Modern Packet Capture and Analysis: Multi-Core, Multi-Gigabit, and Beyond

Luca Deri <deri@ntop.org>







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.

IM 2009 - June 2009

• Strong Multicore NIC: Tilera Tile64





Accelerating Packet Capture and Analysis: PF_RING







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.
- Basic kernel (no low-level programming) knowledge 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	Mbps	% CPU Idle	Wire-Speed
250	259.23	518	> 90%	Yes
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







Socket Packet Ring: Packet Capture Evaluation

- Ability to capture over 1.1 Mpps on commodity hardware with minimal packet sizes (64 bytes).
- Available for Linux 2.4 and 2.6 kernel series.
- Hardware independent: runs on i386, 64bit, MIPS.
- Available for PCs and embedded devices (e.g. OpenWrt, MikroTik routers)





PF_RING on Embedded Devices









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 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.





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 multiprocessing/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 "ancient" 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]

#	tcpo	dump	-d	"udp"	
(000)	ldh		[12]	
(001)	jeq		#0x800	
(002)	ldb		[23]	
(003)	jeq		#0x11jt	4
(004)	ret		#96	
(005)	ret		# O	

# tcpd	dump -d	"udp and port 53"		
(000)	ldh	[12]		
(001)	jeq	#0x800	jt 2	jf
(002)	ldb	[23]		
(003)	jeq	#0x11jt 4	jf 12	
(004)	ldh	[20]		
(005)	jset	#0x1fff	jt 12	jf
(006)	ldxb	4*([14]&0xf)		
(007)	ldh	[x + 14]		
(008)	jeq	#0x35jt 11	jf 9	
(009)	ldh	[x + 16]		
(010)	jeq	#0x35jt 11	jf 12	
(011)	ret	#96		
(012)	ret	# O		





open source

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).
- 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 [1/4]

- An empty bloom is a bit array of <u>m</u> bits all set to zero.
- <u>k</u> hash <u>different</u> functions are used to map a key to an array position (0...m-1 hash function range).
- <u>n</u> is the number of elements insert into the dictionary.
- How to <u>add</u> an element: for each k hash function set to 1 the array bit that corresponds to the hash value.
- How to test if an element belongs to the set: for each hash function calculate the hash element value. The element belongs to the set if and only if all the k bits of the hash values are set to 1.
- How to remove an element: fully rebuild the dictionary or use counting blooms. $\left(\left(\left(1 \right)^{kn} \right)^k \right)^k = \left(\left(\left(1 \right)^{kn} \right)^k \right)^k$

- False positive rate: $\left(1 \left(1 \frac{1}{m}\right)^{kn}\right)^k \approx \left(1 e^{-kn/m}\right)^k$
- Optimal number of hash functions: $k = (m/n)^{\log(2)}$





Dynamic Bloom Filtering [2/4]



Check for inclusion







Dynamic Bloom Filtering [3/4]

- 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 independent from the number of filters.
- The "false positive rate" instead depends on the filters number.







Dynamic Bloom Filtering [4/4]



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







PF_RING Packet Parsing [1/2]

- 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.

```
struct pkt_parsing_info {
    /* Core fields (also used by NetFlow) */
    u_int16_t eth_type; /* Ethernet type */
    u_int16_t vlan_id; /* VLAN Id or NO_VLAN */
    u_int8_t l3_proto, ipv4_tos; /* Layer 3 protocol/TOS */
    u_int32_t ipv4_src, ipv4_dst; /* IPv4 src/dst IP addresses */
    u_int16_t l4_src_port, l4_dst_port; /* Layer 4 src/dst ports */
    u_int8_t tcp_flags; /* TCP flags (0 if not available) */
    [...]
};
```





PF_RING Packet Parsing [2/2]

- 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.

```
Struct pfring_pkthdr {
    struct timeval ts; /* time stamp */
    u_int32_t caplen; /* length of portion present */
    u_int32_t len; /* length this packet (off wire) */
    struct pkt_parsing_info parsed_pkt; /* packet parsing info */
    u_int16_t parsed_header_len; /* Extra parsing data before packet */
};
```

Plugin-based Parsing

 ts
 caplen
 len
 parsed_pkt
 parsed len
 17 parsing
 Payload

 (Optional)

 IM 2009 - June 2009

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) wilcard-ed 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 user-space (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: 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.

```
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;
```

- 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.






