High-Speed Traffic Capture and Analysis Using Open-Source Software and Commodity Hardware

Part 1: Packet Capture

Luca Deri <deri@ntop.org>



1





UNIVERSITÀ DI PISA

Overview

- Accelerating packet capture and analysis: PF_RING.
- Layer 7 kernel packet filtering and processing.
- Direct NIC Access: PF_RING DNA.
- Towards 10 Gbit packet capture using commodity hardware.
- Strong Multicore NIC: Tilera Tile64







Accelerating Packet Capture and Analysis: PF_RING





2nd TMA PhD School - June 2011



Packet Capture: Open Issues

- Monitoring low speed (100 Mbit) networks is already possible using commodity hardware and tools based on libpcap.
- Sometimes even at 100 Mbit there is some (severe) packet loss: we have to shift from thinking in term of speed to number of packets/ second that can be captured analyzed.
- Problem statement: monitor high speed (1 Gbit and above) networks with common PCs (64 bit/66 Mhz PCI/X/Express bus) without the need to purchase custom capture cards or measurement boxes.
- Challenge: how to improve packet capture performance without having to buy dedicated/costly network cards?







Packet Capture Goals

- Use commodity hardware for capturing packets at wire speed with no loss under any traffic condition.
- Be able to have spare CPU cycles for analyzing packets for various purposes (e.g. traffic monitoring and security).
- Enable the creation of software probes that sport the same performance of hardware probes at a fraction of cost.







Socket Packet Ring (PF_RING)



PF_RING Internals



http://en.wikipedia.org/wiki/Circular_buffer







PF_RING Packet Journey [1/2]









PF_RING Packet Journey [2/2]









PF_RING: Benefits

- It creates a straight path for incoming packets in order to make them first-class citizens.
- No need to use custom network cards: any card is supported.
- Transparent to applications: legacy applications need to be recompiled in order to use it.
- No kernel or low-level programming is required.
- Developers familiar with network applications can immediately take advantage of it without having to learn new APIs.







PF_RING: Performance Evaluation

Pkt Size	Kpps	Mpps	% CPU Idle	Wire-Speed	
250	259.23	518	518 > 90%		
250	462.9	925.9	88%	Yes	
128	355.1	363.6	86%	Yes	
128	844.6	864.8	82%	Yes	

Test setup: pfcount, full packet size, 3.2 GHz Celeron (single-core) - IXIA 400 Traffic Generator







PF_RING on Embedded Devices



http://nst.sourceforge.net/nst/docs/user/ch09s02.html



mo

12

PF_RING Socket Clustering [1/2]

- In order to exploit modern computer architectures either multiprocessing or threading have to be used.
- Often computer programs are monolithic and hard to split into several concurrent and collaborating elements.
- In other cases (proprietary applications) source code is not available hence the application cannot be modified and split.
- There are hardware products (e.g. see cPacket's cTap) that split/ balance network traffic across network hosts.
- What is lacking at the operating system level is the concept of distributing sockets across applications. This is because network sockets are proprietary to an application/address-space.







PF_RING Socket Clustering [2/2]

- Socket clustering is the ability to federate PF_RING sockets similar, but opposite, to network interface bonding.
- The idea is simple:
 - Run several monitoring applications, each analyzing a portion of the overall traffic.
 and/or
 - Create multithreaded applications that instead of competing for packets coming from the same socket, have private per-thread sockets.







PF_RING Clustering: Threads



Vanilla PF_RING Application

PF_RING Socket Cluster





2nd TMA PhD School - June 2011



PF_RING Clustering: Applications



- Same as clustering with threads, but across address spaces.
- PF_RING allows clustering to be enabled seamlessly both at thread and application level.





2nd TMA PhD School - June 2011



PF_RING Clustering: Code Example

```
if((pd = pfring open(device, promisc, snaplen, 0)) == NULL) {
   printf("pfring open error\n");
    return (-1);
  } else {
   u int32 t version;
   pfring version(pd, &version);
   printf("Using PF RING v.%d.%d.%d\n",
           (version & 0xFFFF0000) >> 16, (version & 0x0000FF00) >>
8,
           version & 0x00000FF);
  }
  if(clusterId > 0) {
      int rc = pfring set cluster(pd, clusterId);
      printf("pfring set cluster returned %d\n", rc);
```



}



PF_RING Clustering: Summary

- Network traffic balancing policy across socket clusters
 - Per-flow (default)
 - Round-Robin
- Advantages:
 - No locking required when threads are used
 - Ability to distribute the load across multiple applications
 - Very fast as clustering happens into the kernel.
- Socket clustering has been the first attempt to make PF_RING more multi-processing/core friendly.







PF_RING: Packet Filtering [1/2]

- PF_RING has addressed the problem of accelerating packet capture.
- Packet filtering instead is still based on the "legacy" BPF code.
- This means that:
 - Each socket can have up to one filter defined.
 - The packet needs to be parsed in order to match the filter, but the parsing information is not passed to user-space.
 - The BPF filter length can change significantly even if the filter is slightly changed.







PF_RING: Packet Filtering [2/2]

#	tcpc	lump	-d	"udp"	
(0	00)	ldh		[12]	
(0	01)	jeq		#0x800	
(0	02)	ldb		[23]	
(0	03)	jeq		#0x11jt	4
(0	04)	ret		#96	
(0	05)	ret		# O	

# tcpd	lump -d "u	udp and port 53"				
(000)	ldh	[12]				
(001)	jeq	#0x800	jt	2	jf	12
(002)	ldb	[23]				
(003)	jeq	#0x11	jt	4	jf	12
(004)	ldh	[20]				
(005)	jset	#0x1fff	jt	12	jf	6
(006)	ldxb	4*([14]&0xf)				
(007)	ldh	[x + 14]				
(008)	jeq	#0x35	jt	11	jf	9
(009)	ldh	[x + 16]				
(010)	jeq	#0x35	jt	11	jf	12
(011)	ret	#96				
(012)	ret	# O				







Beyond BPF Filtering [1/2]

- VoIP and Lawful Interception traffic is usually very little compared to the rest of traffic (i.e. there is a lot of incoming traffic but very few packets match the filters).
- Capture starts from filtering signaling protocols and then intercepting voice payload.
- BPF-like filtering is not effective (one filter only). When multiple filters need to be enforced, each one has to be executed individually.
- It is necessary to add/remove filters on the fly with hundred active filters.







Beyond BPF Filtering [2/2]

Solution

- Filter packets directly on device drivers (initial release) and PF_RING (second release).
- Implement hash/bloom based filtering (limited false positives) but not BPF at all.
- Memory effective (doesn't grow as filters are added).
- Implemented on Linux on Intel GE cards. Great performance (virtually no packet loss at 1 Gbit).
- No much difference between PF_RING and driver filtering hence the code has been moved to PF_RING.







Dynamic Bloom Filtering [1/2]











Bloom Filters [2/2]

- Ability to specify a thousand different IP packet filters
- Ability to dynamically add/remove filters without having to interrupt existing applications.
- Only "precise" filters (e.g. host X and port Y) are supported.
- The filter processing speed and memory being used is proportional to the number of filters but independent from their number and complexity.







Dynamic Bloom Filtering



- Available into PF_RING (in 3.x series up to 3.7.x).
- Ability to set per-socket bloom filters







PF_RING: Bloom Evaluation

- Tests performed using a dual Xeon 3.2 GHz CPU injecting traffic with an IXIA 400 traffic generator with 1:256 match rate.
- Packet loss only above 1.8 Mpps (2 x 1 Gbit NICs).
- Ability to specify thousand of filters with no performance degradation with respect to a single filter (only false positive rate increases).







Bloom Filters Limitations [1/2]

- Bloom filtering has shown to be a very interesting technology for "precise" packet filtering.
- Unfortunately many application require some features that cannot be easily supported by blooms:
 - port ranges
 - negative expressions (not <expression>)
 - IP address/mask (where mask != /32)
 - in case of match, know what rule(s) matched the filter







Bloom Filters Limitations [2/2]

- Possible workarounds
 - Support ranges by calculating the hash on various combinations
 - 5-tuple for perfect matching (proto, ip/port src, ip/port dst)
 - multiple bloom dictionaries for /32, /24, /16, and /8 networks for network match
- Note that as bloom matching is not exact, using a bloom dictionary for storing negative values (i.e. for implementing the not) is not a good idea. This is because not(false positive) means that a packet might be discarded as the filter is not match although this packet passed the filter.
- In a nutshell:
 - Bloom filters are a fantastic technology for exact packet matching
 - PF_RING must also offer support for 'partial' filtering.







Extended PF_RING Filters [1/2]

The author has made a survey of network applications and created a list of desirable features, that have then been implemented into PF_RING:

- "Wildcard-ed" filters (e.g. TCP and port 80). Each rule has a rule-id and rules are evaluated according to it.
- Precise 5-tuple filters (VLAN, protocol, IP src/dst, port src/dst).
- Precise filters (e.g. best match) have priority over (e.g. generic) wilcarded filters.
- Support of filter ranges (IP and port ranges) for reducing the number of filters.
- Support of mono or bi-directional filters, yet for reducing number of filters.
- Ability to filter both on standard 5-tuple fields and on L7 fields (e.g. HTTP method=GET).







Extended PF_RING Filters [2/2]

- Parsing information (including L7 information) need to be returned to userspace (i.e. do not parse the packet twice) and to all PF_RING components that for various reasons (e.g. due to socket clustering) need to have accessed to this information.
- Per-filter policy in case of match:
 - Stop filtering rule evaluation and drop/forward packet to user-space.
 - Update filtering rule status (e.g. statistics) and stop/continue rule evaluation without forwarding packet to user-space.
 - Execute action and continue rule evaluation (via PF_RING plugins).
- Filtering rules can pass to user-space both captured packets or statistics/ packet digests (this for those apps who need pre-computed values and not just raw packets).







PF_RING Packet Parsing [1/4]

- Contrary to BPF that basically does parse packets while filtering them, PF_RING filtering requires packet to be parsed first.
- Parsing information is propagated up to the userland.
- The basic PF_RING engine contains parsing up to TCP/UDP.







PF_RING Packet Parsing [2/4]

```
struct pkt parsing info
 /* Core fields (also used by NetFlow) */
 u int8 t dmac[ETH ALEN], smac[ETH ALEN]; /* MAC src/dst addresses */
 u int16 t eth type; /* Ethernet type */
 u int16 t vlan id; /* VLAN Id or NO VLAN */
 u int8 t ip version;
 u int8 t 13 proto, ip_tos; /* Layer 3 protocol/TOS */
 ip addr ip src, ip dst; /* IPv4 src/dst IP addresses */
 u int16 t 14 src port, 14 dst port; /* Layer 4 src/dst ports */
  struct {
   u int8 t flags; /* TCP flags (0 if not available) */
   u int32 t seq num, ack num; /* TCP sequence number */
 } tcp;
 u int16 t last matched plugin id; /* If > 0 identifies a plugin to that matched the packet */
 u int16 t last matched rule id; /* If > 0 identifies a rule that matched the packet */
 struct pkt offset offset; /* Offsets of L3/L4/payload elements */
 /* Leave it at the end of the structure */
 packet user detail pkt detail;
};
```







PF_RING Packet Parsing [3/4]

- The decision to always parse the packet is motivated as follows:
 - -Packet parsing is very cheap (in terms of computation) and its slow-down is negligible.
 - Beside rare exceptions (e.g. for packet-to-disk applications), user space applications will need to parse packets.
- PF_RING does not natively include layer-7 packet filtering as this is delegated by plugins as shown later in this presentation.







