

Old Man GURU Magazine

*Wychodzi bardzo nieregularnie, kiedy wydaje mi się,
że mam coś ciekawego lub pożytecznego do napisania...*

Numer 20/2011

17 sierpień 2011

Rozpisałem się ostatnio, ale wydaje mi się, że zagadnienia, którym chciałbym poświęcić ten numer są dość istotne w dobie zdecydowanych zmian w podejściu do budowy systemów informatycznych. A więc do rzeczy:

Konsolidacja:

Coraz więcej organizacji decyduje się na konsolidację serwerów w centrach danych. Przeprowadzenie tego procesu umożliwia osiągnięcie wielu korzyści, z których najważniejsze to:

- obniżenie kosztów utrzymania infrastruktury IT,
- wzrost poziomu bezpieczeństwa danych,
- łatwiejsza kontrola nad pracą całego systemu.

Konsolidacja wiąże się jednak z koniecznością użycia sieci WAN. Już dość dawno (od czasu rozpowszechnienia się standardu 100BaseXX) przestano zwracać uwagę na opóźnienia transmisji wynikające z pracy sieci. Spowodowało to rozpowszechnienie się „gadatliwych” protokołów np. CIFS – Common Internet File System (znanego także pod nazwą SMB – Server Message Block). Wspólną cechą gadatliwych protokołów jest uruchamianie wielu podprogramów, które muszą wymienić kilkakrotnie znaczną ilość danych nim zostanie skompletowana pełna transakcja.

Sieci WAN, nawet jeśli zapewniają szybki transfer plikowy charakteryzują się jednak znacznym opóźnieniami „Round Trip Delay” (RTD) dochodzącymi nawet do setek milisekund. W sieci LAN czas ten nie przekracza na ogół 1-2 ms. Jeśli klient i serwer muszą wielokrotnie wymienić informacje aby zakończyć jedną transakcję, to czas RTD będzie się zwielokrotniał i w efekcie użytkownik pracujący w sieci WAN odczuje opóźnienie dochodzące nawet do kilkunastu sekund.

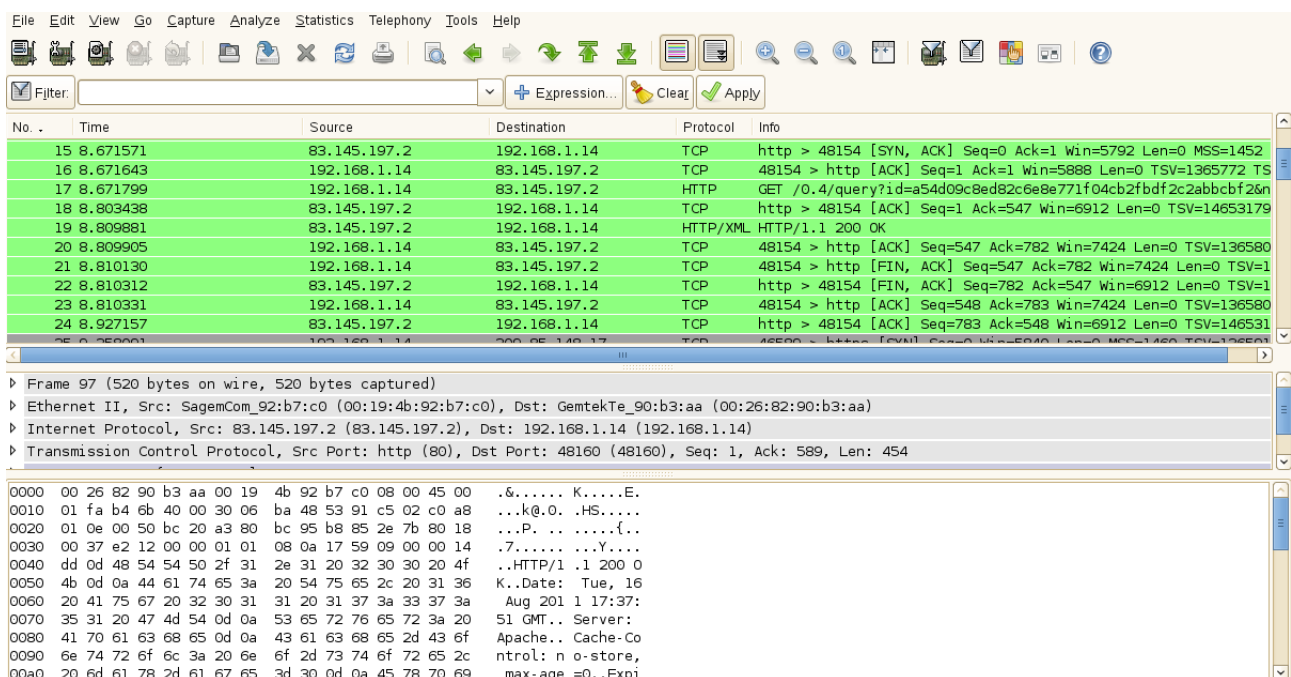
„Webizacja” programów użytkowych:

Przeglądarka WWW staje się coraz częściej uniwersalnym interfejsem użytkownika. Każdy z nas zapewne choć raz skorzystał z poczty Web Mail, coraz większą popularność zyskują również usługi takie jak Google Documents, Picassa i podobne. Jednak „Webizacja” dotyczy nie tylko aplikacji dostępnych w Internecie, ale także systemów klasy „Mission Critical” – przykładem może być choćby wprowadzenie przez Ministerstwo Finansów systemu POLTAX-2, który zastępuje starszą wersję znakową POLTAX-u.

„Webizacja” ma spore zalety, z których największa jest dość oczywista – każdy ma przeglądarkę i umie (jako tako) ją obsługiwać. Nie ma więc potrzeby instalować specjalnych klientów programowych na komputerach PC, w terminalach, tabletach, smartfonach itp.

„Webizacja” oznacza jednak powszechne wykorzystywanie protokołu HTTP i jego pochodnych – np. HTTPS. Protokół ten ma wiele zalet, lecz nie należy bynajmniej do oszczędnych.

Poniżej zamieściłem zrzut fragmentu ekranu prostego analizatora sieciowego rejestrującego transmisje podczas odczytu jednej wiadomości zawierającej jedno słowo „proba” z gmail.com (użytkownik jest zalogowany).



The screenshot shows a Wireshark interface with a list of network packets. The selected packet (Frame 97) is an HTTP response from 83.145.197.2 to 192.168.1.14. The packet details pane shows the following structure:

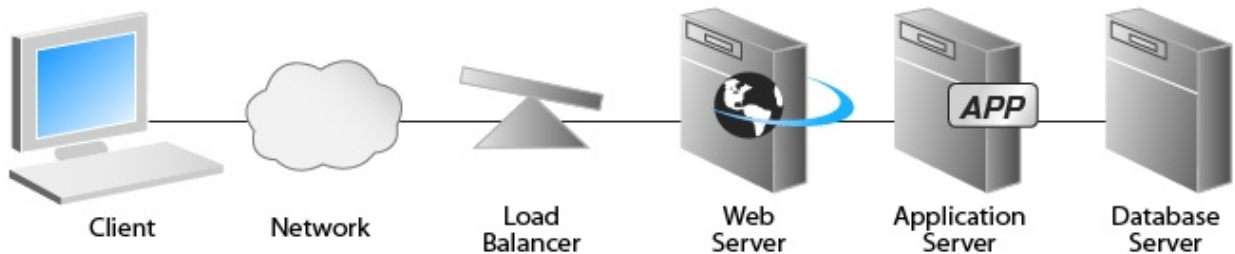
- Ethernet II, Src: SagemCom_92:b7:c0 (00:19:4b:92:b7:c0), Dst: GemtekTe_90:b3:aa (00:26:82:90:b3:aa)
- Internet Protocol, Src: 83.145.197.2 (83.145.197.2), Dst: 192.168.1.14 (192.168.1.14)
- Transmission Control Protocol, Src Port: http (80), Dst Port: 48160 (48160), Seq: 1, Ack: 589, Len: 454
- Application/javascript (application/javascript)
- Raw (454 bytes captured on interface eth0):
 - 0000 00 26 82 90 b3 aa 00 19 4b 92 b7 c0 08 00 45 00 .&.....K....E.
 - 0010 01 fa b4 6b 40 00 30 06 ba 48 53 91 c5 02 c0 a8 ...k@.0..HS....
 - 0020 01 0e 00 50 bc 20 a3 80 bc 95 b8 85 2e 7b 80 18 ...P.{...
 - 0030 00 37 e2 12 00 00 01 01 08 0a 17 59 09 00 00 14 .7.....Y....
 - 0040 dd 0d 48 54 54 50 2f 31 2e 31 20 32 30 30 20 4f ..HTTP/1.1 200 0
 - 0050 4b 0d 0a 44 61 74 65 3a 20 54 75 65 2c 20 31 36 K..Date: Tue, 16
 - 0060 20 41 75 67 20 32 30 31 31 20 31 37 3a 33 37 3a Aug 2011 17:37:
 - 0070 35 31 20 47 4d 54 0d 0a 53 65 72 76 65 72 3a 20 51 GMT.. Server:
 - 0080 41 70 61 63 68 65 0d 0a 43 61 63 68 65 2d 43 6f Apache.. Cache-Co
 - 0090 6e 74 72 6f 6c 3a 20 6e 6f 2d 73 74 6f 72 65 2c ntrol: no-store,
 - 00a0 20 6d 61 78 2d 61 67 65 3d 30 0d 0a 45 78 70 69 max-age=0..Expi