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.

```
typedef struct {
    u_int8_t proto;    /* Use 0 for 'any' protocol */
    u_int16_t vlan_id;    /* Use '0' for any vlan */
    u_int32_t host_low, host_high;    /* User '0' for any host. This is applied to both source
    and destination. */
    u_int16_t port_low, port_high;    /* All ports between 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;
```

- 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.







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).

IM 2009 - June 2009





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]



IM 2009 - June 2009





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 */
                                     IM 2009 - June 2009
                  UNIVERSITÀ DI PISA
```



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.

- 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, then 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







49

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.

IM 2009 - June 2009





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);
```

IM 2009 - June 2009





Università di Pisa



Using PF_RING Filters: HTTP Monitoring [4/5]

Plugin Packet Parsing

/* Fill PF_RING parsing information datastructure just allocated */

```
if((hdr->caplen > offset) && !memcmp(payload, "OPTIONS", 7))
                                                                    packet parsed->the method = method options;
                                                                      packet parsed->the method = method get;
  else if((hdr->caplen > offset) && !memcmp(payload, "GET", 3))
                                                                     packet parsed->the method = method head;
  else if((hdr->caplen > offset) && !memcmp(payload, "HEAD", 4))
  else if((hdr->caplen > offset) && !memcmp(payload, "POST", 4))
                                                                     packet parsed->the method = method post;
                                                                     packet parsed->the method = method put;
  else if((hdr->caplen > offset) && !memcmp(payload, "PUT", 3))
  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 http status code;
  else packet parsed->the method = method other;
}
```

IM 2009 - June 2009







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) {
    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);
```

IM 2009 - June 2009





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.





Advanced PF_RING Filtering: NetFlow [1/5]

- Goal
 - Using PF_RING for packet <u>capture</u> and processing in <u>user space</u>, the target performance (just packet capture, <u>not</u> flow generation) is:
 - Standard Intel driver: 550 Mpps
 - Enhanced Intel driver (see later in this presentation): 950 Mpps
 - Ability to compute NetFlow/IPFIX flows at wire speed at 1 Gbit regardless of the CPU being used and packet size.
- Solution
 - Use PF_RING plugin to "pack" packets belonging to the same flow. This acts as level-1 NetFlow cache.
 - Periodically (e.g. once every 1-5 sec) flush cache flows by forwarding packet digest to userspace via PF_RING.
 - Forwarded packets contains a header, used for computing flows, but not the packet as this is unnecessary. Each PF_RING slot can host several packets/ flows if needed.







Advanced PF_RING Filtering: NetFlow [2/5]

Each PF_RING cache entry contains exactly the same information necessary to generate flows.



 NetFlow cache is walked (for searching expired flows) by user-space application through a dummy call to pfring_purge_idle_hash_rules() that allows to keep kernel code simple as there's no need to spawn a kernel thread for walking the cache.







Advanced PF_RING Filtering: NetFlow [3/5]

- The PF_RING cache has (by default) 4096 entries and it is implemented as an array (i.e. hash buckets are not a linked list) for keeping code simple.
- User-space application can modify cache policy/size when PF_RING is instrumented.
- The plugin is activated with a wildcard-ed rule of 'any' so that any IP packet matching the filter can be computed.
- Modus Operandi
 - When an incoming packet is received, PF_RING parses it, and then it is passes to the plugin.
 - Using parsing information the packet is searched in the cache
 - If found the cache entry is updated
 - if not found the packet is added to the cache (i.e. a filtering rule is added). In case the cache slot where the packet is supposed to be stored is already occupied, the slot is flushed (i.e. the entry is forwarded to the userland) and the packet is accommodated.







Advanced PF_RING Filtering: NetFlow [4/5]

- Using the kernel cache, packets are "merged" in kernel without any userland application intervention.
- In-kernel packet merging does not require any memory/packet copy and it's very fast as the packet is already in the CPU cache (thanks to Intel RSS/DCA, see later in this presentation).
- The "merging rate" increases (in efficiency) with flows speed. In other terms the cache is more efficient as flows are faster. Example:
 - 1 Gbit (1.48 Mpps) flow with minimal packets.
 - Kernel cache duration of 3 sec (i.e. flows older than this duration are exported)
 - "Vanilla" PF_RING: in 3 sec the application receives 4.44 Million packets (3 x 1.48 Mpps).
 - In-kernel cache generates 1 flow for the same amount of traffic.







Advanced PF_RING Filtering: NetFlow [5/5]

- Performance Evaluation
 - Testbed: 1.86 GHz Intel CoreDuo (cost < 100 Euro), IXIA 400 Traffic generator, minimal packet size (64 bytes), Intel e1000 driver, Full 1 Gbit stream, with packet rotation, nProbe (home grown NetFlow probe) used as probe.
- Vanilla PF_RING + nProbe: 100% CPU, ~600 Kpps processed with no loss.
- Kernel NetFlow PF_RING plugin + nProbe: ~60-70% CPU used, wire-rate with no packet-loss.
- Comparison:
 - spare CPU cycles compared to vanilla PF_RING.
 - wire-speed.
 - not suitable (yet) for generating flows with packet payload information (e.g. HTTP URL).







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.

IM 2009 - June 2009





Direct Access to NICs







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.

IM 2009 - June 2009

- Cards falling in this category include:
 - Endace DAG
 - Napatech
 - NetFPGA





Direct NIC Access: Comparison [1/2]





Hardware Acceleration







PF_RING
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 packets 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 Evaluation

- It currently supports Intel 1 GE copper/fiber cards.
- GE Wire speed (1.48 Mpps) full packet capture starting from P4 HT 3 GHz.
- Better results (multiple NICs on the same PC) can be achieved using Opteron machines (HyperTransport makes the difference).
- The nCap speed is limited by the speed applications fetch packets from the NIC, and the PCI bus.

IM 2009 - June 2009





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)



Further nCap Features

- High-speed traffic generation: cheap trafgen as fast as a hardware trafgen (>> 25'000 Euro)
- Precise packet generation.
- Precise packet time-stamping on transmission (no kernel interaction): suitable for precise active monitoring.
- Enhanced driver currently supports Intel cards (1 Gb Ethernet).





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);
```

