# Cutting Long-Tail Latency of Routing Response in Software Defined Networks

Junjie Xie, Deke Guo, Xiaozhou Li, Yulong Shen, and Xiaohong Jiang

Abstract—To enable the network softwarization, network func-1 2 tion virtualization (NFV) and software defined networking (SDN) are integrated to jointly manage and utilize the network resource 3 and virtualized network functions (VNFs). For a network flow 4 resulting from any NFV application, an associated switch would 5 send a routing request to the controller in SDN. The controller then generates and configures a routing path to dynamically steer the flow across appropriate VNFs or service function chains. 8 This process, however, exhibits a skew distribution of response 9 latency with a long tail. Cutting the long-tail latency of response 10 is critical to enable the network softwarization, yet difficult to 11 achieve due to many factors, such as the limited capacities and 12 the load imbalance among controllers. In this paper, we reveal 13 that such flow requests still experience the long-tail response 14 latency, even using the up-to-date controller-to-switch assignment 15 mechanism. To tackle this essential problem, we first propose a 16 light-weight and load-aware switch-to-controller selection scheme 17 to cut the long-tail response latency under the simple scenario 18 of homogeneous controllers, and then design a general delay-19 aware switch-to-controller selection scheme to fundamentally 20 cut the long-tail response latency for the more complicated 21 heterogeneous controller scenario with performance fluctuations. 22 The comprehensive evaluations indicate that our two new switch-23 to-controller selection schemes can significantly reduce the long-24 tail latency and provide higher system throughput. 25

Index Terms—Network softwarization, software defined
 networks, controller selection, long-tail latency.

Manuscript received September 30, 2017; revised February 5, 2018; accepted February 27, 2018. This work was supported in part by the National Natural Science Foundation for Outstanding Excellent Young Scholars of China under Grant 61422214, in part by the National Natural Science Foundation of China under Grant 61772544 and Grant U1536202, in part by the National Basic Research Program (973 program) under Grant 2014CB347800, in part by the Hunan Provincial Natural Science Fund for Distinguished Young Scholars under Grant 2016JJ1002, and in part by the Guangxi Cooperative Innovation Center of cloud computing and Big Data under Grant YD16507 and Grant YD17X11. (*Corresponding authors: Deke Guo; Yulong Shen.*)

J. Xie is with the Science and Technology on Information Systems Engineering Laboratory, National University of Defense Technology, Changsha 410073, China (e-mail: xiejunjie06@gmail.com).

D. Guo is with the College of System Engineering, National University of Defense Technology, Changsha 410073, China. He is also with the School of Computer Science and Technology, Tianjin University, Tianjin 300072, China (e-mail: guodeke@gmail.com).

X. Li is with the Department of Computer Science, Princeton University, Princeton, NJ 08544 USA (e-mail: xl@cs.princeton.edu).

Y. Shen is with the School of Computer Science and Technology, Xidian University, Xi'an 710071, China (e-mail: ylshen@mail.xidian.edu.cn).

X. Jiang is with the School of Systems Information Science, Future University Hakodate, Hokkaido 041-8655, Japan (e-mail: jiang@fun.ac.jp).

Color versions of one or more of the figures in this paper are available online at http://ieeexplore.ieee.org.

Digital Object Identifier 10.1109/JSAC.2018.2815358

#### I. INTRODUCTION

TETWORK softwarization is a transformation trend for designing, implementing, and managing the next generation networks. It exploits the benefits of software to enable the redesign of network and service architectures, optimize the expenditure and operational costs, and bring added values. The key enablers consist of the network function virtualization (NFV), software-defined networking (SDN) and cloud computing, etc [1]. Moreover, 5G systems will also rely on these technologies to attain system's flexibility and true elasticity [2], [3]. Network functions (NFs) are crucial for improving network security by examining and modifying network flows using special-purpose hardware. Recently, NFV has been proposed to execute virtual network functions (VNFs) on generic compute resources [4], such as commodity servers and VMs. Normally, a flow goes through specific VNFs in a particular order to meet its required processing, following the service function chain (SFC) [5], [6] along a routing path.

Additionally, SDN offers the freedom to refactor the control plane and flexibly enables the network softwarization [7]. More precisely, NFV and SDN can jointly manage the network resource and VNFs, and dynamically steer network flows across appropriate VNFs or SFCs. SDN centralizes the network control plane to a programmable software component, i.e., a controller running on a generic server, such as NOX [8]. The controller maintains a global network view and optimizes the forwarding decisions of network flows. For a flow from any NFV application, an associated switch would send a routing request to the controller. It is the controller that generates and configures a routing path to a specific VNF instance or to traverse a SFC on demand. The above process between a pair of switch and controller brings the response latency. Many factors would skew the tail of the latency distribution. For example, a single controller that lacks sufficient capacity to tackle received routing requests quickly and inevitably becomes a performance bottleneck [9]. Thus, such routing requests experience long-tail latency of response, as evaluated in Section II. Cutting the long-tail latency of routing response is critical to enable the network softwarization, yet difficult to achieve due to many factors.

To improve the scalability of SDN, the distributed control plane consisting of multiple controllers has been proposed recently [10], such as ONOS and OpenDaylight. To reduce the long-tail latency of response, they resort to the controllerto-switch assignment mechanism. That is, the control plane

0733-8716 © 2018 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.

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

137

138

168

178

proactively assigns a controller to each switch such that each 74 controller manages the same amount of switches. In reality, 75 the quantities of routing requests coming from switches per 76 unit time are different and dynamic. Consequently, controllers 77 still differ in the amount of received routing requests per time 78 unit. This load imbalance among controllers leads to the long-79 tail latency of response. Additionally, the controller-to-switch 80 assignment requires coordination among controllers, which 81 further aggravates the loads of controllers. 82

In this paper, to cut the long-tail latency of response and 83 lighten the load of controllers, we propose conducting the 84 selection of controllers at the side of each switch instead of the 85 controller-to-switch assignment. This means that each switch 86 actively chooses one controller from multiple available ones, 87 decoupling the static binding between switches and controllers. 88 More precisely, each switch prefers to adaptively select the 89 controller with low response latency for routing requests. This 90 would move the partial intelligences of the network to switches 91 and efficiently reduces the loads of controllers. 92

Despite those potential benefits, the selection of controllers 93 still faces many challenges. First, the switches need to probe 94 the state of controllers via the secure channel between them. 95 The secure channel is one kind of rare resource and affects the 96 performance of the whole network. To save the bandwidth of 97 the secure channel, the selection process of controllers should 98 be light-weight and use a few of the feedbacks from the 99 controllers. Second, the selection scheme needs to be scalable, 100 irrespective of the network size and the number of controllers. 101 Third, the selection scheme needs to accommodate the bursty 102 and skew routing requests from switches. Last, the selection 103 scheme should adapt to the heterogeneous controllers and the 104 performance fluctuation across controllers. 105

To tackle such challenges, we design a *load-aware* selec-106 tion scheme of controllers, which is simple but effective to 107 achieve the load balance among controllers. The load of a 108 controller refers to the number of routing requests waited 109 to be processed. The basic idea of our scheme is that each 110 switch sends routing requests to the controller with the lowest 111 load. In this way, all controllers process the similar number 112 of routing requests per time slot. This is very helpful to cut 113 the long-tail latency of routing response, when all controllers 114 have the same processing capabilities. This method alone, 115 however, is insufficient to deal with more general settings 116 of heterogenous controllers and the performance fluctuation. 117 Those controllers with lower processing capabilities still incur 118 the long-tail latency of response for routing requests, when all 119 controllers achieve the load balance. For this reason, we further 120 present a general delay-aware selection scheme of controllers 121 to fundamentally cut the long-tail latency of routing response. 122

Our delay-aware selection scheme includes two key com-123 ponents. The first one is the controller selection model of each 124 switch, which uses simple and inexpensive probing feedbacks 125 from a few controllers. It is still effective if each switch just 126 randomly probes two controllers and sends upcoming routing 127 requests to the controller with the shorter response delay. This 128 model is scalable and light-weight since it is not affected by 129 the network scale and the number of controllers. The second 130 component is the queue management mechanism of each 131

controller. It could estimate the response delay of a routing request and hence improve the performance predictability of controllers. The evaluation results reveal that our *delay-aware* selection scheme can efficiently reduce the long-tail latency of routing responses and improve the system throughput.

In summary, the major contributions of this paper are as follows.

- We reveal that routing requests experience the long-tail response latency, even using the up-to-date controllerto-switch assignment mechanism in SDN. Therefore, we propose an adaptive selection mechanism of controllers for switches to cut the long-tail latency.
- 2) We first design an efficient *load-aware* selection method 144 of homogeneous controllers for each switch. For more 145 general scenarios, we further propose a general delay-146 aware selection method, which is adaptive to the bursty 147 routing requests, the heterogenous controllers and the 148 performance fluctuations. Our two methods are light-149 weight as they limit the additional overhead caused by 150 probing two controllers. 151
- We further develop a queue management mechanism for each controller, which can efficiently manage the queue length and estimate the response delay of routing requests. The evaluation results reveal that our controller selection methods can accommodate system environment variations and efficiently reduce long-tail latency of routing response.

The paper is organized as follows. In Section II, we present 159 the observation of long-tail latency of routing response. 160 Section III depicts the framework of our controller selection 161 mechanism and the load-aware selection scheme of controllers. 162 We present the delay-aware selection scheme of controllers in 163 Section IV. We conduct massive experiments to evaluate the 164 performance of our controller selection schemes under various 165 system environment in Section V. Section VI introduces the 166 related work. In Section VII, we conclude this paper. 167

# II. LONG TAIL OF RESPONSE LATENCIES

In a SDN, when a switch receives a new flow, the switch 169 sends a routing request to its controller. The controller then 170 computes a route for the flow and inserts flow rules to related 171 switches in the route. Thus, the new flow would be forwarded 172 according to the flow rules in switches [11], [12]. Such an 173 interaction between the switch and the controller causes the 174 response latency. For a routing request, the latency of routing 175 response denotes the time interval from sending the routing 176 request to receiving the flow rules generated by the controller. 177

#### A. Long-Tail Observations of Response Latencies

Fig. 1(a) plots an observation about the long-tail distribu-179 tion under a single instance of ONOS controller. We build 180 a SDN testbed with one controller, running in a virtual 181 machine with 2 CPU cores and 2G RAM. Note that the 182 testbed forms a typical Fat-tree datacenter topology [13]. 183 We record the response latencies of 12,000 routing requests. 184 As shown in Fig. 1(a), the response latencies of 50% of routing 185 requests are lower than 5ms, and 90% of routing requests are 186



(a) A single controller with 12,000 routing(b) Multiple controllers with 150,000 routrequests. ing requests.

Fig. 1. Long-tail distributions of routing responses under a single controller as well as multiple controllers.

served within 30ms. However, there still exist some routing
 requests whose response latencies are more than 50ms. That
 is, the response latencies exhibit a long-tail distribution.

Furthermore, we observe the response latencies of rout-190 ing requests under multiple controllers [14]. We employ 191 300 switches and 40 controllers where each controller 192 manages 7 or 8 switches. Each switch generates routing 193 requests according to a Poisson arrival process with  $\lambda = 0.5$ 194 during 1ms. The processing time of each request in each 195 controller is drawn from an exponential distribution where 196  $\mu^{-1}$  = 2ms. Each controller can process 10 requests in 197 parallel. We run the system 1000ms and record the response 198 latencies of routing requests. The number of arrival routing 199 requests is about 150,000. In Fig. 1(b), 89% of routing 200 requests can be served in 5ms, and the response latencies 201 of 96% of routing requests are lower than 10ms. However, 202 some response latencies are still larger than 20ms. That is, 203 Fig. 1(b) shows that there still exists a long-tail distribution 204 under multiple ONOS controllers. 205

# 206 B. Analysis About the Long-Tail Latencies

Fig. 1(a) results from that the network only employs one 207 controller. Due to the limited capability of the single controller, 208 a large amount of requests have to queue in the controller. 209 Therefore, there is a long-tail latency of routig response caused 210 by the long queueing delays. In Fig. 1(b), the number of 211 routing requests that each switch generates is different, even 212 though they obey the same Poisson distribution. The skew-213 flow requests will make that some controllers are overload, but 214 other controllers would be underutilized. As a consequence, 215 there will be a long-tail latency. We illustrate the problem 216 in Fig. 2. Two controllers are assigned to four switches and 217 have the same processing time of 4ms. Assume  $switch_1$  and 218 switch<sub>2</sub> receive 4 requests each and that switch<sub>3</sub> and switch<sub>4</sub> 219 receive 2 requests each. The requests received by switch<sub>1</sub> 220 and  $switch_2$  can only be processed by  $controller_1$ , which is 221 assigned to manage them. This leads to a maximum latency of 222 32ms, but a load-aware selection obtains a maximum latency 223 of 24ms. Fig. 2 shows that a quantity-based assignment strat-224 egy leads to long-tail latencies because it fails to accommodate 225 the skew-flow requests. 226

227 Quantity-based allocation strategy is commonly employed 228 by many controllers to balance the loads of controllers, such as



Fig. 2. The performances of quantity-based controller allocation and loadaware controller selection.

ONOS [14]. That is a controller-to-switch assignment mech-229 anism, which is abbreviated as the assignment mechanism 230 of controllers. In this case, controllers coordinate to manage 231 switches, and each controller manages an approximately equal 232 quantity of switches. When deploying multiple instances of 233 ONOS in a SDN, the bursty flows from a switch are sent to the 234 same controller. Consequentially, a large number of requests 235 have to queue in the controller. These queueing requests tend 236 to incur long response latencies. However, those controllers, 237 which do not receive bursty routing requests, may even be 238 underloaded. In conclusion, the assignment mechanism fails 239 to efficiently reduce the tail latencies of responses. 240

# III. FRAMEWORK OF CONTROLLER SELECTION MECHANISM

241 242

250

To overcome the drawback of assignment mechanism, we design a switch-to-controller selection scheme, which is abbreviated as the selection scheme of controllers. The selection scheme moves partial intelligences of the network to switches and relieves the loads of controllers. Meanwhile, the selection scheme can efficiently reduce tail latencies of responses. 249

#### A. Overview of Controller Selection Mechanism

For existing designs of control plane, the assignment mech-251 anism of controllers is a static binding between switches and 252 controllers, which fails to deal with the bursty and skew 253 routing requests, and further incurs long response delays. 254 To reduce the tail latencies of responses, routing requests from 255 the same switch need to be processed by appropriate con-256 trollers. This means that the static binding between switches 257 and controllers needs to be decoupled. A better mechanism is 258 to enable switches to select controllers for routing requests. 259

For the assignment mechanism, the load balance among 260 controllers means that each controller manages the same 261 amount of switches. The load is denoted by the number 262 of switches. Moreover, the load balance among controllers 263 requires the coordination of controllers. The coordination 264 will further aggravate the computing and communication 265 overhead. However, in this paper, the selection scheme of 266 controllers achieves the mapping between switches and con-267 trollers through switches conduct simple and actively probing. 268 Therefore, the selection scheme relieves the overhead of coordination and assignment in controllers.

In SDN, the controller periodically probes the switch state, 271 and switches would respond to those probes and forward any 272 switch notifications (e.g., link failures or recovery) to all live 273 controllers [15]. When a switch suffers a failure scenario, 274 it fails to respond to the probe. After that, the controller 275 promptly updates its local topology state and replicates that 276 update to all other controllers in the cluster. In addition, there 277 have been lots of researches on how to achieve the consistency 278 among distributed controllers [15]. These techniques are com-279 plementary to our approach. 280

We design the controller selection mechanism keeping in mind these four goals:

 Light-weight: A light-weight probing method is needed to save the bandwidth of the secure channel. The probing for the controller selection uses the secure channel, which is the communication channel between the control plane and the data plane. The bandwidth of the channel can affect the performance of the whole network.

289
 2) Scalable: The selection scheme of controllers should be irrelevant to the increasing number of controllers. The expansion of network size is common. The method of probing controller should accommodate the increase of deployed controllers and avoid to incur the communica-tion overhead and the computing overhead.

 Burst-immunity: The selection scheme should be burstimmune. There are bursty and skew-flow requests, which
 can lead to long response delays. To shorten the tail
 latency of responses, the selection scheme needs to
 accommodate the bursty and skew-flow requests.

4) Adaptive: The selection scheme of controllers should be adaptive. The capabilities of controllers may be heterogeneous and time-varying. To deal with the general situation, the selection scheme must cope and quickly react to heterogeneous and time-varying processing capabilities across controllers.

#### 306 B. Load-Aware Selection Scheme of Controllers

Accommodating skew-flow requests across controllers 307 necessitates a selection strategy of controllers. The strategy 308 can make switches select faster controllers for routing requests. 309 The controller with fewer unfinished requests can respond to 310 routing requests faster when controllers have the same process-311 ing capability. In this paper, we first present a load-aware 312 selection scheme. To realize this framework, the selection 313 strategy needs to take into account the loads across multiple 314 controllers in the network. The load means the number of 315 unfinished requests. Our load-aware selection scheme is to 316 select a controller with the fewest unfinished requests for 317 newly generated requests. 318

Under the load-aware selection scheme of controllers, the switch needs to send a probing request to each controller after a switch receives a new flow. When the controller receives the probing request, it will return the number of unfinished requests to the switch. After the switch receives all probing results, it then sends the new routing request to the controller with the lightest load since that controller can respond the routing request fastest. The load-aware selection scheme aims to reduce tail latencies of responses by selecting the controller with the lightest load.