Zrzut ekranu (sniffer Wireshark) wykonany przez Autora

Proszę zauważyć, że serwer i klient wymieniły wiele pakietów pomocniczych niezbędnych do zestawienia i zakończenia połączenia. RTD pomiędzy komputerem klienta o adresie 192.168.1.14 (Neostrada 2 Mbps) i serwerem 83.145.197.2 (Finlandia) wynosi około 100 ms, a więc samo nawiązanie połączenia zajmuje około 0,5 s. Prosty analizator nie umożliwia dokładnego wyznaczenia czasów opóźnień, jednak są one łatwo zauważalne przez użytkowników.

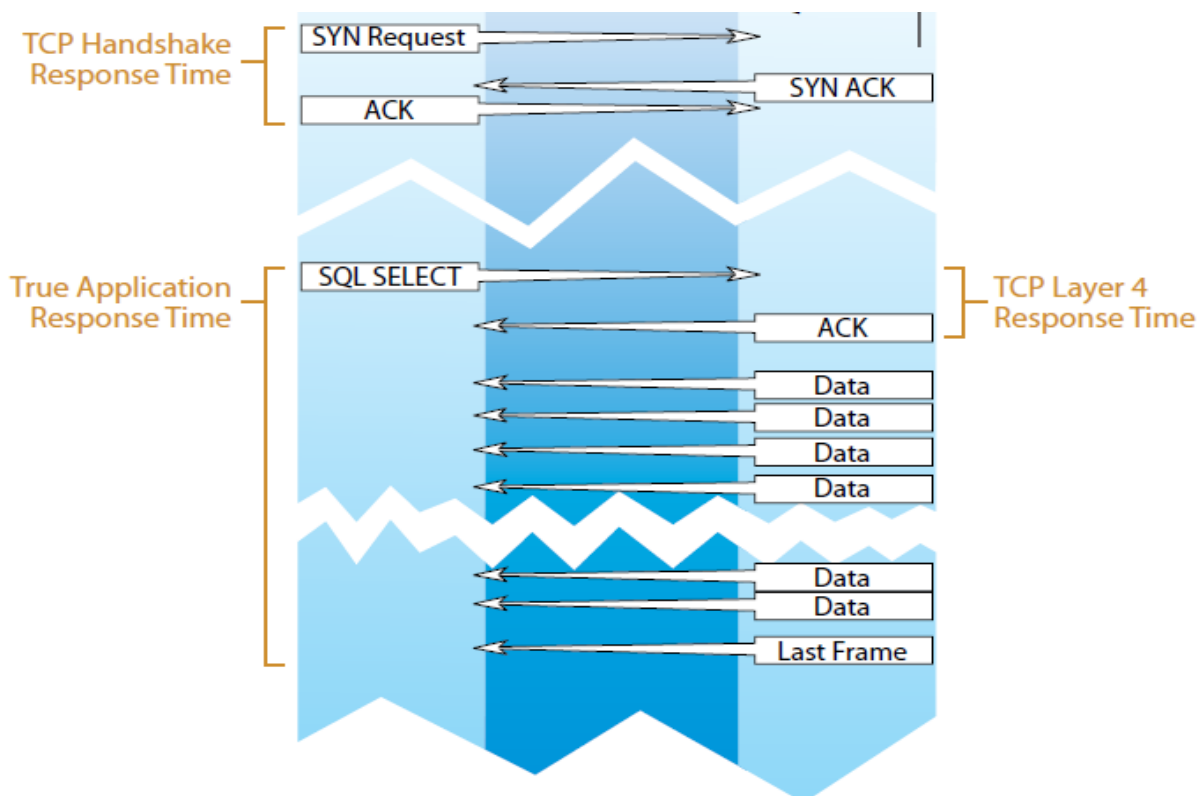
Kolejnym problemem jest coraz powszechniejsze korzystanie z XML. Ogromną zaletą tego języka znaczników jest formalna prostota oraz możliwość kontroli składni (czyżby wróżyło to śmierć wirusów?), jednak obsługa XML wymaga dość znacznego nakładu mocy obliczeniowej, a i sam protokół do oszczędnych nie należy.

