

Improving Ranking Quality and Fairness in Swiss-System Chess Tournaments

Pascal Führllich,¹ Ágnes Cseh,^{2,3} Pascal Lenzner³

¹ Potsdam Institute for Climate Impact Research, Potsdam, Germany

² Institute of Economics, Centre for Economic and Regional Studies, Budapest, Hungary

³ Hasso Plattner Institute, University of Potsdam, Potsdam, Germany

pascal.fuehrlich@pik-potsdam.de, cseh.agnes@krtk.hu, pascal.lenzner@hpi.de

Abstract

The International Chess Federation (FIDE) imposes a voluminous and complex set of player pairing criteria in Swiss-system chess tournaments and endorses computer programs that are able to calculate the prescribed pairings. The purpose of these formalities is to ensure that players are paired fairly during the tournament and that the final ranking corresponds to the players' true strength order.

We contest the official FIDE player pairing routine by presenting alternative pairing rules. These can be enforced by computing maximum weight matchings in a carefully designed graph. We demonstrate by extensive experiments that a tournament format using our mechanism 1) yields fairer pairings in the rounds of the tournament and 2) produces a final ranking that reflects the players' true strengths better than the state-of-the-art FIDE pairing system.

Introduction

How can one determine the relative strength of players who engage in a one-on-one competitive game? This is easy to find out for a group of two players: just let them play a match. For more players, tournaments solve this problem by ranking the players after a limited number of pairwise matches among the participants. The *tournament format* defines a general structure of matches to be played and the method for deriving a ranking from the results of those matches.

Tournament Formats Most tournaments follow an elimination, a round-robin, or a Swiss-system format. In each round of an *elimination* tournament, such as the second stage of the FIFA World Cup, only players who won their match in the previous round are paired again. The last player standing wins the tournament, and the remaining players' strength can only be estimated very roughly from the round they were eliminated in. *Round-robin* tournaments are also called all-play-all tournaments, because each player plays against each other player once. The player with the highest score at the end of the tournament is declared the winner. The pool stage of the FIFA World Cup consists of round-robin tournaments.

The *Swiss-system* tournament format is widely used in competitive games like most e-sports, badminton, and chess,

the last of which this paper focuses on. In such tournaments, the number of rounds is predefined, but the pairing of players in each of these rounds depends on the results of previous rounds. This format offers a convenient golden middle way between the earlier mentioned two tournament formats. However, the features of the Swiss system challenge organizers greatly. Firstly, unlike in elimination tournaments, the goal is to determine a whole ranking of the players and not only to declare the winner. Secondly, the final ranking of each player is greatly influenced by her assigned opponents, which is not an issue in round-robin tournaments.

Therefore, a mechanism that computes suitable player pairings for Swiss-system tournaments is crucially important. However, designing such a system is a challenging task as it boils down to solving a complex combinatorial optimization problem. Interestingly, the state-of-the-art solution to this problem in chess tournaments relies on a complex set of declarative rules and not on a combinatorial optimization algorithm. In this paper we provide an algorithmic approach and we demonstrate that it outperforms the declarative state-of-the-art solution. For this, we do not try to mimic the FIDE solution but instead focus on the most important features of the Swiss system and derive a maximum weight matching formulation that enforces them.

The Swiss-System in Chess In Swiss-system chess tournaments, there are two well-defined and rigid *absolute* and two milder *quality* pairing criteria.

- (A1) No two players play against each other more than once.
- (A2) In each round before the last one, the difference of matches played with white and matches played with black pieces is between -2 and 2 for every player.
- (Q1) Opponents have equal or similar score.
- (Q2) Each player has a balanced color distribution.

Criterion (A1) ensures variety, while criterion (A2) ensures fairness, since the player with white pieces starts the game, and thus has an advantage over her opponent. These absolute criteria must be obeyed at any cost, which often enforces the relaxation of the two quality criteria.

