

REFERENCES

- [1] I.J.B.F. Adan, G.J. van Houtum, and J. van der Wal. Upper and lower bounds for the waiting time in the symmetric shortest queue system. *Annals of Operations Research*, 48:197–217, 1994.
- [2] Kunal Agrawal, I-Ting Angelina Lee, Jing Li, Kefu Lu, and Benjamin Moseley. Practically efficient scheduler for minimizing average flow time of parallel jobs. In *2019 IEEE International Parallel and Distributed Processing Symposium, IPDPS 2019, Rio de Janeiro, Brazil, May 20-24, 2019*, pages 134–144. IEEE, 2019.
- [3] Kunal Agrawal, Jing Li, Kefu Lu, and Benjamin Moseley. Scheduling parallel DAG jobs online to minimize average flow time. In Robert Krauthgamer, editor, *Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2016, Arlington, VA, USA, January 10-12, 2016*, pages 176–189. SIAM, 2016.
- [4] S. Anand, Naveen Garg, and Amit Kumar. Resource augmentation for weighted flow-time explained by dual fitting. In *Proceedings of the Twenty-Third Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2012, Kyoto, Japan, January 17-19, 2012*, pages 1228–1241, 2012.
- [5] Spyros Angelopoulos, Giorgio Lucarelli, and Nguyen Kim Thang. Primal-dual and dual-fitting analysis of online scheduling algorithms for generalized flow-time problems. *Algorithmica*, 81(9):3391–3421, 2019.
- [6] Eitan Bachmat and Hagit Sarfati. Analysis of size interval task assignment policies. *Performance Evaluation Review*, 36(2):107–109, 2008.
- [7] Benjamin Berg, Jan-Pieter Dorsman, and Mor Harchol-Balder. Towards optimality in parallel scheduling. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 1(2):1–30, 2018.
- [8] Benjamin Berg, Mor Harchol-Balder, Benjamin Moseley, Weina Wang, and Justin Whitehouse. Optimal resource allocation for elastic and inelastic jobs. <https://arxiv.org/abs/2005.09745>.
- [9] Carl Bussema and Eric Torng. Greedy multiprocessor server scheduling. *Oper. Res. Lett.*, 34(4):451–458, 2006.
- [10] Jivitej S. Chadha, Naveen Garg, Amit Kumar, and V. N. Muralidhara. A competitive algorithm for minimizing weighted flow time on unrelated machines with speed augmentation. In Michael Mitzenmacher, editor, *Proceedings of the 41st Annual ACM Symposium on Theory of Computing, STOC 2009, Bethesda, MD, USA, May 31 - June 2, 2009*, pages 679–684. ACM, 2009.
- [11] Richard W Conway, Louis W Miller, and William L Maxwell. *Theory of scheduling*. Dover, 2003.
- [12] Jeffrey Dean and Sanjay Ghemawat. Mapreduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1):107–113, 2008.
- [13] Christina Delimitrou and Christos Kozyrakis. Quasar: resource-efficient and qos-aware cluster management. *ACM SIGPLAN Notices*, 49(4):127–144, 2014.
- [14] Jeff Edmonds. Scheduling in the dark. *Theor. Comput. Sci.*, 235(1):109–141, 2000.
- [15] Jeff Edmonds, Sungjin Im, and Benjamin Moseley. Online scalable scheduling for the k-norms of flow time without conservation of work. In Dana Randall, editor, *Proceedings of the Twenty-Second Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2011, San Francisco, California, USA, January 23-25, 2011*, pages 109–119. SIAM, 2011.
- [16] Jeff Edmonds and Kirk Pruhs. Scalably scheduling processes with arbitrary speedup curves. In *Proceedings of the twentieth annual ACM-SIAM symposium on Discrete algorithms*, pages 685–692. SIAM, 2009.
