User-Centric Interference-Aware Load Balancing for Cloud-Deployed Applications

Seyyed Ahmad Javadi[®] and Anshul Gandhi[®]

Abstract—VMs deployed in cloud environments are prone to performance interference due to dynamic and unpredictable contention for shared physical resources among colocated tenants. Current provider-centric solutions, such as careful co-scheduling of VMs and/ or VM migration, require a priori profiling of customer VMs, which is infeasible in public clouds. Further, such solutions are not always aware of the user's SLO requirements or application bottlenecks. This paper presents DIAL, an interference-aware load balancing framework that can directly be employed by cloud users without requiring any assistance from the provider. The key idea behind DIAL is to infer the demand for contended resources on the physical hosts, which is otherwise hidden from users. Estimates of the colocated load are then used to dynamically shift load away from compromised VMs without violating the application's tail latency SLOs. We implement DIAL for web and online analytical processing applications, and show, via experimental results on OpenStack and AWS clouds, that DIAL can reduce tail latencies by as much as 70 percent compared to existing solutions.

13 Index Terms—Cloud computing, performance interference, load balancing

14 **1** INTRODUCTION

5

6

7

g

10

11

12

THE benefits of cloud computing are undeniable – low 15 cost, elasticity, and the ability to pay-as-you-go. Not 16 surprisingly, many online services and applications are 17 now hosted on the cloud, on virtual machines (VMs). 18 Despite its popularity, however, applications deployed 19 in the cloud can experience undesirable performance 20 effects, the most severe of which is interference. Perfor-21 mance interference is caused by contention for physical 22 resources, such as CPU or last-level cache, among colo-23 cated VM users/tenants. 24

Interference is an undesirable side-effect of a fundamen-25 tal design principle of the cloud, namely, multi-tenancy 26 27 (sharing of a physical server among users). While certain resources, such as CPU, can be partitioned among colocated 28 29 VMs by cloud providers, other resources, such as processor caches, are notoriously hard to partition [1]. Nonetheless, 30 partitioning of resources among tenants can adversely 31 impact cloud resource utilization. Further, resource conten-32 tion depends on the workload of all colocated tenant VMs, 33 and is thus dynamic and unpredictable [2]; as a result, static 34 partitioning is not a useful solution. 35

Prior work on interference mitigation has typically focused on provider-centric solutions. A popular approach is to profile applications and co-schedule VMs that do not contend on the same resource(s) [3], [4], [5]. However, since interference is dynamic and can emerge unpredictably, statically co-scheduling VMs will not suffice. VM migration can help in this case, but interference is volatile and short-lived,

Manuscript received 28 Aug. 2018; revised 12 June 2019; accepted 10 Sept. 2019. Date of publication 0 . 0000; date of current version 0 . 0000. (Corresponding author: Seyyed Ahmad Javadi.) Recommended for acceptance by A. Deshpande. For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below.

Digital Object Identifier no. 10.1109/TCC.2019.2943560

often lasting for only a couple minutes [2]; by contrast, 43 migration can take several minutes [6] and can incur over- 44 heads [7], especially for stateful applications [8].

A key challenge that has not been addressed with 46 regards to interference is the *lack of visibility and control* 47 *between the provider and the tenant,* especially in public 48 clouds [9]. Specifically, tenant VMs in a public can not, or 49 should not, be profiled a priori by the provider due to privacy concerns [10]. Further, providers are not always aware 51 of the cloud user's Service Level Objective (SLO) requirements or the user application's bottleneck resources. 53

In this paper, we make the case for a *user-centric* interference 54 mitigation approach which can be employed by the tenant 55 without requiring any assistance from the provider or hypervi-56 sor. Such user-centric solutions empower the tenant to have 57 greater control on their application performance, which is often 58 the most important criteria for users. Further, user-centric solu-59 tions are, by definition, aware of the user application and its 60 SLOs, and can appropriately react to the onset or termination 61 of interference.

We present DIAL, a dynamic solution for mitigating 63 interference in load-balanced cloud deployments. We con- 64 sider a generic cloud-deployed application that has a tier of 65 worker nodes hosted on multiple VMs and experiencing 66 unpredictable interference from colocated VMs (owned by 67 other cloud tenants), as shown in Fig. 1. The incoming load 68 is distributed among the worker nodes via one or more load 69 balancers. This generic model is widely applicable, for 70 example, for web applications (where workers are web or 71 application servers), online analytical processing (OLAP) 72 systems like Pinot [11], etc. 73

The key idea behind DIAL is to *infer* contention in colo- 74 cated VMs. Specifically, by monitoring its own application 75 performance, the user can estimate the colocated load that 76 can induce the observed level of performance degradation, 77 without requiring any assistance from colocated users or 78 the hypervisor. We find that, in addition to estimating the 79 colocated load, it is also important to determine the resource 80

The authors are with the PACE Lab; Stony Brook University, Stony Brook, NY, 11794. E-mail: {sjavadi, anshul}@cs.stonybrook.edu.

^{2168-7161 © 2019} 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.



Fig. 1. Illustration of a generic cloud application containing LBC and processing tier deployed on multiple foreground (fg) VMs experiencing interference from background (bg) VMs.

that is under contention, as this dictates the impact of inter-81 ference on performance. 82

To address the dynamic nature of interference, DIAL 83 adapts the load distribution of incoming requests among 84 user VMs. We introduce a model for interference, based in 85 queueing theory [12], [13], to understand the impact on per-86 formance of contention at shared physical resources. DIAL 87 then optimizes the time-varying load distribution among 88 worker VMs to reduce tail latency. 89

90 We implement and experimentally evaluate DIAL for 91 two specific application classes:

- 92 Web applications: We implement DIAL on HAProxy [14], and evaluate DIAL's benefits using two 93 popular web applications with varying workload 94 under CPU, network, disk, and cache interference on 95 OpenStack and AWS clouds. Our experimental 96 results show that DIAL reduces 90% ile response 97 times by as much as 70 percent compared to interfer-98 ence-oblivious load balancers. Further, compared to 99 100 existing interference-aware solutions, DIAL reduces tail response times by as much as 48 percent. 101
- 2) OLAP Systems: We implement DIAL for a popular and 102 open-source OLAP system called Pinot that has being 103 used in production clusters at LinkedIn and Uber. Our 104 experimental results on a KVM cluster show that 105DIAL can reduce 95% ile query completion times by 106 16-40 percent under CPU and LLC contention. 107

A preliminary version of this paper appeared in the 108 ICAC 2017 conference, but only focused on web applica-109 110 tions [15]. This version extends the performance modeling, optimization, implementation and evaluation of DIAL for 111 the case of OLAP systems, such as Pinot [11]. Further, this 112 version provides the analytical proof sketches for the theo-113 retical results that guide DIAL. 114

2 DIAL SYSTEM DESIGN 115

DIAL is a user-centric interference mitigation solution 116 117 designed for clouds that directly empowers the users. DIAL can complement provider-centric solutions, especially when 118 119 provider efforts to mitigate interference are not enough to avoid specific SLO violations for user applications. 120

2.1 Problem Statement and Scope 121

Fig. 1 illustrates a typical multi-tier cloud deployed applica-122 tion. The worker nodes process the incoming requests, and 123 are illustrated as a tier of VMs. Incoming requests are distrib-124 uted among the worker nodes using an application-specific 125



scheduler or dispatcher or load balancer; we abstract this 134

entity as a Load Balancing Tier (LBT). Our focus in this paper 136 is on the worker nodes and the LBT; specifically, we propose 137 a new technique to dynamically infer the interference on 138 worker nodes and adjust the load balancing weights for the 139 worker VMs in the LBT. We assert that the LBT is ideally 140 suited to mitigate the effects of volatile interference on worker 141 VMs as the LBT acts at the front-end for the worker tier. 142

The worker tier is hosted on multiple foreground (fg) 143 VMs, each of which is hosted on a physical machine (PM); 144 we highlight the worker tier fg VMs in Fig. 1. Each PM may 145 also host background (bg) VMs that do not belong to the fg 146 user, as shown in Fig. 1. The fg and bg VMs on a PM can 147 contend for shared physical resources, such as CPU, net- 148 work bandwidth (NET), disk I/O bandwidth (DISK), and 149 last-level-cache (LLC), resulting in interference. Note that 150 the fg user does not have visibility into the bg VMs; in fact, 151 the fg user is unaware of bg VMs. 152

2.2 DIAL Overview

Our solution, DIAL, is a user-centric dynamic Interference- 154 Aware Load Balancing framework. The design of DIAL 155 addresses two key questions: 156

- (i) How can users estimate the interference that their VMs 157 are experiencing without any assistance from the pro- 158 vider, hypervisor, or colocated users? (Section 2.3) 159
- Given this information, how should users dyna- 160 (ii) mically distribute load among their VMs to mini- 161 mize tail latencies in the presence of interference? 162 (Section 2.4) 163

The key idea in DIAL is to estimate, from within a user VM, 164 the amount of interference being induced by colocated VMs, 165 and then adapt the incoming load intensity for each user VM 166 accordingly. Fig. 2 shows a high-level overview of DIAL's 167 control flow. DIAL monitors performance metrics from 168 within the VMs and signals interference if tail latency goes 169 above a certain threshold (detection, see Section 2.3.1). DIAL 170 then determines if the detected event is a load change for the 171 application or a resource contention event (classification, see 172 Section 2.3.2). Depending on the source of contention, DIAL 173 quantifies, or infers, the severity of resource contention (esti- 174 mation, see Section 2.3.3). Based on this quantification, DIAL 175 determines the theoretically-optimal load balancing weights 176 that the user application should employ to mitigate the impact 177 of contention (see Section 2.4). The above steps are continually 178 employed at runtime, enabling DIAL to respond dynamically 179 to contention. 180

2.3 User-Centric Estimation of Interference

We define amount of interference as the fraction of available 182 physical resources that are in use by colocated background 183 VMs. In the context of Fig. 1, the amount of interference is 184 the fraction of physical resources on a PM that are in use by 185

153

181



Fig. 3. Performance of an OpenStack-deployed Apache web server under interference from colocated VMs running microbenchmarks.

the colocated bg VMs, and are thus unavailable to the fg VM
on that PM. As we show below, estimating the amount of
interference is non-trivial as it requires classification and
modeling of interference.

