# Predicting Application Performance Before Moving to the Cloud

Moving servers and applications from an existing data centre into the cloud (or a new data centre) is a project that comes with a long list of risks and unknowns.

How applications will perform – from a user's perspective – is always an unknown when systems are moved to an alternative data centre.

This presentation discusses a process and technique that measures the exact behaviour of applications.

Early analysis of the exact behaviour of applications can remove the "performance risk" and can reveal:

- Which applications/servers can be moved with no (or little) expected performance impact?

- Which applications/servers will have an expected performance impact – and how much?

- Exactly why are applications predicted to have a large performance impact?

- Workable remediation options that can be proposed and implemented.

- Remediation options that can be implemented and tested well in advance of the physical server moves.

A similar traffic analysis can be used to determine which servers should be moved as groups.

# What Affects Application Performance

Network propagation delay is one of the main and often overlooked variables that can have a large performance impact. If new servers are further away from users, then the minimum network round trip time (RTT) to each server is increased. This will affect individual application transactions as well as "user functions" that involve many transactions, possibly to many servers.

Efficient applications (with a minimal number of transactions to back-end servers) can probably be moved with little performance impact.

However, inefficient "chatty" applications (that perform many transactions) can be severely impacted when the minimum per-server RTT is increased and each transaction time is correspondingly increased.

It is very common for application teams to not fully appreciate how their applications actually work "under the covers" and at the network level.

Discussed here is a mechanism that was used to evaluate key applications and user functions to:

- Predict the user impact of proposed server moves.

- Understand exactly why there would be an impact.

- Confidently provide remediation options for applications.

- Provide a general assessment of application efficiency.

This presentation can be downloaded (as a PDF) from:
 http://www.networkdetective.com.au/downloads

# Terminology

| Term | Definition |
|------|------------|
| Server Transaction | A single paired request and response message from a user's workstation (the *client*) to a server. Is at least one round trip (loop) - but may also involve multiple loops for network delivery of the request and response. |
| User Function | All the activity between a user clicking on a button (or similar event) and the eventual rendering of a new display on the user's screen. It may involve a large number of server transactions, to a variety of different servers. Several server transactions may run concurrently. |

Example "User Functions" for different industries:

| Insurance | Airline | Bank | Real Estate |
|-----------|---------|------|-------------|
| • Launch "Teller" application.<br>• Logon to application.<br>• List all accounts.<br>• Perform a cash withdrawal.<br>• Perform a cheque deposit.<br>• Transfer between accounts.<br>• Foreign currency transaction | • Launch "Check-In" app.<br>• Check-in one passenger.<br>• Check-in multiple passengers.<br>• Add a bag to booking.<br>• List upgrade prices.<br>• List alternate flights.<br>• Print list of passengers. | • Launch "Teller" application.<br>• Logon to application.<br>• List all accounts.<br>• Perform a cash withdrawal.<br>• Perform a cheque deposit.<br>• Transfer between accounts.<br>• Foreign currency transaction. | • Create new property for sale.<br>• Create new tenant record.<br>• Add new tenant to property.<br>• Process a rental payment.<br>• Print payment receipt.<br>• List available properties.<br>• Print monthly statement. |

# Timing of Server Transactions

There are four major components that make up the timing of individual server transactions:

| Time Category | Definition | Display Colour |
|---|---|---|
| Client time | Time spent internally to prepare the next server request | Grey |
| Request message transfer time | Time between the first and last packets of the request. (Zero for single packet requests.) | Green |
| Server time | Time the server takes to process the request and prepare the response. Measured as time of last request packet to first server response packet. Includes one RTT when captures are taken at the user end. | Yellow Orange |
| Response message transfer time | Time between the first and last packets of the response. (Zero for single packet responses.) | Blue |

Note:
Application overheads, such as user authentication, are treated as server transactions within the overall user function.

Likewise, overheads such as TCP 3-way handshakes and SSL handshakes are also included as "transactions" within the user functions (since they also add loop delays that can be "felt" by the user).

# Method Used to Measure Applications

Experienced users were chosen to demonstrate various user functions for each application.

Users were asked to start significant user transactions at given times so that each function's start time and approximate duration could be noted. For example, "press enter when the time is at 00 or 30 seconds". This helps to visually separate network traffic related to each user function.

Network traffic conveying user functions was captured, either on or close to the user's workstation, using Wireshark. Australian software, "NetData Pro," was then used to identify all the server transactions involved in the user functions and measure the four major components of response times.

NetData calculates the minimum round trip time (designated as "loop delay") to each server – as well as the number of transactions and other TCP loops involved for each server. Each user function can usually be characterised and documented with a single chart. The charts also display the various transaction types, showing any network or application inefficiencies that can form the basis for proposed remediation ahead of the server moves.

When there are concurrent TCP connections to one or more servers, NetData chooses a "critical path" through the transactions to count total loop delay (so that round trips that occur in parallel are not double counted).

The expected change of overall timing of the user function is then calculated by informing NetData of the proposed new RTT for each server. Each server can have a different "new RTT". The calculated overall performance impact is displayed in red.