In Fig. 2, the load-aware selection scheme will work as 329 follows. When a switch receives a new flow, it first probes the 330 loads of controllers A and B. Based on the loads of controllers 331 A and B, the switch sends the routing request to the controller 332 with lightest load. This scheme can balance the loads of 333 controllers A and B and can achieve better selection. When the 334 controller finishes processing the routing request, it inserts the 335 flow rules to related switches. Lastly, those switches will deal 336 with the flow according to the actions of matched flow rules. 337 In addition, the new arrival routing requests need to queue in 338 controllers, when controllers are busy. However, the infinite 339 length of queue will incur infinite response delays of routing 340 requests. To cut the tail latency of response, it is necessary to 341 ensure that the length of queue is finite in each controller. 342

# C. The Condition to Finite Queue Length

We give the condition to achieve that the expected number 344 of requests in per controller remains finite for all time. Con-345 sider the following model: requests arrive as a Poisson stream 346 of rate  $\lambda$  at each switch. Requests are processed according 347 to the first-in first-out (FIFO) protocol by controllers. The 348 processing time for a request is exponentially distributed with 349 mean  $\mu$ . When there are *m* switches and *n* controllers in a 350 network, requests arrive as a Poisson stream of rate  $\frac{\lambda m}{n}$  at 351 each controller. We obtain the following theorem. Note that 352  $\frac{\lambda m}{m} < \mu$ , and then the system will be stable, which means that 353 the expected number of requests per controller remains finite 354 in equilibrium. Theorem 1 shows that the system is stable for 355 every  $\frac{\lambda m}{\mu n} < 1$ ; that is, the expected number of requests in each 356 controller remains finite for all time. 357

343

372

373

*Theorem 1:* The system is stable for every  $\frac{\lambda m}{\mu n} < 1$ ; that is, the expected number of requests in each controller remains finite for all time.

Proof: When we treat all controllers as a whole and 361 all switches as a whole, then the system can be seen as 362 a M/M/1 system with Poisson arrival rate  $\lambda m$  and average 363 service rate  $\mu n$ . Let  $P_k(t)$  denote the probability of that there 364 are k requests in the whole system in time t. Accordingly, 365  $P_{k+1}(t)$  denotes the probability where k + 1 requests exist in 366 the system in time t, and  $P_k(t + \Delta t)$  denotes the probability of 367 that there are k requests in the whole system in time  $t + \Delta t$ . 368 Now we consider the evolution of the system. In the time 369  $[t, t + \Delta t]$ , the process of evolution has some attributes just 370 as follows: 371

- One request comes with the probability  $\lambda m \Delta t$ , and the probability of no request comes is  $1 \lambda m \Delta t$ .
- One request departures with the probability  $\mu n \Delta t$ , and the probability of no request departure is  $1 - \mu n \Delta t$ . 375
- The situation of more than one request comes and departures in  $\Delta t$  is a small probability event, and it can be ignored. 376

There are 4 types of evolution process, and we list them <sup>379</sup> in Table I. Take type B for an example. Since *k* is the number <sup>380</sup>

TABLE I Four Types of Evolution Process

Туре	state in time t	come	departure	state in time $(t + \triangle t)$
(A)	k	0	0	k
(B)	k+1	0	1	k
(C)	k-1	1	0	k
(D)	k	1	1	k

of requests in the whole system, so k + 1 means that there 381 are k + 1 requests in the system in time t. After that, when 382 a request departs during  $\Delta t$ , there would be k requests in 383 the system in time  $t + \Delta t$ . Accordingly, we can get the 384 possibility of type A is  $P_k(t)(1-\lambda m \Delta t)(1-\mu n \Delta t)$ , possibility 385 of type B is  $P_{k+1}(t)(1 - \lambda m \triangle t) \mu n \triangle t$ , possibility of type 386 C is  $P_{k-1}(t)\lambda m \Delta t (1 - \mu n \Delta t)$  and possibility of type D is 387  $P_k(t)\lambda m \Delta t \mu n \Delta t$ . Then,  $P_k(t + \Delta t)$  should be the sum of all 388 4 types, shown in Equation (1). 389

$$P_k(t + \Delta t) = P_k(t)(1 - \lambda m \Delta t - \mu n \Delta t)$$
  

$$P_{k+1}(t)\mu n \Delta t + P_{k-1}(t)\lambda m \Delta t + o(\Delta t) \quad (1)$$

And let  $\Delta t \rightarrow 0$ , we can get a differential equation, shown in Equation (2).

<sup>394</sup> 
$$\frac{dP_k(t)}{dt} = \lambda m P_{k-1}(t) + \mu n P_{k+1}(t) - (\lambda m + \mu n) P_k(t) \quad k = 1, 2, \dots$$
 (2)

Noted that if k = 0, there will exist only type A and type B, shown in Equations (3) and (4).

<sup>398</sup> 
$$P_0(t + \Delta t) = P_0(t)(1 - \lambda m \Delta t) + P_1(t)(1 - \lambda m \Delta t)\mu n \Delta t$$
 (3)

<sup>399</sup> 
$$\frac{dP_0(t)}{dt} = -\lambda m P_0(t) + \mu n P_1(t)$$
 (4)

We just have interest in the equilibrium point, and the derivative is 0 in fixed point. Thus, we get Equation (5).

402

$$\begin{cases} -\lambda m P_0 + \mu n P_1 = 0\\ \lambda m P_{k-1} + \mu n P_{k+1} - (\lambda m + \mu n) P_k = 0 \quad k \ge 1 \end{cases}$$
(5)

Resolve equation (5), we can get  $P_k = (\lambda m/\mu n)^k P_0$ . If  $\frac{\lambda m}{\mu n} < 1$ , then the sequence  $P_k$  will be decrease. And we know probability is non-negative, that means  $P_k \ge 0$ . If a sequence is bounded and monotone, it converges [16]. So there exist *K*, when k > K,  $P_k = 0$ . Then the expected total number of requests in all controllers remains finite.

Theorem 1 shows that the expected total number of requests 409 in each controller remains finite, when  $\frac{\lambda m}{\mu n}$  < 1. Therefore, 410 to achieve the finite queue length, it is essential to ensure that 411  $\frac{\lambda m}{\mu n}$  < 1. When the network size *m* increases, if  $\frac{\lambda m}{\mu n} \ge 1$ , it is 412 necessary to increase the number of deployed controllers n, 413 otherwise, the queue length of some controllers will be infinite, 414 and that will incur infinite response delays. Another method 415 to limit the queue length of controllers is to drop some routing 416 requests, which can limit the value of  $\lambda$  and achieve  $\frac{\lambda m}{un} < 1$ . 417

#### 418 D. Limitations of Load-Aware Selection Scheme

Fig. 3 shows the Cumulative Distribution Function (CDF) of response latencies of routing requests under different schemes.



Fig. 3. Response latencies of routing requests under different schemes.



Fig. 4. Distinct selection schemes incur different response latencies.

In Fig. 3, the response latencies of 94% of routing requests 421 are lower than 5ms after adopting the load-aware selection 422 scheme. However, for the quantity-based assignment, only 423 86% of routing requests can be responded in 5ms. In Fig. 3, 424 all routing requests can be responded in 10ms under the load-425 aware scheme. Therefore, Fig. 3 indicates that our load-aware 426 selection scheme can reduce the tail latency of responses than 427 the prior quantity-based allocation method when controllers 428 are homogeneous and exhibit the same processing capabilities. 429 However, we can see that both curves (Quantity-based and 430 Load-aware) are very close to each other, which means that 431 the load-aware selection strategy narrowly reduce the long-432 tail latency. Furthermore, the load-aware selection scheme 433 faces three challenges. First, controllers are heterogeneous. 434 Second, the processing capabilities of controllers are dynam-435 ically changing. Third, the cost of probing is too huge. 436 Therefore, the load-aware selection scheme is still insufficient 437 to completely cut the tail latency. 438

Fig. 4 plots an illustrative example of the limitation. For two 439 controllers, the processing time per request in controller A and 440 controller B are 5ms and 12ms, respectively. Assume all four 441 switches receive a burst of 3 requests each. Each request needs 442 to be forwarded to a single controller. If every switch selects 443 a controller using the load-aware scheme, it will result in each 444 controller receiving an equal share of the requests. This leads 445 to a maximum latency of 72ms, whereas an ideal selection in 446 this case obtains a maximum latency of 45ms. We note that 447 the load-aware scheme will prefer faster controllers over time, 448 but purely relies on the load information. Therefore, when 449 controllers are heterogeneous, the load-aware selection scheme 450 can not efficiently shorten the tail latency. 451

Controllers are commonly heterogeneous for primarily three reasons. First, the hardware is heterogeneous. Controllers run

488

489

in commercial servers. These servers can be heterogeneous 454 due to different hardware configurations, such as CPU and 455 memory. Second, the software is heterogenous. There are 456 multiple different controllers developed by different organiza-457 tions [10], such as NOX, Beacon, Floodlight, Ryu, ONOS and 458 OpenDaylight, etc. Those controllers themselves have different 459 performances. Third, the function is heterogeneous. There 460 are some management applications running in controllers for 461 achieving different functions [12], [17], and these applications 462 will consume some resources of controllers. Consequentially, 463 controllers have different remaining capabilities for processing 464 routing requests, even if the controllers run in servers with 465 the same setting. In this case, queueing routing requests in 466 controllers with low processing capabilities will lead to long 467 response latencies. 468

Additionally, the load-aware scheme probes the loads of all 469 controllers, and then selects the controller that has the lightest 470 load. However, this probing will incur the overhead of com-471 munication and aggravate the loads of controllers when there 472 are a large number of controllers in a large-scale network. For 473 example, there are m switches and n controllers in a network. 474 Suppose that each switch receives  $\lambda$  routing requests in 1ms. 475 There are  $2\lambda \times m \times n$  times communications between switches 476 and controllers during 1ms. Meanwhile, each controller needs 477 to evaluate its own load  $\lambda m$  times in 1ms. The cost of probing 478 is too huge for the load-aware selection scheme. 479

The load-aware controller selection scheme can only reduce 480 tail latencies of responses under the homogeneous controllers. 481 However, the network environment is time-varying in real 482 situations, not only in the processing capabilities of controllers 483 but also in the number of routing requests from switches. 484 We further propose a delay-aware selection scheme of con-485 trollers, which can adapt to the variations of the network 486 environment. 487

# **IV. DELAY-AWARE SELECTION SCHEME** OF CONTROLLERS

We design the delay-aware selection scheme of controllers 490 while keeping the design goals of controller selection mecha-491 nism in mind. We first show an overview of the delay-aware 492 selection scheme. Then, we present two major components of 493 the delay-aware selection scheme, the selection models of con-494 trollers and the queue management mechanism of controller. 495

#### A. Overview of Delay-Aware Selection Scheme 496

To address those problems faced by the load-aware selection 497 scheme, we further design the delay-aware selection scheme of 498 controllers, which is adaptive to the heterogeneous controllers 499 as well as to the dynamic behaviours of flows. The delay-500 aware selection scheme needs to probe the response delays of 501 controllers for routing requests and send the routing requests 502 to the controller with the smallest response delay. The latency 503 of routing response denotes the time interval from sending 504 the routing request to receiving the flow rules generated by 505 the controller and is composed of the queueing delay and 506 the processing delay, as shown in Definition 1. The response 507 delay is an approximate evaluation of response time. Through 508



Fig. 5. Overview of controller selection scheme. CS: Controller Selection scheduler, QM: Queue Management of controller.

probing response delays, the delay-aware selection scheme 500 can accommodate the heterogeneous controllers, while fewer 510 routing requests will be sent to the controllers with low 511 processing capabilities. 512

Definition 1: The latency of routing response denotes the 513 time interval from sending the routing request to receiving the 514 flow rules generated by the controller.

515

Furthermore, the capabilities of controllers are time-varying. 516 With the development of SDN, there are more and more 517 applications running in controllers. When switches send vast 518 requests to the controller that has fast capability of response 519 at before, a large number of requests have to queue in 520 controllers if the capabilities of controllers decrease due to 52 other applications' overconsumption of resources. 522

Our delay-aware selection scheme includes two major com-523 ponents, controller selection (CS) and queue management 524 (QM). Recall the design goal of the selection scheme in 525 Section III-A, CS can achieve that the selection scheme is 526 light-weight, scalable and burst-immune, and QM achieves 527 the goal of adaptivity. First, we design a selection scheme 528 of controllers, which can select the controller based on a 529 little feedback from the controllers, and thus, is light-weight. 530 Second, instead of probing all controllers, the switch randomly 531 probes d controllers where  $d \ge 1$ . The probing is scalable 532 and independent of the network size. Third, the selection 533 scheme can make switches conduct once controller selection 534 for each flow or a batch of flows. It can make those requests 535 be processed by different controllers, and thus can avoid the 536 influence of the bursty and skew-flow requests. Last, through 537 estimating response delays of routing requests, switches can 538 send routing requests to the controller that has the smallest 539 response delay. Based on this estimation, the selection of 540 controllers can be adapted to the heterogeneous and time-541 varying processing capabilities. 542

Fig. 5 depicts the framework of the adaptive switch-to-543 controller selection scheme. When a request is issued at a 544 switch, the switch will work based on Algorithm 1. The switch 545 randomly probes d controllers, where  $d \ge 1$ . The d controllers 546 then estimate their response delays for the routing request 547 based on Algorithm 2 and return the response delays  $\psi$  to 548 the switch. If  $\psi_i$  of controller *i* exceeds the max response 549 delay  $RD_{max}$  limit, then controller *i* will return the response 550 delay  $\psi_i = -1$ . When the switch receives response delays 551 of d controllers, it will select the controller that has the 552 smallest response delay. If all response delays are lower than 0, 553 the switch will reselect a controller whose queue length does 554

	- 7	

Algorithm 1 The Selection of Controllers
<b>Require:</b> Controller set C, d.
1: randomly probe d controllers from C;
2: send estimating request to the $d$ controllers;
3: $\psi \leftarrow$ response delays of <i>d</i> controllers;
4: if there exists $\psi_x \ge 0$ then
5: $id \leftarrow arg \min\{\psi_x \ge 0\};$
6: <b>else</b>
7: for $i = 1$ to C.length do
8: send estimating request to controllers $C[i]$ ;
9: $\psi_0 \leftarrow$ response time of $C[i]$ ;
10: <b>if</b> $\psi_0 >= 0$ <b>then</b>
11: $id = i;$
12: <b>break</b> ;
13: send the routing request to controller $C[id]$ .

Algorithm 2 Queue Management of Controller
<b>Require:</b> the max response delay, $RD_{max}$ .
1: receive an estimating request from switch s;
2: calculate the average processing time of requests $\bar{\nu}$ ;
3: if $r_i < \gamma_i$ then
4: $\psi_i \leftarrow \frac{r_i}{v_i} \bar{v};$
5: else
6: $\psi_i \leftarrow C[q_i mod \gamma_i] + (\lfloor q_i / \gamma_i \rfloor + 1) \times \overline{\nu};$
7: if $\psi_i > RD_{max}$ then
8: $\psi_i = -1;$
9: send $\psi_i$ to switch s.

not exceed limit for the routing request. Last, the switch willsend the request to the selected controller.

# 557 B. The Selection Models of Controllers

When there are only a few controllers in the network, it is feasible to probe all controllers. Considering that this probing could occupy extra bandwidth of the secure link, it is essential to design a per-flow light-weight probing method.

1) Active Per-Flow Selection of Controllers: To deal with 562 the skew-flow requests and reduce response tail latencies, 563 the switches need to select controllers for each flow. Instead 564 of the controller-to-switch assignment, the active per-flow 565 selection of controllers makes that the routing requests from 566 the same switch can be processed by different controllers. 567 This selection scheme can fully exploit the capabilities of 568 controllers and efficiently reduce response tail latencies. 569

To reduce the bandwidth consumption of probing con-570 trollers, one method is to reduce the number of probed 571 controllers. There is a tradeoff between response tail latencies 572 and the number of probed controllers. Probing more controllers 573 can achieve fewer response tail latencies. However, that also 574 means more bandwidth consumption and computing overhead 575 in controllers. The number of controllers increases as the 576 network scale grows. In this case, checking all controllers has 577 a huge cost. To achieve the light-weight probing, we randomly 578 probe d controllers instead of checking all controllers, where 579  $d \geq 1$ . Furthermore, Azar *et al.* [18] have shown that having 580

just two random choices (i.e., d = 2) yields a large reduction in the maximum load over having one choice. This method has been widely studied and applied [19]. Inspired by this fact, our active per-flow selection is to probe two controllers and is thus scalable. Meanwhile, the active per-flow selection of controllers can efficiently reduce the bandwidth consumption of the secure link and the computing load of controllers. 581

Instead of probing the loads of the controllers, probing 588 response delays of controllers can better reduce response tail 589 latencies. The probing of response delays requires that these 590 controllers evaluate their own response delays for routing 591 requests. Since the selection of controllers only needs to get a 592 numerical value of response delay, the selection of controllers 593 is light-weight. Moreover, the heterogeneous and time-varying 594 processing capabilities of controllers increase the complexity 595 of evaluation for response delays of controllers. The response 596 delay estimate model will be introduced in Section IV-C. 597

2) Active Selection of Controllers for a Batch of Flows: 598 When switches meet bursty-flow requests or when the arrival 599 of routing requests is frequent, conducting a controller selec-600 tion for each request still aggravates the bandwidth consump-601 tion and the loads of controllers even if only two controllers 602 need to be probed in one controller selection. Conducting the 603 controller selection for a batch of arrival routing requests is 604 needed to increase the scalability of the controller selection 605 mechanism, when switches suffer bursty flows. 606

For active selection of controllers for a batch of flows, 607 the switch conducts one controller selection after it receives 608 the first flow request. When we set the batch size as  $\delta$ , 609 it means that the following  $\delta - 1$  requests will be sent to 610 the same controller with the first request. That is, the result 611 of controller selection for the first request will be shared by  $\delta$ 612 requests. In addition, the batch selection is irrelevant to the 613 rate of requests because  $\delta$  denotes the number of requests. 614 Therefore, although there would be the high rate of requests 615 in the beginning when the switch changes its controller, those 616 requests would not be sent to the same controller. Given that 617 the request arrival process at each switch is a Poisson process 618 with rate  $\lambda$ , the arrival duration for a flow is exponentially 619 distributed with mean  $1/\lambda$ . Therefore, the arrival duration of  $\delta$ 620 flows is also exponentially distributed with mean  $\delta/\lambda$ . 621

There is a tradeoff between the rounds of controller selection and the performance of controller selection. If  $\delta$  is too small, it is obvious that controller selection should be frequently conducted. However, it will decrease the performance of controller selection when  $\delta$  is too large. It is worth noting that it is unnecessary to adopt the batch selection when the arrival of flows is scattered.

# C. Queue Management Mechanism of Controller

The queue management of controller makes it so that the selection of controllers can cope and quickly react to heterogeneous and time-varying processing time across controllers. Our queue management mechanism of controller includes response delay estimate and queue length bound.

1) Response Delay Estimate Model: As depicted in 635 Section III-D, the load-aware selection scheme of controllers 636

can not accommodate the heterogeneity of controller. To efficiently reduce the long-tail latency of routing response,
switches should select controllers with lower response delays
for each routing request.

