# Zabbix Trending:

# Investigating Production Problems

By Erik Skytthe, DBC, Denmark



DBC: http://www.dbc.dk/english/about_dbc

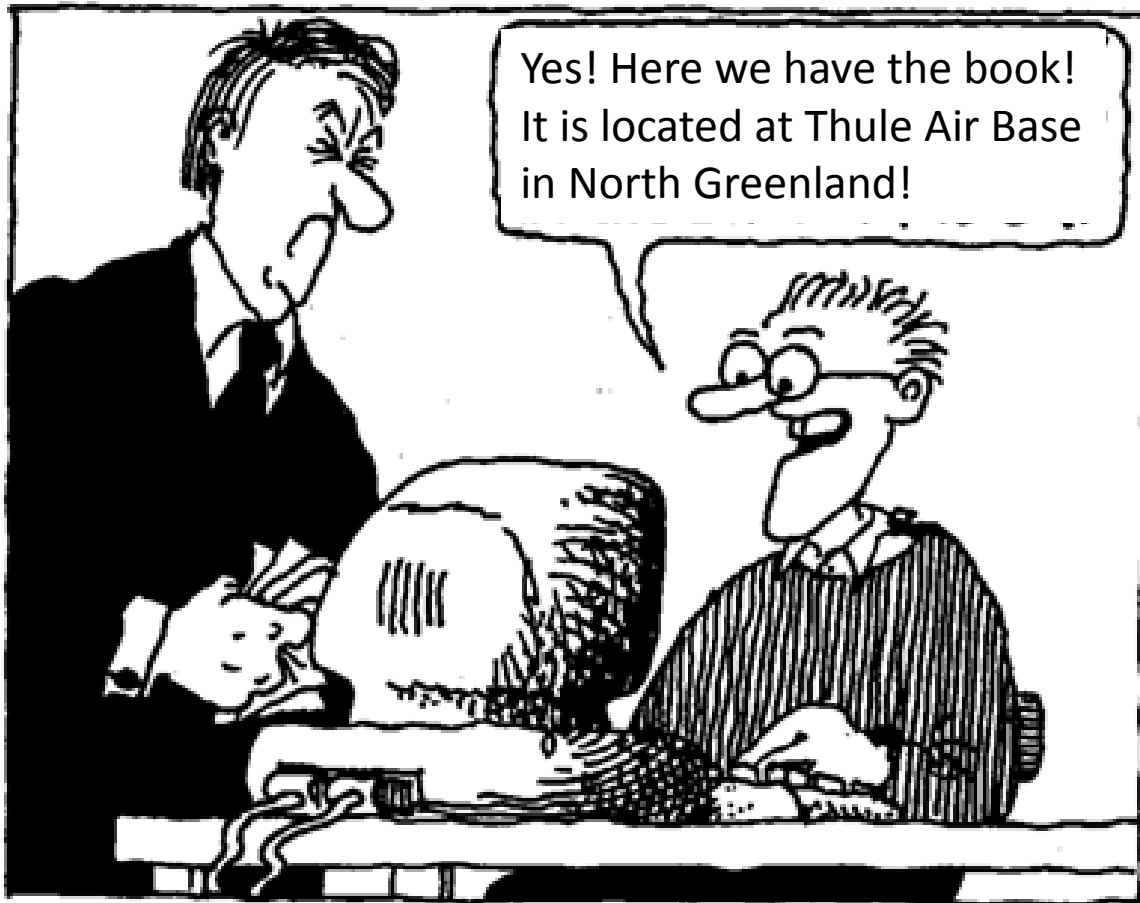Erik: http://www.linkedin.com/pub/erik-skytthe/20/a77/444

Zabbix IRC/Forum: eskytthe

**DBC**

DBC's main task is to develop and maintain the bibliographic and IT-infrastructure in **Danish libraries** by:

- Producing the **Danish National Bibliography** and user oriented cataloguing

- Developing and maintaining **Dan*Bib** as Union Catalog for public, educational and academic libraries as part of an automated ILL

- Developing and maintaining **library.dk** as the common access for citizens to Danish publications and Danish library holdings.

- DBC's IT systems are based on **open source** and service oriented architecture

- … E.g https://github.com/ding2/ding2

# A specialized world



And something completely different:

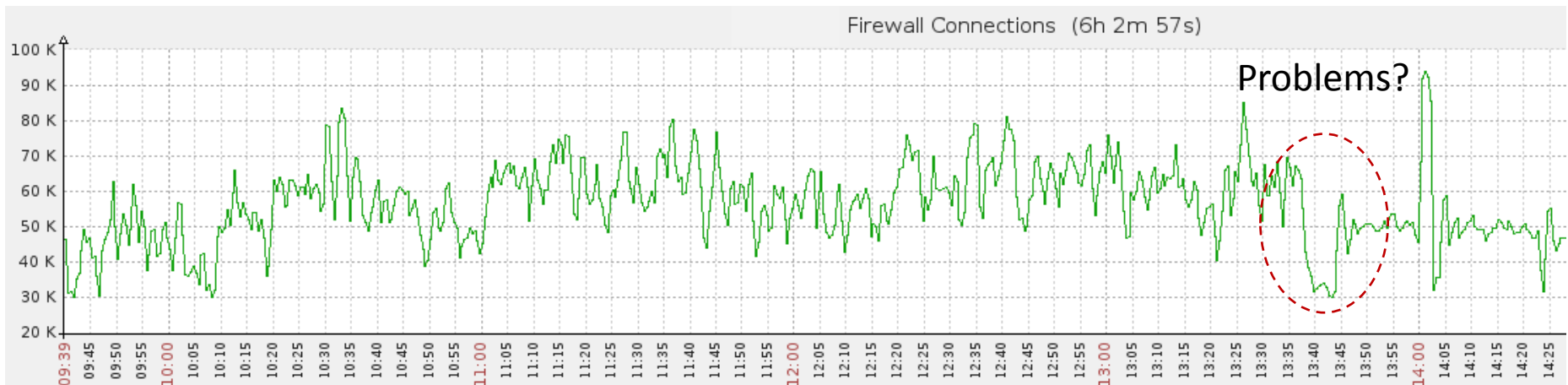*"If you have 50 search fields, a librarian will use them all!"*

# The Ghost case

- Victims:      Dead or unconscious services
- Detail:       Spontaneous "ghost" outages in multiple services
-                … and in non related services and systems !
- Suspect:     Network?
- Evidence:   Trading graphs, monitoring alarms

# The Ghost case

- **Central Firewall was found guilty!**

- **Root cause:** Setup of max connections was to low!

- Background: New Ajax based web services have over time increased the demand of connections allowed (and a bug)

- Tools: Central syslog, Hobbit and Zabbix



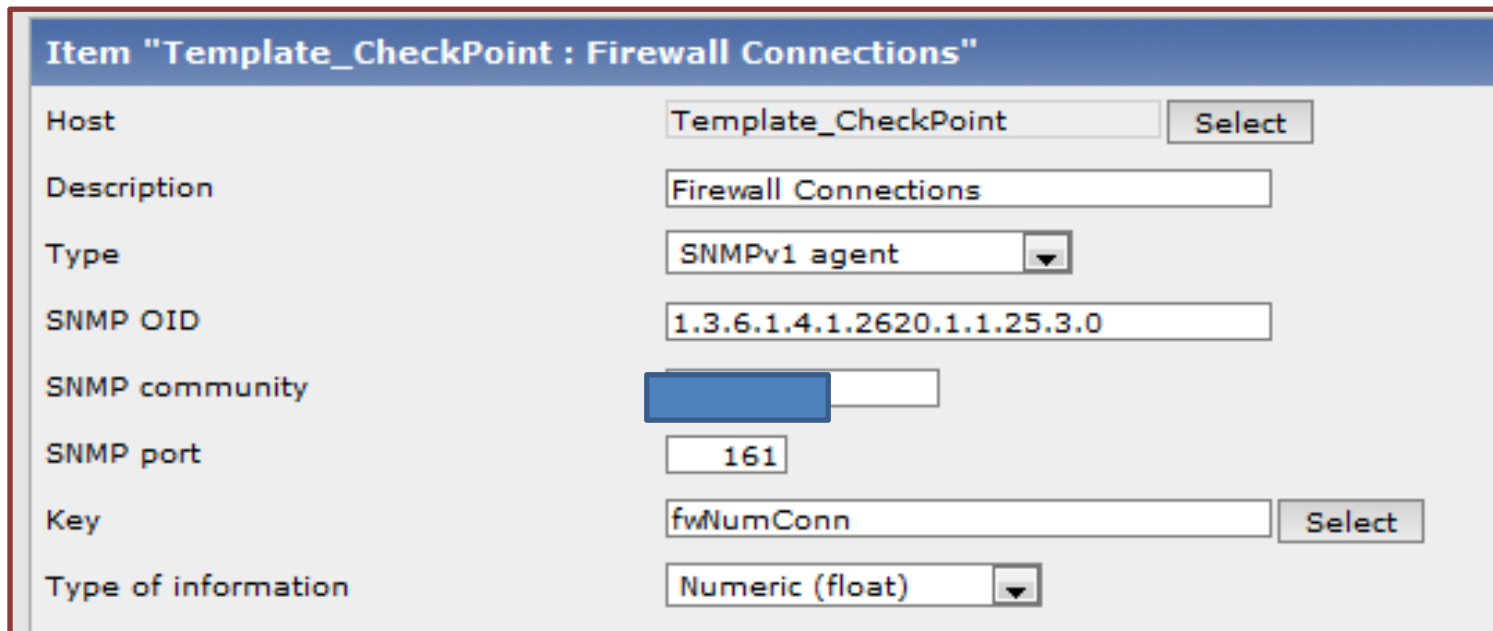Firewall Connections (6h 2m 57s)

Problems?