In order to implement criterion (Q1), players with equal score are grouped into *score groups*. In each round, a chosen *pairing system* allocates each player an opponent from the same score group. If a complete pairing is not possible

within a score group, then one or more players are moved to another score group. Criterion (Q2) requires that after each round of the tournament, the difference between matches played with black and white pieces is small for each player.

Adhering to these four criteria makes pairing design truly challenging. Pairings at FIDE tournaments were traditionally calculated manually by so-called arbiters, often using trial-and-error. Today, pairings are computed by decision-making software, but the FIDE pairing criteria are still written for human instead of computer execution. Over the years, more and more criteria were added to resolve ambiguities, which increased the complexity to a level at which pairing decisions are very challenging to comprehend for most players and even arbiters.

Related Literature

Novel algorithms that assist tournament scheduling regularly evoke interest in the AI community (Larson, Johansson, and Carlsson 2014; Kim and Williams 2015; Chatterjee, Ibsen-Jensen, and Tkadlec 2016; Gupta et al. 2018; Hoshino 2018). We first elaborate on existing work on comparing tournament formats, and then turn to approaches that utilize matchings for scheduling tournaments.

Comparing Tournament Formats Appleton (1995) gives an overview of tournament formats and compares them with respect to how often the best player wins. Scarf, Yusof, and Bilbao (2009) simulate different tournament formats using team data from the UEFA champions league. Elmenreich, Ibounig, and Fehérvári (2009) compare several sorting algorithms, including one based on a Swiss-system tournament, with respect to their robustness, which is defined as the degree of similarity between the resulting ranking and the true strength order of players. They find round-robin sort, merge sort, and Swiss-system sort to be the most robust overall.

Automated Matching Approaches A tournament schedule can be seen as a set of matchings—one for each round. Glickman and Jensen (2005) propose an algorithm based on maximum weight perfect matchings to find the schedule. This algorithm maximizes the information gain about players’ skill. The authors’ approach compares favorably against random and Swiss-system pairing if at least 16 rounds are played. However, almost all real-world Swiss-system chess tournaments have less than 10 rounds according to chess-results.com (Herzog 2020a).

Kujansuu, Lindberg, and Mäkinen (1999) use the stable roommates problem, see (Irving 1985), to model a Swiss-system tournament pairing decision. Each player p has a preference list, which ranks the other players by how desirable a match between player p and each other player would be. The desirability depends on score difference and color balance. In comparison to the official FIDE pairing, this approach produces pairings with slightly better color balance but higher score differences between paired players, or, in other words, clearly favors criterion (Q2) over (Q1).

Weighted Matching Models for Chess Tournaments

The two papers closest to ours focus on modeling the exact

FIDE pairing criteria and computing the prescribed pairings.

Ólafsson (1990) pairs players using a maximum weight matching algorithm on a graph, where players and possible matches are represented by vertices and edges. Edge weights are set so that they model the 1985 FIDE pairing criteria. At that time, pairing criteria were more ambiguous than today, and pairing was done by hand, which sometimes took several hours. In contrast, using Ólafsson’s method, pairings could be calculated fast. Pairings calculated with the commercial software built by Ólafsson are claimed to be preferred by experts to manually calculated pairings. However, Ólafsson only provides examples and does not present any comparison based on formal criteria.

A more recent attempt to convert the FIDE pairing criteria into a weighted matching instance was undertaken by Biró, Fleiner, and Palincza (2017). Due to the extensive criterion system, only a subset of the criteria were modeled. The authors show that a pairing respecting these selected criteria can be calculated in polynomial time, and leave it as a challenging open question whether the other FIDE criteria can also be integrated into a single weighted matching model. The contribution appears to be purely theoretical, since neither a comparison with other pairing programs, nor implementation details are provided.

Our work breaks the line of research that attempts to implement the declarative FIDE pairing criteria via weighted matchings. Instead, we design new pairing rules along with a different mechanism to compute the pairings, and demonstrate their superiority compared to the FIDE pairing criteria and engine. This clearly differentiates our approach from the one in (Ólafsson 1990; Biró, Fleiner, and Palincza 2017).

