

---

# Algorithms with Prediction Portfolios

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The research area of algorithms with predictions has seen recent success showing  
2 how to incorporate machine learning into algorithm design to improve performance  
3 when the predictions are correct, while retaining worst-case guarantees when they  
4 are not. Most previous work has assumed that the algorithm has access to a single  
5 predictor. However, in practice, there are many machine learning methods available,  
6 often with incomparable generalization guarantees, making it hard to pick a best  
7 method a priori. In this work we consider scenarios where multiple predictors are  
8 available to the algorithm and the question is how to best utilize them.

9 Ideally, we would like the algorithm’s performance to depend on the quality of the  
10 *best* predictor. However, utilizing more predictions comes with a cost, since we now  
11 have to identify which prediction is best. We study the use of multiple predictors for  
12 a number of fundamental problems, including matching, load balancing, and non-  
13 clairvoyant scheduling, which have been well-studied in the single predictor setting.  
14 For each of these problems we introduce new algorithms that take advantage of  
15 multiple predictors, and prove bounds on the resulting performance.

## 16 1 Introduction

17 An exciting recent line of research attempts to go beyond traditional worst-case analysis of algorithms  
18 by equipping algorithms with *machine-learned predictions*. The hope is that these predictions  
19 allow the algorithm to circumvent worst case lower bounds when the predictions are good, and  
20 approximately match them otherwise. The precise definitions and guarantees vary with different  
21 settings, but there have been significant successes in applying this framework for many different  
22 algorithmic problems, ranging from very general online problems to classical graph algorithms (see  
23 Section 1.2 for a more detailed discussion of related work, and [29] for a survey). In all of these  
24 settings it turns out to be possible to define a “prediction” where the “quality” of the algorithm  
25 (competitive ratio, running time, etc.) depends the “error” of the prediction. Moreover, in at least  
26 some of these settings, it has been further shown that this prediction is actually learnable with a small  
27 number of samples, usually via standard ERM methods [15].

28 Previous work has shown the power of accurate predictions, and there are numerous examples  
29 showing improved performance in both theory and practice. However, developing accurate predictors  
30 remains an art, and a single predictor may not capture all of the subtleties of the instance space.  
31 Recently, researchers have turned to working with *portfolios of predictors*: instead of training a single  
32 model, we instead train multiple models, with the hope that one of them will give good guarantees.

33 It is easy to see why the best predictor in a portfolio may be *significantly* better than a one-size fits all  
34 predictor. First, many of the modern machine learning methods come with a slew of hyperparameters  
35 that require tuning. Learning rate, mini-batch size, optimizer choice, all of these have significant  
36 impact on the quality of the final solution. Instead of committing to a single setting, one can instead  
37 try to cover the parameter space, with the hope that some of the predictors will generalize better than

38 others. Second, problem instances themselves may come from complex distributions, consisting of  
39 many latent groups or clusters. A single predictor is forced to perform well on average, whereas  
40 multiple predictors can be made to “specialize” to each cluster.

41 In order to take advantage of the increased accuracy provided by the portfolio approach, we must  
42 adapt algorithms with predictions to take advantage of multiple predictions. To recoup the gains in  
43 performance, the algorithm must perform as if equipped with the best predictor, auto-tuning to use the  
44 best one available in the portfolio. However, it is easy to see that there should be a cost as the size of  
45 the portfolio grows. In the extreme, one can add every possible prediction to the portfolio, providing  
46 no additional information, yet now requiring high performance from the algorithm. Therefore, we  
47 must aim to minimize the dependence on the number of predictions in the portfolio.

48 We remark that the high level set up may be reminiscent of expert- or bandit-learning literature.  
49 However, there is a critical distinction. In expert and bandit learning, we are given a sequence of  
50 problem instances, and the goal is to compete (minimize regret) with respect to the best prediction  
51 *averaged* over the whole sequence. On the other hand, in our setup, we aim to compete with the best  
52 predictor on a *per-instance* basis.

53 **Previous work on multiple predictions.** Anand, Ge, Kumar, and Panigrahi recently initiated the  
54 study of algorithms with multiple learned predictions in [6], proving strong bounds for important  
55 online covering problems including online set cover, weighted caching, and online facility location.  
56 However, their techniques and results are limited to online covering problems. Moreover, they do not  
57 discuss the learning aspects at all: they simply assume that they are given  $k$  predictions, and their  
58 goal is to have competitive ratios that are based on the minimum error of any of the  $k$  predictions.  
59 (They actually compete against a stronger dynamic benchmark, but for our purposes this distinction  
60 is not important).

61 On the other hand Balcan et al. [12] look at this problem through a data driven algorithm lens and  
62 study the sample complexity and generalization error of working with  $k$  (as opposed to 1) parameter  
63 settings. The main difference from our work is that they also aim learn a selector, which selects one  
64 of the  $k$  parameters *prior* to beginning to solve the problem instance. In contrast, in this work we  
65 make the selection during the course of the algorithm, and sometimes switch back and forth while  
66 honing in on the best predictor.

## 67 1.1 Our Results and Contributions

68 In this paper we study three fundamental problems, min-cost perfect matching, online load balancing,  
69 and non-clairvoyant scheduling for total completion time, in this new setting. Each of these has seen  
70 significant success in the single-prediction model but is not covered by previous multiple-prediction  
71 frameworks.

72 For each of these we develop algorithms whose performance depends on the error of the *best*  
73 prediction, and explore the effect of the number of predictions,  $k$ . Surprisingly, in the case of  
74 matching and scheduling we show that using a limited number of multiple predictions is essentially  
75 free, and has *no* asymptotic impact on the algorithm’s performance. For load balancing, on the other  
76 hand, we show that the cost of multiple predictions grows *logarithmically* with  $k$ , again implying a  
77 tangible benefit of using multiple predictions. We now describe these in more detail.

78 **Min-Cost Perfect Matching.** We warm up to showcase our approach with the classical min-cost  
79 perfect matching problem in Section 3. This problem was recently studied by [14, 15] to show that it  
80 is possible to use learned predictions to improve *running times* of classical optimization problems.  
81 In particular, [15] showed it is possible to speed up the classical Hungarian algorithm by predicting  
82 dual values, and moreover that it is possible to efficiently (PAC-)learn the best duals. We show that  
83 simple modifications of their ideas lead to similar results for multiple predictions. Interestingly, we  
84 show that as long as  $k \leq O(\sqrt{n})$ , the extra “cost” (running time) of using  $k$  predictions is negligible  
85 compared to the cost of using a single prediction, so we can use up to  $\sqrt{n}$  predictions “for free” while  
86 still getting running time depending on the best of these predictions. Moreover, since in this setting  
87 running time is paramount, we go beyond sample complexity to show that it is also computationally  
88 efficient to learn the best  $k$  predictions.

