

Dynamic Core Allocation and Packet Scheduling in Multicore Network Processors

Muhammad Faisal Iqbal*, Jim Holt[†], Jee Ho Ryoo[‡], Gustavo de Veciana[§], Lizy K. John[¶]

*^{‡§¶}University of Texas at Austin

[†] Freescale Semiconductor Inc. & MIT Computer Science and Artificial Intelligence Laboratory

{*faisaliqbal,[‡]jhr45842}@utexas.edu, {[§]gustavo, [¶]ljohn}@ece.utexas.edu

[†]jim.holt@freescale.com, jholt@csail.mit.edu

Abstract—With ever increasing network traffic rates, multicore architectures for network processors have successfully provided performance improvements through high parallelism. However, naively allocating the network traffic to multiple cores without considering diversified applications and flow locality results in issues such as packet reordering, load imbalance and inefficient cache usage. Consequently, these issues degrade the performance of latency sensitive network processors by dropping packets or delivering packets out of order. In this paper, we propose a packet scheduling scheme that considers the multiple dimensions of locality to improve the throughput of a network processor while minimizing out of order packets. Our scheduling policy tries to maintain packet order by maintaining the flow locality, minimizes the migration of flows from one core to another by identifying the aggressive flows, and partitions the cores among multiple services to gain instruction cache locality.

Our light weight hardware implementation shows improvement of 60% in the number of packets dropped and 80% in the number of out-of-order packet deliveries over previously proposed techniques.

Keywords—Network Processor, Load Balancing, Resource Management.

I. INTRODUCTION

A *Network Processor* is a special-purpose, programmable device that is optimized for network operations. A network processor is generally a multicore processor that can process network packets at wire-speeds of multi-Gbps. Network processors are employed in many demanding network processing environments like core and edge routers. While the main requirement in core routers is high capacity to handle huge amounts of traffic, edge routers require programmability and flexibility in order to support multiple complex applications like intrusion detection, firewalls, protocol gateways, etc. A network processor provides the performance of custom silicon and programming flexibility of general purpose cores. The ability of a network processor to perform complex and flexible processing and its programmability make it an excellent solution for core and edge routers.

A number of network processors exist in the market. These processors can be classified into two categories. The first category includes general purpose multicore processors that are adapted to perform networking functions. Examples of

such processors are the ThunderX [1], Sun Niagara [2] and Tileria [3] processors. The second category includes processors which are specifically designed for networking applications. These processors are equipped with hardware accelerators and co-processors in addition to a large number of general-purpose cores. Examples include the Freescale T4240 [4], Broadcom XLP [5], EZChip [6], Cisco nPowerX1 [7] and IBM PowerNP [8]. Both of these categories have a common attribute: they utilize a large number of cores to achieve desirable performance by exploiting parallelism. Networking applications have abundant parallelism because multiple packets can be processed by different cores in parallel. This packet level parallelism makes multicore architectures well suited for networking applications [9]. Network processors with 64 cores or more have been announced by vendors to handle 100 Gbps network speed [10], [11]. With increasing traffic rates and processing demands, the number and complexity of cores in these processors are on the rise and efficiently managing these cores has become very challenging. In this work we focus on dynamic adaptations based on run time traffic behavior in order to optimize performance and make following contributions.

First, design of a hash based packet scheduler and load balancer is presented in order to achieve the goals of preserving flow locality and packet order. A hash based packet scheduler performs very well in order to achieve these goals because it schedules packet at the flow level and thereby maintains packet order and flow locality inherently. A serious impediment to performance of hash based scheduler is the presence of skewed flow sizes in network traffic. Such skewed distribution of flow sizes can result in overloading some cores and may result in packet loss. To avoid packet loss, a load balancer is designed that migrates some flows from the overloaded cores to under-utilized cores. Flow migrations are undesirable because they result in bad data locality and can result in out of order packets. The load balancer proposed in this study minimizes the number of flow migrations by restricting migrations only to the aggressive flows. We present a low cost hardware to identify aggressive flows. Second, the design of the scheduler is extended to support multiple applications in a router where cores can be dynamically allocated to applications. Furthermore, use of incremental hashing is proposed which is low cost method that minimizes number of flow migrations when cores are allocated or deallocated to the services.

In the next section we present architecture of a network

Manuscript received April 19, 2005; revised January 11, 2007.

processor and discuss issues related to packet scheduling in these processors. The design of the packet scheduler is presented in Sections III and III-C.

II. BACKGROUND AND MOTIVATION

A. Architecture of a Network Processor

Many architectural variations of network processors exist in the market. Although vendors differ in specific implementations, they generally share three main features: 1) Multiple cores to exploit packet level parallelism, 2) Accelerators for networking functions and 3) Optimized path for movement of packet data.

Architecture for a typical network processor is shown in Figure 1. An incoming packet is received by a Frame Manager (FM). FM Places the packet payload in a buffer allocated by the Buffer Manager and places the header, a pointer to the buffer and some meta data as command descriptors in the input queue to a processing core. The general purpose core processes the packet and can offload some of the work to accelerators e.g., it can put some of the work in the queue for security accelerator (SEC). SEC performs the required processing and puts the packet back to the return queue. Eventually, the general purpose core sends the packet back to the FM via an enqueue after finishing the processing.

