# Capacity Planning: A Critical Core Anti-Spam Competency

## MAAWG 9th General Meeting, San Francisco CA
## January 29th-31st, 2007

Joe St Sauver, Ph.D. (joe@uoregon.edu)

MAAWG Senior Technical Advisor

http://www.uoregon.edu/~joe/capacity/

# Why Pay Attention to Capacity Planning?

- **"Logistics sets the campaign's operational limits."**
  -- Joint Pub 1: Joint Warfare of the Armed Forces of
     the United States


- **"He conquers who endures."**
  --Aulus Persius Flaccus, 34-62 A.D.

# We <u>Are</u> In An Arms Race With The Spammers

- It is common to hear the phrase, "We're in an arms race with the spammers," typically after the spammers unleash some new type of tricky-to-filter spam, such as the recent image-based stock and pillz spam. However, that sort of "victory through technical innovation" strategy is subordinate to the <u>real</u> spam "arms race."

- The <u>real</u> spam arms race is one of <u>sheer volume</u>, where the ultimate outcome of the war will be driven by the volume of traffic the spammer can send, versus the volume of traffic that you, the ISP, can successfully deliver/filter.

- If you want to prevail in the "war on spam," your ability to do so will be, at root, a function of your ability to correctly plan for and provision the capacity you'll need to deal with what spammers try to throw at you and your customers.

4

# The Spammer Response To Most Difficulties They Face: <u>Send More</u>

- ISPs deploying better filtering? Send more spam (at least some will probably end up getting through).

- Users less interested in what's being spammed? Send more spam (eventually maybe they'll get curious).

- Not clearing the sort of money you wanted to be making as a spammer? That's easy to fix: send more spam.

- New product available for spamvertising? Great (says the spammer), I'll just send more spam and promote it, too!

- **KEY POINT: When in doubt or facing a problem, spammers will <u>send more spam</u>. Spammers have very scalable delivery infrastructures, and they are generally NOT capacity constrained, so when in doubt, their first instinct is to "hit the gas."**

5

# Spammers <u>Could</u> Also Use Their Excess Delivery Capacity Very Aggressively

- On one level, if spammers have excess delivery capacity, you might imagine they'd just send **more of what they're currently sending, and in pretty much the same way**, and that's often the case. Spammers could, however, use their excess capacity far more "aggressively." For example:

  -- many content based filters don't scan content larger than some threshold size (such as 250KB/message); **what if spammers suddenly began sending jumbo messages** instead of their current comparatively small messages?

  -- currently spammers smear their sending capacity out over multiple ISP targets; but **what if a spammer began concentrating ALL their capacity upon just <u>one</u> ISP target**, eventually moving (on a rolling basis) to another solitary target, etc. Would each ISP be able to keep up?

# Why Pay Attention to Capacity Planning NOW?

# Spam Volumes Are Surging And The Type of Spam Being Sent is Changing

- According to data from Barracuda Networks, an enterprise security appliance vendor in Mountain View, Calif., **there has been a 67 percent increase in overall spam volume and a 500 percent increase in image spam since Aug. 2006**.

  *'Pump-and-Dump' Spam Surge Linked to Russian Bot Herders*

  http://www.eweek.com/article2/0,1895,2060241,00.asp

  November 16, 2006

- Who could have anticipated a surge of that size, or such a radical change in the type of spam being received? No one!

- But what's driving that dramatic jump in spam volumes?

# "Template-Based" Spam…

- "Inside of acting as a proxy for spam senders, each SpamThru client is its own spam engine, downloading a template containing the spam, random phrases to use as hash-busters, random "from" names, and a list of several hundred email addresses to send to."

  Joe Stewart,
  http://www.secureworks.com/research/
  threats/view.html?threat=spamthru

- Thus, spam "pipelines" (e.g., traditional spam zombies functioning  as anonymizing proxies) are being replaced with spam "factories."

- So what's the big deal about that, eh? Well…

# <u>B</u>illion Spam/Day Botnets…

- "[…] SpamThru acts as massive distributed engine for sending spam, but without the cost of maintaining static servers. Total spam capacity is fairly high - with 73,000 bots, given an average SMTP transaction time of 5 seconds, **the botnet is theoretically capable of sending a billion spams in a single day**. This number assumes one recipient per message, however in reality, most spams are delivered in a single message with multiple recipients at the same domain, so **the actual number of separate spams landing in different inboxes could be even higher**, assuming the spammer possesses that many email addresses."