89 **Online Load Balancing with Restricted Assignments.** We continue in Section 4 with the funda-  
90 mental load balancing problem. In this problem there are  $m$  machines, and  $n$  jobs which appear in  
91 online fashion. Each job has a size, and a subset of machines that it can be assigned to. The goal is to  
92 minimize the maximum machine load (i.e., the makespan). This problem has been studied extensively  
93 in the traditional scheduling and online algorithms literature, and recently it has also been the subject  
94 of significant study given a single prediction [22–24]. In particular, Lattanzi, Lavastida, Moseley,  
95 and Vassilvitskii [22] showed that there exist per machine “weights” and an allocation function so  
96 that the competitive ratio of the algorithm depends logarithmically on the maximum error of the  
97 predictions. We show that one can use  $k$  predictions and incur an additional  $O(\log k)$  factor in the  
98 competitive ratio, while being competitive with the error of the *best* prediction. Additionally, we  
99 show that learning the best  $k$  predicted weights (in a PAC sense) can be done efficiently.

100 **Non Clairvoyant Scheduling** Finally, in Section 5 we move to the most technically complex  
101 part of this paper. We study the problem of scheduling  $n$  jobs on a single machine, where all jobs  
102 are released at time 0, but where we do not learn the length of a job until it actually completes (the  
103 *non-clairvoyant* model). Our objective is to minimize the sum of completion times. This problem  
104 has been studied extensively, both with and without predictions [20, 26, 31, 32]. Most recently,  
105 Lindermayr and Megow [26] suggested that we use an *ordering* as the prediction (as opposed to  
106 the more obvious prediction of job sizes), and use the difference between the cost induced by the  
107 predicted ordering and the cost induced by the instance-optimal ordering as the notion of “error”.  
108 In this case, simply following the predicted ordering yields an algorithm with error equal to the  
109 prediction error.

110 We extend this to the multiple prediction setting, which turns out to be surprisingly challenging. The  
111 algorithm of [26] is quite simple: follow the ordering given by the prediction (and run a 2-competitive  
112 algorithm in parallel to obtain a worst-case backstop). But we obviously cannot do this when we are  
113 given multiple orderings! So we must design an algorithm which considers all  $k$  predictions to build  
114 a schedule that has error comparable to the error of the *best* one. Slightly more formally, we prove  
115 that we can get sum of completion times bounded by at most  $(1 + \epsilon)\text{OPT}$  plus  $\text{poly}(1/\epsilon)$  times the  
116 error of the best prediction, under the mild assumption that no set of at most  $\log \log n$  jobs has a large  
117 contribution to  $\text{OPT}$ .

118 To do this, we first use sampling techniques similar to those of [20] to estimate the size of the  
119 approximately  $\epsilon n$ ’th smallest job without incurring much cost. We then use even more sampling and  
120 partial processing to determine for each prediction whether its  $\epsilon n$  prefix has many jobs that should  
121 appear later (a bad sequence) or has very few jobs that should not be in the prefix (a good sequence).  
122 If all sequences are bad then every prediction has large error, so we can use a round robin schedule  
123 and charge the cost to the prediction error. Otherwise, we choose one of the good orderings and  
124 follow it for its  $\epsilon n$  prefix (being careful to handle outliers). We then recurse on the remaining jobs.

## 125 1.2 Related Work

126 As discussed, the most directly related papers are Anand et al. [6] and Balcan, Sandholm, and  
127 Vitercik [12]; these give the two approaches (multiple predictions and portfolio-based algorithm  
128 selection) that are most similar to our setting. The single prediction version of min-cost bipartite  
129 matching was studied in [14, 15], the single prediction version of our load balancing problem was  
130 considered by [22–24] (and a different though related load balancing problem was considered by [4]),  
131 and the single prediction version of our scheduling problem was considered by [26] with the same  
132 prediction that we use (an ordering) and earlier with different predictions by [20, 32]. Online  
133 scheduling with estimates of the true processing times was considered in [9, 10].

134 More generally, there has been an enormous amount of recent progress on algorithms with predictions.  
135 This is particularly true for online algorithms, where the basic setup was formalized by [27] in the  
136 context of caching. For example, the problems considered include caching [21, 27, 33], secretary  
137 problems [7, 17], ski rental [5, 32], and set cover [13]. There has also been recent work on going  
138 beyond traditional online algorithms, including work on running times [14, 15], algorithmic game  
139 theory [2, 18, 28], and streaming algorithms [1, 16, 19].

140 The above is only a small sample of the work on algorithms with predictions. We refer the interested  
141 reader to a recent survey [29], as well as a recently set up website which maintains a list of papers in  
142 the area [25].

143 **2 Learnability**

144 When designing new methods in the algorithms with predictions setting, the predictions under  
 145 consideration must satisfy two constraints. First, they should be useful to the algorithm, so that  
 146 using the predictions allows the algorithm to achieve better running time, competitive ratio, or some  
 147 other performance measure. Second, they must be *learnable*; that is it must be feasible to find good  
 148 predictions given a set of problem instances.

149 To rigorously prove learnability, we follow previous work [11, 15, 30] and focus on proving a bound  
 150 on the sample complexity of finding the best predictions that generalize.

151 Our main result shows that for a given problem, the pseudo-dimension of finding  $k$  predictions is  
 152  $\tilde{O}(k)^1$  factor larger than that for finding a single best predictor. We state the formal Theorem below,  
 153 but defer the proof to the supplementary material.

154 **Theorem 2.1.** *Let  $\mathcal{F}$  be a class of functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  with pseudo-dimension  $d$  and let  $\mathcal{F}^k :=$   
 155  $\{F(x) = \min_{\ell \in [k]} f^\ell(x) \mid f^1, f^2, \dots, f^k \in \mathcal{F}\}$ . Then the pseudo-dimension of  $\mathcal{F}^k$  is at most  
 156  $\tilde{O}(dk)$ .*

157 Note that this directly implies that the sample complexity when looking for  $k$  predictions is a factor  
 158 of  $k$  larger than that of a single predictor.

159 **3 Minimum Cost Bipartite Matching with Predicted Duals**

160 In this section we study the minimum cost bipartite matching problem with multiple predictions. The  
 161 case of a single prediction has been considered recently [14, 15], where they used dual values as a  
 162 prediction and showed that the classical Hungarian algorithm could be sped up by using appropriately  
 163 learned dual values. Our goal in this section is to extend these results to multiple predictions, i.e.,  
 164 multiple duals. In particular, in Section 3.2 we show that we can use  $k$  duals and get running time  
 165 comparable to the time we would have spent if we used the single best of them in the algorithm of [15],  
 166 with no loss whatsoever if  $k$  is at most  $O(\sqrt{n})$ . Then in Section 3.3 we show that  $k$  predictions can  
 167 be learned with not too many more samples (or running time) than learning a single prediction.

168 **3.1 Problem Definition and Predicted Dual Variables**

169 In minimum cost bipartite matching there is a bipartite graph  $G = (V, E)$  (letting  $n = |V|$  and  
 170  $m = |E|$ ) with edge costs  $c \in \mathbb{Z}^E$ . The objective is to output a perfect matching  $M \subseteq E$  which  
 171 minimizes the cost  $c(M) := \sum_{e \in E} c_e$ . This problem is exactly captured by the following primal and  
 172 dual linear programming formulations.