# Example Outputs from the Analysis

The following slides show some example "user function" summary charts from a real-life project.
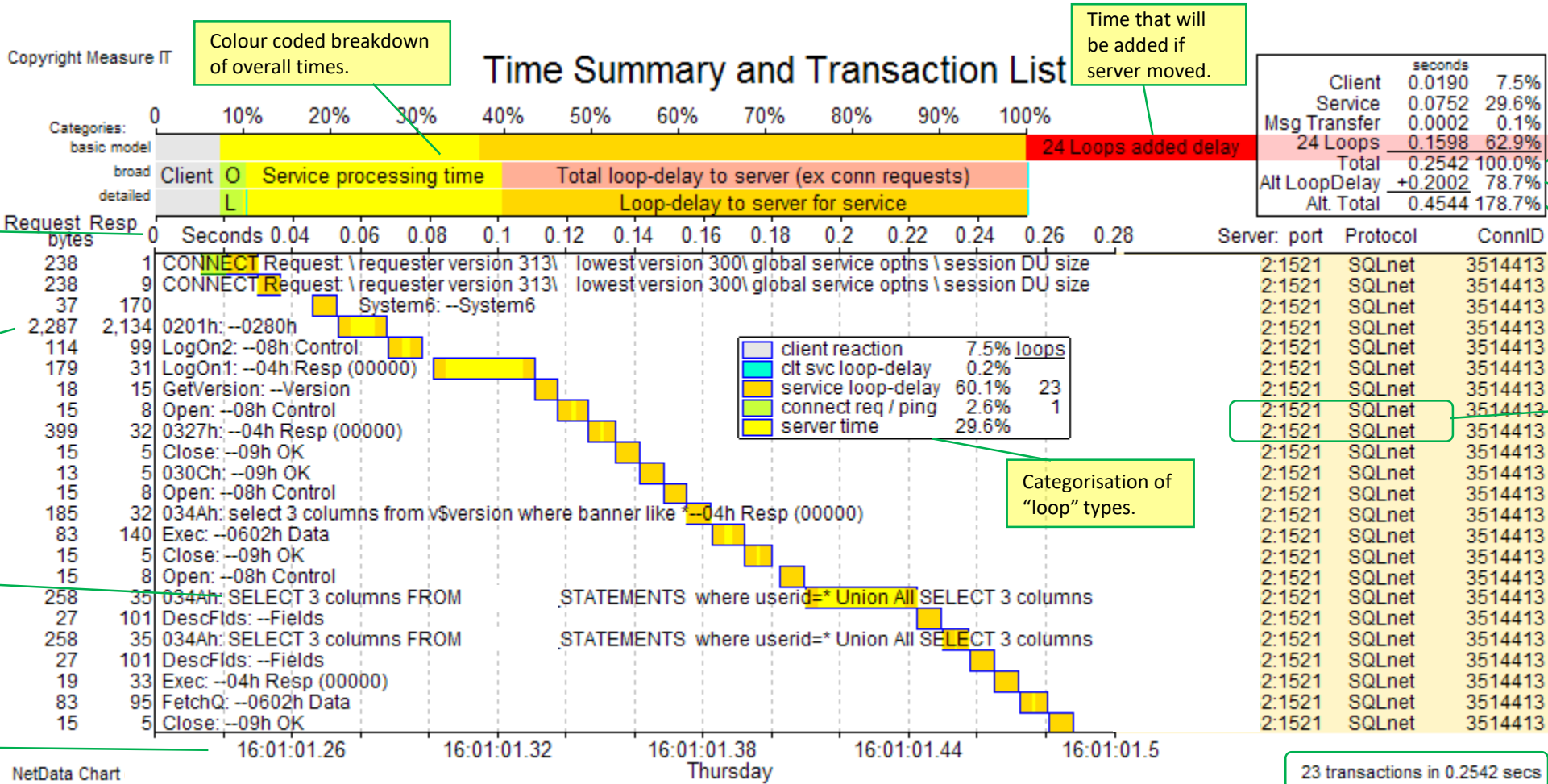
Benefits of these charts:

- Show the expected performance impact to each user function.

- Show the exact reasons for the performance impact.

- Explaining it all to an application team is made easier and clearer by presenting charts rather than just text.

- The application team may learn things about their application behaviour that they hadn't previously known.

# How to Interpret the Charts

We'll begin with a relatively simple chart – which shows the 23 Oracle database transactions that were invoked for a particular end-user task that took about a quarter of a second to complete. In this case, there is only one SQL transaction per row.
Colours are used to differentiate client/server processing times, total network RTT and network time to deliver packet data.
In this case, the original 254 ms function will become 454 ms due to the added 200 ms of loop delay.



Colour coded breakdown of overall times.

Time that will be added if server moved.

Current function timing.

Timing change due to extra loop delay.

Overall timing after server moves.

Elapsed time.

Request / Response byte count for all transactions in row.

Server name/IP, port, application protocol.

Transaction descriptions.

Categorisation of "loop" types.

TCP Connection ID. We can tell if multiple conns are used.

Horizontal axis is time-of-day.

Transaction count is always here.