  Joe Stewart, http://www.secureworks.com/research/threats/view.html?threat=spamthru-stats

# Another Report Regarding
# The "Spam Template" Method

- "[…] since August 16th, a well-known Russian spammer is also suspected of using the MS06-040 exploit technique to target unpatched corporate email servers in order to gain control and use them for distributing spam via Pro Mailer DMS, the notorious spam sending software. DMS is potentially more devastating than most spam-sending software since it is able to use the newer "spam cannon" technique that employs a powerful mail-merge of addresses with pre-prepared spam templates. This approach enables the spammer to maximize throughput and distribute **millions of spam messages per hour through a single compromised computer**, as mentioned in MessageLabs Intelligence Report May 2006."
www.londonactionplan.org/files/messagelabs/MessageLabs%20Intelligence%20Report%20-%20August%202006.pdf

# Expect Higher Volumes
# From Real Mail Servers, Too

- Traditionally, many mail servers software products have offered relatively high throughput, but the bar may be about to be raised still higher.

- Based on testing of the Postfix 2.4 queue manager done by Victor Duchovni (which he shared with me on December 26, 2006 and which he agreed to my sharing with you today), it looks like Postfix 2.4 running on a single well-tuned and well-connected Opteron-class system will be able to deliver well over **150 million messages/day**. That's **ONE** single relatively inexpensive system running an MTA that does things "by the book," and doesn't "cut corners"

- If you don't immediately find that sort of throughput to be impressive, let me provide some numerical context…

# What's Currently a "Lot" of Messages From a Single Domain or a Single Host?

- Checking Senderbase.com on 27 January 2007, the single hottest-running domain, as an entire domain, "only" did an estimated **320 million messages per day**

- The hottest single IP listed by Senderbase? That domain "only" did an estimated **16.2 million messages per day**

- Clearly, if a mail sender -- whether good or bad -- wanted to move to a higher performance MTA such as Postfix, there's still room for **an order of magnitude (10X) increase** over the delivery levels of even the hottest-running current single IP traffic source reported by Senderbase… and a problematic sender might have scores of servers available.

- Template spam and hotter-running real mail servers are not the only things driving the increase in mail volume, though.

# More Servers (Instead of Consumer PCs) Are Going To End Up Compromised

- While traditionally most spam has been sent from compromised consumer PCs running Windows, increasing pressure has recently been put on Unix/Linux-based servers.

- Why the change in miscreant focus? A number of reasons. It is becoming increasingly hard to send mail direct-to-MX from dynamic space because of things like port 25 filters, SPF records, descriptive rDNS nomenclature, restrictive client side port filters integrated into some popular antivirus products, the new Spamhaus PBL, and now even Windows Vista -- all those reduce the chance that a bad guy will be able to successfully get his message out via spam zombies.

- The bad news? Many Unix/Linux servers (up to 70%?) have Web 2.0-related vulnerabilities. Server class systems may also have more horsepower and better network connectivity.

14

# There Will Always Be Newer and "Better" [e.g., worse] Worms, etc.

- As if regular traffic and spam traffic isn't enough, you should also recognize and plan for the fact that there will also almost certainly be one or more newer and "better" worms or rapidly propagating viruses in the days ahead.

- Most miscreants have come to recognize that overwhelming the Internet with overly aggressive worms and viruses is bad for their own personal profitability, but mistakes still get made, and miscalculations still occur… and your servers and network links need to be ready to stay up and functional regardless.

# New International Broadband Markets Continue to Come Online

- For example, I am particularly interested in the effect we may see when consumer broadband grows in popularity in India and Pakistan -- there are tremendous numbers of potential users there, and many of them are fluent english speakers eager to engage western family members and business partners -- that's more real mail for your customers

- Given the volume of potential broadband subscribers in India and other parts of central Asia, I would be surprised if there aren't also at least a couple new major spammers who emerge along with all those legitimate new customers, too.

- Finally there will also probably be a whole pile of new spam zombies, too, at least until the users in these new markets learn to harden their systems, just as has been the case in the United States, Europe, the Far East, South America, etc.