173 Recently, Dinitz et al. [15] studied initializing the Hungarian  
 174 algorithm with a prediction  $\hat{y}$  of the optimal dual solution  $y^*$ .  
 175 They propose an algorithm which operates in two steps. First,  
 176 the predicted dual solution  $\hat{y}$  may not be feasible, so they  
 177 give an  $O(n + m)$  time algorithm which recovers feasibility  
 178 (which we refer to as Make-Feasible). Second, the now-  
 179 feasible dual solution is used in a primal-dual algorithm such  
 180 as the Hungarian algorithm (which we refer to as Primal-  
 181 Dual) and they show that the running time depends on the  $\ell_1$   
 182 error in the predicted solution. In addition to this they show  
 183 that learning a good initial dual solution is computationally  
 184 efficient with low sample complexity. More formally, they  
 185 proved the following theorems.

$$\begin{aligned} \min \quad & \sum_{e \in E} c_e x_e && \text{(MWPM-P)} \\ & \sum_{e \in N(i)} x_e = 1 && \forall i \in V \\ & x_e \geq 0 && \forall e \in E \\ \max \quad & \sum_{i \in V} y_i && \text{(MWPM-D)} \\ & y_i + y_j \leq c_e && \forall e = ij \in E \end{aligned}$$

186 **Theorem 3.1** (Dinitz et al. [15]). *Let  $(G, c)$  be an instance of minimum cost bipartite matching and  
 187  $\hat{y}$  be a prediction of an optimal dual solution  $y^*$ . There exists an algorithm which returns an optimal  
 188 solution and runs in time  $O(m\sqrt{n} \cdot \|y^* - \hat{y}\|_1)$ .*

189 **Theorem 3.2** (Dinitz et al. [15]). *Let  $\mathcal{D}$  be an unknown distribution over instances  $(G, c)$  on  $n$   
 190 vertices and let  $y^*(G, c)$  be an optimal dual solution for the given instance. Given  $S$  independent*

<sup>1</sup> $\tilde{O}(\cdot)$  suppresses logarithmic factors

191 samples from  $\mathcal{D}$ , there is a polynomial time algorithm that outputs a solution  $\hat{y}$  such that

$$\mathbb{E}_{(G,c) \sim \mathcal{D}} [\|y^*(G,c) - \hat{y}\|_1] \leq \min_y \mathbb{E}_{(G,c) \sim \mathcal{D}} [\|y^*(G,c) - y\|_1] + \epsilon$$

192 with probability  $1 - \delta$  where  $S = \text{poly}(n, \frac{1}{\epsilon}, \frac{1}{\delta})$ .

### 193 3.2 Using $k$ Predicted Dual Solutions Efficiently

194 Given  $k$  predicted dual solutions  $\hat{y}^1, \hat{y}^2, \dots, \hat{y}^k$ , we would like to efficiently determine which solution  
 195 has the minimum error for the given problem instance. Note that the predicted solutions may still be  
 196 infeasible and that we do not know the target optimal dual solution  $y^*$ . We propose the following  
 197 simple algorithm which takes as input  $k$  predicted solutions and whose running time depends only  
 198 on the  $\ell_1$  error of the *best* predicted solution. First, we make each predicted solution feasible, just  
 199 as before. Next, we select the (now-feasible) dual solution with highest dual objective value and  
 200 proceed running the primal-dual algorithm with only that solution. See Algorithm 1 for pseudo-code.

201 We have the following  
 202 result concerning Algo-  
 203 rithm 1. To interpret this  
 204 result, note that the cost for  
 205 increasing the number of  
 206 predictions is  $O(k(n+m))$ ,  
 207 which will be dominated  
 208 by the  $m\sqrt{n}$  term we pay  
 209 for running the Hungarian  
 210 algorithm unless  $k$  is  
 211 extremely large (certainly  
 212 larger than  $\sqrt{n}$ ) or there  
 213 is a prediction with 0 error

---

#### Algorithm 1 Minimum cost matching with $k$ predicted dual solutions

---

```

1: procedure  $k$ -PREDICTEDPRIMAL-DUAL( $G, c, \hat{y}^1, \hat{y}^2, \dots, \hat{y}^k$ )
2:   for  $\ell \in [k]$  do
3:      $y^\ell \leftarrow \text{MakeFeasible}(G, c, \hat{y}^\ell)$ 
4:   end for
5:    $\ell' \leftarrow \arg \max_{\ell \in [k]} \sum_{i \in V} y_i^\ell$ 
6:    $M \leftarrow \text{Primal-Dual}(G, c, y^{\ell'})$ 
7:   Return  $M$ 
8: end procedure

```

---

(which is highly unlikely). Hence we can reap the benefit of a large number of predictions “for free”.

215 **Theorem 3.3.** *Let  $(G, c)$  be a minimum cost bipartite matching instance and let  $\hat{y}^1, \hat{y}^2, \dots, \hat{y}^k$  be*  
 216 *predicted dual solutions. Algorithm 1 returns an optimal solution and runs in time  $O(k(n+m)) +$*   
 217  *$m\sqrt{n} \cdot \min_{\ell \in [k]} \|y^* - \hat{y}^\ell\|_1$ .*

218 We defer the proof to the supplementary material. But correctness is essentially direct from [15], and  
 219 the running time requires just a simple modification of the analysis of [15].

### 220 3.3 Learning $k$ Predicted Dual Solutions

221 Next we extend Theorem 3.2 to the setting where we output  $k$  predictions. Let  $\mathcal{D}$  be a distribution  
 222 over problem instances  $(G, c)$  on  $n$  vertices. We show that we can find the best set of  $k$  predictions.  
 223 More formally, we prove the following theorem.

224 **Theorem 3.4.** *Let  $\mathcal{D}$  be an unknown distribution over instances  $(G, c)$  on  $n$  vertices and let  $y^*(G, c)$*   
 225 *be an optimal dual solution for the given instance. Given  $S$  independent samples from  $\mathcal{D}$ , there is a*  
 226 *polynomial time algorithm that outputs a  $k$  solutions  $\hat{y}^1, \hat{y}^2, \dots, \hat{y}^k$  such that*

$$\mathbb{E}_{(G,c) \sim \mathcal{D}} \left[ \min_{\ell \in [k]} \|y^*(G, c) - \hat{y}^\ell\|_1 \right] \leq O(1) \cdot \min_{y^1, y^2, \dots, y^k} \mathbb{E}_{(G,c) \sim \mathcal{D}} \left[ \min_{\ell \in [k]} \|y^*(G, c) - y^\ell\|_1 \right] + \epsilon$$

227 with probability  $1 - \delta$  where  $S = \text{poly}(n, k, \frac{1}{\epsilon}, \frac{1}{\delta})$ .

228 The proof of this theorem can be found in the supplementary material, but it is straightforward.  
 229 The sample complexity is due to combining Theorem 2.1 with Theorem 3.2 (or more precisely,  
 230 with the pseudo-dimension bound which implies Theorem 3.2). The  $O(1)$ -approximation factor  
 231 and polynomial running time is from the observation that the ERM problem in this setting is just  
 232  $k$ -median.