# The Ghost case

- Made a trigger on max connections
- Had to increase the setting several times
- Periodic peeks because of testing from development team

$ snmpwalk -v1 -c community aaa-fw-1 1.3.6.1.4.1.2620.1.1.25.3

SNMPv2-SMI::enterprises.2620.1.1.25.3.0 = INTEGER: 49776

**Item "Template_CheckPoint : Firewall Connections"**

| | |
|---|---|
| Host | Template_CheckPoint [Select] |
| Description | Firewall Connections |
| Type | SNMPv1 agent |
| SNMP OID | 1.3.6.1.4.1.2620.1.1.25.3.0 |
| SNMP community | |
| SNMP port | 161 |
| Key | fwNumConn [Select] |
| Type of information | Numeric (float) |

# The Ghost case … continued

- Started to setup central syslog …
- This message dropped in from local iptable firewall on host1
- Needed to increase the iptable connection tracking table on the host

Aug  2 13:09:47 host1 kernel: [8453153.139095] nf_conntrack: table full, dropping packet.

Aug  2 13:09:47 host1 kernel: [8453153.141819] nf_conntrack: table full, dropping packet.

# The Ghost case … continued

- Have "integrated" syslog files in Zabbix screens
- URL's to files on central syslog server
- syslog-today,  syslog-yesterday, syslog-weekly
- Will filter syslog on central syslog server and send critical lines to Zabbix server
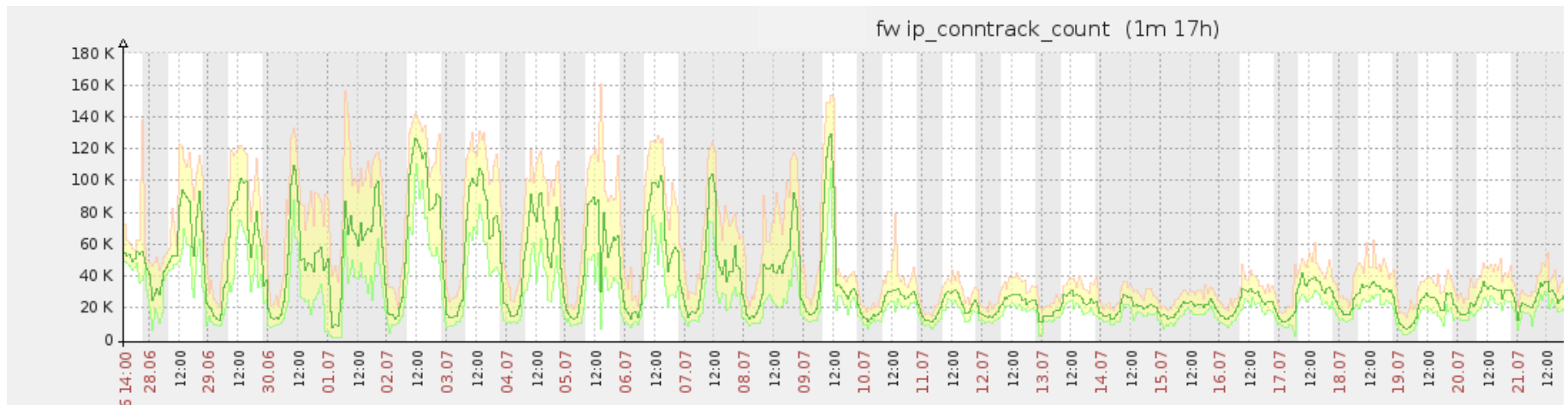- Rsyslog has some nice filter options

**Syslog**

⌄ Filter ⌄

```
daemon.err<27>: Aug 31 13:15:53          ntpd[21385]: couldn't unlink /var/log/ntpstats/peerstats: Permission denied
daemon.err<27>: Aug 31 13:15:53          ntpd[21385]: couldn't unlink /var/log/ntpstats/peerstats: Permission denied
daemon.err<27>: Aug 31 13:15:53          ntpd[21385]: can't open /var/log/ntpstats/peerstats.20120831: Permission denied
daemon.err<27>: Aug 31 13:15:53          ntpd[21385]: can't open /var/log/ntpstats/peerstats.20120831: Permission denied
kern.err<3>: Aug 31 13:52:38      kernel: [    2.226018] sd 2:0:0:0: [sda] Assuming drive cache: write through
kern.err<3>: Aug 31 13:52:38      kernel: [    2.226018] sd 2:0:0:0: [sda] Assuming drive cache: write through
kern.err<3>: Aug 31 13:52:38      kernel: [    2.226333] sd 2:0:0:0: [sda] Assuming drive cache: write through
kern.err<3>: Aug 31 13:52:38      kernel: [    2.226333] sd 2:0:0:0: [sda] Assuming drive cache: write through
kern.err<3>: Aug 31 13:52:38      kernel: [    2.231033] sd 2:0:0:0: [sda] Assuming drive cache: write through
kern.err<3>: Aug 31 13:52:38      kernel: [    2.231033] sd 2:0:0:0: [sda] Assuming drive cache: write through
kern.err<3>: Aug 31 13:52:38      kernel: [    3.679140] ACPI: I/O resource piix4_smbus [0x1040-0x1047] conflicts with
kern.err<3>: Aug 31 13:52:38      kernel: [    3.679140] ACPI: I/O resource piix4_smbus [0x1040-0x1047] conflicts with
```

# The Ghost case … continued … cont.

- Later having major problems on host after linux upgrade
- Have to optimize some TCP kernel parameters

- Bug: TCP sockets were not reused in applications (old bug)
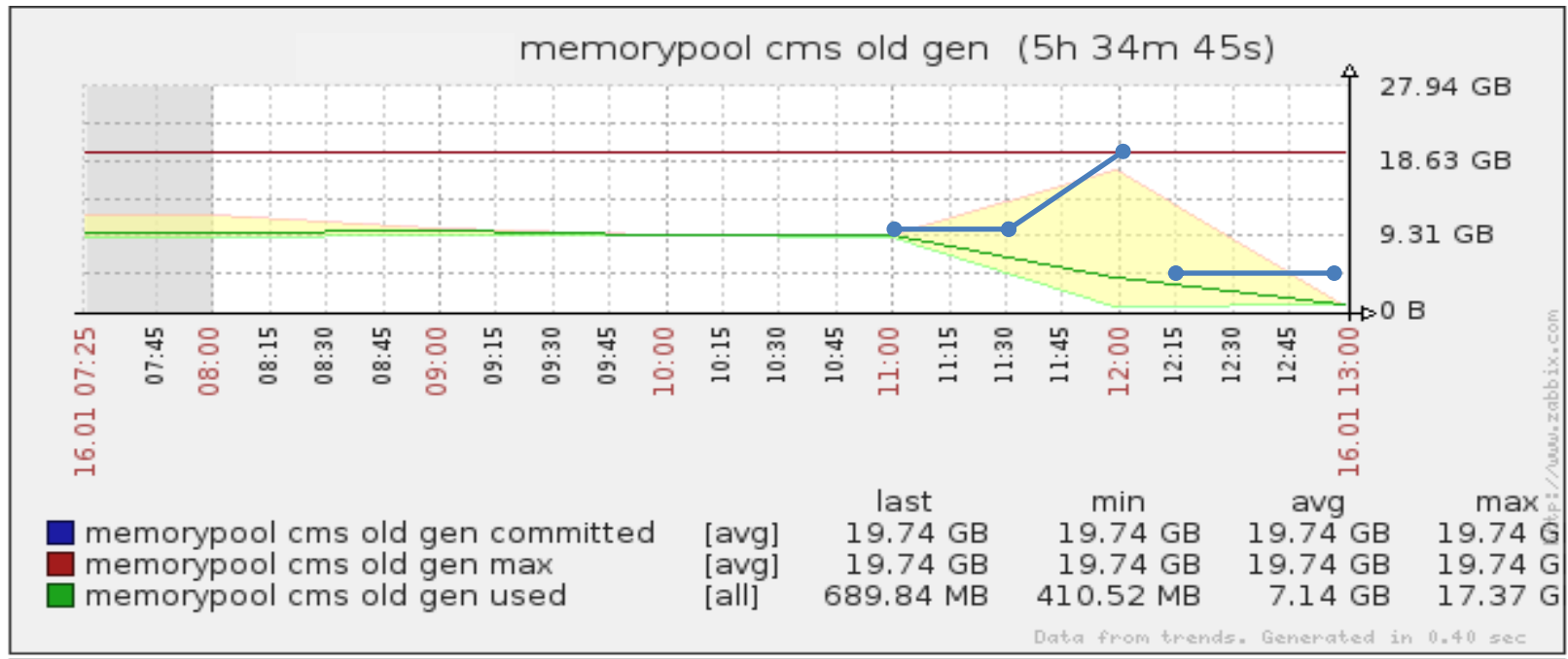- Effect of fixed bug can be seen in graph

# The Java case

- Victim:        Dead java application
- Suspect:       ??
- Evidence:    BIG BIG java log files
- Starting monitoring of java (Tomcat etc)
- Tool:           Zapcat JMX Zabbix Bridge



© Universal Press Syndicate

# The Java case



**Observations and actions:**

- Short before crime: "cms old gen" went unusually high

- Making trigger on 70% of cms old gen used - alert/early warning

- Using graph to point out time for problem start – problem reporting to development team

# The Java case

- **Internal search update routine was found guilty!**
- **Route Cause:** Did not handle timeout from data store system ok
- – start accumulated search over and over again!
- Still using graphs to point out start time for problems

- **Extension of java monitoring and templates**
- Solr and other production data
- Own developed attributes used in new templates

