

Technical Report : Study and Implementation of 20 Scheduling and Shaping Policies

Mostafa Elbediwy, Bill Pontikakis, Jean-Pierre David, and Yvon Savaria, *Fellow, IEEE*

Abstract—In this technical report, the pseudocodes for different scheduling policies, listed by the P4 working group, are presented. We are targeting to implement these algorithms in the ingress (enqueue) and egress (dequeue) pipelines of the P4 target switches. Between the Enqueue and the Dequeue levels, there are the DR-PIFO scheduler [1] [2], which always schedules the lowest rank eligible packet. The DR-PIFO scheduler is successfully integrated into target switches that are programmed by the P4 programming language. So, the following algorithms calculating the rank of packets/flows and/or the eligible time of packets. In case of shaping policies, a packet should not be scheduled before its eligible time.

In this work, we present 20 scheduling and shaping policies. 15 of them already have a well-defined rank-based version which is compatible with the DR-PIFO scheduler. the other 5 policies which are : Deficit Round Robin (DRR), Weighted Round Robin (WRR), Weighted Deficit Round Robin (WDRR), Slytherin, and the rate-limited Strictly Priority (non-work-conserving), require novel rank-based versions to utilize the DR-PIFO scheduler. Slytherin is the only scheduling policies out of these 5 policies that have an approximated rank-based version which seems to provide an unfair service to the network users. Accordingly, in a companion paper, we propose novel rank-based versions for these 5 scheduling policies while, in this report, we only present the original algorithms of these policies in Algorithms 1 2 3 4 5.

Algorithm 1: (P1) Original Deficit Round Robin [3]

```

/* In DRR, there are multiple "FIFO" queues,
   one dedicated to each flow.
   "Quantum" is the maximum number of bytes,
   that a flow can send every round.
   N : is the number of the supported flows. */
1 At Enqueue of each new packet (P) :
   /* Insert the new packet at the tail of the
      FIFO queue of its flow. */
2   FIFO[P.flow_id].enqueue(P);

Dequeueing Module :
3 while true do
4   for  $i \leftarrow 1$  to  $N$  do
5     if  $FIFO[i]$  not empty then
6        $deficit\_count[i] += Quantum$ 
7       while  $deficit\_count[i] \geq$ 
          $Size(FIFO[i].head)$  and  $FIFO[i]$  not empty
         do
8          $deficit\_count[i] -=$ 
            $Size(FIFO[i].head);$ 
9          $FIFO[i].dequeue();$ 
10      if  $FIFO[i]$  is empty then
11         $deficit\_count[i] = 0$ 

```

Algorithm 2: (P2) Original Weighted Round Robin [4]

```

/* "Weight" is the maximum number of packets,
   that a flow can send every round. */
1 At Enqueue of each new packet (P) :
2   FIFO[P.flow_id].enqueue(P);

Dequeueing Module :
3 while true do
4   for  $i \leftarrow 1$  to  $N$  do
5      $sent\_pkts = 0$ 
6     while  $sent\_pkts <$ 
        $weight[i]$  and  $FIFO[i]$  not empty do
7        $++ sent\_pkts;$ 
8        $FIFO[i].dequeue();$ 

```

Algorithm 3: (P3) Original Weighted Deficit Round Robin [5][6]

```

1 At Enqueue of each new packet (P) :
2   FIFO[P.flow_id].enqueue(P);
   Dequeuing Module :
3 while true do
4   for  $i \leftarrow 1$  to  $N$  do
5     if FIFO[ $i$ ] not empty then
6       deficit_count[ $i$ ] +=  $Quantum[i]$ 
7       while deficit_count[ $i$ ] >=
          $Size(FIFO[i].head)$  and FIFO[ $i$ ] not empty
         do
8         deficit_count[ $k$ ] -=
            $Size(FIFO[i].head)$ ;
9         FIFO[ $i$ ].dequeue();
10      if FIFO[ $i$ ] is empty then
11      deficit_count[ $i$ ] = 0

```

Algorithm 4: (P4) Original Slytherin [7]

```

1 At Enqueue of each new packet (P) :
2 if CE bit is set then then
3   P.rank = low_rank; // high priority
4 else
5   P.rank = high_rank; // low priority

