

```
1 -----
2 Morgan GangwereGSOC 2016
3
4             *** FINAL DRAFT ***
5
6             A proposal for the NetBSD project
7             Inetd improvements
8
9 -----
10
11
12 *** DRAFT HISTORY ***
13
14 2016-03-21: Initial proposal presented to Google Summer of Code 2016
15 2016-03-21: Fix some wording, make the whole document fit in 80cols
16 2016-03-21: Add better contact information
17 2016-03-22: Refine information about goals (prefork) and bonuses.
18 2016-03-24: Clarify certain points
19 2015-03-25: Final edits + argument for proposed path.
20
21 *** PROVIDING FEEDBACK / CONTACTING THE AUTHOR **
22
23 Feedback can be provided directly to me via email < morgan.gangwere@gmail.com >
24 as well as via the NetBSD tech-userlevel mailing list (where this will be
25 posted for public review)
26
27 I may be reached over XMPP via indrora@rows.io or via email (see above). I can
28 also be found on Freenode as indrora (responses may be delayed due to ZNC/idle)
29
30 A copy of my resume is available at http://tsunami.zaibatsutel.net/cv.pdf
31
32 This document is available in plain text at
33 http://tsunami.zaibatsutel.net/gsoc16\_netbsd\_proposal.txt
34
35 ---
36
37 1. Introduction (About the author)
38
39 I'm Morgan Gangwere, a student from the University of New Mexico in Albuquerque,
40 New Mexico. I've worked on lots of things, but some of my favorites are open
41 source contributions. I've made contributions in my past to the libmtp project,
42 fixed bugs in Travis (in the end, the patch wasn't accepted -- my solution was
43 against the documented details, but it sparked the conversation), worked around
44 Android's limitations in OpenKeychain, then helped people communicate when I
45 forked yaaic as Atomic and started making it better.
46
47 I've built my own kernel, booted UEFI by hand, forged raw TCP sockets from raw
48 hot bits with a hex editor and a sledge, helped port Android to a new phone,
49 pushed new binaries to devices over TFTP. I've worked out what a "bad cast to
50 std::Allocator<_T, alloc<T>>(std::stream_t<T>)" means. I've created ext2 fs
51 superblocks by hand, beaten U-Boot into compliance, slimmed down FAT filesystems
52 and run rsync over amateur radio.
53
54 My editor is vim, my shell is zsh, my work? it makes the internet happen. My
55 goal is to make the internet a good place for data to live.
56
57 I got a degree in network administration (focusing on Cisco networks) back when
58 I was in high school. I automated my work in my sophomore year of classes, often
59 spending the second half of the class playing quake. I learned how to build
60 networks that are reliable and consistent in those classes. Since then, my work
61 has focused on making things reliable and easy to manage.
62
63 I'm familiar with NetBSD; I typically run it in virtual machines but also use it
64 on the odd 'I need a lightweight not-linux' system. I'm not afraid of digging
65 through manpages. I'm definitely not afraid of digging into /usr/src to find the
66 function I'm looking for.
67
68 I believe in Systems Engineering -- the idea that all things should be done with
69 a requirements document, a plan and a good debate as to whether this is the
70 right direction to go. I believe in change requests, in code that is built well
```

71 and comments that explain what is being done. I believe in code quality, secure
72 design, portability and expected execution. I believe in clever design over
73 clever implementation, a value I know the NetBSD project holds dearly.
74
75 2. The problems with inetd
76
77 Inetd is plagued with a few problems that the NetBSD projects wants to address
78 through Google Summer of Code, as well as one that I personally wish to tackle
79 early on to tackle a maintainability issue that plagues the NetBSD project.
80
81 There's a few bonuses here and there that I'd like to implement. Implementing
82 a new per-service configuration format makes adding features to inetd simple
83 and paves the way for other tools to do the same.