# We Should Also Not Discount the Possibility of Online Attack Traffic

- Computer network operations ("CNO") are now an accepted part of US and foreign cyberwar doctrine, and because many military and governmental networks are blocking large parts of the regular Internet, major commercial service providers may be the only "soft" online strategic target left (heck, they need to have SOMETHING they can still attack, right?)

- Similarly, if we believe terrorists might target tangible western interests at home and abroad, is there any reason to believe that western online activities are not also at risk?

- Even if you discount both military computer network operations and cyber terrorism, cyber crime (other than spam) may still end up targeting your systems or you network, with the same result: you'd need more capacity

# Bottom Line, It's Time to Review Your Current Capacity Levels

- Regardless of whether it is new template-based spam cannons, more capable legitimate mailers used by both the good and the bad guys, more capable machines getting popped, a virus/worm that may have gotten out of control, new markets (and new spam sources) coming on line, or attack traffic targeting your servers as part of cyberwar activity, it will be just a matter of time, in my opinion, before things end up getting pretty hot out there….

- **For all these reasons, you may want to review your current capacity, and decide if it might be prudent to increase your system and network capacity now or in the immediate future.**

- Let's talk a little more about that.

# Capacity and the Service Provider

# Sending's Easy But Receiving (And Filtering Spam) Is Hard!

- Each message that's sent by the spammer needs to be received (or at least filtered) by the ISP, and receiving and filtering messages is much more demanding than just sourcing that same traffic -- this is an asymmetric issue.

- One source of that asymmetry is the ISP's need to meet customer expectations (perhaps expressed in the form of a rigid SLA). Customers may be <u>very</u> sensitive to any delays in the receipt of legitimate mail, and any outages caused by a failure to handle surges in volume can have immediate and ongoing financial impacts.

- Spammers, on the other hand, can decide both <u>when</u> and <u>how much</u> to spam, and spammers are generally <u>not</u> under any time or performance pressures (except perhaps in the case of some spammer-for-hire scenarios).

# Capacity Thus Becomes A Critical ISP Defensive Strength

- Need to deploy additional computationally intensive filters to address some new type of spam that's being seen? If you have "lots of capacity" that <u>will</u> be an option for you. If you don't have enough capacity you'll just need to muddle along.

- New email worm or email born virus? If you have "lots of capacity" you can soak up the associated surge in volume without having that unexpected increase in volume impact delivery of legitimate mail.

- Competitive environment change? For example, are your competitors offering customers twice as much (or ten times as much!) disk space as they used to? If you have sufficient capacity, you can directly match that competitive threat.

# Upgrades, Safety Margins and the Unexpected

- Many companies do equipment upgrades or equipment replacement on a scheduled basis, or when utilization reaches a defined threshold. Sometimes though, tight budgets might lead decision makers to "stretch" server replacement plans or to defer network upgrade plans just a *little* bit further. Who among us hasn't heard…

  "Well, the servers and bandwidth have held up *so far,* and *surely* there's a little engineering safety margin built in… I bet that they'll be okay for another quarter or two (or year or two!)… won't they?"

- Of course, if you run through any safety margins, you're not going to be in very good shape when the unexpected occurs.

# Additional Capacity Is Cheap Insurance for Uncertain Times

- If there is one message you take away from this somewhat noisy talk, make it be this one:

  **We live in uncertain times when it comes to spam, and additional capacity is cheap insurance that prevents or ameliorates a lot of ills. Don't run out of gas in the middle of the war. Insure you have sufficient capacity!**

# One Slight Problem: Capacity Costs Money

- If money were no object, we'd all just deploy an "infinite" amount of capacity, and voila, life would be grand.

- Unfortunately:
  -- capital is constrained at most companies
  -- an investment in hardware capacity in one area may result in an a reduction in funds available for use in another area (for example, you may have a choice between purchasing more server capacity or more disk storage, and hiring more staff to respond to complaints or to help users clean up their compromised systems)
  -- truly excess capacity accomplishes nothing, and may hurt some performance metrics (such as ROI)
  -- deployed technology can quickly become obsolete; you don't want to deploy capacity until you need it

# Personnel Capacity Constraints

- It is tempting to focus just on just hard capacity constraints, things such as CPU, memory, I/O, and network capacity, but you must also keep personnel-related capacity constraints in mind.

- The servers you add do need to be administered by someone, and even with the use of modern, highly scalable approaches to server management, you still need to make sure that you're not stretching system administration staff to the breaking point.