190 2.3.1 Impact of Interference on Tail Latencies

Interference is known to impact application response times [3], [16], [17]. DIAL leverages this fact to estimate the amount of interference that an fg VM is experiencing because of resource contention created by colocated bg VMs. Specifically, DIAL aims to infer the amount of interference, or resource contention, that the bg VMs must be creating to effect the observed rise in fg response times.

Fig. 3 shows the impact of different types of resource 198 contention on the 90% ile response time of an OpenStack 199 cloud-deployed Apache web server VM hosting files and 200 driven by the httperf load generator. We create contention 201 for this fg VM by running various microbenchmarks in colo-202 cated bg VMs. The x-axis denotes the percentage of total 203 resource usage, which is the sum of resource usage by the 204 fg VM and all colocated bg VMs, normalized by peak 205 resource capacity or bandwidth. For example, if the total 206 network bandwidth usage is 80MB/s, and the peak network 207 bandwidth is about 115MB/s, then the resource utilization 208 is $80/115 \approx 0.7$. 209

We make two observations from this figure:

210

(i) response time increases considerably under interference, (ii)
 the relationship between total resource usage and response time
 depends on the exact resource under contention.

Detecting Interference. DIAL uses the first observation to 214 *detect* when the fg application VM is under interference. 215 216 Specifically, from Fig. 3, we see that application response times, or T_x , are initially low and stable (left of the graph). 217 However, once the total resource usage increases (right of 218 219 the graph), because of the increased resource demand from bg VMs, the fg response times rise sharply. Thus, DIAL sig-220 nals interference when T_x goes beyond the 95 percent *confi*-221 dence intervals (around the mean of periodically monitored 222 223 tail latencies) observed during no or low interference.

Need for Identifying the Source of Interference: The second observation suggests that using tail response times to estimate interference will require knowledge of the specific resource that is under contention.

228 2.3.2 Classifying Interference using Decision Trees

Our next task is to classify the source of interference, which is defined as the *dominant resource under contention*. Note that it is possible for several resources to be simultaneously under contention; however, we only consider dominant ²³² resource contention. Our key idea in classification is to ²³³ observe the impact of interference on user metrics, such as ²³⁴ CPU utilization and I/O wait time, which can be easily ²³⁵ obtained from within the VM via the /proc subsystem. ²³⁶

DIAL uses decision trees to classify contention. The decision tree classifier is trained by running controlled interference experiments using microbenchmarks and monitoring the metrics in each case. After training, the decision tree can classify the source of interference, even for unseen workloads, based on the observed metric values (Section 4.3.2). 242

Distinguishing Interference from Workload Variations. An 243 application's response time can degrade for various reasons, such as workload surges, in addition to interference. 245 While our detection methodology detailed in Section 2.3.1 246 does not distinguish between interference and workload 247 variations. DIAL makes this distinction at the classification 248 stage by again leveraging the decision tree classifier. Specifi- 249 cally, to distinguish interference from workload variations, 250 DIAL normalizes the observed metric values with *predicted* 251 values based on monitored workload intensity. Prior work 252 has shown that linear models can accurately predict CPU 253 usage based on workload intensities [18]. We thus use linear regression to predict the metric values as a linear 255 function of the number of requests seen in the past monitor- 256 ing interval. 257

The intuition behind this approach is that, in the absence 258 of interference, the normalized values will be close to 1 259 under workload variations. The decision tree can thus use 260 the deviation of the observed metrics from the normalized 261 metrics to distinguish workload changes from interference. 262

2.3.3 Queueing-based Model for Interference

The final step is to use the classification information to estimate the amount of interference, which is the *fraction of* 265 *resources that are in use by colocated bg VMs*. Once we have 266 these estimates, DIAL can redistribute incoming load accordingly to mitigate the impact of interference (Section 2.4). 268

From Fig. 3, we see that tail response times increase non-269 linearly with the total usage of the resource under conten-270 tion. Recall that the total resource usage is the sum of 271 resource usage of the fg VM (can be monitored by the fg 272 user) and all colocated bg VMs (cannot be monitored by the 273 fg user). Our key idea is to model this non-linear relationship for each resource; this allows inferring the resource 275 usage of the colocated bg VMs based on observed fg tail 276 latencies, which in turn gives us the amount of interference. 277

Modeling Interference. We employ queueing theory to 278 model the non-linear relationship between resource usage 279 and tail latencies. Queueing models suggest that the tail 280 response time for an application is inversely proportional 281 to the unused capacity of the VM [12]. Mathematically, 282 $T_x \sim 1/(1 - \rho_{fg})^{\alpha}$, for some parameter α , where ρ_{fg} is the 283 resource load of the fg application (such as CPU utilization 284 or I/O bandwidth utilization), normalized to peak resource 285 usage; that is, $0 \le \rho_{fg} \le 1$. Prior work [19] has shown that 286 $\alpha = 2$ works well for practical settings given the high variability in real workloads. Prior theoretical work has also 288 shown that a quadratic term in the denominator can result 289 in better predictability under high loads [20]. However, 290 such models do *not* account for interference.

Under interference, the fg application experiences con- 292 gested resources due to colocated bg VMs. As a result, the 293

303

294 application experiences higher load than it would in the 295 absence of interference. We model this effect by adding the resource usage of colocated bg VMs to that of the fg VM, 296 resulting in fg response times being inversely proportional 297 to $(1 - (\rho_{fg} + \rho_{bg}))$. The sum of loads exerted by the fg 298 and bg VMs, $(\rho_{fg} + \rho_{bq})$, represents the normalized total 299 resource utilization. We thus approximate x%ile response 300 time as: 301

$$T_x = c_0 + c_1 / (1 - \rho_{fg} - \rho_{bg}) + c_2 / (1 - \rho_{fg} - \rho_{bg})^2, \tag{1}$$

where \vec{c} is the coefficient vector that depends on the *specific resource under contention*. The polynomial function in Eq. (1) is inspired by prior work on queueing systems [20], [21] to interpolate between low load (linear term in denominator) and high load (quadratic term in denominator) regimes.

To determine the coefficients, we train the model in 309 Eq. (1) by creating different levels of resource usage and 310 monitoring the T_x of fg VMs (see Section 2.5.1). We then use 311 312 multiple linear regression over this training data to derive the resource-specific coefficients. While Eq. (1) is inspired 313 by queueing models, it can accurately track the relationship 314 315 between tail response times and resource usage for realistic 316 web applications, as we show in Section 4.3.

Applying the Model to Estimate Interference: Eq. (1) can be easily employed to estimate the amount of interference. After detection and classification, the fg user can estimate ρ_{bg} by monitoring T_x and ρ_{fg} , and solving Eq. (1) for ρ_{bg} .