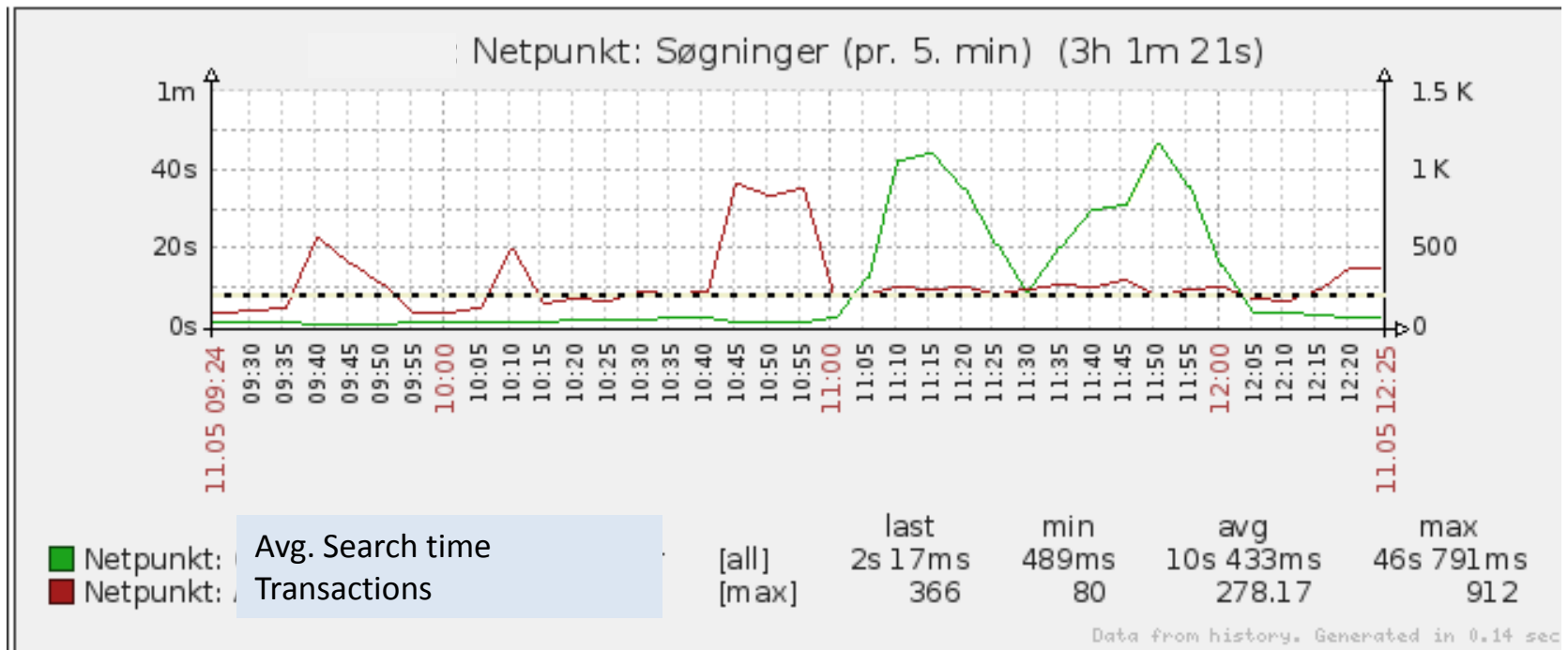| Description | Triggers | Key ↕ |
|---|---|---|
| QueryCache - Hitratio | | jmx[solr/:type=queryResultCache,id=org.apache.solr.search.LRUCache][hitratio] |
| QueryCache - Hits | | jmx[solr/:type=queryResultCache,id=org.apache.solr.search.LRUCache][hits] |
| QueryCache - Inserts | | jmx[solr/:type=queryResultCache,id=org.apache.solr.search.LRUCache][inserts] |
| QueryCache - Lookups | | jmx[solr/:type=queryResultCache,id=org.apache.solr.search.LRUCache][lookups] |
| QueryCache - Size | | jmx[solr/:type=queryResultCache,id=org.apache.solr.search.LRUCache][size] |
| SolrSearch - avgTimePerRequest | | jmx[solr/:type=search,id=org.apache.solr.handler.component.SearchHandler][avgTimePerRequest] |
| SolrSearch - Errors | | jmx[solr/:type=search,id=org.apache.solr.handler.component.SearchHandler][errors] |
| SolrSearch - Requests | | jmx[solr/:type=search,id=org.apache.solr.handler.component.SearchHandler][requests] |
| SolrSearcher - MaxDoc | | jmx[solr/:type=searcher,id=org.apache.solr.search.SolrIndexSearcher][maxDoc] |
| SolrSearcher - NumDocs | | jmx[solr/:type=searcher,id=org.apache.solr.search.SolrIndexSearcher][numDocs] |
| SolrSearcher - openedAt | | jmx[solr/:type=searcher,id=org.apache.solr.search.SolrIndexSearcher][openedAt] |

# Getting Production data
## (Enter the crime scene)

- Discovered that many DBC services were using nearly the same log format

- Python script to parse the logfiles

- Running from cron locally now.

- Moving to central solution, with syslog etc.

- Parse logfile from EOF, because of BIG logfiles ( back to *–B sec.* )

- Counting transactions, errors and avg. transactions time

- Sending items to Zabbix server with zabbix_sender

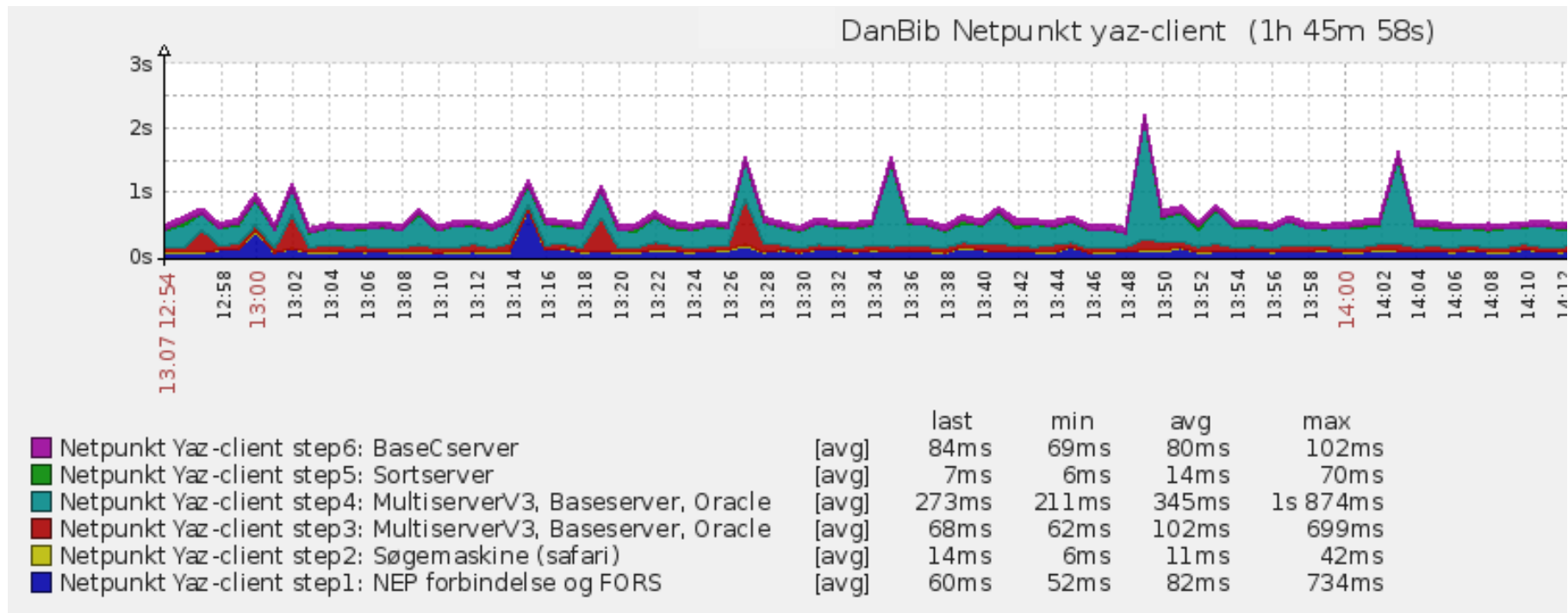- *parsetracelog.py -r --tag=Sog --sysname=systemX -B 300 -z file.log*

# Getting Production data

- Big moment when first real life production data was showed up in the Zabbix GUI !!

# Getting Production data

- Discovered test utility at a workshop

- Also running from cron – every 1 min.

- Test all main backend subsystems and return times for steps

- *yaz-client –f inputfile | grep Elapsed ….*



DanBib Netpunkt yaz-client (1h 45m 58s)

| | | | last | min | avg | max |
|---|---|---|---|---|---|---|
| ■ Netpunkt Yaz-client step6: BaseCserver | | [avg] | 84ms | 69ms | 80ms | 102ms |
| ■ Netpunkt Yaz-client step5: Sortserver | | [avg] | 7ms | 6ms | 14ms | 70ms |
| ■ Netpunkt Yaz-client step4: MultiserverV3, Baseserver, Oracle | | [avg] | 273ms | 211ms | 345ms | 1s 874ms |
| ■ Netpunkt Yaz-client step3: MultiserverV3, Baseserver, Oracle | | [avg] | 68ms | 62ms | 102ms | 699ms |
| ■ Netpunkt Yaz-client step2: Søgemaskine (safari) | | [avg] | 14ms | 6ms | 11ms | 42ms |
| ■ Netpunkt Yaz-client step1: NEP forbindelse og FORS | | [avg] | 60ms | 52ms | 82ms | 734ms |

# Looking for traces

- **Talk with people!**

- Get different people's views and opinions

- ... the programmer, the customer, the super user, the application people, the database admins, the end user support team ...

