

```
110 * checked comprehensively. Further, these variables should not be tuned
111 * dynamically via "mdb -kw" or other means; they should only be tuned via
112 * /etc/system.
113 */
114 int          dtrace_destructive_disallow = 0;
115 dtrace_optval_t dtrace_nonroot_maxsize = (16 * 1024 * 1024);
116 size_t        dtrace_difo_maxsize = (256 * 1024);
117 dtrace_optval_t dtrace_dof_maxsize = (256 * 1024);
118 size_t        dtrace_global_maxsize = (16 * 1024);
119 size_t        dtrace_actions_max = (16 * 1024);
120 size_t        dtrace_retain_max = 1024;
121 dtrace_optval_t dtrace_helper_actions_max = 32;
122 dtrace_optval_t dtrace_helper_providers_max = 32;
123 dtrace_optval_t dtrace_dstate_defsize = (1 * 1024 * 1024);
124 size_t        dtrace_strsize_default = 256;
125 dtrace_optval_t dtrace_cleanrate_default = 9900990;          /* 101 hz */
126 dtrace_optval_t dtrace_cleanrate_min = 200000;             /* 5000 hz */
127 dtrace_optval_t dtrace_cleanrate_max = (uint64_t)60 * NANOSEC; /* 1/minute */
128 dtrace_optval_t dtrace_aggrate_default = NANOSEC;          /* 1 hz */
129 dtrace_optval_t dtrace_statusrate_default = NANOSEC;        /* 1 hz */
130 dtrace_optval_t dtrace_statusrate_max = (hrtime_t)10 * NANOSEC; /* 6/minute */
131 dtrace_optval_t dtrace_switchrate_default = NANOSEC;        /* 1 hz */
132 dtrace_optval_t dtrace_nspec_default = 1;
133 dtrace_optval_t dtrace_specsize_default = 32 * 1024;
134 dtrace_optval_t dtrace_stackframes_default = 20;
135 dtrace_optval_t dtrace_ustackframes_default = 20;
136 dtrace_optval_t dtrace_jstackframes_default = 50;
137 int          dtrace_msgdsz_max = 120;
138 int          dtrace_chill_max = 500 * (NANOSEC / MILLISEC); /* 500 ms */
139 hrtime_t     dtrace_chill_interval = NANOSEC;                /* 1000 ms */
140 int          dtrace_devdepth_max = 32;
141 int          dtrace_err_verbosity;
142 int          dtrace_man_interval = NANOSEC;
143 hrtime_t     dtrace_deadman_timeout = (hrtime_t)10 * NANOSEC;
144 hrtime_t     dtrace_deadman_per = (hrtime_t)30 * NANOSEC;
145
146
147 /*
```

# Hibakeresés – Mennyország és Pokol

Fischer Erik

Principal Engineer

[erik.fischer@augmentit.eu](mailto:erik.fischer@augmentit.eu)

# A hibakeresés sohasem egyszerű

- Hibák pedig léteznek
  - > Nem csak programozási hibák vannak...
  - > Konfigurációs hibák
  - > Architektúrális hibák
  - > Tervezési hibák
- A jó hibákat könnyű megtalálni
- A rosszakat meg nagyon nehéz
- Az undorítóakat meg egyszerűen lehetetlen...

# Hibakeresés – in vitro

- Végzetes hibáknál postmortem hibakeresés
  - > Core dump és debugger [mdb(1), dbx(1)]
- Tranziens hibáknál (programozási hibák, vélt vagy valós teljesítmény problémák)
  - > Ad hoc megoldások vagy szinte semmi
  - > truss(1), mdb(1), \*stat(1M), sar(1)
  - > Ha forráskód is elérhető, akkor pl. Sun Studio Performance Analyzer

# Hibakeresés – in vivo

- Szinte minden esetben csak invazív technikák
  - > Bináris instrumentálás
  - > Forrás szintű instrumentálás
  - > Interposer library-k
  - > Debug library-k
  - > Debug kernel
- Általában
  - > Durva beavatkozások
  - > Lassú
  - > Nagy az additív hibák injekciójának esélye