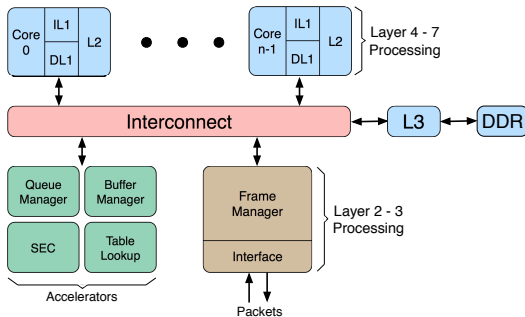


Fig. 1. Typical Architecture of a Network Processor

Network processing can be classified as either *Control Plane* or *Data Plane*. *Control Plane* is responsible for control and management processing e.g., maintaining and updating the routing tables. Control Plane processing may involve executing routing protocols like RIP, OSPF, and BGP, or control and signalling protocols such as RSVP or LDP. *Data Plane* deals with actual processing involved in packet forwarding. The data plane execution involves compression, encryption, address searches, address prefix matching for forwarding, classification, traffic shaping, network address translation and so on. In many network processors, the general purpose cores are responsible for processing both data and control plane packets. However, in majority of modern high speed network processors, control plane processing is separated from data plane processing [4], [12]. When a packet arrives, a packet classifier in the FM decides whether it is a control or a data plane packet. Control plane packets take the *slow path* through general purpose cores. The data plane packets (Layer 2 or

possibly Layer 3) take the fast path and are not offloaded to general purpose cores. Fast path processing is handled by the FM itself.

The FM is equipped with a large number (32-120) of small cores called I/O Processors (IOP). These IOPs are in-order, dual issue cores with non coherent memory, and generally do not have an operating system. When a packet arrives the packet classifier first identifies whether it is a control plane or a data plane packet. If it is a L2 or possibly L3 data plane packet, it is handled in the FM autonomously by IOPs, otherwise it goes to general purpose cores. This configuration describes a notional system that represents a class of chips as they look today and moving into the next 3-5 years. In this work we are interested in scheduling of data plane packets on IOPs. Since these packets arrive at a very high rate (100Gbps and even higher in future), an efficient scheduling of packets on IOPs is needed in order to gain good performance. The term IOP and core are used interchangeably in this work.

B. Challenges in Packet Scheduling

The design of scheduler for these applications is very challenging. First, the scheduler is in the data path and therefore it should be as efficient as possible. Second, it should meet the requirements of packet ordering, flow locality and cache locality.

a) *Packet Ordering*: Although the internet is designed to tolerate out-of-order packets, performance of upper layer protocols, such as Transmission Control Protocol (TCP), greatly depends on packet ordering. Out of order packets can falsely trigger congestion control mechanisms and degrade throughput unnecessarily [13]. Also, applications like Voice Over IP (VOIP) and multimedia transcoding require that packets arrive in order because the receiver might not be able to easily reorder the packets. Hence, it is important to preserve the order among the packets of a flow. In this work, a flow is defined as a set of packets that have the same source address, destination address, source port, destination port and protocol. If packets from the same flow are processed by different cores, they can experience different queuing and processing delays, and consequently, the probability of out of order delivery of packets increases. Careful scheduling of packets is needed in the network processors to minimize out of order departure of packets.

b) *Load Balancing*: Load balancing is an important technique to efficiently utilize multiple cores in a network processor. Packets arriving at the input should be distributed uniformly to the available processing cores to maximize performance. An unbalanced allocation of load can swamp some cores. As a result, incoming packets assigned to overloaded cores will experience large delays and may even result in packet loss due to limited storage in the network processor.

c) *Data Cache Locality*: If different cores process packets of the same flow, the data cache will be used inefficiently as the same data is copied to multiple caches. Packet processing needs to access per flow data (state, statistics), as well as more global data (routing table). If packets of a flow always go to the same core, locality can be preserved for both local

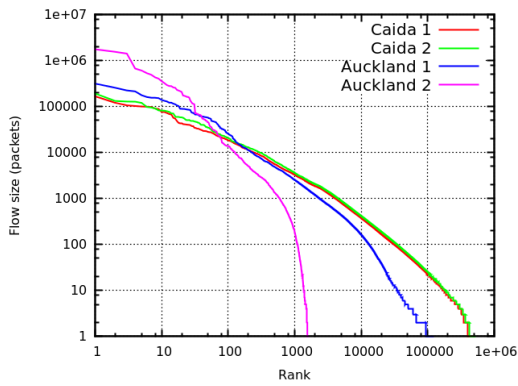


Fig. 2. Distribution of flow sizes in real network traces. Rank 1 is the flow with the highest flow size.

flow migrations. Flow migrations result in out-of-order packets and also badly affect data cache performance.

In order to minimize the number of flow migrations, previous research has made the observation that migrations should be limited to only top aggressive flows [23]. In this way load balance can be achieved by minimum number of flow migrations. However, previous research based its study purely on offline analysis and kept per flow statistics to identify aggressive flows. Maintaining per flow statistics has a lot of overhead and is not possible in realistic designs. Although many per flow statistics are maintained by software, accessing those software statistics is very time consuming for a scheduler that is trying to schedule data plane packets. The data plane packet scheduler needs to function with minimum software intervention for good performance.

This research presents the design of a hardware scheduler for data plane packets. A novel low-overhead hardware technique to identify aggressive flows is presented. The aggressive flow detection scheme is based on two-level caching idea of annex-cache [27] used in general purpose applications. The caching based aggressive flow detector integrates readily with a hash based packet scheduler. The complete design and evaluation of the scheduler is presented in this section.

B. Packet Scheduler Design

The proposed packet scheduler uses a hash based design which is a natural way of maintaining flow locality and packet order. The scheduler is called Locality Aware Packet Scheduler (LAPS). When a packet arrives, its flow identifier is extracted from the header. Flow identifier is a five tuple consisting of source and destination IP addresses, source and destination ports and protocol ID. This five tuple is hashed using CRC16 to get an index into a map table. The map table² stores target core ID where the packet is eventually forwarded. In the presence of skewed flow size distribution as shown in Figure 2, the scheduler identifies and migrates the aggressive flows from the overloaded core to achieve load balance. An efficient scheme for identifying and migrating aggressive flows is presented.