Request response time consists of the queueing time and 641 the processing time. Furthermore, the queueing time is related 642 to the length of queue, which is equal to the number of 643 queueing requests. Meanwhile, to estimate the queueing time, 644 it is essential to estimate the processing time of each request. 645 In our design, the controller records  $v_i$ , which is the processing 646 time of the latest responded *j*th requests. Given the number 647 of the latest finished requests s, we calculate  $\bar{\nu}$ , which denotes 648 the average processing time of s requests in controller i. Thus, 649  $\bar{\nu} = \frac{1}{s} \sum_{i=1}^{s} \nu_i$ . We use  $\bar{\nu}$  to estimate the processing time of 650 requests. 651

Consider that the controller can process multiple routing 652 requests simultaneously. We use  $\gamma_i$  to denote the number of 653 requests that *controller*<sub>i</sub> can simultaneously process.  $q_i$  and  $r_i$ 654 are the number of queueing and running requests in con-655 troller *i*, respectively. To improve the system utilization, 656 the controller that has idle running slots should have a lower 657 estimated response delay. Therefore, the estimated response 658 delay  $\psi_i = \frac{r_i}{\gamma_i} \bar{\nu}$  when  $r_i < \gamma_i$ . When there are requests 659 queueing in a controller, the controller records the running 660 duration A[k] of the running request in the kth slot where 661  $1 \le k \le \gamma_i$ . The controller then estimates that the queueing 662 request will run in which slot. To achieve this goal, we use 663  $B[k] = |\bar{v} - A[k]|$  and then sort B[k] as non-decreasing order. 664 Then, the controller estimates that the request will run in 665  $[(q_i \mod \gamma_i) + 1]$ th slot. We can get the queueing time 666  $wt_i = B[(q_i \mod \gamma_i) + 1] + (\lfloor q_i / \gamma_i \rfloor) \times \overline{\nu}$ . At last,  $\psi_i = wt_i + \overline{\nu}$ 667 when  $r_i = \gamma_i$ . 668

In summary,  $controller_i$  uses the following estimation function for response delay:

$$\psi_{i} = \begin{cases} B[(q_{i}mod\gamma_{i})+1] + (\lfloor q_{i}/\gamma_{i}\rfloor+1) \times \bar{\nu} & r_{i} = \gamma_{i} \\ \frac{r_{i}}{\gamma_{i}}\bar{\nu} & r_{i} < \gamma_{i} \end{cases}$$
(6)

When a switch sends an estimating request to  $controller_i$ , 672 the *controller<sub>i</sub>* adopts the formula (6) to estimate the response 673 delay. In general,  $\gamma_i = 1$  means that the controller only can 674 process one request once. In this case,  $\psi_i = B[1] + (q_i + 1) \times \overline{\nu}$ . 675 We suppose that the controller is empty at the beginning. After 676 that,  $\bar{\nu}$  is equal to the average processing time of finished 677 requests when the number of finished messages is lower than 678 the given threshold s. 679

2) Cutting Tail Latencies: Since switches conduct con-680 troller selection simultaneously, there may be "herd behav-681 iors," wherein multiple switches are coaxed to direct requests 682 towards the best controller. There are many requests queueing 683 in a controller under herd behavior that could leads to long-tail 684 latencies of routing responses. Moreover, it is possible that the 685 probed controllers all have low processing capabilities or long 686 queues. In this case, it is not suitable to select a controller 687 from those probed controllers. 688

To cut long-tail latencies and reduce the influence of herd behavior, the controller necessitates to bound its queue length. Determining the length of queues at controllers is crucial. Queues that are too short lead to lower controller utilization, as resources may remain idle between allocations. Queues that are too long may incur excessive queuing delays.

When fewer requests are sent to a controller, this may incur 695 under-utilization of its resources, whereas significant delays 696 may occur when requests need longer processing time. Hence, 697 after estimating request response delay, we further design 698 a delay-aware bounding mechanism to bound queue length, 699 which can accommodate the heterogeneous and time-varying 700 capabilities of controllers. Meanwhile, bounding queue length 701 can efficiently weaken the influence of herd behaviors. At one 702 point, a controller receives a burst of flows, and that exceeds 703 the limit of queue length. After that, the following flows will 704 not queue in the controller until the queue length is lower than 705 the limit. This delay-aware bounding mechanism relies on the 706 response delay estimation of request, which is reported by the 707 controller. 708

In particular, we specify the maximum response delay 709  $RD_{max}$  that a request is allowed to wait in a queue. When 710 we are about to place a request at the queue of  $controller_i$ , 711 we first check the estimated response delay  $\psi_i$  reported by 712 *controller<sub>i</sub>*. Only if  $\psi_i < RD_{max}$  is the request queued at 713 that controller. We sample d controllers while conducting the 714 controller selection. If the d selected controllers all do not sat-715 isfy the maximum response delay constraint. The switch needs 716 to reselect a new controller. Using this method, the number of 717 requests in each controller gets dynamically adapted based on 718 the current capability of the controller. 719

 $RD_{max}$  is set to make requests prefer faster controllers. 720 Furthermore,  $RD_{max}$  can be dynamically regulated to fit the 721 variation of controllers' capabilities. For controllers that have 722 low processing capabilities,  $RD_{max}$  can limit the number 723 of queueing requests in these controllers. After that, these 724 requests can be sent to faster controllers. However, if  $RD_{max}$ 725 is too small, most of controllers refuse to receive new requests 726 because their response delays exceed the limit of  $RD_{max}$ . 727 In this case, it is necessary to extend the value of  $RD_{max}$ . 728

#### V. PERFORMANCE EVALUATION

729

736

We start with the evaluation methodology and scenarios. This section, we evaluate the selection schemes of controller and the assignment mechanism of controller under the general settings of controllers, the queue length bound  $RD_{max}$ , The heavy request-skews, the time-varying service rates and the batch selection. The setting of the batch selection. The setting setting service rates are setting setting. The setting settin

#### A. Experimental Setup

We build a discrete-event simulator wherein workload 737 generators create flow requests at a set of switches. These 738 switches then employ the controller selection scheme to select 739 a controller for each request. Unless otherwise specified, 740 the network consists of 300 switches and 40 controllers. 741 However, when a large amount of flow requests emerges, 742 we need to employ more controllers to deal with the burstly 743 flows. The workload generators create flow requests according 744 to a Poisson arrival process to mimic arrival of user requests 745 at web servers [20]. Unless otherwise specified, the Poisson 746

arrival process is with  $\lambda = 0.5$  during 1ms. At the beginning, 747 the system is empty, and there are no requests. With the system 748 running, the switches start to produce flow requests and select 749 controllers for flow requests. We run the system 1000ms, 750 and the number of arrival flow requests is about 150,000. 751 Each controller maintains a FIFO request queue. Moreover, 752 in the settings of controllers, each controller can service a 753 tunable number of requests in parallel (10 in our settings). 754 The processing time each request experiences is drawn from 755 an exponential distribution (as in [21]) with a mean processing 756 time  $\mu^{-1} = 2$ ms. Furthermore, we incorporate controller 757 heterogeneity into the network as follows: each controller, 758 independently and with a uniform probability, sets its service 759 rate to a different  $\mu$ , where  $\mu = 0.5$  or  $\mu = 0.1$  in our settings. 760 To estimate the response delay of each request, we set s = 100, 761 That is,  $\bar{\nu}$  denotes the average processing time of the latest 762 finished 100 requests. We repeat every experiment 20 times 763 using different random seeds, and then get the average result. 764

<sup>765</sup> We compare our design against three strategies:

 Quantity-based allocation: Controllers achieve load balance through balancing the number of switches, which each controller manages. Currently, ONOS controller utilizes this strategy to balance the loads of distributed controllers.

- 2) Load-aware selection: The switch probes two controllers for each request and sends the request to the controller with fewer number of requests because probing all controllers is not scalable.
- 3) Delay-aware selection: The switch probes two controllers for each request and gets request response delays, which rely on the feedbacks of probed controllers. Then, the switch sends the request to the controller with smaller response delay.

### 780 B. The Impact of d

Azar *et al.* [18] have shown that the situation of d = 2 yields 781 a large reduction in the maximum load over d = 1, while 782 each additional choice beyond two decreases the maximum 783 load by just a constant factor. Further, to verify the theory 784 and determine the value of d, we evaluate the impact of 785 d on the performance of the delay-aware selection strategy 786 under heterogeneous controllers where d denotes the number 787 of probed controllers. The processing time of each request 788 in each controller is drawn from an exponential distribution 789 where  $\mu$  was randomly set as 0.5 or 0.1. Other parameters are 790 the same as Section V-A. 791

Fig. 6 shows the delay-aware selection strategy significantly 792 reduces the mean response time as the value of d increase 793 from 1 to 2. However, probing more controllers just incur a 794 little bit reduction on the mean response time after d exceeds 795 2. In addition, multiple switches may simultaneously select 796 the same controller if each of them randomly probes more 797 controllers. This would in turn aggravates the load of the 798 controller. Thus, we set d = 2 in the next experiments. 799

#### 800 C. General Settings of Controllers

We evaluate the performances of different schemes with heterogeneous controllers. We employ 60 controllers where the



Fig. 6. The impact of d on the performance of the delay-aware selection strategy.



(a) Request response time with different (b) Throughput under different schemes.

Fig. 7. The performances of different schemes where controllers are heterogeneous.

bound of maximal response delay  $RD_{max} = 20$ ms. The other settings of experiments are consistent with that of Section V-B.

Fig. 7(a) shows that our delay-aware scheme can signifi-805 cantly reduce the response tail latencies. Basically, all requests 806 can be finished in 50ms while adopting our delay-aware 807 scheme. The load-aware scheme achieves better performance 808 than the quantity-based scheme in Fig. 7(a). Over 90% of 809 requests can be processed during 150ms based on the load-810 aware scheme. The quantity-based scheme leads to long 811 response delays, and there are over 20% of requests whose 812 response delays are more than 200ms in Fig. 7(a). This is 813 because a large of requests queue in controllers that have lower 814 processing capabilities. Meanwhile, the load-aware scheme 815 also failed to respond to requests quickly. The processing 816 delays of flow requests in different controllers are different 817 when controllers are heterogeneous. As a consequence, select-818 ing a controller by the number of requests is not efficient. 819 Fig. 7(a) reveals that our delay-aware scheme achieved the 820 lowest response duration due to not only estimating response 821 delay but also cutting tail latencies. Fig. 7(a) also reveals that 822 our delay-aware scheme can be adapted to the system where 823 controllers have heterogeneous capabilities. 824

Fig. 7(b) shows that our delay-aware scheme can respond<br/>to more requests than the load-aware scheme and the quantity-<br/>based scheme can in the same time period. Meanwhile,<br/>the throughput difference among schemes grows as the system<br/>runs. In summary, with heterogeneous controllers, our delay-<br/>aware scheme can efficiently reduce response tail latencies and<br/>improve the throughput of controllers.825826827

## D. Impact of Queue Length Bound $RD_{max}$

We evaluate the impact of  $RD_{max}$  on the performance of the delay-aware scheme. We set  $RD_{max} = 20$ ms,  $RD_{max} = 834$ 100ms and  $RD_{max} = \infty$  respectively.  $RD_{max} = \infty$  means



(a) The performance of delay-aware(b) The difference of responded requests scheme under different  $RD_{max}$  settings.under different  $RD_{max}$  settings.

Fig. 8. The impact of  $RD_{max}$  on the performance and throughput of the delay-aware scheme.

that there is no limit on the queue length of controllers. Other parameters are the same with Section V-C.

Fig. 8(a) reveals that our delay-aware scheme has bet-838 ter performance when  $RD_{max}$  has a smaller value. Under 839  $RD_{max} = 20$ ms, all requests can be finished in 20ms. How-840 ever, The maximal response delay is 40ms under  $RD_{max}$  = 841 100ms. Therefore, the performance of delay-aware scheme 842 under  $RD_{max} = 20$ ms is better than that of  $RD_{max} = 100$ ms. 843 Comparing with  $RD_{max} = \infty$ , delay-aware scheme can sig-844 nificantly reduce response tail latencies when  $RD_{max} = 20$ ms. 845 Furthermore, we compare the throughput of controllers under 846 different  $RD_{max}$  settings. Fig. 8(b) shows that the system can 847 respond to 300 more requests under  $RD_{max} = 20$  ms than 848  $RD_{max} = 100$ ms. It was obvious that controllers have higher 849 throughput when  $RD_{max} = 20$ ms than when  $RD_{max} = 100$ ms 850 and  $RD_{max} = \infty$ . The difference of responded requests 851 between  $RD_{max} = 20$ ms and  $RD_{max} = 100$ ms remains 852 stable. However, the difference of responded requests between 853  $RD_{max} = 20$ ms and  $RD_{max} = \infty$  grows as the system runs. 854 Fig. 8 reveals that the setting of  $RD_{max}$  can make flow 855 requests prefer the controllers that have faster processing 856 capabilities. Additionally, it is noteworthy that  $RD_{max}$  can not 857

860 E. Performance Under Heavy Request-Skews

while conducting the controller selection.

858

859

In this section, we study the effect of heavy demand skews 861 on the observed latencies where controllers are homogeneous 862 with average server rate  $\mu = 0.5$ . We set request-skew= 20% 863 and request-skew= 50%. That is, 20% and 50% of switches 864 generated 80% of the total requests towards the controllers. 865 Most of parameters are inherited from Section V-A. To enable 866 20% of switches to generate 80% of the total requests, 867 we randomly select 60 switches and set the arrival rate of flow 868 requests  $\lambda = 2$ . Other switches set  $\lambda = 0.125$ . Under request-869 skew= 50%, half of the switches set  $\lambda = 0.8$ , and the other 870 half of the switches set  $\lambda = 0.2$ . We set  $RD_{max} = 150$ ms 871 because there are too many requests in a short time and these 872 requests have to queue in controllers. 873

be set too small, otherwise, there are no available controllers

Fig. 9(a) shows that over 5% requests have response delays of more than 200ms for the quantity-based scheme. The quantity-based scheme suffers decreased performance due to the request-skews. However, the load-aware and delay-aware schemes can significantly reduce response tail latencies. Based



Fig. 9. Request response time with different schemes under the heavy request-skews.

![](_page_9_Figure_11.jpeg)

(a) Request response time with different(b) The difference of responded requests schemes. between different schemes.

Fig. 10. The performances of different schemes under the time-varying service rates.

on the load-aware and delay-aware schemes, all requests can 879 be finished in 10ms in Fig. 9(a). Fig. 9 reveals that the load-880 aware and delay-aware schemes achieve very similar perfor-881 mances since controllers are homogeneous in this section. 882 Meanwhile, the quantity-based scheme suffers decreased per-883 formance due to the request-skews. Under request-skews, 884 a part of switches generate a large number of the requests, 885 which incur long queues in some controllers for the quantity-886 based scheme. 887

Comparing Fig. 9(a) and Fig. 9(b), we can find that the quantity-based scheme under request-skew= 50% achieves a lower response latency than under request-skew= 20%. It is because the load balance among controllers where requestskew= 50% is better than that of request-skew= 20%. Moreover, Fig. 9 also reveals that our delay-aware scheme can accommodate the heavy request-skews.

# F. Impact of Time-Varying Service Rates

In this section, we study the effect of the service rate fluctuation on the tail latency of response. We change the average service rates of controllers in the system every 50ms, and all controllers randomly set  $\mu = 0.5$  or  $\mu = 0.1$ . Other parameters are inherited from Section V-C.

895

Fig. 10(a) reveals that our delay-aware scheme can respond to all requests in 60ms, the load-aware scheme can respond to all requests during 80ms, and the response tail latency of the quantity-based scheme is more than 200ms. Therefore, our delay-aware scheme can efficiently reduce response tail latencies. For the quantity-based scheme, it could not exploit the feedbacks of controllers to select the controller and further

![](_page_10_Figure_1.jpeg)

(a) Request response time with different (b) Throughput under different schemes.

Fig. 11. The performances of different schemes under the batch selection.

![](_page_10_Figure_4.jpeg)

Fig. 10(b) shows that our delay-aware scheme can respond 912 to more requests than the quantity-based scheme and the load-913 aware scheme. With the increase of system running time, 914 the advantage of the delay-aware scheme is more obvious 915 than the quantity-based scheme in Fig. 10(b). Meanwhile, 916 the difference of responded requests between the delay-917 aware scheme and the load-aware scheme cyclically fluctu-918 ates because the service rates of controllers are periodically 919 changed. In summary, our delay-aware scheme can be better 920 adapted to the time-varying service rate and can efficiently 921 reduce tail latencies of responses. 922

## 923 G. Evaluation for Batch Selection

In this section, we evaluate the performance of different schemes for batch arrival requests. We set the arrival rate of flow requests  $\lambda = 5$ . There are about 300,000 requests during 200ms. To deal with the burst requests, we employ 150 controllers, and each controller can process 40 requests in parallel. We set the batch size  $\delta = 10$  and  $RD_{max} = 15$ ms, and other parameters are set as Section V-C.

Fig. 11(a) reveals that our delay-aware scheme can still 931 reduce tail latencies of responses under the batch selection. 932 Our delay-aware scheme can respond to all requests in 40ms, 933 and the load-aware scheme can respond to all requests during 934 60ms. However, the response tail latency of the quantity-935 based scheme is more than 130ms. The performance of the 936 quantity-based scheme is mainly affected by controllers that 937 have lower processing capabilities. The tail latency of response 938 939 under the load-aware scheme is generated because it fails to accommodate the heterogeneous controllers. Meanwhile, 940 Fig. 11(b) shows that our delay-aware scheme can respond to 941 more flow requests than either of the other two schemes. This 942 means that our delay-aware scheme can not only reduce the 943 tail latency, but can also improve the throughput of controllers 944 under the batch selection. 945

<sup>946</sup> *Impact of Batch Size*  $\delta$ : Furthermore, we evaluate the impact <sup>947</sup> of the batch size  $\delta$  on the performance of the delay-aware <sup>948</sup> scheme. Fig. 12(a) depicts the performances of the delay-<sup>949</sup> aware scheme under different batch sizes. Fig. 12(a) shows

![](_page_10_Figure_10.jpeg)

(a) The performance of delay-aware(b) The number of CS operations. CS: the scheme.

