# Demographic Bias in Data Center Scheduling Systems

Sara Mahdizadeh Shahri, smahdiz@cmu.edu

## I. INTRODUCTION

Modern web services such as web search, video streaming, and online healthcare run across hundreds of thousands of servers in a data center to prioritize sending quick responses to the end user [22], [39]. Since web services are user-facing, they must often meet stringent tail ($99^{th}\%$) latency constraints expressed as soft real-time deadlines called Service Level Objectives (SLOs) [41].

To send quick responses and meet SLOs, web systems adopt a "performance-first" approach [13], [31], [40] by often introducing request priorities [47], [51]. For example, scheduling systems that implement the Shortest Job First policy [25] prioritize processing shorter requests over longer ones. Such prioritization-driven web systems improve performance metrics such as response latency [41], resource utilization [13], and the user's Quality of Experience (QoE) [31].

We posit that web systems that adopt a "performance-first" approach via request prioritization can potentially implicitly prioritize requests from one user demographic over another, and may thereby introduce *demographic bias*. Our hypothesis is motivated by prior works such as Zhang et al. [51], who propose a system that prioritizes faster processing of some requests that encounter lower network delays over requests that face greater network delays. In a similar vein, prior work proposes a scheduler that prioritizes requests from users who are closer to the server [33]. Hence, we posit that in such cases, a system might implicitly prioritize requests that originate from urban areas over rural ones, as urban areas typically face lower network delays, thereby causing demographic bias based on the user's geographical location.

To validate our hypothesis, we must first qualitatively and quantitatively define demographic bias. Using this definition, we must investigate whether existing prioritization-driven systems can cause demographic bias. If such systems indeed cause bias, we must develop solutions to control such bias. Thus, in this work and our future extension we systematically answer three questions: (1) How do we quantify a system's demographic bias in a way that is similar to today's norm of quantifying its performance, power, or fairness [19] metrics? (2) Can existing (open-source) systems, especially scheduling systems, introduce demographic bias? (3) How do we detect and control demographic bias?

Answering these three questions to define, identify, and control bias is important and challenging due to several reasons. First, while there are metrics to measure fairness [14], [44], [45], [50], such metrics do not reveal whether a user is discriminated against based on their demographic. Thus, there is a lack of a clear metric to quantify a web system's demographic bias. Second, if we do not identify and control bias, we can precipitate poor latencies for users for certain demographics.

We first define and quantify demographic bias for web systems by extending the Earth Mover's Distance (EMD) that is commonly used to quantify bias in ML systems [11]. Next, via a case study, we demonstrate that an open-source scheduling system is susceptible to demographic bias. We show that an existing open-source SJF-based scheduler [25] that prioritizes processing shorter requests implicitly introduces demographic bias in a web search service by prioritizing requests from male users over female users. This bias occurs since male users typically send shorter search queries that require shorter processing times (i.e., a shorter job size) [7], [23], [38], [49]. We limit our case study to open-source web systems since (1) most real-world, closed-source web systems are difficult to study outside industry and (2) web system operators are unlikely to be amenable to disclosing induced biases for fear of losing their user base [16], [21], [26], [46].

In our future work, we plan to develop a system called Bias-Controller, which features a framework to detect and control demographic bias in scheduling systems. Bias-Controller's goal is to achieve latency SLOs while minimizing bias.

## II. STUDYING DEMOGRAPHIC BIAS

### A. *How* should we study and define bias in web systems?

Given that services have access to demographic data, we investigate how demographic-driven bias is studied today. We consider the end-to-end pipeline of user data flow in an ML-driven data center system (e.g., scheduling system): (1) the user inputs their data, (2) data mining/collection systems collect this data, (3) ML algorithms are run using this data, (4) ML models predict using this data, (5) web system design decisions are made using these ML models, (6) application behaviors are driven by these systems decisions, and (7) the user receives a response from the application. From this end-to-end pipeline, we can clearly see that the introduction of bias in any stage will propagate to successive stages. Therefore, demographic biases caused by: user inputs (e.g., a user's implict bias in their data entry [12]); data collection algorithms [15], [18], [28]; or ML algorithms/models [5], [8], [48] can introduce biases in a web system's decisions.