321 2.4 Interference-Aware Load Balancing

Interference-aware load balancing is the key component of 322 323 DIAL. When there is no interference, balancing the load among VMs works well to provide low response times. 324 However, if one of the VMs is facing interference (can be 325 estimated via the above-described interference modeling), 326 then its share of the load must be adjusted accordingly. One 327 might think that reducing the share of load in proportion to 328 the available capacity at the compromised VM, $(1 - \rho_{ba})$, 329 should work well. Unfortunately, this approach can be far 330 from optimal, as we show via experiments in Section 4.3. 331

2.4.1 Minimizing Tail Response Times for Web Applications

334 To minimize application tail response times under interference, we again employ queueing theory. We first consider 335 the case where any VM can serve an incoming request, as in 336 the case of a web application tier. Consider a cluster of n337 VMs, with VM *i* facing interference of $\rho_{bq,i}$. Let the fraction 338 of total incoming load that is directed to VM i be p_i ; we refer 339 to p_i as the weight assigned by the load balancer (LB) to VM 340 *i*. If the total arrival rate for the application is *a*, the arrival 341 rate for VM *i* is $a \cdot p_i$. Our goal is to determine the p_i s that 342 minimize the x%ile response time, T_x . 343

To obtain a simple closed-form expression for the theo-344 345 retically optimal p_i s, we model each VM as an M/M/1 346 system. By focusing on the dominant resource that is caus-347 ing interference, as classified using the decision tree, we employ the M/M/1 model to represent the contention at 348 349 the dominant resource. While this is an oversimplification, the resulting closed-form tail latency expression enables the 350 optimization and determination of theoretically-optimal 351 load balancer weights. We note that the resulting p_i s are 352

only optimal under the M/M/1 model; we refer to these as 353 the "theoretically optimal" weights in the rest of the paper. 354

For the M/M/1 model, the response time is known to fol- 355 low an Exponential distribution [12]. We can thus obtain 356 any tail probability of response time by using the CDF of 357 the Exponential distribution. Under the M/M/1 assump- 358 tion, T_x for a cluster of *n* VMs is approximated as: 359

$$T_x \approx \sum_{i=1}^n p_i \cdot \frac{-\ln(1 - x/100)}{r_i - a \cdot p_i},$$
 (2)

where r_i represents the throughput of VM *i* (with contention). Since interference reduces the throughput of the compromised VM, we set $r_i = r \cdot (1 - \rho_{bg,i})$, where *r* is the peak 364 throughput of an application VM. For example, if the peak 365 throughput of our Apache server is r = 1000 req/sec, and it 366 is experiencing an estimated interference of $\rho_{bg} = 0.6$, then 367 we set $r = 1000 \times 0.4 = 400$ req/sec. 368

Eq. (2) above works for all percentiles of response time. 369 For example, if x = 90, meaning we focus on the 90%ile 370 response time, then the term in the numerator becomes 371 $-\ln(1-0.9) = \ln 10$. For 95%ile response times, the numerator becomes $\ln 20$. Interestingly, the optimization for p_i s 373 discussed below does *not* depend on the numerator value 374 (since it is independent of p_i), and thus *our results apply, asis, for any percentile of response times,* including the median. 376

Given *a* (monitored at the LB) and r_i (derived as discussed above using interference estimation from Section 378 2.3), T_x can be expressed as a function of p_i via Eq. (2). We 379 can now derive the theoretically optimal weights, p_i s, that 380 minimize T_x in Eq. (2) via calculus, as presented below. 381

Lemma 1. The theoretically optimal load split for minimizing T_x 382 for a cluster of n VMs with total arrival rate a and individual 383 VM throughputs r_i is given by: 384

$$p_{i}^{*} = \left(r_{i}\sum_{j=1}^{n}\sqrt{r_{j}} - \sqrt{r_{i}}\sum_{j=1}^{n}r_{j} + a\sqrt{r_{i}}\right) / \left(a\sum_{j=1}^{n}\sqrt{r_{j}}\right).$$
(3) 386
387

Proof. The proof proceeds via mathematical induction on *n*. 388 We first prove the base case for n = 2. Let the probability 389 of sending a request to VM 1 (with throughput r_1) be *p*; 390 thus, arrival rate into VM 1 is $a \cdot p$. Then, under the M/ 391 M/1 queueing model [12], the response time for VM 1 is 392 distributed as $Exp(r_1 - a \cdot p)$. Based on this, the x%ile 393 response time is $\frac{-\ln(1-x/100)}{r_1 - a \cdot p}$. Likewise, the x%ile response 394 time for VM 2 (with arrival rate $a \cdot (1 - p)$) is $\frac{-\ln(1-x/100)}{r_2 - a \cdot (1-p)}$. 395 We now approximate T_x for the 2-VM system as: 396

$$T_x \approx p \cdot \frac{-ln(1-x/100)}{r_1 - a \cdot p} + (1-p) \cdot \frac{-ln(1-x/100)}{r_2 - a \cdot (1-p)}.$$

We now derive the optimal value of $0 \le p \le 1$ that minimizes T_x . Taking the derivative of T_x w.r.t. p, we get:

$$p_1^* = p^* = \frac{r_1\sqrt{r_2} - r_2\sqrt{r_1} + a\sqrt{r_1}}{a(\sqrt{r_1} + \sqrt{r_2})}.$$
402
403

Now assume that the above expression for p^* is true for 404 n = k. Then, for n = (k + 1), we partition the (k + 1) VM 405 system into a single VM with request probability p_n and a 406 k-VM system with request probability $(1 - p_n)$. For the 407 k-VM system (with primed variables) with request rate 408

409

411

415

420

452

454

455

 $a' = a \cdot (1 - p_n)$, by the inductive hypothesis, we have:

$$p_i'^* = \frac{r_i \sum_{j=1}^k \sqrt{r_j} - \sqrt{r_i} \sum_{j=1}^k r_j + a' \sqrt{r_i}}{a' \sum_{j=1}^k \sqrt{r_j}}$$

412 The approximate x%ile response time for the (k + 1)-VM 413 system can then be written as:

$$T_x \approx p_n \cdot \frac{-ln(1 - x/100)}{r_n - a \cdot p_n} + (1 - p_n) \cdot \sum_{j=1}^k p_i'^* \cdot \frac{-ln(1 - x/100)}{r_i - a' \cdot p_i'^*}$$
(4)

⁴¹⁶ Note that T_x is itself a function of p_n since request rate for ⁴¹⁷ the *k*-VM system is $a' = a(1 - p_n)$). We now derive the ⁴¹⁸ theoretically optimal p_n^* by differentiating Eq. (4) to get:

$$p_n^* = \frac{r_n \sum_{j=1}^n \sqrt{r_j} - \sqrt{r_n} \sum_{j=1}^n r_j + a\sqrt{r_n}}{a \sum_{j=1}^n \sqrt{r_j}}.$$
 (5)

The remaining *k* theoretically optimal probabilities can then be derived by noting that $p_i^* = (1 - p_n^*) \cdot p_i^{**}$.

Note that p_i^* depends on the estimates of r_i , thus necessitat-423 ing the interference estimation of Section 2.3. Also note that p_i^* 424 depends on the total arrival rate, a. This is to be expected since, 425 for example, if the arrival rate is very low, we can send all 426 requests to the VM with the highest throughput to minimize 427 response times; however, if the arrival rate is very high, then a 428 single VM cannot handle all requests, and we have to distrib-429 ute the load. Importantly, both r_i and a can change unpredict-430 ably at any time (r_i due to interference and a due to variable 431 customer traffic), motivating the need for a *dynamic* solution 432 instead of existing static solutions. 433

434 2.4.2 Minimizing Tail Response Times for OLAP 435 Applications

We now extend the above analysis to the case where only a subset of workers (replicas) can serve an incoming request due to data locality, as in the case of OLAP systems. Let the number of replicas be c. Let us first consider the case where one worker, say w, out of n, is under interference. Let the non-interference throughput be r and that of w be $r_w < r$.

442 In the absence of interference, 1/c fraction of requests that have a replica on *w* would be sent there by the LB; further, the 443 arrival rate to w would be a/n, assuming a fair distribution of 444 445 replicas among workers. In the presence of interference, let the fraction sent to w be p, and so the fraction sent to the remaining 446 (c-1) replicas is $\frac{1-p}{c-1}$. Thus, the arrival rate into w is now 447 $\frac{a}{n} \cdot \frac{p}{1/c} = \frac{a}{n} \cdot p \ c$, and the fraction of *all* requests that go to *w* is $q = \frac{p}{n} \cdot \frac{c}{n}$. Likewise, arrival rate into *each* non-interference worker 448 449 is $\frac{a}{n} + \frac{a}{n} \cdot \frac{1-p}{(n-1)} = a \cdot \frac{1-q}{n-1}$, and the fraction of all requests that go 450 to each non-interference worker is $\frac{1-q}{n-1}$ 451

Using the M/M/1 model [12], we have, similar to Eq. (2):

$$T_x \approx q \cdot \frac{-\ln(1 - x/100)}{r_w - a \cdot q} + (n - 1) \cdot \frac{1 - q}{n - 1} \cdot \frac{-\ln(1 - x/100)}{r - a \cdot \frac{1 - q}{n - 1}}.$$
(6)

Observe that Eq. (6) is exactly the same as Eq. (2) when 456 $r_1 = r_w$ and $r_i = r$ for $i \neq 1$, except that p_1 is replaced by q. 457 Thus, the theoretically optimal solution, via Eq. (3), is $q^* = p_1^*$, 458 and thus the theoretically optimal split for the worker under 459 interference is $q^* \cdot \frac{n}{c} = p_1^* \cdot \frac{n}{c}$. Intuitively, this result says that 460 more load needs to be placed on the interference worker in 461 case of OLAP when compared to web applications; this makes 462 sense as there are fewer alternative workers in case of OLAP 463 applications as opposed to web applications, i.e., (c-1) as 464 opposed to (n-1). We can similarly obtain the theoretically 465 optimal split when more than one worker is under interference.