Fig. 12. The impact of the batch size  $\delta$ .

that the performance of the delay-aware scheme has a modest 950 decrease with the increase of the batch size  $\delta$ . When the 951 batch size  $\delta$  increases, it means that more requests will be 952 sent to the same controller, even though the controller has a 953 low processing capacity. As a consequence, the delay-aware 954 scheme suffers a little of decreased performance. we can 955 find that experiments show a similar performance under the 956 different settings of  $\delta$  in Fig. 12(a). Although more requests 957 would be directed to the controller with a low processing 958 capacity at some time when the value of  $\delta$  is larger, after that, 959 the following requests would have less chance to select the 960 controller. Therefore, the decrease of performance is modest. 961 Furthermore, Fig. 12(b) reveals that the number of controller 962 selection operations dramatically decreases when the batch 963 size  $\delta$  goes up. The fewer controller selection operations in 964 switch end will incur less bandwidth consumption in the secure 965 links between switches and controllers and less computing 966 load in controllers. In conclusion, active selection of con-967 trollers for a batch of flows can efficiently reduce the resourse 968 consumption of communication and computing with a little bit 969 of performance reduction. 970

#### VI. RELATED WORK

### A. Network Softwarization

Network softwarization is a transformation trend for design-973 ing, implementing, and managing the 5G and next generation 974 networks. It exploits the benefits of software to enable the 975 redesign of network and service architectures, optimize the 976 expenditure and operational costs, and bring new values in 977 the infrastructures. The key enablers consist of the network 978 function virtualization (NFV), software-defined networking 979 (SDN) and cloud computing, etc. Meanwhile, 5G systems will 980 also rely on these technologies to attain the systems flexibility 981 and elasticity [2]. 982

Along with recent and ongoing advances in cloud com-983 puting, it has become promising to design flexible, scalable, 984 and elastic 5G systems benefiting from advanced virtualization 985 techniques of cloud computing [1]. Taleb and Ksentini [22] 986 introduces the Follow-Me Cloud concept and proposes its 987 framework. There has been research on the possibility of 988 extending cloud computing beyond data centers toward the 989 mobile end user, providing end-to-end mobile connectivity 990 as a cloud service [23]. Software defined networking (SDN) 991

971

acts as a promising enabler for network softwarization and 992 plays a crucial role in the design of 5G wireless networks [1]. 993 SDN has been also utilized in the virtualization of mobile 994 network functions [24]. Kempf et al. [25] describe an evolution 995 of the mobile Evolved Packet Core (EPC) utilizing SDN 996 that allows the EPC control plane to be moved into a data 997 center. Taleb et al. [1] introduce the concept of "Anything 998 as a Service" (ANYaaS), which relies on the reference ETSI 999 NFV architecture to orchestrate and manage important ser-1000 vices. NFV aims at offering network services using network 1001 functions implemented in software and deployed in an on-1002 demand and elastic manner on the cloud [26]. Medhat *et al.* [5] 1003 introduce a service function chaining taxonomy as the basis 1004 for the subsequent state-of-the-art analysis. 1005

### 1006 B. SDN Scalability

To tackle the problem of scaling SDN, Xie et al. [27] 1007 resort to design the distributed controllers, such as Onix [28] 1008 and ONOS [14], which try to distribute the control plane 1009 while maintaining a logically centralized management. These 1010 approaches balance the load of controllers based on the 1011 number of switches, which can not efficiently reduce the tail 1012 latencies of responses. Aissioui *et al.* [29] propose a two-level 1013 hierarchical controller platform to address these scalability and 1014 performance issues in the context of 5G mobile networks. 1015

One key limitation of the distributed controllers is that 1016 the mapping between a switch and a controller is statically 1017 configured. The static configuration results in the uneven 1018 load distribution among the controllers. Bari et al. [30] 1019 propose a management framework, which periodically eval-1020 uates the current controller-to-switch assignment. After that, 1021 it needs to decide whether to perform a reassignment. If a 1022 reassignment is performed, the management framework also 1023 changes the controller-to-switch assignment in the network. 1024 Dixit et al. [31] propose ElastiCon, an elastic distributed 1025 controller architecture in which the controller pool is dynam-1026 ically grown or shrunk according to traffic conditions. In this 1027 case, the load is dynamically shifted across controllers, which 1028 similarly relieves the static mapping between a switch and a 1029 controller. Zhou et al. [32] propose a dynamic and adaptive 1030 1031 algorithm (DALB), which is running as a module of SDN controller. The controllers in distributed environment can 1032 cooperate with each other to keep load balancing. Overloaded 1033 controller can be detected, and high-load switches mapped to 1034 this controller can be smoothly migrated to the under-load 1035 controllers. However, these dynamic frameworks require that 1036 the control plane monitors the state of the whole network and 1037 conducts the reassignments, which aggravate the computing 1038 load of the control plane. 1039

In contrast to these works, our approach relies on simple and 1040 inexpensive feedback of controllers and efficiently relieve the 1041 load of the control plane. Mao and Shen [33] use the principles 1042 of SDN to achieve the server load balancing by setting the 1043 SDN flow table, which does not aim to solve the load balance 1044 of the distributed controllers of SDN. Palma et al. [34] develop 1045 the QueuePusher, which is a queue management extension to 1046 OpenFlow controllers supporting the Open vSwitch Database 1047

Management Protocol (OVSDB) standard. QueuePusher can generate the appropriate queue configuration messages for switches. In addition, there have been lots of researches on how to achieve the consistency among distributed controllers [15]. These techniques are complementary to our approach.

#### VII. CONCLUSION

NFV and SDN can dynamically redistribute the flow across 1055 appropriate VNFs or service function chains if the controller 1056 configures a desired routing path for each network flow 1057 resulting from NFV applications. In this paper, we present 1058 the long-tail observations of the routing response latencies 1059 while using the up-to-date controller-to-switch assignment 1060 mechanism. To tackle this essential problem, we first propose 1061 a light-weight and load-aware switch-to-controller selection 1062 scheme to cut the long-tail response latency under the sim-1063 ple scenario of homogeneous controllers, and then design 1064 a general delay-aware switch-to-controller selection scheme 1065 to fundamentally cut the long-tail response latency for the 1066 more complicated heterogeneous controller scenario with per-1067 formance fluctuations. Through comprehensive performance 1068 evaluation, we demonstrate that our adaptive controller selec-1069 tion schemes can efficiently reduce response long-tail latencies 1070 and accommodate various system environments including the 1071 request-skews, the fluctuation of service rates and so on. 1072

#### REFERENCES

- [1] T. Taleb, A. Ksentini, and R. Jantti, "Anything as a service' for 5G mobile systems," *IEEE Netw.*, vol. 30, no. 6, pp. 84–91, Nov. 2016. 1074
- T. Taleb, S. Dutta, A. Ksentini, M. Iqbal, and H. Flinck, "Mobile edge computing potential in making cities smarter," *IEEE Commun. Mag.*, vol. 55, no. 3, pp. 38–43, Mar. 2017.
- Y. Zhou, D. Zhang, and N. Xiong, "Post-cloud computing paradigms: 1079 A survey and comparison," *Tsinghua Sci. Technol.*, vol. 22, no. 6, 1080 pp. 714–732, 2017.
- M. Kablan, A. Alsudais, E. Keller, and F. Le, "Stateless network functions: Breaking the tight coupling of state and processing," in *Proc.* 1083 14th USENIX NSDI, Boston, MA, USA, Mar. 2017, pp. 97–112.
- [5] A. M. Medhat, T. Taleb, A. Elmangoush, G. A. Carella, S. Covaci, and T. Magedanz, "Service function chaining in next generation networks: 1086 State of the art and research challenges," *IEEE Commun. Mag.*, vol. 55, 1087 no. 2, pp. 216–223, Feb. 2017.
- [6] P. Zave, R. A. Ferreira, X. K. Zou, M. Morimoto, and J. Rexford, 1089
   "Dynamic service chaining with dysco," in *Proc. ACM SIGCOMM*, 1090
   Los Angeles, CA, USA, 2017, pp. 57–70.
- [7] D. L. C. Dutra, M. Bagaa, T. Taleb, and K. Samdanis, "Ensuring end-toend QoS based on multi-paths routing using SDN technology," in *Proc.* 1093 *IEEE Global Commun. Conf.*, Dec. 2017, pp. 1–6.
- [8] N. Gude et al., "NOX: Towards an operating system for networks," ACM 1095 SIGCOMM Comput. Commun. Rev., vol. 38, no. 3, pp. 105–110, 2008. 1096
- [9] A. Ksentini, M. Bagaa, T. Taleb, and I. Balasingham, "On using bargaining game for optimal placement of SDN controllers," in *Proc. IEEE Int. Conf. Commun. (ICC)*, May 2016, pp. 1–6.
- [10] J. Xie, D. Guo, Z. Hu, T. Qu, and P. Lv, "Control plane of software defined networks: A survey," *Comput. Commun.*, vol. 67, pp. 1–10, 1101 Aug. 2015.
- D. Kreutz, F. Ramos, P. E. Veríssimo, C. E. Rothenberg, S. Azodolmolky, and S. Uhlig, "Software-defined networking: A comprehensive survey," *Proc. IEEE*, vol. 103, no. 1, pp. 14–76, Jan. 2015.
- [12] A. Ksentini, M. Bagaa, and T. Taleb, "On using SDN in 5G: 1106 The controller placement problem," in *Proc. IEEE Global Commun.* 1107 *Conf. (GLOBECOM)*, Dec. 2016, pp. 1–6. 1108
- M. Al-Fares, A. Loukissas, and A. Vahdat, "A scalable, commodity data center network architecture," ACM SIGCOMM Comput. Commun. Rev., vol. 38, no. 4, pp. 63–74, 2008.

1073

- [14] P. Berde et al., "ONOS: Towards an open, distributed SDN OS," in Proc. 1112 ACM HotSDN, 2014, pp. 1-6. 1113
- [15] A. Panda, W. Zheng, X. Hu, A. Krishnamurthy, and S. Shenker, "SCL: 1114 Simplifying distributed SDN control planes," in Proc. 14th USENIX 1115 1116 NSDI, Mar. 2017, pp. 329-345.
- [16] D. Hugheshallett, A. M. Gleason, and W. G. Mccallum, Calculus: Single 1117 1118 and Multivariable, Student Solutions Manual, 6th ed. Hoboken, NJ, USA: Wiley, 2013. 1119
- A. Bremler-Barr, Y. Harchol, and D. Hay, "OpenBox: A software-defined 1120 [17] 1121 framework for developing, deploying, and managing network functions," in Proc. ACM SIGCOMM, Salvador, Brazil, Aug. 2016, pp. 511-524. 1122
- 1123 Y. Azar, A. Z. Broder, A. R. Karlin, and E. Upfal, "Balanced alloca-[18] tions," SIAM J. Comput., vol. 29, no. 1, pp. 180-200, 1999. 1124
- M. Mitzenmacher, A. W. Richa, and R. Sitaraman, "The power of two [19] 1125 1126 random choices: A survey of techniques and results," in Handbook of Randomized Computing, vol. 11. Norwell, MA, USA: Kluwer, 2000, 1127 pp. 255-312. 1128
- [20] R. Nishtala et al., "Scaling memcache at Facebook," in Proc. USENIX 1129 NSDI, 2013, pp. 385-398. 1130
- 1131 [21] A. Vulimiri, P. B. Godfrey, R. Mittal, J. Sherry, S. Ratnasamy, and S. Shenker, "Low latency via redundancy," in Proc. ACM Conf. Emerg. 1132 Netw. Experim. Technol., 2013, pp. 283-294. 1133
- 1134 [22] T. Taleb and A. Ksentini, "Follow me cloud: Interworking federated clouds and distributed mobile networks," IEEE Netw., vol. 27, no. 5, 1135 pp. 12-19, Sep./Oct. 2013. 1136
- T. Taleb, "Toward carrier cloud: Potential, challenges, and solutions," [23] 1137 IEEE Wireless Commun., vol. 21, no. 3, pp. 80-91, Jun. 2014. 1138
- 1139 [24] T. Taleb, A. Ksentini, and A. Kobbane, "Lightweight mobile core networks for machine type communications," IEEE Access, vol. 2, 1140 pp. 1128-1137, 2014. 1141
- [25] J. Kempf, B. Johansson, S. Pettersson, H. Lüning, and T. Nilsson, 1142 "Moving the mobile evolved packet core to the cloud," in Proc. IEEE 1143 8th Int. Conf. Wireless Mobile Comput., Netw. Commun. (WiMob), 1144 Oct. 2012, pp. 784-791. 1145
- T. Taleb et al., "EASE: EPC as a service to ease mobile core net-[26] 1146 1147 work deployment over cloud," IEEE Netw., vol. 29, no. 2, pp. 78-88, 1148 Mar./Apr. 2015.
- J. Xie, D. Guo, X. Zhu, B. Ren, and H. Chen, "Minimal fault-tolerant 1149 [27] coverage of controllers in IaaS datacenters," IEEE Trans. Services 1150 Comput., to be published, doi: 10.1109/TSC.2017.2753260. 1151
- [28] T. Koponen et al., "Onix: A distributed control platform for large-scale 1152 production networks," in Proc. USENIX OSDI, 2010, pp. 351-364. 1153
- [29] A. Aissioui, A. Ksentini, A. M. Gueroui, and T. Taleb, "Toward elastic 1154 distributed SDN/NFV controller for 5G mobile cloud management 1155 systems," IEEE Access, vol. 3, pp. 2055-2064, 2015. 1156
- [30] M. F. Bari et al., "Dynamic controller provisioning in software defined 1157 networks," in Proc. Int. Conf. Netw. Service Manage., Oct. 2013, 1158 pp. 18-25. 1159
- A. Dixit, F. Hao, S. Mukherjee, T. Lakshman, and R. Kompella, [31] 1160 1161 "Towards an elastic distributed SDN controller," in Proc. ACM HotSDN, Hong Kong, Aug. 2013, pp. 7-12. 1162
- Y. Zhou et al., "A load balancing strategy of sdn controller based on 1163 [32] distributed decision," in Proc. IEEE 13th Int. Conf. Trust, Secur. Privacy 1164 Comput. Commun., Sep. 2014, pp. 851-856. 1165
- Q. Mao and W. Shen, "A load balancing method based on SDN," in 1166 [33] Proc. 7th Int. Conf. Meas. Technol. Mechatronics Autom., Jun. 2015, 1167 pp. 18–21. 1168
- D. Palma et al., "The queuepusher: Enabling queue management in 1169 [34] openflow," in Proc. 3rd Eur. Workshop Softw. Defined Netw., Sep. 2014, 1170 pp. 125-126. 1171

![](_page_12_Picture_22.jpeg)

Deke Guo received the B.S. degree in industry engi-1184 neering from the Beijing University of Aeronautics 1185 and Astronautics, Beijing, China, in 2001, and the 1186 Ph.D. degree in management science and engineer-1187 ing from the National University of Defense Tech-1188 nology, Changsha, China, in 2008. He is currently a 1189 Professor with the College of Information System 1190 and Management, National University of Defense 1191 Technology, and a Professor with the School of 1192 Computer Science and Technology, Tianjin Uni-1193 versity. His research interests include distributed 1194

systems, software-defined networking, data center networking, wireless and 1195 mobile systems, and interconnection networks. He is a member of the ACM. 1196

![](_page_12_Picture_25.jpeg)

Xiaozhou Li received the Ph.D. degree in com-1197 puter science from Princeton University in 2016. 1198 He is currently a Software Engineer with Barefoot 1199 Networks, where he is building new networking 1200 systems with fast programmable network chips. His 1201 research improves the performance, scalability, and 1202 efficiency of datacenter services, with a particular 1203 focus on combining new hardware and infrastructure 1204 capabilities with careful architectural design and 1205 algorithm engineering. 1206

![](_page_12_Picture_27.jpeg)

Yulong Shen received the B.S. and M.S. degrees in 1207 computer science and the Ph.D. degree in cryptog-1208 raphy from Xidian University, Xian, China, in 2002, 1209 2005, and 2008, respectively. He is currently a 1210 Professor with the School of Computer Science and 1211 Technology, Xidian University, and also an Asso-1212 ciate Director of the Shaanxi Key Laboratory of 1213 Network and System Security. He has also served on 1214 the technical program committees of several interna-1215 tional conferences, including the NANA, the ICEBE, 1216 the INCoS, the CIS, and the SOWN. His research 1217 1218

interests include wireless network security and cloud computing security.

![](_page_12_Picture_30.jpeg)

science and technology from the Beijing Institute of Technology, Beijing, China, in 2013, and the M.S. degree in management science and engineering from the National University of Defense Technology, Changsha, China, in 2015, where he is currently pursuing the Ph.D. degree. He is also a joint Ph.D. student with the University of California, Santa Cruz (UCSC), CA, USA, from 2017. His study in the UCSC is supported by the China Scholarship Council. His research interests include distributed

systems, data-center networks, and software-defined networks. 1183

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

![](_page_12_Picture_33.jpeg)

Xiaohong Jiang received the B.S., M.S., and Ph.D. 1219 degrees from Xidian University, China, in 1989, 1220 1992, and 1999 respectively. From 2005 to 2010, he 1221 was an Associate professor, Tohoku University. He 1222 is currently a Full Professor with Future University 1223 Hakodate, Japan. His research interests include com-1224 puter communications networks, mainly wireless 1225 networks and optical networks, network security, 1226 and routers/switches design. He has authored or co-1227 authored over 300 technical papers at premium inter-1228 national journals and conferences, which include 1229

over 70 papers published in top IEEE journals and top IEEE conferences. He 1230 was a recipient of the Best Paper Award of IEEE HPCC 2014, IEEE WCNC 1231 2012, IEEE WCNC 2008, IEEE ICC 2005-Optical Networking Symposium, 1232 and IEEE/IEICE HPSR 2002. He is a Member of IEICE. 1233

# Cutting Long-Tail Latency of Routing Response in Software Defined Networks

Junjie Xie, Deke Guo, Xiaozhou Li, Yulong Shen, and Xiaohong Jiang

Abstract— To enable the network softwarization, network func-1 2 tion virtualization (NFV) and software defined networking (SDN) are integrated to jointly manage and utilize the network resource 3 and virtualized network functions (VNFs). For a network flow 4 resulting from any NFV application, an associated switch would 5 send a routing request to the controller in SDN. The controller then generates and configures a routing path to dynamically steer the flow across appropriate VNFs or service function chains. 8 This process, however, exhibits a skew distribution of response 9 latency with a long tail. Cutting the long-tail latency of response 10 is critical to enable the network softwarization, yet difficult to 11 achieve due to many factors, such as the limited capacities and 12 the load imbalance among controllers. In this paper, we reveal 13 that such flow requests still experience the long-tail response 14 latency, even using the up-to-date controller-to-switch assignment 15 mechanism. To tackle this essential problem, we first propose a 16 light-weight and load-aware switch-to-controller selection scheme 17 to cut the long-tail response latency under the simple scenario 18 of homogeneous controllers, and then design a general delay-19 aware switch-to-controller selection scheme to fundamentally 20 cut the long-tail response latency for the more complicated 21 heterogeneous controller scenario with performance fluctuations. 22 The comprehensive evaluations indicate that our two new switch-23 to-controller selection schemes can significantly reduce the long-24 tail latency and provide higher system throughput. 25

Index Terms—Network softwarization, software defined
 networks, controller selection, long-tail latency.