- **Key point: ongoing personnel costs may likely dwarf one time capital expenditures for systems or ongoing expenditures for networks!**

- Nonetheless, you <u>must</u> scale server administration staff up along with the hardware you deploy.

# A Second Slight Problem: We Need <u>Lead Time</u> To Add Capacity

- If we could add capacity on demand with zero advance notice (or zero "lead time") we'd have a far easier time when it comes to managing our capacity requirements. If there was no need to plan ahead, we could just add (or remove) capacity as required, adjusting our capacity to meet our empirically observed requirements.

- In Real Life™, however, we need to "pull the trigger" on orders for additional capacity in advance of the time we actually need that capacity. It takes time for orders to be approved, and for purchase orders to be cut, and for systems to be built or customized, and for gear to be shipped and installed and burned in and configured and integrated into production. We need to plan ahead.

# A Third Slight Problem: We Can Only Add Capacity in "Chunks"

- A "chunk" may be the amount of capacity a new server delivers, or a "chunk" may be going from a fast ethernet to a gigabit ethernet or from a gigabit ethernet to a ten gigabit ethernet link, or a "chunk" may be some other unit of incremental capacity.

- Because of this "chunking" phenomena, capacity will typically not always get added in a smooth, continuous way, but in a series of discontinuous increments.

- Chunking means we may have a choice between too much (which can be bad), and too little (which can be worse).

# A Final Slight Problem: Our Crystal Ball Is Still a Little Cloudy

- If we could forecast the future with perfect certainty, we'd know (with perfect certainty) the amount of capacity we'd need, and we could then very efficiently plan to deploy just that much and no more.

- In this world, though, none of us are "Svengalis" and we need to do the best we can with the information we have available to us.

- One thing we need to decide is what we need to measure: what are our fundamental capacity constraints?

# Multiple Potential Capacity Constraints

- If we think about a typical mail server, multiple potential capacity constraints may exist. For example, a server may be:
  -- CPU bound,
  -- Memory bound,
  -- Disk subsystem bound (disk space, inodes, I/O's, etc.),
  -- Network throughput limited (due to the TCP/IP stack),
  -- Network throughput limited (due to the speed of its
      local connection)
  -- Limited by how the server was (mis)-configured, etc.

- Sets of servers may also be subject to aggregate constraints that may not apply to any individual server, such as:
  -- limits to the rack space available in the data center
  -- limits to available power,
  -- limits to available cooling,
  -- limits to available aggregate wide area bandwidth, etc. 29

# All Applicable Constraints May Not Be Immediately/Simultaneously Knowable

- Eliminating one constraint may expose a new constraint, in a fashion that's often referred to as uncovering "rolling bottlenecks."

- Because of the "rolling bottlenecks" phenomena, you may need to do multiple successive capacity upgrades before you truly all the capacity constraints a given system is laboring under.

- In some cases, you may even run into situations where by fixing one capacity constraint, you make another capacity constrained situation more dire still. For example, fixing a network bottleneck may put more I/O load on your disks.

- Capacity constraints do not respect neat little departmental boundaries, either…

# Strict Division of Duties Can Hurt Your Ability to Do Capacity Planning

- How much capacity do you currently have? "A lot?" "Just the amount you currently need?" "Too little even to meet even your current requirements?"

- I'm not really going to ask, but at least some providers here today may not even KNOW the current capacity of some critical assets. For example, if you're a "server guy" you may not know if your network bandwidth is currently sufficient, and if you're a router guy you may not have the slightest idea if your company's mail servers have enough physical memory, and if you're doing MTA administration you may have no idea whether or not DNS servers are struggling or coasting. That needs to change.To properly do capacity planning, someone needs to understand the impact of capacity limitations in ALL relevant areas.

# Capacity Planning
# Needs to Be Data Driven

- Your job is to be, or to become, a data driven son-of-a-gun -- you can't plan without data, so your first objective is (or should be!) to start collecting data.

- There are many different ways you can potentially gather data for tracking and planning purposes, a few of which include:

  -- Analysis of log files
  -- SNMP
  -- Agent-based reporting

- Because each of those approaches has limitations, most providers will want to gather data via multiple methods.

# Analysis of Log Files

- Many providers start their capacity review by analyzing their server's log files:
  - -- How many messages per day are being accepted?
  - -- How many messages per day are being rejected at connection time?
  - -- How many messages per day are being subsequently identified as spam after being accepted for delivery?
  - -- etc., etc., etc.