```

Algorithm 5: (P5) Original Rate Limited Strict Priority (Non-Work-Conserving)

```

1 Dequeuing Module :
2 while true do
3   for  $Q \leftarrow high\_queues$  to  $low\_queues$  do
4     if  $Q$  not empty and
        $current\_rate(Q) < rate\_limit(Q)$  then
5      $Q$ .dequeue();

```

Algorithm 6: (P6) Least Attained Service (LAS) [8]

```

1 // array to count the service received by each flow.
2   service_count = array[number_of_flows];

   At Enqueue of each new packet (P) :
3   service_count[P.flow_id] =
4     service_count[P.flow_id] + P.length;
5   P.rank = service_count[P.flow_id];

```

Algorithm 7: (P7) pFabric with starvation prevention

```

1 current_rank = array[N];

   At Enqueue of each new packet (P) :
2 P.rank = P.remaining_flow_size;
3 if  $P.rank < current\_rank[P.flow\_id]$  or
    $current\_rank[P.flow\_id] == 0$  then
4   current_rank[P.flow_id] = P.rank;
5   update_rank = current_rank[P.flow_id];
6   update_id = P.flow_id;

```

Algorithm 8: (P8) Weighted Fair Queueing (WFQ) [9] [10]

```

1 // array for the weights assigned to each flow.
2   weights = array[number_of_flows];

   // array for the finish time for each flow
3   finish_time = array[number_of_flows];

   At Enqueue of each new packet (P) :
4   length_over_weight = P.length / weights[P.flow_id];

   if  $finish\_time[P.flow\_id] > P.arrival\_time$  then
5     finish_time[P.flow_id] =  $finish\_time[P.flow\_id]$ 
6       + length_over_weight;
7 else
8   finish_time[P.flow_id] = P.arrival_time
9     + length_over_weight;
10
11   P.rank = finish_time[P.flow_id];

```

Algorithm 9: (P9) Least Slack Time First (LSTF) [11]

```

1 At Enqueue of each new packet (P) :
2   P.rank = P.slack + P.arrival_time;

   At Dequeue of each scheduled packet (P) :
3   P.slack = P.rank - P.departure_time;
4 // if a packet queued for more than its slack, it can be
5 // dropped.

```

Algorithm 10: (P10) Least Attained Recent Service (LARS) [12]

```

1 // array to count the service received by each flow.
2   attained_service = array[number_of_flows];

   At Enqueue of each new packet (P) :
3 if current_time_unit >= decay_time_unit then
4   |   attained_service[P.flow_ID] =
5   |   |   attained_service[P.flow_ID] * decay_factor;
6   ++ current_time_unit;
7   attained_service[P.flow_ID] =
8   |   attained_service[P.flow_ID] + P.length;
9   P.rank = attained_service[P.flow_ID];

   // To avoid starvation, the ranks of previous packets
10 // may need to be updated if their ranks are higher
11 // than the new packet's rank.

```

Algorithm 11: (P11) Stop-and-Go [13]

```

1 // T : the length of the time frames
   At Enqueue of each new packet (P) :
2 if current_time >= frame_end_time then
3   |   frame_begin_time = frame_end_time;
4   |   frame_end_time = frame_begin_time + T;
5 P.rank = frame_end_time;

```

Algorithm 12: (P12) Rate Controlled Service Disciplines (RCSD), (Jitter-EDD)

```

1 At Enqueue of each new packet (P) :
2   P.rank = P.local_deadline + P.arrival_time;
3   P.eligible_time = P.Ahead + current_time;

   At Dequeue of each scheduled packet (P) :
4 if P.rank < P.departure_time then
5   |   P.Ahead = P.departure_time - P.rank;
6 else
7   |   P.Ahead = 0;