84
85 2.1 No early service availability
86
87 In conversation with Riastradh on IRC, this is sort of a hot-button topic. The
88 GSoC project page specifies the following:
89
90 > Prefork: Support pre-forking multiple children and keeping them alive
91 > for multiple invocations.
92
93 How this should be implemented or otherwise handled is left very much as an
94 "exercise to the reader".
95
96 As it stands, inetd waits until a service is requested (through kqueue) to
97 fork() and exec(). It's simple, straightforward and is a direct approach to
98 the problem.
99
100 FreeBSD's inetd has a "max children" but has nothing to create children early,
101 only a mechanism to limit the number of sessions at any one given time.
102 [FREEBSD-INETD] However, no
103
104 In 1999, we faced the c10k problem; cdrom.com saw 10,000 concurrent connections
105 to its FTP server. As more devices hit the net, we saw the c100k problem, then
106 in 2014 or so, we began seeing the c10m problem: 10 million active connections
107 in software like Facebook, Twitter, etc. [C10K] [C10M]
108
109 kqueue is a competent system to deliver messages among processes and inetd can
110 use this in a manner similar to how nginx's event-driven worker system. This
111 path works well for Nginx and for fairly direct, computationally inexpensive
112 services such as those which inetd is built for.
113
114 Several years ago (2014), a set of informal studies were done comparing nginx
115 and Apache against each other with various configurations. One such found that
116 nginx's pure event system was better for the static content that was being
117 served in a specific case. [APACHE-PREFORK-NGINX]
118
119 2.2 A configuration format from ages ago
120
121 inetd's configuration format is right out of 1980. Whitespace-delimited, it's a
122 definite relic of how we used to do software. And it was good back when people
123 didn't want to put their entire system's configuration into version control. But
124 now we do want our system to be held inside a version control system and oh man,
125 how are we going to do that?
126
127 The biggest question becomes "what happens when someone upgrades from NetBSD-N
128 to NetBSD-N+1?"
129
130 FreeBSD just crammed more options on top of the wait/nowait column. This isn't
131 a good design decision as it really encourages shoehorning features on top of
132 a format. [FREEBSD-INETD]
133
134 ... Choices need to be made ...
135
136 There's a *lot* of things inside the inetd configuration file that make it a
137 hard to parse format that needs to be migrated away from in the long run. A
138 format for setting bindhosts, specifying ipsec rules -- all these have been
139 shoehorned onto the inetd configuration through special keywords that change the
140 state of the parser as it's going on.

```
141
142 2.3 A solid log of a codebase
143
144 The codebase is pretty much one single file. Modern compilers and modern build
145 systems encourage smaller compilation units for massively parallel compiles.
146
147 I'd personally like to start a trend where NetBSD moves from where the tools
148 are maintained as only a few large files and broken into multiple compilation
149 units.
150
151 This has one upshot: Multiple authors can work on the same tool and not have
152 merge conflicts
153
154 3. Changes that should be made
155
156 3.1 preforking: A change to the fork()/execve() model:
157
158 We must take note of our fellow web daemons, study their actions. Their attack
159 is that of preforking. Implementing preforking in inetd requires a certain
160 restructuring of the codebase.
161
162 This restructuring means that we need to reconsider how inetd handles its child
163 pids -- notably, it currently just keeps a linked list of them and cleans up
164 some things after the children are done.