- Let them show you how they work with the systems

- Get overview drawings etc.

- Make logs / notes about what you see and find

- **Get data!**

- From logs

- From scripts

- From polling

- From Zabbix database – SQL , event etc.

Take care !!
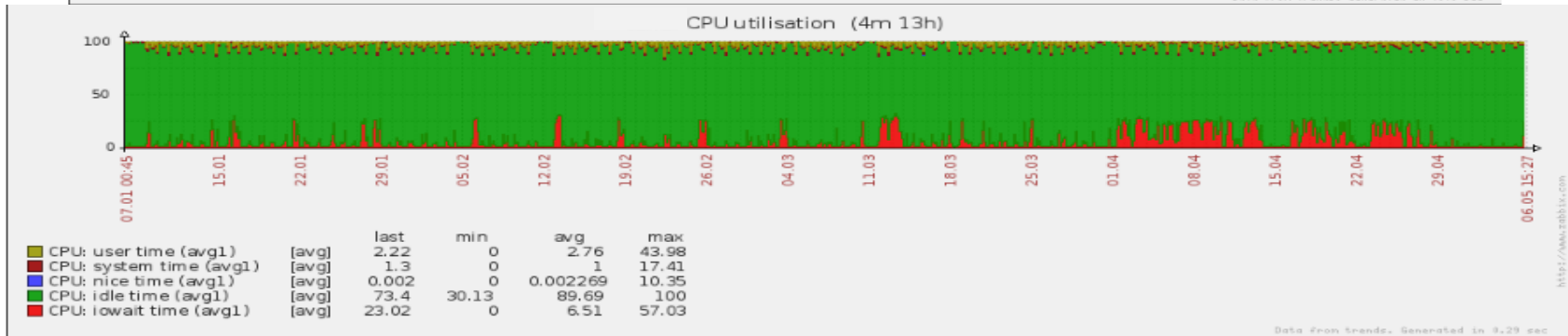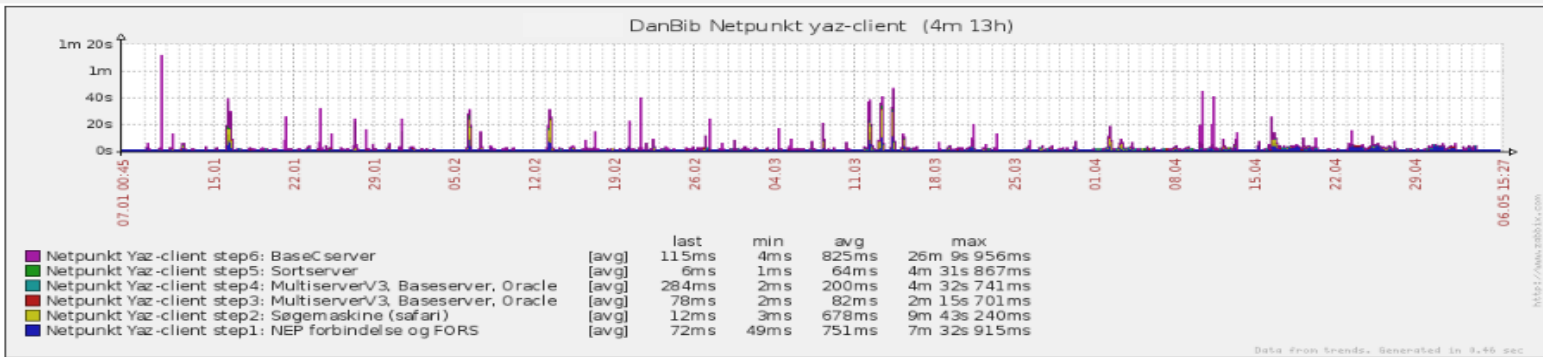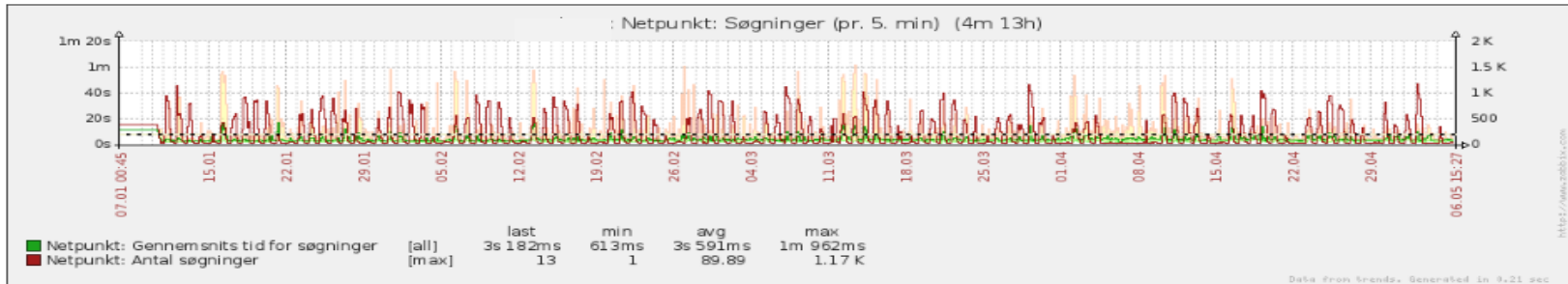- With avg. in graphs
- Your statistics

# Looking for traces

- My main work board / Tool is **Zabbix Screens**!
- I combine different graphs
- Move around with them, play with them
- Zoom in and out
- Try to make overviews
- Make several different screens, looking at different aspects
- Go down in details at other times

- **In this way I:**
- Try to analyse data and looking for trends
- Try to find patterns and correlate graphs
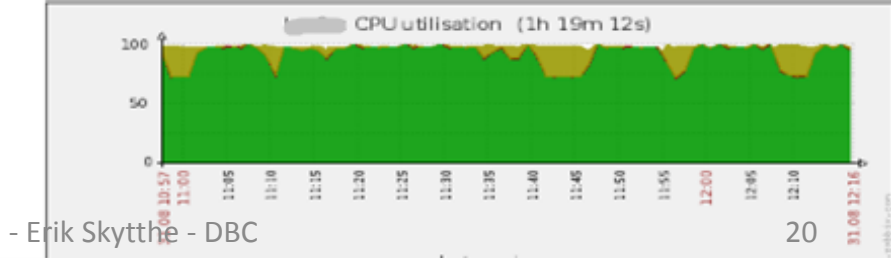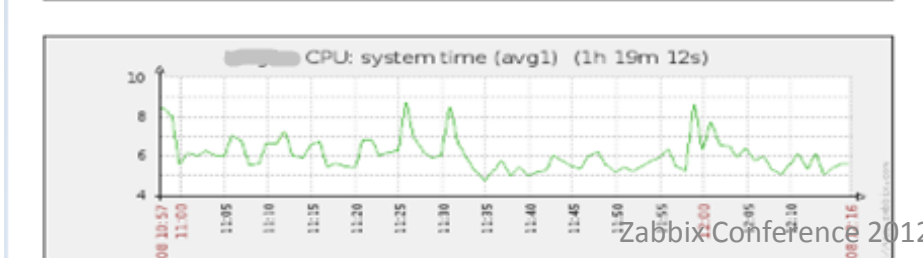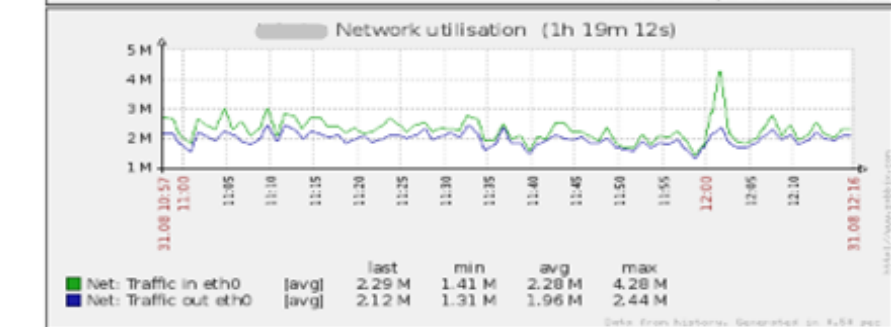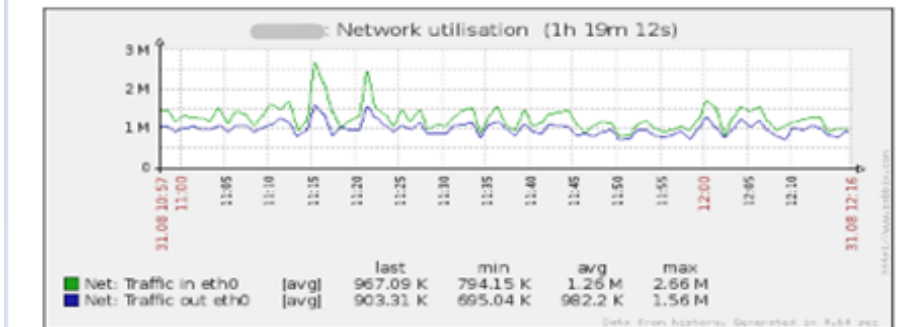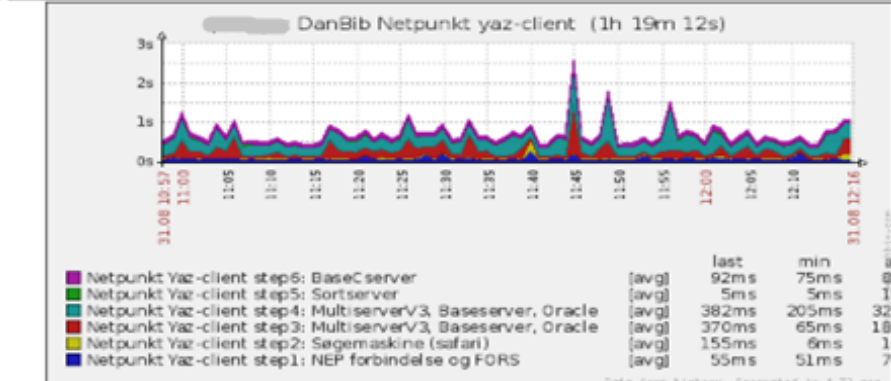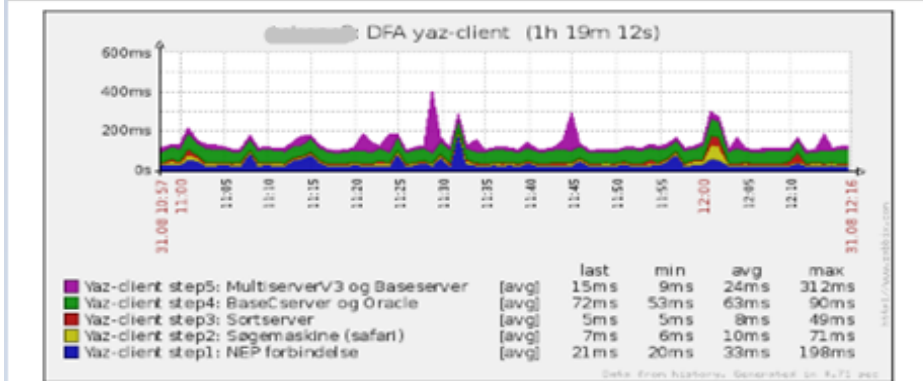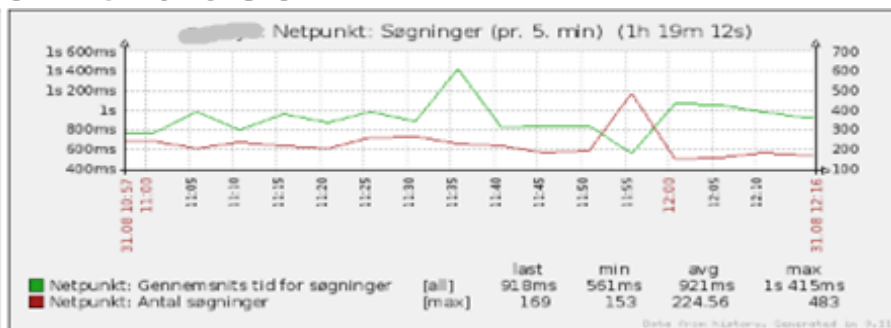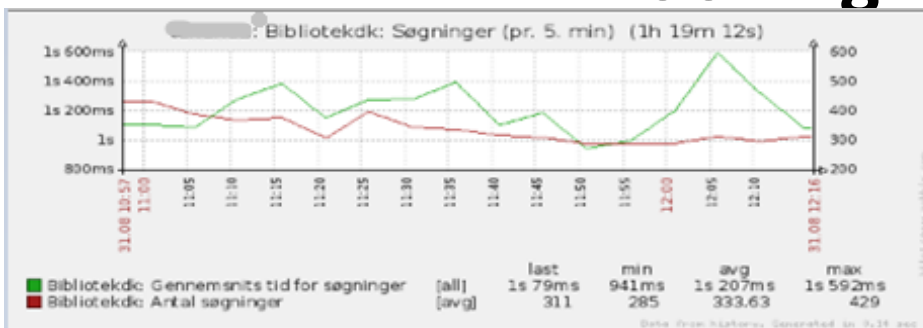- Try to isolate the problems