```

Algorithm 13: (P13) Window Constrained Scheduling - Virtual Deadline Scheduling (VDS) [14] [15]

```

1 // Constraints parameters.
2   M = array[number_of_flows];
3   K = array[number_of_flows];
4   T = array[number_of_flows];
5   M' = array[number_of_flows];
6   K' = array[number_of_flows];
7 // the arrival time of the last packet from each flow.
8   arrival_times = array[number_of_flows];
9
10 At Enqueue of each new packet (P) :
11   arrival_times[P.flow_id] = P.arrival_time;
12   -- K'[P.flow_id];
13 if K'[P.flow_id] == 0 then
14   |   // if M' is above zero, that could be a violation.
15   |   K'[P.flow_id] = K[P.flow_id];
16   |   M'[P.flow_id] = M[P.flow_id];
17
18   P.rank = ( ( T[P.flow_id] * K'[P.flow_id] ) /
19   |   M'[P.flow_id] ) + P.arrival_time;
20 for all previous packets (PP) with the same P.flow_id
21 do
22   |   PP.rank = P.rank;
23
24 At Dequeue of each scheduled packet (P) :
25 // assume synchronization reset of M' and K'
26   -- M'[P.flow_id];
27 if M'[P.flow_id] == 0 then
28   |   M'[P.flow_id] = M[P.flow_id];
29   new_rank = ( ( T[P.flow_id] * K'[P.flow_id] ) /
30   |   M'[P.flow_id] ) + arrival_times[P.flow_id];
31 for all previous packets (PP) with the same P.flow_id
32 do
33   |   PP.rank = new_rank;
34
35 // Note : in case of violation, packets may be dropped.

```

Algorithm 14: (P14) Approximate Fair Queueing [16]

```

1 "srv_cntr" : array to count the service received by each
  flow.
  "Q" : "Quantum" the maximum number of bytes, that
  a flow can send every round.
   At Enqueue of each new packet (P) :
2 srv_cntr[P.flow_id] += P.length;
3 round_id = ⌊(srv_cntr[P.flow_id] - 1)/Q⌋;
4 P.rank = round_id;

```

Algorithm 15: (P15) Hierarchical Policies (FIFO-DRR-SP)

```

1 // As an example :
2 // Assume 3-level hierarchy,
3 // level-2 (root) Strict-Priority (SP) scheduling,
4 // level-1 (middle) Weighted Fair Queueing (WFQ),
5 // level-0 (leaf) First-In-First-Out (FIFO) scheduling.
6
7 // array for the weights assigned to each flow.
8     weights = array[number_of_flows];

    // array for the finish time for each flow
9     finish_time = array[number_of_flows];

    At Enqueue of each new packet (P) :
10 // FIFO ordering
11     P.rank[0] = P.arrival_time;
12 // WFQ
13     length_over_weight = P.length / weights[P.flow_id];
14 if finish_time[P.flow_id] > P.arrival_time then
15     |   finish_time[P.flow_id] = finish_time[P.flow_id]
16     |   + length_over_weight;
17 else
18     |   finish_time[P.flow_id] = P.arrival_time
19     |   + length_over_weight;
20     P.rank[1] = finish_time[P.flow_id]; // SP (ToS, is
    the priority class for this packet).
21     P.rank[2] = P.ToS; // assign high or low priority

```

Algorithm 16: (P16) Penalize Heavy Hitters (PHH) [17]

```

1 // array to count the number of received packets
2 // from each flow in the current time window.
3     received_packets = array[number_of_flows];
4
5 // the time of the last reset triggered for each flow id
6     last_reset_time = array[number_of_flows];
7
8 At Enqueue of each new packet (P) :
9 if (current_time - last_reset_time[P.flow_id]) >=
    window_period then
10 |   last_reset_time[P.flow_id] = current_time;
11 |   received_packets[P.flow_id] = 0;
12
13     ++ received_packets[P.flow_id];
14 if received_packets[P.flow_id] >= Threshold then
15 |   // Low priority
16 |   P.rank = Higher_Rank_Value;
17 else
18 |   // High priority
19 |   P.rank = Lower_Rank_Value;

```

Algorithm 17: (P17) WFQ with weights determined by queue occupancy (WFQ-QO) [18]

```

1 // array for the expected ratio between the delay of
2 // each flow versus the delays of the rest of the flows.
3     delay_ratio = array[number_of_flows];
4
5 // array to store the weights of each flow.
6     weight = array[number_of_flows];
7
8 // array to calculate the number of stored packets
9 // from each flow.
10    queue_occupancy = array[number_of_flows];
11
12 // array for the finish time for each flow
13    finish_time = array[number_of_flows];

    // the time of the last reset triggered for each flow id
14    last_reset_time = array[number_of_flows];
15
16 At Enqueue of each new packet (P) :
17 if (current_time - last_reset_time[P.flow_id]) >=
    window_period then
18 |   weight[P.flow_id] = queue_occupancy[P.flow_id] /
19 |   delay_ratio[P.flow_id];
20 |   last_reset_time[P.flow_id] = current_time;
21
22    length_over_weight = P.length / weights[P.flow_id];
23
24 if finish_time[P.flow_id] > P.arrival_time then
25 |   finish_time[P.flow_id] = finish_time[P.flow_id]
26 |   + length_over_weight;
27 else
28 |   finish_time[P.flow_id] = P.arrival_time
29 |   + length_over_weight;
30
31     ++ queue_occupancy[P.flow_id];
32     P.rank = finish_time[P.flow_id];
33
34
35 At Dequeue of each scheduled packet (P) :
36     -- queue_occupancy[P.flow_id];
37
38 // Note : to implement the "Dynamic WFQ"
39 // scheduling algorithm, instead of only
40 // (queue_occupancy), it should be replaced by
41 // ( queue_occupancy + arrival_rate[P.flow_ID] ).

