#### Samba and the road to 100,000 users

Presented by Andrew Bartlett Samba Team - Catalyst // 16 Jan 2017

### catalyst 🞑

open source technologists





Samba is a member project of the Software Freedom Conservancy



#### **Andrew Bartlett**

- Samba developer since 2001
- Working on the AD DC since soon after the start of the 4.0 branch, since 2004!
  - Driven to work on the AD DC after being a high school Systems Administrator
- Working for Catalyst in Wellington since 2013
  - Now leading a team of 5 Catalyst Samba Engineers
- These views are mine alone
- Please ask questions during the talk

![](_page_1_Picture_8.jpeg)

![](_page_1_Picture_10.jpeg)

![](_page_1_Picture_11.jpeg)

#### In 2017, Samba is released Fast!

- New release cycle
  - **Strict** 6 months cycle
  - Supported for 18 months
- Support status
  - Samba 4.0 is 4 years ago old now
  - Samba 4.6 about to be released
  - Samba 3.6 and 4.0, 4.1, 4.2 already out of security support
    - Your distributor may be providing extended support for Samba

![](_page_2_Picture_10.jpeg)

![](_page_2_Picture_11.jpeg)

#### But also much faster as an AD DC

- In a two-hour benchmark adding users and adding to four groups:
  - Samba 4.4: 26,000 users
  - Samba 4.5: 48,000 users
  - Samba 4.6: 55,000 users
  - With pending patches: 85,000 users!
    - The first 55,000 added in just 50mins

![](_page_3_Picture_8.jpeg)

![](_page_3_Picture_9.jpeg)

#### Scale is important to us

- Every user account implies a computer account also
  - Computers are domain joined and get 'user' objects
- Samba3 deployed widely using OpenLDAP for the hard work
  - Ideally we need to match that scale!
- We really want to remove barriers, both real and perceived to Samba's use
  - Not reasonable to ask that Samba be deployed on the very edge of its capability

![](_page_4_Picture_8.jpeg)

![](_page_4_Picture_9.jpeg)

#### **Rebuilding Samba for performance**

- Once we started looking at performance, we quickly found things to fix
- Performance issues, not bugs, are now the biggest area of work!
  - Customers deploying Samba at scale

- Customers growing and very keen to keep Samba
- Very glad to be the backbone of some multi-national corporate networks!

![](_page_5_Picture_6.jpeg)

![](_page_5_Picture_7.jpeg)

#### **Replication as a performance bottleneck**

- So what if it takes time to add 10,000 users or so?
  - Companies can't hire that fast anyway
  - Even students don't enrol that fast!
- Biggest bottleneck is adding new DCs to Samba domains
  - e.g. opening a new office

![](_page_6_Picture_7.jpeg)

#### **Improvements in Samba 4.5**

- Samba 4.5 addressed major issues with the client-side of replication
  - 3 of the 4 O(n<sup>2</sup>) loops removed
  - Critical as these were under the transaction lock
- Turned on graph (rather than all to all) replication by default
  - Previously every Samba DC would notify every other Samba DC about changes
  - This could trigger a short replication storm

![](_page_7_Picture_8.jpeg)

![](_page_7_Picture_9.jpeg)

#### Some improvement in 4.6

- Samba 4.6 will avoid over-replication of links
  - When replicating from server A, we also ask is what changes it got from B
  - That means we don't need to ask B for changes directly
  - We did this for attributes, but didn't do this for links previously
- Faster parsing of links also improved performance around 20% for some tasks
  - Avoid sscanf() and malloc()

![](_page_8_Picture_7.jpeg)

![](_page_8_Picture_8.jpeg)

![](_page_8_Picture_9.jpeg)

#### Two steps forward, one back

- But there appears to be a regression with Samba 4.5 and 4.6
  - When joining new Samba 4.5 and 4.6 DCs we request 'parents before children'
  - Sorting the DB (for initial replication) appears to take longer than a timeout period
  - Stefan Metzmacher from SerNet has a likely fix
  - My team is looking to isolate and confirm or improve that fix