²Map table is used instead of direct hashing because it allows dynamic core allocation presented in the next section.

When a core becomes overloaded, i.e., its queue size reaches a threshold, the scheduler needs to migrate some of the incoming traffic from that core to a less loaded core. This migration of flows has two drawbacks: One, it makes some cached data in the source core useless and triggers some cold misses in the cache of newly allocated core. Two, flow migration makes it harder to maintain the order among packets of the flow. The new incoming packets will potentially experience less queuing delay as compared to older packets that are waiting in the overloaded core's queue.

To avoid the above two situations, it is desirable to minimize the number of flow migrations. If only the most aggressive flows can be identified and migrated, load balance can be achieved with minimum disruption, i.e., only a few flows need to be migrated to achieve load balance. In order to achieve this, a low cost mechanism is needed to identify top aggressive flows. This research proposes a novel cache based hardware called Aggressive Flow Detector (AFD) to identify the top flows. The hardware consists of a small fully associative cache called Aggressive Flow Cache (AFC). AFC is augmented with a cache assist called annex cache. Detailed architecture of annex cache and AFC is presented in Section V-C.

Figure 3 presents the scheduler design. The incoming packets are hashed to get index into a map table that stores the target core IDs. On load imbalance, the incoming packet flow to the overloaded core is migrated to the least loaded core if the flow is identified as an aggressive flow by AFD. The decision is recorded in the Migration table. So the future packets of the same flow are always migrated to the newly allocated core. The scheduler gives priority to the output of migration table over the default hash table. If the input queue indicated by the scheduler is filled up, the incoming packet is dropped.

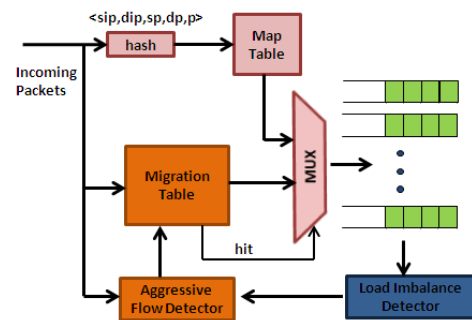


Fig. 3. Load balancer design

1) *Aggressive Flow Detection*: The design of Aggressive Flow Detector is based on *annex cache*. Annex Cache was proposed by John [27] to exploit locality in the memory references in general purpose processor workloads. This study shows that such a structure can be very useful in to identify aggressive flows.

The AFD has two main components as shown in Figure 4. One component is a small fully associative cache called Aggressive Flow Cache (AFC). AFC holds the IDs of top aggressive flows. All entries into AFC come via annex cache. Items referenced only rarely will be filtered out by annex cache

and will never enter AFC. The basic premise is that a flow deserves to enter AFC only if it proves its right to be in AFC by showing locality in the annex cache. Annex cache also serves as a victim cache and provides some inertia before a flow is excluded from the AFC. Both AFC and annex cache use Least Frequently Used (LFU) replacement policy.

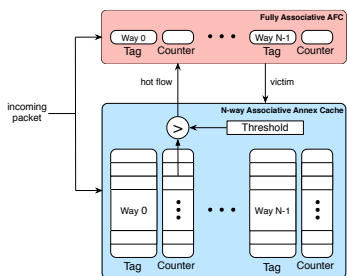


Fig. 4. Structure of Aggressive Flow Detector (AFD)

The design of AFD is slightly different from the one presented in [27] because in AFD annex cache is bigger than AFC. A larger AFC is undesirable because the proposed scheme wants to limit the number of monitored aggressive flows. Annex cache is a bigger structure that serves as a qualifying station for large number of flows to demonstrate their eligibility to be cached into the AFC. When a packet arrives, its flow ID is checked in both AFC and annex cache. If it is a hit in AFC, the hit counter is incremented. On a hit in the annex cache, flow counter is incremented and the value is compared with a pre-defined threshold. The threshold for promotion to AFC is the LFU count in AFC. If the hit count in annex cache exceeds the threshold, the flow is promoted to AFC. The victim flow from AFC is then placed in the annex cache. Finally on a miss in annex cache, a flow replaces the LFU flow of the annex cache.

2) *Load Imbalance Detection:* Length of the longest queue is used to detect the load imbalance in the system, i.e., when the length of the longest queue in the system reaches a predefined threshold, load imbalance signal is asserted. As long as the load imbalance signal is asserted, all the aggressive flows are migrated to the least loaded core. The migrated flows are forwarded to the new core even after the load imbalance signal is de-asserted as a result of flow migration.

Most modern network processors have dedicated hardware units for management of packet queues [4], [28], [8] and a lot of research has been done on design of these hardware queue managers [29], [30], [31]. These queue managers implement different active queue management algorithms (e.g., Random Early Detection RED) and monitor the queue length as part of their normal operation. This queue length information can easily be used by the load balancer to detect the need for low migration, i.e., it can easily be reported to the packet scheduler when the queue reaches a threshold. Hence, additional hardware resources are not needed to monitor queue length, because queues are already monitored for congestion control purposes. In this work, it is assumed that hardware queue manager monitors the queue state and generates the load imbalance signal.

C. Packet Scheduler For Multiservice Routers

A simple hash based design as presented in III can result in inefficient i-cache usage. In order to exploit i-cache locality, LAPS divides the pool of cores among all active applications or services. In effect, there is a separate map table for each service. All cores in a single map table always get packets that require same processing. Hence, i-cache locality is preserved. The main question is how to allocate cores to applications? If cores are allocated on compile time, we need provisioning for worst case traffic requirements of each service requiring a huge number of cores. Fortunately, all services do not experience their worst case traffic simultaneously and hence cores can be multiplexed dynamically among services to keep the total number of cores reasonable. Many dynamic allocation schemes have been proposed in the past. In this work we adapt the policies presented in [32] to integrate it with hash load balancer. We further make it flow aware so that the number of flow migrations are minimized on dynamic allocation and deallocation of cores. LAPS utilizes incremental hashing (also known as Linear hashing) to minimize number of flow migrations on dynamic adaptation.

