 2 s et e pr em i I/O  
a e pe te i /preempti ser e e t rea s  
a i er e s pp r t r ti i t e ser e e s e  
er t ese e e ts. T is apr a e s res ea i t e rati  
t et rea i rar a t e perati s stem e er  
t e ar e am t er e a es i e se em t a e  
pre e i e a pti . A t er p t e tia pr em it  
t is apr a ist att ere i e e s e era ti ati  
r ea t sta i I/O perati i a m er i  
t et e s t sa s r I t er et ser ers. T is res t is  
tr ar t t e r i a a re i t e m er  
er e t rea s e e . T is pr em appare t stems r m  
t e a t t at s e era ti ati s are e e pe primari  
r i per rma e mp ti e ir me ts ere is  
a ast et r I/O are mi a t. Ne er t e e s s e er  
a ti ati s a e a ia e apr a t ea i it pa e  
a ts a preempti s i apr i i. mp i s e er  
a ti ati s a s a t e ser e e s e r t i  
e e t e r e s e i s i a t i er e t rea t  
preempt. T is s eme a e se t s e i ffi t pr em s  
i e *priority inversion* a t e *convoy eno enon* .  
pp r t r ser e e preempti a :Nt rea i (i.e.  
r i ser e e t rea s t p N er e t rea s) is  
tri . T e i q es s as ptimisti rre tr  
a a i s r stea i 7 a e se e e t i e t ma  
a e t rea a s e er ata str t res. r i a prese ts  
a i e es ripti t ese a t er te i q es i t e  
te t Li 12. e e pe t t emp ma t ese te  
i q es i apr i i e ea s pp r t r :Nt rea i .

## K l Th d

T e N TL pr e t r Li as ma e reat stri es t  
ar impr i t e e ffi e Li er e t rea s. T ese  
a a esi ea m er er e e e impr eme t s s  
as etter ata str t res er mem r er ea a t e  
se  $O(1)$  t rea ma a eme t perati s. N TL is q ite  
e a is sti er a ti e e e pme t. He e e pe t  
t at s me t e per rma e era ati e  
i er m ers t rea s ma eres e ast e e e pers  
s a reate aster a rit ms.

## A l c o - c f i c z o

er rma e ptimizati t r appiati spe i  
tr s stem res res is a imp rta t t eme i O re  
sear . a 21 a e appiati st spe i t e ir  
V pa i s eme i impr e per rma e r ap  
pi ati s t at e a t t e ir p mi mem r ee s  
a is a ess patter s. UN T 37 i simi ar ti s  
r et r I/O impr i e i i t a re i er  
ea it t mpr misi sa et . T e IN perati  
s stem 5 a t e VINO perati s stem 29 pr i e ser  
st mizati a i appiati et em e i t  
t e er e. T e er e 13 t t e pp site apr a  
a m e m st t e O t ser e e . A t ese s stems  
a appiati spe i ptimizati ear a aspe ts  
t e s stem.

T ese te i q es req ire pr rammers t tai r t e ir ap  
pi ati t ma a er es res rit se t ist pe t i  
is t e i ffi t a ritte. A iti a t e tie pr rams  
t sta ar A Is re i t e ir p rta i it. apr i  
ta es a e apr a t appiati spe i ptimizati  
e a i a t mati mpi er ire te a ee a ase  
t i t et rea pa a e. e e i e t att is apr a  
i ma e t ese te i q es m re pra ti a a i a a  
i er ra e appiati st e e t r m t em.

## A ch o ou I/

A m er a t rs pr p se impr e er e i ter a es  
t at t a ea imp rta t impa t ser e e t rea i .  
As r s I/O primiti es s as Li s **epoll** 20  
is AIO 17 a reeB s q e e i ter a e 19 are e  
tra t reati s a a e ser e e t rea pa a e. apr i  
i ta es a a ta e t ese i ter a es a e e t  
r m impr eme ts s as re i t e m er er e  
r ssi s.

## ck M

T ere are a m er reate apr a est t e pr em  
prea ati ar e sta s. me ti a a a es  
s as ta ar L Ne Jerse 3 t se a a  
sta at a rat er t e a ate a a ti ati re r s  
t e eap. T is apr a is reas a e i t e t e t  
a a a e t at ses a ar a e e tra t at s pp rts  
i er r er ti sa rst ass ti ati s . H  
e er t ese eat res are t pr i e t e pr rammi  
a a e i mea s t at ma t ear me ts i a r  
eap a ate a ti ati re r s t app i r ase.  
rt erm re e t is t i r t e er ea ass i  
ate it a i a ar a e e t r t r s stem pre i s  
r ass t at Ja a s e era p r p se ar a e e  
t r i s i apr priate r i per rma es stems 30.