![](_page_9_Picture_6.jpeg)

![](_page_9_Picture_7.jpeg)

#### Before optimisation: Samba 4.4

 Adding a user and adding that user to four groups in a two-hour limit

![](_page_10_Figure_2.jpeg)

open source technologists

,S'AMBA

![](_page_10_Picture_5.jpeg)

#### Much improved scale factors: two-hour limit

![](_page_11_Figure_1.jpeg)

SAMBA

![](_page_11_Picture_3.jpeg)

#### Supporting more users on each DC

- Hoping to avoid needing to run extra DCs to spread the load
- Samba 4.6 removes single-process restrictions on NETLOGON
  - Really important for 802.1x backed wireless authentication
  - Unbreak the WiFi and watch the DC melt instead :-(
- Samba 4.7 will support a multi-process LDAP server
  - Easy to turn on in the code
  - Currently fork() and cleanup for exit() costs are too high

![](_page_12_Picture_8.jpeg)

![](_page_12_Picture_9.jpeg)

![](_page_12_Picture_10.jpeg)

#### Finding a Bottleneck: Number of group members

- The slowest part of the code was not user or group objects
- Main cost was each group member (a link and backlink)
- As discussed, quite painful during replication
  - Client-side link processing is during the transaction lock
  - 4 different O(n<sup>2</sup>) loops found!
  - talloc() and talloc\_free() quite expensive

![](_page_13_Picture_7.jpeg)

![](_page_13_Picture_8.jpeg)

![](_page_13_Picture_9.jpeg)

#### Flame graphs

- We used linux perf and Brendan D. Gregg's Flame Graphs
- If you have performance issues with Samba:
  - Install the linux-perf tools
  - Clone https://github.com/brendangregg/FlameGraph
  - Follow http://www.brendangregg.com/FlameGraphs/cpuflamegraphs.html
  - Send us the Interactive SVG
    - Sensitive user data not included, just function names!

![](_page_14_Picture_8.jpeg)

![](_page_14_Picture_9.jpeg)

![](_page_14_Picture_10.jpeg)

talloc_free_in talloc_free_in talloc_free_int talloc_free_int talloc_free_int talloc_free_int talloc_free_int talloc_free_int talloc_free_int talloc_free_int talloc_free_int	asanm Idb_val_equal_exact Idb_msg_find_val Idb_msg_find_val Idb_callback tevent_common_loop_timer epoll_event_loop_once std_event_loop_once std_event_loop_once std_event_loop_once tevent_loop_once Idb_wait dsdb_module_modify commit _commit _commit _commit	Idif. 1. I.	GUID_comp GUID_equal repImd_add		i dn l ltd ltdb ltdb ltdb b ltdb b ltdb ca tevent epoll_e std_eve tevent ldb_wait ldb_ext dsdb_r libnet
descriptor_prepare_commit					dsdb_r
ldb_transaction_prepare_commit					libnet
ldb_transaction_commit					py_net
py_Idb_transaction_commit					PyEval
PyEval_EvalFrameEx				[unknown]	_dl_mak

[unknown] python all

#### Flame graphs are interactive

- When used in a web browser
- I blogged about this for catalyst:
  - Burning Samba with perf and FlameGraph
  - https://catalyst.net.nz/blog/burning-samba-perf-and-flamegraph

![](_page_16_Picture_6.jpeg)

![](_page_16_Picture_7.jpeg)

## **Performance** graphs from <sup>10</sup> March 2016

![](_page_17_Figure_1.jpeg)

#### The difference a sorted list makes!

- Our code needs to find group members to support add/delete/modfy
- Previously, we had to parse every link

• Now we sort by GUID, and so can do a binary search

![](_page_18_Picture_6.jpeg)

![](_page_18_Picture_7.jpeg)

# Pending changes for sorting links

• Over a 60% drop in time for some tests

![](_page_19_Figure_2.jpeg)