- [17] Kyle Fox and Benjamin Moseley. Online scheduling on identical machines using SRPT. In Dana Randall, editor, *Proceedings of the Twenty-Second Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2011, San Francisco, California, USA, January 23-25, 2011*, pages 120–128. SIAM, 2011.
- [18] Anshul Gandhi and Mor Harchol-Balder. How data center size impacts the effectiveness of dynamic power management. In *2011 49th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 1164–1169. IEEE, 2011.
- [19] Isaac Grosf, Ziv Scully, and Mor Harchol-Balder. Srpt for multiserver systems. *Performance Evaluation*, 127:154–175, 2018.
- [20] Abhishek Gupta, Bilge Acun, Osman Sarood, and Laxmikant V Kalé. Towards realizing the potential of malleable jobs. In *2014 21st International Conference on High Performance Computing (HiPC)*, pages 1–10. IEEE, 2014.
- [21] Anupam Gupta, Sungjin Im, Ravishankar Krishnaswamy, Benjamin Moseley, and Kirk Pruhs. Scheduling heterogeneous processors isn't as easy as you think. In Yuval Rabani, editor, *Proceedings of the Twenty-Third Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2012, Kyoto, Japan, January 17-19, 2012*, pages 1242–1253. SIAM, 2012.
- [22] Varun Gupta, Benjamin Moseley, Marc Uetz, and Qiaomin Xie. Stochastic online scheduling on unrelated machines. In *Integer Programming and Combinatorial Optimization - 19th International Conference, IPCO 2017, Waterloo, ON, Canada, June 26-28, 2017, Proceedings*, pages 228–240, 2017.
- [23] Varun Gupta, Karl Sigman, Mor Harchol-Balder, and Ward Whitt. Insensitivity for ps server farms with jsq routing. *ACM SIGMETRICS Performance Evaluation Review*, 35(2):24–26, 2007.
- [24] Mor Harchol-Balder. Task assignment with unknown duration. *Journal of the ACM*, 49(2):260–288, March 2002.
- [25] Mor Harchol-Balder. *Performance modeling and design of computer systems: queueing theory in action*. Cambridge University Press, 2013.
- [26] Mor Harchol-Balder, Cuihong Li, Takayuki Osogami, Alan Scheller-Wolf, and Mark Squillante. Cycle stealing under immediate dispatch task assignment. In *Proceedings of the 15th ACM Symposium on Parallel Algorithms and Architectures*, pages 274–285, San Diego, CA, June 2003.
- [27] Mor Harchol-Balder, Cuihong Li, Takayuki Osogami, Alan Scheller-Wolf, and Mark Squillante. Task assignment with cycle stealing under central queue. In *Proceedings of the 23rd International Conference on Distributed Computing Systems*, pages 628–637, Providence, RI, May 2003.
- [28] Mor Harchol-Balder, Takayuki Osogami, Alan Scheller-Wolf, and Adam Wierman. Multi-server queueing systems with multiple priority classes. *Queueing Systems: Theory and Applications*, 51(3-4):331–360, 2005.
- [29] Mor Harchol-Balder, Alan Scheller-Wolf, and Andrew Young. Surprising results on task assignment in server farms with high-variability workloads. In *ACM Sigmetrics 2009 Conference on Measurement and Modeling of Computer Systems*, pages 287–298, 2009.
- [30] Benjamin Hindman, Andy Konwinski, Matei Zaharia, Ali Ghodsi, Anthony D Joseph, Randy H Katz, Scott Shenker, and Ion Stoica. Mesos: A platform for fine-grained resource sharing in the data center. In *NSDI*, volume 11, pages 22–22, 2011.
- [31] Sungjin Im, Benjamin Moseley, Kirk Pruhs, and Eric Torng. Competitively scheduling tasks with intermediate parallelizability. *ACM Transactions on Parallel Computing (TOPC)*, 3(1):1–19, 2016.
- [32] Cheeha Kim and Ashok K Agrawala. Analysis of the fork-join queue. *IEEE Transactions on computers*, 38(2):250–255, 1989.