```

Algorithm 18: (P18) Rate Limited Strict Priority and Work-Conserving (RL-SP-WC)

```

1 // array to count the number of received packets
2 // from each flow in the current time window.
3     received_packets = array[number_of_flows];
4
5 // the current window ID of each flow.
6     window_ID = array[number_of_flows];
7
8 // array contains the maximum packets that each
9 // flow can send during a single time window.
10    maximum_packets = array[number_of_flows];
11
12 // number of possible ranks at a single time window.
13    possible_ranks = number_of_possible_SP_ranks *
14         $\sum_{n=0}^{n < \text{number\_of\_flows}}$  maximum_packets[n]
15
16 At Enqueue of each new packet (P) :
17 if (received_packets[P.flow_ID] ==
    maximum_packets[P.flow_ID]) then
18     |   received_packets[P.flow_ID] = 0;
19     |   ++ window_ID[P.flow_id];
20
21     ++ received_packets[P.flow_id];
22 // P.ToS is the strict priority value for each packet.
23     P.rank = P.ToS + window_ID[P.flow_id] *
24         possible_ranks;

```

Algorithm 19: (P19) Variant of Weighted Fair Queueing (NUMFabric) [19]

```

1 // array for the finish time for each flow
2     finish_time = array[number_of_flows];

At Enqueue of each new packet (P) :
3 // weights are dynamically updated by the received
4 // packets. maybe a register is needed, if weights are
5 // usually updated once in a while.
6     length_over_weight = P.length / P.weight;

if finish_time[P.flow_id] > P.arrival_time then
7     |   finish_time[P.flow_id] = finish_time[P.flow_id]
8     |   + length_over_weight;
9 else
10    |   finish_time[P.flow_id] = P.arrival_time
11    |   + length_over_weight;
12
13    P.rank = finish_time[P.flow_id];

```

Algorithm 20: (P20) Independent Scheduling and Shaping Policies (ISSP)

```

/* Example : 2-level hierarchy, level-1 (root,
SP) and level-2 (leaf, WFQ). In addition,
rate-limiting is applied at level-2. */
1 weights = array[N];
2 finish_time = array[N];
3 srv_cntr = array[N];
4 pkts_per_rnd = array[N];
5 window_ID = array[N];
At Enqueue of each new packet (P) :
// rate-limiting
6 if srv_cntr[P.flow_ID] == pkts_per_rnd[P.flow_ID]
7 then
8     |   ++ window_ID[P.flow_ID];
9     |   srv_cntr[P.flow_ID] = 0;
10 ++ srv_cntr[P.flow_ID];
11 P.eligible_time = window_ID[P.flow_ID] *
12     window_period;
13 length_over_weight = P.length / weights[P.flow_id];
14 if finish_time[P.flow_id] > P.arrival_time then
15     |   finish_time[P.flow_id] = finish_time[P.flow_id]
16     |   + length_over_weight;
17 else
18     |   finish_time[P.flow_id] = P.arrival_time
19     |   + length_over_weight;
20 P.rank[0] = finish_time[P.flow_id]; // WFQ
21 P.rank[1] = P.ToS; // SP