Wydajność „Webizowanych” aplikacji zależy jednak nie tylko od transmisji pomiędzy przeglądarką, a serwerem WWW. We współczesnych rozwiązaniach serwer WWW jest jedynie częścią bardziej złożonego systemu:



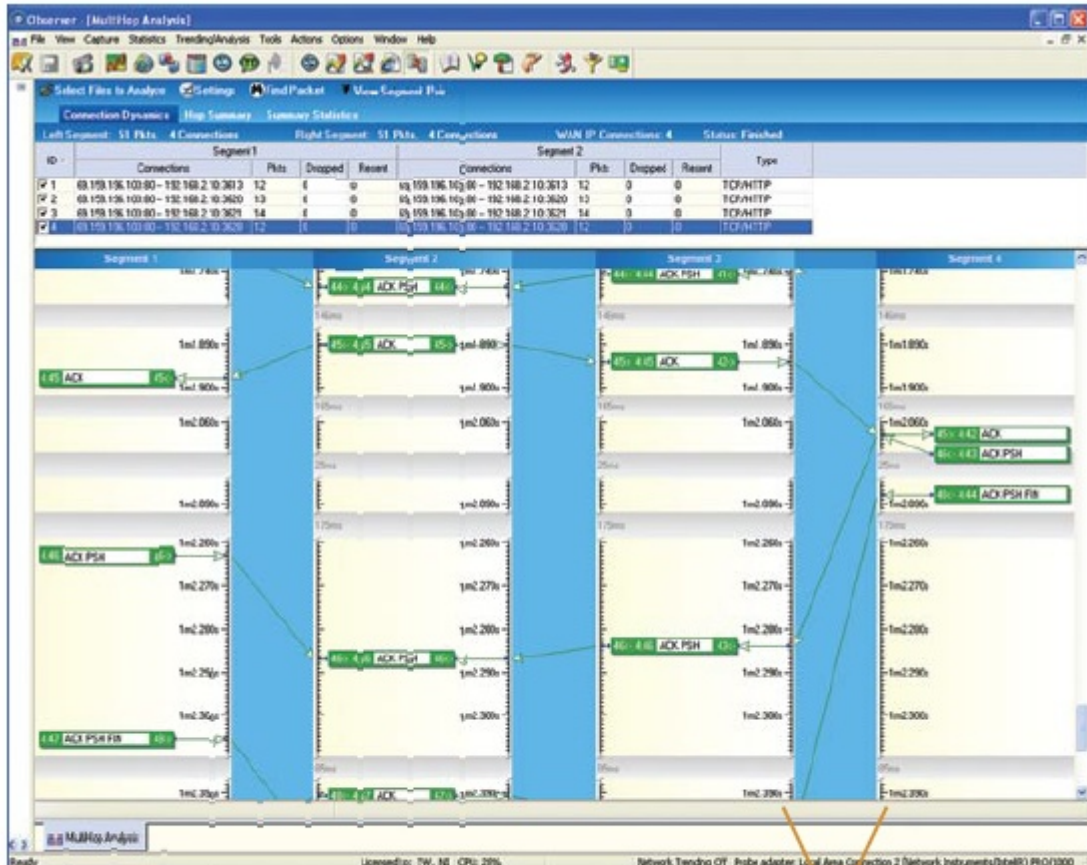
Źródło: Jim Metzler, IT Advisory

Powyższa, dość powszechnie stosowana infrastruktura nosi nazwę „czteropoziomowej” (4-tier). Wydajność całego systemu wynika oczywiście z wydajności wszystkich jego elementów, czasów odpowiedzi serwerów i przepływu danych pomiędzy przeglądarką (Client) a serwerem bazy danych i z powrotem. Optymalizacja pracy takiego systemu bez dysponowania narzędziami umożliwiającymi pomiar tych interwałów czasowych – np. odpowiedzi serwera bazy SQL.



Rysunek z materiałów technicznych Network Instruments