## 233 4 Online Load Balancing with Predicted Machine Weights

234 We now apply our framework to online load balancing with restricted assignments. In particular, we  
 235 consider proportional weights, which have been considered in prior work [22–24]. Informally, we

236 show in Section 4.2 that if  $\beta$  is the cost of the *best* of the  $k$  predictions, then even without knowing a  
 237 priori which prediction is best, we get cost of  $O(\beta \log k)$ . Then in Section 4.3 we show that it does  
 238 not take many samples to actually learn the best  $k$  predictions.

#### 239 4.1 Problem Definition and Proportional Weights

240 In online load balancing with restricted assignments there is a sequence of  $n$  jobs which must be  
 241 assigned to  $m$  machines in an online fashion. Upon seeing job  $j$ , the online algorithm observes  
 242 its size  $p_j > 0$  and a neighborhood  $N(j) \subseteq [m]$  of *feasible* machines. The algorithm must then  
 243 choose some feasible machine  $i \in N(j)$  to irrevocably assign the job to before seeing any more  
 244 jobs in the sequence. We also consider fractional assignments, i.e. vectors belonging to the set  
 245  $X = \{x \in \mathbb{R}_+^{m \times n} \mid \forall j \in [n], \sum_i x_{ij} = 1, \text{ and } x_{ij} = 0 \iff i \notin N(j)\}$ .

246 Prior work studied the application of proportional weights[3, 22–24]. Intuitively, a prediction in this  
 247 setting is a weighting of *machines*, which then implies an online assignment, which is shown to be  
 248 near-optimal. Slightly more formally, suppose that we are given weights  $w_i$  for each machine  $i$ . Then  
 249 we say that each job  $j$  is *fractionally* assigned to machine  $i$  to a fractional amount of  $\frac{w_i}{\sum_{i' \in N(j)} w_{i'}}$ .

250 Notice that given weights, this also gives an online assignment. It is known that there exist weights for  
 251 any instance where the fractional solution has a near optimal makespan, even though there are only  
 252  $m$  “degree of freedom” in the weights compared to  $mn$  in an assignment. That is, for all machines  $i$ ,  
 253  $\sum_{j \in [n]} p_j \cdot \frac{w_i}{\sum_{i' \in N(j)} w_{i'}}$  is at most a  $(1 + \epsilon)$  factor larger than the optimal makespan for any constant  
 254  $\epsilon > 0$  [3, 22].

255 Let  $w^*$  be a set of near optimal weights for a given instance. Lattanzi et al. [22] showed the following  
 256 theorem:

257 **Theorem 4.1.** *Given predicted weights  $\hat{w}$ , there is an online fractional algorithm which has makespan*  
 258  *$O(\log(\eta(\hat{w}, w^*) \text{OPT}))$ , where  $\eta(\hat{w}, w^*) := \max_{i \in [m]} \max(\frac{\hat{w}_i}{w_i^*}, \frac{w_i^*}{\hat{w}_i})$  to be the error in the prediction.*

259 Moreover, this fractional assignment can be converted online to an integral assignment while losing  
 260 only an  $O(\log \log m)$  factor in the makespan [22, 24].

#### 261 4.2 Combining Fractional Solutions Online

262 Given  $k$  different predicted weight vectors  $\hat{w}^1, \hat{w}^2, \dots, \hat{w}^k$ , we want to give an algorithm which is  
 263 competitive against the *minimum* error among the predicted weights, i.e. we want the competitiveness  
 264 to depend upon  $\eta_{\min} := \min_{\ell \in [k]} \eta(\hat{w}^\ell, w^*)$ .

265 The challenge is that we do not know up front which  $\ell \in [k]$  will yield the smallest error, but instead  
 266 learn this in hindsight. We consider the following approach to address this. For each  $\ell \in [k]$ , let  $x^\ell$ ,  
 267 be the resulting fractional assignment from applying the fractional online algorithm due to [22] with  
 268 weights  $\hat{w}^\ell$ . This fractional assignment is revealed one job at a time.

269 We give an algorithm which is  $O(\log k)$ -competitive against any collection of  $k$  fractional assignments  
 270 which are revealed online. Moreover, our result applies to the unrelated machines setting, in which  
 271 each job has a collection of machine-dependent sizes  $\{p_{ij}\}_{i \in [m]}$ . The algorithm is based on the  
 272 doubling trick and is similar to results due to [8] which apply to metrical task systems. Let  $\beta :=$   
 273  $\min_{\ell \in [k]} \max_i \sum_j p_{ij} x_{ij}^\ell$  be the best fractional makespan in hindsight. As in previous work, our  
 274 algorithm is assumed to know  $\beta$ , an assumption that can be removed [22]. At a high level, our  
 275 algorithm maintains a set  $A \subseteq [k]$  of solutions which are good with respect to the current value of  $\beta$ ,  
 276 averaging among these. See Algorithm 2 for a detailed description. We have the following theorem.

277 **Theorem 4.2.** *Let  $x^1, x^2, \dots, x^k$  be fractional assignments which are revealed online. If Algorithm 2*  
 278 *is run with  $\beta := \min_{\ell \in [k]} \max_i \sum_j p_{ij} x_{ij}^\ell$ , then it yields a solution of cost  $O(\log k) \cdot \beta$  and never*  
 279 *reaches the fail state (line 7 in Algorithm 2).*

280 Let  $\beta_\ell = \max_i \sum_j p_{ij} x_{ij}^\ell$  and OPT be the optimal makespan. Theorem 4.1 shows that  $\beta_\ell \leq$   
 281  $O(\log \eta_\ell) \text{OPT}$ . The following corollary is then immediate:

282 **Corollary 4.3.** *Let  $w^1, w^2, \dots, w^k$  be the predicted weights with errors  $\eta^1, \eta^2, \dots, \eta^k$ . Then Algo-*  
 283 *gorithm 2 returns a fractional assignment with makespan at most  $\text{OPT} \cdot O(\log k) \cdot \min_{\ell \in [k]} \log(\eta^\ell)$ .*

---

**Algorithm 2** Algorithm for combining fractional solutions online for load balancing.

---

```
1: procedure COMBINE-LOADBALANCING( $\beta$ )
2:    $A \leftarrow [k]$  ▷ Initially all solutions are good
3:   for each job  $j$  do
4:     Receive the assignments  $x^1, x^2, \dots, x^k$ 
5:      $A(j, \beta) \leftarrow \{\ell \in A \mid \forall i \in [m], x_{ij}^\ell > 0 \implies p_{ij} x_{ij}^\ell \leq \beta\}$ 
6:     if  $A = \emptyset$  or  $A(j, \beta) = \emptyset$  then
7:       Return “Fail”
8:     end if
9:      $\forall i \in [m], x_{ij} \leftarrow \frac{1}{|A(j, \beta)|} \sum_{\ell \in S(j, \beta)} x_{ij}^\ell$ 
10:     $A \leftarrow \{\ell \in S \mid \max_{i \in [m]} \sum_{j' \leq j} p_{ij'} x_{ij'}^\ell > \beta\}$  ▷ Bad solutions w.r.t.  $\beta$ 
11:     $A \leftarrow A \setminus B$ 
12:  end for
13: end procedure
```

---

284 We defer the proof of Theorem 4.2 to the Supplementary material.

### 285 4.3 Learning $k$ Predicted Weight Vectors

286 We now turn to the question of showing how to learn  $k$  different predicted weight vectors  
287  $\hat{w}^1, \hat{w}^2, \dots, \hat{w}^k$ . Recall that there is an unknown distribution  $\mathcal{D}$  over sets of  $n$  jobs from which  
288 we receive independent samples  $J_1, J_2, \dots, J_S$ . Our goal is to show that we can efficiently learn (in  
289 terms of sample complexity)  $k$  predicted weight vectors to minimize the expected minimum error.  
290 Let  $w^*(J)$  be the correct weight vector for instance  $J$  and let  $\eta(w, w') = \max_{i \in [m]} \max(\frac{w_i}{w'_i}, \frac{w'_i}{w_i})$  be  
291 the error between a pair of weight vectors. We have the following result.