**Time Summary and Transaction List**

| | seconds | |
|---|---|---|
| Client | 0.0190 | 7.5% |
| Service | 0.0752 | 29.6% |
| Msg Transfer | 0.0002 | 0.1% |
| 24 Loops | 0.1598 | 62.9% |
| Total | 0.2542 | 100.0% |
| Alt LoopDelay | +0.2002 | 78.7% |
| Alt. Total | 0.4544 | 178.7% |

Copyright Measure IT

Categories:
basic model
broad — Client  O  Service processing time — Total loop-delay to server (ex conn requests) — 24 Loops added delay
detailed — L — Loop-delay to server for service

| client reaction | 7.5% | loops |
|---|---|---|
| clt svc loop-delay | 0.2% | |
| service loop-delay | 60.1% | 23 |
| connect req / ping | 2.6% | 1 |
| server time | 29.6% | |

| Request Resp bytes | | | Server: port | Protocol | ConnID |
|---|---|---|---|---|---|
| 238 | 1 | CONNECT Request \ requester version 313\ lowest version 300\ global service optns \ session DU size | 2:1521 | SQLnet | 3514413 |
| 238 | 9 | CONNECT Request \ requester version 313\ lowest version 300\ global service optns \ session DU size | 2:1521 | SQLnet | 3514413 |
| 37 | 170 | System6: --System6 | 2:1521 | SQLnet | 3514413 |
| 2,287 | 2,134 | 0201h: --0280h | 2:1521 | SQLnet | 3514413 |
| 114 | 99 | LogOn2: --08h Control | 2:1521 | SQLnet | 3514413 |
| 179 | 31 | LogOn1: --04h Resp (00000) | 2:1521 | SQLnet | 3514413 |
| 18 | 15 | GetVersion: --Version | 2:1521 | SQLnet | 3514413 |
| 15 | 8 | Open: --08h Control | 2:1521 | SQLnet | 3514413 |
| 399 | 32 | 0327h: --04h Resp (00000) | 2:1521 | SQLnet | 3514413 |
| 15 | 5 | Close: --09h OK | 2:1521 | SQLnet | 3514413 |
| 13 | 5 | 030Ch: --09h OK | 2:1521 | SQLnet | 3514413 |
| 15 | 8 | Open: --08h Control | 2:1521 | SQLnet | 3514413 |
| 185 | 32 | 034Ah: select 3 columns from v$version where banner like *--04h Resp (00000) | 2:1521 | SQLnet | 3514413 |
| 83 | 140 | Exec: --0602h Data | 2:1521 | SQLnet | 3514413 |
| 15 | 5 | Close: --09h OK | 2:1521 | SQLnet | 3514413 |
| 15 | 8 | Open: --08h Control | 2:1521 | SQLnet | 3514413 |
| 258 | 35 | 034Ah: SELECT 3 columns FROM STATEMENTS where userid=* Union All SELECT 3 columns | 2:1521 | SQLnet | 3514413 |
| 27 | 101 | DescFlds: --Fields | 2:1521 | SQLnet | 3514413 |
| 258 | 35 | 034Ah: SELECT 3 columns FROM STATEMENTS where userid=* Union All SELECT 3 columns | 2:1521 | SQLnet | 3514413 |
| 27 | 101 | DescFlds: --Fields | 2:1521 | SQLnet | 3514413 |
| 19 | 33 | Exec: --04h Resp (00000) | 2:1521 | SQLnet | 3514413 |
| 83 | 95 | FetchQ: --0602h Data | 2:1521 | SQLnet | 3514413 |
| 15 | 5 | Close: --09h OK | 2:1521 | SQLnet | 3514413 |

16:01:01.26   16:01:01.32   16:01:01.38   16:01:01.44   16:01:01.5
Thursday

NetData Chart

23 transactions in 0.2542 secs

# An Oracle Application

In this example, there is still only one TCP connection to a single Oracle database server. There are 32 transactions that took 178 ms to complete. In this case, each row contains multiple SQL transactions of a particular type.

The extra loop-delay would add 746 ms to the 178 ms function – to become 924 ms.  Perhaps not noticeable to a human?
Could the time be reduced by combining the 27 SQL requests into one?  We also need to ensure the new environment can cater for the large data transfers.

Colour coded breakup of various times. There's a lot of network time (blue) due to the large responses.

Copyright Measure IT

## Time Summary and Transaction List

Elapsed time.

| | seconds | |
|---|---|---|
| Client | 0.0466 | 26.2% |
| Service | 0.0447 | 25.1% |
| Msg Transfer | 0.0829 | 46.5% |
| 50 Loops | 0.0040 | 2.2% |
| Total | 0.1782 | 100.0% |
| Alt LoopDelay | +0.7460 | 418.5% |
| Alt. Total | 0.9242 | x 5.2 |

This is always the most critical table to examine.

Categories: 0   10%   20%   30%   40%   50%   60%   70%   80%   90%   100%

basic model — Total message-transfer time ex loops — 50 Loops added delay

broad — Client reaction time — Service processing — L Service response-message time ex loops

detailed — max. 0.0050 secs — L

Request Resp bytes   0   Seconds   0.04   0.06   0.08   0.1   0.12   0.14   0.16   0.18   0.2   Server:  port   Protocol   ConnID

| 109 | 246 | SELECT: select * from attachment where att_sn = *; Close--1017h Format Control (00000) | Oracle:1521 | SQLnet | 2985397 |
| 51 | 194 | FetchQ: --0602h Data (*****) | Oracle:1521 | SQLnet | 2985397 |
| 231 | 258 | 2 x 0360h: --08h | Oracle:1521 | SQLnet | 2985397 |
| 117 | 129 | 0360h: --08h | Oracle:1521 | SQLnet | 2985397 |
| 3,210 | 891K | 27 x 0360h: --0EFEh Data | Oracle:1521 | SQLnet | 2985397 |

Only one TCP connection.