# Looking for traces

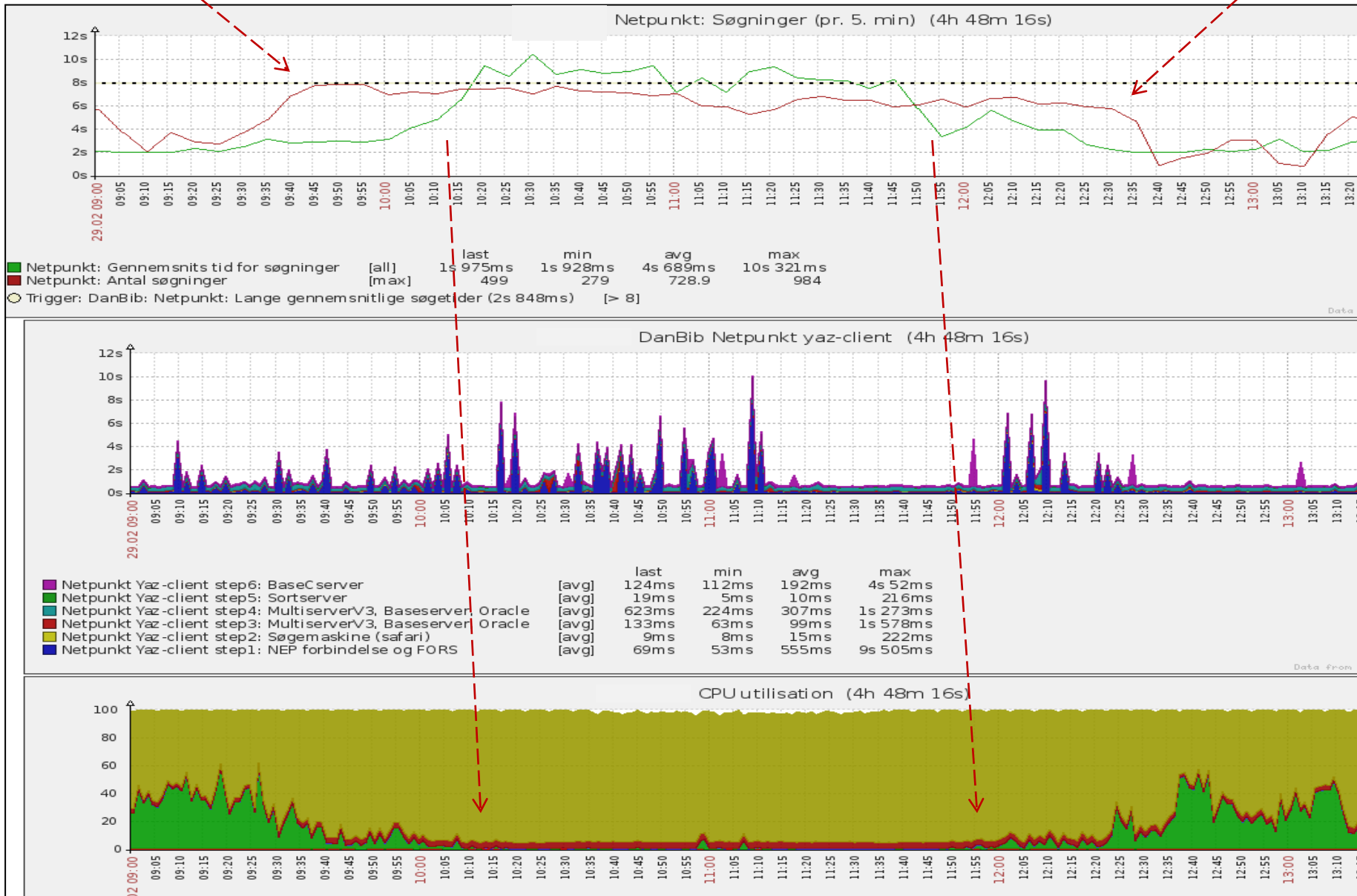# Looking for traces

# Looking for traces

# The "Wednesday monster" case

- Victims:      Users of application.

- **Details:**

- On most Wednesdays, major performance problems, or even service breakdown are seen. Sometimes  it is Thursday instead.

- Suspects:    Backend, frontend systems?

- Evidence:    Trading graphs, monitoring alarms
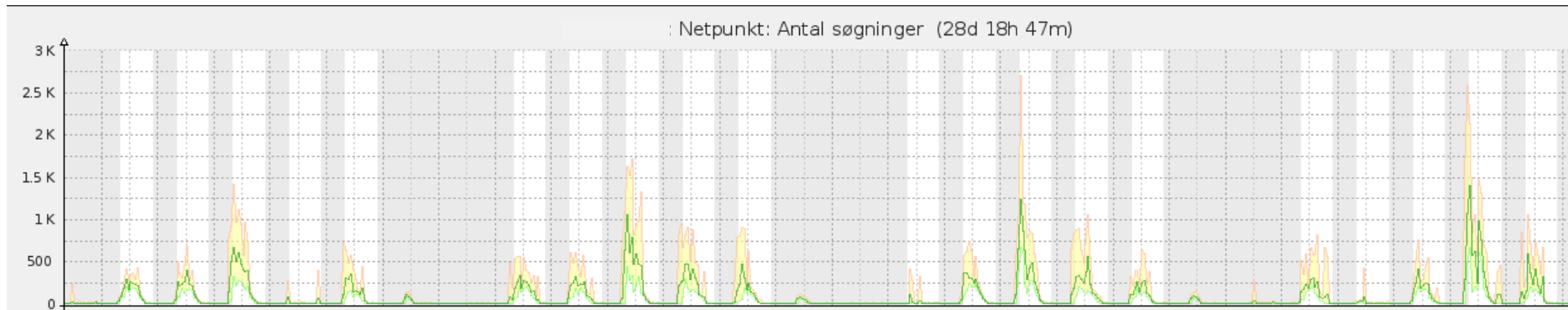
# The "Wednesday monster" case