292 **Theorem 4.4.** *Let  $\mathcal{D}$  be an unknown distribution over restricted assignment instances on  $n$  jobs*  
293 *and let  $w^*(J)$  be a set of good weights for instance  $J$ . Given  $S$  independent samples from  $\mathcal{D}$ ,*  
294 *there is a polynomial time algorithm that outputs  $k$  weight vectors  $\hat{w}^1, \hat{w}^2, \dots, \hat{w}^k$  such that*  
295  $\mathbb{E}_{J \sim \mathcal{D}} [\min_{\ell \in [k]} \log(\eta(\hat{w}^\ell, w^*(J)))] \leq O(1) \cdot \min_{w^1, w^2, \dots, w^k} \mathbb{E} [\min_{\ell \in [k]} \log(\eta(w^\ell, w^*(J)))] + \epsilon$   
296 *with probability  $1 - \delta$ , where  $S = \text{poly}(m, k, \frac{1}{\epsilon}, \frac{1}{\delta})$*

297 The proof of Theorem 4.4 is deferred to the Supplementary material, but we note that to get a  
298 polynomial time algorithm we carry out an interesting reduction to  $k$ -median clustering. Namely, we  
299 show that the function  $d(w, w') := \log(\eta(w, w'))$  satisfies the triangle inequality and thus forms a  
300 metric space.

## 301 5 Scheduling with Predicted Permutations

302 In this problem there are  $n$  jobs, indexed by  $1, 2, \dots, n$ , to be scheduled on a single machine. We  
303 assume that they are all available at time 0. Job  $j$  has size  $p_j$  and needs to get processed for  $p_j$  time  
304 units to complete. If all job sizes are known a priori, Shortest Job First (or equivalently Shortest  
305 Remaining Time First), which processes jobs in non-decreasing order of their size, is known to be  
306 optimal for minimizing total completion time. We assume that the true value of  $p_j$  is unknown and is  
307 revealed only when the job completes, i.e. the *non-clairvoyant* setting. In the non-clairvoyant setting,  
308 it is known that Round-Robin (which processes all alive jobs equally) is 2-competitive and that this is  
309 the best competitive ratio one can hope for [31].

310 We study this basic scheduling problem assuming certain predictions are available for use. Fol-  
311 lowing the recent work by Lindermayr and Megow [26], we will assume that we are given  $k$   
312 orderings/sequences as prediction,  $\{\sigma_\ell\}_{\ell \in [k]}$ . Each  $\sigma_\ell$  is a permutation of  $J := [n]$ . Intuitively, it  
313 suggests an ordering in which jobs should be processed. This prediction is inspired by the afore-  
314 mentioned Shortest Job First (SJF) as an optimum schedule can be described as an ordering of jobs,  
315 specifically increasing order of job sizes.

316 For each  $\sigma_\ell$ , its error is measured as  $\eta(J, \sigma_\ell) := \text{COST}(J, \sigma_\ell) - \text{OPT}(J)$ , where  $\text{COST}(J, \sigma_\ell)$  denotes  
317 the objective of the schedule where jobs are processed in the order of  $\sigma_\ell$  and  $\text{OPT}(J)$  denotes the  
318 optimum objective. We may drop  $J$  from notation when it is clear from the context.

319 As observed in [26], the error can be expressed as  $\eta(J, \sigma_\ell) = \sum_{i < j \in J} I_{i,j}^\ell \cdot |p_i - p_j|$ , where  $I_{i,j}^\ell$   
320 is an indicator variable that has value 1 if and only if  $\sigma_\ell$  predicts the pairwise ordering of  $i$  and  $j$   
321 incorrectly. That is, if  $p_i < p_j$ , then the optimal schedule would process  $i$  before  $j$ ; here  $I_{i,j}^\ell = 1$  iff  
322  $i \succ_{\sigma_\ell} j$ .

323 As discussed in [26], this error measure satisfies two desired properties, monotonicity and Lipschitz-  
324 ness, which were formalized in [20].

325 Our main result is the following.

326 **Theorem 5.1.** *Consider a constant  $\epsilon > 0$ . Suppose that for any  $S \subseteq J$  with  $|S| = \Theta(\frac{1}{\epsilon^3}(\log \log n +$   
327  $\log k + \log(1/\epsilon)))$ , we have  $\text{OPT}(S) \leq c\epsilon \cdot \text{OPT}(J)$  for some small absolute constant  $c$ . Then,  
328 there exists a randomized algorithm that yields a schedule whose expected total completion time is at  
329 most  $(1 + \epsilon)\text{OPT} + (1 + \epsilon)\frac{1}{\epsilon^3}\eta(J, \sigma_\ell)$  for all  $\ell \in [k]$ .*

330 As a corollary, by running our algorithm with  $1 - \epsilon$  processing speed and simultaneously running  
331 Round-Robin with the remaining  $\epsilon$  of the speed, the cost increases by a factor of at most  $\frac{1}{1-\epsilon}$  while  
332 the resulting hybrid algorithm is  $2/\epsilon$ -competitive.<sup>2</sup>

## 333 5.1 Algorithm

334 To make our presentation more transparent we will first round job sizes. Formally, we choose  $\rho$   
335 uniformly at random from  $[0, 1)$ . Then, round up each job  $j$ 's size to the closest number of the form  
336  $(1 + \epsilon)^{\rho+t}$  for some integer  $t$ . Then, we scale down all job sizes by  $(1 + \epsilon)^\rho$  factor. We will present  
337 our algorithm and analysis assuming that every job has a size equal to a power of  $(1 + \epsilon)$ . In the  
338 supplementary we show how to remove this assumption without increasing our algorithm's objective  
339 by more than  $1 + \epsilon$  factor in expectation.

340 We first present the following algorithm that achieves Theorem 5.1 with  $|S| = \Theta(\frac{1}{\epsilon^3}(\log n + \log k))$ .  
341 The improved bound claimed in the theorem needs minor tweaks of the algorithm and analysis and  
342 they are deferred to the supplementary material.

343 Our algorithm runs in rounds. Let  $J_r$  be the jobs that complete in round  $r \geq 1$ . For any subset  $S$  of  
344 rounds,  $J_S := \cup_{r \in S} J_r$ . For example,  $J_{\leq r} := J_1 \cup \dots \cup J_r$ . Let  $n_r := |J_{\geq r}| = n - |J_{< r}|$  denote  
345 the number of alive jobs at the beginning of round  $r$ .