PF_RING Packet Parsing [4/4]



Plugin-based Parsing

ts	caplen	len	parsed_pkt	parsed len	l7 parsing	Payload
					(Optional)	







PF_RING: Exact Filters [1/2]

• Exact filters (called hash filtering rules) are used whenever all the filtering criteria are present in the filter.

```
typedef struct {
    u_int16_t vlan_id;
    u_int8_t proto;
    u_int32_t host_peer_a, host_peer_b;
    u_int16_t port_peer_a, port_peer_b;
    [...]
    hash_filtering_rule;
```

- Exact filters are stored in a hash table whose key is calculated on the filter values.
- When a packet is received, the key is calculated and searched into the filter hash.







PF_RING: Exact Filters [2/2]

- Filters can have a rule associated to it such as:
 - Pass packet to userland in case of match.
 - Drop packet in case of match.
 - Execute the action associated with the packet.
- Actions are implemented into plugins. Typical action include:
 - Add/delete filtering rule
 - Increment specific traffic counters.
 - Interact with the Linux kernel for performing specific actions.

```
typedef struct {
  [...]
  rule_action_behaviour rule_action; /* What to do in case of match */
  filtering_rule_plugin_action plugin_action;
  unsigned long jiffies_last_match;
} hash_filtering_rule;
```




PF_RING: Wildcard-ed Filters [1/2]

- This filter family has to be used whenever:
 - Not all filter elements are set to a specific value.
 - The filter contains value ranges.
- Filters are bi-directional (i.e. they are checked on both source and destinations fields.
- Filtering rules have a unique (in the PF_RING socket) numeric identifier that also identifies the rule evaluation order.

```
typedef struct {
 u int8 t dmac[ETH ALEN], smac[ETH ALEN]; /* Use '0' (zero-ed MAC address) for any MAC address.
                                              This is applied to both source and destination. */
                                     /* Use '0' for any vlan */
 u int16 t vlan id;
 u int8 t proto;
                                     /* Use 0 for 'any' protocol */
 ip addr host low, host high;
                                     /* User '0' for any host. This is applied to both source
                                        and destination. */
                                     /* All ports between port low...port high
 u int16 t port low, port high;
                                        0 means 'any' port. This is applied to both source
                                        and destination. This means that
                                        (proto, sip, sport, dip, dport) matches the rule if
                                        one in "sip & sport", "sip & dport" "dip & sport"
                                        match. */
} filtering rule core fields;
```

ntop



2nd TMA PhD School - June 2011



PF_RING: Wildcard-ed Filters [2/2]

- Filters can optionally contain some extended fields used for:
 - Matching packet payload
 - Implementing more complex packet filtering by means of plugins (see later).
- User-space PF_RING library allows plugins to specify some parameters to be passed to filters (e.g. pass only HTTP packets with method POST).



2nd TMA PhD School - June 2011



Combining Filtering with Balancing [1/4]

- PF_RING clustering allows socket to be grouped so that they be used for effectively sharing load across threads and processes.
- Clustering works at PF_RING socket level and it's basically a mechanism for balancing traffic across packet consumers.
- PF_RING filtering rules combine the best of these technologies by implementing traffic balancing for those packets that match a certain filter.
- The idea is to have the same filter specified for various sockets that are the grouped together. Packets matching the filter are then forwarded only to one of the sockets.







Combining Filtering with Balancing [2/4]









Combining Filtering with Balancing [3/4]

Filtered packets are balanced across sockets as follows







Combining Filtering with Balancing [4/4]



 Using balancing for distributing load across applications/ threads is very effective for exploiting multi-processor/core architectures.







PF_RING Packet Reflection [1/3]

- Often, monitoring applications need to forward filtered packets to remote systems or applications.
- Traffic balancers for instance are basically a "filter & forward" application.
- Moving packets from the kernel to userland and then back to the kernel (for packet forwarding) is not very efficient as:
 - Too many actors are involved.
 - The packet journey is definitively too long.
- PF_RING packet reflection is a way to forward packets that matched a certain filter towards a remote destination on a specific NIC (that can be different from the one on which the packet has been received).
- Packet reflection is configured from userland at startup.
- All forwarding is performed inside the kernel without any application intervention at all.







PF_RING Packet Reflection [2/3]

```
/* open devices */
if((pd = pfring open(in dev, promisc, 1500, 0)) == NULL)
  printf("pfring open error for %s\n", in dev);
  return -1;
} else
  pfring set application name(pd, "forwarder");
if ((td = pfring open(out dev, promisc, 1500, 0)) == NULL) {
  printf("pfring open error for %s\n", out dev);
  return -1;
} else
  pfring set application name(td, "forwarder");
/* set reflector */
if (pfring set reflector(pd, out dev) != 0)
  printf("pfring set reflector error for %s\n", out dev);
  return -1;
/* Enable rings */
pfring enable ring(pd);
pfring enable ring(td);
while(1) sleep(60); /* Loop forever */
```







PF_RING Packet Reflection [3/3]

- PF_RING packet reflection allows easily and efficiently to implement:
 - Filtering packet balancers
 - (Filtering) Network bridges
- In a nutshell this technique allows to easily implement the "divide and conquer" principle and to combine it with techniques just presented.







PF_RING Kernel Plugins [1/3]

- Implementing into the kernel is usually more efficient than doing the same in userland because:
 - Packets do not need to travel from kernel to userland.
 - If a packet is supposed to be received by multiple applications it is not duplicated on the various sockets, but processed once into the kernel
- For packet filtering, it is important to filter as low as possible in the networking stack, as this prevents packet not matching the filter to be propagated and the discarded later on.
- PF_RING plugins allow developers to code small software modules that are executed by PF_RING when incoming packets are received.
- Plugins can be loaded and unloaded dynamically via insmod/ rmmod commands.







PF_RING Kernel Plugins [2/3]

• Each plugin need to declare a data structure according to the format below.

```
struct pfring plugin registration {
 u int16 t plugin id;
                       /* Unique plugin name (e.g. sip, udp) */
 char name[16];
 char description[64];
                       /* Short plugin description */
                      pfring_plugin_filter skb; /* Filter skb: 1=match, 0=no match */
 plugin filter skb
 plugin handle skb
                      pfring plugin handle skb;
                      pfring plugin get stats;
 plugin get stats
 plugin free ring mem pfring plugin free ring mem;
 plugin add rule pfring plugin add rule;
 plugin register pfring plugin register;
 kernel packet start pfring packet start;
 kernel packet reader pfring packet reader;
 kernel packet term pfring packet term;
};
```

- The various pfring_plugin_* variables are pointers to functions that are called by PF_RING when:
 - A packet has to be filtered.
 - An incoming packet has been received and needs to be processed.
 - A userland application wants to know stats about this plugin.
 - A filtering rule will be removed and the memory allocated by the plugin needs to be released.







PF_RING Kernel Plugins [3/3]

- Plugins are associated with filtering rules and are triggered whenever a packet matches the rule.
- If the plugin has a filter function, the this function is called in order to check whether a packet passing the header filter will also pass other criteria. For instance:
 - 'tcp and port 80' is a rule filter used to select http traffic
 - The HTTP plugin can check the packet payload (via DPI) to verify that the packet is really http and it's not another protocol that hides itself on the http port.
- In order to perform complex checks, rules need to be stateful hence to allocate some memory, private to the plugin, that is used to keep the state.
- PF_RING delegates to the plugin the duty of managing this opaque memory that is released by PF_RING when a rule is deleted, by calling the plugin callback.







Efficient Layer 7 Packet Analysis





mo



Using PF_RING Filters: HTTP Monitoring [1/5]

- Goal
 - Passively produce HTTP traffic logs similar to those produced by Apache/Squid/W3C.
- Solution
 - Implement plugin that filters HTTP traffic.
 - Forward to userspace only those packets containing HTTP requests for all known methods (e.g. GET, POST, HEAD) and responses (e.g. HTTP 200 OK).
 - All other HTTP packets beside those listed above are filtered and not returned to userspace.
 - HTTP response length is computed based on the "Content-Length" HTTP response header attribute.







Using PF_RING Filters: HTTP Monitoring [2/5]

Plugin Registration

```
static int __init http_plugin_init(void)
{
    int rc;
    printk("Welcome to HTTP plugin for PF_RING\n");
    reg.plugin_id = HTTP_PLUGIN_ID;
    reg.pfring_plugin_filter_skb = http_plugin_plugin_filter_skb;
    reg.pfring_plugin_handle_skb = NULL;
    reg.pfring_plugin_get_stats = NULL;
    snprintf(reg.name, sizeof(reg.name)-1, "http");
    snprintf(reg.description, sizeof(reg.description)-1, "HTTP protocol analyzer");
    rc = do_register_pfring_plugin(&reg);
    printk("HTTP plugin registered [id=%d][rc=%d]\n", reg.plugin_id, rc);
    return(0);
}
```







Using PF_RING Filters: HTTP Monitoring [3/5]

Plugin Packet Filtering

```
static int http plugin plugin filter skb(filtering rule element *rule,
     struct pfring pkthdr *hdr, struct sk buff *skb,
     struct parse buffer **parse memory)
  struct http filter *rule filter = (struct http filter*)rule-
>rule.extended fields.filter plugin data;
  struct http parse *packet parsed filter;
  if((*parse memory) == NULL) {
   /* Allocate (contiguous) parsing information memory */
    (*parse memory) = kmalloc(sizeof(struct parse buffer*), GFP KERNEL);
   if(*parse memory) {
      (*parse memory) ->mem len = sizeof(struct http parse);
      (*parse memory) ->mem = kcalloc(1, (*parse memory) ->mem len, GFP KERNEL);
      if((*parse memory)->mem == NULL) return(0); /* no match */
   packet parsed filter = (struct http parse*)((*parse memory)->mem);
   parse http packet(packet parsed filter, hdr, skb);
  } else {
   /* Packet already parsed: multiple HTTP rules, parse packet once */
   packet parsed filter = (struct http parse*)((*parse memory)->mem);
  }
```

return((rule_filter->the_method & packet_parsed_filter->the_method) ? 1 /* match */ : 0);







Using PF_RING Filters: HTTP Monitoring [4/5]

Plugin Packet Parsing

```
else if((hdr->caplen > offset) && !memcmp(payload, "HEAD", 4)) packet_parsed->the_method = method_head;
else if((hdr->caplen > offset) && !memcmp(payload, "POST", 4)) packet_parsed->the_method = method_post;
else if((hdr->caplen > offset) && !memcmp(payload, "PUT", 3)) packet_parsed->the_method = method_put;
else if((hdr->caplen > offset) && !memcmp(payload, "DELETE", 6)) packet_parsed->the_method = method_delete;
else if((hdr->caplen > offset) && !memcmp(payload, "TRACE", 5)) packet_parsed->the_method = method_trace;
else if((hdr->caplen > offset) && !memcmp(payload, "CONNECT", 7)) packet_parsed->the_method = method_connect;
else if((hdr->caplen > offset) && !memcmp(payload, "HTTP ", 4)) packet_parsed->the_method = method_connect;
method_http_status_code;
```

```
else packet_parsed->the_method = method_other;
```