Preliminaries and FIDE Criteria

Players are entities participating in a Swiss-system tournament. Each player has an *Elo rating*, which is a measure designed to capture her current playing strength from the outcome of her earlier matches (Elo 1978). In a *match* two players, a and b , play against each other. The three possible *match results* are: a wins and b loses, a and b draw, a loses and b wins. The winner receives 1 point, the loser 0 points, while a draw is worth 0.5 points. A Swiss-system tournament consists of multiple *rounds*, each of which is defined by a *pairing*: a set of disjoint pairs of players, where each pair plays a match. At the end of the tournament, a strict ranking of the players is derived from the match results.

Bye Allocation In general, each player plays exactly one match per round. For an odd number of players, one of them receives a so-called ‘bye’, which is a point rewarded without a match. This is always the player currently ranked last among those who have not yet received a bye.

Color Balance The FIDE Handbook (FIDE 2020, Chapter C.04.1) states that ‘For each player the difference between the number of black and the number of white games shall not be greater than 2 or less than -2.’ This criterion may only be relaxed in the last round. This corresponds to our criterion (A2). Also, a ban on a color that is assigned to a player three times consecutively, and further milder criteria

are phrased to ensure a color assignment as close to an alternating white-black sequence as possible (FIDE 2020, Chapters C.04.3.A.6 and C.04.3.C).

Pairing Systems Players are always ranked by their current tournament score. Furthermore, within each score group the players are ranked by their Elo rank. The score groups and this ranking are the input of the *pairing system*, which assigns an opponent to each player. Three main pairing systems are defined for chess tournaments. Table 1 shows an example pairing for each of them.

- **Dutch:** Each score group is cut into an upper and a lower half. The upper half is then paired against the lower half so that the i th ranked player in the upper half plays against the i th ranked player in the lower half. Dutch is the de facto standard for major chess tournaments.
- **Burstein:** For each score group, the highest ranked unpaired player is paired against the lowest ranked unpaired player repeatedly until all players are paired.
- **Monrad:** In ascending rank order each unpaired player in a score group is paired against the next highest ranked player in that score group.

Dutch	Burstein	Monrad
1–5	1–8	1–2
2–6	2–7	3–4
3–7	3–6	5–6
4–8	4–5	7–8

Table 1: Example pairing for each pairing system in a score group of 8 players. Players are referenced by rank within the score group, i.e., player 1 has the highest Elo rank.

For comparison, we propose two additional pairing systems based on randomness.

- **Random:** Every player within a score group is paired against a random player from her score group.
- **Random2:** Every player from the top half of her score group is paired against a random player from the bottom half of her score group.

Floating Players Players who are paired outside of their own score group are called *floaters*. To ensure that opponents are of similar strength—our criterion (Q1)—, the FIDE criteria require to minimize the number of such floaters and aim to float them to a score group of similar score. However, floating is unavoidable, e.g., in score groups with an odd number of players, and also in score groups where the first or second criterion eliminates too many possible matches.

The BBP Pairing Engine A *pairing engine* is used to calculate the pairing for each round, based on the results of previous rounds. The BBP pairing engine was developed by Bierema (2017). It implements the FIDE criteria strictly (FIDE 2020, C.04.3 and C.04.4.2) for the Dutch and Burstein pairing systems and outputs the unique pairing adhering to each of them. BBP uses a weighted matching algorithm, similarly as the approaches in (Ólafsson 1990; Biró,

Fleiner, and Palincza 2017). The main difference to our algorithm is that while the weighted model of BBP was designed to follow the declarative criteria of FIDE and output the prescribed pairings, our pairing engine relies on a different weighted model, computes completely different pairs, and while doing so, it is able to reach a better ranking quality and a higher degree of fairness. The output of Dutch BBP will serve as a base for our comparisons throughout the paper, because Dutch is the pairing system implemented by all 8 pairing programs currently endorsed by the FIDE (FIDE 2020, C.04.A.10.Annex-3).