27 similar SQL transactions totalled 3 KB of request data and 891 KB of response data.

| | | loops |
|---|---|---|
| client reaction | 26.2% | |
| clt svc loop-delay | 0.2% | |
| clt loop-d in msg | 0.1% | |
| service loop-delay | 1.3% | 32 |
| svr loop-d in msg | 0.7% | 18 |
| server time | 25.1% | |
| respnse msg time | 46.5% | |

Categorisation of "loop" types

15:40:01.4   15:40:01.44   15:40:01.48   15:40:01.52   15:40:01.56

NetData Chart   Tuesday

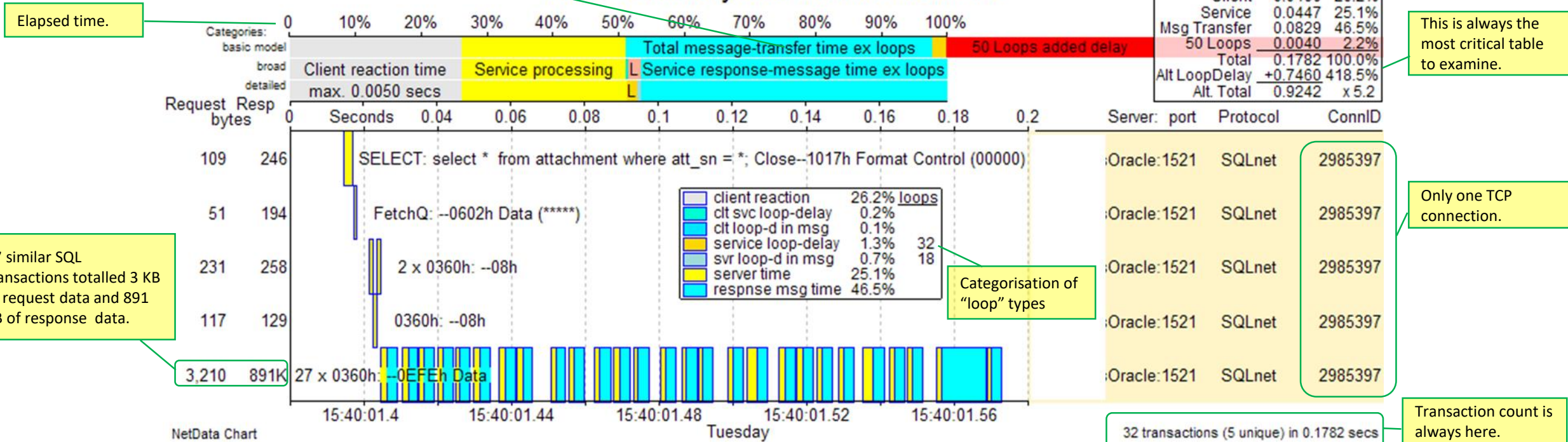32 transactions (5 unique) in 0.1782 secs

Transaction count is always here.

# Inefficient Oracle Application

In this example, there is still only one TCP connection to a single Oracle database server.  However, there are 21,948 transactions that took 33 seconds to complete. Retrieving thousands of rows of data either one or eleven at a time is "expensive" in both client and server processing - but extremely expensive in network round trip times.  This current 33 second task would become 6 minutes.

Simple application coding changes (or driver settings) could easily eliminate 21,939 round trips and use the network more efficiently.

Copyright Measure IT

## Time Summary and Transaction List

|  | seconds |  |
|---|---|---|
| Client | 13.6776 | 41.5% |
| Service | 10.4743 | 31.8% |
| Msg Transfer | 0.0000 | 0.0% |
| 21948 Loops | 8.8304 | 26.8% |
| Total | 32.9822 | 100.0% |
| Alt LoopDelay | 320.3896 | x 9.7 |
| Alt. Total | 353.3718 | x 10.7 |

Categories:
- basic model
- broad — Client reaction time max. 0.0178 secs — Service processing time — Loops / Total loops
- detailed — Loops / Service loops

21948 Loops added delay

This is much more serious!

| Request bytes | Resp | | Server: port | Protocol | ConnID |
|---|---|---|---|---|---|
| 33 | 1,536 | Describe: GREPDATAFILE--Description | :1521 | SQLn312V | 3122236 |
| 234 | 1,103 | BEGIN: BEGIN GREPDATAFILE(*,*,*,Null,*,Null,Null,*,*,Null,Null,*,Null,*,*,*,Null,Null,:RCUR ); | :1521 | SQLn312V | 3122236 |
| 52 | 228 | FetchQ: --0602h Data Control (00000) | :1521 | SQLn312V | 3122236 |
| 353K | 2,636K | 20781 x Fetch Continuation 1 | :1521 | SQLn312V | 3122236 |
| 96 | 18 | 3 x Commit | :1521 | SQLn312V | 3122236 |
| 177 | 370 | SELECT: SELECT 2 columns FROM _DATAFILE_LOG"--1017h Format Control | :1521 | SQLn312V | 3122236 |
| 19,686 | 537K | 1158 x Fetch Continuation 11 | | | 3122236 |
| 166 | 242 | 035Eh: SELECT 5 columns FROM _DATAFILE_LOG"; Close--1017h Format Control | | | 3122236 |
| 381 | 72 | INSERT: INSERT INTO _DATAFILE_LOG"  (5 columns) VALUES (*); Close--08h | :1521 | SQLn312V | 3122236 |

20,782 rows of data (2.6 MB) were retrieved one at a time, resulting in 20,782 loops.