}







Using PF_RING Filters: HTTP Monitoring [5/5]

Userland application

```
if((pd = pfring open(device, promisc, 0)) == NULL) { printf("pfring open error\n"); return(-1); }
pfring toggle filtering policy(pd, 0); /* Default to drop */
memset(&rule, 0, sizeof(rule));
rule.rule id = 5, rule.rule action = forward packet and stop rule evaluation;
rule.core fields.proto = 6 /* tcp */;
rule.core fields.port low = 80, rule.core fields.port high = 80;
rule.plugin_action.plugin_id = HTTP_PLUGIN ID; /* HTTP plugin */
rule.extended fields.filter plugin id = HTTP PLUGIN ID; /* Enable packet parsing/filtering */
filter = (struct http filter*)rule.extended fields.filter plugin data;
filter->the method = method get | method http status code;
 if(pfring add filtering rule(pd, &rule) < 0) {</pre>
   printf("pfring_add filtering rule() failed\n");
   return(-1); }
 while(1) {
  u char buffer[2048];
  struct pfring pkthdr hdr;
  if(pfring recv(pd, (char*)buffer, sizeof(buffer), &hdr, 1) > 0)
    dummyProcesssPacket(&hdr, buffer);
}
pfring close(pd);
```





YouTube Monitoring [1/2]

- YouTube monitoring is an extension of the HTTP plugin.
- HTTP is used by YouTube to transport videos usually encoded in H.264 or Flash Video.
- The HTTP plugin can be used for monitoring, from the network point of view, the YouTube traffic and detecting whether the network quality is adequate or if the user should have experienced unstable playback.
- Video streams are tracked by checking the URL (e.g. GET /get_video? video_id=...) and the server host (www.youtube.com).
- Whenever a YouTube video stream is detected, the HTTP plugin adds an exact matching rule on the hash, used to track the stream, with the YouTube plugin specified as rule action.







YouTube Monitoring [2/2]

• The YouTube plugin is able to measure some stream statistics such as throughput, jitter, bandwidth used.

```
struct youtube_http_stats {
    u_int32_t initialTimestamp, lastTimestamp, lastSample; /* Packet Timestamps [jiffies] */
    struct timeval initial_tv;
    u_int32_t tot_pkts, tot_bytes, cur_bytes;
    u_int32_t num_samples;
    u_int8_t signaling_stream; /* 1=signaling, 2=real video stream */
    char url[URL_LEN];
    char video_id[VIDEO_ID_LEN], video_playback_id[VIDEO_ID_LEN];
    u_int32_t min_thpt, avg_thpt, max_thpt; /* bps */
    u_int32_t min_jitter, avg_jitter, max_jitter; /* jiffies */
    u_int32_t duration_ms;
    char content_type[CONTENT_TYPE_LEN];
    u_int32_t tot_jitter, num_jitter_samples;
};
```

- When a stream is over, the plugin return to userland a packet with the stream statistics.
- Note that all stream packets are not returned to userland, but just the statistics, that contributes to reduce load on the probe and improve performance.







Dynamic PF_RING Filtering: VoIP [1/6]

• Goal

- Track VoIP (SIP+RTP) calls at any rate on a Gbit link using commodity hardware.
- Track RTP streams and calculate call quality information such as jitter, packet loss, without having to handle packets in userland.

• Solution

- Code a PF_RING plugin for tracking SIP methods and filter-out:
 - Uninteresting (e.g. SIP Options) SIP methods
 - Not well-formed SIP packets
 - Dummy/self calls (i.e. calls used to keep the line open but that do not result in a real call).
- Code a RTP plugin for computing in-kernel call statistics (no pkt forwarding).
- The SIP plugin adds/removes a precise RTP PF_RING filtering rule whenever a call starts/ends.







Dynamic PF_RING Filtering: VoIP [2/6]

- Before removing the RTP rule though PF_RING library calls, call information is read and then the rule is deleted.
- Keeping the call state in userland and do not forwarding RTP packets, allows the code of VoIP monitoring applications to be greatly simplified.
- Furthermore as SIP packets are very few compared to RTP packets, the outcome is that most packets are not forwarded to userland contributing to reduce the overall system load.









Dynamic PF_RING Filtering: VoIP [3/6]

- SIP Plugin
 - It allows to set filters based on SIP fields (e.g. From, To, Via, CalIID)
 - Some fields are not parsed but the plugin returns an offset inside the SIP packet (e.g. SDP offset, used to find out the IP:port that will be used for carrying the RTP/RTCP streams).
 - Forwarded packets contain parsing information in addition to SIP payload.
- RTP Plugin
 - It tracks RTP (mono/by-directional) flows.
 - The following, per-flow, statistics are computed: jitter, packet loss, malformed packets, out of order, transit time, max packet delta.
 - Developers can decide not to forward packets (this is the default behavior) or to forward them (usually not needed unless activities like lawful interception need to be carried on).







Dynamic PF_RING Filtering: VoIP [4/6]

- Validation
 - A SIP test tool and traffic generator (sipp) is used to create synthetic SIP/RTP traffic.
 - A test application has been developed: it receives SIP packets (signaling) and based on them it computes RTP stats.
 - A traffic generator (IXIA 400) is used to generate noise in the line and fill it up. As RTP packets are 100 bytes in average, all tests are run with 128 bytes packets.
 - The test code runs on a cheap single-core Celeron 3.2 GHz (cost < 40 Euro).
 - In order to evaluate the speed gain due to PF_RING kernel modules, the same test application code is tested:
 - Forwarding SIP/RTP packets to userland without exploiting kernel modules (i.e. the code uses the standard PF_RING).
 - RTP packets are not forwarded, SIP packets are parsed/filtered in kernel.







Dynamic PF_RING Filtering: VoIP [5/6]

% Idle CPU [128 bytes packets]



Dynamic PF_RING Filtering: VoIP [6/6]

- Validation Evaluation
 - In-kernel acceleration has demonstrated that until 40K rules, kernel plugins are faster than a dummy application that simply captures packets without any processing.
 - On a Gbit link it is possible to have up to ~10K concurrent calls with G.711 (872 Mbit) or ~30K calls with G.729 (936 Mbit). This means that with the current setup and a slow processor, it is basically possible to monitor a medium/large ISP.
- Future Work Items
 - The plugins are currently used as building blocks glued together by means of the user-space applications.
 - The SIP plugin can dynamically add/remove RTP rules, so that it is possible to avoid (even for SIP) packet forwarding and send to userland just VoIP statistics for even better performance figures.







PF_RING Content Inspection

- PF_RING allows filtering to be combined with packet inspection.
- Ability to (in kernel) search multiple string patterns into packet payload.
- Algorithm based on Aho-Corasick work.
- Ideal for fields like lawful interception and security (including IDSs).
- Major performance improvement with respect to conventional pcap-based applications.







L7 Analysis: Summary

- The use of kernel plugins allows packets to have a short journey towards the application.
- In-kernel processing is very efficient and it avoids the bottleneck of several userland application threads competing for packets.
- As PF_RING requires minimal locking (when the filtering rule is accessed and updated), packets are processed concurrently without any intervention from userland applications.
- As the Linux kernel concurrently fetches packets from adapters, this is a simple way to exploit multi-processing/core without having to code specific (multithreaded) userland applications and serialize packets on (PF_RING) sockets.







Direct Access to NICs





2nd TMA PhD School - June 2011



Direct NIC Access: Introduction

- Commercial accelerated NICs are accelerated either using ASIC (rare) or FPGAs (often) chips.
- Accelerators improve common activities such as packet filtering and are also responsible of pushing packets to memory with very limited (< 1%) load on the main CPU.
- Applications access packets directly without any kernel intervention at all.
- A kernel-mapped DMA memory allows the application to manipulate card registers and to read packets from this memory where incoming packets are copied by the hardware accelerators.
- Cards falling in this category include:
 - Endace DAG
 - Napatech
 - NetFPGA







Direct NIC Access: Comparison [1/2]





Hardware Acceleration



UNIVERSITÀ DI PISA

ntop

Direct NIC Access: Comparison [2/2]

- The reason why accelerated cards are so efficient are:
 - The FPGA polls packets as fast as possible without any intervention from the main CPU. In Linux the main CPU has to periodically read packets through NAPI from the NIC.
 - Received packets are copied on a pre-allocated large memory buffer so no per-packet allocation/deallocation is necessary at all, as it happens in vanilla Linux.
 - Similar to PF_RING, packets are read from circular buffer without any kernel interaction (beside packet polling).
- Limitations
 - As applications access packets directly, if they improperly manipulate card's memory the whole system might crash.
 - FPGA filtering is very limited and not as rich as PF_RING.
 - Contrary to PF_RING, only one application at time can read packers from the ring.







Welcome to nCap (Circa 2003)









nCap Features

	Packet Capture Acceleration	Wire Speed Packet Capture	Number of Applications per Adapter
Standard TCP/IP Stack with accelerated driver	Limited	No	Unlimited
PF_RING with accelerated driver	Great	Almost	Unlimited
Straight Capture	Extreme	Yes	One







nCap Internals

- nCap maps at userland the card registers and memory.
- The card is accessed by means of a device /dev/ncap/ethX
- If the device is closed it behaves as a "normal" NIC.
- When the device is open, it is completely controlled by userland the application.
- A packet is sent by copying it to the TX ring.
- A packet is received by reading it from the RX ring.
- Interrupts are disabled unless the userland application wait for packets (poll()).
- On NIC packet filtering (MAC Address/VLAN only).







nCap Comparison (1 Gbit)

	Maximum	Estimated	Manufacturer
	Packet Loss	Card	
	at Wire Speed	Price	
DAG	0%	> 5-7 K Euro	Endace.com
nCap	0.8%	100 Euro	
Combo 6 (Xilinx)	5%	> 7-10 K Euro	Liberouter.com

Source Cesnet (http://luca.ntop.org/ncap-evaluation.pdf)






Beyond PF_RING

- PF_RING has shown to be an excellent packet capture acceleration technology compared to vanilla Linux.
- It has reduced the cost of packet capture and forward to userland.
- However it has some design limitations as it requires two actors for capturing packets that result in sub-optimal performance:
 - kernel: copy packet from NIC to ring.
 - userland: read packet from ring and process it.
- PF_RING kernel modules demonstrated that limiting packet processing in user-space by moving it to kernel results in major performance improvements.
- A possible solution is to map a NIC to user-space and prevent the kernel from using it.







PF_RING DNA (Direct NIC Access)

- PF_RING DNA is an extension for PF_RING that allows NICs to be accessed in direct mode fully bypassing Linux NAPI.
- Based on the lessons learnt while developing nCap, DNA is a technology developed in clean-room that has been designed to be NIC-neutral in order to allows various NICs to be supported.
- The NIC mapping is driver dependent hence it requires some driver modifications in order to:
 - Disable NAPI when the NIC is accessed in DNA mode.
 - Contiguously allocate RX card's memory in one shot (and not per packet).
 - Register the NIC with PF_RING so the card is accessed seamlessly from PF_RING applications without the need to know the NIC internals and its memory layout.