IM 2009 - June 2009





PF_RING DNA: Current Status

- As of today, DNA is available for Intel 1 Gbit NICs (e1000 driver).
- it is planned to support more modern 1G NICs later this year.
- 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

IM 2009 - June 2009





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: 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 = E1000 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);
```

IM 2009 - June 2009

[....]







Memory Allocation Life Cycle [4/5]







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.

IM 2009 - June 2009





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

UNIVERSITÀ DI PISA

- 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).
- 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.

```
while(<polling packets from RX queue X>) {
    /* Avoid CPU monopolization */
    if(rx_budget > 0)
        rx_budget--;
    else {
        rx_budget = DEFAULT_RX_BUDGET;
        yield();
    }
}
M 2009-June 2009
```





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								
	CPU0	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	pts								
	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 000000 # echo # cat 000000 # cat 198:</pre>	/proc/irq/198 008 5 2 > /proc/ir /proc/irq/198 002 /proc/interru 0	/smp_affini q/198/smp_a /smp_affini pts grep ' 67	affinity [00 ity "eth7:v6-Rx" 0	0000010 whe 5309	re 1 = CPU : 1	1] 2	0	1	PCI-MSI-edge	eth7:v6-Rx

IM 2009 - June 2009

- How to map a process to a CPU/core

```
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>
```





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)





Enhanced NIC Drivers: TNAPI [8/8]

- Additional Performance results:
 - 10 Gbit
 - The testbed is a 4 x 1 G ports IXIA 400 traffic generator that are mixed into a 10G stream using a HP ProCurve 3400cl-24 switch.
 - A dual 4-core 3 GHz Xeon has been used for testing.
 - Using the accelerated driver it is possible to driver-filter 512 bytes packets at 7 Gbps with a 1:256 packet forward rate to user-space with no loss.
 - 1 Gbit
 - The same testbed for 10G has been used.
 - The same packet filtering policy applied to 2 x 1 Gbit ports works with no loss and with minimal (~10%) CPU load.
 - The performance improvement also affects packet capture. For instance with a Core 2 Duo 1.86 GHz, packet capture improved from 580 Kpps to around 900 Kpps.







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).
- 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





-

-



IM 2009 - June 2009



Strong Multicore NICs: Tilera Tile64







108
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
 - Use communication libraries to optimize
 - Use tools to tune performance











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...









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.









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 1 Gbit 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, <u>http://people.redhat.com/drepper/cpumemory.pdf</u>, RedHat 2007.







Programming for Multicore [3/4]





Programming for Multicore [4/4]



Great Application DesignExploit Native MulticoreFully Lockless Hash

Lockeless hashes:

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







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.







pen source

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 <u>http://www.hoard.org</u>/





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>>