12,749 rows of data (0.5 MB) were retrieved 11 at a time, resulting in 1,159 loops.

This is not one blue rectangle, it is 20,781 thin blue lines.

Still only one TCP connection.

| client reaction | 41.5% loops |
|---|---|
| clt svc loop-delay | 8.0% |
| service loop-delay | 18.7%21948 |
| server time | 31.8% |

NetData Chart

10:26:54   10:27:00   10:27:06   10:27:12   10:27:18   10:27:24   10:27:30
Friday

21,948 transactions in 32.9822 secs
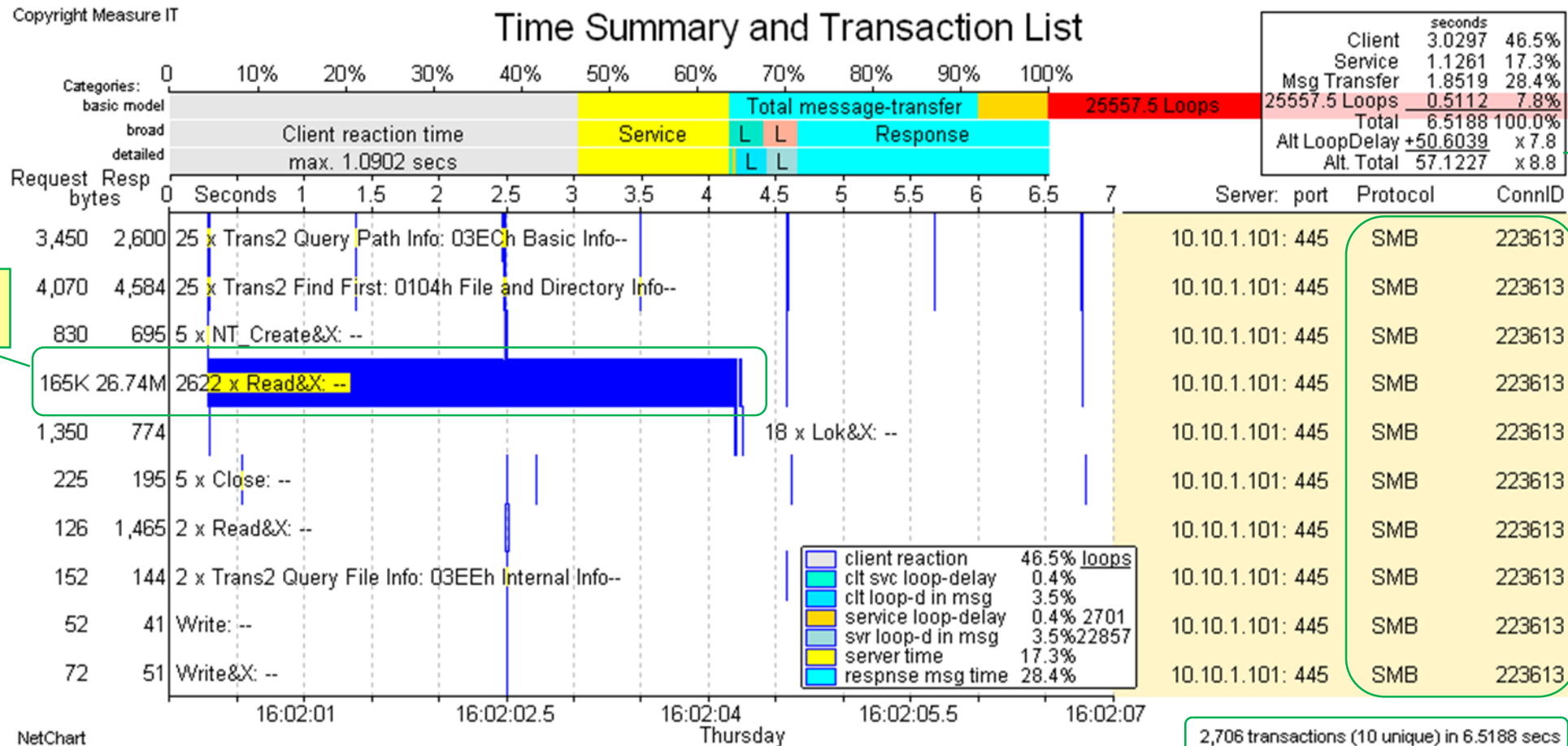
Transaction count is always here.

# SMB File Server

Here there is one TCP connection to a single SMB file server.

A total of 26.7 MB of data is downloaded via 2,622 Read requests (only 10 KB each) taking 4 seconds for overall completion.

This 6.5 second overall task would become 57 seconds without remediation.

Notice the large amount of grey "client" time, where the user PC has to handle all those SMB requests.  It also takes 1.8 seconds of network time to transfer 26.7 MB.  Could users live with this function taking nearly a minute?

Copyright Measure IT

## Time Summary and Transaction List

| | seconds | |
|---|---|---|
| Client | 3.0297 | 46.5% |
| Service | 1.1261 | 17.3% |
| Msg Transfer | 1.8519 | 28.4% |
| 25557.5 Loops | 0.5112 | 7.8% |
| Total | 6.5188 | 100.0% |
| Alt LoopDelay | +50.6039 | x 7.8 |
| Alt. Total | 57.1227 | x 8.8 |

Categories:

basic model — Total message-transfer — 25557.5 Loops

broad — Client reaction time — Service — L L — Response

detailed — max. 1.0902 secs — L L

Very serious!