PF_RING DNA (De)Registration

```
/* Register with PF RING */
            do ring dna device handler(add device mapping,
                adapter->tnapi.dma mem.packet_memory,
                adapter->tnapi.dma mem.packet num slots,
                                                                  NIC Memory
                adapter->tnapi.dma mem.packet slot len,
                                                                   Pointers
                adapter->tnapi.dma mem.tot packet memory,
                rx ring->desc,
                rx ring->count, /* # of items */
NIC DMA Ring
                sizeof(struct e1000 rx desc),
                rx ring->size, /* tot len (bytes) */
                0, /* Channel Id */
                (void*)netdev->mem start,
                                                                 NIC Registers
                netdev->mem end,
                                                                   Memory
                netdev,
                intel e1000,
                &adapter->tnapi.packet waitqueue,
Packet Polling
                &adapter->tnapi.interrupt received,
                (void*)adapter,
                wait packet function ptr);
```





PF_RING DNA: Current Status

- As of today, DNA is available for Intel-based 1 Gbit (e1000 driver) and 10 Gbit (ixgbe) NICs.
- Any modern dual-core (or better) system can achieve wire rate packet capture at any packet size using DNA.
- A userland library used to manipulate card registers has been integrated into PF_RING.
- Applications do not need to do anything different from standard PF_RING with the exception that the ring memory has to be open using pfring_open_dna() instead of the standard pfring_open().
- When an application opens the adapter in DNA mode, other applications using the same adapter in non-DNA mode will stop receiving packets until the application quits.







Towards 10 Gbit Packet Capture Using Commodity Hardware







Enhanced NIC Drivers [1/5]

- The current trend in computer architecture is towards multi-core systems.
- Currently 4-core CPUs are relatively cheap, some manufacturers announced a 64-core x86 CPU by the end of 2008.
- Exploiting multi-core in userland applications is relatively simple by using threads.
- Exploiting multi-core in kernel networking code is much more complex.
- Linux kernel networking drivers are single-threaded and the model is still the same since many years.
- It's not possible to achieve good networking performance unless NIC drivers are also accelerated and exploit multi-core.







Enhanced NIC Drivers [1/4]

- The current trend in computer architecture is towards multi-core systems.
- Currently 4-core CPUs are relatively cheap and rather common on the market. Intel announced Xeon Nehalem-EX with 16 threads (8 cores) for late 2009. The core rush is not yet over.
- Exploiting multi-core in userland applications is relatively simple by using threads.
- Exploiting multi-core in kernel networking code is much more complex.
- Linux kernel networking drivers are single-threaded and the model is still the same since many years.
- It's not possible to achieve good networking performance unless NIC drivers are also accelerated and exploit multi-core.







Enhanced NIC Drivers [2/4]

Intel has recently introduced a few innovations in the Xeon 5000 chipset series that have been designed to accelerate networking applications:

- I/O Acceleration Technology (I/OAT)
 - Direct Cache Access (DCA) asynchronously move packets from NIC directly on CPU's cache in DMA.
 - Multiple TX/RX queues (one per core) that improve system throughput and utilization.
 - MSI-X, low latency interrupts and load balancing across multiple RX queues.
 - RSS (Receive-Side Scaling) balances (network flow affinity) packets across RX queue/ cores.
 - Low-latency with adaptive and flexible interrupt moderation.

In a nutshell: increase performance by distributing workloads across available CPU cores.







Enhanced NIC Drivers [3/4]









Enhanced NIC Drivers: Linux NAPI [4/4]









Linux NAPI Limitations [1/2]









Linux NAPI Limitations [2/2]

- Multiple-RX queues are not fully exploited by Linux as NAPI polls them in sequence and not concurrently
- Interrupts are enabled/disabled globally (i.e. for all queues at the same time) whereas they should be managed queue-per-queue as not all queues have the same amount of traffic (it depends on how balance-able is the ingress traffic).
- Original queue index (that can be used as flow identifier) is lost when the packet is propagated inside the kernel and then to userland.
- Userland applications see the NIC as a single entity and not as a collection of queues as it should be. This is a problem as the software could take advantage of multiple queues by avoiding threads competing for incoming packets all coming from the same NIC but from different queues.







Example of Multi-Queue NIC Statistics

```
# ethtool -S eth5
NIC statistics:
     rx packets: 161216
     tx packets: 0
     rx bytes: 11606251
     tx bytes: 0
     lsc int: 1
     tx busy: 0
     non eop descs: 0
     rx errors: 0
     tx errors: 0
     rx dropped: 0
     tx dropped: 0
     multicast: 4
     broadcast: 1
     rx no buffer count: 2
     collisions: 0
     rx over errors: 0
     rx crc errors: 0
     rx frame errors: 0
     rx fifo errors: 0
     rx missed errors: 0
     tx aborted errors: 0
     tx carrier errors: 0
     tx fifo errors: 0
```

tx heartbeat errors: 0 tx timeout count: 0 tx restart queue: 0 rx long length errors: 0 rx short length errors: 0 tx tcp4 seg ctxt: 0 tx tcp6 seg ctxt: 0 tx flow control xon: 0 rx flow control xon: 0 tx flow control xoff: 0 rx flow control xoff: 0 rx csum offload good: 153902 rx csum offload errors: 79 tx csum offload ctxt: 0 rx header split: 73914 low latency interrupt: 0 alloc rx page failed: 0 alloc rx buff failed: 0 lro flushed: $\overline{0}$ lro coal: 0 tx queue 0 packets: 0 tx queue 0 bytes: 0 rx queue 0 packets: 79589 rx queue 0 bytes: 5721731 rx queue 1 packets: 81627 rx queue 1 bytes: 5884520







Memory Allocation Life Cycle [1/5]

• Incoming packets are stored into kernel's memory that has been previously allocated by the driver.



- As soon that a packet is received, the NIC NPU (Network Process Unit) checks if there's an empty slot and if so it copies the packet in the slot.
- The slot is removed from the RX buffer and propagated through the kernel.
- A new bucket is allocated and places on the same position of the old slot.







Memory Allocation Life Cycle [2/5]

- The consequence of this allocation policy is that:
 - Every new packet requires one slow allocation (and later-on a free).
 - As the traffic rate increases, increasing allocations/free will happen.
 - In particular at 10 Gbit, if there's a traffic spike or a traffic shot, the system may run out of memory as incoming packets:
 - require memory hence the memory allocator does its best to allocate new slots.
 - are stuck in the network kernel queue because the packet consumers cannot keep-up with the ingress traffic rate.
 - When the system runs in low memory it tries to free cached memory in order to free some space.
 - Unfortunately when the ingress rate is very high, the memory recover process does not have enough time hence the system runs out of memory and the result is that Linux's OOM (Out Of Memory) killer has to kill some processes in order to recover some memory.







Memory Allocation Life Cycle [3/5]

```
if(rx desc->status & E1000 RXD STAT DD) {
       /* A packet has been received */
#if defined (CONFIG RING) || defined(CONFIG RING MODULE)
          handle ring skb ring handler = get skb ring handler();
          if (ring handler && adapter->soncap.soncap enabled) {
            ring handler(skb, 0, 1, (hash value % MAX NUM CHANNELS));
          } else {
#endif
          [....]
              if (++i == rx ring -> count) i = 0;
              next rxd = E1\overline{0}00 RX DESC(*rx ring, i);
              prefetch(next rxd);
              next buffer = &rx ring->buffer info[i];
              cleaned = TRUE;
              cleaned count++;
              pci unmap single(pdev, buffer info->dma, PAGE SIZE, PCI DMA FROMDEVICE);
              [....]
              skb = netdev alloc skb(netdev, bufsz);
              buffer info->dma = pci map single(pdev,
                                                  skb->data,
                                                  adapter->rx buffer len,
                                                  PCI DMA FROMDEVICE);
```

[....]







Memory Allocation Life Cycle [4/5]







2nd TMA PhD School - June 2011



Memory Allocation Life Cycle [5/5]

- Avoiding memory allocation/deallocation has several advantages:
 - No need to allocate/free buffers
 - No need to map memory though the PCI bus
 - In case of too much incoming traffic, as the kernel has more priority than userland applications, there's no risk to run out of memory as it happens with standard NAPI.
- The last advantage of doing a packet copy to the PF_RING buffer is the speed. Depending on the setup, the packet capture performance can be increased of 10-20% with respect to standard NAPI.







Enhanced NIC Drivers: TNAPI [1/8]

- In order to enhance and accelerate packet capture under Linux, a new Linux driver for Intel 1 and 10 Gbit cards has been developed. Main features are:
 - Multithreaded capture (one thread per RX queue, per NIC adapter).
 The number of rings is the number of cores (i.e. a 4 core system has 4 RX rings)
 - RX packet balancing across cores based on RSS: one core, one RX ring.
 - Driver-based packet filtering (PF_RING filters port into the driver) for stopping unwanted packets at the source.
 - Development drivers for Intel 82598/9 (10G) and 82575/6 (1G) ethernet controllers.
- For this reason the driver has been called TNAPI (Threaded NAPI).







Enhanced NIC Drivers: TNAPI [2/8]











Enhanced NIC Drivers: TNAPI [3/8]

- Packet capture has been greatly accelerated thanks to TNAPI as:
 - Each RX queue is finally independent (interrupts are turned on/off per queue and not per card)
 - Each RX queue has a thread associated and mapped on the same CPU core as the one used for RSS (i.e. cache is not invalidated)
 - The kernel thread pushes packets as fast as possible up on the networking stack.
 - Packets are copied from the NIC directly into PF_RING (allocation/ deallocation of skbuffers is avoided).
 - Userland applications can capture packets from a virtual ethernet NIC that maps the RX ring directly into userspace via PF_RING.







Enhanced NIC Drivers: TNAPI [4/8]

- TNAPI Issues: CPU Monopolization
 - As the thread pushes packets onto PF_RING, it should be avoided that this thread monopolizes. This is because of the all CPU is used by the kernel for receiving packets, then packet loss won't happen in kernel but in userspace (i.e. the packet loss problem is not solved, but just moved).

```
while(<polling packets from RX queue X>) {
    /* Avoid CPU monopolization */
    if(rx_budget > 0)
        rx_budget--;
    else {
        rx_budget = DEFAULT_RX_BUDGET;
        yield();
    }
}
```

 Solution: every X polling cycles, the thread has to give away some CPU cycles. This is implemented as follow rx_budget that's consumed whenever a packet is received and sent to PF_RING.





Enhanced NIC Drivers: TNAPI [5/8]

- TNAPI Issues: Interrupts and Cores Allocation
 - RX ring interrupts must be sent to the right core that's manipulating the queue in order to preserve cache coherency.
 - The userland application that's fetching packets from queue X, should also be mapped to core X.
 - As interrupts are now sent per-queue (and not per-nic as it used to be) we must make sure that they are sent to the same core that's fetching packets.