Final Ranking The major organizing principle for the final ranking of players is obviously the final score. Players with the same final score are sorted by tiebreakers. The FIDE (FIDE 2020, Chapter C.02.13) defines 14 types of tiebreakers, and the tournament organizer lists some of them to be used at the specific tournament. If all tiebreaks fail, the tie is required to be broken by drawing of lots.

Our Contribution

In this paper, we present a novel mechanism for calculating pairings in Swiss-system chess tournaments. With this, we contest the state-of-the-art mechanism endorsed by FIDE. We compare the two systems by three measures: ranking quality, number of floaters, and color balance quality, in accordance with the FIDE tournament schedule goals. Our main findings are summarized in the following list.

1. We implemented the pairing systems Dutch, Burstein, Monrad, Random, and Random2 with an extensible and easy-to-understand approach that uses maximum weight matchings.
2. The pairing systems in descending order by expected **ranking quality** are: Burstein > Random2 > Dutch = Dutch BBP > Random > Monrad. In particular, our implementations of Burstein and Random2 both yield higher ranking quality, while our implementation of Dutch yields similar ranking quality as the one reached by the Dutch BBP pairing engine.
3. We utilize our weighted matching model to define a novel measure called ‘normalized strength difference’, which we identify as the main reason for a good ranking quality.
4. The pairing systems in ascending order by expected **number of floaters** are: Burstein < Random2 = Dutch = Monrad < Dutch BBP < Random. Compared to Dutch BBP, our mechanism is fairer in terms of matching more players within their own score group.
5. All our pairing systems ensure the same **color balance quality** as Dutch BBP, with Random even reaching a better color balance. Moreover, we show that our approach can easily be modified to enforce an even stronger color balance. This does not significantly affect the ranking quality—only the number of floaters increases slightly.
6. As the previous points demonstrate, our implementations of Burstein and Random2 either outperform or are on a par with Dutch BBP. Our implementation of Dutch leads to pairings that perform just as well or even better than

the ones prescribed by the official FIDE (Dutch) criteria and computed by Dutch BBP.

Pairings via Maximum Weight Matching

Our novel mechanism is based on computing a maximum weight matching (MWM) in an auxiliary, suitably weighted graph. The MWM engine is optimized for simplicity: score groups, color balances, and the employed pairing system are modeled by weights, so only a single computation of a MWM is needed in each round. We now describe the MWM engine.

Input

Each tournament has n players $P = \{p_1, \dots, p_n\}$, a chosen pairing system (Dutch, Burstein, Monrad, Random, or Random2), and a maximum allowed color difference β . As criterion (A2) states, FIDE aims for $\beta = 2$. If n is odd, the weakest performing player who has not received a bye yet is given one, in accordance with the FIDE rules. In the MWM engine we will exclude the same player while constructing the auxiliary graph. Hence, from this point on we can assume that n is even.

Before each round of the tournament, the following input parameters are defined for each player $p_i \in P$:

- $Elo(p_i)$: the Elo rating of p_i prior to the tournament. This remains unchanged for all rounds.
- $s(p_i)$: the current score of p_i , defined as the sum of points player p_i collected so far.
- $r(p_i)$: the current rank of p_i , calculated from ordering all players in decreasing order according to their scores and their Elo ratings. Higher score and higher Elo rating yield better rank. Players with equal Elo rating are ordered randomly at the beginning, and their order is kept for all rounds.
- $cd(p_i)$: the current color difference of p_i , defined as the number of matches played with white minus the number of matches played with black pieces.

Graph Construction

With these parameters as input, we construct the corresponding auxiliary weighted graph $G_r = (V, E, w)$ for round r as follows. Let $V := P$ and for all pairs of players $p_i \neq p_j$, let the edge set E contain the edge $\{p_i, p_j\}$ if

- (1) p_i and p_j have not yet played against each other, and
- (2) $|cd(p_i) + cd(p_j)| < 2\beta$.