Request Resp bytes / Seconds

| Request bytes | Resp | Transaction | Server: port | Protocol | ConnID |
|---|---|---|---|---|---|
| 3,450 | 2,600 | 25 x Trans2 Query Path Info: 03ECh Basic Info-- | 10.10.1.101: 445 | SMB | 223613 |
| 4,070 | 4,584 | 25 x Trans2 Find First: 0104h File and Directory Info-- | 10.10.1.101: 445 | SMB | 223613 |
| 830 | 695 | 5 x NT_Create&X: -- | 10.10.1.101: 445 | SMB | 223613 |
| 165K | 26.74M | 2622 x Read&X: -- | 10.10.1.101: 445 | SMB | 223613 |
| 1,350 | 774 | 18 x Lok&X: -- | 10.10.1.101: 445 | SMB | 223613 |
| 225 | 195 | 5 x Close: -- | 10.10.1.101: 445 | SMB | 223613 |
| 126 | 1,465 | 2 x Read&X: -- | 10.10.1.101: 445 | SMB | 223613 |
| 152 | 144 | 2 x Trans2 Query File Info: 03EEh Internal Info-- | 10.10.1.101: 445 | SMB | 223613 |
| 52 | 41 | Write: -- | 10.10.1.101: 445 | SMB | 223613 |
| 72 | 51 | Write&X: -- | 10.10.1.101: 445 | SMB | 223613 |

2622 SMB Reads, retrieving 26.7 MB.

One TCP connection. This time, SMB.

| Legend | % | |
|---|---|---|
| client reaction | 46.5% | loops |
| clt svc loop-delay | 0.4% | |
| clt loop-d in msg | 3.5% | |
| service loop-delay | 0.4% | 2701 |
| svr loop-d in msg | 3.5% | 22857 |
| server time | 17.3% | |
| respnse msg time | 28.4% | |

16:02:01   16:02:02.5   16:02:04   16:02:05.5   16:02:07

Thursday

NetChart

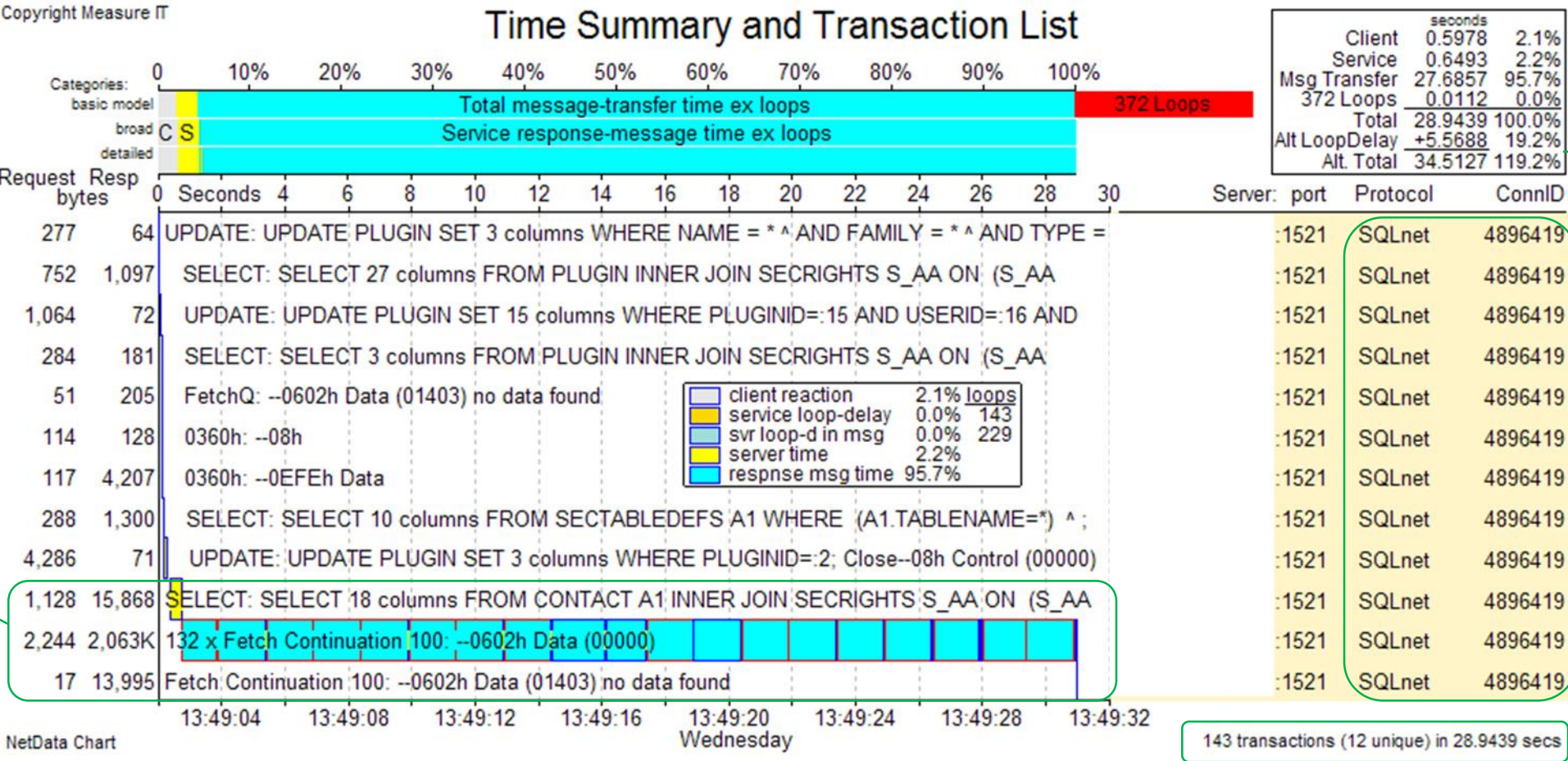2,706 transactions (10 unique) in 6.5188 secs