- Data derived from log files is usually summarized over a period of time, such as a day, or perhaps hour-by-hour.

- Summarization is required for anyone to be able to glean usable nuggets from the otherwise overwhelming volume of data that's often found in log files, but aggregation or summarization in that fashion inescapably leads to "loss" or "condensation" of some information.
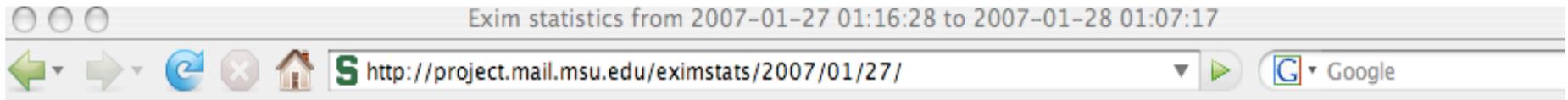
33

# Log File Data Needs to Be Accessible

- In order for log file data to be analyzed, it needs to be accessible, but at least in some cases, log files live only on each physical machine and access to them is then often limited. If you're not currently doing centralized syslogging, you may want to consider doing so.

- Centralized syslogging also reduces the chance that :
  -- too-small local log partitions will overflow
  -- log data will end up getting lost or overlooked,
  -- logs will get tampered with in the event of an intrusion or other security event, etc.

- A popular replacement central syslogger is syslog-ng, see http://www.balabit.com/products/syslog_ng/
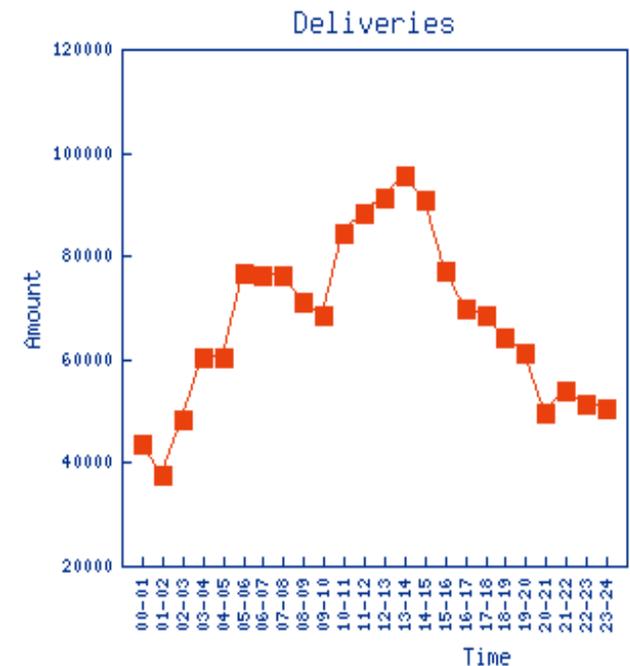
# Log File Summarization/Analysis Tools

- Ultimately, many people find they want custom summaries and end up doing their own custom log file summarization and reporting, but there are many free and commercial log analysis tools one can try first, instead, such as Analog (see http://www.analog.cx/ ), or even log analyzers which may come with your MTA

# One Part of Exim's Log Analysis

http://project.mail.msu.edu/eximstats/2007/01/27/

Google

## Deliveries per hour (each dot is 1912 deliveries)

```
00-01   43588   ......................
01-02   37652   ...................
02-03   48532   .........................
03-04   60132   ...............................
04-05   60206   ...............................
05-06   76859   .......................................
06-07   76266   .......................................
07-08   76236   .......................................
08-09   70923   ....................................
09-10   68703   ...................................
10-11   84184   ...........................................
11-12   88104   ..............................................
12-13   91062   ...............................................
13-14   95594   ..................................................
14-15   90943   ...............................................
15-16   76949   .......................................
16-17   69931   ....................................
17-18   68345   ...................................
18-19   64408   .................................
19-20   61064   ...............................
20-21   49485   .........................
21-22   53939   ............................
22-23   51380   ..........................
23-24   50322   ..........................
```

Deliveries

Amount

Time

# A Warning About Synchronous Logging

- On some systems, logging is done synchronously, and can materially slow system performance or increase system load.