1) *Allocation of Cores to Services:* LAPS keeps a list of cores that are marked as surplus cores by other services (Section III-C2). When a service becomes overloaded, LAPS looks through the list of surplus cores and finds the core that has been marked extra for the longest period of time and allocates this core to fulfill the demands of requesting service. This policy makes sure that the deallocated core has the least utility for the victim service. The core ID is added to the list of allocated cores for the requesting service.

2) *Release of Cores by Services:* When input queue to a core becomes empty, a timer starts. When the timer reaches idleth, the core is marked surplus by adding it to a list of extra cores. The core still remains allocated to the same service. In case, the same service needs more resources in near future (before the core is put to deep sleep state by the power management scheme), this core can be unmarked and removed from the list of surplus cores without incurring the overhead of context switch. If the core is actually allocated to another service, it is removed from the bucket list of the victim service. Other core IDs will be shifted to take the place of this ID. The bucket size b is decremented by 1 and the hash function is also changed accordingly. This may result in some flow migrations but the performance overhead is tolerable because this service is only lightly loaded anyway. The value of idleth is set to $10\mu s$ based on previous research [33].

3) *Load Redistribution on Core Allocation:* When an additional core is allocated to a service, the resource manager appends the core ID to the end of the hash table for that service, i.e., hash table size grows by 1. *Linear Hashing* scheme allows a hash table to grow one bucket at a time and does not require rehashing of all flows currently allocated. This makes it useful for load balancing because it is desirable to minimize the flow disruption when an additional core is allocated to a service. The Linear Hashing scheme was introduced by [34] and has been described in [35]. Following is a brief introduction of how this scheme works.

and the dynamic size b_i changes with traffic variations. The hash function for each service is decided based on the size of its bucket list. Following steps are taken when a packet arrives:

- 1) If the flow ID hits in the migration table, the packet is forwarded to the core ID indicated by the migration table.
- 2) If the flow ID does not hit the migration table, the map table is searched using the hash function and the packet is forwarded to the core indicated by the mapping table.
- 3) Under load imbalance, the aggressive flows (flows that hit in AFC) are migrated to the least loaded core allocated to that service similar to the load balancing scheme of Section III.
- 4) When number of cores allocated to a service become insufficient, the bucket lists are updated. An idle core is removed from the bucket list of donor service and is added to the bucket list of overloaded service.

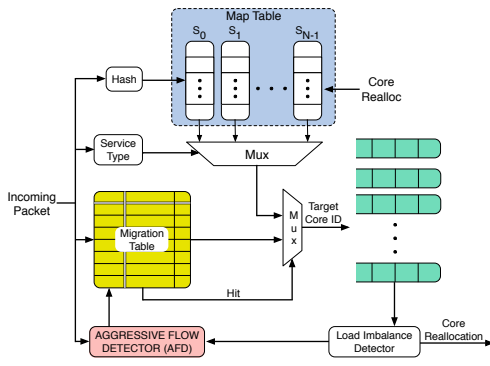


Fig. 8. Locality Aware Packet Scheduler

IV. EVALUATION INFRASTRUCTURE

A. Traffic Traces

In this work we used real network traces to evaluate the performance of packet scheduler. Following is a small description of set of traces used in this study.

1) *CAIDA Traces*: This dataset contains anonymized traffic traces from CAIDA's equinix-sanjose monitor [36]. This monitor is connect to OC-192 link. These set of traces are captured in year 2011 and are of duration of 1 minute each.

Trace	Name
Caida 1	20110120-125905.UTC.anon.pcap.gz
Caida 2	20110120-130000.UTC.anon.pcap.gz
Caida 3	20110120-130100.UTC.anon.pcap.gz
Caida 4	20110120-130200.UTC.anon.pcap.gz

TABLE I. LIST OF CAIDA TRACES USED IN THE STUDY

2) *University of Auckland Traces*: This set of traces, also known as AUCK-II, is captured at University of Auckland and captures the traffic between the university and its ISP [37]. All connections from the university to external world pass through this measurement point. These traces are of one hour long duration each.

Trace	Name
Auckland 1	20000614-181539-0.gz
Auckland 2	20000614-181539-1.gz
Auckland 3	20000619-183717-1.gz
Auckland 4	20000621-105006-0.gz
Auckland 5	20000621-105006-1.gz
Auckland 6	20000630-175712-0.gz
Auckland 7	20000630-175712-1.gz
Auckland 7	20000703-152100-0.gz

TABLE II. LIST OF AUCKLAND-II TRACES USED IN THE STUDY

B. Workload Model

In order to model the different services running on a multi-service router we consider a workload similar to the one presented in Figure 9. This model is based on methodology presented in [16]. In modern network processors, all tasks of the same path are scheduled on the same core to reduce the communication overhead. Hence, in this study we consider all the tasks on the same path as a single service. Thus our simulations have four active services in the processors. A packet is tied to a single core for the life time of its processing. The incoming packets can be serviced by one of the four

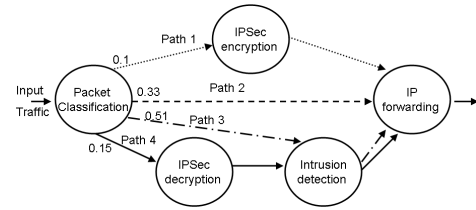


Fig. 9. Example Task graph for an edge router

services represented by different paths of Figure 9. *Path 1* describes the path of outgoing packets which are tunnelled via VPN. *Path 2* represents the default handling of packets. *Path 3* is the path of incoming packets on edge router that are scanned for malware and *Path 4* is for incoming VPN packets which are decrypted and scanned for malware.