These rules ensure criteria (A1) and (A2). The second condition in our model will enforce $-2 \leq cd(p_i) \leq 2$ together with our color assignment rule in Section . In the appendix we additionally consider a variant where $-1 \leq cd(p_i) \leq 1$ is enforced. This implements FIDE’s criterion that the color assignment should be as close to an alternating white-black sequence as possible and that no player can be assigned the same color three times in a row.

The weight of an edge $\{p_i, p_j\} \in E$ is defined as the tuple $w(p_i, p_j) := (-|s(p_i) - s(p_j)|, -|cd(p_i) + cd(p_j)|, \pi(p_i, p_j))$, where the value of $\pi(p_i, p_j)$ depends on the pairing system as follows.

- Monrad: $\pi(p_i, p_j) := -|r(p_i) - r(p_j)|$.
- Burstein: $\pi(p_i, p_j) := |r(p_i) - r(p_j)|^{1.01}$.
- Dutch: $\pi(p_i, p_j) := -\left|\frac{\text{sg size}}{2} - |r(p_i) - r(p_j)|\right|^{1.01}$, where sg size is set to 0 if p_i and p_j belong to different score groups, and it is the size of the score group of p_i and p_j otherwise.
- Random: $\pi(p_i, p_j) :=$ random number in the interval $(0, 1)$.
- Random2: $\pi(p_i, p_j)$ is set to a random number in the interval $(0, 1)$ if p_i and p_j belong to different halves of the same score group, otherwise it is set to a random number in the interval $(-1, 0)$.

The exponent 1.01 in the function for Burstein rewards a larger rank difference, i.e., the Burstein pairing in Table 1 indeed carries a larger weight than the Dutch pairing, which has the same sum of rank differences. Similarly, the exponent for Dutch penalizes a larger distance from $\frac{\text{sg size}}{2}$. Notice that this exponent could be an arbitrary number as long as it is larger than 1.

Algorithm

The edge weights of G_r are compared lexicographically and a maximum weight matching is sought for. This implies that pairing players within their score groups has the highest priority, optimizing color balance is second, and adhering to the pairing system is last. The comprehensive rules of our framework consist of our two absolute rules for including an edge in the graph G_r , and this priority ordering serving as our quality rule. See Figure 1 for an illustration.

Before round r , we compute a maximum weight matching M in graph G_r and derive the player pairing from the edges in M . If $\{p_i, p_j\} \in M$ then the players p_i and p_j will play against each other in round r . Between them, the respective player with the lower color difference will play white. If they have the same color difference, then colors are assigned randomly.

Assumptions and Experimental Setup

In our simulations we assume that each player $p_i \in P$ has true playing strength $str(p_i)$ that is approximated by her Elo rating $Elo(p_i)$ and we treat both values as constant throughout the tournament. The probabilities of match results and optimal rankings are defined by the playing strength. More precisely, each player’s playing strength is a random number drawn from a uniform distribution of values between 1400 and 2200. We also justified our claims on ranking quality using other realistic player strength distributions. We elaborate on these in the appendix. The results are in line with the results for the uniform distribution.

Elo ratings are used for computing $r(p_i)$ and for breaking ties in the final order. The Elo rating of player p_i is randomly drawn from a normal distribution with mean $str(p_i)$ and standard deviation $\frac{3000 - str(p_i)}{20}$. This function mirrors the assumption that a higher Elo rating estimates the strength more accurately.