- [33] Leonard Kleinrock. *Queueing systems, volume 2: Computer applications*, volume 66. Wiley New York, 1976.
- [34] Guy Latouche and V. Ramaswami. *Introduction to Matrix Analytic Methods in Stochastic Modeling*. ASA-SIAM, Philadelphia, 1999.
- [35] Stefano Leonardi and Danny Raz. Approximating total flow time on parallel machines. *Journal of Computer and System Sciences*, 73(6):875–891, 2007.
- [36] Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. In *Advances in Neural Information Processing Systems*, pages 5330–5340, 2017.
- [37] Richard Liaw, Romil Bhardwaj, Lisa Dunlap, Yitian Zou, Joseph E Gonzalez, Ion Stoica, and Alexey Tumanov. Hypersched: Dynamic resource reallocation for model development on a deadline. In *Proceedings of the ACM Symposium on Cloud Computing*, pages 61–73, 2019.
- [38] David Lo, Liqun Cheng, Rama Govindaraju, Parthasarathy Ranganathan, and Christos Kozyrakis. Heracles: Improving resource efficiency at scale. In *Proceedings of the 42nd Annual International Symposium on Computer Architecture*, pages 450–462, 2015.
- [39] Jason Mars, Lingjia Tang, Robert Hundt, Kevin Skadron, and Mary Lou Soffa. Bubble-up: Increasing utilization in modern warehouse scale computers via sensible co-locations. In *Proceedings of the 44th annual IEEE/ACM International Symposium on Microarchitecture*, pages 248–259, 2011.
- [40] Robert McNaughton. Scheduling with deadlines and loss functions. *Management Science*, 6(1):1–12, 1959.
- [41] Philipp Moritz, Robert Nishihara, Stephanie Wang, Alexey Tumanov, Richard Liaw, Eric Liang, Melih Elibol, Zongheng Yang, William Paul, Michael I Jordan, et al. Ray: A distributed framework for emerging {AI} applications. In *13th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 18)*, pages 561–577, 2018.
- [42] Randolph Nelson and Asser Tantawi. Approximate analysis of fork/join synchronization in parallel queues. *Transactions on Computers*, 37(6):739–743, 1988.
- [43] Marcel F. Neuts. *Matrix-Geometric Solutions in Stochastic Models*. Johns Hopkins University Press, 1981.
- [44] Marcel F. Neuts. *Structured Stochastic Matrices of M/G/1 Type and Their Applications*. Marcel Dekker, 1989.
- [45] Takayuki Osogami and Mor Harchol-Balder. Closed form solutions for mapping general distributions to quasi-minimal PH distributions. *Performance Evaluation*, 63(6):524–552, 2006.
- [46] Takayuki Osogami, Mor Harchol-Balder, and Alan Scheller-Wolf. Analysis of cycle stealing with switching times and thresholds. In *Proceedings of ACM Sigmetrics*, pages 184–195, San Diego, CA, June 2003.
- [47] Yanghua Peng, Yixin Bao, Yangrui Chen, Chuan Wu, and Chuanxing Guo. Optimus: an efficient dynamic resource scheduler for deep learning clusters. In *Proceedings of the Thirteenth EuroSys Conference*, pages 1–14, 2018.
- [48] Donald R Smith. A new proof of the optimality of the shortest remaining processing time discipline. *Operations Research*, 26(1):197–199, 1978.
- [49] Abhishek Verma, Luis Pedrosa, Madhukar Korupolu, David Oppenheimer, Eric Tune, and John Wilkes. Large-scale cluster management at google with borg. In *Proceedings of the European Conference on Computer Systems*, pages 1–17, 2015.
- [50] Weina Wang, Mor Harchol-Balder, Haotian Jiang, Alan Scheller-Wolf, and Rayadurgam Srikant. Delay asymptotics and bounds for multitask parallel jobs. *Queueing Systems*, 91(3-4):207–239, 2019.