Netpunkt: Søgninger (pr. 5. min)  (4h 48m 16s)

| | | last | min | avg | max |
|---|---|---|---|---|---|
| ▬ Netpunkt: Gennemsnits tid for søgninger | [all] | 1s 975ms | 1s 928ms | 4s 689ms | 10s 321ms |
| ▬ Netpunkt: Antal søgninger | [max] | 499 | 279 | 728.9 | 984 |
| ○ Trigger: DanBib: Netpunkt: Lange gennemsnitlige søgetider (2s 848ms) | [> 8] | | | | |

DanBib Netpunkt yaz-client  (4h 48m 16s)

| | | last | min | avg | max |
|---|---|---|---|---|---|
| ▪ Netpunkt Yaz-client step6: BaseC server | [avg] | 124ms | 112ms | 192ms | 4s 52ms |
| ▪ Netpunkt Yaz-client step5: Sortserver | [avg] | 19ms | 5ms | 10ms | 216ms |
| ▪ Netpunkt Yaz-client step4: MultiserverV3, Baseserver Oracle | [avg] | 623ms | 224ms | 307ms | 1s 273ms |
| ▪ Netpunkt Yaz-client step3: MultiserverV3, Baseserver Oracle | [avg] | 133ms | 63ms | 99ms | 1s 578ms |
| ▪ Netpunkt Yaz-client step2: Søgemaskine (safari) | [avg] | 9ms | 8ms | 15ms | 222ms |
| ▪ Netpunkt Yaz-client step1: NEP forbindelse og FORS | [avg] | 69ms | 53ms | 555ms | 9s 505ms |

CPU utilisation  (4h 48m 16s)

# The "Wednesday monster" case

- **Graph conclusions:**
- Major increase in amount of transactions for period
- Correlation between "bad" avg. search time and user CPU time at frontend host
- User CPU time consumption is HEAVY …
- … But is it because of slow backend systems or ???
- … does not seem so – the middle graph is showing backend subsystems. No clear correlation there.
- (Well, middle graph is polled test values, but should be more clear if correlations – Upper graph have higher weight)
- **Other conclusions:**
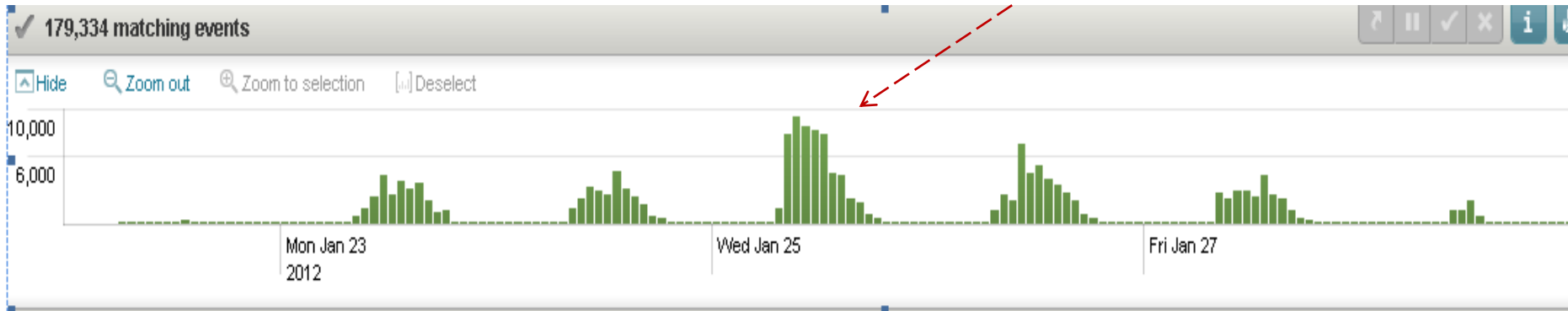- Twin system does not show any problems … so no general infrastructure problems … (problem isolation)

# The "Wednesday monster" case

- Zoomed out 4 weeks – only looking at transactions
- Seems like Wednesday hang out

# The "Wednesday monster" case

# The "Wednesday monster" case

- Did a deeper analysis of the transaction file
- Using "Splunk" - a not so free log analyzer tool (other can do)

- Looking into Wednesday …
- … Top10 source ip ~ more then 40 % of transactions
- Looking at one Top10 source ip …
- … Most transactions on Wednesday!

# The "Wednesday monster" case

- **After some talk with application people:**

- Every Tuesday night there is a update of the central database

- Some local library systems do a data harvest the day after, via DBC web frontend system (to update their local databases)

- Yes Wednesday !!!

- … and sometime Thursday instead

- … if DBC central update have been delayed !


- **Case closed:**

- Customers were found guilty ☺ or ?

- No **root cause was old slow server hardware** at frontend host

- … problem solved after hardware upgrade

# The "Search monster" case

- Victims:      Users of application.
- **Details:**
- General performance problems in intervals of 5 – 30 min
- Some days or even weeks there is no problems at all.

- Suspects:    Backend or frontend systems or?
- Evidence:   Trading graphs, monitoring alarms

# The "Search monster" case



: Netpunkt: Søgninger (pr. 5. min)  (5d 8h 58m)

| | | last | min | avg | max |
|---|---|---|---|---|---|
| ■ Netpunkt: Gennemsnits tid for søgninger | [all] | 3s 140ms | 953ms | 5s 214ms | 1m 962ms |
| ■ Netpunkt: Antal søgninger | [max] | 2 | 1 | 117.64 | 1 K |



: DanBib Netpunkt yaz-client  (5d 8h 58m)

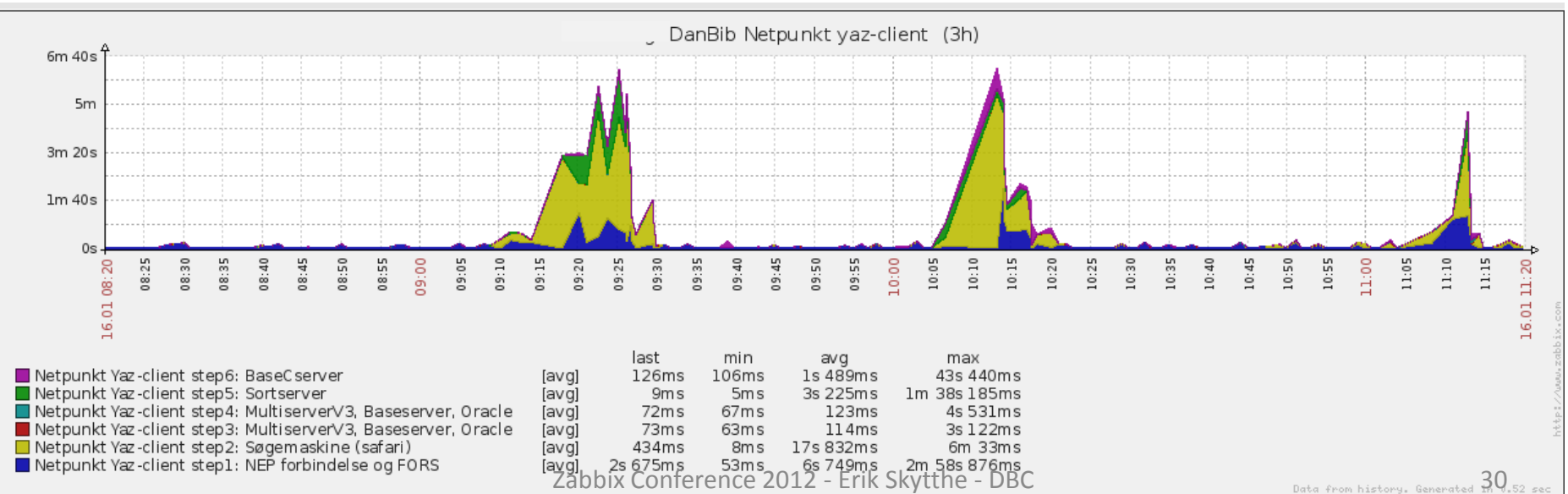| | | last | min | avg | max |
|---|---|---|---|---|---|
| ■ Netpunkt Yaz-client step6: BaseCserver | [avg] | 109ms | 98ms | 811ms | 7m 22s 34ms |
| ■ Netpunkt Yaz-client step5: Sortserver | [avg] | 10ms | 5ms | 522ms | 2m 20s 874ms |
| ■ Netpunkt Yaz-client step4: MultiserverV3, Baseserver, Oracle | [avg] | 221ms | 2ms | 255ms | 2s 930ms |
| ■ Netpunkt Yaz-client step3: MultiserverV3, Baseserver, Oracle | [avg] | 66ms | 2ms | 99ms | 2m 15s 701ms |
| ■ Netpunkt Yaz-client step2: Søgemaskine (safari) | [avg] | 12ms | 7ms | 4s 618ms | 9m 43s 240ms |
| ■ Netpunkt Yaz-client step1: NEP forbindelse og FORS | [avg] | 67ms | 49ms | 2s 402ms | 5m 5s 785ms |

# The "Search monster" case



Netpunkt: Søgninger (pr. 5. min)  (3h)

| | last | min | avg | max |
|---|---|---|---|---|
| Netpunkt: Gennemsnits tid for søgninger  [avg] | 28s 686ms | 2s 977ms | 15s 389ms | 56s 161ms |
| Netpunkt: Antal søgninger  [avg] | 319 | 33 | 155.5 | 404 |

Data from history. Generated in 0.40 s



DanBib Netpunkt yaz-client  (3h)