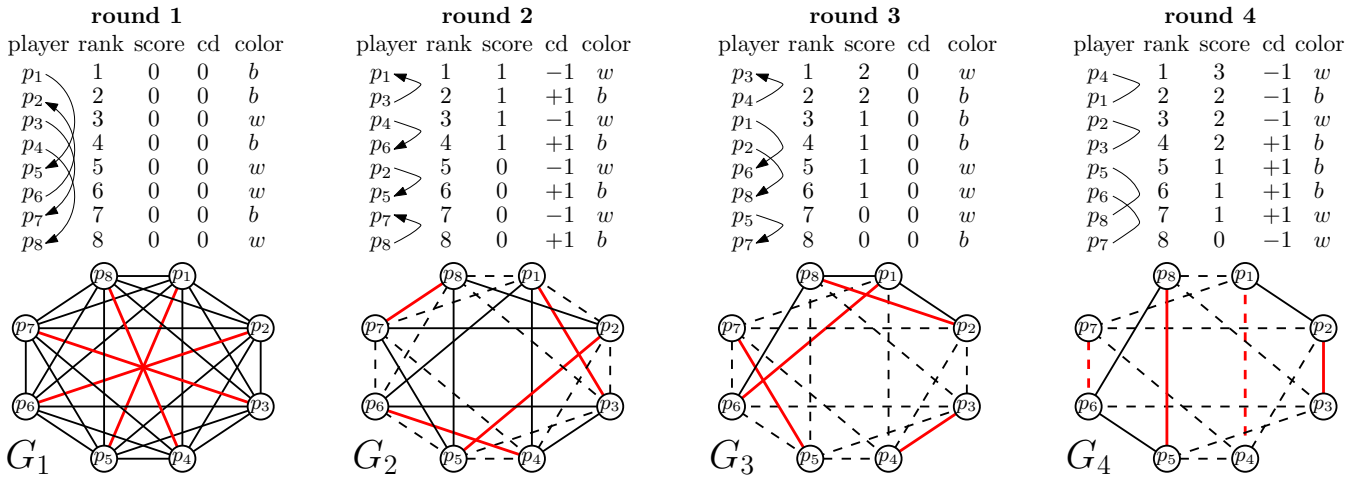


Figure 1: Example pairings of a 4-round tournament with 8 players generated via the MWM engine using the Dutch pairing system. Initially players are sorted decreasingly according to their Elo rating. Bold edges are possible matches within the same score group whereas dashed edges are other possible matches. The maximum weight matchings are shown in red. Arrows indicate the match outcomes (winner points to loser, no draws), and the color column shows the corresponding color distribution. The table for round $i + 1$ is based on the table of round i . As score and color difference are equal, the pairing in round 1 is enforced by the Dutch pairing system. The pairing in round 2 is the outcome of optimizing first for criterion (Q1) and then for criterion (Q2), e.g., in G_2 we have $w(p_1, p_3) = w(p_4, p_6) = (0, 0, -1)$ and $w(p_1, p_4) = w(p_3, p_6) = (0, -2, 0)$ so the MWM picks the edges $\{p_1, p_3\}$ and $\{p_4, p_6\}$. In G_3 players p_3 and p_4 are paired since $w(p_3, p_4) = (0, 0, 0)$ whereas the weight of any other incident edge of both p_3 and p_4 has lexicographically lower weight. The matching in G_4 is enforced by maximizing the number of matches within score groups. If p_1 and p_2 would be paired, then, since p_3 and p_4 already played, player p_4 would float to a match with a player with score 1, which implies that no match within the group with score 1 is possible.

To avoid the noise introduced by byes, we assume that the number of players n is even. The number of rounds is chosen to lie between $\lceil \log_2 n \rceil$ and $\frac{n}{2}$, as at least $\lceil \log_2 n \rceil$ rounds ensure that a player who wins all matches is the sole winner and at most $\frac{n}{2}$ rounds ensures that, according to Dirac’s theorem (Dirac 1952), a perfect matching always exists. The tiebreakers used for obtaining the final tournament ranking are based on the FIDE recommendation (FIDE 2020, C.02.13.16.5).