We find that today, bias is primarily studied in data mining research [15], [18], [28] and ML algorithms/models research [4], [9], [17], [29], [32], [37], [42]

In comparison, such efforts to define and mitigate bias is woefully lacking in systems research, despite modern web systems being driven by these data collection and ML algorithms [52]. In short, while bias awareness has promoted more responsible data mining or ML algorithms research, systems research is yet to prioritize building bias-free web systems. Hence, we motivate how mitigating demographic bias is an endeavour that must be undertaken by not just data mining or ML algorithms research, but by low-level systems research as well. In short, everyone, all the way from the user inputting data parameters to ML algorithms/models, web systems, and applications, must work in coalition to ensure that there is no bias in the end-to-end data flow path, for modern web services to truly be bias-free. To this end, we make a case for systems research to treat demographic bias as a first-order concern.

Defining and mitigating demographic bias in web systems is challenging due to the lack of quantitative systems metrics to measure demographic bias. The closest, widely-used systems metric is fairness. Fairness is interpreted and quantified differently in different contexts, but mostly does not consider user demographics. For

example, a widely-used fairness concept is proportional sharing, where resource allocation is in proportion to the user priority [14], [44], [45], [50]. Similarly, prior work [19] proposes max-min fairness to guarantee that the network throughput for every user is at least as large as another user, when they both face the same bottleneck. In contrast to fairness, measuring demographic bias reveals whether a user is discriminated against (e.g., granted less resources) on the basis of their demographic. In short, while a system can be unfair, it might not suffer from demographic bias by discriminating against certain user demographics. Hence, we need a better metric to define and measure demographic bias.

*Defining bias.* To design bias-free data center systems, we must define demographic bias for each web system. Just like how a longer latency to a deprived demographic indicates a biased scheduling system, similarly, we must define what indicates a " demographic bias" for different data center systems (e.g., power management systems).

Additionally, intuitively we want the bias definition to holistically compare how users of different demographics experience a service. To define bias, similar to Dwork et. al. [11], we use the statistical distance. Assume $u$ and $v$ are the latency distributions of demographic A and demographic B respectively. We define bias as the Earth Mover's Distance (EMD) between $u$ and $v$ [36]. At high level, this distance shows the minimum amount of work required to transform one distribution into the other. If $U$ and $V$ are the respective CDFs of $u$ and $v$, then the minimum work to transform one distribution into the other is the area between the graphs of $U$ and $V$. As a result, we define the demographic bias as follows:

$$Demographic\ Bias(A, B) = (\int_{\Re} |F_A(x) - F_B(x)|^p\, dx)^{\frac{1}{p}} \quad (1)$$

**B. Why should we study and eliminate bias in web systems?**
Existing web systems can be susceptible to bias when they make decisions that prioritize performance [20], [43]. For example, prior work finds that web search result quality varies across diverse users, with certain results prioritized based on user profiles [47].

Web systems' susceptibility to bias is possibly even more likely today, since systems decisions today are often driven by ML, which is prone to being biased [27], [30]). Such ML algorithm-driven biases might percolate into web system and application behaviors. In this work, we will show that even without ML-driven scheduling, classical web system scheduling algorithms, such as "Shortest Job First," can inherently induce demographic bias.

Designing demographic bias-free data center systems is important since the success of several web service categories depends on minimal latency—an increase can have severe consequences. In the FinTech sector, millisecond-scale delays can cost profits on the order of millions of dollars [2], [35]. Similarly, online retailers stand to lose 7% in conversions for every 100 ms of delay; moreover, over 50% of visitors to a mobile retail site will leave if the page takes over three seconds to load [3].