# A tökéletes hibakereső rendszer

- Dinamikus
- Információt gyűjt
- Éles rendszereken is biztonságosan használható
- “Teljesen” biztonságosan
- Teljesítmény veszteség “nélkül” használható

# DTrace

- Interpretált nyelv (valójában D)
  - > Probe-ok
  - > Prédikátumok
  - > Akciók
- Dinamikus (nem csak) függvény be- és kilépési pontokat instrumentáló rendszer
- Kernel modul
- Unix parancs: `dtrace(1M)`

# A DTrace jellemzői

- User és Kernel rétegben is működik
- Ismeri és érti a syscall-okat
- 50000+ probe egy átlagos Solaris 10 rendszeren (update és hardver architektúra függő)
- Alapállapotban csak root-ként működik
- 3 privilégiumot igényel: `dtrace_kernel`, `dtrace_proc` és `dtrace_user`
- Bő 800 oldalas dokumentáció

# Probe

- Az instrumentáció pontos helye egy hierarchiában (leírója egy 'n'-es)  
`<provider, module, function, probe_name>`
- A hierarchiában egy vagy akár szint is üresen hagyható
- Egy provider bocsájtja a rendelkezésünkre
- Modulokra és függvényekre bomlik
- Van neve (általában entry és return)
- Lehet minimalista módon reguláris kifejezéseket is használni: \* ? [ ]
- `dtrace -l`



# Provider

- Alapvetően az instrumentáció metodológiáját nyújtja
- Szoros kapcsolatban áll a DTrace keretrendszerrel
  - > Probe-okat ad a keretnek
  - > A keret értesíti a provider-t az egyes probe-ok aktiválásáról
  - > A provider egy probe tüzelésekor átadja a vezérlést a DTrace keretnek

# Provider-ek I.

- fbt
  - > Function Boundary Tracing
  - > Minden kernel függvény be- és kilépési pontja
  - > ~45000
- syscall
  - > syscall tábla be- és kilépési pontjai
  - > ~450
- lockstat
  - > Kernel szinkronizációs primitívek
    - Contention-event
    - Hold-event

# Provider-ek II.

- vminfo
  - > Kernel vm statisztika
- sysinfo
  - > Kernel sys statisztikák
- pid
  - > Alkalmazások függvényeinek be- és kilépési pontjainak instrumentálása
  - > Akár minden utasítást is képes instrumentálni...
- io
  - > Diszk I/O események instrumentálása