165
166 The proposed mechanism for preforking is restructuring the `servtab` such that
167 the number of fields is reduced. This would create (roughly) the following
168 C-style structures:
169
170 struct service {
171     char    *name;           /* service name */
172     uint8_t  type;          /* type of service */
173     int      proto;         /* protocol */
174
175     char    *progpath;      /* program path (not name) */
176     char    *argv;          /* arguments to pass to the program */
177     int      port;          /* tcp/udp port to run on (allows for override
178                             from /etc/services?) */
179
180     struct {
181         int  children;      /* number of children to keep at hand at any time */
182         int  overload;     /* Allow overflow (fallback to fork()/exec() on
183                             demand after pool is exhausted) */
184     } child_opts;
185     struct {
186         int  hits;          /* Number of times a service can be used ... */
187         int  time;          /* ... over a certain amount of time ... */
188         int  cooldown;     /* ... before we wait for some time
189                             and let it cool */
190     } load_opts;
191
192     /* ... */
193
194 }
195
196 struct child_service {
197     char    *name;          /* tied to service name */
198     int      fd;            /* file descriptor we have to read/write (stdin
199                             and stdout bound to this) */
200     /* ... other control stuff */
201 }
202
203 (this isn't final by any means)
204
205 3.2 Per-Service configuration
206
207 This is by far the hardest part of what needs to happen to inetd. Inetd's
208 configuration hasn't been made with backwards (or really forwards) compatibility
209 in mind, and as a result we're having to collect on the technical debt.
210
```

211 3.2.1: The options on the table:
212
213 3.2.1.1: Fork inetd, creating inetd-legacy (default) and inetd-ng (not default).
214
215 In this situation, we can ignore the past and rebuild as necessary, totally
216 revamping the inner workings of inetd. However, this comes at a bit of a cost
217 in that we've now got **two** versions of inetd to maintain! (I seriously doubt
218 that the netbsd team wants to maintain two versions of inetd.)
219
220 Pros:
221
222 - Erasure of the past's technical debt with regards to the configuration
223 - keeps upgrades simple: Not ready? Keep the old behavior as it is, get no
224 new features
225
226 Cons:
227 - Two versions of inetd to maintain (with a consistently diverging codebase
228 - -ng has the very real chance of just rotting and being left in the rain.
229 - Your daily dose of technical debt at twice the cost: Fixing bugs in the old
230 (-legacy) variant **plus** any bugs in the new (-ng) variant.
231
232 3.2.1.2: Want new features? Use the new format. Otherwise, fine whatever.
233
234 In this situation, we keep all the old parsing of /etc/inetd.conf as it stood.
235 Any usage of new features is dependant on the new configuration format.
236
237 Pros:
238
239 - Existing installations can upgrade and not have to worry about it.
240 - Adding new features simple (new features are only in the old format)
241 - Zero upgrade cost
242
243 Cons:
244
245 - Support **nightmare** - you have to ask "are you using the old or the new
246 format? What's your setup?" to see if there's a service collision. Ewww no.
247 - All the technical debt of the old version is still there, plus new, added
248 complexity from the two formats.
249 - Encourages old users to keep doing what they're doing, "but it's always
250 been done like this"
251 - Someone, somewhere is going to abuse some aspect of this combination to
252 do something nasty.
253 - No incentive to move to the new format
254 - Documentation nightmare
255
256 christos@ suggests this as the best route.
257
258 3.2.1.3 Drop the old format, introduce the new format, keep a tool to spit
259 out configuration files in the new format.
260
261 Here, we totally drop the old table-driven style configuration and move to a
262 format that's tolerable for everyone.
263
264 As part of this, include a tool (using awk, python, whatever) that approximately
265 creates the new format to let people migrate a service.
266
267 Pros:
268
269 - Erasure of technical debt (inetd configuration can be gutted and rebuilt)
270 - Super-duper simple addition of configuration options
271 - Backwards compatible through the future (bring a newer service file in and
272 all known configuration knobs will be brought in with it)
273 - a new, consistent format means easier documentation.
274
275 Cons:
276
277 - Documentation needs to be rewritten
278 - no backwards compatibility from the past **except** through a tool or manual
279 configuration
280 - Need to maintain the tool to convert inetd lines into config lines.

```
281     - Future needs to be told that the past changed the configuration format:
282     Documentation needs to be rewritten but include a mechanism to find the
283     old format's documentation.
284
285     This is my preferred route.
286
287 3.2.1.4 shoehorn new options on top of the old ones as FreeBSD did
288
289 This is what several of our siblings have done. FreeBSD in particular has done
290 several things to try and cram more into the fields of the original inetd config
291 format.