In some cases, web applications can be life critical, intensifying the vitality of dependable latency bounds. Healthcare applications for Electronic Health Records can be a direct factor in patient care quality. Without high availability, incomplete health record data can lead to misdiagnoses or errors in prescribing treatments [6], [24]. More recently, social media has become essential to emergency response, both in mitigating primary deaths (resulting directly from the emergency) and secondary deaths (resulting from infrastructure breakdown) [34].
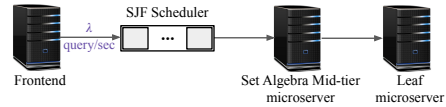


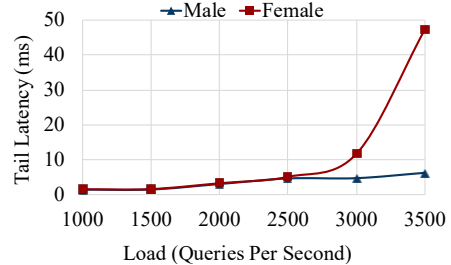Fig. 1: `Set Algebra` pipeline with integrated SJF scheduler.



Fig. 2: Throughput vs. tail latency curves for male and female demographics for `Set Algebra` service under SJF scheduling policy.

### III. BIASES IN SCHEDULING SYSTEMS

We study demographic-driven bias in one class of data center systems—scheduling systems. We investigate whether prioritization-driven scheduling algorithms can introduce bias to achieve performance gains. We conduct a case study to demonstrate demographic bias in the Shortest Job First (SJF)-driven scheduler that is integrated with a document search service [40]. We now detail the experimental setup and the case study's results.

**Experimental setup.** We evaluate our study on c6420 Xeon servers on CloudLab [10] using a web service from the $\mu$Suite benchmark suite [40]—`Set Algebra`, shown in Figure 1. `Set Algebra` performs document search by intersecting posting lists. It searches a corpus of 4.3 million WikiText documents in Wikipedia [1] sharded uniformly across leaf microservers, to identify documents containing all search terms. The mid-tier microserver forwards client queries containing search terms to the leaf microservers, which then return intersected posting lists to the mid-tier for their respective shards. The mid-tier aggregates the per-shard posting lists and returns their union to the client. Leaf servers look up posting lists for all search terms and then intersect the sorted lists. The resulting intersection is returned to the mid-tier. We query `Set Algebra` using search queries based on prior works that detail the type, distribution, and length of queries from male and female web search users [49].

**Experimental results.** Fig. 2 shows the tail ($99^{th}\%$) latency achieved across different load conditions measured in Queries Per Second (QPS) for both male and female demographics. This graph demonstrates that, although the SJF scheduler is not intentionally prioritizing a demographic over another, due to the job characteristics associated with each demographic, the scheduler can sustain a larger load for male users.

### IV. FUTURE WORK: BIAS-CONTROLLER

As part of our future work, we design a feedback-based controller that dynamically prioritizes incoming requests with the objective of meeting the application' QoS constraints while minimizing bias in the system. Here we describe the design of the Bias-Controller in detail.

Bias-Controller consists of a monitoring and priority component. The former monitors the application latency distribution of requests and consequently the SLO violation and the corresponding bias, while the latter uses the information in the monitoring component to determine appropriate priority assignment for the request, and enforces them using priority queues.

REFERENCES

[1] "Wikipedia, the free encyclopedia," https://en.wikipedia.org/w/index.php?title=Plagiarism&oldid=5139350.

[2] "Concept release on equity market structure," *Securities and Exchange Commission*, 2010. [Online]. Available: https://www.sec.gov/rules/concept/2010/34-61358.pdf

[3] "Milliseconds are critical," *Akamai online retail performance report*, 2017. [Online]. Available: https://www.ir.akamai.com/news-releases/news-release-details/akamai-online-retail-performance-report-milliseconds-are

[4] J. Angwin, J. Larson, S. Mattu, and L. Kirchner, *Machine bias: There's software used across the country to predict future criminals. and it's biased against blacks*, ProPublica, May 2016, https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing.

[5] S. Biswas and H. Rajan, "Do the machine learning models on a crowd sourced platform exhibit bias? an empirical study on model fairness," in *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ser. ESEC/FSE 2020. New York, NY, USA: Association for Computing Machinery, 2020, p. 642–653. [Online]. Available: https://doi.org/10.1145/3368089.3409704

[6] S. Bowman, "Impact of electronic health record systems on information integrity: Quality and safety implications," *Perspectives in Health Information Management*, 2013.