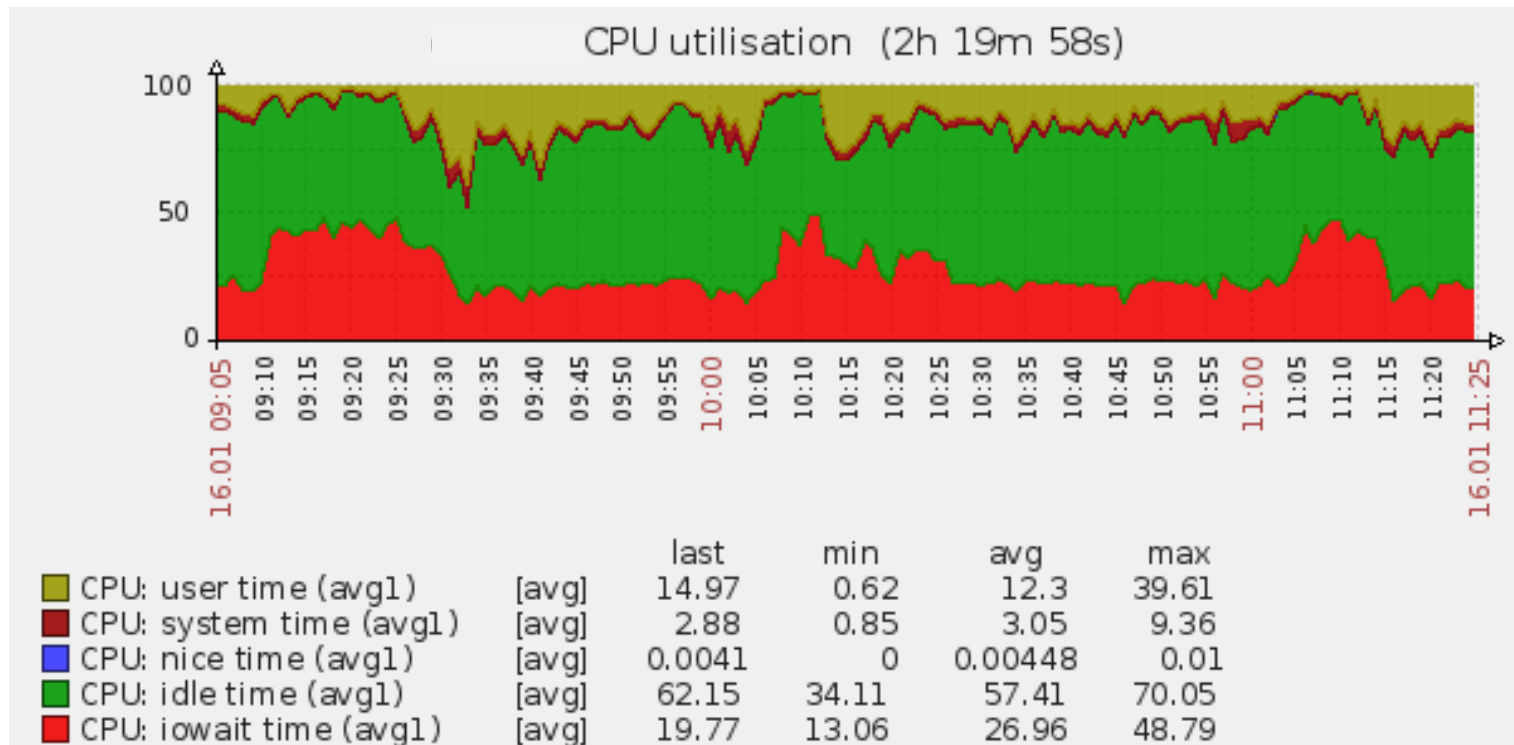| | last | min | avg | max |
|---|---|---|---|---|
| Netpunkt Yaz-client step6: BaseCserver  [avg] | 126ms | 106ms | 1s 489ms | 43s 440ms |
| Netpunkt Yaz-client step5: Sortserver  [avg] | 9ms | 5ms | 3s 225ms | 1m 38s 185ms |
| Netpunkt Yaz-client step4: MultiserverV3, Baseserver, Oracle  [avg] | 72ms | 67ms | 123ms | 4s 531ms |
| Netpunkt Yaz-client step3: MultiserverV3, Baseserver, Oracle  [avg] | 73ms | 63ms | 114ms | 3s 122ms |
| Netpunkt Yaz-client step2: Søgemaskine (safari)  [avg] | 434ms | 8ms | 17s 832ms | 6m 33ms |
| Netpunkt Yaz-client step1: NEP forbindelse og FORS  [avg] | 2s 675ms | 53ms | 6s 749ms | 2m 58s 876ms |

Data from history. Generated in 0.52 sec

# The "Search monster" case

- **Graph week overview:**
- Clear correlation
- Search subsystem hangout
- Small BaseCserver subsystem peak Monday
- Not as bad as it looked like – not whole day … do a zoom in

- **Graph hour:**
- Search subsystem peeks of 10 – 15 min
- Clear correlation with bad performance (avg. transaction)

# The "Search monster" case

- **Update:**
- **Commit:**



CPU utilisation (2h 19m 58s)

| | | last | min | avg | max |
|---|---|---|---|---|---|
| CPU: user time (avg1) | [avg] | 14.97 | 0.62 | 12.3 | 39.61 |
| CPU: system time (avg1) | [avg] | 2.88 | 0.85 | 3.05 | 9.36 |
| CPU: nice time (avg1) | [avg] | 0.0041 | 0 | 0.00448 | 0.01 |
| CPU: idle time (avg1) | [avg] | 62.15 | 34.11 | 57.41 | 70.05 |
| CPU: iowait time (avg1) | [avg] | 19.77 | 13.06 | 26.96 | 48.79 |

# The "Search monster" case

- **Added CPU graph from a search server**
- I/O wait peeks correlate with application performance problems
- Application colleague saw the correlation immediately:

- **Explanation / Root cause:**
- Blue lines is updates of search registers in search subsystem
- Red lines is commits of these updates!
- … and these commits are HEAVY … at disk I/O
- Updates are running in periods (sometimes with weeks in between)
- New search engine did also make more load on servers

- **Case closed:**
- **Search server hardware was found guilty!**
- … problem solved after hardware upgrade
- … AND heavy Disk system upgrade (read: faaaast disks)

# The "Monday monster" case

- Victims: Users of application.



© Universal Press Syndicate

- **Details:**

- A performance problem. Even that it is called the Monday problem, it do not have special connections with Mondays. It is seen at all days and all times of days. Even at weekends. It is not seen that often … 1 – 3 times a week … but can block activity for some time

- Suspects: Backend systems, database system, heavy batch jobs?

- Evidence: Historical data, trading graphs, monitoring alarms

The "Monday monster" case

# The "Monday monster" case

- **BaseCserver (backend subsystem)**

- Did early see correlation with baseCserver tests

- However Logs and statistics for baseCserver did not indicate problems

- Did not have logs and statistics for systems in between frontend and baseCserver

- **Analyses of frontend log files:**

- Only transaction times from 0 – 20 sec and > 1 min (timeout) ??

- 25 % of transactions was timeouts  ?

# The "Monday monster" case

- **Database server**

- Did later find correlation with systemtime on DB server

- But response time statistics etc. from DB server was fine

- Even specialized DB monitor tool did not show problems


- Further search for correlations did not succeed

- Eyes was turned onto batch jobs and other systems doing heavy database searches etc.

# The "Monday monster" case

- **Other systems – the long search**

- Did get list of "heavy" systems

- Did go for the logs again

- Different log formats, but a matter of scripting

- Did not want to wait for things (problems) to happen this time

- So pulled out historical data of system logs and transferred it into Zabbix

# The "Monday monster" case

- **Other systems – the long search … continued**
- Struggle with how to visualize data from log files
  - Time span
  - Records / transactions processed
- Miss bar graphs
- Tried to "simulate" bar graphs

- **Did not find any correlation with heavy systems**
- … at least not the systems I did investigate

- Did however find jobs that should not run in day time …

**The "Monday monster" case**

# The "Monday monster" case

- **Zabbix_sender**
- **$ awk -f adhl.awk adhl-log > outputfile**
- **Outputfile:** **(*simulation of bar graphs)***
- adhl danbib.adhlimport.time 1343180920  0  *(edge start)*
- adhl danbib.adhlimport.time 1343180930 3783
- adhl danbib.adhlimport.time 1343184704 3783
- adhl danbib.adhlimport.time 1343184714  0  *(edge end)*
- adhl danbib.adhlimport.time 1343184725  0
- adhl danbib.adhlimport.time 1343184735 877
- adhl danbib.adhlimport.time 1343185603 877
- adhl danbib.adhlimport.time 1343185613  0
-
- **$ Zabbix_sender -z Zabbix -i outputfile -T**
- Info from server: "Processed 250 Failed 0 Total 250 Seconds spent 0.001821"
- …
- Info from server: "Processed 250 Failed 0 Total 250 Seconds spent 0.001739"
- Info from server: "Processed 80 Failed 0 Total 80 Seconds spent 0.000577"
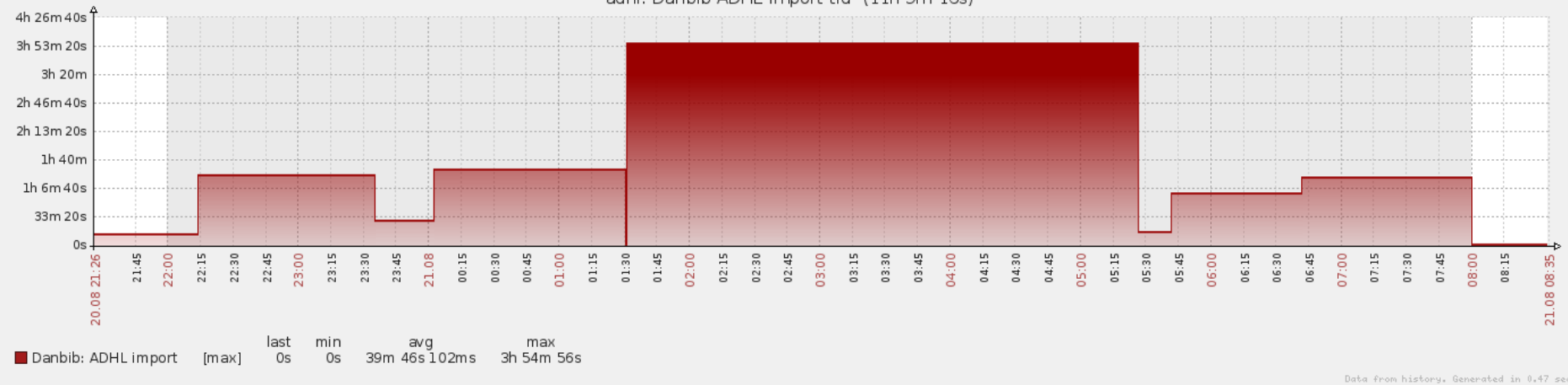- sent: 1580; skipped: 0; total: 1580