- If asynchronous logging is supported by your system, you may want to test it as an option (if it is supported, enabling it may be as simple a matter as adding a dash prefix to the logfile name in the configuration file).
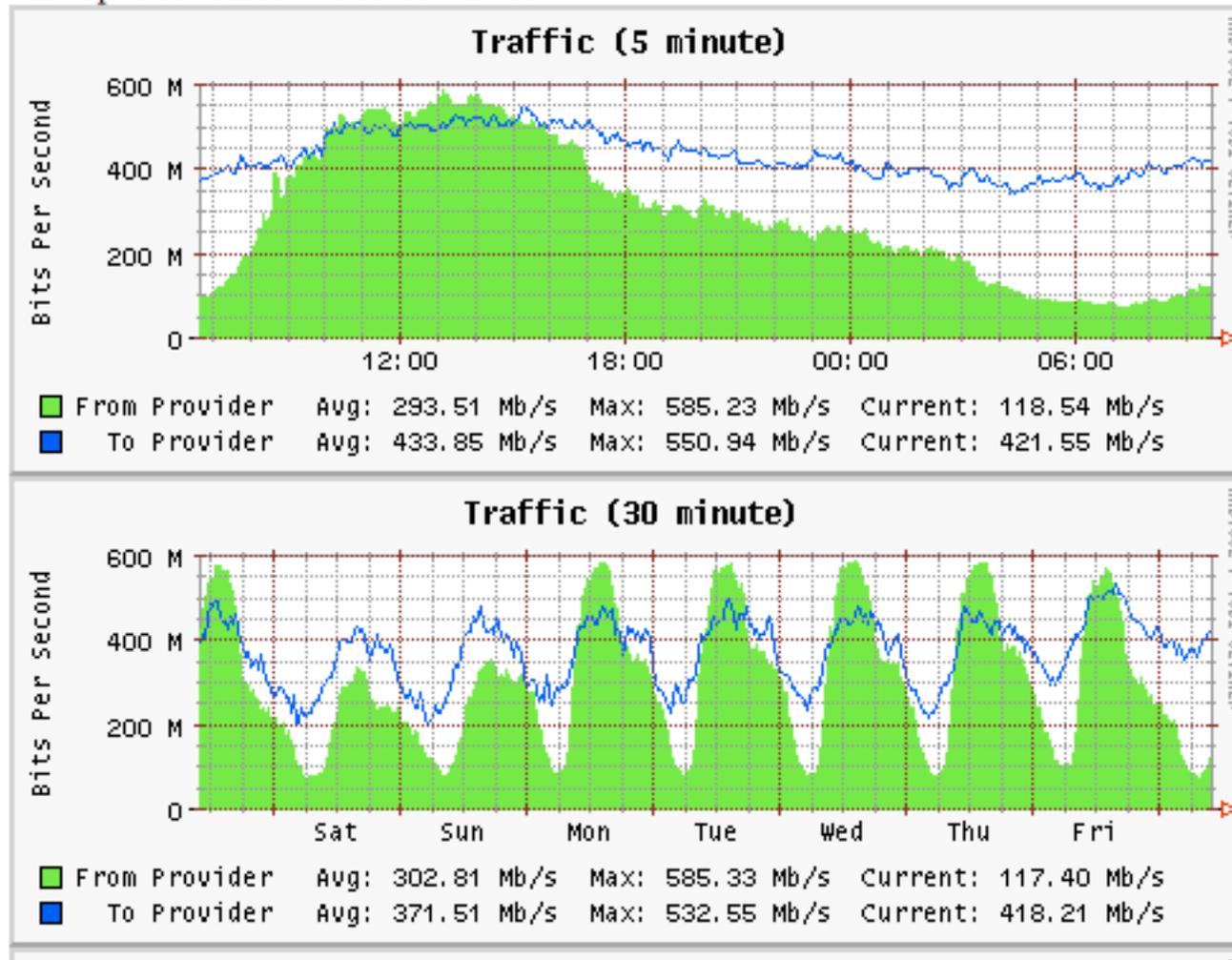
# Watching "Now" Data Instead of Looking Backward

- Log file analysis is retrospective, or "backwards looking" -- the "rearview mirror" on your server.

- Most of us are more interested in what's happening *now* or what's going to happen *in the future*.

- SNMP data is an excellent way of focusing on the "now" data.

# A Quick & Incomplete Overview of SNMP

- SNMP is the "Simple Network Management Protocol" and has traditionally been associated with monitoring and managing things like routers and switches, but SNMP can also be used as a way to collect data from host systems.