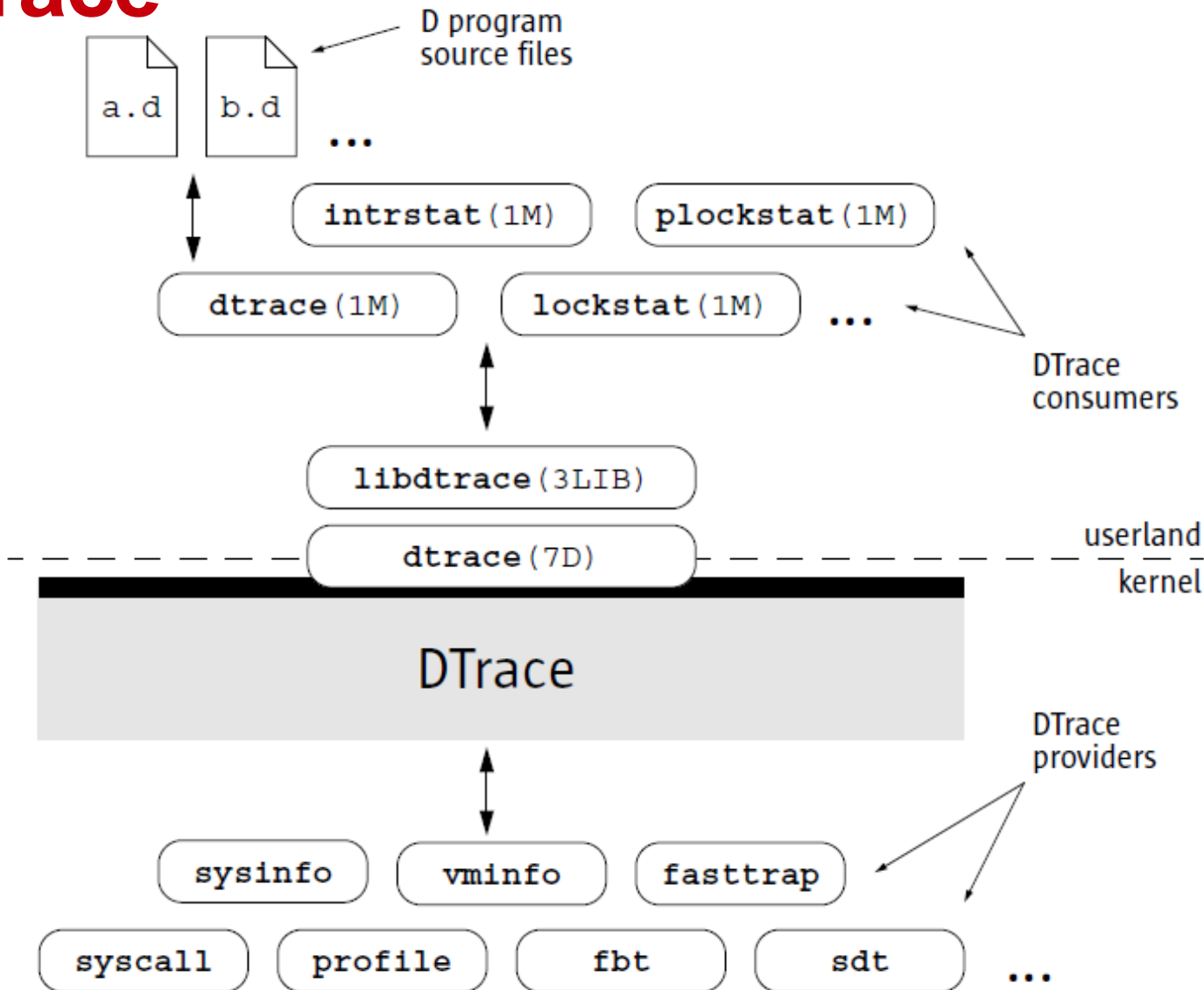
# Provider-ek III.

- proc
  - > Processz életciklus események
  - > create, start, fault, exit, exec, signal
- sched
  - > CPU ütemezési események
  - > enqueue, dequeue, sleep, wakeup, on-cpu, off-cpu
- sdt
- mib
- cpc

# Provider-ek IV.

- profile, tick
  - > Konfigurálható rátájú profiling interrupt
- fpuinfo
- plockstat
  - > User land szinkronizációs primitívek
- dtrace
  - > BEGIN
  - > END

# DTrace



# Szkriptek

- dtrace(1M) támogatja a szkripteket

```
#!/usr/sbin/dtrace -s
```

- Szkript szintaktika

```
[  
provider:module:function:name[,  
provider:module:function:name]*  
[/predicate/]  
{  
  [action;]*  
}  
]+
```

# Változók

- Skalár
  - > Szokásos típusok (char, short, int, long, long long, float, double, long double)
- Asszociatív array-ek
  - > Továbbfejlesztett Perl hash
  - > A kulcs egy tetszőleges 'n'-es
  - > Pl. array[execname, timestamp]
- String



# Beépített változók

- pid – az aktuális processz ID-je
  - ppid – a szülő processz ID-je
  - tid – az aktuális thread ID-je
  - execname – az aktuális programnév
  - \*curcpu – az aktuális CPU információs struktúrája
  - \*curthread – az aktuális thread kernel struktúrája
  - \*curpsinfo – az aktuális processz állapotleíró struktúrája
  - \*curlwpsinfo – az aktuális lwp állapotleíró struktúrája
  - lgrp – az adott CPU lgroup ID-je
  - pset – az adott CPU processor set ID-je
  - timestamp – a bootolás óta eltelt idő [ns]
  - vtimestamp – az adott processz virtualizált órája [ns]
  - probprov, probemod, probefunc, probename – a probe adatai
- ...és még sokan mások