2.5 The DIAL Control Flow

467

491

503

The control flow for our DIAL implementation (for web and 468 OLAP applications) is as follows: 469

- 0) Monitoring: DIAL monitors the fg application's T_{x} , 470 arrival rate, *a*, load, $\rho_{fg,i}$, and classification metrics 471 (e.g., connection time), averaged every interval, for 472 all fg VMs. 473
- 1) Detection: DIAL signals interference if T_x exceeds its 474 95 percent confidence bounds for successive moni- 475 toring intervals. 476
- Classification: DIAL next employs the decision tree 477 to identify the dominant resource under contention. 478
- 3) Estimation: DIAL then uses the T_x and $\rho_{fg,i}$ values in 479 Eq. (1), with the dominant resource-specific coeffi-480 cients, to estimate the interference, $\rho_{bg,i}$. The interfer-481 ence-aware throughput for fg VM *i* is adjusted by 482 $(1 - \rho_{bg,i})$.
- 4) Interference-aware load balancing: Given these estimates, and the monitored *a* value, DIAL derives the 485 LB weights, $\vec{p^*}$, via Eq. (3), and inputs them to the LB. 486

We continue monitoring the VMs' performance to detect 487 further changes in interference and to detect the end of 488 interference. When T_x returns to normal (for successive 489 intervals), we reset the LB weights. 490

2.5.1 Training the DIAL Controller

DIAL requires some model training to build the decision tree 492 (Section 2.3.2) and derive the coefficients of the estimation 493 model (Eq. (1)). The above training tasks can be performed 494 offline on a dedicated server in a cloud environment by con-495 trolling the bg VMs to run microbenchmarks at different 496 intensities while monitoring relevant metrics. In a private 497 cloud environment, such as OpenStack, we can set aside a 498 dedicated host using Availability Zones. In some public 499 clouds, such as Amazon, dedicated hosts can be rented. We 500 use these options for training the DIAL controller, as dis-501 cussed in Sections 4.3 and 5.4.

2.5.2 Assumptions for DIAL Controller Training

The above-described DIAL control flow and training makes 504 certain implicit assumptions about the incoming workload. 505 Specifically, by training at different load intensities, DIAL 506 assumes that (i) the workload request mix does not change significantly at runtime, and (ii) the distribution of inter-arrival 508 times does not change significantly at runtime. When the 509 request mix changes, for example, to a more database-heavy 510 request mix, then the workload will have a greater sensitivity 511 to disk I/O contention. This will thus require a retaining of the 512 DIAL controller to infer the new model parameters. Note that 513

the mean arrival rate and/or the number of workers may change dynamically, and this is already monitored by DIAL and is taken into account when determining the theoretically optimal load balancing weights via the *a* and *n* parameters, respectively. We show, in Section 5.4.3, that DIAL works well even under an abrupt request rate change and a change in the number of workers.

521 **3 EVALUATION METHODOLOGY**

To evaluate the efficacy of DIAL, we implement it for realistic applications and study the reduction in tail latency when worker nodes face interference. This section describes the fg and bg application setup we employ, and our resource monitoring approach for worker nodes.

527 3.1 Foreground (fg) Applications

528 3.1.1 Web Applications

We consider multi-tier web applications where the application server tier is treated as the worker tier. The incoming requests are distributed among application servers via a load balancer.

We employ two multi-tier web benchmarks as our fg application, CloudSuite [22] and WikiBench [23]. Unless specified otherwise, we use CloudSuite in foreground.

CloudSuite. The CloudSuite 2.0 Web Serving benchmark is
 a multi-tier, multi-request class, PHP-MySQL based social
 networking application. The benchmark uses several request
 classes, e.g., HomePage, TagSearch, EventDetail, etc.

540 Our CloudSuite setup consists of: (i) Faban workload generator for creating realistic session-based web requests. 541 We set the number of users to 1000 for OpenStack and 5000 542 for AWS; the think time is 5s (default). (ii) HAProxy LB dis-543 tributes incoming http requests (from Faban) among the 544 back-end application tier VMs. We use the default Round 545 Robin policy, unless stated otherwise. (iii) Application VMs 546 installed with Apache, PHP, Memcached, and an NFS-Cli-547 548 ent. We employ 3 application VMs in OpenStack and 10 in AWS. (iv) A MySQL server and an NFS server, hosting the 549 file store, are installed on separate, large VMs (to avoid 550 being the bottleneck). 551

WikiBench. WikiBench is a Web hosting benchmark that mimics wikipedia.org. Our WikiBench setup consists of: (i) wikijector load generator to replay real traffic from past traces of requests to Wikipedia, (ii) HAProxy LB, and (iii) three VMs running the MediaWiki application (the same application that hosts wikipedia.org), and (iv) a MySQL database to store the Wikipedia database dump.

559 3.1.2 OLAP Applications

We use the open-source Pinot [11] system as our represen-560 tative OLAP application. Pinot is a low-latency, scalable, 561 distributed OLAP data store that is used at LinkedIn and 562 Uber for various user-facing functions and internal analy-563 sis. The Pinot architecture consists of three main compo-564 565 nents: (1) controller, (2) broker, and (3) historical worker 566 nodes. The controller is responsible for cluster-wide coor-567 dination and segment (data shard) assignment to worker nodes. The broker(s) receives queries from clients, distrib-568 569 utes them among workers, and integrates the results from the workers and sends the final result back to clients. The 570 brokers act as our load balancing tier (LBT), see Section 571 2.1. The historical worker nodes host data segments and 572

respond to queries that originate from the broker. The 573 worker nodes constitute our worker tier. Historical work- 574 ers store data in the form of an index called *segment*; every 575 table has its own segments. 576

For all the above fg applications, we use the suggested 577 default configuration values, resulting in average CPU utili-578 zation of about 25 percent for CloudSuite, 34 percent for 579 WikiBench, and 62 percent for Pinot, without interference. 580 Recent studies, including those at Azure [24] and Alibaba 581 [25], reported average CPU utilizations of about 20 percent 582 for fg VMs. 583

3.2 Background (bg) Workloads

In our experiments, we emulate interference by employing 585 several bg workloads to create contention for the fg application. The bg workloads are hosted on VMs colocated with 587 the fg application layer VMs. Each fg VM under interference 588 is hosted separately from other fg VMs, and is colocated 589 with bg VMs. We first employ *microbenchmarks* to stress 590 individual resources for analyzing fg interference. We then 591 employ *test workloads* to evaluate DIAL for fg applications 592 under realistic cloud workloads. 593

584

606

Microbenchmarks. We employ: (i) stress-ng tool on bg 594 VMs to create controlled CPU contention; (ii) httperf load 595 generator (on a separate VM and PM) to retrieve hosted files 596 from the colocated bg VMs at different, controllable request 597 rates to create NET contention; (iii) dcopy benchmark on bg 598 VMs to create LLC contention; and (iv) stress on bg VMs to 599 create DISK contention. 600

Test Workloads. We employ: (i) SPEC CPU to create CPU 601 contention, (ii) Memcache server (driven by mutilate client) 602 to create NET contention, (iii) STREAM to create LLC contention, and (iv) Hadoop running TeraSort with a large data 604 set to create DISK contention. 605

3.3 Resource usage Monitoring

We study resource contention for four resources: (i) network 607 (NET), CPU, Last-Level-Cache (LLC), and disk. We now 608 explain how we monitor fg and bg resource usage, from 609 within the VMs, for our model training. 610

- NET: We use the dstat Linux tool to monitor the used 611 network bandwidth for bg and fg VMs. We then nor-612 malize their sum by the peak bandwidth. 613
- CPU: We consider fair-sharing of the possibly over- 614 committed PM cores among VMs. If a PM has n 615 cores available and all VMs together require m cores, 616 then the CPU usage of each VM is normalized by 617 max{m,n}.
- LLC: Since memory bandwidth for a VM cannot be 619
 easily monitored, we employ the RAMspeed bench- 620
 mark to measure the available memory bandwidth. 621
 We obtain this bandwidth for each experiment and 622
 then estimate the LLC usage by computing the dif- 623
 ference between peak bandwidth and experiment 624
 bandwidth. Finally, we normalize this difference by 625
 peak bandwidth to estimate LLC usage. 626
- DISK: Disk usage typically depends on the access 627 pattern (sequential versus random). We thus use the 628 same approach as for LLC, but with sysbench 629 instead of RAMspeed, for estimating DISK usage. 630



Fig. 4. Illustration of our OpenStack cloud setup.

631 4 DIAL FOR WEB APPLICATIONS

We first explain our DIAL implementation, and then present evaluation results for CloudSuite and WikiBench.

634 4.1 DIAL Implementation

For DIAL web application deployment, we implement the 635 DIAL controller logic using: (i) a C program to execute the 636 detection, classification, and estimation tasks, and (ii) a set 637 of bash scripts to monitor metrics from the /proc subsys-638 tem (from within the VM) and the LB logs, and to communi-639 cate with the LB to reconfigure the weights. The overhead 640 641 of the DIAL controller is negligible in practice since the decision tree building, response time modeling, and LB 642 weights optimization are performed offline, and are only 643 leveraged periodically during run time using the monitored 644 metrics. Our evaluation results show that the average 645 increase in CPU utilization of the LB VM under DIAL is 646 about 2 percent. ough interval. 647

648 4.2 Cloud Environments

We set up two cloud environments for our evaluation, an
OpenStack based private cloud environment and an AWSbased public cloud environment. Unless specified otherwise, we use the OpenStack environment.