C. Simulation Infrastructure

For evaluating different scheduling strategies, we developed a simulation model in SpecC [38]. SpecC is similar to systemC [39] in its design and philosophy. Different components of the simulator are shown in Figure 10.

1) *Packet Generator*: *Packet Generator* generates traffic with programmable traffic rates. To generate packets, it reads the real packet traces. We govern the traffic for each path based on Holt-Winterz forecasting as suggested in [40]. The traffic rate is governed by the equation 1.

$$x_i(t) = a + b.t + C.S(t \bmod m) + n(\sigma) \quad (1)$$

where $x_i(t)$ is the traffic rate for service i , a is the baseline traffic component, b is the trend component, C is the magnitude of seasonal component S , m is the period of seasonal component, n is random noise with a standard deviation of σ . Total incoming traffic is the sum of traffic of each individual service. The header for each generated packet is taken from

real network traces. We use a separate packet trace for each path of the flow graph. The use of real network traces ensures that realistic flow scenarios are created.

2) *Scheduler*: The *Scheduler* module implements the different scheduling strategies. Once a decision has been made the input packets are enqueued into the input queue of the target core. The queue size is set to 32 packet descriptors for each queue based on pervious research [19]. A packet is lost when it is assigned to a queue which is already full.

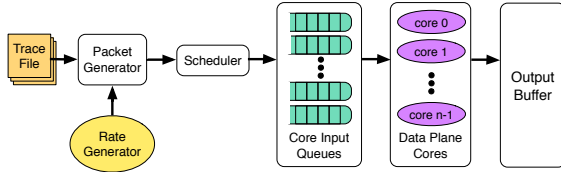


Fig. 10. Simulation infrastructure

3) *Processing Latencies*: Each packet of a service i , experiences a Processing Delay (PD_i) in the core based on the following equation

$$PD_i = T_{proc,i} + FM_{penalty} + CC_{penalty} \quad (2)$$

Where $T_{proc,i}$ is the processing time, $FM_{penalty}$ is the penalty due to flow migration and $CC_{penalty}$ is the cold cache penalty which occurs when subsequent packet needs different processing than the previous packet. $T_{proc,i}$ is derived from real delays seen by the packets when the packet processing is implemented in software on a full system GEMS [41] simulator. The configuration of in-order cores is shown in Table III.

Frequency	Pipeline	Branch Predictor	I-Cache	D-Cache
1GHz	7 stage 2-issue	gshare/BTB 128 entry each	16KB 2 way	32KB 4 way

TABLE III. DATA PLANE CORE CONFIGURATION

We executed these packet processing applications and derived a packet processing delay model for each service. T_{Proc} is measured to be $0.5\mu s$ for path 2 i.e., IP forwarding. For path 3 it is measured to be $3.53\mu s$. For Path 1 it also depends on the packet size and is given as

$$T_{proc,path1} = 3.7\mu s + \frac{PacketSize}{64byte} \times 0.23\mu s \quad (3)$$

Similarly the processing time for path 4 is given as

$$T_{proc,path4} = 5.8\mu s + \frac{PacketSize}{64byte} \times 0.21\mu s \quad (4)$$

$FM_{penalty}$ is set to four cache misses ($0.8\mu s$) conservatively (two for routing data and two for per flow data). In reality, a flow migration can cause a lot more misses depending on how much per flow data is maintained. Becasue of small I-cache, these cores can hold instructions of only the last executed program (e.g., AES encryption used in IPSec requires 16KB). So whenever a packet of different service arrives at a core,

it will experience cold cache penalty. We set the cold cache penalty to $10\mu s$ which is the cold cache penalty for the smallest service i.e., IP Forwarding as observed in GEMS simulator. In practice this penalty will be higher because many services are larger and a context switch can result in some D-Cache misses too. For simplicity, we ignore the D-cache misses due to context switch in this work.

V. RESULTS AND DISCUSSION

A. Throughput Improvement with LAPS

LAPS aims to improve throughput by exploiting locality in instruction and data caches. Figure 12 shows effectiveness of LAPS in improving throughput of a sixteen core system. In this experiment, all four services of Figure 9 are active. Simulation infrastructure of Figure 10 is used. The traffic rate generator is configured to increase the traffic gradually to measure the maximum throughput supported by the system. Traffic is equally divided among the four services, i.e., Path 1 through 4 of Figure 9. Caida 1, Caida 2, Caida 3 and Caida 4 traces are used for generating packets for Path 1, Path 2, Path 3 and Path 4. Figure 12 compares throughput of LAPS with a First Come First Served (FCFS) and Arbitrary Flow Shift (AFS) scheduler. X-axis shows combined input traffic rate which is equally divided among all the services and y-axis shows the traffic rate observed at the output.