346 Fix the beginning of round  $r$ . The algorithm processes the job in the following way for this round. If  
347  $n_r \leq \frac{1}{\epsilon^3}(\log n + \log k)$ , we run Round-Robin to complete all the remaining jobs,  $J_{\geq r}$ . This is the  
348 last round and it is denoted as round  $L + 1$ . Otherwise, we do the following Steps 1-4:

349 **Step 1. Estimating  $\epsilon$ -percentile.** Roughly speaking, the goal is to estimate the  $\epsilon$ -percentile of job  
350 sizes among the remaining jobs. For a job  $j \in J_{\geq r}$ , define its rank among  $J_{\geq r}$  as the number of jobs  
351 no smaller than  $j$  in  $J_{\geq r}$  breaking ties in an arbitrary yet fixed way. Ideally, we would like to estimate  
352 the size of job of rank  $\epsilon n_r$ , but do so only approximately.

353 The algorithm will find  $\tilde{q}_r$  that is the size of a job whose rank lies in  $[\epsilon(1 - \epsilon)n_r, \epsilon(1 + \epsilon)n_r]$ . To  
354 handle the case that there are many jobs of the same size  $\tilde{q}_r$ , we estimate  $y_r$  the number of jobs no  
355 bigger than  $\tilde{q}_r$ ; let  $\tilde{y}_r$  denote our estimate of  $y_r$ . We will show how we can do these estimations  
356 without spending much time by sampling some jobs and partially processing them in Round-Robin  
357 manner (the proof of the following lemma can be found in the supplementary material.)

358 **Lemma 5.2.** *W.h.p. the algorithm can construct estimates  $\tilde{q}_r$  and  $\tilde{y}_r$  in time at most  $O(\tilde{q}_r \frac{1}{\epsilon^2} \log n)$   
359 such that there is a job of size  $\tilde{q}_r$  whose rank lies in  $[\epsilon(1 - \epsilon)n_r, \epsilon(1 + \epsilon)n_r]$  and  $|\tilde{y}_r - y_r| \leq \epsilon^2 n_r$ .*

360 **Step 2. Determining Good and Bad Sequences.** Let  $\sigma_\ell^r$  denote  $\sigma_\ell$  with all jobs completed in the  
361 previous rounds removed and with the relative ordering of the remaining jobs fixed. Let  $\sigma_{\ell,\epsilon}^r$  denote  
362 the first  $\tilde{y}_r$  jobs in the ordering. We say a job  $j$  is big if  $p_j > \tilde{q}_r$ ; middle if  $p_j = \tilde{q}_r$ ; small otherwise.  
363 Using sampling and partial processing we will approximately distinguish good and bad sequences.

<sup>2</sup>This hybrid algorithm is essentially the preferential time sharing [20, 26, 32]. Formally, we run our algorithm ignoring RR's processing and also run RR ignoring our algorithm; this can be done by a simple simulation. Thus, we construct two schedules concurrently and each job completes at the time when it does in either schedule. This type of algorithms was first used in [32].

364 Informally  $\sigma_\ell^r$  is good if  $\sigma_{\ell,\epsilon}^r$  has few big jobs and bad if it does many big jobs. The proof of the  
 365 following lemma can be found in the supplementary material.

366 **Lemma 5.3.** *For all  $\ell \in [k]$ , we can label sequence  $\sigma_\ell^r$  either good or bad in time at most*  
 367  *$O(\tilde{q}_r \frac{1}{\epsilon^2} (\log n + \log k))$  that satisfies the following with high probability: If it is good,  $\sigma_{\ell,\epsilon}^r$  has*  
 368 *at most  $3\epsilon^2 n_r$  big jobs; otherwise  $\sigma_{\ell,\epsilon}^r$  has at least  $\epsilon^2 n_r$  big jobs.*

369 **Step 3. Job Processing.** If all sequences are bad, then we process all jobs, each up to  $\tilde{q}_r$  units in an  
 370 arbitrary order. Otherwise, we process the first  $\tilde{y}_r$  jobs in an arbitrary good sequence, in an arbitrary  
 371 order, each up to  $\tilde{q}_r$  units.

372 **Step 4. Updating Sequences.** The jobs completed in this round drop from the sequences but the  
 373 remaining jobs' relative ordering remains fixed in each (sub-)sequence. For simplicity, we assume that  
 374 partially processed jobs were never processed—this is without loss of generality as this assumption  
 375 only increases our schedule's objective.

## 376 5.2 Analysis of the Algorithm's Performance

377 We defer the analysis of the above algorithm (the proof of Theorem 5.1) to the supplementary material,  
 378 as it is quite technical and complex. At a very high level, though, we use the fact that the error in each  
 379 prediction can be decomposed into pair-wise inversions, and moreover we can partition the inversions  
 380 into the rounds of the algorithm in which they appear. Then we look at each round, and split into  
 381 two cases. First, if all sequences are bad then every prediction has large error, so we can simply use  
 382 Round Robin (which is 2-competitive against OPT) and the cost can be charged to the error of any  
 383 prediction. Second, if there is a good sequence, then in any good sequence the number of big jobs  
 384 is small (so we do not spend much time processing them), and we therefore complete almost all of  
 385 the non-big jobs. Here, we crucially use the fact that we can process the first  $\epsilon$  fraction of jobs in a  
 386 sequence in an arbitrary order remaining competitive against the sequence. Finally, we show that all  
 387 of the additional assumptions and costs (e.g., rounding processing times and the cost due to sampling)  
 388 only change our performance by a  $1 + \epsilon$  factor. Getting all of these details right requires much care.

## 389 5.3 Learning $k$ Predicted Permutations

390 Now we show that learning the best  $k$  permutations has polynomial sample complexity.

391 **Theorem 5.4.** *Let  $\mathcal{D}$  be an unknown distribution of instances on  $n$  jobs. Given  $S$  indepen-*  
 392 *dent samples from  $\mathcal{D}$ , there is an algorithm that outputs  $k$  permutations  $\hat{\sigma}_1, \hat{\sigma}_2, \dots, \hat{\sigma}_k$  such that*  
 393  *$\mathbb{E}_{J \sim \mathcal{D}} [\min_{\ell \in [k]} \eta(J, \hat{\sigma}_\ell)] \leq \min_{\sigma_1, \sigma_2, \dots, \sigma_k} \mathbb{E}_{J \sim \mathcal{D}} [\min_{\ell \in [k]} \eta(J, \sigma_\ell)] + \epsilon$  with probability  $1 - \delta$ ,*  
 394 *where  $S = \text{poly}(n, k, \frac{1}{\epsilon}, \frac{1}{\delta})$ .*

395 *Proof.* The algorithm is basic ERM, and the polynomial sample complexity follows from Theorem 2.1  
 396 and Theorem 20 in Lindermayr and Megow [26].  $\square$

## 397 6 Conclusion

398 Despite the explosive recent work in algorithms with predictions, almost all of this work has assumed  
 399 only a single prediction. In this paper we study algorithms with *multiple* machine-learned predictions,  
 400 rather than just one. We study three different problems that have been well-studied in the single  
 401 prediction setting but not with multiple predictions: faster algorithms for min-cost bipartite matching  
 402 using learned duals, online load balancing with learned machine weights, and non-clairvoyant  
 403 scheduling with order predictions. For all of the problems we design algorithms that can utilize  
 404 multiple predictions, and show sample complexity bounds for learning the best set of  $k$  predictions.