```

References

- [1] M. Elbediwy, B. Pontikakis, A. Ghaffari, J.-P. David, and Y. Savaria, "Dr-pifo: A dynamic ranking packet scheduler using a push-in-first-out queue," *IEEE Transactions on Network and Service Management*, pp. 1–1, 2023.
- [2] M. Elbediwy, B. Pontikakis, J.-P. David, and Y. Savaria, "A hardware architecture of a dynamic ranking packet scheduler for programmable network devices," *IEEE Access*, 2023.
- [3] M. Shreedhar and G. Varghese, "Efficient fair queuing using deficit round-robin," *IEEE/ACM Transactions on networking*, vol. 4, no. 3, pp. 375–385, 1996.
- [4] M. Katevenis, S. Sidiropoulos, and C. Courcoubetis, "Weighted round-robin cell multiplexing in a general-purpose atm switch chip," *IEEE Journal on selected Areas in Communications*, vol. 9, no. 8, pp. 1265–1279, 1991.
- [5] J. Lakkakorpi, A. Sayenko, and J. Moilanen, "Comparison of different scheduling algorithms for wimax base station: Deficit round-robin vs. proportional fair vs. weighted deficit round-robin," in *2008 IEEE Wireless Communications and Networking Conference*. IEEE, 2008, pp. 1991–1996.
- [6] R. Sharmal, S. S. Sehra, and S. K. Sehra, "Review of different queuing disciplines in voip, video conferencing and file transfer," *International Journal of Advanced Research in Computer and Communication Engineering*, vol. 4, no. 3, pp. 264–267, 2015.
- [7] H. Rezaei, M. Malekpourshahraki, and B. Vamanan, "Slytherin: Dynamic, network-assisted prioritization of tail packets in datacenter networks," in *2018 27th International Conference on Computer Communication and Networks (ICCCN)*. IEEE, 2018, pp. 1–9.
- [8] I. A. Rai, E. W. Biersack, and G. Urvoy-Keller, "Size-based scheduling to improve the performance of short tcp flows," *IEEE network*, vol. 19, no. 1, pp. 12–17, 2005.
- [9] A. Demers, S. Keshav, and S. Shenker, "Analysis and simulation of a fair queueing algorithm," *ACM SIGCOMM Computer Communication Review*, vol. 19, no. 4, pp. 1–12, 1989.
- [10] K. Parekh, Abhay and G. Robert, G., "A generalized processor sharing approach to flow control in integrated services networks: the single-node case," *IEEE/ACM transactions on networking* 1, no. 3, pp. 344–357, 1993.
- [11] J. Y.-T. Leung, "A new algorithm for scheduling periodic, real-time tasks," *Algorithmica*, vol. 4, no. 1-4, p. 209, 1989.
- [12] M. Heusse, G. Urvoy-Keller, A. Duda, and T. X. Brown, "Least attained recent service for packet scheduling over wireless lans," in *2010 IEEE International Symposium on "A World of Wireless, Mobile and Multimedia Networks" (WoWMoM)*. IEEE, 2010, pp. 1–9.
- [13] S. J. Golestani, "A stop-and-go queueing framework for congestion management," in *Proceedings of the ACM symposium on Communications architectures & protocols*, 1990, pp. 8–18.
- [14] Y. Zhang and R. West, "End-to-end window-constrained scheduling for real-time communication," in *Proceedings of the 10th International Conference on Real-Time and Embedded Computing Systems and Applications (RTCSA'04)*, 2004.
- [15] Y. Zhang, R. West, and X. Qi, "A virtual deadline scheduler for window-constrained service guarantees," in *25th IEEE International Real-Time Systems Symposium*. IEEE, 2004, pp. 151–160.
- [16] N. K. Sharma, M. Liu, K. Atreya, and A. Krishnamurthy, "Approximating fair queueing on reconfigurable switches," in *15th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 18)*, 2018, pp. 1–16.
- [17] L. Yang, B. Ng, and W. K. Seah, "Heavy hitter detection and identification in software defined networking," in *2016 25th International Conference on Computer Communication and Networks (ICCCN)*. IEEE, 2016, pp. 1–10.
- [18] C.-C. Li, S.-L. Tsao, M. C. Chen, Y. Sun, and Y.-M. Huang, "Proportional delay differentiation service based on weighted fair queuing," in *Proceedings Ninth International Conference on Computer Communications and Networks (Cat. No. 00EX440)*. IEEE, 2000, pp. 418–423.
- [19] K. Nagaraj, D. Bharadia, H. Mao, S. Chinchali, M. Alizadeh, and S. Katti, "Numfabric: Fast and flexible bandwidth allocation in datacenters," in *Proceedings of the 2016 ACM SIGCOMM Conference*, 2016, pp. 188–201.