Analiza czasów odpowiedzi pozwala na dokładne poznanie przepływu strumienia danych na całej ich drodze z uwzględnieniem wszystkich serwerów, urządzeń sieciowych, protokołów – a więc stwierdzenie, czy wydajności jego elementów są poprawnie dobrane. Pozwala to na optymalizację kosztów rozwiązania.

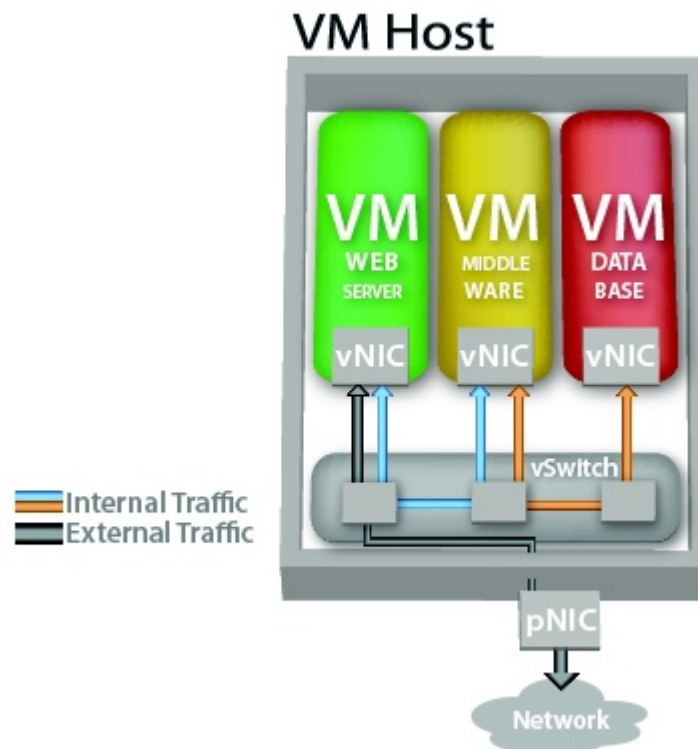


MPLS portion of route

Rysunek z materiałów technicznych Network Instruments

Wirtualizacja:

Środowiska wielopoziomowe są bardzo często realizowane w oparciu o serwery wirtualne. W takim rozwiązaniu zarówno serwer bazy danych, serwer aplikacji (middleware) i serwer WWW wykorzystują jedną fizyczną maszynę. Należy jednak zauważyć, że w takiej sytuacji znaczna część transmisji jest realizowana przez system gospodarza. Często zapominamy o tym, że także te transmisje wpływają na wydajność całego systemu. Zmienne obciążenie zasobów sprzętowych przez maszyny wirtualne wpływa na przepustowość wirtualnego przełącznika – a tym samym na wydajność całego systemu.