# Thread lokális változók

- Egy név, de minden thread-nek saját adat
- Jelölése: `self->`
- A nem definiált változóknak nulla az értékük
- Ha nullázunk egy thread lokális változót, avval deallokáljuk

# Makrók

- `#[1-9]` – az adott parancssori paraméterek
- `$egid`, `$euid`, `$uid`, `$gid`, `$pid`, `$pgid`, `$ppid`
- `$projid`
- `$sid`
- `$taskid`

# Prédikátumok

- Akciók csak bizonyos feltételek megléte esetén történő érvényesítése
- Egy logikai D kifejezés
- Ha a kifejezés igaz, akkor és csak akkor hat az akció

# Az akciókról általában I.

- Ez a valódi program
- Általában információt rögzít
- Ritkán információt változtat meg
- Akciók
  - > Műveletek változókkal és array-ekkel
  - > **NINCS FLOATING POINT ARITHMETIKA**
  - > Függvényhívások
- Bizonyos programozási szerkezetek nem léteznek
  - > Elágazás (if)
  - > Ciklus (for, while, until)

# Az akciókról általában II.

- Arithmetikai műveletek
- Logikai műveletek
- Komparálás
- Értékadás
- Feltételes kifejezés

```
x = i == 0 ? "zero" : "non-zero";
```

# Destruktív akciók

- A `-w` flag-gel engedélyezni kell
- `stop()` – megállítja az aktuális processzt
- `raise()` – szignált küld az adott processznek
- `breakpoint()` – kernel breakpoint triggerelése
- `panic()` – kernel pánik
- `chill()` – adott időre leállítja a processzt
- `copyout()` és `copyoutstr()` – adott változóból adott számú adatot egy memória helyre másol

# Aggregációk

- DTrace függvény
- Egy aggregáció olyan  $f(x)$  függvény, ahol  $x$  egy tetszőleges  $n$  elemű adathalmaz és:  
$$f(f(x_0) \cup f(x_1) \cup \dots \cup f(x_n)) = f(x_0 \cup x_1 \cup \dots \cup x_n)$$
- Ilyenek a count, a sum, az avg, a max és a min (pl. a median nem az, éppen ezért nincs is ilyen)
- Speciális aggregáció a quantize és az lquantize



# Normalizáció és egyéb műveletek

- `normalize()`
  - > Nem valódi normálás, csak egy faktor tárolása
  - > Alapvetően megjelenítéshez
- `denormalize()`
  - > Egy meglévő norma faktor törlése
- Aggregáció törlése
  - > `clear()` - csak az értéket törli, a kulcsok maradnak
- Aggregáció csonkolás
  - > `trunc()` - az aggregáció csonkolása

# Spekulatív követés

- Bizonyos esetekben csak speciális feltételek együttállása esetén kellenek adatok
- Viszont minden más esetben minden begyűjtött adatot el kell dobni
- Nem szűrhető prédikátumokkal
- Tipikusan hibák előfordulásakor gyűjtendő adatok, speciális program ágak monitorozása, etc.

# Anonymous követés

- Fogasztóhoz nem kötött nyomkövetés
- Pl. A boot folyamat követése
- /kernel/drv/dtrace.conf fájlba bekerülő direktívák
  - > A dtrace(7D) betöltésekor automatikusan aktiválódik

# Kernel instrumentáció

- A belépési utasítás branch always utasításra cserélése
- Ugrás a DTrace keretbe
- A feladatok végrehajtása
- Az elcserélt utasítás végrehajtása
- Ugrás vissza oda, ahonnan jöttünk