[7] D. Carmel, L. Lewin-Eytan, A. Libov, Y. Maarek, and A. Raviv, "The demographics of mail search and their application to query suggestion," in *International conference on world wide web*, 2017.

[8] L. Dixon, J. Li, J. Sorensen, N. Thain, and L. Vasserman, "Measuring and mitigating unintended bias in text classification," in *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, ser. AIES '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 67–73. [Online]. Available: https://doi.org/10.1145/3278721.3278729

[9] J. Dressel and H. Farid, "The accuracy, fairness, and limits of predicting recidivism," *Science advances*, vol. 4, no. 1, p. eaao5580, 2018.

[10] D. Duplyakin, R. Ricci, A. Maricq, G. Wong, J. Duerig, E. Eide, L. Stoller, M. Hibler, D. Johnson, K. Webb, A. Akella, K. Wang, G. Ricart, L. Landweber, C. Elliott, M. Zink, E. Cecchet, S. Kar, and P. Mishra, "The design and operation of cloudlab," in *2019 USENIX Annual Technical Conference (USENIX ATC 19)*. Renton, WA: USENIX Association, Jul. 2019, pp. 1–14.

[11] C. Dwork, M. Hardt, T. Pitassi, O. Reingold, and R. Zemel, "Fairness through awareness," in *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, ser. ITCS '12. New York, NY, USA: Association for Computing Machinery, 2012, p. 214–226. [Online]. Available: https://doi.org/10.1145/2090236.2090255

[12] A. Esmaieeli Sikaroudi, G. Rouzbeh, and A. Sikaroudi, "A data mining approach to employee turnover prediction (case study: Arak automotive parts manufacturing)," *Journal of Industrial and Systems Engineering*, vol. 8, 11 2015.

[13] Y. Gan, Y. Zhang, D. Cheng, A. Shetty, P. Rathi, N. Katarki, A. Bruno, J. Hu, B. Ritchken, B. Jackson, K. Hu, M. Pancholi, Y. He, B. Clancy, C. Colen, F. Wen, C. Leung, S. Wang, L. Zaruvinsky, M. Espinosa, R. Lin, Z. Liu, J. Padilla, and C. Delimitrou, "An open-source benchmark suite for microservices and their hardware-software implications for cloud & edge systems," in *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 3–18. [Online]. Available: https://doi.org/10.1145/3297858.3304013

[14] A. Gulati, A. Merchant, and P. J. Varman, "Pclock: An arrival curve based approach for qos guarantees in shared storage systems," in *Proceedings of the 2007 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems*, ser. SIGMETRICS '07. New York, NY, USA: Association for Computing Machinery, 2007, p. 13–24. [Online]. Available: https://doi.org/10.1145/1254882.1254885

[15] S. Hajian, F. Bonchi, and C. Castillo, "Algorithmic bias: From discrimination discovery to fairness-aware data mining," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ser. KDD '16. New York, NY, USA: Association for Computing Machinery, 2016, p. 2125–2126. [Online]. Available: https://doi.org/10.1145/2939672.2945386

[16] A. Heath, "Facebook's lost generation," https://www.theverge.com/22743744/facebook-teen-usage-decline-frances-haugen-leaks, 2021.

[17] A. Howard and J. Borenstein, "The ugly truth about ourselves and our robot creations: the problem of bias and social inequity," *Science and engineering ethics*, vol. 24, pp. 1521–1536, 2018.

[18] M. Z. Islam and L. Brankovic, "A framework for privacy preserving classification in data mining," in *Proceedings of the Second Workshop on Australasian Information Security, Data Mining and Web Intelligence, and Software Internationalisation - Volume 32*, ser. ACSW Frontiers '04. AUS: Australian Computer Society, Inc., 2004, p. 163–168.

[19] J. Jaffe, "Bottleneck flow control," *IEEE Transactions on Communications*, vol. 29, no. 7, pp. 954–962, 1981.

[20] T. Karr, "Digital denied: Free press report exposes the impact of systemic racism on internet adoption," https://www.freepress.net/news/press-releases/digital-denied-free-press-report-exposes-impact-systemic-racism-internet, 2016.