Rysunek z materiałów technicznych Network Instruments

Monitorowanie ruchu „sietciowego” pomiędzy maszynami wirtualnymi nie jest zadaniem prostym. Jednym z możliwych rozwiązań jest utworzenie dodatkowej maszyny wirtualnej, która zostanie wykorzystana jako baza dla oprogramowania monitorującego oraz analizującego wewnętrzny ruch sieciowy. Innym możliwym rozwiązaniem jest zainstalowanie w serwerze dodatkowej karty sieciowej, która będzie wykorzystywana do przekazywania wewnętrznego ruchu sieciowego w celu udostępnienia go zewnętrznym urządzeniom analizującym. Rozwiązanie to nosi nazwę „Virtual TAP” (Testing Access Point). Zaletą vTAP jest możliwość integracji analizy wewnętrznego ruchu sieciowego z systemem analizy ruchu w sieci LAN i WAN.

Problemy związane z analizą i nadzorem pracy złożonych systemów wykorzystujących nowoczesne techniki konsolidacji, wirtualizacji oraz udostępniania aplikacji w prywatnych oraz publicznych chmurach obliczeniowych zostały jedynie w tym numerze GURU zasygnalizowane. Warto jednak się im bliżej przyjrzeć, bo uzyskanie oczekiwanych efektów i tym samym zadowolenia użytkowników decyduje o powodzeniu całego przedsięwzięcia. W tym celu musimy dysponować odpowiednim narzędziem. Proste sniffery nie sprawdzą się w złożonych wielopoziomowych (n-Tier) systemach. Wymagania, które musi spełniać nowoczesny analizator pracy sieci są bardzo wygórowane.

Zapis pakietów sieciowych (Packet Capture):

Musi być realizowany bez straty informacji. Jeśli stosowane są sieci o prędkości 10 Gbps (a mamy z tym do czynienia coraz częściej) analizator powinien móc je przechwytywać z pełną prędkością Wire Speed nawet przy pełnym obciążeniu sieci. Należy też uwzględnić ilość danych, która w takim przypadku bardzo szybko może osiągnąć wiele Terabajtów. Wymaga to (nawet dla sieci 1 Gbps) użycia wydajnego sprzętu i 64 bitowego systemu operacyjnego.

Jeśli stosowane są łącza zwielokrotniane (np. 2x1 Gbps) analizator powinien umożliwiać zapis obu kanałów oraz późniejsze korelowanie zapisanych informacji.

Ponieważ analizowane są na ogół złożone sieci o skomplikowanej topologii oraz wykorzystujące różne protokoły współczesne analizatory są wyposażane w system próbników (Network Probes), z których dane są pobierane, analizowane i prezentowane na stanowisku zarządzającym (konsoli systemu). Dane mogą być pobierane zarówno z portów przełączników i innych urządzeń sieciowych lub bezpośrednio z mediów transmisyjnych – kabli miedzianych, światłowodowych (za pomocą rozgałęźników nTAP) lub z sieci bezprzewodowych. Umożliwia to również wykrywanie błędów w warstwie fizycznej (np. generowanych przez uszkodzone urządzenia), które powodują poważne zakłócenia w pracy sieci, a trudne są do wykrycia innymi metodami. W praktyce udało mi się to kilkakrotnie.

Należy z całym naciskiem podkreślić, że nawet wydajny komputer przenośny typu laptop może w wielu przypadkach nie sprostać wymaganiom wydajnościowym. Potrzebny jest silny komputer wieloprocessorowy (quad core) z 64 bitowym systemem operacyjnym, sporą pamięcią i szybkim podsystemem dyskowym. Przeanalizować przecież możemy tylko te dane (pakiety), które zdołamy zapisać.

Parsing i prezentacja wyników:

Analizując pracę rzeczywistych sieci mamy na ogół do czynienia z ogromną ilością danych. Wymaga to odfiltrowania informacji, które nas interesują. Jest to złożone zadanie. Administratora sieci interesuje jak najszybsze dotarcie do źródeł zakłóceń, opóźnień w pracy aplikacji, na które skarżą się użytkownicy systemu a nawet otrzymywanie ostrzeżeń o pojawiających się problemach – na przykład okresowych spadkach wydajności systemu w warunkach zwiększonego obciążenia - przekroczenia dopuszczalnego limitu czasowego realizacji transakcji.