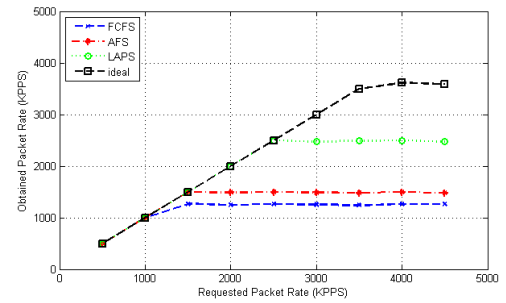


Fig. 12. Throughput comparison of different schedulers

FCFS scheduler services packets in their arrival order and does not consider flow or instruction locality. As a result, it causes many data and instruction cache misses and results in the worst throughput among the three schedulers. As compared to FCFS, AFS reduces some flow migrations and is able to improve throughput a little. But AFS is still unaware of instruction locality and results in suboptimal performance. In comparison to these two schemes, LAPS improves both flow and instruction locality and results in substantially better throughput (56% more than AFS and about 100% more than FCFS). Ideal throughput represents a system with no cache miss penalties. The plot is obtained by setting the cache miss penalties to zero. Although such a system is infeasible, it represents a theoretical maximum which can be achieved if the system has full knowledge of everything and is able to move data and instructions into caches before they are needed.

Note that the throughput supported by the simulated sixteen core system is much less than the industrial system. There are

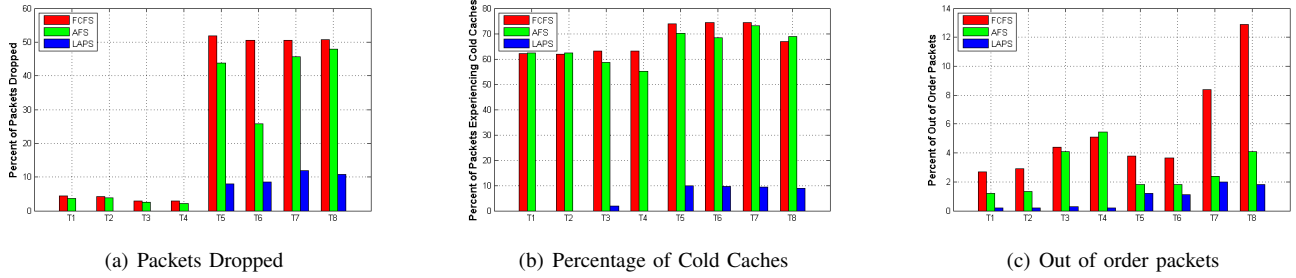


Fig. 11. Comparison of LAP with FCFS and AFS with different traffic scenarios listed in Table VI

two reasons for this: First, the software implementations of services are taken from open source benchmark suites where as companies use highly optimized implementations. Second, the experiments are based on software only implementations and do not use hardware accelerations.

B. Performance Improvement with LAPS

In this section, we compare the performance of LAPS with a First Come First Served (FCFS) scheduler and the scheduler presented in [20]. This scheduler migrates arbitrary flows when load imbalance is detected. We call this scheme Arbitrary Flow Shift (AFS). For this set of experiments, traffic rate is governed by equation 1. We experimented with different sets of parameters for equation 1 and LAPS outperforms other schemes in all scenarios. We present results with two sets of parameters listed in Table IV. Set 1 represent the under-load scenario i.e., the aggregate traffic rate is less than the ideal capacity of 16 cores. Set 2 represents an overload scenario i.e., the data rate is more than the capacity of the 16 core system.

	Service	a	b	C	m	σ
Set 1	S1	1	0.03	0.3	40	0.1
	S2	1.8	-.025	0.1	25	0.05
	S3	0.5	0.01	0.07	60	0.25
	S4	0.3	0.005	0.09	600	0.3
Set 2	S1	1.5	0.002	0.3	100	0.3
	S2	1.3	-.02	0.15	25	0.05
	S3	1	0.004	0.25	30	0.25
	S4	0.7	0.01	0.18	200	0.3

TABLE IV. PARAMETERS GOVERNING TRAFFIC RATE. RATE IS IN MPPTS AND PERIOD IS IN SECONDS

For each service, we use real network traces to generate the packet. We used different sets of traces listed in Table V. The combination of sets of equation in Table IV and traces in Table V creates difference traffic scenarios listed in Table VI. Figure

Group	S1	S2	S3	S4
G1	Caida1	Caida2	Caida3	Caida4
G2	Caida5	Caida6	Caida2	Caida3
G3	auck1	auck2	auck3	auck4
G4	auck5	auck6	auck7	auck8

TABLE V. TRACES USED IN EXPERIMENT FOR PACKETS OF INDIVIDUAL SERVICES

Scenario	Parameter Set	Trace Group
T1	Set 1	G1
T2	Set 1	G2
T3	Set 1	G3
T4	Set 1	G4
T5	Set 2	G1
T6	Set 2	G2
T7	Set 2	G3
T8	Set 2	G3

TABLE VI. TRAFFIC SCENARIOS USED IN FIGURE 11

11(a) shows packets dropped with three schemes under the traffic scenarios shown in Table VI. LAPS outperforms FCFS and AFS in both underload and overload conditions. FCFS and AFS distribute packets of different services arbitrarily to cores and suffer from poor I-cache locality (Figure 11(b)). These schemes drop packets even in underload conditions because almost 60% of packets suffer from cold cache penalties. On the other hand, LAPS partitions the cores among services effectively and enjoys good I-Cache performance. Under overload scenarios (T5 through T8), LAPS also suffers from some cold caches because cores are dynamically switched between services based on traffic variations.

LAPS maintains data and instruction cache locality and is able to sustain higher traffic input rates. Figure 11(c) shows the effectiveness of LAPS in preserving packet order under traffic scenarios of Table VI. FCFS does not care for packet ordering and hence results in most out of order packets. AFS reduces these out of order packets but still there are considerable amount of out of order packets due to large number of flow migrations. LAPS minimizes the flow migrations by only migrating the top flows and hence result in very few packets being delivered out of order. Next, we show how our proposed Aggressive Flow Detector (AFD) helps in identifying the top flows and helps to achieve load balance with minimum flow migrations.