653 OpenStack-based Private Cloud. Fig. 4 depicts our experi-654 mental setup. We use an OpenStack Icehouse-based private 655 cloud with several dedicated Dell C6100 physical machines, referred to as PMs. Each PM has 2 sockets with 6 cores each, 656 and 48GB memory. The host OS is Ubuntu 14.04. All PMs 657 are connected to a network switch via a 1Gb Ethernet cable. 658 659 Our experiments reveal that the maximum achievable network bandwidth is about 115 MB/sec (we flood the net-660 661 work using a simple load generator, httperf, and measure the peak observed bandwidth under various request rates 662 and request sizes). Likewise, we find that the maximum 663 achievable memory and (sequential) hard disk drive I/O 664 bandwidths are about 11GB/sec (using RAMspeed) and 665 50 MB/sec (using sysbench), respectively. 666

AWS-based Public Cloud. We rent 10 c4.large instances (2
vCPUs and 3.75GB of memory) in AWS EC2's US East (N.
Virginia) region. We also rent a c4 dedicated server (PM) for
hosting one of the instances colocated with bg VMs.

671 4.3 Evaluation

We first present results for classification and estimation of test workloads. We then present results for performance improvement (reduction in T_{90}) under DIAL for OpenStack and AWS setups for CloudSuite and WikiBench. Unless mentioned otherwise, we compare performance under DIAL with performance without DIAL, referred to as



Fig. 5. Observed and modeled response times for CloudSuite under resource contention via microbenchmarks. Average modeling error: 6.1 percent.

baseline. In Section 4.3.6 we compare DIAL against existing 678 interference-aware techniques that are popularly employed. 679

4.3.1 Evaluating Detection, Classification, and Estimation

Detection. The crosses in Fig. 5 show the impact of different 682 resource contentions, created by microbenchmarks, on 683 CloudSuite's HomePage request class response time under 684 the OpenStack setup. Every data point (cross) in Fig. 5 is 685 obtained by averaging the 90%ile of response times in every 686 monitoring interval over three different experiments, each of 687 which takes 300s. To detect contention, we use the 95 percent 688 confidence intervals around the mean (see Section 2.3.1) to 689 obtain the following detection rule for both the OpenStack 690 and AWS setups: $T_{90} > 5ms$, for HomePage; similar rules 691 can be derived for other request classes. We run several 692 experiments using the bg test workloads and find that our 693 detection rule results in a low false positive rate of 5.7 percent.

Prior work has employed similar techniques to detect 695 and analyze interference using hardware performance 696 counters such as CPI [27], MIPS [5], cache miss rate [1], [2], 697 etc.; such values are visible to the hypervisor, but are diffi-698 cult and often infeasible to obtain from within the VM. We 699 tried accessing such counters through VMs hosted by AWS 700 EC2, Google Cloud Platform, and our OpenStack environ-701 ment, but the values were either not supported or were 702 incorrectly reported as all zeros; similar observations were made for AWS EC2 VMs in prior work [2]. 704

Classification. We monitor the user space CPU utilization, 705 usr, the system space CPU utilization, sys, the I/O wait 706 time, wai, the rate of segments retransmitted, seg_ret, and 707 the 90%ile time taken to establish a connection to the 708



Fig. 6. Our trained decision tree. Leaves represent the contention classification with numbers in the leaves representing the total classification instances (left) and the number of misclassified ones, if any (right).

680



Fig. 7. Performance comparison between DIAL and baseline for test background workloads. The red, blue, green, and gray regions represent NET, CPU, DISK, and LLC contention, respectively. DIAL reduces 90% ile response time during these contentions by 39.1, 56.3, 16.2, and 59.2 percent.



Fig. 8. *T_c*, *wai*, *sys*, and *usr* metrics for apache1 application layer VM for the experiments in Fig. 7. apache1 VM experiences NET, DISK, and LLC contention, and shows an increase in relevant metrics under those contentions.



Fig. 9. T_c, wai, sys, and usr metrics for apache2 application layer VM for the experiments in Fig. 7. apache2 VM experiences only CPU contention, and consequently shows an increase in relevant metrics under CPU contention.

application VM, T_c (via HAProxy logs). Note that all metrics 709 are monitored from within the VMs, to comply with the 710 user-centric design of DIAL. We normalize usr and sys 711 using predicted values to distinguish from workload varia-712 tions, as discussed in Section 2.3.2. The usr and sys metrics 713 can help detect CPU and LLC contention as the processor 714 might have to do more work under these contentions. wai 715 716 could potentially help classify DISK contention. Finally, seg and T_c could help classify NET contention because of the 717 reduced available network bandwidth. 718

Our decision tree for CloudSuite, trained using microbe-719 nchmarks, is shown in Fig. 6. The decision tree is generated 720 using WEKA [27]; in particular, WEKA determines the 721 nodes and cutoff values using the J48 algorithm. The tree 722 structure may be different for different applications. How-723 ever, we expect the high level rules to be the same, as illus-724 trated by our classification results for the Pinot OLAP 725 application in Section 5.4.1. For example, we expect that 726 LLC interference will lead to an increase in CPU usage. 727

Our 10-fold cross-validation error is 7.8 percent. Our clas-728 729 sifier shows that high (normalized to predicted contention) 730 usr signals LLC contention, possibly because more work has to be done to service the LLC misses. A high T_c signals 731 NET contention, which seems intuitive. A moderate drop in 732 usr and moderate rise in sys signals CPU contention; we 733 734 believe this is because throughput decreases under contention, resulting in lower *usr*, and thus exhibiting a relative 735

rise in *sys*. A high *wai* suggests DISK contention. Finally, a 736 moderate rise in *seg_ret* and T_c signals workload variations 737 (denoted as Δ _load in Fig. 6). 738

We also evaluate our classifier using test workloads that 739 were *not* seen during classifier training. We run 50 total 740 experiments using 10 experiments each for Memcache (NET 741 contention), SPEC (CPU contention), Hadoop (DISK conten-742 tion) and STREAM (LLC contention), in addition to 10 743 experiments under varying CloudSuite application load. 744 Our decision tree successfully classifies 44 of the 50 test 745 instances; the misclassifications are observed for change in 746 workload and DISK contention. The "misclassifications" for 747 DISK contention (as LLC) under Hadoop are because of the 748 numerous memory accesses made by the colocated Slave 749 VMs; we believe that Hadoop interference cannot always be 750 classified as a single resource due to its complex and 751 dynamic resource needs.

Estimation. The solid lines in Fig. 5 show our modeling 753 results for CloudSuite interference estimation (see Section 754 2.3.3) under different resource contentions via training. Our 755 average modeling error across all contentions is 6.1 percent. 756 If we instead use $\alpha = 1$ in Eq. (1) to model T_{90} simply as 757 $c_0 + c_1/(1 - \rho_{fg} - \rho_{bg})$, the modeling error increases to about 758 15 percent. However, when we increase the value of α 759 beyond 2, we find only modest improvements in accuracy. 760

Effect of Monitoring Interval Length. We use a metrics mon- 761 itoring interval length of 10s for the above evaluation. 762



Fig. 10. DIAL reduces the response time of all request classes by 37 and 56 percent under CPU and combined CPU + NET contention, respectively.

Experimentally, we find that shorter interval lengths, such as 1s or 5s, lead to inaccurate classification and estimation due to system noise and load fluctuations. On the other hand, intervals larger than 10s do not significantly improve accuracy. For these reasons, we choose an interval length of 10s; prior work has also reported such reaction times to avoid rash decisions [2], [26], [28].

770 4.3.2 Evaluating DIAL under Real Workloads

Fig. 7 shows our experimental results for CloudSuite under
OpenStack for various time-varying contentions created
using test workloads in bg VMs. The y-axis shows the tail
latency for CloudSuite across all request classes. We create
NET, DISK, and LLC contention for apache1 VM using
Memcache, Hadoop (TeraSort), and STREAM, respectively.
We use SPEC to create CPU contention for apache2.

We see that DIAL significantly reduces tail response times,
when compared to the baseline, under all contentions; the
reduction ranges from 16 percent under DISK contention to
59 percent under LLC contention. The relatively low improvement under DISK contention is because Hadoop intermittently utilizes disk I/O bandwidth; further, not all CloudSuite
request classes require (or contend for) disk access.

Without DIAL, the tail response time can be as high as 20-30ms; with DIAL, the tail response time is almost always around 4-5ms. Note that DIAL requires some time (at least two successive intervals of high response time) for interference detection during which response time continues to be high, as seen at the start of each contention.