Wyniki powinny być prezentowane w sposób umożliwiający szybką i jednoznaczną ich interpretację oraz generowane w postaci raportów. Oczywiście przy tworzeniu raportu powinny być uwzględnione dane ze wszystkich (lub wybranych) próbników umieszczonych w różnych punktach sieci. Cenną właściwością jest również możliwość analizy „historycznej” (np. pracy sieci w określonych dniach i godzinach) lecz wymaga to retencji dużej ilości danych, co wiąże się z określonymi kosztami.

Praca z analizatorem sieciowym nie może polegać na przeglądaniu zapisanych ciągów pakietów – nawet jeśli zastosujemy ich filtrację. W taki sposób można oczywiście rozwiązać wiele problemów, lecz jest to pracochłonne dla administratora. Użytkowników (a co za tym idzie Szefów) nie interesują szczegóły. Stwierdzają po prostu: „dzisiaj rano znów miałem trudności z zalogowaniem się i wszystko działało bardzo wolno”, „Nie mogłem skorzystać z dostępu do systemu POLTAX-2 (lub innego)”, „Nie miałem dostępu do serwera aplikacji biurowych” - a po 10:15 wszystko zaczęło działać normalnie.

Administrator musi na takie pytania możliwie jak najszybciej odpowiedzieć oraz przygotować odpowiednie propozycje, które zapobiegą takim sytuacjom w przyszłości. Wiele osób zarzuca mi, że zbyt często stosuję porównania „motoryzacyjne”, jednak narzucają się one same – samochód z lat 1970 mógł naprawić prawie każdy. Dziś niezbędny jest dostęp do urządzenia umożliwiającego odczyt kodów OBD II. Oto małe fragmenty tabeli tych kodów (w praktyce stosuje się prawie 1000!):

DTC Codes - P0100-P0199 – Fuel and Air Metering

- * P0100 Mass or Volume Air Flow Circuit Malfunction
- * P0101 Mass or Volume Air Flow Circuit Range/Performance Problem
- * P0102 Mass or Volume Air Flow Circuit Low Input
- * P0103 Mass or Volume Air Flow Circuit High Input
- * P0104 Mass or Volume Air Flow Circuit Intermittent
- * P0105 Manifold Absolute Pressure/Barometric Pressure Circuit Malfunction

Profesjonalne urządzenie diagnostyczne nie tylko powinno odczytać kod błędu (wartość kodu DTC) ale także przekazać diagnoście informację o rodzaju usterki – a nawet możliwych sposobach jej usunięcia. Trudno bowiem marnować czas na przeszukiwanie wielostronicowych tabeli – znacznie efektywniejsze jest skorzystanie z systemu ekspertowego.

Współczesne systemy IT na pewno nie są mniej skomplikowane niż nowoczesny samochód i ich obsługa wymaga również profesjonalnych narzędzi. Przygotowanie takich narzędzi wymaga bardzo dużej wiedzy specjalistycznej i często są one oferowane przez wysoko wyspecjalizowane firmy. W moim przekonaniu na narzędziach nie należy oszczędzać. Posiadanie systemu ekspertowego do analizy pracy sieciowego systemu komputerowego jest pomocne nie tylko przy wykrywaniu źródeł zakłóceń w jego pracy, lecz także przy projektowaniu jego rozbudowy lub modernizacji dzięki narzędziom umożliwiającym prowadzenie analizy „Co jeżeli” (What If Analysis).

Analiza prowadzi do dokładnego określenia kierunków rozwoju systemu i jej przeprowadzenie umożliwia uniknięcie zbędnych kosztów, które powoduje dość popularny „oversizing” zasobów „na wszelki wypadek”. Efektywny nadzór nad pracą sieci zwiększa także efektywność pracy organizacji, co wpływa korzystnie na osiągnięte przez nią wyniki biznesowe i ocenę jej działalności.

Korzyści związane z posiadaniem dobrych, profesjonalnych narzędzi zawsze przewyższają ich koszt zakupu.

Dziękuję IT Advisory i firmie Network Instruments za zezwolenie na wykorzystanie rysunków i diagramów.