- SNMP data is collected by polling SNMP-enabled devices via a network management system, snmpget, or other tool.

- To retrieve data, a user typically needs:
  -- the FQDN or IP address of the SNMP managed device
  -- the "community string" (or password) for SNMP access (all too often this is just "public" for read-only access)
  -- the object ID (or variable name) of the MIB (management information base) of interest, normally a series of numeric values separated by dots

- The value of that OID can then be periodically polled, and will often be graphed using MRTG or RRDtool.

# Sample SNMP-Derived Graphs

# SNMP Limitations

- SNMP is truly a very SIMPLE ("primitive?") protocol. E.G.:

- SNMPv1 was not very secure (commands, data and community strings were all passed "in the clear," and thus were easy to eavesdrop upon)

- It was tedious to walk all subelements of a MIB branch

- Counters would often roll over rapidly due to their limited range

- General access to SNMP data would often have to be limited via firewalls or router ACLs, because SNMP-using devices would often not have the ability to control access themselves

- Subsequent versions of the SNMP protocol addressed these issues -- but at a price of additional device complexity. As a result you may still see many simple network devices that only "speak" SNMPv1

- SNMP is still hugely popular and useful for collecting data.

41

# Warning

- There has recently been increased security interest in network monitoring and management software, with material vulnerabilities found in some popular packages. It is extremely important that you keep all software you use on your network management station up to date, and you should harden and shelter your network monitoring and management station from miscreant attention.

- An example of the sort of thing that's being found is shown on the following slide…

## Cacti Command Execution and SQL Injection V

| | |
|---|---|
| **Secunia Advisory:** | SA23528 |
| **Release Date:** | 2006-12-28 |
| **Last Update:** | 2007-01-18 |
| **Critical:** | Highly critical |
| **Impact:** | Security Bypass<br>Manipulation of data<br>System access |
| **Where:** | From remote |
| **Solution Status:** | Vendor Patch |
| **Software:** | Cacti 0.x |
| **CVE reference:** | CVE-2006-6799 (Se |

**Description**:
rgod has discovered four vulnerabilities in Cacti, which can be exploited by malicious people to bypass certain security restrictions, manipulate data and compromise vulnerable systems.

1) The "cmd.php" and "copy_cacti_user.php" scripts do not properly restrict access to command line usage and are installed in a web-accessible location.

Successful exploitation requires that "register_argc_argv" is enabled.