791 Figs. 8 and 9 show our classification metrics for apache1 and apache2, respectively; we only show T_c , wai, sys, and 792 usr (and not seg_ret) for ease of presentation. Note that the 793 y-axis range in Fig. 9 is intentionally smaller to focus on the 794 rise in the sys metric. For apache1, under NET contention, 795 T_c is high while the other metrics are unaffected. For DISK 796 and LLC contentions, sys is high, especially for LLC; fur-797 ther, usr is also high under LLC contention. Finally, the wai 798 metric, though noisy, is higher under DISK contention. By 799 contrast, these metrics are unaffected for the corresponding 800 time periods under apache2. 801

Likewise, for apache2, for CPU contention, sys is moder-802 803 ately high but not as high as that under DISK and LLC con-804 tention under apache1. Again, the metrics are unaffected for 805 the CPU contention period under apache1. This shows that the relevant metrics on the compromised VM change under 806 contention, but are unaffected for uncompromised VMs. Fur-807 ther, the change in metric values under the contention peri-808 ods are in agreement with the rules of the decision tree 809 classifier in Fig. 6, even though the classifier was trained on 810



Fig. 11. Performance under LLC contention for fg WikiBench. DIAL reduces response times by ${\sim}23.6$ percent during contention (gray regions).

microbenchmarks and not on these test workloads. This 811 highlights the efficacy of our classifier. 812

For Memcache, the server is hosted on a bg VM and is 813 driven by mutilate clients (on different hosts) issuing a high 814 request rate for a small set of key-value pairs, resulting in 815 NET being the dominant resource. DIAL correctly classifies 816 this Memcache bg VM as creating NET contention. For 817 Hadoop, there is significant demand for disk and memory 818 bandwidth; our classifier suggests DISK contention. 819

4.3.3 Evaluating DIAL under Multiple Contentions

DIAL is capable of dynamically responding to multiple 821 compromised VMs. The optimization in Section 2.4.1 pro- 822 vides estimates for LB weights, via Eq. (3), for all VMs. This 823 is different from the case of multiple resource contentions 824 on the *same* VM, which is beyond the scope of this paper. 825

Fig. 10 shows our experimental results for CloudSuite 826 where initially apache2 VM is under CPU contention, but 827 then, after about 5 mins, apache1 (on a different host) also 828 starts experiencing NET contention, resulting in very high 829 interference for the application. After an additional 5 mins, 830 both contentions are terminated. We see that DIAL substan- 831 tially reduces T_{90} under interference. This example highlights 832 the dynamic nature of DIAL. Compared to existing techniques 833 that employ (static) VM placement to mitigate interference, 834 DIAL is able to adapt to variations in interference by constantly 835 updating its estimates and re-distributing load accordingly. 836 For the above experiment, for CPU contention, the DIAL 837 weights are $\{0.45, 0.1, 0.45\}$ (apache2 under contention), and 838 for combined CPU and NET contention, the weights are 839 $\{0, 0.27, 0.73\}$ (apache1 under severe NET contention). 840

4.3.4 Evaluating DIAL for the WikiBench fg Application 841 Fig. 11 shows our results for WikiBench under LLC conten- 842

tion created by the dcopy microbenchmark. Here, we have 843 two application VMs and one of them is under contention. 844 The figure shows the response time for baseline and DIAL 845 for all request classes. We create three different contention 846 levels for this experiment, shown in gray. DIAL reduces 847 response time by about 23 percent when compared to the 848 baseline. We also measure the *usr* and *sys* metrics for classification and find that both increase considerably, by about 850 62 and 41 percent, respectively, under interference; this is in 851 agreement with our decision tree classifier. 852

4.3.5 Evaluating DIAL in the AWS Setup

Fig. 12 shows our results for CloudSuite under LLC conten- 854 tion created by the dcopy microbenchmark in the AWS 855 setup. Here, we have 10 application VMs and only one of 856 them is under contention. The figure shows the response 857 time for baseline and DIAL for all request classes served by 858

820



Fig. 12. Performance under LLC contention for AWS setup. DIAL reduces response times by around 22.3 percent.



Fig. 13. *usr* and *sys* metrics for the gray region in Fig. 12. These metrics clearly increase during contention.

all VMs in the AWS setup. We create several different contention levels for this experiment. We see that DIAL reduces
response time by about 22 percent when compared to the
baseline. This shows that *even one* compromised VM (out of
10) can considerably impact the overall response time.

Fig. 13 shows the *usr* and *sys* metrics for the shaded region in Fig. 12 to assess classification. Clearly, both the *usr* and *sys* metrics increase considerably during contention when compared to the low, flat lines during no contention. Further, the regions of contention can be easily discerned from the figure, resulting in good detection accuracy.

4.3.6 Comparison with Existing user-Centric Techniques

Utilization-based Strategies.Fig. 14 shows our experimental 872 results for high CPU contention under DIAL and under 873 ICE [1]. Similar to DIAL, ICE is an interference-aware load 874 875 balancer that adjusts the traffic directed towards compromised VMs. However, instead of using LB weights, ICE 876 ensures that the CPU utilization for the compromised VMs 877 878 stays below a certain threshold. The authors do not mention this threshold value in the paper, and so we experi-879 mentally determine the best threshold value across 880 experiments. Unfortunately, we find that the optimal CPU 881 utilization threshold varies with the amount and type of 882 interference. For example, we find that 15 percent CPU 883 utilization works well for moderate CPU interference 884 under ICE, but does not work well for high CPU interfer-885 ence, as shown in Fig. 14. Under DIAL, response time is 886 significantly lower, and the observed CPU usage at the 887 compromised VM is about 8-10 percent; results are similar 888 889 for other contentions.

Queue-Length based Strategies. Queue-length or load-based
 strategies send traffic to the VM that has the lowest load. We
 consider the Least Connections (LC) strategy that directs the
 next incoming request to the VM that has the least number
 of active connections. Under interference, the outstanding
 requests for the compromised VM will be higher, resulting
 in fewer additional requests being sent to it under LC.



Fig. 14. Comparison of DIAL with ICE under CPU contention. DIAL $_{905}$ reduces T_{30} by 25-48 percent for all request classes.



(a) TagSearch for LLC contention. (b) TagSearch for CPU contention.

Fig. 15. Comparison of DIAL with other LB heuristics.

Fig. 15 shows the reduction in T_{90} afforded by DIAL over 905 LC (and other heuristics that we discuss next) for the 908 TagSearch request class under CPU and LLC conten-909 tions; results are similar for other classes and for NET 910 and DISK contention. We see that DIAL lowers response 911 time significantly, by as much as 70-80 percent, when 912 compared to LC (red dashed line). The improvement is 913 greater at higher contentions. The reason for this 914 improvement is that the compromised VM does not just 915 have lower capacity, but also requires (non-linearly) *more* 916 *time* to serve each request. The weights under DIAL take 917 both these into consideration, as opposed to LC that 918 only addresses the former.

Weighted Load Balancing Strategies. We now compare 920 DIAL with other weighted load balancing heuristics, such 921 as Weighted Round Robin (WRR) and Weighted Least Con-922 nections (WLC). For WRR and WLC, we use proportional 923 interference-aware weights, as discussed in Section 2.4. 924 Fig. 15 shows the reduction in T_{90} afforded by DIAL over 925 WRR (blue solid line) and WLC (black dotted line). We see 926 that DIAL lowers response time considerably when com-927 pared to these heuristics. It is interesting to note that WRR 928 is typically worse than WLC under CPU contention, but bet-929 ter than WLC under LLC contention; this observation reaffirms the fact that the impact of interference depends on the 931 type of resource under contention. 932

5 DIAL FOR PINOT

We now present our implementation of DIAL and its evaluation for a widely used OLAP solution, Pinot [11]. 935

5.1 DIAL Implementation

The Load Balancing Tier (LBT) for Pinot consists of the Broker nodes (see Section 3.1.2), that distribute queries to the back-end workers nodes. The Brokers rely on routing tables, stored in Broker memory, to determine which worker nodes host the data segments that are needed to serve the incoming query. Each routing table is a map from every segment 942

936



Fig. 16. Interference modeling and classification for Pinot.

to one worker node; since each segment is stored on several
replicas, numerous unique routing tables can be generated.
By randomly selecting a routing table for each query, the
Brokers balance load among the worker nodes.

We implement DIAL on the Broker using \sim 300 lines of Java code. Once interference is detected, DIAL updates the routing tables to remap segments that were initially assigned to the worker(s) under interference to other replicas, based on the theoretically-derived optimal fractions, q^* (see Section 2.4.2). Likewise, once interference ceases, the routing tables are updated to the default balanced weights.

954 5.2 Cloud Environment

We use several blade servers (PMs) from a HP Proliant 955 956 C7000 Chassis. Each PM has 2 sockets with 4-core CPUs 957 each, and 32 GB memory. The host OS is Ubuntu 16.04. The servers are connected through 1Gb/s network links. We use 958 KVM (on top of Ubuntu 16.04) to deploy VMs on these 959 PMs. We deploy 6 Pinot worker nodes on 1 vCPU, 16GB 960 memory VMs; each VM is on a separate PM. We deploy the 961 Pinot Controller and 2 Pinot Brokers using VMs with 8 962 vCPUs and 16GB memory, on different PMs. 963

We experiment with CPU and LLC contention for Pinot. For CPU contention, we use a 1 vCPU bg VM that is statically pinned to the same core as the fg VM (via hyperthreading). For LLC contention, we use a 3 vCPU bg VM that is pinned to the remaining 3 cores of the 4-core socket that hosts the fg VM; in this way, we do not share the same core as the fg VM to avoid CPU contention.