# Oracle Server – Network Constraint

A total of 2 MB of data is downloaded from an Oracle server via 134 requests (about 16 KB each). 95% of the overall time here was network data transfer time (blue). 28 seconds to transfer 2 MB is around 72 KB/s or 0.5 Mbps.
This 29 sec task would only become 35 secs after moving. However, network transfer speed is also an issue here, not just RTT.

Therefore, we also need to ensure that our connection to the cloud (or new DC) is provisioned with enough capacity.
It could still be worth a conversation with the application developer to retrieve larger chunks of data and reduce the 134 loops.



Copyright Measure IT

## Time Summary and Transaction List

| | seconds | |
|---|---|---|
| Client | 0.5978 | 2.1% |
| Service | 0.6493 | 2.2% |
| Msg Transfer | 27.6857 | 95.7% |
| 372 Loops | 0.0112 | 0.0% |
| Total | 28.9439 | 100.0% |
| Alt LoopDelay | +5.5688 | 19.2% |
| Alt. Total | 34.5127 | 119.2% |

An extra 5.5 secs may be OK.

| Request bytes | Resp | | Server: port | Protocol | ConnID |
|---|---|---|---|---|---|
| 277 | 64 | UPDATE: UPDATE PLUGIN SET 3 columns WHERE NAME = * ^ AND FAMILY = * ^ AND TYPE = | :1521 | SQLnet | 4896419 |
| 752 | 1,097 | SELECT: SELECT 27 columns FROM PLUGIN INNER JOIN SECRIGHTS S_AA ON (S_AA | :1521 | SQLnet | 4896419 |
| 1,064 | 72 | UPDATE: UPDATE PLUGIN SET 15 columns WHERE PLUGINID=:15 AND USERID=:16 AND | :1521 | SQLnet | 4896419 |
| 284 | 181 | SELECT: SELECT 3 columns FROM PLUGIN INNER JOIN SECRIGHTS S_AA ON (S_AA | :1521 | SQLnet | 4896419 |
| 51 | 205 | FetchQ: --0602h Data (01403) no data found | :1521 | SQLnet | 4896419 |
| 114 | 128 | 0360h: --08h | :1521 | SQLnet | 4896419 |
| 117 | 4,207 | 0360h: --0EFEh Data | :1521 | SQLnet | 4896419 |
| 288 | 1,300 | SELECT: SELECT 10 columns FROM SECTABLEDEFS A1 WHERE (A1.TABLENAME=*) ^ ; | :1521 | SQLnet | 4896419 |
| 4,286 | 71 | UPDATE: UPDATE PLUGIN SET 3 columns WHERE PLUGINID=:2; Close--08h Control (00000) | :1521 | SQLnet | 4896419 |
| 1,128 | 15,868 | SELECT: SELECT 18 columns FROM CONTACT A1 INNER JOIN SECRIGHTS S_AA ON (S_AA | :1521 | SQLnet | 4896419 |
| 2,244 | 2,063K | 132 x Fetch Continuation 100: --0602h Data (00000) | :1521 | SQLnet | 4896419 |
| 17 | 13,995 | Fetch Continuation 100: --0602h Data (01403) no data found | :1521 | SQLnet | 4896419 |

One TCP connection to an Oracle server.

| | | |
|---|---|---|
| client reaction | 2.1% | loops |
| service loop-delay | 0.0% | 143 |
| svr loop-d in msg | 0.0% | 229 |
| server time | 2.2% | |
| respnse msg time | 95.7% | |

13,400 rows of data (2 MB) were retrieved 100 at a time, resulting in 134 loops.

13:49:04  13:49:08  13:49:12  13:49:16  13:49:20  13:49:24  13:49:28  13:49:32
Wednesday

NetData Chart

143 transactions (12 unique) in 28.9439 secs

# Three Different Server Types

In this example, the overall function takes 1.3 seconds.
With 42 round trips, it would have just 0.6 seconds added to it (everything else being equal).

Users may not notice the extra 625 ms. However, given that server and network transfer times dominate here, we need to ensure that our new servers are capable and that our connection to the cloud (or new DC) is provisioned with enough capacity.

Colour coded breakup of times. There's a lot of server (yellow) and network time (blue) due to the large responses.

5 requests returning 182 KB.

An extra 0.6 secs may be OK.

8 TCP conns to 3 different servers.

# Four Different Protocol Types

In this example, the overall function takes 796 seconds. It has a mix of protocols, including 95 TN3270 requests to a mainframe. With 2442 round trips, it would have 35 seconds added to it. A 4.5% increase may be acceptable for this function.

Of interest here is that 90% of the time for this function is grey, meaning that it is internal processing (no network activity). The place to start looking for performance improvements is within the application code on the user's PC.