Manuscript received September 30, 2017; revised February 5, 2018; accepted February 27, 2018. This work was supported in part by the National Natural Science Foundation for Outstanding Excellent Young Scholars of China under Grant 61422214, in part by the National Natural Science Foundation of China under Grant 61772544 and Grant U1536202, in part by the National Basic Research Program (973 program) under Grant 2014CB347800, in part by the Hunan Provincial Natural Science Fund for Distinguished Young Scholars under Grant 2016J11002, and in part by the Guangxi Cooperative Innovation Center of cloud computing and Big Data under Grant YD16507 and Grant YD17X11. (*Corresponding authors: Deke Gua; Yulong Shen.*)

J. Xie is with the Science and Technology on Information Systems Engineering Laboratory, National University of Defense Technology, Changsha 410073, China (e-mail: xiejunjie06@gmail.com).

D. Guo is with the College of System Engineering, National University of Defense Technology, Changsha 410073, China. He is also with the School of Computer Science and Technology, Tianjin University, Tianjin 300072, China (e-mail: guodeke@gmail.com).

X. Li is with the Department of Computer Science, Princeton University, Princeton, NJ 08544 USA (e-mail: xl@cs.princeton.edu).

Y. Shen is with the School of Computer Science and Technology, Xidian University, Xi'an 710071, China (e-mail: ylshen@mail.xidian.edu.cn).

X. Jiang is with the School of Systems Information Science, Future University Hakodate, Hokkaido 041-8655, Japan (e-mail: jiang@fun.ac.jp).

Color versions of one or more of the figures in this paper are available online at http://ieeexplore.ieee.org.

Digital Object Identifier 10.1109/JSAC.2018.2815358

#### I. INTRODUCTION

TETWORK softwarization is a transformation trend for designing, implementing, and managing the next generation networks. It exploits the benefits of software to enable the redesign of network and service architectures, optimize the expenditure and operational costs, and bring added values. The key enablers consist of the network function virtualization (NFV), software-defined networking (SDN) and cloud computing, etc [1]. Moreover, 5G systems will also rely on these technologies to attain system's flexibility and true elasticity [2], [3]. Network functions (NFs) are crucial for improving network security by examining and modifying network flows using special-purpose hardware. Recently, NFV has been proposed to execute virtual network functions (VNFs) on generic compute resources [4], such as commodity servers and VMs. Normally, a flow goes through specific VNFs in a particular order to meet its required processing, following the service function chain (SFC) [5], [6] along a routing path.

Additionally, SDN offers the freedom to refactor the control plane and flexibly enables the network softwarization [7]. More precisely, NFV and SDN can jointly manage the network resource and VNFs, and dynamically steer network flows across appropriate VNFs or SFCs. SDN centralizes the network control plane to a programmable software component, i.e., a controller running on a generic server, such as NOX [8]. The controller maintains a global network view and optimizes the forwarding decisions of network flows. For a flow from any NFV application, an associated switch would send a routing request to the controller. It is the controller that generates and configures a routing path to a specific VNF instance or to traverse a SFC on demand. The above process between a pair of switch and controller brings the response latency. Many factors would skew the tail of the latency distribution. For example, a single controller that lacks sufficient capacity to tackle received routing requests quickly and inevitably becomes a performance bottleneck [9]. Thus, such routing requests experience long-tail latency of response, as evaluated in Section II. Cutting the long-tail latency of routing response is critical to enable the network softwarization, yet difficult to achieve due to many factors.

To improve the scalability of SDN, the distributed control plane consisting of multiple controllers has been proposed recently [10], such as ONOS and OpenDaylight. To reduce the long-tail latency of response, they resort to the controllerto-switch assignment mechanism. That is, the control plane

0733-8716 © 2018 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.

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

proactively assigns a controller to each switch such that each 74 controller manages the same amount of switches. In reality, 75 the quantities of routing requests coming from switches per 76 unit time are different and dynamic. Consequently, controllers 77 still differ in the amount of received routing requests per time 78 unit. This load imbalance among controllers leads to the long-79 tail latency of response. Additionally, the controller-to-switch 80 assignment requires coordination among controllers, which 81 further aggravates the loads of controllers. 82

In this paper, to cut the long-tail latency of response and 83 lighten the load of controllers, we propose conducting the 84 selection of controllers at the side of each switch instead of the 85 controller-to-switch assignment. This means that each switch 86 actively chooses one controller from multiple available ones, 87 decoupling the static binding between switches and controllers. 88 More precisely, each switch prefers to adaptively select the 89 controller with low response latency for routing requests. This 90 would move the partial intelligences of the network to switches 91 and efficiently reduces the loads of controllers. 92

Despite those potential benefits, the selection of controllers 93 still faces many challenges. First, the switches need to probe 94 the state of controllers via the secure channel between them. 95 The secure channel is one kind of rare resource and affects the performance of the whole network. To save the bandwidth of 97 the secure channel, the selection process of controllers should be light-weight and use a few of the feedbacks from the 99 controllers. Second, the selection scheme needs to be scalable, 100 irrespective of the network size and the number of controllers. 101 Third, the selection scheme needs to accommodate the bursty 102 and skew routing requests from switches. Last, the selection 103 scheme should adapt to the heterogeneous controllers and the 104 performance fluctuation across controllers. 105

To tackle such challenges, we design a *load-aware* selec-106 tion scheme of controllers, which is simple but effective to 107 achieve the load balance among controllers. The load of a 108 controller refers to the number of routing requests waited 109 to be processed. The basic idea of our scheme is that each 110 switch sends routing requests to the controller with the lowest 111 load. In this way, all controllers process the similar number 112 of routing requests per time slot. This is very helpful to cut 113 the long-tail latency of routing response, when all controllers 114 have the same processing capabilities. This method alone, 115 however, is insufficient to deal with more general settings 116 of heterogenous controllers and the performance fluctuation. 117 Those controllers with lower processing capabilities still incur 118 the long-tail latency of response for routing requests, when all 119 controllers achieve the load balance. For this reason, we further 120 present a general delay-aware selection scheme of controllers 121 to fundamentally cut the long-tail latency of routing response. 122

Our delay-aware selection scheme includes two key com-123 ponents. The first one is the controller selection model of each 124 switch, which uses simple and inexpensive probing feedbacks 125 from a few controllers. It is still effective if each switch just 126 randomly probes two controllers and sends upcoming routing 127 requests to the controller with the shorter response delay. This 128 model is scalable and light-weight since it is not affected by 129 the network scale and the number of controllers. The second 130 component is the queue management mechanism of each 131

controller. It could estimate the response delay of a routing request and hence improve the performance predictability of controllers. The evaluation results reveal that our *delay-aware* selection scheme can efficiently reduce the long-tail latency of routing responses and improve the system throughput.

In summary, the major contributions of this paper are as follows.

- We reveal that routing requests experience the long-tail response latency, even using the up-to-date controllerto-switch assignment mechanism in SDN. Therefore, we propose an adaptive selection mechanism of controllers for switches to cut the long-tail latency.
- 2) We first design an efficient *load-aware* selection method 144 of homogeneous controllers for each switch. For more 145 general scenarios, we further propose a general delay-146 aware selection method, which is adaptive to the bursty 147 routing requests, the heterogenous controllers and the 148 performance fluctuations. Our two methods are light-149 weight as they limit the additional overhead caused by 150 probing two controllers. 151
- We further develop a queue management mechanism for each controller, which can efficiently manage the queue length and estimate the response delay of routing requests. The evaluation results reveal that our controller selection methods can accommodate system environment variations and efficiently reduce long-tail latency of routing response.

The paper is organized as follows. In Section II, we present 159 the observation of long-tail latency of routing response. 160 Section III depicts the framework of our controller selection 161 mechanism and the load-aware selection scheme of controllers. 162 We present the delay-aware selection scheme of controllers in 163 Section IV. We conduct massive experiments to evaluate the 164 performance of our controller selection schemes under various 165 system environment in Section V. Section VI introduces the 166 related work. In Section VII, we conclude this paper. 167

# II. LONG TAIL OF RESPONSE LATENCIES

In a SDN, when a switch receives a new flow, the switch 169 sends a routing request to its controller. The controller then 170 computes a route for the flow and inserts flow rules to related 171 switches in the route. Thus, the new flow would be forwarded 172 according to the flow rules in switches [11], [12]. Such an 173 interaction between the switch and the controller causes the 174 response latency. For a routing request, the latency of routing 175 response denotes the time interval from sending the routing 176 request to receiving the flow rules generated by the controller. 177

#### A. Long-Tail Observations of Response Latencies

Fig. 1(a) plots an observation about the long-tail distribu-179 tion under a single instance of ONOS controller. We build 180 a SDN testbed with one controller, running in a virtual 181 machine with 2 CPU cores and 2G RAM. Note that the 182 testbed forms a typical Fat-tree datacenter topology [13]. 183 We record the response latencies of 12,000 routing requests. 184 As shown in Fig. 1(a), the response latencies of 50% of routing 185 requests are lower than 5ms, and 90% of routing requests are 186

168

![](_page_15_Figure_1.jpeg)

(a) A single controller with 12,000 routing(b) Multiple controllers with 150,000 routrequests. ing requests.

Fig. 1. Long-tail distributions of routing responses under a single controller as well as multiple controllers.

served within 30ms. However, there still exist some routing
 requests whose response latencies are more than 50ms. That
 is, the response latencies exhibit a long-tail distribution.

Furthermore, we observe the response latencies of rout-190 ing requests under multiple controllers [14]. We employ 191 300 switches and 40 controllers where each controller 192 manages 7 or 8 switches. Each switch generates routing 193 requests according to a Poisson arrival process with  $\lambda = 0.5$ 194 during 1ms. The processing time of each request in each 195 controller is drawn from an exponential distribution where 196  $\mu^{-1}$  = 2ms. Each controller can process 10 requests in 197 parallel. We run the system 1000ms and record the response 198 latencies of routing requests. The number of arrival routing 199 requests is about 150,000. In Fig. 1(b), 89% of routing 200 requests can be served in 5ms, and the response latencies 201 of 96% of routing requests are lower than 10ms. However, 202 some response latencies are still larger than 20ms. That is, 203 Fig. 1(b) shows that there still exists a long-tail distribution 204 under multiple ONOS controllers. 205

# 206 B. Analysis About the Long-Tail Latencies

Fig. 1(a) results from that the network only employs one 207 controller. Due to the limited capability of the single controller, 208 a large amount of requests have to queue in the controller. 209 Therefore, there is a long-tail latency of routig response caused 210 by the long queueing delays. In Fig. 1(b), the number of 211 212 routing requests that each switch generates is different, even though they obey the same Poisson distribution. The skew-213 flow requests will make that some controllers are overload, but 214 other controllers would be underutilized. As a consequence, 215 there will be a long-tail latency. We illustrate the problem 216 in Fig. 2. Two controllers are assigned to four switches and 217 have the same processing time of 4ms. Assume  $switch_1$  and 218 switch<sub>2</sub> receive 4 requests each and that switch<sub>3</sub> and switch<sub>4</sub> 219 receive 2 requests each. The requests received by  $switch_1$ 220 and  $switch_2$  can only be processed by  $controller_1$ , which is 221 assigned to manage them. This leads to a maximum latency of 222 32ms, but a load-aware selection obtains a maximum latency 223 of 24ms. Fig. 2 shows that a quantity-based assignment strat-224 egy leads to long-tail latencies because it fails to accommodate 225 the skew-flow requests. 226

227 Quantity-based allocation strategy is commonly employed 228 by many controllers to balance the loads of controllers, such as

![](_page_15_Figure_9.jpeg)

Fig. 2. The performances of quantity-based controller allocation and loadaware controller selection.

ONOS [14]. That is a controller-to-switch assignment mech-229 anism, which is abbreviated as the assignment mechanism 230 of controllers. In this case, controllers coordinate to manage 231 switches, and each controller manages an approximately equal 232 quantity of switches. When deploying multiple instances of 233 ONOS in a SDN, the bursty flows from a switch are sent to the 234 same controller. Consequentially, a large number of requests 235 have to queue in the controller. These queueing requests tend 236 to incur long response latencies. However, those controllers, 237 which do not receive bursty routing requests, may even be 238 underloaded. In conclusion, the assignment mechanism fails 239 to efficiently reduce the tail latencies of responses. 240

# III. FRAMEWORK OF CONTROLLER SELECTION MECHANISM

241 242

250

To overcome the drawback of assignment mechanism, we design a switch-to-controller selection scheme, which is abbreviated as the selection scheme of controllers. The selection scheme moves partial intelligences of the network to switches and relieves the loads of controllers. Meanwhile, the selection scheme can efficiently reduce tail latencies of responses. 249

#### A. Overview of Controller Selection Mechanism

For existing designs of control plane, the assignment mech-251 anism of controllers is a static binding between switches and 252 controllers, which fails to deal with the bursty and skew 253 routing requests, and further incurs long response delays. 254 To reduce the tail latencies of responses, routing requests from 255 the same switch need to be processed by appropriate con-256 trollers. This means that the static binding between switches 257 and controllers needs to be decoupled. A better mechanism is 258 to enable switches to select controllers for routing requests. 259

For the assignment mechanism, the load balance among 260 controllers means that each controller manages the same 261 amount of switches. The load is denoted by the number 262 of switches. Moreover, the load balance among controllers 263 requires the coordination of controllers. The coordination 264 will further aggravate the computing and communication 265 overhead. However, in this paper, the selection scheme of 266 controllers achieves the mapping between switches and con-267 trollers through switches conduct simple and actively probing. 268

Therefore, the selection scheme relieves the overhead of 269 coordination and assignment in controllers. 270

In SDN, the controller periodically probes the switch state, 271 and switches would respond to those probes and forward any 272 switch notifications (e.g., link failures or recovery) to all live 273 controllers [15]. When a switch suffers a failure scenario, 274 it fails to respond to the probe. After that, the controller 275 promptly updates its local topology state and replicates that 276 update to all other controllers in the cluster. In addition, there 277 have been lots of researches on how to achieve the consistency 278 among distributed controllers [15]. These techniques are com-279 plementary to our approach. 280

We design the controller selection mechanism keeping in 281 mind these four goals: 282

1) Light-weight: A light-weight probing method is needed 283 to save the bandwidth of the secure channel. The probing 284 for the controller selection uses the secure channel, 285 which is the communication channel between the control 286 plane and the data plane. The bandwidth of the channel 287 can affect the performance of the whole network. 288

Scalable: The selection scheme of controllers should be 2) 289 irrelevant to the increasing number of controllers. The 290 expansion of network size is common. The method of 291 probing controller should accommodate the increase of 292 deployed controllers and avoid to incur the communica-293 tion overhead and the computing overhead. 294

Burst-immunity: The selection scheme should be burst-295 immune. There are bursty and skew-flow requests, which 296 can lead to long response delays. To shorten the tail 297 latency of responses, the selection scheme needs to 298 accommodate the bursty and skew-flow requests. 299

Adaptive: The selection scheme of controllers should be 300 adaptive. The capabilities of controllers may be hetero-301 geneous and time-varying. To deal with the general situ-302 ation, the selection scheme must cope and quickly react 303 to heterogeneous and time-varying processing capabili-304 ties across controllers. 305

#### B. Load-Aware Selection Scheme of Controllers 306

Accommodating skew-flow requests across controllers 307 necessitates a selection strategy of controllers. The strategy 308 can make switches select faster controllers for routing requests. 309 The controller with fewer unfinished requests can respond to 310 routing requests faster when controllers have the same process-311 ing capability. In this paper, we first present a load-aware 312 selection scheme. To realize this framework, the selection 313 strategy needs to take into account the loads across multiple 314 controllers in the network. The load means the number of 315 unfinished requests. Our load-aware selection scheme is to 316 select a controller with the fewest unfinished requests for 317 newly generated requests. 318

Under the load-aware selection scheme of controllers, 319 the switch needs to send a probing request to each controller 320 after a switch receives a new flow. When the controller receives 321 the probing request, it will return the number of unfinished 322 requests to the switch. After the switch receives all probing 323 results, it then sends the new routing request to the controller 324

with the lightest load since that controller can respond the 325 routing request fastest. The load-aware selection scheme aims 326 to reduce tail latencies of responses by selecting the controller 327 with the lightest load. 328

In Fig. 2, the load-aware selection scheme will work as 329 follows. When a switch receives a new flow, it first probes the 330 loads of controllers A and B. Based on the loads of controllers 331 A and B, the switch sends the routing request to the controller 332 with lightest load. This scheme can balance the loads of 333 controllers A and B and can achieve better selection. When the 334 controller finishes processing the routing request, it inserts the 335 flow rules to related switches. Lastly, those switches will deal 336 with the flow according to the actions of matched flow rules. 337 In addition, the new arrival routing requests need to queue in 338 controllers, when controllers are busy. However, the infinite 339 length of queue will incur infinite response delays of routing 340 requests. To cut the tail latency of response, it is necessary to 341 ensure that the length of queue is finite in each controller. 342

# C. The Condition to Finite Queue Length

We give the condition to achieve that the expected number 344 of requests in per controller remains finite for all time. Con-345 sider the following model: requests arrive as a Poisson stream 346 of rate  $\lambda$  at each switch. Requests are processed according 347 to the first-in first-out (FIFO) protocol by controllers. The 348 processing time for a request is exponentially distributed with 349 mean  $\mu$ . When there are *m* switches and *n* controllers in a 350 network, requests arrive as a Poisson stream of rate  $\frac{\lambda m}{n}$  at 351 each controller. We obtain the following theorem. Note that 352  $\frac{\lambda m}{m} < \mu$ , and then the system will be stable, which means that 353 the expected number of requests per controller remains finite 354 in equilibrium. Theorem 1 shows that the system is stable for 355 every  $\frac{\lambda m}{\mu n}$  < 1; that is, the expected number of requests in each 356 controller remains finite for all time. 357

Theorem 1: The system is stable for every  $\frac{\lambda m}{\mu n}$  < 1; that 358 is, the expected number of requests in each controller remains finite for all time.