# cat ,	/proc/interr	rupts								
	CPUO	CPU1	CPU2	CPU3	CPU4	CPU5	CPU6	CPU7		
191:	1	1	2656	1	2	2	1	2	PCI-MSI-edge	eth3
192:	1	4	0	0	2655	3	1	2	PCI-MSI-edge	eth2
193:	78634	14	7	13	9	13	13	18	PCI-MSI-edge	eth1
194:	3	15964	6	3	3	5	3	5	PCI-MSI-edge	eth0
195:	0	0	0	0	0	0	0	0	PCI-MSI-edge	eth7:lsc
196:	1	2	2	0	0	2658	1	0	PCI-MSI-edge	eth7:v8-Tx
197:	1	0	2	0	1	0	1	5309	PCI-MSI-edge	eth7:v7-Rx
198:	1	0	0	5309	1	2	0	1	PCI-MSI-edge	eth7:v6-Rx
199:	0	1	0	1	0	1	2	5309	PCI-MSI-edge	eth7:v5-Rx
200:	0	1	1	5307	2	2	1	0	PCI-MSI-edge	eth7:v4-Rx
201:	1	0	1	2	1	5307	2	0	PCI-MSI-edge	eth7:v3-Rx
202:	2	2	0	1	1	0	5307	1	PCI-MSI-edge	eth7:v2-Rx
203:	0	1	5309	1	1	1	0	1	PCI-MSI-edge	eth7:v1-Rx
204:	2	2	1	0	5307	1	1	0	PCI-MSI-edge	eth7:v0-Rx







Enhanced NIC Drivers: TNAPI [6/8]

• Example:

-RX ring 6 and 4 use the same CPU 3.

-We want to move RX ring 6 to CPU 1

# cat /	proc/interru	ıpts								
	CPU0	CPU1	CPU2	CPU3	CPU4	CPU5	CPU6	CPU7		
198:	1	0	0	5309	1	2	0	1	PCI-MSI-edge	eth7:v6-Rx
200:	0	1	1	5307	2	2	1	0	PCI-MSI-edge	eth7:v4-Rx
<pre># cat / 0000000 # echo # cat / 0000000</pre>	proc/irq/198 8 2 > /proc/ir proc/irq/198 2	8/smp_affin: cq/198/smp_a 8/smp_affin:	ity affinity [0 ity	0000010 whe:	re 1 = CPU 3	L]				
# cat /	proc/interru	upts grep '	"eth7:v6-Rx	11						
100.	0	67	0	5309	1	2	0	1	DCT MCT odao	
198:	0	07	0	5505	Ţ	2	0	T	PCI-MSI-edge	eth7:v6-Rx
198:	0	07	U	5505	Ţ	2	0	1	rci-Msi-eage	eth7:v6-Rx

```
unsigned long mask = 7; /* processors 0, 1, and 2 */
unsigned int len = sizeof(mask);
if (sched_setaffinity(0, len, &mask) < 0) {
    perror("sched_setaffinity");
}</pre>
```

-How to map a process to a CPU/core







Enhanced NIC Drivers: TNAPI [7/8]

Test Type	Max Packet Capture Speed				
PF_RING	300K pps	560K pps			
PF_RING+TNAPI Mono RX queue	750K pps	920K pps			
PF_RING+TNAPI Multi RX queue	860K pps	Wire Rate (1 Gbit) ~ 3 Million pps (10 Gbit) ~ 5 Million pps (10 Gbit - 2 x Xeon)			
	Intel Core2Duo 1.86 GHz (Dual Core) No Intel I/OAT	CPU Intel Xeon 2.4 GHz (Quad Core) Intel 5000 chipset (I/OAT support)			





UNIVERSITÀ DI PISA

Enhanced NIC Drivers: TNAPI [8/8]



Testbed: Xeon X3450 @ 2.67GHz







RX Multi-Queue and DNA

- As previously explained, DNA is an excellent technology for those application developers who need wire speed packet capture, but that do not need features such as:
 - packet filtering
 - multiple application packet consumers.
- DNA so far has been ported to the Intel mono-queue 1 Gbit driver (e1000) and multi-queue 10 Gbit driver (ixgbe).
- Currently the port of DNA to 1 Gbit RX multi-queue driver (igb) is in progress and it will be available later this year.
- Combining DNA with multi-queue allows applications to be split into concurrent execution threads that enables multicore architectures to be further exploited.
- Additionally by exploiting hardware traffic balancing, it allows flow-based applications such as netflow probes, to be further accelerated.







Multi-Queue on Accelerated NICs









Exploiting PF_RING Multi-Queue: nProbe



- Packet balancing across cores.
- Peak nProbe performance: 1.48 Mpps (packet rate) x 2 Cores.







Strong Multicore NICs: Tilera Tile64









Towards Strong Multicore [1/2]

General Perception is that people usually think that multicore is a good idea, although difficult to implement.

- General PC market
 - Input data is unstructured, sequential
 - Billions of lines of sequential applications
 - Hard to migrate it to parallel code
- Embedded market
 - Data is inherently parallel
 - Engineers have designed parallel applications
 - Their main challenge is complexity of design







Towards Strong Multicore [2/2]

• Some applications are naturally parallel as in networking where a network pipe is a multiplex of many "flows" or distinct streams.



- The only barriers towards adopting strong multicore are:
 - Design the application program so that it can take advantage of multicore without sequentially performing activities that could be carried on in parallel.
 - Entry ticket for learning multicore development tools.
 - Low-level programming required to take advantage of the technology.







Programming Paradigms

- Run to completion model
 - Sequential C/C++ applications
 - Run multiple application instances one/core
 - Use load balancer library for distribution
 - Use tools to tune performance
- Parallel programming
 - Parallelize application with pthreads shared memory
 - Run on multiple cores

UNIVERSITÀ DI PIS.

- Use communication libraries to optimize
- Use tools to tune performance







105

Parallel Processing Without Parallel Programming

- Standard model in the embedded world
 - Facilitates immediate results using off-the-self code
- Simple architecture
 - Each core runs complete application and handles one or multiple flows or channels
 - I/O management and load distribution
 - Most embedded applications fit this category
 - Large numbers of flows, frames channels, streams, etc...





UNIVERSITÀ DI PISA



Tilera TILExpress64

- 64-core CPU.
- Linux-based 2.6 operating system running on board.
- Programmable in C/C++.
- Eclipse Integration for easing software development and debugging.







2nd TMA PhD School - June 2011



TILE64 Architecture [1/2]








TILE64 Architecture [2/2]



Each tile is a complete processor

2 Dimensional iMesh connects tiles



38 terabits of on-chip bandwidth







Tilera Advantages

- No need to capture packets as it happens with PCs.
- 12 x 1 Gbit, or 6 x 1Gbit and 1 x 10 Gbit Interfaces (XAUI connector).
- Ability to boot from flash for creating stand-alone products.
- Standard Linux development tools available including libpcap for packet capture.
- Application porting is very quick and simple: less than 100 lines of code changed in nProbe.







Porting Exiting Applications to Tile64: nProbe









nProbe Performance on Tile64









Final Remarks









Programming for Multicore [1/4]

- Multicore is not the solution to all performance and scalability problems.
- Actually it can decrease the performance of poorly designed applications.
- Like it or not, multicore is the future of CPUs, and programmers have to face with it.
- From author's experience before adding threads and semaphores to parallelize an existing program, it's worth to think if instead the basic algorithm used are compatible with multicore.







Programming for Multicore [2/4]

- When multiple cores are used, efficient memory caching is the way to improve application performance.
- Hardware CPU caches are rather sophisticated, however they cannot work optimally without programmer's assistance.
- Cache coherence can be rather costly if programs invalidate it when not necessary.
- False sharing (when a system participant attempts to periodically access data that will never be altered by another party, but that data shares a cache block with data that is altered, the caching protocol may force the first participant to reload the whole unit despite a lack of logical necessity) is just an example of performance degrading due to poor programming.
- Reference
 - U. Drepped, What Every Programmer Should Know About Memory, http://people.redhat.com/drepper/cpumemory.pdf, RedHat 2007.





Programming for Multicore [3/4]





UNIVERSITÀ DI PISA

open source

Programming for Multicore [4/4]



Great Application DesignExploit Native MulticoreFully Lockless Hash

Lockeless hashes:

http://video.google.com/videoplay?docid=2139967204534450862





2nd TMA PhD School - June 2011



Memory Allocation [1/2]

Limit Memory Allocation (if not necessary)

- Multithreaded programs often do not scale because the heap is a bottleneck.
- When multiple threads simultaneously allocate or deallocate memory from the allocator, the allocator will serialize them.
- Programs making intensive use of the allocator actually slow down as the number of processors increases.







Memory Allocation [2/2]

- Programs should avoid, if possible, allocating/deallocations memory too often and in particular whenever a packet is received.
- In the Linux kernel there are available kernel/driver patches for recycling skbuff (kernel memory used to store incoming/outgoing packets).
- Using PF_RING (into the driver) for copying packets from the NIC to the circular buffer without any memory allocation increases the capture performance (around 10%) and reduces congestion issues.

References:

- A Comparison of Memory Allocators
 <u>http://developers.sun.com/solaris/articles/multiproc/multiproc.html</u>
- The Hoard Memory Allocator http://www.hoard.org/







PF_RING on VMs [1/4]



Open Issues

- Long packet journey from NIC to the VM.
- Various packet copies are involved.
- Packets replicated on all VMs.
- Overhead due to the abstraction level.

Goal

• Implement straight capture to the VM.







PF_RING on VMs [2/4]



vNPlug

- It implements a shared memory area (host <->VM) that is mapped as a dummy PCI device
- Host->Guest signaling by emulating interrupts.
- Guest->Host signaling by writing on PCI registers that are monitored via ioeventfd.





11(0)

PF_RING on VMs [3/4]



UNIVERSITÀ DI PISA

122

open source

PF_RING on VMs [4/4]



UNIVERSITÀ DI PISA

123

References

- http://www.ntop.org/
- http://www.intel.com/cd/network/connectivity/emea/eng/226275.htm
- http://www.tilera.com

Email: Luca Deri <<u>deri@ntop.org</u>>







High-Speed Traffic Capture and Analysis Using Open-Source Software and Commodity Hardware

Part 2: Traffic Monitoring

Luca Deri <<u>deri@ntop.org</u>>







Monitoring Goals

- Analysis of LAN and WAN Traffic
- Unaggregated raw data storage for the near past (-3 days) and long-term data aggregation on selected network traffic metrics (limit: available disk space)
- Data navigation by means of a web 2.0 GUI
- Geolocation of network flows and their aggregation based on their geographical source.
- Integration with routing information in order to provide accurate traffic path analysis.







Traffic Collection Architecture [1/2]

- Available Options
 - 1.Exploit network equipment (routers and switches)
 - -Advantages:
 - Maximize investment.
 - Avoid adding extra network equipment/complexity in the network.
 - •No additional point of Failure
 - Disadvantages:
 - •Often is necessary to buy costly netflow engines
 - •Have to survive with bugs (e.g. Juniper have issues with AS information)







Traffic Collection Architecture [2/2]

- 2.Custom Network Probes
- Advantages
 - -Ability to avoid limitations of commercial equipment
 - (Often) Faster and more flexible than hw probes









Introduction to Cisco NetFlow

- Flow: "Set of network packets with some properties in common". Typically (IP src/dst, Port src/dst, Proto, TOS, VLAN).
- Network Flows contain:
 - –Peers: flow source and destination.
 - -Counters: packets, bytes, time.
 - Routing information: AS, network mask, interfaces.