971 5.3 Workload and Benchmark

We implement a query generator for Pinot based on our 972 tables. For each table, we create several realistic queries. An 973 example query for the ProfileView table is "SELECT COUNT 974 (*) FROM ProfileView WHERE ViewedProfileID=ID", where 975 *ID* is a (randomized) query parameter. Our table and query 976 design is based on LinkedIn's Pinot deployment [11]. Our 977 978 benchmark is implemented in ~2000 lines of code, and is open-sourced [30]. 979

980 5.4 Evaluation

We use a warm-up time of 120s for all our experiments in this section. We focus on 95%ile response times for Pinot.

983 5.4.1 Evaluating Detection, Classification, 984 and Estimation

Detection. The crosses in Fig. 16a show the impact of CPU and LLC contention on Pinot response times. The detection rule of $T_{95} > 61ms$ is obtained based on the discussion in Section 2.3.1. We run several experiments using the bg test workloads and find that our detection rule results in a low false positive rate of 3.3 percent.



Fig. 17. Comparison of DIAL with other heuristics for Pinot.

Classification. We monitor total CPU usage, *cpu*, and the 991 system space CPU utilization, *sys*; we normalize these val-992 ues using predicted values to distinguish from workload 993 variations, as discussed in Section 2.3.2. Our decision tree 994 for Pinot is shown in Fig. 16b. The classification rules in 995 Fig. 16b closely resemble those for the web application in 996 Fig. 6. Our 10-fold cross-validation error is 2.3 percent. 997

We also evaluate our classifier using test workloads 998 that were *not* seen during training. We run 10 experi- 999 ments each for SPEC (CPU contention) and STREAM 1000 (LLC contention), and 10 experiments under varying 1001 Pinot workload. Our decision tree classifier is able to 1002 accurately classify all instances, except one CPU interfer- 1003 ence instance which is misclassified as LLC interference. 1004 Our classification accuracy based on these 30 experiments 1005 is 96.7 percent.

Estimation. The solid lines in Fig. 16a show our modeling 1007 results for Pinot interference estimation (as discussed in 1008 Section 2.3.3) under different resource contentions via training. Our average modeling error is 11.5 percent. 1010

5.4.2 Evaluating DIAL for Pinot under Real bg Workloads

Fig. 17 shows our experimental results for Pinot under our 1013 KVM setup for CPU and LLC contentions created using test 1014 workloads SPEC and STREAM, respectively. Here, the 1015 request rate for Pinot is set to 200 queries/sec, which results 1016 in a CPU load of about 60 percent. We consider 6 worker 1017 nodes with a replication factor of 3. The contention is cre- 1018 ated in bg VMs on one of the six PMs hosting the Pinot 1019 worker VMs. We show the tail response time values for no 1020 contention, baseline (with contention), DIAL, using theoreti- 1021 cally optimal interference-aware weights from Section 2.4.2, 1022 and Weighted Round Robin (WRR), which uses propor- 1023 tional interference-aware weights, as discussed in Section 1024 2.4. The response time is the query completion time monitored at the Broker, and depends on the performance of all 1026 workers. 1027

We see that DIAL significantly improves tail response 1028 times when compared to baseline; the average reduction in 1029 95% ile response times for CPU and LLC contention is 40.5 1030 and 25.8 percent, respectively. Compared to WRR, DIAL provides an average reduction in 95% ile response times for CPU 1032 and LLC contention of 16.1 and 16.5 percent, respectively. 1033

5.4.3 Evaluating DIAL for Pinot under Dynamic Conditions

Fig. 18 shows the 95% ile response time (tail latency) for 1036 Pinot under DIAL and baseline for our dynamic workload 1037 experiment. Here, we start with a load of 200 queries/s 1038 (or, qps) and no interference; as before, we have 6 Pinot 1039 workers and a replication factor of 3. Then, in the next 1040

11

1011

1012

1034



Fig. 18. Performance of DIAL and baseline using Pinot under CPU contention and under dynamic workload and cloud conditions.

phase (yellow shaded region), one of the fg worker VMs 1041 1042 experiences CPU interference due to a colocated VM running SPEC. DIAL responds, after monitoring and detec-1043 tion, by setting the theoretically optimal load balancing 1044 weights for the 3 replicas of segments hosted on the under-1045 interference worker (see Section 5.1). For this experiment, 1046 the theoretically optimal weights in this phase are 1047 $\{0.16, 0.42, 0.42\}$, obtained via the analysis discussed in 1048 1049 Section 2.4.2. By setting these weights, the tail latency low-1050 ers from about 147ms under the baseline to 88ms (39.9 per-1051 cent improvement).

1052 In the next phase (green shaded region), Pinot experiences an increase in load to 300 qps, severely impacting 1053 tail latency. DIAL detects this load change via request 1054 rate monitoring (see Section 2.5), and updates the load 1055 balancing weights to $\{0.2, 0.4, 0.4\}$, thus lowering tail 1056 latency from 266ms under the baseline to 189ms (28.2 1057 percent improvement). Our theoretically derived weights 1058 from Section 2.4.2 already take request rate into account 1059 (via the *a* parameter), and thus the updated weights can 1060 be easily obtained. 1061

To handle the increased load, Pinot eventually scales-out 1062 1063 by adding 3 new workers and redistributing data segments 1064 across all workers. We assume that the scale-out and data 1065 segment mapping is handled by an external autoscaling entity (e.g., MLscale [30] or other similar works [31], [32]). 1066 1067 With the additional workers, the tail latency of Pinot decreases, as seen in the last phase (gray shaded region). 1068 DIAL again updates the weights for this new configuration, 1069 by updating the *n* parameter (that represents the number of 1070 workers), resulting in a further lowering of tail latency from 1071 about 123ms under the baseline to 87ms (28.8 percent 1072 improvement). 1073

We repeated the experiments for a total of 5 runs. The 1074 results were qualitatively similar to Fig. 18, with the aver-1075 age improvement in 95% ile response time across all runs 1076 afforded by DIAL in the three shaded phases being about 1077 1078 33.1, 29.8, and 30.7 percent. We also repeated the experi-1079 ment with LLC contention, and obtained qualitatively similar results with an average improvement of up to 1080 20 percent. 1081

1082 6 PRIOR WORK IN THE CONTEXT OF DIAL

Interference Detection. Recent work has emphasized the need for user-centric interference detection [1], [2], [33], [34]. IC² [2] employs decision trees using VM-level statistics to detect interference at the cache; this information is then used to tune the configuration of web servers in co-located environments. Casale et al. [33] focus on CPU interference and present a user-centric technique to detect contention 1089 by analyzing the CPU steal metric. CRE [35] makes use of 1090 collaborative filtering to detect interference in web services 1091 by monitoring response times. While we also monitor 1092 response time, we go beyond detection and also estimate 1093 the amount of interference. CPI² [26] employs statistical 1094 approaches to analyze an application's CPI metric to detect 1095 and mitigate processor interference between different jobs. 1096 While CPI² can be used in virtual environments, public 1097 cloud VMs (e.g., AWS) do not always expose performance 1098 counters. 1099

There have also been prior works on hypervisor-centric 1100 interference detection (e.g., ILA [3]). While effective, such 1101 techniques require hypervisor access for monitoring hostlevel metrics, making them infeasible for cloud users. 1103

Interference-Aware Performance Management. ICE [1] pro- 1104 poses interference-aware load balancing by limiting the 1105 CPU utilization of the affected VM below a certain thresh- 1106 old. While effective, we find, via experiments (see Section 1107 4.3.6), that this strategy is not adaptive to different levels of 1108 interference. Mukherjee et al. [34] propose a tenant-centric 1109 interference estimation technique that employs a software 1110 probe periodically on all tenant VMs, and compares the per- 1111 formance of the probe at runtime versus that in isolation to 1112 quantify interference. The authors later extended this work 1113 to PRIMA [37], which is an interference-aware load balanc- 1114 ing and auto-scaling technique that leverages the above- 1115 described probing technique to make load balancing deci- 1116 sions. However, PRIMA only focuses on mean response 1117 time (as opposed to the more practical tail response time 1118 metric) and limits itself to network interference. Bubble- 1119 Up [37], Tarcil [7], Quasar [8], and ESP [38] profile workload 1120 classes and carefully colocate workloads that do not signifi- 1121 cantly impact each others' performance due to their specific 1122 resource requirements. By contrast, DIAL does not control 1123 colocation (VM placement is not in the user's control); 1124 instead, DIAL globally adjusts the fg LB policy to reroute 1125 some of the requests directed at affected VMs. 1126

7 CONCLUSION

We presented DIAL, a user-centric dynamic Interference-Aware Load Balancing framework that can be employed 1129 directly by cloud users without requiring any assistance 1130 from the hypervisor or cloud provider to reduce tail 1131 response times during interference. DIAL works by leverag-1132 ing two important components: (i) An accurate user-centric, 1133 response time-monitoring based interference detector, clas-1134 sifier, and estimator, and (ii) A framework for deriving theo-1135 retically optimal load balancer weights under interference. 1136 Our experimental results for web and OLAP applications 1137 on several cloud platforms, under interference from realistic 1138 benchmarks, demonstrate the benefits of DIAL.