405 Surprisingly, we have shown that in some cases, using multiple predictions is essentially “free.” For  
 406 instance, in the case of min-cost perfect matching examining  $k = O(\sqrt{n})$  predictions takes the same  
 407 amount of time as one round of the Hungarian algorithm, but the number of rounds is determined by  
 408 the quality of the *best* prediction. In contrast, for load balancing, using  $k$  predictions always incurs  
 409 an  $O(\log k)$  cost, so using a constant number of predictions may be best. More generally, studying  
 410 this trade-off between the cost and the benefit of multiple predictions for other problems remains an  
 411 interesting and challenging open problem.

412 **References**

- 413 [1] Anders Aamand, Piotr Indyk, and Ali Vakilian. (learned) frequency estimation algorithms under  
414 zipfian distribution. *arXiv preprint arXiv:1908.05198*, 2019.
- 415 [2] Priyank Agrawal, Eric Balkanski, Vasilis Gkatzelis, Tingting Ou, and Xizhi Tan. Learning-  
416 augmented mechanism design: Leveraging predictions for facility location, 2022. URL <https://arxiv.org/abs/2204.01120>.  
417
- 418 [3] Shipra Agrawal, Morteza Zadimoghaddam, and Vahab S. Mirrokni. Proportional allocation:  
419 Simple, distributed, and diverse matching with high entropy. In Jennifer G. Dy and Andreas  
420 Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML  
421 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of  
422 Machine Learning Research*, pages 99–108. PMLR, 2018. URL [http://proceedings.mlr.  
423 press/v80/agrawal18b.html](http://proceedings.mlr.press/v80/agrawal18b.html).
- 424 [4] Sara Ahmadian, Hossein Esfandiari, Vahab Mirrokni, and Binghui Peng. *Robust Load Balancing  
425 with Machine Learned Advice*, pages 20–34. 2022. doi: 10.1137/1.9781611977073.2. URL  
426 <https://epubs.siam.org/doi/abs/10.1137/1.9781611977073.2>.
- 427 [5] Keerti Anand, Rong Ge, and Debmalya Panigrahi. Customizing ML predictions for online  
428 algorithms. In *Proceedings of the 37th International Conference on Machine Learning, ICML  
429 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*,  
430 pages 303–313. PMLR, 2020. URL [http://proceedings.mlr.press/v119/anand20a.  
431 html](http://proceedings.mlr.press/v119/anand20a.html).
- 432 [6] Keerti Anand, Rong Ge, Amit Kumar, and Debmalya Panigrahi. Online algorithms with multiple  
433 predictions. In *Proceedings of the 39th International Conference on Machine Learning, ICML  
434 2022, 2022*.
- 435 [7] Antonios Antoniadis, Themis Gouleakis, Pieter Kleer, and Pavel Kolev. Secretary and online  
436 matching problems with machine learned advice. In Hugo Larochelle, Marc’Aurelio Ranzato,  
437 Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Informa-  
438 tion Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020,  
439 NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL [https://proceedings.neurips.  
440 cc/paper/2020/hash/5a378f8490c8d6af8647a753812f6e31-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/5a378f8490c8d6af8647a753812f6e31-Abstract.html).
- 441 [8] Yossi Azar, Andrei Z. Broder, and Mark S. Manasse. On-line choice of on-line algorithms. In Vi-  
442 jaya Ramachandran, editor, *Proceedings of the Fourth Annual ACM/SIGACT-SIAM Symposium  
443 on Discrete Algorithms, 25-27 January 1993, Austin, Texas, USA*, pages 432–440. ACM/SIAM,  
444 1993. URL <http://dl.acm.org/citation.cfm?id=313559.313847>.
- 445 [9] Yossi Azar, Stefano Leonardi, and Noam Touitou. Flow time scheduling with uncertain  
446 processing time. In Samir Khuller and Virginia Vassilevska Williams, editors, *STOC ’21:  
447 53rd Annual ACM SIGACT Symposium on Theory of Computing, Virtual Event, Italy, June  
448 21-25, 2021*, pages 1070–1080. ACM, 2021. doi: 10.1145/3406325.3451023. URL <https://doi.org/10.1145/3406325.3451023>.  
449
- 450 [10] Yossi Azar, Stefano Leonardi, and Noam Touitou. Distortion-oblivious algorithms for mini-  
451 mizing flow time. In Joseph (Seffi) Naor and Niv Buchbinder, editors, *Proceedings of the 2022  
452 ACM-SIAM Symposium on Discrete Algorithms, SODA 2022, Virtual Conference / Alexandria,  
453 VA, USA, January 9 - 12, 2022*, pages 252–274. SIAM, 2022. doi: 10.1137/1.9781611977073.13.  
454 URL <https://doi.org/10.1137/1.9781611977073.13>.
- 455 [11] Maria-Florina Balcan, Dan F. DeBlasio, Travis Dick, Carl Kingsford, Tuomas Sandholm, and  
456 Ellen Vitercik. How much data is sufficient to learn high-performing algorithms? *CoRR*,  
457 <abs/1908.02894>, 2019. URL <http://arxiv.org/abs/1908.02894>.
- 458 [12] Maria-Florina Balcan, Tuomas Sandholm, and Ellen Vitercik. Generalization in portfolio-based  
459 algorithm selection. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021,  
460 Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The  
461 Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual  
462 Event, February 2-9, 2021*, pages 12225–12232. AAAI Press, 2021. URL [https://ojs.aaai.  
463 org/index.php/AAAI/article/view/17451](https://ojs.aaai.org/index.php/AAAI/article/view/17451).