292
293 FreeBSD has made the following the columns for /etc/inetd.conf:
294
295 service-name
296 socket-type
297 protocol
298 {wait|nowait}[/max-child[/max-connections-per-ip-per-minute[/max-child-per-ip]]]
299 user[:group][[/login-class]]
300 server-program
301 server-program-arguments
302
303 This makes sense in FreeBSD's case. Not so much in the case of NetBSD.
304
305 Pros:
306
307     - Familiar format
308     - Does not take a lot of time to implement similar extensions
309
310 Cons:
311
312     - Shoehorning is only sustainable for so long
313     - *not* future proof nor backwards compatible
314     - Doesn't allow for a lot of future changes (e.g. chroot, etc)
315
316 3.2.2: Suggested new format
317
318 Up for suggestion is the following:
319
320 * Key-value pairs in the style of sysctl
321 * A set of directories within /etc, /usr possibly in /usr/pkg/etc and others
322 * A master configuration in /etc/rc.inetd
323
324 Grammar for /etc/rc.inetd would look like this:
325
326   ``
327   # comment
328   IncludeDir /usr/inetd/
329   IncludeDir /usr/pkg/inetd/
330
331   # to specify ssh over ipv4 but not ipv6
332   ssh.tcp4=yes
333   ssh.tcp6=no
334
335   # dns over defaults (ip4/ip6)
336   dns.udp=yes
337
338
339   ### alternately
340
341   dns=tcp4:dgram,bind=xx.xx.xx.xx
342   dns=tcp6:dgram,bind=dead:b33f:caf3:f00d::101
343
344   ``
345
346 within /usr/inetd/ or /usr/pkg/inetd would be a file named "http" containing
347 something like
348
349   ``
350   executable=/usr/pkg/bin/muhttpd
```

```
351 progname=muhttpd
352 arguments=--config muhttpd.conf
353 chroot=/srv/http/
354 ```
355
356 If a configuration file exists in multiple directories, the *last scanned* wins.
357 That is, if the include order is
358
359 * /usr/inetd/
360 * /usr/pkg/inetd/
361 * /usr/local/inetd/
362
363 and the files /usr/pkg/inetd/ssh and /usr/local/inetd/ssh exist, whatever is in
364 /usr/local/inetd/ssh overrides any settings from /usr/pkg/inetd/ssh
365
366 This means that a service on one system can override specific details (such as
367 environment variables or command line options) for a service whose binaries are
368 hosted over NFS.
369
370 Valid options within a service definition file would be, at the start
371
372 name          type          desc
373 type          enum          Socket type, one of (stream|dgram|raw|rdm|seqpacket)
374 wait         bool          Should inetd wait on this process to finish before
375               accepting another connection?
376 user         string         User to run as. Default gid is default gid of user,
377               unless specified with :group
378 server       string         full path to executable
379 progname     string         argv[0]
380 arguments    string         argv[1] and beyond
381 env         string         environment declaration
382 chroot       string         path to chroot() to before calling the executable
383 workers      int           Number of prefork threads to keep at any one given time
384 overflow     bool          allow the prefork pool to overflow.
385
386 *** NOT FINAL ***
387
388 if progname is empty or unspecified, the filename specified in server is used
389 in place.
390
391 Environment variables stack and overwrite each other. If /usr/pkg/inetd/ssh
392 has "env=foo=bar;baz=whatever" and later on ../ssh says "env=foo=wonka",
393 the final environment for the executable will be "foo=wonka;baz=whatever"
394
395 Other
396
397 3.3 a codebase in need of a breakup
398
399 This is actually the easiest part of the task.
400
401 At the moment, inetd.c is a huge 2ksloc file. Making each part of inetd's
402 component parts a separate compilation unit makes compilation easier on big
403 systems (modern systems can handle 2-3 compiler threads at a time) and makes
404 maintenance easier in future (less chance to trample something accidentally.)