A m er t er s stems a e se ist s sma sta  
s i pa e ti s sta s. B r a e reit  
es ri e ate i q et at ses a si esta r m tipe e  
ir me ts e e t i e i i i t e sta i t s sta s  
e er t e t a a ze t e pr ram t attempt t  
re e t e am t r time e s req ire . O e  
i is a a a e a r time s stem r para e izi  
pr rams se a simp i e ersi B r a e  
reit s te i q e a e spa etti sta s 9. I t is te  
i q e a ti ati re r s r i ere t t rea sare i ter ea e  
a si esta e er ea a ti ati re r s i t e  
mi e t e sta a t e re a ime i i e a ti ati  
re r s sti e ist rt er t e sta i a a  
t e am t aste sta spa e t r it t .

re re e t t e Laz T rea s pr e t i tr e sta  
 ets i are i e sta s r se i mpi i  
 para e a a es 15. T is me a ism pr i es r time  
 sta er e s a it ses a mpi e r a a sis t  
 e imi ate e s e sta sa e a e e e er  
 t is a a sis t at es t a ere rsi as apri i  
 es a it es t pr i et i parameters. e a  
 Be as se e size sta ets t pr i e s  
 pr essi time i a para e rea time ar a e e t r 11.

## 7 F T EW

e are i t e pr ess e te i apri i t r it  
 m ti U ma i es. T e ame ta a e e pr i e  
 m tipe Us ist at e a er re t e p  
 erati et rea i m e t pr i e at mi it . H e er e  
 e i e et at i rmati pr e t e mpi e r a assist  
 t e s e er i ma i e isi s t at ara tee at mi it  
 ertai s e at t e app i ati e e .

T ere are a m er aspe ts apri i s imp eme ta  
 ti e e i e t e pr e . e e i e e e ramati  
 a re e er e rssi s er ea et r a it  
 a at i i ter a e ras r s et r I/O . e a s  
 e pe t t ere are ma a s t impr e r res r e a are  
 s e er s a s tra i t e a r i a ei t er e s r e sa e  
 i rap es a impr i r ete ti  
 t ras i .

T ere are se era a s i i r sta a a sis a  
 e impr e . As me ti e ear ier e se a ser ati e  
 appr imati t e a rap i t e prese e ti  
 p i ters r t er a a e eat res t at req ire i ire t  
 a s (e . . i er r er ti s irt a met ispat  
 a e epti s). Impr eme ts t t is appr imati  
 s sta tia impr e r res ts. I parti ar e pa  
 t a apt t e ata a a sis re 2 i r er t  
 isam i ate ma t e ti p i ter a sites. e  
 mpi i t er a a es e start it simi ar  
 ser ati e a rap sa t e emp e isti tr  
 a a ses (e . . t e 0 A a a ses 31 r ti a  
 a e t r i e t e a a es a a es r irt a ti  
 res ti a a ses 27 r e t r i e t e a a es).

I a iti e pa t pr e pr i t s t at a  
 assist t e pr rammer a t e mpi e r i t i apri  
 i s sta parameters t t e app i ati s ee s. I parti  
 ar e a re r i rmati a t i ter a a e t er a  
 aste spa e a e a at er statisti s a t i  
 ti a s a se e sta s t e i e . B  
 ser i t i s i rmati r a r a e parameter a es  
 e a a t mate parameter t i . e a a s s est  
 p te tia ptimizati s t t e pr rammer i i ati  
 i ti s are m st t e resp si e r i reasi  
 sta size a sta aste.

I e era e e i e et at mpi e r t e i pa  
 a imp rta t r e i t e e ti t e t e iq es e  
 s ri e i t is paper. r e ampe e are i t e pr ess  
 e isi a mpi e r a a sis t at is apa e e erati  
 a i rap at mpi e time t ese res ts i impr e  
 t e effi e t er times stem (si e a tra es are  
 req ire t e erate t e rap ) a t e i a s t  
 et at mi it r ree ara teei stati a t at ertai  
 riti a se ti s t tai i p i ts. I a iti  
 e pa t i esti ate strate ies r i serti i p i ts  
 i t t e e at mpi e time i r er t e r e air ess.