C. Performance of Aggressive Flow Detector (AFD)

The proposed AFD has two components: An aggressive flow cache (AFC), and an annex cache. An annex cache can be viewed as a preliminary filter where non-aggressive flows are filtered out from entering the small AFC. Therefore, any entry in AFC is a considered an aggressive flow. We evaluate the effectiveness of AFD by varying annex cache size while setting the size of AFC constant at 16 entries. Since our

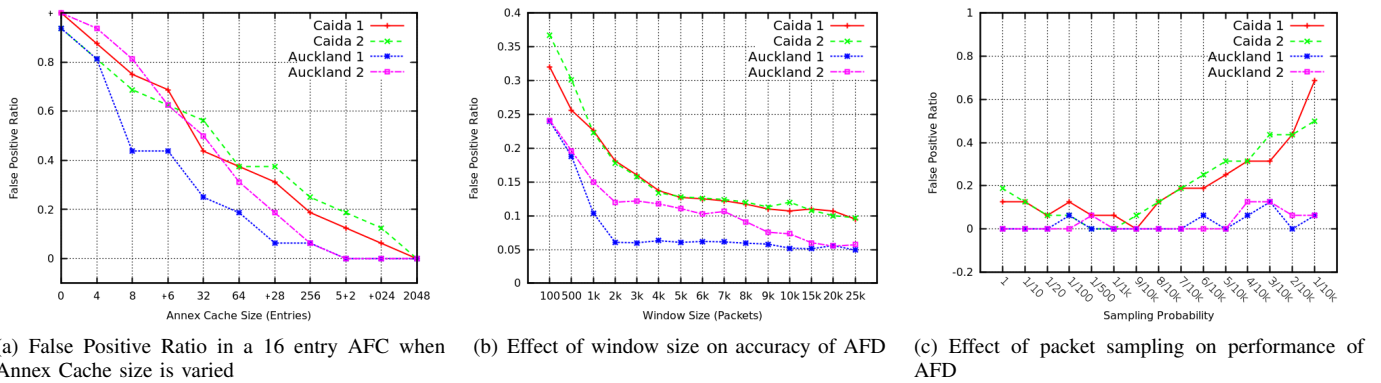


Fig. 13. Effectiveness of AFD in identifying aggressive flows

AFC size is fixed, we can only detect up to the maximum of 16 top aggressive flows. A perfectly accurate AFC will hold the IDs of top 16 aggressive flows. A flow found in AFC, which is not among the top 16 flows identified by offline analysis is considered a false positive. Figure 13(a) shows the false positive ratio (false positives/total entries) in AFC when annex cache size is varied. As the size increases, the annex cache can hold more flows to choose a possible candidate for a promotion to the AFC. In other words, the pool of aggressive flow candidates increases and the chances of aggressive flows residing in the cache for the AFC promotion becomes higher. For Auckland traces, AFC can identify all top 16 flows with 100% accuracy with a 512 entry annex cache. Caida traces have much more flows active and thus require a larger annex cache. In Caida 1 and 2 respectively, only 14 and 13 most aggressive flows are correctly identified with a 512 entry annex cache. When we double the size to 1024 entries, accuracy improved an average of 6.25%. Although there are 2 or 3 false positives in Caida 1 and 2 cases, they are not random flows that are promoted to the AFC. In fact, when we consider 20 most aggressive flows as our area of interest, these false positives fall into the aggressive flow category. Yet, for consistency of our work, we treat those flows as false positives. We only looked at the accuracy of our mechanism at the end of our simulations until now. Since LAPS needs to peek into the AFC whenever load balancing is required, we performed another experiment where the accuracy is checked at every fixed intervals. In Figure 13(b), we performed the same accuracy evaluation with varying interval steps. In this experiment, we assumed the fixed 512 entries for the size of the annex cache. Our mechanism shows above 90% accuracy from a small step size such as every 1000 packets to large step sizes. This implies that our AFC will contain the most aggressive active flows regardless of when it is accessed. In dynamic scheduling schemes like ours, it is key to maintain a high level of accuracy across the entire execution. Figure 13(c) shows the false positive ratio when packets are sampled with a probability p and not all packets access the AFD. It is interesting to note that FPR improves initially with sampling. This is because sampling acts as a filter i.e., the probability of large flows being sampled is higher than the smaller flows. However, the

performance deteriorates for Caida traces at larger sampling intervals. Sampling up to 1/1k probability gives better or equal performance than sampling all packets for all traces. Caida traces have generally large number of high data rate flows and hence their performance deteriorates if sampling is increased too much. Sampling not only improves the accuracy but also reduces power consumption because now each packet does not have to access the AFD.

D. Dynamic Behavior of the System

In order to observe the effectiveness of dynamic resource allocation scheme, temporal behavior of number of cores allocated to each service is plotted. Figure 14 shows the dynamic behavior when two services are active in the system. Service 1 is the same as Path 1 of Figure 9, i.e., the outgoing VPN traffic is encrypted using IPSEC encryption. Service 2 is the Path 3 of Figure 9 which corresponds to processing incoming packets through a firewall. The traffic requirements of each service are varied over time and the response of the resource allocation system is observed. Figure 14 shows that the system is very effective in following the changing traffic requirements and changes the core allocations to match the demands of each service very effectively.

E. Opportunities of Migration Without Reordering

The results presented in the previous section show significant improvements in minimizing number of flow migrations and percentage of packets which leave the system out of order. However, there is still a small percentage of packets which are out of order (1-2%). Ideally, any out of order packets should be eliminated. In this section we study, the opportunities of migrating flows without any risk of packet reorder. Such opportunities exist because the gap between the packets may be large enough to allow flows to migrate without any risk of reordering [42].