# User processz instrumentáció

- Alapvetően a kernel instrumentáció logikájával azonos
- A probléma, hogy user land-ből kernel land-be kell átkerülni
- A megoldás egy trap utasítás inzertálása
- Következménye egy speciális DTRACE trap entry

# Syscall instrumentáció

- A syscall-ok valójában trap-eken keresztül hívódnak meg
- A syscall-ok táblákban vannak és ide indexelünk a trap "argumentumával"
- A megoldás az adott syscall bejegyzés cseréje egy DTRACE\_SYSTRACE\_SYSCALL bejegyzésre

# DTrace a nagyvilágban

- Portok
  - > FreeBSD
  - > OS X
  - > Linux (Oracle Enterprise Linux), folyamatban...
- Az OpenSolaris megszűnésével keletkezett *illumos* kernelre épülő rendszerekben:
  - > Nexenta NexentaStor és illumian
  - > OpenIndiana
  - > Joyent SmartOS
  - > OmniOS
  - > DEYOS (egyelőre még nem publikus)
- Jelentős változások az Oracle Solaris 11-ben

# DTrace kérdések

- Az eltérő rendszerek kezelés
- Univerzális Dtrace szkriptek
- Univerzális és dinamikus transzlátorok



# Példák

```
110 * checked comprehensively. Further, these variables should not be tuned
111 * dynamically via "mdb -kw" or other means; they should only be tuned via
112 * /etc/system.
113 */
114 int          dtrace_destructive_disallow = 0;
115 dtrace_optval_t dtrace_nonroot_maxsize = (16 * 1024 * 1024);
116 size_t        dtrace_difo_maxsize = (256 * 1024);
117 dtrace_optval_t dtrace_dof_maxsize = (256 * 1024);
118 size_t        dtrace_global_maxsize = (16 * 1024);
119 size_t        dtrace_actions_max = (16 * 1024);
120 size_t        dtrace_retain_max = 1024;
121 dtrace_optval_t dtrace_helper_actions_max = 32;
122 dtrace_optval_t dtrace_helper_providers_max = 32;
123 dtrace_optval_t dtrace_dstate_defsize = (1 * 1024 * 1024);
124 size_t        dtrace_strsize_default = 256;
125 dtrace_optval_t dtrace_cleanrate_default = 9900990;          /* 101 hz */
126 dtrace_optval_t dtrace_cleanrate_min = 200000;             /* 5000 hz */
127 dtrace_optval_t dtrace_cleanrate_max = (uint64_t)60 * NANOSEC; /* 1/minute */
128 dtrace_optval_t dtrace_aggrate_default = NANOSEC;          /* 1 hz */
129 dtrace_optval_t dtrace_statusrate_default = NANOSEC;        /* 1 hz */
130 dtrace_optval_t dtrace_statusrate_max = (hrtime_t)10 * NANOSEC; /* 6/minute */
131 dtrace_optval_t dtrace_switchrate_default = NANOSEC;        /* 1 hz */
132 dtrace_optval_t dtrace_nspec_default = 1;
133 dtrace_optval_t dtrace_specsize_default = 32 * 1024;
134 dtrace_optval_t dtrace_stackframes_default = 20;
135 dtrace_optval_t dtrace_ustackframes_default = 20;
136 dtrace_optval_t dtrace_jstackframes_default = 50;
137 dtrace_optval_t dtrace_jstackstrsize_default = 512;
138 int          dtrace_msgdsize_max = 128;
139 hrtime_t     dtrace_chill_max = 500 * (NANOSEC / MILLISEC); /* 500 ms */
140 hrtime_t     dtrace_chill_interval = NANOSEC;                /* 1000 ms */
141 int          dtrace_devdepth_max = 32;
142 int          dtrace_err_verbose;
143 int          dtrace_deadman_interval = NANOSEC;
144 hrtime_t     dtrace_deadman_timeout = (hrtime_t)10 * NANOSEC;
145 hrtime_t     dtrace_deadman_min = (hrtime_t)30 * NANOSEC;
146
147 /*
```

Fischer Erik  
Principal Engineer  
erik.fischer@augmentit.eu