405
406 Breaking up inetd.c into its constituent parts (builtins, configuration, inetd
407 itself, some kqueue stuff, structures, etc) improves code readability and
408 maintenance in the future.
409
410 3.4 integration, configuration, etc.
411
412 There are a few per-service configuration callouts in the project listing. These
413 are things like per-service rate limiting, blacklistd integration, etc.
414
415 I'd like to hit four major things:
416
417 * per-service ratelimiting
418 * blacklistd integration
419 * per-service logging configuration
420 * chroot() support
```

```
421
422 4. Timeline
423
424 Week 21 -> week 1 of GSOC work
425
426
427 WEEK (real) MONDAY          SUNDAY
428 Week 21      May 23          2016 May 29, 2016
429
430 Work pinning-down: A requirements document to define the grammar and specific
431 details of new configuration format.
432
433 Week 22      May 30          2016 June 5, 2016
434
435 Break up inetd into multiple compilation units, including
436 disabling of builtin services at runtime.
437
438 Implement each builtin as "inetd.<name>" and call depending on argv[0]
439
440 Week 23      June 6          2016 June 12, 2016
441 Week 24      June 13         2016 June 19, 2016
442 Week 25      June 20         2016 June 26, 2016
443
444 Implement per-service configuration:
445
446 * parser for service configuration
447 * parser for main service configuration
448 * new service structures to separate instances of a service from
449 configuration.
450 * Starting parts for logging on a per-service basis
451 * Configuration parity to current design (incl. bindhosts & ipsec config)
452
453 Deliverable: inetd that reads from sparse configuration files
454 Deliverable: documentation for inetd configuration
455
456 At this point, inetd would be otherwise "at parity"
457
458 Week 26      June 27         2016 July 3, 2016  ** FOURTH OF JULY WEEKEND
459 Week 27      July 4          2016 July 10, 2016  FOURTH OF JULY WEEKEND **
460 Week 28      July 11         2016 July 17, 2016
461
462 Implement service availability + pool sizing
463
464 Deliverable: inetd that handles early service availability
465 Deliverable: documentation covering specifics and caveats of early service
466 availability
467
468 ** FULL IMPLEMENTATION SHOULD BE COVERED IN REQUIREMENTS DOC **
469
470 Week 29      July 18         2016 July 24, 2016
471
472 Documentation week: Full documentation + conversion of a variety
473 of pkgsrc configurations
474
475 * openSSH
476 * bozohttpd
477 * BIND?
478
479 Also, a means to auto-generate from inetd classic lines into new-world
480 inetd service files. Suggesting awk (as it's common and kindof built for this)
481 or python (as stream parsers are easy to write in it)
482
483 Deliverable: A set of ready-to-run inetd service files for pkgsrc
484
485 Week 30      July 25         2016 July 31, 2016
486
487 * Logging configuration (log flags, etc)
488 * chroot() on a per-process basis
489
490 *** BONUS GOALS ***
```

491
492 * performance writeup
493 * blacklistd integration
494
495 Week 31 August 1 2016 August 7, 2016
496
497 ~ buffer week and writeups ~
498
499 ** DEFCON THIS WEEK / 4-7th IN LAS VEGAS, NV **
500
501 Week 32 August 8 2016 August 14, 2016
502
503 ~ buffer week and writeups ~
504
505 5. Prior work in the area
506
507 Some work has been done in the "overhaul inetd and friends" camp.
508
509 * Systemd: A total overhaul of Linux's init system. It has a complex, but
510 very comprehensive system for socket activation (see system.socket(5))
511 which can do wonderful things. Systemd's support includes adding pre-
512 and post-exec commands to setup/teardown a service, as well as dependencies
513 (X service must be running for Y server to run).
514 * xinetd: A *major* overhaul of inetd complete with a new configuration format.
515 this format is... Kindof JSON, Kindof not. Configuration is powerful, but
516 is poorly documented. There hasn't been a huge amount of information on what
517 the future of xinetd looks like.