- 464 [13] Étienne Bamas, Andreas Maggiori, and Ola Svensson. The primal-dual method for learning  
465 augmented algorithms. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina  
466 Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33:  
467 Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December  
468 6-12, 2020, virtual*, 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/  
469 e834cb114d33f729dbc9c7fb0c6bb607-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/e834cb114d33f729dbc9c7fb0c6bb607-Abstract.html).
- 470 [14] Justin Y. Chen, Sandeep Silwal, Ali Vakilian, and Fred Zhang. Faster fundamental graph  
471 algorithms via learned predictions. *CoRR*, abs/2204.12055, 2022. doi: 10.48550/arXiv.2204.  
472 12055. URL <https://doi.org/10.48550/arXiv.2204.12055>.
- 473 [15] Michael Dinitz, Sungjin Im, Thomas Lavastida, Benjamin Moseley, and Sergei Vassil-  
474 vitskii. Faster matchings via learned duals. In Marc’Aurelio Ranzato, Alina Beygelz-  
475 zimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan, editors, *Ad-  
476 vances in Neural Information Processing Systems 34: Annual Conference on Neural  
477 Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*,  
478 pages 10393–10406, 2021. URL [https://proceedings.neurips.cc/paper/2021/hash/  
479 5616060fb8ae85d93f334e7267307664-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/5616060fb8ae85d93f334e7267307664-Abstract.html).
- 480 [16] Elbert Du, Franklyn Wang, and Michael Mitzenmacher. Putting the “learning” into learning-  
481 augmented algorithms for frequency estimation. In Marina Meila and Tong Zhang, editors,  
482 *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Pro-  
483 ceedings of Machine Learning Research*, pages 2860–2869. PMLR, 18–24 Jul 2021. URL  
484 <https://proceedings.mlr.press/v139/du21d.html>.
- 485 [17] Paul Dütting, Silvio Lattanzi, Renato Paes Leme, and Sergei Vassilvitskii. Secretaries with  
486 advice. *CoRR*, abs/2011.06726, 2020. URL <https://arxiv.org/abs/2011.06726>.
- 487 [18] Vasilis Gkatzelis, Kostas Kollias, Alkmini Sgouritsa, and Xizhi Tan. Improved price of anarchy  
488 via predictions, 2022. URL <https://arxiv.org/abs/2205.04252>.
- 489 [19] Chen-Yu Hsu, Piotr Indyk, Dina Katabi, and Ali Vakilian. Learning-based frequency estimation  
490 algorithms. In *7th International Conference on Learning Representations*, 2019.
- 491 [20] Sungjin Im, Ravi Kumar, Mahshid Montazer Qaem, and Manish Purohit. Non-clairvoyant  
492 scheduling with predictions. In Kunal Agrawal and Yossi Azar, editors, *SPAA ’21: 33rd  
493 ACM Symposium on Parallelism in Algorithms and Architectures, Virtual Event, USA, 6-8  
494 July, 2021*, pages 285–294. ACM, 2021. doi: 10.1145/3409964.3461790. URL <https://doi.org/10.1145/3409964.3461790>.
- 495 [21] Zhihao Jiang, Debmalya Panigrahi, and Kevin Sun. Online algorithms for weighted paging  
496 with predictions. In Artur Czumaj, Anuj Dawar, and Emanuela Merelli, editors, *47th Interna-  
497 tional Colloquium on Automata, Languages, and Programming, ICALP 2020, July 8-11, 2020,  
498 Saarbrücken, Germany (Virtual Conference)*, volume 168 of *LIPICs*, pages 69:1–69:18. Schloss  
499 Dagstuhl - Leibniz-Zentrum für Informatik, 2020. doi: 10.4230/LIPICs.ICALP.2020.69. URL  
500 <https://doi.org/10.4230/LIPICs.ICALP.2020.69>.
- 501 [22] Silvio Lattanzi, Thomas Lavastida, Benjamin Moseley, and Sergei Vassilvitskii. Online  
502 scheduling via learned weights. In Shuchi Chawla, editor, *Proceedings of the 2020 ACM-  
503 SIAM Symposium on Discrete Algorithms, SODA 2020, Salt Lake City, UT, USA, January  
504 5-8, 2020*, pages 1859–1877. SIAM, 2020. doi: 10.1137/1.9781611975994.114. URL  
505 <https://doi.org/10.1137/1.9781611975994.114>.
- 506 [23] Thomas Lavastida, Benjamin Moseley, R. Ravi, and Chenyang Xu. Learnable and instance-  
507 robust predictions for online matching, flows and load balancing. In Petra Mutzel, Rasmus  
508 Pagh, and Grzegorz Herman, editors, *29th Annual European Symposium on Algorithms, ESA  
509 2021, September 6-8, 2021, Lisbon, Portugal (Virtual Conference)*, volume 204 of *LIPICs*, pages  
510 59:1–59:17. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2021. doi: 10.4230/LIPICs.  
511 ESA.2021.59. URL <https://doi.org/10.4230/LIPICs.ESA.2021.59>.
- 512 [24] Shi Li and Jiayi Xian. Online unrelated machine load balancing with predictions revisited. In  
513 Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on*  
514

- 515 *Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of*  
516 *Machine Learning Research*, pages 6523–6532. PMLR, 2021. URL <http://proceedings.mlr.press/v139/li21w.html>.
- 518 [25] Alexander Lindermayr and Nicole Megow. Algorithms with predictions. [https://](https://algorithms-with-predictions.github.io/)  
519 [algorithms-with-predictions.github.io/](https://algorithms-with-predictions.github.io/), 2022.
- 520 [26] Alexander Lindermayr and Nicole Megow. Non-clairvoyant scheduling with predictions revisited. *CoRR*, abs/2202.10199, 2022. URL <https://arxiv.org/abs/2202.10199>.
- 522 [27] Thodoris Lykouris and Sergei Vassilvitskii. Competitive caching with machine learned advice. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*, pages 3302–3311, 2018.
- 525 [28] Andrés Muñoz Medina and Sergei Vassilvitskii. Revenue optimization with approximate bid predictions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, page 1856–1864, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- 529 [29] Michael Mitzenmacher and Sergei Vassilvitskii. *Algorithms with Predictions*, page 646–662. Cambridge University Press, 2021. doi: 10.1017/9781108637435.037.
- 531 [30] Jamie H Morgenstern and Tim Roughgarden. On the pseudo-dimension of nearly optimal auctions. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 136–144. Curran Associates, Inc., 2015. URL <http://papers.nips.cc/paper/5766-on-the-pseudo-dimension-of-nearly-optimal-auctions.pdf>.
- 536 [31] Rajeev Motwani, Steven Phillips, and Eric Torng. Nonclairvoyant scheduling. *Theoretical Computer Science*, 130(1):17–47, 1994.
- 538 [32] Manish Purohit, Zoya Svitkina, and Ravi Kumar. Improving online algorithms via ml predictions. In *Advances in Neural Information Processing Systems*, pages 9661–9670, 2018.
- 540 [33] Dhruv Rohatgi. Near-optimal bounds for online caching with machine learned advice. In *Symposium on Discrete Algorithms (SODA)*, 2020.

## 542 Checklist

- 543 1. For all authors...
- 544 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
545 contributions and scope? [Yes]
- 546 (b) Did you describe the limitations of your work? [Yes]
- 547 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 548 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
549 them? [Yes]
- 550 2. If you are including theoretical results...
- 551 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 552 (b) Did you include complete proofs of all theoretical results? [Yes]
- 553 3. If you ran experiments...
- 554 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
555 mental results (either in the supplemental material or as a URL)? [N/A]
- 556 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
557 were chosen)? [N/A]
- 558 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
559 ments multiple times)? [N/A]
- 560 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
561 of GPUs, internal cluster, or cloud provider)? [N/A]

- 562 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 563 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 564 (b) Did you mention the license of the assets? [N/A]
- 565 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 566
- 567 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 568 using/curating? [N/A]
- 569 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 570 information or offensive content? [N/A]
- 571 5. If you used crowdsourcing or conducted research with human subjects...
- 572 (a) Did you include the full text of instructions given to participants and screenshots, if
- 573 applicable? [N/A]
- 574 (b) Did you describe any potential participant risks, with links to Institutional Review
- 575 Board (IRB) approvals, if applicable? [N/A]
- 576 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 577 spent on participant compensation? [N/A]