mpi e time a a sis a a s re e t e rre e  
 s ar i t e pr rammer a t ata ra es. A  
 t stati ete ti ra e iti s is a e i  
 t ere as ee re e t pr res s et mpi e r impr eme ts  
 a tra ta e e pr ram a a ses. I es 1 a a  
 a e r et r e se s r s t ere iss pp rt r at mi se  
 ti s a t e mpi e r er sta s t e rre m e .  
 It ses a mi t re I/O mp eti s a r t mp eti  
 t rea s a t e mpi e r ses a ariati a a rap  
 t at is simi ar t r i rap . T e mpi e r e s res  
 t at at mi se ti s resi e it i e e e t at rap  
 i parti ar a s it i a at mi se ti a t i e  
 r (e e i ire t ). T is i s pp rt e  
 e treme p er r a t r i ser ers. i a e e pe t  
 t at at mi se ti s i a s e a e etter s e i a  
 e e ea ete ti .

## 8 C CL S S

T e apri i t rea pa a e pr i es empiri a e i e e  
 t at i t rea pa a es is a ia es ti t t e pr  
 em i i s a a e i rre I ter et ser ers.  
 O re perie e it riti s pr ram s s est t at t e  
 t rea e pr rammi m e is a m re se a s tra ti  
 t a t e e t ase m e r riti mai tai i a  
 e i t ese ser ers. B e pi t e t rea impe  
 me tati r m t e perati s stem it se e a ta e  
 a a ta e e I/O me a isms a mpi e r s pp rt.  
 As a res t e a se t e iq es s as i e sta s  
 a res r e a res e i i a s ta i e e  
 si i a t s a a i it a per rma e impr eme ts e  
 mp are t e isti t rea ase r e e t ase s stems.

As t is t e mat res e e pe t e e m re t ese  
 t e iq es t e i t e rate it mpi e r t e . B  
 riti pr ram s i t rea e s t e pr rammers pr i e  
 t e mpi e r it m re i rmati a t t e i e e  
 str t r e t e tas s t at t e ser er m st per rm. Usi  
 t is i rmati e pe t at t e mpi e r a e p se e e  
 m re pp rt ities r t stati a ami per rma e  
 t i .

## 9 EFE E CES

- 1 A. A a J. H e . T eimer . J. B s a  
 J. R. e r. perati e tas ma a eme t it t  
 ma a sta ma a eme t. I *Proceedings of t e 2002  
 Usenix ATC J e 2002.*
- 2 T. . A ers B. N. Bers a . . Laz s a a  
 H. . Le . e er A ti ati s: e ti e Ker e  
 pp rt r t e User Le e a a eme t  
 ara e ism. *ACM Transactions on Co uter Syste s*  
 10(1):53-79 e r ar 1992.
- 3 A. . Appe a . B. a Q ee . ta ar L  
 Ne Jerse . I *Proceedings of t e 3rd International  
 Sy osiu on Progra ing Language  
 I le entation and Logic Progra ing* pa es 1-13  
 1991.  
 A. . Appe a Z. a . A empiri a a a a ti  
 st sta s. eap st r a a es it  
 s res. *Journal of Functional Progra ing*  
 (1): 7-7 Ja 199 .
- 5 B. N. Bers a . am ers . J. ers . ae a  
 . Namee . ar a . a a e a . G. irer.  
 IN a e t e si e mi r e r