518 * launchd: OSX's precursor to systemd. There have been a few attempts at porting
519 it over to the BSD world, but that looks to be in vain [LAUNCHD-GODOT]
520
521 6. Specific thoughts on timeline, implementation
522
523 This timeline is intentionally a little... sparse. There is slop allowed for
524 each section of the work as I feel comfortable, with two weeks at the end for
525 the inevitable slippage that is software work.
526
527 I'm basing much of this timeline off conservative estimates off how comfortable
528 I am with each part. I fully expect that I'll be behind or ahead in the middle,
529 however I want to account for any other space that isn't well budgeted enough.
530
531 I'm not sure what the best direction to take is. There's evidence in both
532 directions on what's better. For example, in 2009, one test showed it really
533 matters on what's being done (e.g. static files vs. PHP) in Nginx vs. Apache
534 [APACHE-NGINX-2009] whereas another [APACHE-PREFORK-NGINX] really doesn't have
535 much to say positive for Apache.
536
537 I'd like to do, as part of the writeup, addressing performance before/after to
538 validate or invalidate this approach.
539
540 Basing the configuration on sysctl makes it easy for people to add new options
541 as well as graceful fallback when an unknown/unimplemented option is used. This
542 style is also consistent with other BSDs; notably, OpenBSD which has made many
543 system tunables a sysctl mechanism.
544
545 This also comes with a cheap bonus: adding features is easy.
546
547 7. Argument for preferred configuration change
548
549 My argument for a change in the configuration format, dropping the old format
550 fundamentally comes down to two factors:
551
552 * Documentation
553 * Maintaining
554
555 NetBSD is well documented -- it's a pride of the BSDs that the BSDs ship with a
556 book that defines how they act. Moreso, the documentation is clear and easy
557 to follow. We're not in the business of making things hard.
558
559 In NetBSD (and other BSDs) follow two basic forms of configuration: the table
560 form (see fstab) and the sysctl style. I personally feel that inetd was built

561 using the table form without a lot of thought of what the future was going to
562 look like (especially with regard to SMP, manycore, highly parallel systems)
563 whereas fstab was definitely built for that future oriented style.
564
565 This brings us to maintenance. Trying to keep the old format on life support is
566 inviting two major factors: Poor support for both (see also "do thing well", a
567 typical UNIX philosophy) and potential nasty bugs cropping up in one.
568
569 The table driven form has accrued a certain amount of technical debt. The OpenBSD
570 project has worked to slowly find ways to make the inetd form less bug-ridden,
571 but the problem still exists: Technical debt needs to be paid; I personally will
572 argue that the table-driven format carries too much technical debt and possibly
573 doubling that debt is not the direction that NetBSD should take. It is not my
574 belief that the NetBSD project wants to be known for the type of hacks that are
575 normally associated with the Linux kernel.
576
577 Dropping support for the table driven style of configuration and pushing on a
578 new, well understood form fundamentally makes it easier to maintain in the near
579 future. Including an awk script that ingests the old form and spits out the new
580 form is not a huge problem, and could even be made to understand different kinds
581 of inetd configurations.
582
583
584
585 A. End references
586
587 [FREEBSD-INETD] <https://www.freebsd.org/doc/handbook/network-inetd.html>
588 [C10K] <http://www.kegel.com/c10k.html>
589 [C10M] <http://c10m.robertgraham.com/p/manifesto.html>
590 [APACHE-PREFORK-NGINX] <http://www.eschrade.com/page/performance-of-apache-2-4-with-the-event-mpm-compared-to-nginx/>
591 [APACHE-NGINX-2009] https://blog.a2o.si/2009/06/24/apache-mod_php-compared-to-nginx-php-fpm/
592 [LAUNCHD-GODOT] <http://homepage.ntlworld.com/jonathan.deboynepollard/FGA/launchd-on-bsd.html>
593
594
595