*Proof:* When we treat all controllers as a whole and 361 all switches as a whole, then the system can be seen as 362 a M/M/1 system with Poisson arrival rate  $\lambda m$  and average 363 service rate  $\mu n$ . Let  $P_k(t)$  denote the probability of that there 364 are k requests in the whole system in time t. Accordingly, 365  $P_{k+1}(t)$  denotes the probability where k + 1 requests exist in 366 the system in time t, and  $P_k(t + \Delta t)$  denotes the probability of 367 that there are k requests in the whole system in time  $t + \Delta t$ . 368 Now we consider the evolution of the system. In the time 369  $[t, t + \Delta t]$ , the process of evolution has some attributes just 370 as follows: 371

- One request comes with the probability  $\lambda m \Delta t$ , and the probability of no request comes is  $1 - \lambda m \Delta t$ .
- One request departures with the probability  $\mu n \Delta t$ , and 374 the probability of no request departure is  $1 - \mu n \Delta t$ . 375
- The situation of more than one request comes and depar-376 tures in  $\triangle t$  is a small probability event, and it can be 377 ignored. 378

There are 4 types of evolution process, and we list them 379 in Table I. Take type B for an example. Since k is the number 380

343

359

360

372

TABLE IFour Types of Evolution Process

Туре	state in time t	come	departure	state in time $(t + \Delta t)$
(A)	k	0	0	k
(B)	k+1	0	1	k
(C)	k-1	1	0	k
(D)	k	1	1	k

of requests in the whole system, so k + 1 means that there 381 are k + 1 requests in the system in time t. After that, when 382 a request departs during  $\Delta t$ , there would be k requests in 383 the system in time  $t + \Delta t$ . Accordingly, we can get the 384 possibility of type A is  $P_k(t)(1-\lambda m \Delta t)(1-\mu n \Delta t)$ , possibility 385 of type B is  $P_{k+1}(t)(1 - \lambda m \triangle t) \mu n \triangle t$ , possibility of type 386 C is  $P_{k-1}(t)\lambda m \Delta t (1 - \mu n \Delta t)$  and possibility of type D is 387  $P_k(t)\lambda m \Delta t \mu n \Delta t$ . Then,  $P_k(t + \Delta t)$  should be the sum of all 388 4 types, shown in Equation (1). 389

$$P_{k}(t + \Delta t) = P_{k}(t)(1 - \lambda m \Delta t - \mu n \Delta t)$$

$$P_{k+1}(t)\mu n \Delta t + P_{k-1}(t)\lambda m \Delta t + o(\Delta t) \quad (1)$$

And let  $\Delta t \rightarrow 0$ , we can get a differential equation, shown in Equation (2).

<sup>394</sup> 
$$\frac{dP_k(t)}{dt} = \lambda m P_{k-1}(t) + \mu n P_{k+1}(t)$$
  
<sup>395</sup>  $-(\lambda m + \mu n) P_k(t) \quad k = 1, 2, \dots$  (2)

Noted that if k = 0, there will exist only type A and type B, shown in Equations (3) and (4).

<sup>398</sup> 
$$P_0(t + \Delta t) = P_0(t)(1 - \lambda m \Delta t) + P_1(t)(1 - \lambda m \Delta t)\mu n \Delta t$$
 (3)

<sup>399</sup> 
$$\frac{dP_0(t)}{dt} = -\lambda m P_0(t) + \mu n P_1(t)$$
 (4)

We just have interest in the equilibrium point, and the derivative is 0 in fixed point. Thus, we get Equation (5).

402

$$\begin{cases} -\lambda m P_0 + \mu n P_1 = 0\\ \lambda m P_{k-1} + \mu n P_{k+1} - (\lambda m + \mu n) P_k = 0 \quad k \ge 1 \end{cases}$$
(5)

Resolve equation (5), we can get  $P_k = (\lambda m/\mu n)^k P_0$ . If  $\frac{\lambda m}{\mu n} < 1$ , then the sequence  $P_k$  will be decrease. And we know probability is non-negative, that means  $P_k \ge 0$ . If a sequence is bounded and monotone, it converges [16]. So there exist K, when k > K,  $P_k = 0$ . Then the expected total number of requests in all controllers remains finite.

Theorem 1 shows that the expected total number of requests 409 in each controller remains finite, when  $\frac{\lambda m}{\mu n}$  < 1. Therefore, 410 to achieve the finite queue length, it is essential to ensure that 411  $\frac{\lambda m}{\mu n} < 1$ . When the network size *m* increases, if  $\frac{\lambda m}{\mu n} \ge 1$ , it is 412 necessary to increase the number of deployed controllers n, 413 otherwise, the queue length of some controllers will be infinite, 414 and that will incur infinite response delays. Another method 415 to limit the queue length of controllers is to drop some routing 416 requests, which can limit the value of  $\lambda$  and achieve  $\frac{\lambda m}{un} < 1$ . 417

## 418 D. Limitations of Load-Aware Selection Scheme

Fig. 3 shows the Cumulative Distribution Function (CDF) of response latencies of routing requests under different schemes.

![](_page_17_Figure_16.jpeg)

Fig. 3. Response latencies of routing requests under different schemes.

![](_page_17_Figure_18.jpeg)

Fig. 4. Distinct selection schemes incur different response latencies.

In Fig. 3, the response latencies of 94% of routing requests 421 are lower than 5ms after adopting the load-aware selection 422 scheme. However, for the quantity-based assignment, only 423 86% of routing requests can be responded in 5ms. In Fig. 3, 424 all routing requests can be responded in 10ms under the load-425 aware scheme. Therefore, Fig. 3 indicates that our load-aware 426 selection scheme can reduce the tail latency of responses than 427 the prior quantity-based allocation method when controllers 428 are homogeneous and exhibit the same processing capabilities. 429 However, we can see that both curves (Quantity-based and 430 Load-aware) are very close to each other, which means that 431 the load-aware selection strategy narrowly reduce the long-432 tail latency. Furthermore, the load-aware selection scheme 433 faces three challenges. First, controllers are heterogeneous. 434 Second, the processing capabilities of controllers are dynam-435 ically changing. Third, the cost of probing is too huge. 436 Therefore, the load-aware selection scheme is still insufficient 437 to completely cut the tail latency. 438

Fig. 4 plots an illustrative example of the limitation. For two 439 controllers, the processing time per request in controller A and 440 controller B are 5ms and 12ms, respectively. Assume all four 441 switches receive a burst of 3 requests each. Each request needs 442 to be forwarded to a single controller. If every switch selects 443 a controller using the load-aware scheme, it will result in each 444 controller receiving an equal share of the requests. This leads 445 to a maximum latency of 72ms, whereas an ideal selection in 446 this case obtains a maximum latency of 45ms. We note that 447 the load-aware scheme will prefer faster controllers over time, 448 but purely relies on the load information. Therefore, when 449 controllers are heterogeneous, the load-aware selection scheme 450 can not efficiently shorten the tail latency. 451

Controllers are commonly heterogeneous for primarily three reasons. First, the hardware is heterogeneous. Controllers run

488

489

in commercial servers. These servers can be heterogeneous 454 due to different hardware configurations, such as CPU and 455 memory. Second, the software is heterogenous. There are 456 multiple different controllers developed by different organiza-457 tions [10], such as NOX, Beacon, Floodlight, Ryu, ONOS and 458 OpenDaylight, etc. Those controllers themselves have different 459 performances. Third, the function is heterogeneous. There 460 are some management applications running in controllers for 461 achieving different functions [12], [17], and these applications 462 will consume some resources of controllers. Consequentially, 463 controllers have different remaining capabilities for processing 464 routing requests, even if the controllers run in servers with 465 the same setting. In this case, queueing routing requests in 466 controllers with low processing capabilities will lead to long 467 response latencies. 468

Additionally, the load-aware scheme probes the loads of all 469 controllers, and then selects the controller that has the lightest 470 load. However, this probing will incur the overhead of com-471 munication and aggravate the loads of controllers when there 472 are a large number of controllers in a large-scale network. For 473 example, there are m switches and n controllers in a network. 474 Suppose that each switch receives  $\lambda$  routing requests in 1ms. 475 There are  $2\lambda \times m \times n$  times communications between switches 476 and controllers during 1ms. Meanwhile, each controller needs 477 to evaluate its own load  $\lambda m$  times in 1ms. The cost of probing 478 is too huge for the load-aware selection scheme. 479

The load-aware controller selection scheme can only reduce 480 tail latencies of responses under the homogeneous controllers. 481 However, the network environment is time-varying in real 482 situations, not only in the processing capabilities of controllers 483 but also in the number of routing requests from switches. 484 We further propose a delay-aware selection scheme of con-485 trollers, which can adapt to the variations of the network 486 environment. 487

# **IV. DELAY-AWARE SELECTION SCHEME** OF CONTROLLERS

We design the delay-aware selection scheme of controllers 490 while keeping the design goals of controller selection mecha-491 nism in mind. We first show an overview of the delay-aware 492 selection scheme. Then, we present two major components of 493 the delay-aware selection scheme, the selection models of con-494 trollers and the queue management mechanism of controller. 495

#### A. Overview of Delay-Aware Selection Scheme 496

To address those problems faced by the load-aware selection 497 scheme, we further design the delay-aware selection scheme of 498 controllers, which is adaptive to the heterogeneous controllers 499 as well as to the dynamic behaviours of flows. The delay-500 aware selection scheme needs to probe the response delays of 501 controllers for routing requests and send the routing requests 502 to the controller with the smallest response delay. The latency 503 of routing response denotes the time interval from sending 504 the routing request to receiving the flow rules generated by 505 the controller and is composed of the queueing delay and 506 the processing delay, as shown in Definition 1. The response 507 delay is an approximate evaluation of response time. Through 508

![](_page_18_Figure_9.jpeg)

Fig. 5. Overview of controller selection scheme. CS: Controller Selection scheduler, QM: Queue Management of controller.

probing response delays, the delay-aware selection scheme 500 can accommodate the heterogeneous controllers, while fewer 510 routing requests will be sent to the controllers with low 511 processing capabilities.

Definition 1: The latency of routing response denotes the time interval from sending the routing request to receiving the flow rules generated by the controller.

Furthermore, the capabilities of controllers are time-varying. With the development of SDN, there are more and more applications running in controllers. When switches send vast requests to the controller that has fast capability of response at before, a large number of requests have to queue in controllers if the capabilities of controllers decrease due to 52 other applications' overconsumption of resources. 522

Our delay-aware selection scheme includes two major com-523 ponents, controller selection (CS) and queue management 524 (QM). Recall the design goal of the selection scheme in 525 Section III-A, CS can achieve that the selection scheme is 526 light-weight, scalable and burst-immune, and QM achieves 527 the goal of adaptivity. First, we design a selection scheme 528 of controllers, which can select the controller based on a 529 little feedback from the controllers, and thus, is light-weight. 530 Second, instead of probing all controllers, the switch randomly 531 probes d controllers where  $d \ge 1$ . The probing is scalable 532 and independent of the network size. Third, the selection 533 scheme can make switches conduct once controller selection 534 for each flow or a batch of flows. It can make those requests 535 be processed by different controllers, and thus can avoid the 536 influence of the bursty and skew-flow requests. Last, through 537 estimating response delays of routing requests, switches can 538 send routing requests to the controller that has the smallest 539 response delay. Based on this estimation, the selection of 540 controllers can be adapted to the heterogeneous and time-541 varying processing capabilities. 542

Fig. 5 depicts the framework of the adaptive switch-to-543 controller selection scheme. When a request is issued at a 544 switch, the switch will work based on Algorithm 1. The switch 545 randomly probes d controllers, where  $d \ge 1$ . The d controllers 546 then estimate their response delays for the routing request 547 based on Algorithm 2 and return the response delays  $\psi$  to 548 the switch. If  $\psi_i$  of controller *i* exceeds the max response 549 delay  $RD_{max}$  limit, then controller *i* will return the response 550 delay  $\psi_i = -1$ . When the switch receives response delays 551 of d controllers, it will select the controller that has the 552 smallest response delay. If all response delays are lower than 0, 553 the switch will reselect a controller whose queue length does 554

517 518 519

Algorithm 1 The Selection of Controllers
<b>Require:</b> Controller set <i>C</i> , <i>d</i> .
1: randomly probe <i>d</i> controllers from <i>C</i> ;
2: send estimating request to the $d$ controllers;
3: $\psi \leftarrow$ response delays of <i>d</i> controllers;
4: if there exists $\psi_x \ge 0$ then
5: $id \leftarrow arg \min\{\psi_x \ge 0\};$
6: else
7: for $i = 1$ to C.length do
8: send estimating request to controllers $C[i]$ ;
9: $\psi_0 \leftarrow$ response time of $C[i]$ ;
10: <b>if</b> $\psi_0 >= 0$ <b>then</b>
11: $id = i;$
12: <b>break</b> ;
13: send the routing request to controller $C[id]$ .

Algorithm 2 Queue Management of Controller
<b>Require:</b> the max response delay, $RD_{max}$ .
1: receive an estimating request from switch s;
2: calculate the average processing time of requests $\bar{\nu}$ ;
3: if $r_i < \gamma_i$ then
4: $\psi_i \leftarrow \frac{r_i}{v_i} \bar{v};$
5: else
6: $\psi_i \leftarrow C[q_i mod \gamma_i] + (\lfloor q_i / \gamma_i \rfloor + 1) \times \overline{\nu};$
7: if $\psi_i > RD_{max}$ then
8: $\psi_i = -1;$
9: send $\psi_i$ to switch s.

not exceed limit for the routing request. Last, the switch willsend the request to the selected controller.

#### 557 B. The Selection Models of Controllers

When there are only a few controllers in the network, it is feasible to probe all controllers. Considering that this probing could occupy extra bandwidth of the secure link, it is essential to design a per-flow light-weight probing method.

1) Active Per-Flow Selection of Controllers: To deal with 562 the skew-flow requests and reduce response tail latencies. 563 the switches need to select controllers for each flow. Instead 564 of the controller-to-switch assignment, the active per-flow 565 selection of controllers makes that the routing requests from 566 the same switch can be processed by different controllers. 567 This selection scheme can fully exploit the capabilities of 568 controllers and efficiently reduce response tail latencies. 569

To reduce the bandwidth consumption of probing con-570 trollers, one method is to reduce the number of probed 571 controllers. There is a tradeoff between response tail latencies 572 and the number of probed controllers. Probing more controllers 573 can achieve fewer response tail latencies. However, that also 574 means more bandwidth consumption and computing overhead 575 in controllers. The number of controllers increases as the 576 network scale grows. In this case, checking all controllers has 577 a huge cost. To achieve the light-weight probing, we randomly 578 probe d controllers instead of checking all controllers, where 579  $d \ge 1$ . Furthermore, Azar *et al.* [18] have shown that having 580

just two random choices (i.e., d = 2) yields a large reduction in the maximum load over having one choice. This method has been widely studied and applied [19]. Inspired by this fact, our active per-flow selection is to probe two controllers and is thus scalable. Meanwhile, the active per-flow selection of controllers can efficiently reduce the bandwidth consumption of the secure link and the computing load of controllers. 581

Instead of probing the loads of the controllers, probing 588 response delays of controllers can better reduce response tail 589 latencies. The probing of response delays requires that these 590 controllers evaluate their own response delays for routing 591 requests. Since the selection of controllers only needs to get a 592 numerical value of response delay, the selection of controllers 593 is light-weight. Moreover, the heterogeneous and time-varying 594 processing capabilities of controllers increase the complexity 595 of evaluation for response delays of controllers. The response 596 delay estimate model will be introduced in Section IV-C. 597

2) Active Selection of Controllers for a Batch of Flows: 598 When switches meet bursty-flow requests or when the arrival 599 of routing requests is frequent, conducting a controller selec-600 tion for each request still aggravates the bandwidth consump-601 tion and the loads of controllers even if only two controllers 602 need to be probed in one controller selection. Conducting the 603 controller selection for a batch of arrival routing requests is 604 needed to increase the scalability of the controller selection 605 mechanism, when switches suffer bursty flows. 606

For active selection of controllers for a batch of flows, 607 the switch conducts one controller selection after it receives 608 the first flow request. When we set the batch size as  $\delta$ , 609 it means that the following  $\delta - 1$  requests will be sent to 610 the same controller with the first request. That is, the result 611 of controller selection for the first request will be shared by  $\delta$ 612 requests. In addition, the batch selection is irrelevant to the 613 rate of requests because  $\delta$  denotes the number of requests. 614 Therefore, although there would be the high rate of requests 615 in the beginning when the switch changes its controller, those 616 requests would not be sent to the same controller. Given that 617 the request arrival process at each switch is a Poisson process 618 with rate  $\lambda$ , the arrival duration for a flow is exponentially 619 distributed with mean  $1/\lambda$ . Therefore, the arrival duration of  $\delta$ 620 flows is also exponentially distributed with mean  $\delta/\lambda$ . 621

There is a tradeoff between the rounds of controller selection and the performance of controller selection. If  $\delta$  is too small, it is obvious that controller selection should be frequently conducted. However, it will decrease the performance of controller selection when  $\delta$  is too large. It is worth noting that it is unnecessary to adopt the batch selection when the arrival of flows is scattered.

#### C. Queue Management Mechanism of Controller

The queue management of controller makes it so that the selection of controllers can cope and quickly react to heterogeneous and time-varying processing time across controllers. Our queue management mechanism of controller includes response delay estimate and queue length bound.

1) Response Delay Estimate Model: As depicted in 635 Section III-D, the load-aware selection scheme of controllers 636

can not accommodate the heterogeneity of controller. To efficiently reduce the long-tail latency of routing response,
switches should select controllers with lower response delays
for each routing request.