ACKNOWLEDGMENTS

We thank the anonymous reviewers for their helpful sug- 1141 gestions. This work was supported by NSF grants CNS- 1142 1750109, CNS-1717588, and CNS-1617046. 1143

REFERENCES

 [1] A. Maji, et al., "ICE: An integrated configuration engine for interference mitigation in cloud services," in *Proc. IEEE Int. Conf. Auto-* 1146 *nomic Comput.*, 2015, pp. 91–100.

1140

1144

- 1148 A. Maji, et al., "Mitigating interference in cloud services by mid-[2] dleware reconfiguration," in Proc. 15th Int. Middleware Conf., 2014, 1149 pp. 277–288 1150
- X. Bu, et al., "Interference and locality-aware task scheduling for 1151 [3] MapReduce applications in virtual clusters," in Proc. 22nd Int. 1152 Symp. High-Perform. Parallel Distrib. Comput., 2013, pp. 227–238. 1153
- 1154 R. Nathuji, et al., "Q-clouds: Managing performance interference [4] effects for QoS-aware clouds," in Proc. 5th Eur. Conf. Comput. Syst., 1155 2010, pp. 237-250. 1156
- 1157 [5] C. Delimitrou and C. Kozyrakis, "QoS-aware scheduling in het-1158 erogeneous datacenters with paragon," ACM Trans. Comput. Syst., 1159 vol. 31, no. 4, 2013, Art. no. 12.
- S. Nathan, et al., "Towards a comprehensive performance model 1160 [6] of virtual machine live migration," in Proc. 6th ACM Symp. Cloud 1161 1162 Comput., 2015, pp. 288-301.
- C. Delimitrou, et al., "Tarcil: High quality and low latency sched-[7] 1163 1164 uling in large, shared clusters," in Proc. 6th ACM Symp. Cloud Comput., 2015, pp. 97-110. 1165
- 1166 [8] C. Delimitrou, et al., "Quasar: Resource-efficient and QoS-aware cluster management," in Proc. 19th Int. Conf. Architectural Support 1167 1168 Programm. Languages Operating Syst., 2014, pp. 127-144.
- [9] A. Gandhi, et al., "The unobservability problem in clouds," in 1169 1170 Proc. Int. Conf. Cloud Autonomic Comput., 2015, pp. 13-20.
- D. Novakovic, et al., "Deepdive: Transparently identifying and [10] 1171 1172 managing performance interference in virtualized environments," in Proc. USENIX Conf. Annu. Tech. Conf., 2013, pp. 219-230. 1173
- 1174 J.-F. Im, et al., "Pinot: Realtime OLAP for 530 million users," in Proc. Int. Conf. Manage. Data, 2018, pp. 583-594. 1175
- [12] M. Harchol-Balter, Performance Modeling and Design of Computer 1176 Systems: Queueing Theory in Action. Cambridge, U.K.: Cambridge 1177 1178 Univ. Press
- 1179 [13] K. Leonard, Queueing Systems, vol. 2, Hoboken, NJ, USA: Wiley, 1976 1180
- [14] "HAProxy: The reliable, high performance TCP/HTTP load bal-1181 1182 ancer," 2019. [Online]. Available: http://www.haproxy.org
- [15] A. Javadi and A. Gandhi, "DIAL: Reducing tail latencies for cloud 1183 1184 applications via dynamic interference-aware load balancing," in 1185 Proc. IEEE Int. Conf. Autonomic Comput., 2017, pp. 135-144.
- C. Wang, et al., "Effective capacity modulation as an explicit con-1186 [16] trol knob for public cloud profitability," in Proc. IEEE Int. Conf. 1187 1188 Autonomic Comput., 2016, pp. 95–104.
- 1189 Y. Xu, et al., "Small is better: Avoiding latency traps in virtualized [17] data centers," in Proc. 4th Annu. Symp. Cloud Comput., 2013, 1190 рр. 7–16. 1191
- Q. Zhang, et al., "A regression-based analytic model for capacity 1192 [18] 1193 planning of multi-tier applications," Cluster Comput., vol. 11, no. 3, 1194 pp. 197-211, 2008.
 - T. Horvath, et al., "Multi-mode energy management for multi-tier [19] server clusters," in Proc. Int. Conf. Parallel Architectures Compilation Techn., 2008, pp. 270-279.

1195

1196

1197

1211 1212

1213

1214

- [20] M. I. Reiman and B. Simon, "An interpolation approximation for 1198 queueing systems with poisson input," Operations Res., vol. 36, no. 1199 1200 3, pp. 454-469, 1988.
- 1201 [21] M. Boon, E. Winands, I. Adan, and A. van Wijk, "Closed-form waiting time approximations for polling systems," Perform. Eval., 1202 vol. 68, no. 3, pp. 290-306, 2011. 1203
- [22] M. Ferdman, et al., "Clearing the clouds: A study of emerging 1204 scale-out workloads on modern hardware," in Proc. 17th Int. Conf. 1205 1206 Architectural Support Programm. Languages Operating Syst., 2012, 1207 pp. 37-48.
- [23] E.J. Van Baaren, "WikiBench: A distributed, Wikipedia based web application benchmark," Master's thesis, Department of Computer 1208 1209 1210 Science, Vrije Univesiteit Amsterdam, the Netherlands, 2009.
 - [24] E. Cortez, et al., "Resource central: Understanding and predicting workloads for improved resource management in large cloud platforms," in Proc. 26th Symp. Operating Syst. Principles, 2017, pp. 153-167.
- [25] "Cluster data collected from production clusters in Alibaba for 1215 cluster management research," 2018, [Online]. Available: https:// 1216 github.com/alibaba/clusterdata 1217
- X. Zhang, et al., "CPI²: CPU performance isolation for shared [26] 1218 compute clusters," in Proc. 8th ACM Eur. Conf. Comput. Syst., 2013, 1219 pp. 379-391. 1220
- 1221 M. Hall, et al., "The WEKA data mining software: An update," [27] SIGKDD Explorations Newsletter, vol. 11, no. 1, pp. 10–18, 2009. 1222

- [28] David Lo, et al., "Heracles: Improving resource efficiency at 1223 scale," ACM SIGARCH Comput. Archit. News, vol. 43, no. 3, 2015, 1224 pp. 450–462. "PACELab/pinot," 2019. [Online]. Available: https://github.com/ 1225
- [29] 1226 PACELab/pinot 1227
- [30] M. Wajahat, A. Karve, A. Kochut, and A. Gandhi, "Using machine 1228 learning for black-box autoscaling," in Proc. 7th Int. Green Sustain-1229 able Comput. Conf., 2016, pp. 1-8. 1230
- [31] A. Gandhi, T. Zhu, M. Harchol-Balter, and M. Kozuch, 1231 "SOFTScale: Scaling opportunistically for transient scaling," in 1232 Proc. 13th Int. Middleware Conf., 2012, pp. 142-163. 1233
- [32] A. Gandhi, P. Dube, A. Kochut, L. Zhang, and S. Thota, 1234 "Autoscaling for hadoop clusters," in Proc. IEEE Int. Conf. Cloud 1235 *Eng.*, 2016, pp. 109–118. 1236
- [33] G. Casale, et al., "A feasibility study of host-level contention 1237 detection by guest virtual machines," in Proc. IEEE Int. Conf. Cloud 1238 Comput. Technol. Sci., 2013, pp. 152-157. 1239
- [34] J. Mukherjee, et al., "Subscriber-driven interference detection 1240 for cloud-based web services," IEEE Trans. Netw. Service Manag., 1241 vol. 14, no. 1, pp. 48-62, Mar. 2017. 1242
- [35] Y. Amannejad, et al., "Detecting performance interference in 1243 cloud-based web services," in Proc. IFIP/IEEE Int. Symp. Integrated 1244 Netw. Manage., 2015, pp. 423–431. 1245
- J. Mukherjee and D. Krishnamurthy, "Subscriber-driven cloud [36] 1246 interference mitigation for network services," in Proc. IEEE/ACM 1247 Int. Symp. Quality Service, 2018, pp. 1-6. 1248
- J. Mars, et al., "Bubble-Up: Increasing utilization in modern ware-1249 [37] house scale computers via sensible co-locations," in Proc. 44th 1250 Annu. IEEE/ACM Int. Symp. Microarchitecture, 2011, pp. 248–259. 1251
- N. Mishra, et al., "ESP: A machine learning approach to predict-[38] 1252 ing application interference," in Proc. IEEE Int. Conf. Autonomic 1253 Comput., 2017, pp. 125-134. 1254



Seyyed Ahmad Javadi received the PhD degree 1255 in computer science from Stony Brook University, 1256 in 2019. He is a researcher with the Department 1257 of Computer Science and Technology, University 1258 of Cambridge. He is currenly interested in to do 1259 research on performance challenges and privacy 1260 concerns in cloud computing environments. 1261



Anshul Gandhi received the PhD degree in 1262 computer science from Carnegie Mellon Univer- 1263 sity, in 2013. He is an assistant professor of Com- 1264 puter Science at Stony Brook University. His 1265 current research interests are in Performance 1266 Modeling and Cloud Computing. 1267

▷ For more information on this or any other computing topic, 1268 please visit our Digital Library at www.computer.org/publications/dlib. 1269