# The "Monday monster" case

# The "Monday monster" case

# The "Monday monster" case

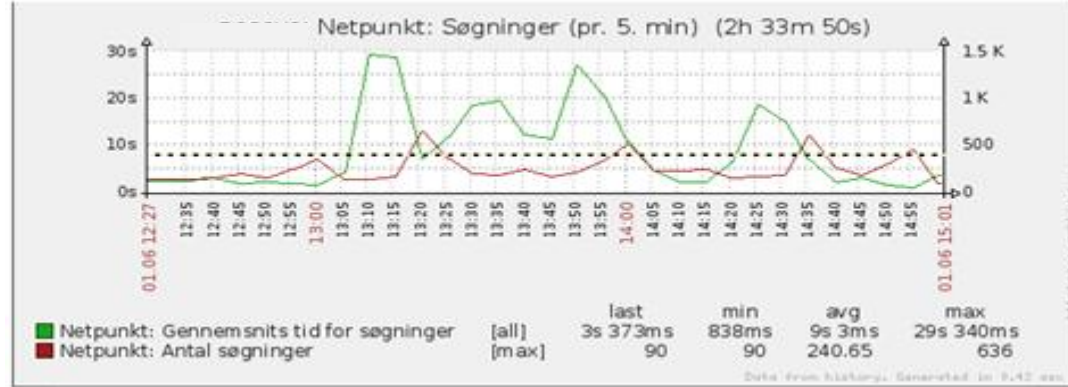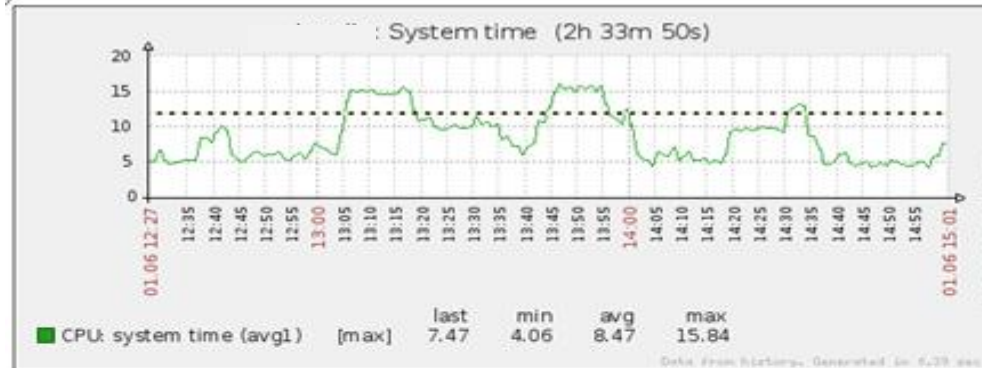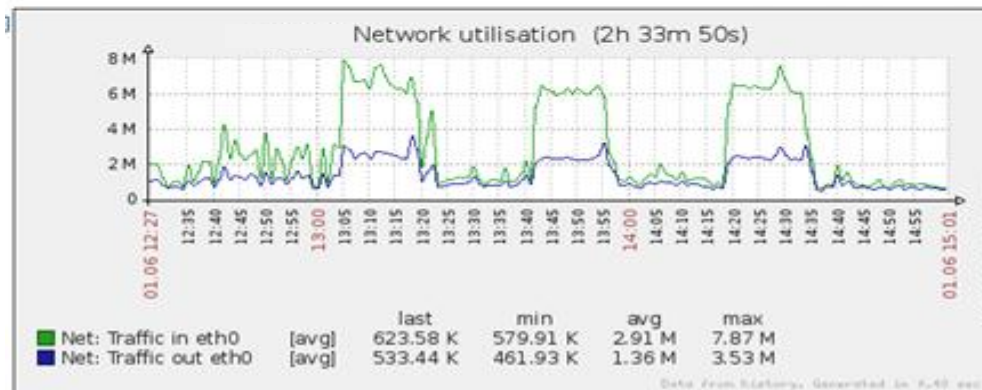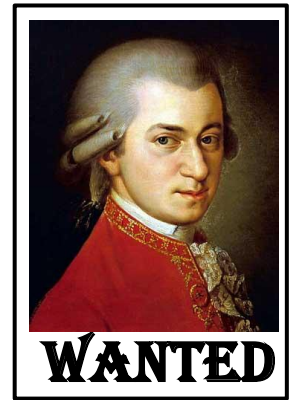Frontend
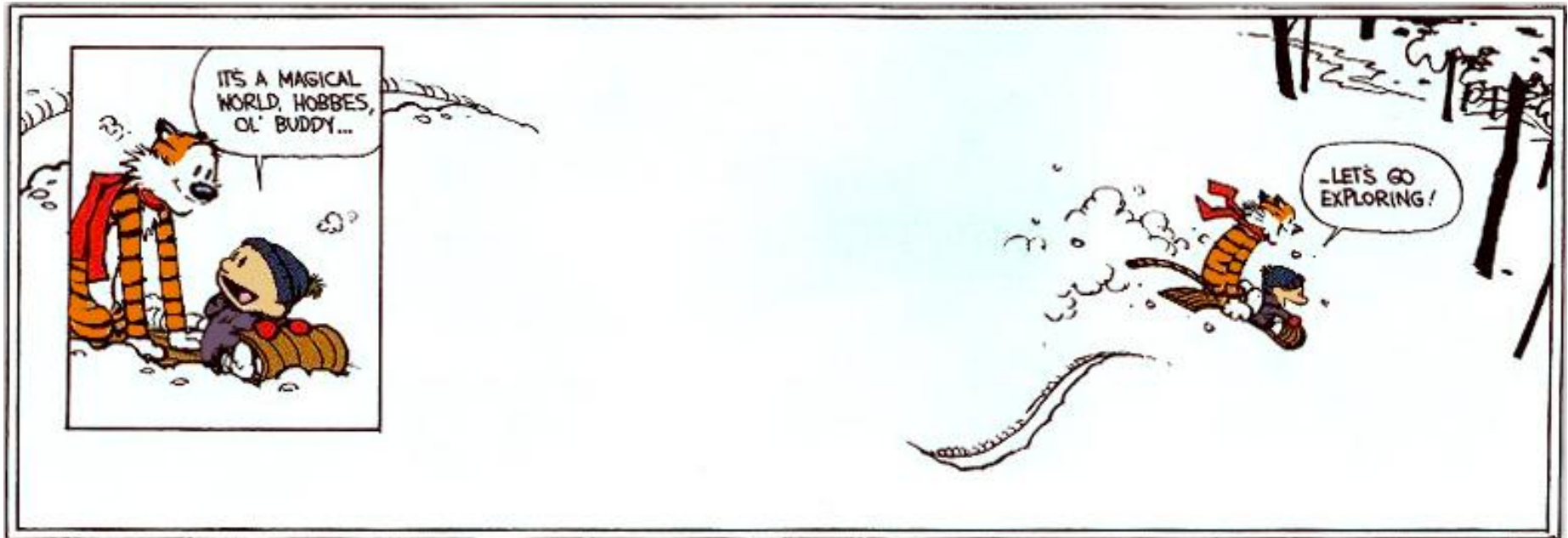
DB server

BaseCserver
host

# The "Monday monster" case

- **Finally a breakthrough !**

- One morning … looking at other graphs … I went to check the baseCserver hosts OS statistic screen …

- Peaks with high Network I/O caught my attention

- Some zoom in and out and …

- … **correlation with DB server system time found !!**

- From time period provided by the correlation, my database admin colleague found this in the DB log:

- **<Timing>Total time for request : 399** (sec)

- Development colleague found the request in a baseCserver log …

- … and could explain the whole story …

# The "Monday monster" case


WANTED

- **The guilty was:          Wolfgang Amadeus Mozart** ☺

- **Explanation / Root cause:**
- A subtask in a search, is to do a dynamic matching of all references
- The problem appeared when someone did a search on media with very many references
- **BaseCserver (was the real guilty)**
- The baseCserver does this matching
- The baseCserver has to match on a big amount of data (records) …
- … that's why we see network peaks on baseCserver hosts
- … and higher systemtime on DB server
- 4 instances of the baseCserver were running
- 1 will be blocked – explaining 25 % transactions timeouts!
- **Solution**
- The "simple" solution was to increase the number of instances of the baseCserver !!!

# The Start …



Calvin and Hobbes

- Zabbix 1.8 Tutorials:
- http://www.zabbix.com/forum/showthread.php?t=22211