Request response time consists of the queueing time and 641 the processing time. Furthermore, the queueing time is related 642 to the length of queue, which is equal to the number of 643 queueing requests. Meanwhile, to estimate the queueing time, 644 it is essential to estimate the processing time of each request. 645 In our design, the controller records  $v_i$ , which is the processing 646 time of the latest responded *j*th requests. Given the number 647 of the latest finished requests s, we calculate  $\bar{\nu}$ , which denotes 648 the average processing time of s requests in controller i. Thus, 649  $\bar{\nu} = \frac{1}{s} \sum_{i=1}^{s} \nu_i$ . We use  $\bar{\nu}$  to estimate the processing time of 650 requests. 651

Consider that the controller can process multiple routing 652 requests simultaneously. We use  $\gamma_i$  to denote the number of 653 requests that *controller*<sub>i</sub> can simultaneously process.  $q_i$  and  $r_i$ 654 are the number of queueing and running requests in con-655 troller *i*, respectively. To improve the system utilization, 656 the controller that has idle running slots should have a lower 657 estimated response delay. Therefore, the estimated response 658 delay  $\psi_i = \frac{r_i}{\gamma_i} \bar{\nu}$  when  $r_i < \gamma_i$ . When there are requests 659 queueing in a controller, the controller records the running 660 duration A[k] of the running request in the kth slot where 661  $1 \le k \le \gamma_i$ . The controller then estimates that the queueing 662 request will run in which slot. To achieve this goal, we use 663  $B[k] = |\bar{v} - A[k]|$  and then sort B[k] as non-decreasing order. 664 Then, the controller estimates that the request will run in 665  $[(q_i \mod \gamma_i) + 1]$ th slot. We can get the queueing time 666  $wt_i = B[(q_i \mod \gamma_i) + 1] + (\lfloor q_i / \gamma_i \rfloor) \times \overline{\nu}$ . At last,  $\psi_i = wt_i + \overline{\nu}$ 667 when  $r_i = \gamma_i$ . 668

In summary,  $controller_i$  uses the following estimation function for response delay:

$$\psi_{i} = \begin{cases} B[(q_{i}mod\gamma_{i})+1] + (\lfloor q_{i}/\gamma_{i} \rfloor+1) \times \bar{\nu} & r_{i} = \gamma_{i} \\ \frac{r_{i}}{\gamma_{i}}\bar{\nu} & r_{i} < \gamma_{i} \end{cases}$$
(6)

When a switch sends an estimating request to  $controller_i$ , 672 the *controller<sub>i</sub>* adopts the formula (6) to estimate the response 673 delay. In general,  $\gamma_i = 1$  means that the controller only can 674 process one request once. In this case,  $\psi_i = B[1] + (q_i + 1) \times \overline{\nu}$ . 675 We suppose that the controller is empty at the beginning. After 676 that,  $\bar{\nu}$  is equal to the average processing time of finished 677 requests when the number of finished messages is lower than 678 the given threshold s. 679

2) Cutting Tail Latencies: Since switches conduct con-680 troller selection simultaneously, there may be "herd behav-681 iors," wherein multiple switches are coaxed to direct requests 682 towards the best controller. There are many requests queueing 683 in a controller under herd behavior that could leads to long-tail 684 latencies of routing responses. Moreover, it is possible that the 685 probed controllers all have low processing capabilities or long 686 queues. In this case, it is not suitable to select a controller 687 from those probed controllers. 688

To cut long-tail latencies and reduce the influence of herd behavior, the controller necessitates to bound its queue length. Determining the length of queues at controllers is crucial. Queues that are too short lead to lower controller utilization, as resources may remain idle between allocations. Queues that are too long may incur excessive queuing delays.

When fewer requests are sent to a controller, this may incur 695 under-utilization of its resources, whereas significant delays 696 may occur when requests need longer processing time. Hence, 697 after estimating request response delay, we further design 698 a delay-aware bounding mechanism to bound queue length, 699 which can accommodate the heterogeneous and time-varying 700 capabilities of controllers. Meanwhile, bounding queue length 701 can efficiently weaken the influence of herd behaviors. At one 702 point, a controller receives a burst of flows, and that exceeds 703 the limit of queue length. After that, the following flows will 704 not queue in the controller until the queue length is lower than 705 the limit. This delay-aware bounding mechanism relies on the 706 response delay estimation of request, which is reported by the 707 controller. 708

In particular, we specify the maximum response delay 709  $RD_{max}$  that a request is allowed to wait in a queue. When 710 we are about to place a request at the queue of  $controller_i$ , 711 we first check the estimated response delay  $\psi_i$  reported by 712 controller<sub>i</sub>. Only if  $\psi_i < RD_{max}$  is the request queued at 713 that controller. We sample d controllers while conducting the 714 controller selection. If the d selected controllers all do not sat-715 isfy the maximum response delay constraint. The switch needs 716 to reselect a new controller. Using this method, the number of 717 requests in each controller gets dynamically adapted based on 718 the current capability of the controller. 719

 $RD_{max}$  is set to make requests prefer faster controllers. 720 Furthermore,  $RD_{max}$  can be dynamically regulated to fit the 721 variation of controllers' capabilities. For controllers that have 722 low processing capabilities,  $RD_{max}$  can limit the number 723 of queueing requests in these controllers. After that, these 724 requests can be sent to faster controllers. However, if  $RD_{max}$ 725 is too small, most of controllers refuse to receive new requests 726 because their response delays exceed the limit of  $RD_{max}$ . 727 In this case, it is necessary to extend the value of  $RD_{max}$ . 728

#### V. PERFORMANCE EVALUATION

729

736

We start with the evaluation methodology and scenarios. <sup>730</sup> In this section, we evaluate the selection schemes of controller <sup>731</sup> and the assignment mechanism of controller under the general settings of controllers, the queue length bound  $RD_{max}$ , <sup>733</sup> the heavy request-skews, the time-varying service rates and the batch selection. <sup>735</sup>

# A. Experimental Setup

We build a discrete-event simulator wherein workload 737 generators create flow requests at a set of switches. These 738 switches then employ the controller selection scheme to select 739 a controller for each request. Unless otherwise specified, 740 the network consists of 300 switches and 40 controllers. 741 However, when a large amount of flow requests emerges, 742 we need to employ more controllers to deal with the burstly 743 flows. The workload generators create flow requests according 744 to a Poisson arrival process to mimic arrival of user requests 745 at web servers [20]. Unless otherwise specified, the Poisson 746

arrival process is with  $\lambda = 0.5$  during 1ms. At the beginning, 747 the system is empty, and there are no requests. With the system 748 running, the switches start to produce flow requests and select 749 controllers for flow requests. We run the system 1000ms, 750 and the number of arrival flow requests is about 150,000. 751 Each controller maintains a FIFO request queue. Moreover, 752 in the settings of controllers, each controller can service a 753 tunable number of requests in parallel (10 in our settings). 754 The processing time each request experiences is drawn from 755 an exponential distribution (as in [21]) with a mean processing 756 time  $\mu^{-1} = 2$ ms. Furthermore, we incorporate controller 757 heterogeneity into the network as follows: each controller, 758 independently and with a uniform probability, sets its service 759 rate to a different  $\mu$ , where  $\mu = 0.5$  or  $\mu = 0.1$  in our settings. 760 To estimate the response delay of each request, we set s = 100, 761 That is,  $\bar{\nu}$  denotes the average processing time of the latest 762 finished 100 requests. We repeat every experiment 20 times 763 using different random seeds, and then get the average result. 764

<sup>765</sup> We compare our design against three strategies:

 Quantity-based allocation: Controllers achieve load balance through balancing the number of switches, which each controller manages. Currently, ONOS controller utilizes this strategy to balance the loads of distributed controllers.

2) Load-aware selection: The switch probes two controllers for each request and sends the request to the controller with fewer number of requests because probing all controllers is not scalable.

3) Delay-aware selection: The switch probes two controllers for each request and gets request response delays, which rely on the feedbacks of probed controllers. Then, the switch sends the request to the controller with smaller response delay.

#### 780 B. The Impact of d

Azar *et al.* [18] have shown that the situation of d = 2 yields 781 a large reduction in the maximum load over d = 1, while 782 each additional choice beyond two decreases the maximum 783 load by just a constant factor. Further, to verify the theory 784 and determine the value of d, we evaluate the impact of 785 d on the performance of the delay-aware selection strategy 786 under heterogeneous controllers where d denotes the number 787 of probed controllers. The processing time of each request 788 in each controller is drawn from an exponential distribution 789 where  $\mu$  was randomly set as 0.5 or 0.1. Other parameters are 790 the same as Section V-A. 791

Fig. 6 shows the delay-aware selection strategy significantly 792 reduces the mean response time as the value of d increase 793 from 1 to 2. However, probing more controllers just incur a 794 little bit reduction on the mean response time after d exceeds 795 2. In addition, multiple switches may simultaneously select 796 the same controller if each of them randomly probes more 797 controllers. This would in turn aggravates the load of the 798 controller. Thus, we set d = 2 in the next experiments. 799

#### 800 C. General Settings of Controllers

We evaluate the performances of different schemes with heterogeneous controllers. We employ 60 controllers where the

![](_page_21_Figure_11.jpeg)

Fig. 6. The impact of d on the performance of the delay-aware selection strategy.

![](_page_21_Figure_13.jpeg)

(a) Request response time with different (b) Throughput under different schemes.

![](_page_21_Figure_15.jpeg)

bound of maximal response delay  $RD_{max} = 20$ ms. The other settings of experiments are consistent with that of Section V-B.

Fig. 7(a) shows that our delay-aware scheme can signifi-805 cantly reduce the response tail latencies. Basically, all requests 806 can be finished in 50ms while adopting our delay-aware 807 scheme. The load-aware scheme achieves better performance 808 than the quantity-based scheme in Fig. 7(a). Over 90% of 809 requests can be processed during 150ms based on the load-810 aware scheme. The quantity-based scheme leads to long 811 response delays, and there are over 20% of requests whose 812 response delays are more than 200ms in Fig. 7(a). This is 813 because a large of requests queue in controllers that have lower 814 processing capabilities. Meanwhile, the load-aware scheme 815 also failed to respond to requests quickly. The processing 816 delays of flow requests in different controllers are different 817 when controllers are heterogeneous. As a consequence, select-818 ing a controller by the number of requests is not efficient. 819 Fig. 7(a) reveals that our delay-aware scheme achieved the 820 lowest response duration due to not only estimating response 821 delay but also cutting tail latencies. Fig. 7(a) also reveals that 822 our delay-aware scheme can be adapted to the system where 823 controllers have heterogeneous capabilities. 824

Fig. 7(b) shows that our delay-aware scheme can respond<br/>to more requests than the load-aware scheme and the quantity-<br/>based scheme can in the same time period. Meanwhile,<br/>the throughput difference among schemes grows as the system<br/>runs. In summary, with heterogeneous controllers, our delay-<br/>aware scheme can efficiently reduce response tail latencies and<br/>improve the throughput of controllers.825826827

# D. Impact of Queue Length Bound RD<sub>max</sub>

We evaluate the impact of  $RD_{max}$  on the performance of the delay-aware scheme. We set  $RD_{max} = 20$ ms,  $RD_{max} = 834$ 100ms and  $RD_{max} = \infty$  respectively.  $RD_{max} = \infty$  means

![](_page_22_Figure_1.jpeg)

(a) The performance of delay-aware(b) The difference of responded requests scheme under different  $RD_{max}$  settings.under different  $RD_{max}$  settings.

Fig. 8. The impact of  $RD_{max}$  on the performance and throughput of the delay-aware scheme.

that there is no limit on the queue length of controllers. Other 836 parameters are the same with Section V-C. 837

Fig. 8(a) reveals that our delay-aware scheme has bet-838 ter performance when  $RD_{max}$  has a smaller value. Under 839  $RD_{max} = 20$ ms, all requests can be finished in 20ms. How-840 ever, The maximal response delay is 40ms under  $RD_{max}$  = 841 100ms. Therefore, the performance of delay-aware scheme 842 under  $RD_{max} = 20$ ms is better than that of  $RD_{max} = 100$ ms. 843 Comparing with  $RD_{max} = \infty$ , delay-aware scheme can sig-844 nificantly reduce response tail latencies when  $RD_{max} = 20$ ms. 845 Furthermore, we compare the throughput of controllers under 846 different  $RD_{max}$  settings. Fig. 8(b) shows that the system can 847 respond to 300 more requests under  $RD_{max} = 20$  ms than 848  $RD_{max} = 100$ ms. It was obvious that controllers have higher 849 throughput when  $RD_{max} = 20$ ms than when  $RD_{max} = 100$ ms 850 and  $RD_{max} = \infty$ . The difference of responded requests 851 between  $RD_{max} = 20$ ms and  $RD_{max} = 100$ ms remains 852 stable. However, the difference of responded requests between 853  $RD_{max} = 20$ ms and  $RD_{max} = \infty$  grows as the system runs. 854 Fig. 8 reveals that the setting of  $RD_{max}$  can make flow 855 requests prefer the controllers that have faster processing 856 capabilities. Additionally, it is noteworthy that  $RD_{max}$  can not 857

be set too small, otherwise, there are no available controllers

#### E. Performance Under Heavy Request-Skews 860

while conducting the controller selection.

858

859

In this section, we study the effect of heavy demand skews 861 on the observed latencies where controllers are homogeneous 862 with average server rate  $\mu = 0.5$ . We set request-skew= 20% 863 and request-skew= 50%. That is, 20% and 50% of switches 864 generated 80% of the total requests towards the controllers. 865 Most of parameters are inherited from Section V-A. To enable 866 20% of switches to generate 80% of the total requests, 867 we randomly select 60 switches and set the arrival rate of flow 868 requests  $\lambda = 2$ . Other switches set  $\lambda = 0.125$ . Under request-869 skew= 50%, half of the switches set  $\lambda = 0.8$ , and the other 870 half of the switches set  $\lambda = 0.2$ . We set  $RD_{max} = 150$ ms 871 because there are too many requests in a short time and these 872 requests have to queue in controllers. 873

Fig. 9(a) shows that over 5% requests have response delays 874 of more than 200ms for the quantity-based scheme. The 875 quantity-based scheme suffers decreased performance due to 876 the request-skews. However, the load-aware and delay-aware 877 schemes can significantly reduce response tail latencies. Based 878

![](_page_22_Figure_9.jpeg)

Fig. 9. Request response time with different schemes under the heavy requestskews

![](_page_22_Figure_11.jpeg)

(a) Request response time with different(b) The difference of responded requests schemes between different schemes.

Fig. 10. The performances of different schemes under the time-varying service rates

on the load-aware and delay-aware schemes, all requests can 879 be finished in 10ms in Fig. 9(a). Fig. 9 reveals that the load-880 aware and delay-aware schemes achieve very similar perfor-881 mances since controllers are homogeneous in this section. 882 Meanwhile, the quantity-based scheme suffers decreased per-883 formance due to the request-skews. Under request-skews, 884 a part of switches generate a large number of the requests, 885 which incur long queues in some controllers for the quantity-886 based scheme. 887

Comparing Fig. 9(a) and Fig. 9(b), we can find that the 888 quantity-based scheme under request-skew= 50% achieves a 889 lower response latency than under request-skew= 20%. It is 890 because the load balance among controllers where request-891 skew = 50% is better than that of request-skew = 20%. 892 Moreover, Fig. 9 also reveals that our delay-aware scheme 893 can accommodate the heavy request-skews. 894

## F. Impact of Time-Varying Service Rates

In this section, we study the effect of the service rate 896 fluctuation on the tail latency of response. We change the 897 average service rates of controllers in the system every 50ms, 898 and all controllers randomly set  $\mu = 0.5$  or  $\mu = 0.1$ . Other 899 parameters are inherited from Section V-C.

895

900

Fig. 10(a) reveals that our delay-aware scheme can respond 901 to all requests in 60ms, the load-aware scheme can respond 902 to all requests during 80ms, and the response tail latency of 903 the quantity-based scheme is more than 200ms. Therefore, 904 our delay-aware scheme can efficiently reduce response tail 905 latencies. For the quantity-based scheme, it could not exploit 906 the feedbacks of controllers to select the controller and further 907

![](_page_23_Figure_1.jpeg)

(a) Request response time with different (b) Throughput under different schemes.

Fig. 11. The performances of different schemes under the batch selection.

suffers lower system utilization. The load-based scheme also
suffers lower performance because it fails to consider the timevarying service rate. Fig. 10(a) reveals that our delay-based
scheme can accommodate time-varying service rate.

Fig. 10(b) shows that our delay-aware scheme can respond 912 to more requests than the quantity-based scheme and the load-913 aware scheme. With the increase of system running time, 914 the advantage of the delay-aware scheme is more obvious 915 than the quantity-based scheme in Fig. 10(b). Meanwhile, 916 the difference of responded requests between the delay-917 aware scheme and the load-aware scheme cyclically fluctu-918 ates because the service rates of controllers are periodically 919 changed. In summary, our delay-aware scheme can be better 920 adapted to the time-varying service rate and can efficiently 921 reduce tail latencies of responses. 922

## 923 G. Evaluation for Batch Selection

In this section, we evaluate the performance of different schemes for batch arrival requests. We set the arrival rate of flow requests  $\lambda = 5$ . There are about 300, 000 requests during 200ms. To deal with the burst requests, we employ 150 controllers, and each controller can process 40 requests in parallel. We set the batch size  $\delta = 10$  and  $RD_{max} = 15$ ms, and other parameters are set as Section V-C.

Fig. 11(a) reveals that our delay-aware scheme can still 931 reduce tail latencies of responses under the batch selection. 932 933 Our delay-aware scheme can respond to all requests in 40ms, and the load-aware scheme can respond to all requests during 934 60ms. However, the response tail latency of the quantity-935 based scheme is more than 130ms. The performance of the 936 quantity-based scheme is mainly affected by controllers that 937 have lower processing capabilities. The tail latency of response 938 under the load-aware scheme is generated because it fails 939 to accommodate the heterogeneous controllers. Meanwhile, 940 Fig. 11(b) shows that our delay-aware scheme can respond to 941 more flow requests than either of the other two schemes. This 942 means that our delay-aware scheme can not only reduce the 943 tail latency, but can also improve the throughput of controllers 944 under the batch selection. 945

<sup>946</sup> *Impact of Batch Size*  $\delta$ : Furthermore, we evaluate the impact <sup>947</sup> of the batch size  $\delta$  on the performance of the delay-aware <sup>948</sup> scheme. Fig. 12(a) depicts the performances of the delay-<sup>949</sup> aware scheme under different batch sizes. Fig. 12(a) shows

![](_page_23_Figure_10.jpeg)

(a) The performance of delay-aware(b) The number of CS operations. CS: the scheme.