Collection Architectures [1/2]



ntop

UNIVERSITÀ DI PISA



Collection Architectures [2/2]









Flow Journey: Creation









Flow Journey: Export

Srclf	SrclPadd	Dstlf	DstlPadd	Protocol	TOS	Flgs	Pkts	Src Port	Src Msk	Src AS	Dst Port	Dst Msk	Dst AS	NextHop	Bytes/ Pkt	Active	Idle
a1/0	173.100.21.2	Fa0/0	10.0.227.12	11	80	10	11000	162	/24	5	163	/24	15	10.0.23.2	1528	1745	4
a1/0	173.100.3.2	Fa0/0	10.0.227.12	6	40	0	2491	15	/26	196	15	/24	15	10.0.23.2	740	41.5	1
a1/0	173.100.20.2	Fa0/0	10.0.227.12	11	80	10	10000	161	/24	180	10	/24	15	10.0.23.2	1428	1145.5	3
a1/0	173.100.6.2	Fa0/0	10.0.227.12	6	40	0	2210	19	/30	180	19	/24	15	10.0.23.2	1040	24.5	14
	2. Flow A	ging T	imers	 Inacti Long Flow 	ive FI Flow ends	ow (* 7 (30) 8 by F	15 sec min (18 8ST or	is de 800 s FIN 1	fault iec) i ICP f) s def ⁼ lag	'ault)						
Srclf	2. Flow A	ging T	imers DstlPadd	Inacti Long Flow	TOS	ow (1 (30) by Figs	15 sec min (18 ST or Pkts	is de 800 s FIN 1 Src Port	fault ec) i CP f) s def Flag Src AS	Dst Port	Dst Msk	Dst	NextHop	Bytes/	Active	Idi
Srclf a1/0	2. Flow A SrcIPadd 173.100.21.2	ging T Dstlf Fa0/0	DstlPadd 10.0227.12	 Inacti Long Flow Protocol 11	TOS 80	ow (1 (30) by F Flgs 10	Pkts 11000	Src Port 00A2	fault ec) i CP f Src Msk /24) s def Elag Src AS 5	Dst Port 00A2	Dst Msk /24	Dst AS 15	NextHop 10.023.2	Bytes/ Pkt 1528	Active 1800	ldi 4







Flow Format: NetFlow v5 vs v9

	v5	v9
Flow Format	Fixed	User Defined
Extensible	No	Yes (Define new FlowSet Fields)
Flow Type	Unidirectional	Bidirectional
Flow Size	48 Bytes (fixed)	It depends on the format
IPv6 Aware	No	IP v4/v6
MPLS/VLAN	No	Yes







Flow Format: NetFlow v9/IPFIX

Header	NetFlow Version 9 Header: 32 bits	→		
First Template FlowSet	Version 9 Count = 4 (FlowS	Sets)		
Template Record	System Uptime			
First Record FlowSet	UNIX Seconds			
(Template ID 256)	Package Sequence			
First Data Record	Source ID			
Second Data Record	A Translate Flow Oct. 10 hits. b			
Third Data Record	<- Template FlowSet: 16 bits->	← Data FlowSet: 32 bits →		
Second Template Flowset	FlowSet ID = 0	FlowSet Length =		
Template Record	Length = 28 bytes	🔻 ID = 256 64 bytes		
Template Record	Template ID = 256	192.168.1.12		
Second Record Flowset	Field Count = 5	*		
(Template ID 257)	IPv4_SRCADDR (0x0008)	10.5.12.254		
Data Record	Length = 4	≭		
Data Record	IPv4_DSTADDR (0x000C)	192.168.1.1		
Data Record	Length = 4			
Data Record	IPv4_NEXT_HOP (0x000E)	5009		
	Length = 4	7		
	PKTS_32 (0x0002)	5344385		
	Length = 4	₮		
	BYTES_32 (0x0001)	192.168.1.27		
	Length = 4	10.5.12.23		
		192.168.1.1		
		748		
		388934		
		192.168.1.56		
		10.5.12.65		
		192.168.1.1		
		5		
		6534		







InMon sFlow



% Sampling Error <= 196 * sqrt(1 / number of samples) [http://www.sflow.org/packetSamplingBasics/]









Integrated Network Monitoring









Traffic Collection: A Real Scenario









Heterogeneous Flow Collection









nProbe: sFlow/NF/IPFIX Probe+Collector



Raw Files / MySQL / SQLite / FastBit







Problem Statement [1/2]

- NetFlow and sFlow are the current state-of-theart standard for network traffic monitoring.
- As the number of generated flows can be quite high, operators often use sampling in order to reduce their number.
- Sampling leads to inaccuracy so it cannot always be used in production networks.
- Thus network operators have to face the problem of collecting and analyzing a large number of flow records.







Problem Statement [2/2]

Where to store collected flows?

- -Relational Databases
 - Pros: Expressiveness of SQL for data search.
 - •Cons: Sacrifice flow collection speed and query response time.
- -Raw Disk Archives
 - Pros: Efficient flow-to-disk collection speed (> 250K flow/s).
 - Cons: Limited query facilities as well search time proportional to the amount of collected data (i.e. no indexing is used).







Towards Column-Oriented Databases [1/3]

- Network flow records are read-only, shouldn't be modified after collection, and several flow fields have very few unique values.
- B-tree/hash indexes used in relational DBs to accelerate queries, encounter performance issues with large tables as:
 - need to be updated whenever a new flow is stored.
 - require a large number of tree-branching operations as they use slow pointer chases in memory and random disk access (seek), thus taking a long time.
- Thus with relational DBs it is not possible to do live flow collection/ import as index update will lead to flow loss.







Towards Column-Oriented Databases [2/3]

- A column-oriented database stores its content by column rather than by row. As each column is stored contiguously, compression ratios are generally better than row-stores because consecutive entries in a column are homogeneous to each other.
- Column-stores are more I/O efficient (than row stores) for readonly queries since they only have to read from disk (or from memory) those attributes accessed by a query.
- Indexes that use bit arrays (called bitmaps) answer queries by performing bitwise logical operations on these bitmaps.






Towards Column-Oriented Databases [3/3]

- Bitmap indexes perform extremely well because the intersection between the search results on each value is a simple AND operation over the resulting bitmaps.
- As column data can be individually sorted, bitmap indexes are also very efficient for range queries (e.g. subnet search) as data is contiguous hence disk seek is reduced.
- As column-oriented databases with bitmap indexes provide better performance compared to relational databases, the authors explored their use in the field of flow monitoring.







nProbe + FastBit

- FastBit is not a database but a C++ library that implements efficient bitmap indexing methods.
- Data is represented as tables with rows and columns.
- A large table may be partitioned into many data partitions and each of them is stored on a distinct directory, with each column stored as a separated file in raw binary form.
- nProbe natively integrates FastBit support and it automatically creates the DB schema according to the flow records template.
- Flows are saved in blocks of 4096 records.
- When a partition is fully dumped, columns to be indexed are first sorted then indexed.







Performance Evaluation: Disk Space

MySQL	No/With Indexes	1.9 / 4.2
FastBit	Daily Partition (no/with Indexes)	1.9 / 3.4
	Hourly Partition (no/with Indexes)	1.9 / 3.9
nfdump	No indexes	1.9

Results are in GB







Performance Evaluation: Query Time [1/2]

nProbe+FastBit vs MySQL

Query	My:	SQL	nProbe Daily Po	+ FastBit artitions	nProbe + FastBit Hourly Partitions		
	No Index	With Indexes	No Cache	Cached	No Cache	Cached	
Ql	20.8	22.6	12.8	5.86	10	5.6	
Q2	23.4	69	0.3	0.29	1.5	0.5	
Q3	796	971	17.6	14.6	32.9	12.5	
Q4	1033	1341	62	57.2	55.7	48.2	
Q5	1754	2257	44.5	28.1	47.3	30.7	

Results are in seconds





Performance Evaluation: Query Time [2/2]

nProbe+FastBit vs nfdump

nProbe+FastBit	45 sec
nfdump	1500 sec

SELECT IPV4_SRC_ADDR, L4_SRC_PORT, IPV4_DST_ADDR, L4_DST_PORT, PROTOCOL FROM NETFLOW WHERE IPV4_SRC_ADDR=X OR IPV4_DST_ADDR=X

worth 19 GB of data (14 hours of collected flows)

nfdump query time = (time to sequentially read the raw data) + (record filtering time)







Host Geolocation [1/3]

- Host geolocation is a known problem (vd http:// en.wikipedia.org/wiki/Geoip)
- Need to handle thousand flows/sec (no inline internet query)
- Requirements: IP -> Location e IP -> ASN



MaxMind, GeoIP and related marks are registered trademarks of MaxMind, Inc. Copyright © 2010 MaxMind, Inc. All Rights Reserved. <u>Terms of use</u>.





2nd TMA PhD School - June 2011



Host Geolocation [2/3]

- Interactive Flash[™] world map, that displays hosts distribution by country and by cities of a selected country
- nProbe + GeoIP + Python + Google Visualization. The script
 - Cycles through all the hosts seen by ntop
 - Gets their GeolP info
 - Counts them based on their location.
- Google GeoMap and Visualization Table
- Ajax/JSON communications with web server for updated data







Host Geolocation [3/3]









How to Add Geolocation Data [1/3]

- Routers are unable to export any geolocation information.
- NetFlow/IPFIX flows do not contain any information about geolocation into standard flow formats.
- Solution:
 - Let the collector add geolocation information to flows received by routers
 - -Let the softprobe export this information to collectors.







How to Add Geolocation Data [2/3]

- nProbe takes advantage of GeoIP library (GPL) to
 - -Add geolocation information to flows
 - Map IP addresses to ASN (Autonomous System Numbers) for adding ASN awareness.
 - -GeoIPASNum.dat (ASN)
 - -GeoLiteCity.dat (GeoLocation)







How to Add Geolocation Data [3/3]

```
if(host->ipVersion == 4)
    return(GeoIP_record_by_ipnum(readOnlyGlobals.geo_ip_city_db, host->ipType.ipv4));
#ifdef INET6
    else if(host->ipVersion == 6)
    return(GeoIP_record_by_ipnum_v6(readOnlyGlobals.geo_ip_city_db, host->ipType.ipv6));
#endif
```

```
char *rsp = NULL;
u_int32_t as;
if(ip.ipVersion == 4)
    rsp = GeoIP_name_by_ipnum(readOnlyGlobals.geo_ip_asn_db, ip.ipType.ipv4);
    else {
    #ifdef INET6
        rsp = GeoIP_name_by_ipnum_v6(readOnlyGlobals.geo_ip_asn_db, ip.ipType.ipv6);
#endif
    }
    as = rsp ? atoi(&rsp[2]) : 0;
    free(rsp);
```







BGP Data Integration [1/2]









BGP Data Integration [2/2]

```
# Constructor
$update = Net::BGP::Update->new(
   NLRI
                  => [qw(10/8 172.168/16)],
                  => [ qw( 192.168.1/24 172.10/16 192.168.2.1/32 ) ],
   Withdraw
   # For Net::BGP::NLRI
   Aggregator => [ 64512, '10.0.0.1' ],
   AsPath
                  => [ 64512, 64513, 64514 ],
   AtomicAggregate => 1,
   Communities
                  = [qw(64512:1000064512:10001)],
   LocalPref
                  => 100,
                  => 200,
   MED
                  => '10.0.0.1',
   NextHop
   Origin
                  => INCOMPLETE,
```

);







What if you have no BGP Router? [1/3]