[21] J. Keeley, "4 reasons why facebook is starting to lose users," https://www.makeuseof.com/why-facebook-is-losing-users/, 2022.

[22] S. Kemp, "6 in 10 people around the world now use the internet," https://datareportal.com/reports/6-in-10-people-around-the-world-now-use-the-internet.

[23] L. Lorigo, B. Pan, H. Hembrooke, T. Joachims, L. Granka, and G. Gay, "The influence of task and gender on search and evaluation behavior using google," *Information processing & management*, 2006.

[24] S. Z. Lowry, M. Ramaiah, S. Taylor, E. S. Patterson, S. Spickard Prettyman, D. Simmons, D. Brick, P. Latkany, and M. C. Gibbons, "Technical evaluation, testing, and validation of the usability of electronic health records: Empirically based use cases for validating safety-enhanced usability and guidelines for standardization," 2015.

[25] H. Mao, M. Schwarzkopf, S. B. Venkatakrishnan, Z. Meng, and M. Alizadeh, "Learning Scheduling Algorithms for Data Processing Clusters," in *Proceedings of the ACM Special Interest Group on Data Communication*, 2019.

[26] A. Marantz, "The moral bankruptcy of facebook," https://www.newyorker.com/news/daily-comment/the-moral-bankruptcy-of-facebook, 2021.

[27] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan, "A survey on bias and fairness in machine learning," *ACM Comput. Surv.*, vol. 54, no. 6, jul 2021. [Online]. Available: https://doi.org/10.1145/3457607

[28] E. Ntoutsi, P. Fafalios, U. Gadiraju, V. Iosifidis, W. Nejdl, M. Vidal, S. Ruggieri, F. Turini, S. Papadopoulos, E. Krasanakis, I. Kompatsiaris, K. Kinder-Kurlanda, C. Wagner, F. Karimi, M. Fernandez, H. Alani, B. Berendt, T. Kruegel, C. Heinze, and S. Staab, "Bias in data-driven artificial intelligence systems—an introductory survey," *WIREs Data Mining and Knowledge Discovery*, vol. 10, 05 2020.

[29] C. O'Neil, *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*. USA: Crown Publishing Group, 2016.

[30] O. A. Osoba and W. Welser, *An Intelligence in Our Image: The Risks of Bias and Errors in Artificial Intelligence*. RAND Corporation, 2017.

[31] Z. peng GAO, J. CHEN, X. song QIU, and L. ming MENG, "Qoe/qos driven simulated annealing-based genetic algorithm for web services selection," *The Journal of China Universities of Posts and Telecommunications*, vol. 16, pp. 102–107, 2009. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S1005888508603477

[32] I. D. Raji and J. Buolamwini, "Actionable auditing: Investigating the impact of publicly naming biased performance results of commercial ai products," in *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, ser. AIES '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 429–435. [Online]. Available: https://doi.org/10.1145/3306618.3314244

[33] M. Rawat and A. Kshemkalyani, "Swift: Scheduling in web servers for fast response time," in *Proceedings of the Second IEEE International Symposium on Network Computing and Applications*, ser. NCA '03. USA: IEEE Computer Society, 2003, p. 51.

[34] P. Rengaraju, K. Sethuramalingam, and C.-H. Lung, "Providing internet access for post-disaster communications using balloon networks," in *Proceedings of the 18th ACM Symposium on Performance Evaluation of Wireless Ad Hoc, Sensor, & Ubiquitous Networks*, ser. PE-WASUN '21. New York, NY, USA: Association for Computing Machinery, 2021, p. 111–117. [Online]. Available: https://doi.org/10.1145/3479240.3488497

[35] V. Reporter, "The value of a millisecond: Finding the optimal speed of a trading," Apr 2008. [Online]. Available: https://research.tabbgroup.com/report/v06-007-value-millisecond-finding-optimal-speed-trading-infrastructure

[36] Y. Rubner, C. Tomasi, and L. J. Guibas, "The earth mover's distance as a metric for image retrieval," *International journal of computer vision*, vol. 40, no. 2, p. 99, 2000.