To investigate such opportunities a simple experiment is conducted. Whenever a packet arrives, the opportunity count is incremented by 1 if the target queue does not contain any packets from the same flow. Table VII shows the opportunity counts for Caida 1 trace. Note that this is a conservative

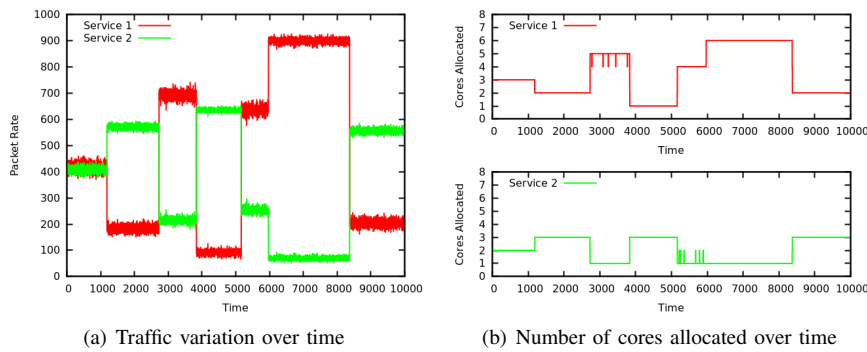


Fig. 14. Temporal behavior of the resource allocator

estimate because even if packets exist in the queue, it is still possible to migrate without reordering if the existing packets will be processed before the new incoming packet.

First row shows, when any packet arrives and it is safe to migrate it without any reordering. The second row shows the same results when the input queue length of the target core is greater than the imbalance threshold. This is more important than because this is the situation where flow migration is needed. The third row represents the situation when the incoming packet belongs to the top flow and the target core is overloaded. This is the situation where our scheme can benefit by migrating the large flow without causing any reordering. This is a great opportunity to further minimize out of order departure of packets. It is possible that current scheme when moves a big flow does not cause reorder. But the scheme in itself does not have any such mechanism. It can be further extended such that when its time to migrate flow, priority could be given to flows with minimum packets in the system to minimize reordering. This will help minimize reorderings but will make the system little more complex. Because now, we need to store the target cores for the large flows and also need to keep track of the number of packets belonging to that flow in the system. Such a system is part of future work.

Flow Type	Qlength	Total Packets	PMO
Any	Any	28219211	20385555
Any	$> qthresh$	14068978	9375921
Top	$> qthresh$	619555	252825

TABLE VII. OPPORTUNITIES FOR FLOW MIGRATION WITHOUT REORDERING

F. Analysis of Flows on Migration

In order to understand the behavior of the system, the flows allocated to the overloaded core at the instance of flow migration are analyzed. Figure 15(a) shows the number of unique flows present in the queue of the overloaded core when a big flow is migrated from that core. Generally, a large number (15-20) of flows are present. This indicates that the overload is caused by a combination of large and small flows and migrating the large flow is expected to mitigate the load imbalance. A small number of flows would indicate that the

overload is caused by small number of large flows and there is a potential that migrating the large flow would result in imbalance even in the newly allocated core. Figure 15(b) shows the number of big flows allocated to the overloaded core at the time of flow migration. From the plot it can be seen that the imbalance is usually caused by multiple big flows and migrating one flow is not expected to cause imbalance in the new core.

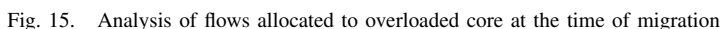
Presence of multiple big flows in the queue of overloaded core opens up the opportunity to further improve the scheme. For example, when migration decision is made preference could be given to the big flow which will cause less distortion in the order of packets, e.g., the flow with less number of packets in the queue could be preferred over the flow with larger number of packets. Such a scheme is likely to increase the complexity of the system because now core association of the flows and their number of packets in the system need to be monitored. Design of such a system is part of future work.

VI. RELATED WORK

A. Existing Work on Packet Ordering

Previous researchers have adopted two different approaches to minimize packet reordering in network processors: order restoration and order preservation.

1) *Order Restoration*: This technique allows packets belonging to a flow to be processed out of order by different cores and restores the order at the output [8], [19], [18]. At the input, each packet is tagged with a sequence number and the packets are allowed to exit the system in strict sequence order. Per flow tagging is needed in order to preserve order among packets of the same flow. This requires keeping per flow information, which is a huge overhead as there can be millions of flows active at a time [43], [44]. Overhead of per flow tagging can be reduced by using global tagging. Global tagging is easy to implement but it is overly restrictive as it forces order even among packets of different flows and results in throughput degradation. Order restoration also requires an expensive synchronization mechanism because multiple cores may be required to update the ordered list of packets of the same flow at the same time. This scheme also results in poor data cache locality because flow locality is not preserved.



Many researches have observed the need for dynamic resource allocation in network processors [14], [51], [52] and there have been proposals for runtime resource allocations in the past [15], [53], but these schemes consider a packet processing application as a graph where different tasks within the application form the nodes of the graph. These schemes consider adjacency between nodes for task scheduling as packets move between different cores in a pipelined manner during processing. In contrast, this research considers each service as a single entity, i.e., a packet is tied to a single core for the whole processing and graph or pipeline scheduling is not considered. Wolf et al. [54] observed that the mix of packets destined for each service varies with time. If packets of different services are sent to the same core, i-cache locality cannot be maintained. This results in huge performance overhead. They attempted to address the issue of i-cache locality through careful packet scheduling. When a core becomes idle, their scheme searches for a packet of

0018-9340 (c) 2016 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications_standards/publications/rights/index.html for more information.

- Lizy Kurian John** is B. N. Gafford Professor of Electrical Engineering at the University of Texas at Austin. She received her Ph.D in computer engineering from The Pennsylvania State University in 1993. Her research interests include microprocessor architecture, performance and power modeling, workload characterization, and low power architecture. She is an IEEE fellow.