![](_page_19_Picture_3.jpeg)

![](_page_19_Picture_4.jpeg)

![](_page_19_Picture_5.jpeg)

#### The future for performance

- Remove other O(n) and O(n<sup>2</sup>) operations
  - Multi-valued attribute handling
- Better index handling
  - Our current index code is still very much a first pass
  - Proposal to move to a GUID based index
- Reaching the limits for the current DB:
  - memcpy() and memmove() from ldb\_tdb transactions are 20% of the time

![](_page_20_Picture_8.jpeg)

![](_page_20_Picture_9.jpeg)

![](_page_20_Picture_10.jpeg)

#### **Lightening Memory-mapped Database from Symas**

- The company behind OpenLDAP
- Built by Howard Chu to make OpenLDAP fly
- LMDB backend prototyped by Red Hat for sssd
  - Appears to be 3 times faster for some operations
- Garming Sam has been working on reimplementation
  - Preparing it in a way that could be submitted
  - Based more tightly on the TDB LDB backend
- As of Friday it successfully ran provision!

![](_page_21_Picture_10.jpeg)

![](_page_21_Picture_11.jpeg)

#### Maintaining Performance and scale

- Large scale operation needs to be part of Samba's autobuild
- Project to develop a new performance metric for Samba domains
  - Currently awaiting client approval
- Ongoing graphing of performance measurements
  - Try to spot regressions before they get too old

![](_page_22_Picture_7.jpeg)

![](_page_22_Picture_8.jpeg)

#### Help wanted!

- For the performance metric tool I need to calibrate it
- I need volunteers running AD willing to run a tshark script
  - Windows or Samba AD welcome
  - What does your busy hour look like?
  - What is the pattern of requests?
- E-mail abartlet@samba.org if you can help

![](_page_23_Picture_7.jpeg)

![](_page_23_Picture_8.jpeg)

#### **Beyond performance**

- Inter-forest trusts
  - Because sometimes sharding the data is really the right approach
  - Initial support in Samba 4.3 but more work needed
- Inter-domain trusts
  - To allow migration from per-department Samba domains
  - Still pending further development
  - Most companies move to one domain, one forest

![](_page_24_Picture_8.jpeg)

![](_page_24_Picture_9.jpeg)

![](_page_24_Picture_10.jpeg)

#### **Beyond just pure AD**

- What would make Samba compelling for your networks?
- Can be integrate better with POSIX systems?
  - Become the natural directory for Linux networks too?
  - Can Mac OS X be better supported?
- Samba 4.5 includes a Samba-specific password sync extension

![](_page_25_Picture_7.jpeg)

![](_page_25_Picture_8.jpeg)

#### **MIT Kerberos**

- Blocking Samba being a part of SLES and RHEL
- Still in progress
- Very important as Heimdal Upstream only just restarted releases
- I'm hoping to update the Heimdal copy as well
  - 5 year old security code is not a great thing

![](_page_26_Picture_7.jpeg)

![](_page_26_Picture_8.jpeg)

#### **OpenLDAP** backend

- The original 'make samba faster' proposal
- Sadly little progress other than a presentations in 2015
- No public code
- I'm hesitant about another lift-and-shift like MIT Kerberos
- Prefer to fix one identified, isolated issue at a time
- Incremental progress can pay off now

![](_page_27_Picture_8.jpeg)

# Become an OFFICIAL CONSERVANCY SUPPORTER!

![](_page_28_Picture_1.jpeg)

![](_page_28_Picture_2.jpeg)

![](_page_28_Picture_3.jpeg)

#### **Catalyst's Open Source Technologies – Questions?**

![](_page_29_Figure_1.jpeg)

Want to join Catalyst? We love Linux-passionate Sysadmins and our Samba dev workload is growing: talk to me in the hallway track!

![](_page_29_Picture_4.jpeg)

![](_page_29_Picture_5.jpeg)