[37] T. Schnabel, A. Swaminathan, A. Singh, N. Chandak, and T. Joachims, "Recommendations as treatments: Debiasing learning and evaluation," in *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*, ser. ICML'16. New York, NY, USA: JMLR.org, 2016, p. 1670–1679.

[38] N. V. Spirin, J. He, M. Develin, K. G. Karahalios, and M. Boucher, "People search within an online social network: Large scale analysis of facebook graph search query logs," in *Proceedings of the 23rd acm international conference on conference on information and knowledge management*, 2014.

[39] A. Sriraman and A. Dhanotia, "Accelerometer: Understanding acceleration opportunities for data center overheads at hyperscale," in *Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS '20. New York, NY, USA: Association for Computing Machinery, 2020, p. 733–750. [Online]. Available: https://doi.org/10.1145/3373376.3378450

[40] A. Sriraman and T. F. Wenisch, "µ suite: a benchmark suite for microservices," in *2018 IEEE International Symposium on Workload Characterization (IISWC)*. Raleigh, NC, USA: IEEE, 2018, pp. 1–12.

[41] A. Sriraman and T. F. Wenisch, "µtune: Auto-tuned threading for oldi microservices," in *Proceedings of the 13th USENIX Conference on Operating Systems Design and Implementation*, ser. OSDI'18. USA: USENIX Association, 2018, p. 177–194.

[42] S. Tolan, M. Miron, E. Gómez, and C. Castillo, "Why machine learning may lead to unfairness: Evidence from risk assessment for juvenile justice in catalonia," in *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Law*, ser. ICAIL '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 83–92. [Online]. Available: https://doi.org/10.1145/3322640.3326705

[43] S. D. Turner, *Digital Denied: The Impact of Systemic Racial Discrimination on Home-Internet Adoption*, 2016. [Online]. Available: https://www.freepress.net/sites/default/files/legacy-policy/digital_denied_free_press_report_december_2016.pdf

[44] H. Wang and P. Varman, "Balancing fairness and efficiency in tiered storage systems with bottleneck-aware allocation," in *Proceedings of the 12th USENIX Conference on File and Storage Technologies*, ser. FAST'14. USA: USENIX Association, 2014, p. 229–242.

[45] Y. Wang and A. Merchant, "Proportional-share scheduling for distributed storage systems." in *FAST*, vol. 7. San Jose, CA, USA: USENIX, 2007, pp. 4–4.

[46] Wired, "What Really Happened When Google Ousted Timnit Gebru," https://www-wired-com.proxy.lib.umich.edu/story/google-timnit-gebru-ai-what-really-happened/.

[47] X. Xing, W. Meng, D. Doozan, N. Feamster, W. Lee, and A. C. Snoeren, "Exposing inconsistent web search results with bobble," in *Proceedings of the 15th International Conference on Passive and Active Measurement*, ser. PAM 2014. Berlin, Heidelberg: Springer-Verlag, 2014, p. 131–140. [Online]. Available: https://doi.org/10.1007/978-3-319-04918-2_13

[48] S. Xue, M. Yurochkin, and Y. Sun, "Auditing ml models for individual bias and unfairness," in *International Conference on Artificial Intelligence and Statistics*, 2020.

[49] E. Yom-Tov, "Demographic differences in search engine use with implications for cohort selection," *Information Retrieval Journal*, 2019.

[50] J. Zhang, A. Sivasubramaniam, Q. Wang, A. Riska, and E. Riedel, "Storage performance virtualization via throughput and latency control," *ACM Trans. Storage*, vol. 2, no. 3, p. 283–308, aug 2006. [Online]. Available: https://doi.org/10.1145/1168910.1168913

[51] X. Zhang, S. Sen, D. Kurniawan, H. Gunawi, and J. Jiang, "E2E: Embracing User Heterogeneity to Improve Quality of Experience on the Web," in *Proceedings of the ACM Special Interest Group on Data Communication*, 2019.

[52] Y. Zhang, W. Hua, Z. Zhou, G. E. Suh, and C. Delimitrou, "Sinan: Ml-based and qos-aware resource management for cloud microservices," in *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS '21. New York, NY, USA: Association for Computing Machinery, 2021, p. 167–181. [Online]. Available: https://doi.org/10.1145/3445814.3446693