About RIPE NCC | Contact | Search | Sitemap RIPE NCC LIR Portal RIPE RIPE \$ 00 Quick Links Routing Information Service NCC you are here: home -> RIPE NCC Projects -> RIS **RELATED TOPICS** RIS: **RIS Raw Data** DNS Monitoring This page links to the raw data collected by the RRCs using Quagga routing software, stored in MRT format. This format is described in an IETF draft. These **RIS Home Page** . DISI files can be read using libbgpdump, a library written in C by Dan Ardelean, currently maintained by the RIPE NCC. A Python library also exists, PyBGPdump **RIS Raw Data** K-root Name Server which provides access to MRT files via Python. Documen***** RRCC Rawdata Please note that the raw data format recently changed to include support for 4-byte ASNs. Analysis using RIS TTM Contact Us For more information see:http://www.ripe.net/projects/ris/docs/asn.html Send Feedback Two sets of files are available for each of the RRC's: • All BGP packets, created with the Zebra command "dump bgp all". The filenames start with updates and are created every five minutes. The entire BGP routing table, created with the Zebra command "dump bgp routes-mrt ...". These files are created every eight hours, the filenames start with byiew. BGP Timer settings since 23 November 2006: Keepalives: 60 seconds Holddown: 180 seconds BGP Timer settings between 17 October 2002 and 23 November 2006: Keepalives: 0 (disabled) Holddown: 0 (disabled) BGP Timer settings before 17 October 2002: Keepalives: 60 seconds Holddown: 180 seconds Please note that this machine is connected to the Internet over a link that is loaded during Amsterdam office hours. If you want to download a large number of files, please do so at a spread-out rate, preferably at a time outside Amsterdam office hours (08:00 - 18:00 UTC). Click on a collector to get a list of available files. Files are grouped collector by collector, then month by month. rcc00.ripe.net at RIPE NCC, Amsterdam, collects default free routing updates from peers. From October 1999. rrc01.ripe.net at LINX, London. Collects route updates announced by LINX members. From July 2000. rrc02.ripe.net at SFINX. Paris. Collected route updates announced by SFINX members from March 2001 until October 2008.







What if you have no BGP Router? [2/3]

Index of /rrc10/2010.06

	Name	Last modifie	ad	<u>Size</u>	Descri	ption
2	Parent Directory			-		
N.	bview.20100601.0759.gz	01-Jun-2010	08:00	8.6M		
N	bview.20100601.1559.gz	01-Jun-2010	16:00	8.6M		
	bview.20100601.2359.gz	02-Jun-2010	00:00	8.6M		
	bview.20100602.0759.gz	02-Jun-2010	08:00	8.7M	. Pa	
	bview.20100602.1559.gz	02-Jun-2010	16:00	8.7M		bview
D	bview.20100602.2359.gz	03-Jun-2010	00:00	8.7M		bview
	bview.20100603.0759.gz	03-Jun-2010	08:00	8.7M		bview
	bview.20100603.1559.gz	03-Jun-2010	16:00	8.7M		updat
	bview.20100603.2359.gz	04-Jun-2010	00:00	8.6M		updat
	byiew.20100604.0759.gz	04-Jun-2010	08:00	8.6M		updat
A	byjew.20100604.1559.gz	04-Jun-2010	16:00	8.6M		updat
A	byjew 20100604 2359 gg	05-Jun-2010	00.00	8 64	1	updat
A	DV16w.20100004.2339.02	03-04H-2010	00:00	0.01	1	updat

	and the second			
	bview.20100613.2359.gz	14-Jun-2010	00:00	8.6M
	bview.20100614.0759.gz	14-Jun-2010	08:00	8.7M
Ū.	bview.20100614.1559.gz	14-Jun-2010	16:00	8.7M
J.	updates.20100601.0000.gz	01-Jun-2010	00:04	15K
D	updates.20100601.0005.gz	01-Jun-2010	00:10	14K
D	updates.20100601.0010.gz	01-Jun-2010	00:15	19K
D	updates.20100601.0015.gz	01-Jun-2010	00:20	17K
Ð	updates.20100601.0020.gz	01-Jun-2010	00:25	12K
D	updates.20100601.0025.gz	01-Jun-2010	00:29	13K
B	updates.20100601.0030.gz	01-Jun-2010	00:35	178
Ā	updates 20100601 0035 gz	01-Jun-2010	00.40	198
Ā	updates 20100601 0040 gg	01_Tup_2010	00.45	112
A	upuaces.20100001.0040.gz	01-Jun-2010	00:45	TIK





L



What if you have no BGP Router? [3/3]

- libbgpdump can be used to read BGP dump and updates.
- Periodically poll the RIPE RIS directory searching for full dumps or updates.

```
TIME: 06/15/10 15:59:58

TYPE: TABLE_DUMP_V2/IPV4_UNICAST

PREFIX: 12.51.167.0/24

SEQUENCE: 1321

FROM: 217.29.66.65 AS12779

ORIGINATED: 06/15/10 13:20:28

ORIGIN: IGP

ASPATH: 12779 1239 3356 19343 19343 19343 19343

NEXT_HOP: 217.29.66.65

COMMUNITY: 12779:1239 12779:65098
```

- Connect to the probe and refresh the routes according to the values being read.
- NOTE: always use the BGP dumps for a location near to you in order to have your view of the Internet.







Implementing a Web 2.0 GUI

- Web server: Lighttpd (easy and fast), avoid Apache.
- Ajax: use established frameworks such as jQuery or Prototype.



- Implement class libraries used to read your monitoring data. Python is used for speed, ease of use and script compilation.
- Use templates (e.g. Mako) for generating (XML-free) pages.



• Web frameworks are perhaps easier to use, but you will be bound to them forever (pros and cons).





Storing Historical Data [1/2]

RRD is the de-facto standard for permanently storing numerical data.

```
$rrd = "$dataDir/$agent-$ifIndex.rrd";
if(! -e $rrd) {
    RRDs::create ($rrd, "--start",$now-1, "--step",20,
        "DS:bytesIn:COUNTER:120:0:10000000",
        "DS:bytesOut:COUNTER:120:0:10000000",
        "RRA:AVERAGE:0.5:3:288");
    $ERROR = RRDs::error;
    die "$0: unable to create `$rrd': $ERROR\n" if $ERROR;
}
    RRDs::update $rrd, "$now:$ifInOctets:$ifOutOctets";
    if ($ERROR = RRDs::error) {
        die "$0: unable to update `$rrd': $ERROR\n";
    }
}
```

	6.0	T						
	5.0	T					-	in des
						100		
	4.0	1			1.1			
	3.0	τ				and the second second	ALC: NOT THE OWNER OF	
	2.0	T -	-					
							THE OWNER WATER	_
	1.0	1				and the second se		
	0.	e Oct	27 Jan 9	2. Bor 20	Jul 99	0::+20	angg Borg	
	0.	e Octi	97 Jan9	9 Apr-90	Jul 98	Oct90	lan99 Apr9	2
=	1.0 0.	e Octi	Jan 9	9 Apr90	Jul 98	Oct90	Ian99 Apr9	0
E	0. 04 1 04 1	e Octi	7 Jan9 101	9 Apr 90 1 If A 1 If H	Jul 99 LD1 118	0ct90 J 1/E 156 J	Ian99 Apr9 NNRI ISC-Mac	9 153 153
100	1.0 0. DH I EK C	e Octi DZ Biwi 0.2	7 Jan 9 104 104 187 00 115	9 Apr-90 1 IfA 1 IfH 1 fH	Jul 99 LEPI 118 07 196	0ct90 J	Am79 Apr9 NREI ISC-Mac	9 151 151
	1.0 0. DH 1 EK C TIK: TIK:	e Octi Dz 8:01 0.2 1.6	7 Jan9 104 187 187 187 187 183 119	9 Apr-90 10A 10A 10A 10A 10A 10A 10A 10	Jul 99 LD1 118 07 196 28 190	0ct90 J	Am79 Apr9 NREI ISC-Mac (TDyte) (TDyte)	9 151 151
	1.0 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0	e Octi DZ Bivi 0.2 1.6 0.9	7 Jan9 104 107 107 101 101 101 101 101 101	9 Apr-90 104 104 104 104 104 104 104 10	Jul 90 LD1 115 07 155 28 195 69 195	0ct90 J I/E 196 1 : 0.32: : 0.53: : 0.450	Anno Apro Neri ISC-flac (TDyte) (TDyte) (TDyte) (TDyte)	9 153 8 TI

		_	_				_								
	100			-							1.000	1000			
	80														
ž	60									۲.					
ŝ	40									1					
-															
-					- /	-		-				-			_
2	20				1	-	-	.	mu i	yel-	1-	-	-		-
2	20	Mar	Apr	Мау	Jan	Ju1	Aug	Sep	OCT	Nov	Dec	Jan	Feb	Mar	Apr
a o -	20	Mar 1	Apr - 10	May x3	Jun 10 -	ји1 100 ка	Aug 10	Sep 0 - 10	OCT OCT	Nou Gr	Dec eater t	Jan than 1	Feb	Rar	Apr
0 - 1	20 0- 1 KB	Mar 1	Apr - 10 CW	May KB	Jun 10 - 21	Jul 100 KB 0.672	Aug 10 Aver	Sep 0 - 10 age:	Oct 00 HB 23.4	Nov E Gr	Dec eater 1 Min:	Jan than 1 0.	Feb HB	Bar Nac	Apr 51.0
0 - 1 1 - 10	20 0- 1 KB KB	Mar 1	Apr - 10 CN CN	May x3	Jun 10 - 21 41	Ju1 100 kB 0.672 0.251	Aug 10 Aver Aver	Sap 0 - 10 age: age:	OCT 00 HB 23.4 50.2	Nov Gr 136 136	Dec eater t Min: Min:	Jan than 1 0.	Feb H0 000	Mar Max: Max:	Apr 51.0 65.6
0 - 1 1 - 10 10 - 1	20 0- 1 KB KB 100 KB	Har 1	Apr - 10 Cu Cu	May x3	Jun 10 - 21 41 21	Ju1 100 kB 3.672 3.251 5.013	Aug 10 Aver Aver Aver	Sep 0 - 10 age: age: age:	0ct 00 KB 23.4 50.2 21.7	Nov Gr 136 104 779	Dec eater 1 Min: Min: Min:	Jan than 1 0. 0.	Feb 140 000 000 000	Mar Mac: Mac: Mac:	Apr 51.0 65.6 72.2
0 - 1 1 - 10 10 - 1 100 -	20 0 1 KB 100 KB 1000 KB	Har 1	Apr - 10 Cu Cu Cu Cu	May KB rreat: rreat: rreat:	Jun 10 - 21 21 21	Jul 100 kB 3.572 3.251 5.013 1.910	Aug 10 Aver Aver Aver Aver	Sep 0 - 10 age: age: age: age: age:	0ct 00 KB 23.4 50.2 21.7 2.0	Nov Gr 136 129 151	Dec eater 1 Min: Min: Min: Min: Min:	Jan than 1 0. 0. 0.	Feb H6 000 000 000 000	Mar Mac Mac Mac Mac	Apr 51.0 65.6 72.2 3.6







Storing Historical Data [2/2]

- RRD has several limitations:
 - -Only one (quantity one) numerical data can be stored at each time interval (e.g. # of bytes received).
 - You must know 'in advance' what you want to store. For instance you can't store anything like 'the name and amount of traffic sent by the top host': the top host changes overtime, so you need an rrd per top host and this is not what you want.
 - Sets or lists of data (e.g. top protocols with bytes on interval X) cannot be stored in RRD.