- appiati spesi perati sistem ser ies. I *ACM SIGOPS Euroean Works o* pa es -71 199 .
- . Bas e J. Gra . . it ma a T. G. ri e. T e p e me . *O erating Syste s Revie* 13(2):20-25 1979.
- 7 R. . B m e . . J er B. . K szma . . . Leisers K. H. Ra a a Y. Z . i : A effie t m tit rea e r times stem. *Journal of Parallel and Distri uted Co uting* 37(1):55- 9 199 .
- . G. B r a B. e reit. A m e a sta imp eme tati m tipe e ir me ts. *Co unications of t e ACM* 1 (10):591- 03 O t 1973.
- 9 . . ar is e A. R ers J. Repp a L. He re . ar e perie es it O e. I *Proceedings of t e 6t International Works o on Languages and Co ilers for Parallel Co uting (LNCS)* 1993.
- 10 A. a t . B. a zi . Neer aes . . artz a K. J. rre . A Hierar ia I ter et O e t a e. I *Proceedings of t e 1996 Usenix Annual Tec nical Conference* Ja ar 199 .
- 11 . e a G. . Be . A para e rea time ar a e e t r. I *Proceedings of t e 2001 ACM SIGPLAN Conference on Progra ing Language Design and I le entation (PLDI '01)* 2001.
- 12 J. r i a. ast m tit rea i sare mem r m tipr ess rs. Te ia rep rt U i ersit ata J e 2000.
- 13 . R. er . . Kaas e a J. O T e. er e: A perati s stem ar ite t r e appiati e e res r e ma a eme t. I *Sy osiu on O erating Syste s Princi les* pa es 251-2 1995.
- 1 . Ga . Le is R. Be re . es . Bre er a . er. T e es a a e: A isti appra t et r e em e e s stems. I *ACM SIGPLAN Conference on Progra ing Language Design and I le entation* 2003.
- 15 . . G stei K. . a ser a . . er. Laz T rea s ta ets a r izers: a i primiti es r mpi i para e a a es. I *T ird Works o on Langauges, Co ilers, and Run-Ti e Syste s for Scala le Co uters* 1995.
- 1 T. H . i ima te t T rea 0.7 ma a . ttp:// . ara et r . m/ s/ m t ma a . p 2002.
- 17 B. LaHaise. Li AIO me pa e. ttp:// . a . r / a / ai / .
- 1 H. . La er a R. . Nee am. O t e ait perati s stem str t res. I *Second Inernational Sy osiu on O erating Syste s, IR1A O t er* 197 .
- 19 J. Lem . Kq e e: A e eri a s a a e e e t ti ati a iit. I *USENIX Tec nical conference* 2001.
- 20 . Li e zi. Li ep pat . ttp:// . mai ser er. r / i pat es/ i impr e. tm .
- 21 . Namee a K. Armstr . te i t e a e ter a pa eri ter a e t a mm ate ser e e pa e rep a eme t p i ies. Te ia Rep rt TR 90 09 05 U i ersit as i t 1990.
- 22 . . i . *Advanced Co iler Design and I le entation*. r a Ka ma a ra is 2000.
- 23 G. . Ne a . ea . . Ra a . eimer. II: I terme iate a a e a t s r a a sis a tra s rmati pr rams. *Lecture Notes in Co uter Science* 230 :213-229 2002.
- 2 G. . Ne a . ea a . eimer. re : T pe sa e retr tti e a e. I *T e 29t Annual ACM Sy osiu on Princi les of Progra ing Languages* pa es 12 -139. A Ja . 2002.
- 25 J. K. O ster t. T rea s Are A Ba I ea ( r m st p rp ses). rese tati i e at t e 199 Use i A a Te ia ere e Ja ar 199 .
- 2 V. . ai . rs e a . Z ae ep e. as : A ffi e t a rta e e er er. I *Proceedings of t e 1999 Annual Usenix Tec nical Conference* J e 1999.
- 27 H. . a e a B. G. R er. ata ase irt a ti res ti . *Lecture Notes in Co uter Science* 11 5:23 -25 199 .
- 2 . a a . . G i. A a rit m rs i strai t satis a ti pr ems.
- 29 . I. etzer Y. . ma a K. A. mit . eai it isaster: r i i mis e a e ere e e te si s. I *Proceedings of t e 2nd Sy osiu on O erating Syste s Design and I le entation* pa es 213-227 eatt e as i t 199 .
- 30 . A. a . a e . J. ra i a J. . He erstei . Ja a s pp rt r ata i te si e s stems: perie es i i t e Tee rap ata s stem. *SIGMOD Record* 30( ):103-11 2001.
- 31 O. i ers. *Control-Flo Analysis of Hig er-Order Languages*. t esis ar e ie e U i ersit a 1991.
- 32 . T er i . r ti e i rar s re. ttp:// . r . e / r ese/ r / .
- 33 U . A e erati Apa e pr e t. ttp:// aap.s r e r e. et/ .
- 3 U . tat e t rea s r I ter et appiati s. ttp:// state t rea s s r e r e. et/ s/ st. tm .
- 35 J. R. Be re . Bre er N. B ris . e . es J. a a J. La . Gri e a . er. Ni a: A rame r r et r ser ies. I *Proceedings of t e 2002 Usenix Annual Tec nical Conference* J e 2002.
- 3 R. Be re J. it a . Bre er. e e ts are a a i ea ( r i rre ser ers). I *Proceedings of t e 2003 HotOS Works o a* 2003.
- 37 T. i e A. Bas V. B a . V es. U Net: A User Le e Net r I ter a e r ara e a istri te mp ti . I *Proceedings of t e 15t ACM Sy osiu on O erating Syste s Princi les* pper tai Res rt O U A e eme er 1995.
- 3 . es . . er a . A. Bre er. A: A ar ite t r e r e iti e s a a e I ter et ser i es. I *Sy osiu on O erating Syste s Princi les* pa es 230-2 3 2001.