2) Input passed in the URL to cmd.php is not properly sanitised before being used in SQL queries. This can be exploited to manipulate
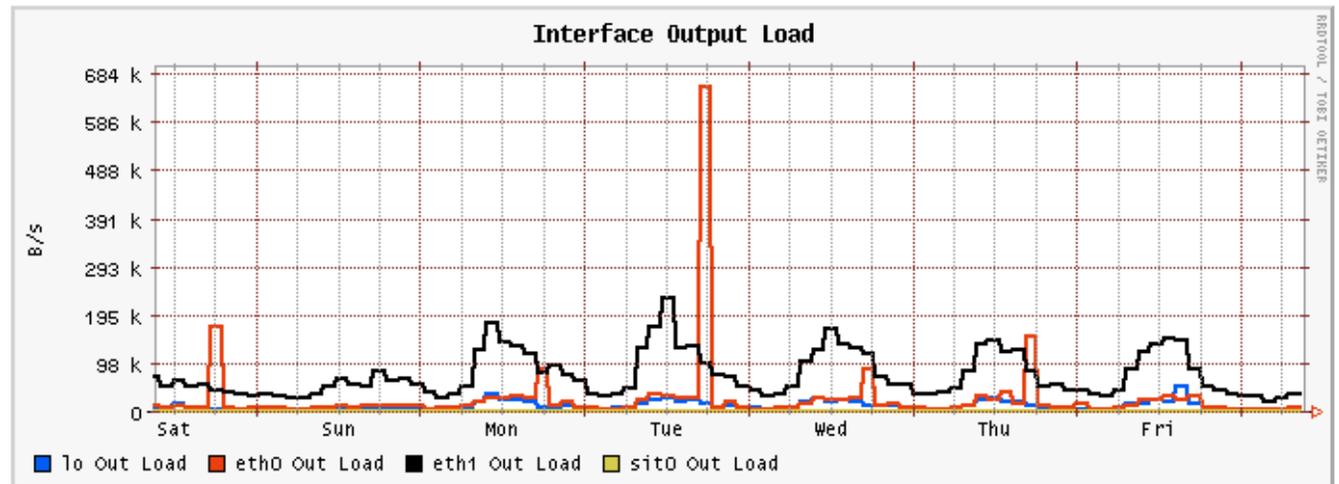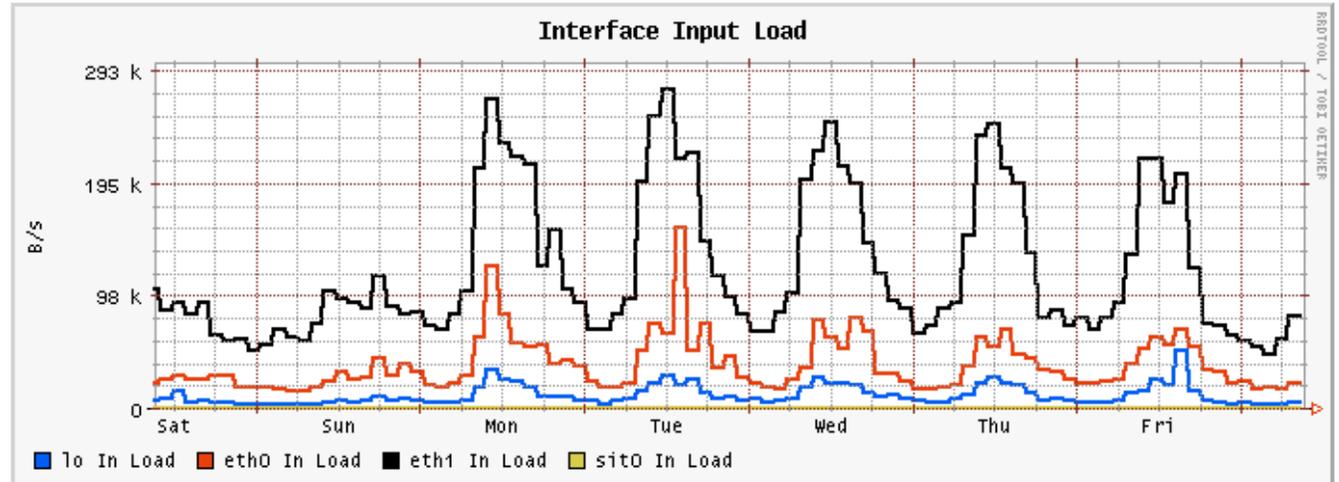
43

# Agent Based Reporting

- An alternative to having an external agent poll systems or devices (SNMP style) is to have systems "phone home" or "report back" with the required data, pushing the data to a central collection host. This can be very convenient for hosts where filter rules do not allow externally instigated connections to be made.

- The process of collecting the required data and then pushing it out is usually coordinated by a software "agent" (or reporting program) running on the system being monitored.

- One example of a free agent-based reporting tool is BigSister (see http://www.bigsister.ch/ ). A list of BigSister agents (or "plugins") is at http://www.bigsister.ch/plugins.html

- BigSister is also able to do direct SNMP queries when that may be more convenient.

# Sample Big Sister Output

# "I Think I <u>AM</u> Seeing Capacity Issues"

# That's Great!

- Well, you know what I mean. It is great that you're paying attention to capacity and it is great that you've identified what may be a bottleneck, not that you HAVE a bottleneck.

- The next step is probably to discuss what you're seeing with all relevant parties. Do not assume that a bottleneck you're seeing will always require you to immediately purchase new gear -- tuning or other steps may allow you to overcome the bottleneck you've identified at no out of pocket cost.

- Be sure to document any changes you make!

- Having found one issue, don't stop there…After you correct that first problem, go back and look again -- you may very well find additional issues.

# What About Formal Capacity Forecasting Models?

# We Can't Cover Statistical Capacity Forecasting Models in Half an Hour

- I considered attempting to include a quick introduction to statistical approaches to capacity requirements planning (e.g., covering regression analysis and/or ARIMA models for example) but then decided that wouldn't make sense given the limited time available.

- Would that be an area of interest for a possible longer future tutorial session?

# Thanks for the Chance to Talk Today!

- Are there any questions?