Copyright Measure IT

## Time Summary and Transaction List

**Large time of no network activity.**

**An extra 35 secs (4.5%) may be OK.**

**507 requests returning 10 MB.**

**181 conns to 4 different servers, incl mainframe.**

|  | seconds |  |
|---|---|---|
| Client | 722.3860 | 90.8% |
| Service | 51.7948 | 6.5% |
| Msg Transfer | 20.4462 | 2.6% |
| 2442.5 Loops | 1.0110 | 0.1% |
| Total | 795.6381 | 100.0% |
| Alt LoopDelay | +35.6265 | 4.5% |
| Alt. Total | 831.2645 | 104.5% |

| Request bytes | Resp | Transaction | Server: port | Protocol | ConnID |
|---|---|---|---|---|---|
| 54,249 | 160K | 542 x SELECT | :4100 | TDS-5 | 3627560 |
| 3,223 | 521 | 3 x BEGIN | :4100 | TDS-5 | 3627560 |
| 7,127 | 1,274 | 22 x INSERT | :4100 | TDS-5 | 3627560 |
| 48 | 67 | PA_RECORDS | :4100 | TDS-5 | 3627560 |
| 50,032 | 179K | 184 x lAnywhere | :2999 | Anywhere | 11 conns |
| 17,088 | 9,707K | 323 x lAnywhere | :2999 | Anywhere | 11 conns |
| 5,742 | 5,394 | 87 x Open Conn | ports | TowerIDM | 87 conns |
| 414 | 168 | 8 x Tower Technology | 2900 | TowerIDM | 5 conns |
| 3,359 | 23,772 | 95 x 3270 | n: 23 | Telnet | 4160775 |
| 1,748 | 3,516 | 85 x Tower Technology | 2900 | TowerIDM | 62 conns |
| 193 | 851 | 5 x ALLOCATENEXTID | :4100 | TDS-5 | 3627560 |
| 93 | 17 | DELETE: DELETE FROM Fo    WHERE CaseID = * and    CaseID = NULL- | :4100 | TDS-5 | 3627560 |
| 6,211 | 547 | 11 x UPDATE | | | 3627560 |
| 144 | 630 | 6 x USE | | | 3627560 |
| 146 | 88 | 2 x Stream (v5) | | | 3627560 |
| 2,206 | 296 | INSERT UPDATE INSERT (2) UPDATE: INSERT INTO Acti    (8 columns) VALUES (*)w | | | 3627560 |
| 9,576 | 1,512 | 6 x BEGIN SELECT UPDATE IF UPDATE IF SELECT IF ELSE UPDATE IF SELECT IF UPDATE | | | 3627560 |

: --result of stored proc

| | | |
|---|---|---|
| client reaction | 90.8% | loops |
| request msg time | 0.1% | |
| service loop-delay | 0.1% | 1294 |
| svr loop-d in msg | 0.0% | 1062 |
| connect req / ping | 0.0% | 87 |
| server time | 6.5% | |
| respnse msg time | 2.4% | |

Client reaction time max. 42.8201 secs

NetChart of 12:49 22/09/09

Tuesday 22/09/09

1,382 transactions in 795.6381 secs

# HTTP Protocol

This is an example showing HTTP transactions. With 257 round trips, the servers being 10 ms away would add 2.5 seconds to this 19 second set of transactions (a 13.4% increase). 87% of the overall time is "Server (yellow)" with 6.4% being "Msg Transfer (blue)" – which is time taken to deliver the packets of the large responses.

It is worth noting that this chart was produced with NetData Lite. Further, that this is just page 1 of a 7 page scrollable chart.

# How to Produce a Waterfall Chart

To produce a chart like the ones in this presentation, first produce a Packet Timing chart in NetData.  When loading data, packets must also be loaded in order to display the timing breakdown table at the top right of the Waterfall chart.

The "Waterfall" button will toggle the chart style. However, more options are available via the "Format" button.  The yellow highlights show the settings used to produce the chart on the previous slide.

First produce a "Packet Timing" chart – which will look something like this.

To get the time summary table, packets must be loaded.

The "Format" button in the top left of the Timing chart will show this window.

Alternatively, use the "Waterfall" button to directly toggle chart modes.

# Conclusions

Performing an analysis of application behaviour before servers are moved between data centres (or to the cloud) can provide an enormous benefit to the project team.

a) A data centre (or cloud) move comes with a long list of risks and unknowns.

b) The exact behaviour of an application is rarely fully known by the team that "owns" the application.

c) This type of analysis reduces unknowns and reduces risks.  More so if done early in the project timeframe.

d) Revealing the exact behaviour of applications can provide early notice of:
   - Which applications/servers can be moved with no (or little) expected performance impact.
   - Which applications/servers will have a performance impact – and how much?
   - Exactly why applications are predicted to have a large performance impact?
   - Workable remediation options that can be proposed and implemented.
   - Remediation options that can be implemented and tested well in advance of the physical server moves.

e) A similar traffic analysis can be used to determine which servers should be moved as groups.

If you'd like to discuss the possibility of performing this type of analysis in advance of your project (or if you've already moved and want to determine why you now have degraded performance), please get in touch via my contact details on the next page.

# Phil Storey

Phil@NetworkDetective.com.au

www.NetworkDetective.com.au

www.linkedin.com/in/philipstorey3/

@PhilStorey24

www.youtube.com/c/NetworkDetective

ask.wireshark.org:    @philst