Fig. 12. The impact of the batch size  $\delta$ .

that the performance of the delay-aware scheme has a modest 950 decrease with the increase of the batch size  $\delta$ . When the 951 batch size  $\delta$  increases, it means that more requests will be 952 sent to the same controller, even though the controller has a 953 low processing capacity. As a consequence, the delay-aware 954 scheme suffers a little of decreased performance. we can 955 find that experiments show a similar performance under the 956 different settings of  $\delta$  in Fig. 12(a). Although more requests 957 would be directed to the controller with a low processing 958 capacity at some time when the value of  $\delta$  is larger, after that, 959 the following requests would have less chance to select the 960 controller. Therefore, the decrease of performance is modest. 961 Furthermore, Fig. 12(b) reveals that the number of controller 962 selection operations dramatically decreases when the batch 963 size  $\delta$  goes up. The fewer controller selection operations in 964 switch end will incur less bandwidth consumption in the secure 965 links between switches and controllers and less computing 966 load in controllers. In conclusion, active selection of con-967 trollers for a batch of flows can efficiently reduce the resourse 968 consumption of communication and computing with a little bit 969 of performance reduction. 970

#### VI. RELATED WORK

#### A. Network Softwarization

Network softwarization is a transformation trend for design-973 ing, implementing, and managing the 5G and next generation 974 networks. It exploits the benefits of software to enable the 975 redesign of network and service architectures, optimize the 976 expenditure and operational costs, and bring new values in 977 the infrastructures. The key enablers consist of the network 978 function virtualization (NFV), software-defined networking 979 (SDN) and cloud computing, etc. Meanwhile, 5G systems will 980 also rely on these technologies to attain the systems flexibility 981 and elasticity [2]. 982

Along with recent and ongoing advances in cloud com-983 puting, it has become promising to design flexible, scalable, 984 and elastic 5G systems benefiting from advanced virtualization 985 techniques of cloud computing [1]. Taleb and Ksentini [22] 986 introduces the Follow-Me Cloud concept and proposes its 987 framework. There has been research on the possibility of 988 extending cloud computing beyond data centers toward the 989 mobile end user, providing end-to-end mobile connectivity 990 as a cloud service [23]. Software defined networking (SDN) 991

971

acts as a promising enabler for network softwarization and 992 plays a crucial role in the design of 5G wireless networks [1]. 993 SDN has been also utilized in the virtualization of mobile 994 network functions [24]. Kempf et al. [25] describe an evolution 995 of the mobile Evolved Packet Core (EPC) utilizing SDN 996 that allows the EPC control plane to be moved into a data 997 center. Taleb et al. [1] introduce the concept of "Anything 998 as a Service" (ANYaaS), which relies on the reference ETSI 999 NFV architecture to orchestrate and manage important ser-1000 vices. NFV aims at offering network services using network 1001 functions implemented in software and deployed in an on-1002 demand and elastic manner on the cloud [26]. Medhat *et al.* [5] 1003 introduce a service function chaining taxonomy as the basis 1004 for the subsequent state-of-the-art analysis. 1005

#### 1006 B. SDN Scalability

To tackle the problem of scaling SDN, Xie et al. [27] 1007 resort to design the distributed controllers, such as Onix [28] 1008 and ONOS [14], which try to distribute the control plane 1009 while maintaining a logically centralized management. These 1010 approaches balance the load of controllers based on the 1011 number of switches, which can not efficiently reduce the tail 1012 latencies of responses. Aissioui et al. [29] propose a two-level 1013 hierarchical controller platform to address these scalability and 1014 performance issues in the context of 5G mobile networks. 1015

One key limitation of the distributed controllers is that 1016 the mapping between a switch and a controller is statically 1017 configured. The static configuration results in the uneven 1018 load distribution among the controllers. Bari et al. [30] 1019 propose a management framework, which periodically eval-1020 uates the current controller-to-switch assignment. After that, 1021 it needs to decide whether to perform a reassignment. If a 1022 reassignment is performed, the management framework also 1023 changes the controller-to-switch assignment in the network. 1024 Dixit et al. [31] propose ElastiCon, an elastic distributed 1025 controller architecture in which the controller pool is dynam-1026 ically grown or shrunk according to traffic conditions. In this 1027 case, the load is dynamically shifted across controllers, which 1028 similarly relieves the static mapping between a switch and a 1029 controller. Zhou et al. [32] propose a dynamic and adaptive 1030 algorithm (DALB), which is running as a module of SDN 1031 controller. The controllers in distributed environment can 1032 cooperate with each other to keep load balancing. Overloaded 1033 controller can be detected, and high-load switches mapped to 1034 this controller can be smoothly migrated to the under-load 1035 controllers. However, these dynamic frameworks require that 1036 the control plane monitors the state of the whole network and 1037 conducts the reassignments, which aggravate the computing 1038 load of the control plane. 1039

In contrast to these works, our approach relies on simple and 1040 inexpensive feedback of controllers and efficiently relieve the 1041 load of the control plane. Mao and Shen [33] use the principles 1042 of SDN to achieve the server load balancing by setting the 1043 SDN flow table, which does not aim to solve the load balance 1044 of the distributed controllers of SDN. Palma et al. [34] develop 1045 the QueuePusher, which is a queue management extension to 1046 OpenFlow controllers supporting the Open vSwitch Database 1047

Management Protocol (OVSDB) standard. QueuePusher can generate the appropriate queue configuration messages for switches. In addition, there have been lots of researches on how to achieve the consistency among distributed controllers [15]. These techniques are complementary to our approach.

#### VII. CONCLUSION

NFV and SDN can dynamically redistribute the flow across 1055 appropriate VNFs or service function chains if the controller 1056 configures a desired routing path for each network flow 1057 resulting from NFV applications. In this paper, we present 1058 the long-tail observations of the routing response latencies 1059 while using the up-to-date controller-to-switch assignment 1060 mechanism. To tackle this essential problem, we first propose 1061 a light-weight and load-aware switch-to-controller selection 1062 scheme to cut the long-tail response latency under the sim-1063 ple scenario of homogeneous controllers, and then design 1064 a general delay-aware switch-to-controller selection scheme 1065 to fundamentally cut the long-tail response latency for the 1066 more complicated heterogeneous controller scenario with per-1067 formance fluctuations. Through comprehensive performance 1068 evaluation, we demonstrate that our adaptive controller selec-1069 tion schemes can efficiently reduce response long-tail latencies 1070 and accommodate various system environments including the 1071 request-skews, the fluctuation of service rates and so on. 1072

#### REFERENCES

- [1] T. Taleb, A. Ksentini, and R. Jantti, "Anything as a service' for 5G mobile systems," *IEEE Netw.*, vol. 30, no. 6, pp. 84–91, Nov. 2016.
- T. Taleb, S. Dutta, A. Ksentini, M. Iqbal, and H. Flinck, "Mobile edge computing potential in making cities smarter," *IEEE Commun. Mag.*, vol. 55, no. 3, pp. 38–43, Mar. 2017.
- Y. Zhou, D. Zhang, and N. Xiong, "Post-cloud computing paradigms: 1079 A survey and comparison," *Tsinghua Sci. Technol.*, vol. 22, no. 6, 1080 pp. 714–732, 2017.
- M. Kablan, A. Alsudais, E. Keller, and F. Le, "Stateless network functions: Breaking the tight coupling of state and processing," in *Proc.* 1083 14th USENIX NSDI, Boston, MA, USA, Mar. 2017, pp. 97–112.
- [5] A. M. Medhat, T. Taleb, A. Elmangoush, G. A. Carella, S. Covaci, and T. Magedanz, "Service function chaining in next generation networks: State of the art and research challenges," *IEEE Commun. Mag.*, vol. 55, no. 2, pp. 216–223, Feb. 2017.
- [6] P. Zave, R. A. Ferreira, X. K. Zou, M. Morimoto, and J. Rexford, 1089
   "Dynamic service chaining with dysco," in *Proc. ACM SIGCOMM*, 1090
   Los Angeles, CA, USA, 2017, pp. 57–70.
- [7] D. L. C. Dutra, M. Bagaa, T. Taleb, and K. Samdanis, "Ensuring end-toend QoS based on multi-paths routing using SDN technology," in *Proc. IEEE Global Commun. Conf.*, Dec. 2017, pp. 1–6.
- [8] N. Gude et al., "NOX: Towards an operating system for networks," ACM 1095 SIGCOMM Comput. Commun. Rev., vol. 38, no. 3, pp. 105–110, 2008. 1096
- [9] A. Ksentini, M. Bagaa, T. Taleb, and I. Balasingham, "On using bargaining game for optimal placement of SDN controllers," in *Proc. IEEE Int. Conf. Commun. (ICC)*, May 2016, pp. 1–6.
- [10] J. Xie, D. Guo, Z. Hu, T. Qu, and P. Lv, "Control plane of software defined networks: A survey," *Comput. Commun.*, vol. 67, pp. 1–10, 1101 Aug. 2015.
- D. Kreutz, F. Ramos, P. E. Veríssimo, C. E. Rothenberg, S. Azodolmolky, and S. Uhlig, "Software-defined networking: A comprehensive survey," *Proc. IEEE*, vol. 103, no. 1, pp. 14–76, Jan. 2015.
- [12] A. Ksentini, M. Bagaa, and T. Taleb, "On using SDN in 5G: 1106 The controller placement problem," in *Proc. IEEE Global Commun.* 1107 *Conf. (GLOBECOM)*, Dec. 2016, pp. 1–6.
- M. Al-Fares, A. Loukissas, and A. Vahdat, "A scalable, commodity data center network architecture," ACM SIGCOMM Comput. Commun. Rev., vol. 38, no. 4, pp. 63–74, 2008.

1073

- [14] P. Berde et al., "ONOS: Towards an open, distributed SDN OS," in Proc. 1112 ACM HotSDN, 2014, pp. 1-6. 1113
- [15] A. Panda, W. Zheng, X. Hu, A. Krishnamurthy, and S. Shenker, "SCL: 1114 Simplifying distributed SDN control planes," in Proc. 14th USENIX 1115 1116 NSDI, Mar. 2017, pp. 329-345.
- [16] D. Hugheshallett, A. M. Gleason, and W. G. Mccallum, Calculus: Single 1117 1118 and Multivariable, Student Solutions Manual, 6th ed. Hoboken, NJ, USA: Wiley, 2013. 1119
- A. Bremler-Barr, Y. Harchol, and D. Hay, "OpenBox: A software-defined 1120 [17] 1121 framework for developing, deploying, and managing network functions," in Proc. ACM SIGCOMM, Salvador, Brazil, Aug. 2016, pp. 511-524. 1122
- 1123 [18] Y. Azar, A. Z. Broder, A. R. Karlin, and E. Upfal, "Balanced allocations," SIAM J. Comput., vol. 29, no. 1, pp. 180-200, 1999. 1124
- M. Mitzenmacher, A. W. Richa, and R. Sitaraman, "The power of two [19] 1125 1126 random choices: A survey of techniques and results," in Handbook of Randomized Computing, vol. 11. Norwell, MA, USA: Kluwer, 2000, 1127 pp. 255-312. 1128
- [20] R. Nishtala et al., "Scaling memcache at Facebook," in Proc. USENIX 1129 NSDI, 2013, pp. 385-398. 1130
- A. Vulimiri, P. B. Godfrey, R. Mittal, J. Sherry, S. Ratnasamy, and 1131 [21] S. Shenker, "Low latency via redundancy," in Proc. ACM Conf. Emerg. 1132 Netw. Experim. Technol., 2013, pp. 283-294. 1133
- 1134 [22] T. Taleb and A. Ksentini, "Follow me cloud: Interworking federated clouds and distributed mobile networks," IEEE Netw., vol. 27, no. 5, 1135 pp. 12-19, Sep./Oct. 2013. 1136
- T. Taleb, "Toward carrier cloud: Potential, challenges, and solutions," [23] 1137 IEEE Wireless Commun., vol. 21, no. 3, pp. 80-91, Jun. 2014. 1138
- 1139 [24] T. Taleb, A. Ksentini, and A. Kobbane, "Lightweight mobile core networks for machine type communications," IEEE Access, vol. 2, 1140 pp. 1128-1137, 2014. 1141
- [25] J. Kempf, B. Johansson, S. Pettersson, H. Lüning, and T. Nilsson, 1142 "Moving the mobile evolved packet core to the cloud," in Proc. IEEE 1143 8th Int. Conf. Wireless Mobile Comput., Netw. Commun. (WiMob), 1144 Oct. 2012, pp. 784-791. 1145
- T. Taleb et al., "EASE: EPC as a service to ease mobile core net-[26] 1146 1147 work deployment over cloud," IEEE Netw., vol. 29, no. 2, pp. 78-88, 1148 Mar./Apr. 2015.
- J. Xie, D. Guo, X. Zhu, B. Ren, and H. Chen, "Minimal fault-tolerant 1149 [27] coverage of controllers in IaaS datacenters," IEEE Trans. Services 1150 Comput., to be published, doi: 10.1109/TSC.2017.2753260. 1151
- [28] T. Koponen et al., "Onix: A distributed control platform for large-scale 1152 production networks," in Proc. USENIX OSDI, 2010, pp. 351-364. 1153
- A. Aissioui, A. Ksentini, A. M. Gueroui, and T. Taleb, "Toward elastic 1154 [29] distributed SDN/NFV controller for 5G mobile cloud management 1155 systems," IEEE Access, vol. 3, pp. 2055-2064, 2015. 1156
- [30] M. F. Bari et al., "Dynamic controller provisioning in software defined 1157 networks," in Proc. Int. Conf. Netw. Service Manage., Oct. 2013, 1158 pp. 18-25. 1159
- A. Dixit, F. Hao, S. Mukherjee, T. Lakshman, and R. Kompella, [31] 1160 1161 "Towards an elastic distributed SDN controller," in Proc. ACM HotSDN, Hong Kong, Aug. 2013, pp. 7-12. 1162
- Y. Zhou et al., "A load balancing strategy of sdn controller based on 1163 [32] distributed decision," in Proc. IEEE 13th Int. Conf. Trust, Secur. Privacy 1164 Comput. Commun., Sep. 2014, pp. 851-856. 1165
- Q. Mao and W. Shen, "A load balancing method based on SDN," in 1166 [33] Proc. 7th Int. Conf. Meas. Technol. Mechatronics Autom., Jun. 2015, 1167 pp. 18-21. 1168
- D. Palma et al., "The queuepusher: Enabling queue management in 1169 [34] openflow," in Proc. 3rd Eur. Workshop Softw. Defined Netw., Sep. 2014, 1170 pp. 125-126. 1171

![](_page_25_Picture_22.jpeg)

Junjie Xie received the B.S. degree in computer science and technology from the Beijing Institute of Technology, Beijing, China, in 2013, and the M.S. degree in management science and engineering from the National University of Defense Technology, Changsha, China, in 2015, where he is currently pursuing the Ph.D. degree. He is also a joint Ph.D. student with the University of California, Santa Cruz (UCSC), CA, USA, from 2017. His study in the UCSC is supported by the China Scholarship Council. His research interests include distributed

systems, data-center networks, and software-defined networks. 1183

![](_page_25_Picture_25.jpeg)

Deke Guo received the B.S. degree in industry engi-1184 neering from the Beijing University of Aeronautics 1185 and Astronautics, Beijing, China, in 2001, and the 1186 Ph.D. degree in management science and engineer-1187 ing from the National University of Defense Tech-1188 nology, Changsha, China, in 2008. He is currently a 1189 Professor with the College of Information System 1190 and Management, National University of Defense 1191 Technology, and a Professor with the School of 1192 Computer Science and Technology, Tianjin Uni-1193 versity. His research interests include distributed 1194

systems, software-defined networking, data center networking, wireless and 1195 mobile systems, and interconnection networks. He is a member of the ACM. 1196

![](_page_25_Picture_28.jpeg)

Xiaozhou Li received the Ph.D. degree in com-1197 puter science from Princeton University in 2016. 1198 He is currently a Software Engineer with Barefoot 1199 Networks, where he is building new networking 1200 systems with fast programmable network chips. His 1201 research improves the performance, scalability, and 1202 efficiency of datacenter services, with a particular 1203 focus on combining new hardware and infrastructure 1204 capabilities with careful architectural design and 1205 algorithm engineering. 1206

![](_page_25_Picture_30.jpeg)

Yulong Shen received the B.S. and M.S. degrees in 1207 computer science and the Ph.D. degree in cryptog-1208 raphy from Xidian University, Xian, China, in 2002, 1209 2005, and 2008, respectively. He is currently a 1210 Professor with the School of Computer Science and 1211 Technology, Xidian University, and also an Asso-1212 ciate Director of the Shaanxi Key Laboratory of 1213 Network and System Security. He has also served on 1214 the technical program committees of several interna-1215 tional conferences, including the NANA, the ICEBE, 1216 the INCoS, the CIS, and the SOWN. His research 1217 1218

interests include wireless network security and cloud computing security.

![](_page_25_Picture_33.jpeg)

Xiaohong Jiang received the B.S., M.S., and Ph.D. 1219 degrees from Xidian University, China, in 1989, 1220 1992, and 1999 respectively. From 2005 to 2010, he 1221 was an Associate professor, Tohoku University. He 1222 is currently a Full Professor with Future University 1223 Hakodate, Japan. His research interests include com-1224 puter communications networks, mainly wireless 1225 networks and optical networks, network security, 1226 and routers/switches design. He has authored or co-1227 authored over 300 technical papers at premium inter-1228 national journals and conferences, which include 1229

over 70 papers published in top IEEE journals and top IEEE conferences. He 1230 was a recipient of the Best Paper Award of IEEE HPCC 2014, IEEE WCNC 1231 2012, IEEE WCNC 2008, IEEE ICC 2005-Optical Networking Symposium, 1232 and IEEE/IEICE HPSR 2002. He is a Member of IEICE. 1233