Beyond RRD

- Requirements:
 - -Store network values are tuples (list of <name>:<value>, where <value> can also be a list).
 - Ability to aggregate tuples using a user-defined function (i.e. not just max/min/average).
 - -Manipulate values as RRD does: create, update, last, export, fetch and graph.
 - -Graph: images are not enough as we have tuples (not just one value) and also because the user must be able to interact with data, not just look at it.







pSWTDB [1/4]

- pSWTDB (Sliding Window Tuple DB).
- python class used to store tuples on disk using data serialization (called pickle on python).

-Pros:

- native in python
- portable across datatypes (i.e. no need to define the type)

-Cons:

- Slow as RRD (deserialize/update/serialize at each update)
- Same principle of RRD with the exception that here we use <u>tuples</u> and <u>not</u> numerical values.







pSWTDB [2/4]

- It comes with aggregation functions such as:
 - -Each time interval has a list of (key, value).
 - -Sum values with same key.
 - -Sort values
 - Discard values ranking after position X (e.g. take the top/bottom X values).
- Examples
 - -Top X protocols (list of <proto>:<value>)
 - -Top X hosts (list of <host>:(<proto>:<value>,...))







pSWTDB [3/4]

- Data are plotted using SVG/JavaScript.
- Users can interact with data (pan, zoom, move).
- Multiple criteria can be plotted at the same time (e.g. top X hosts and Y protocols).
- Clicking on data can be used to trigger GUI updates









pSWTDB [4/4]

```
deri@MacLuca.local 233> cat pcreate.py
#!/usr/bin/python
import pSWTDB
t = pSWTDB.pSWTDB('ptest.pkl')
# Hearbeat is 5 min
t.create(300)
# Keep 60 samples, one per minute
t.add base aggregation('1min', 60, 60)
# Keep 50 samples, each aggregating 5 samples
# of the base aggregation
t.add_aggregation('5min', 5, 50, pSWTDB.sum, '')
# Keep 60 samples, each aggregating 24 samples
# of the 5min aggregation
t.add aggregation('hour', 24, 60, pSWTDB.sum, '5min')
# Keep 30 samples, each aggregating 12 samples
# of the hour aggregation
t.add aggregation('day', 12, 30, pSWTDB.sum, 'hour')
```

```
deri@MacLuca.local 234> cat pupdate.py
#!/usr/bin/python
import pSWTDB
t = pSWTDB.pSWTDB('IT.pkl')
t.update('now',
         { 'keys' : ['APPL PROTOCOL'],
           'values' : ['SUM PKTS'],
           'data' : {
            'das' : ( 4522726 ),
            'domain': (1706286),
            'whois': ( 62838 ),
            'www': (28699),
            'smtp': (16149),
            'https' : ( 10892 ),
            'Unknown': (4934),
           })
deri@MacLuca.local 238> cat pfetch.py
#!/usr/bin/python
import pSWTDB
import pprint
t = pSWTDB.pSWTDB('IT.pkl')
ret = t.fetch('', 'now-1h', 'now')
print t.plot(ret)
```







Traffic Data Analysis [1/4]









Traffic Data Analysis [2/4]

```
deri@anifani 208> ls -l
total 24
16 drwxr-xr-x 3 root root 16384 May 25 08:21 aggregations/
 4 drwxr-xr-x 4 deri deri 4096 Mar 27 12:07 queries/
 4 drwxr-xr-x 6 deri deri 4096 Mar 18 19:37 rrd/
deri@anifani 209> ls -l *
aggregations:
total 34000
20 -rw-r--r-- 1 root root 18768 May 25 16:12 A1.pkl
164 -rw-r--r-- 1 root root 167641 May 25 16:12 A2.pkl
152 -rw-r--r-- 1 root root 154778 May 25 16:12 AD.pkl
216 -rw-r--r-- 1 root root 219872 May 25 16:13 AE.pkl
148 -rw-r--r-- 1 root root 148012 May 25 16:13 AF.pkl
152 -rw-r--r-- 1 root root 152841 May 25 16:13 AG.pkl
100 -rw-r--r-- 1 root root 100615 May 25 16:12 AI.pkl
. . .
152 -rw-r--r-- 1 root root 154259 May 25 16:13 YE.pkl
12 -rw-r--r-- 1 root root 10101 May 25 15:13 YT.pkl
200 -rw-r--r-- 1 root root 201469 May 25 16:12 ZA.pkl
148 -rw-r--r-- 1 root root 151246 May 25 16:12 ZM.pkl
156 -rw-r--r-- 1 root root 156071 May 25 16:12 ZW.pkl
308 -rw-r--r-- 1 root root 315311 May 25 16:13 all countries.pkl
  4 -rw-r--r-- 1 root root 791 May 15 23:55 ne.pkl
 24 drwxr-xr-x 2 root root 20480 May 22 13:57 top hosts/
queries:
total 8
4 drwxr-xr-x 7 deri deri 4096 May 1 00:05 2010/
rrd:
total 144
48 -rw-r--r-- 1 root root 47128 May 25 16:13 bits.rrd
12 drwxr-xr-x 2 root root 12288 May 6 02:06 bytes/
12 drwxr-xr-x 475 root root 12288 May 16 19:26 country/
12 drwxr-xr-x 2 root root 12288 May 24 23:36 flows/
48 -rw-r--r-- 1 root root 47128 May 25 16:13 flows.rrd
12 drwxr-xr-x 2 root root 12288 May 12 20:42 pkts/
```







Traffic Data Analysis [3/4]

rrd/country/CH/mandelspawn.rrd rrd/country/CH/gds db.rrd rrd/country/CH/dircproxy.rrd rrd/country/CH/rmtcfg.rrd rrd/country/CH/ssh.rrd rrd/country/CH/isisd.rrd rrd/country/CH/cfinger.rrd rrd/country/CH/gris.rrd rrd/country/CH/daap.rrd rrd/country/CH/x11.rrd rrd/country/CH/postgresgl.rrd rrd/country/CH/amanda.rrd rrd/country/CH/zephyr-hm.rrd rrd/country/CH/gsigatekeeper.rrd rrd/country/CH/fax.rrd rrd/country/CH/netbios-ssn.rrd rrd/country/CH/afs3-fileserver.rrd rrd/country/CH/cvspserver.rrd rrd/country/CH/ospf6d.rrd rrd/country/CH/bpcd.rrd rrd/country/CH/proofd.rrd rrd/country/CH/afs3-errors.rrd rrd/country/CH/qqz.rrd rrd/country/CH/tproxy.rrd rrd/country/CH/cfengine.rrd rrd/country/CH/x11-6.rrd rrd/country/CH/msp.rrd rrd/country/CH/rje.rrd rrd/country/CH/sane-port.rrd rrd/country/CH/smtp.rrd

deri@anifani 213> ls queries/2010/05/25/16/00/ total 1172 1164 cities.pkl 8 top n 17 protocols.pkl







Traffic Data Analysis [4/4]

```
deri@anifani 215> ~/nProbe/fastbit/python/dump.py cities.pkl |m
{'city': [['SRC COUNTRY',
           'SRC CITY',
           'SRC LATITUDE',
           'SRC LONGITUDE',
           'SRC REGION',
           'COUNT'],
          ['', '', '', '', 15079],
          ['IT', 'Rome', 41.899999999999999, 12.4832, 'Lazio', 1427],
          ['KR',
           'Seoul',
           37.56640000000002,
           126.9997,
           "Seoul-t'ukpyolsi",
           12501,
          ['RU',
           'Moscow',
           55.75220000000002,
           37.61560000000001,
           'Moscow City',
           1243],
          ['IT',
           'Milan',
           45.46670000000003,
           9.1999999999999993,
           'Lombardia',
           936],
```





Remote Probe Deployment [1/2]

- In order to monitor a distributed network it is often necessary to deploy remote probes.
- Exporting flows towards a central location is not always possible:
 - -Limited bandwidth available.
 - -Need to have a separate/secure network/tunnel as flows contain sensitive data.
 - -Interference with other network activities.
 - -Export of raw flows is much more costly than exporting the metrics we're interested in.







Remote Probe Deployment [2/2]

- Exporting data on off-peak times is not an option:
 - -We would introduce latency in data consumption.
 - The amount of data to transfer is not significantly reduced (zip flows) with respect to live data export.
 - -Unable to use the system for near-realtime analysis and alarm generation.
- Better solution
 - -Create a web service for querying data remotely in realtime
 - -Export aggregated metrics (e.g. .pkl files)







Web Interface: Internals [1/3]



(C) 2009-10 - IIT







Web Interface: Internals [2/3]









Web Interface: Internals [2/3]

Partition /home/deri/fastbit/netflow/2010/05/25/16/00

Protocol All Protocols

Live FastBit Query+Aggregation Python Glue Software

	Distance: 1				
ASN	AS Name	Traffic	:	Flows	Path
2597	REGISTRO CCTLD IT	1.0 GB	100.0 %	1367295	
	Distance: 2				
ASN	AS Name	Flows	Path		
3356	Level 3 Communications, LLC	358.3 MB	95.7 %	633432	9
<u>137</u>	GARR Italian academic and research network	12.0 MB	3.2 %	6951	9
<u>12637</u>	Seeweb Srl	3.1 MB	0.8 %	2949	0
21056	Welcome Italia S.p.A.	468.0 KB	0.1 %	1183	9
21309	ACANTHO SPA	434.8 KB	0.1 %	284	0
<u>64862</u>	??	216.0 KB	0.1 %	270	Q
16004	MIX S.r.L.	9.3 KB	0.0 %	20	9
15469	Warinet NOC AS	365.0 bytes	0.0 %	1	0
	Distance: 3				
ASN	AS Name	Traffic		Flows	Path
9035	Wind Telecomunicazioni spa	35.5 MB	14.3 %	34395	9
702	UUNET - Commercial IP service provider in Europe	25.9 MB	10.4 %	31974	9
6762	Telecom Italia international high speed,	21.6 MB	8.7 %	36706	9
3549	Global Crossing	17.7 MB	7.1 %	22630	9
24940	Hetzner Online AG RZ-Nuernberg	16.0 MB	6.4 %	8078	9
6453	Teleglobe Inc.	15.6 MB	6.3 %	22171	Q
12956	Telefonica Data Autonomous System	12.4 MB	5.0 %	28594	Q
1239	Sprint	11.2 MB	4.5%	17075	0







Using Geolocation Data [1/2]







2nd TMA PhD School - June 2011



Using Geolocation Data [2/2]







2nd TMA PhD School - June 2011



Disk and Memory Usage

- Collection of ~5k flows netflow/sec
- Each 5 min partition takes ~150 MB in FastBit format (32 GB/day)
- Partitions with raw data stay 3 days on disk (limited by available disk space)
- Each tuple archive in pickle format takes up to 400 KB (112 MB in total, almost constant).
- BGP patricia tree (inside the probe) of all routing tables takes about ~100 MB






Final Remarks

- NetFlow and sFlow are the two leading monitoring protocols.
- nProbe is an open-source software probe that can efficiently act as a probe/collector/proxy
- Storing and analyzing large volume of data is challenging but there are solutions available for doing it efficiently.
- Geolocation and routing information are useful for mapping traffic with users.