Computing the Maximum Weight Matching First we transform each edge weight given as a tuple to a rational number. In particular, $w(p_i, p_j)$ is transformed to $10000 \cdot (-|s(p_i) - s(p_j)|) + 100 \cdot (-|cd(p_i) + cd(p_j)|) + \pi(p_i, p_j)$. The factors 10000 and 100 ensure that each lexicographically maximum solution corresponds to a maximum weight solution with the new weights and vice versa. We compute pairings using the LEMON Graph Library (Dezső, Jüttner, and Kovács 2011) implementation of the maximum weight perfect matching algorithm, which is based on the blossom algorithm of Edmonds (1965) and has the same time and space complexity (Kolmogorov 2009). The implementation we use has $O(nm \log n)$ time complexity, where n is the number of players and m is the number of edges in the constructed graph G_r .

Realistic Probabilistic Model for Match Results The results of the individual matches are computed via a probabilistic model that is designed to be as realistic as possible. Match results are drawn at random from a suitably chosen

probability distribution based on the players’ strength and on the assigned colors for that match. For this, we use the probability distribution proposed by Milvang (2016), which was featured in a recent news article of the FIDE commission System of Pairings and Programs (FIDE SPP Commission 2020). Milvang’s probability distribution was engineered via a Data Science approach that used real-world data from almost 4 million real chess matches from 50 000 tournaments. It is based on Elo ratings and color information, whereas we use true strength values instead of Elo ratings to get unbiased match result probabilities.

Using Milvang’s approach, the probability for a certain outcome of a match depends on the actual strengths of the involved players, not only on their strength difference. Draw probability increases with mean strength of the players. The probabilities also depend on colors, as the player playing with white pieces has an advantage. See Table 2 for some example values drawn from Milvang’s distribution.

Player Strengths	Win White	Win Black	Draw
1200 (w) vs 1400 (b)	26 %	57 %	17 %
2200 (w) vs 2400 (b)	14 %	55 %	31 %
2400 (w) vs 2200 (b)	63 %	11 %	26 %

Table 2: Example match outcome probabilities drawn from Milvang’s probability distribution (Milvang 2016).

Measuring Ranking Quality Ranking quality measures how similar the tournament’s final ranking is to the ranking

that sorts the players by their strength. One popular measure for the difference between two rankings is the Kendall τ distance (Kendall 1945). It counts the number of discordant pairs: pairs of elements x and y , where $x < y$ in one ranking, but $y < x$ in the other. We use its normalized variant, where $\tau \in [-1, 1]$, and $\tau = 1$ means the rankings are identical, while $\tau = -1$ means one ranking is the inverse of the other. A higher Kendall τ is better, because it indicates a larger degree of similarity between the true and the output ranking.

We also justify our claims on ranking quality using two other well-known and possibly more sophisticated similarity measures, the Spearman ρ distance (Spearman 1904) and normalized discounted cumulative gain (NDCG). We elaborate on these measures and their behavior for our problem in the appendix. The results are in line with the ones derived for the Kendall τ distance.

Measuring Fairness We measure fairness in terms of the two relaxable criteria of Swiss-system chess tournaments: (Q1) on the equal score of opponents and (Q2) on the color distribution balance. Adhering to (Q1) is measured by the number of float pairs, which equals the number of matches with opponents from different score groups throughout the tournament. We measure the absolute color difference of a round as the sum of color differences for all players: $acd = \sum_{p_i \in P} |cd(p_i)|$. Note that as $|cd(p_i)| \geq 1$ for all players after each odd round, $acd \geq n$ in those rounds.

Presentation of the Data Data is presented in the form of *violin plots* (Hintze and Nelson 1998), *letter value plots* (Hofmann, Wickham, and Kafadar 2017), and *scatter plots* (Friendly and Denis 2005). For violin plots, kernel density estimation is used to show a smoothed probability density function of the underlying distribution. Additionally, similar to box plots, quartiles are shown by dashed lines. Letter value plots are enhanced box plots that show more quantiles. Unlike violin plots, they are suitable for discrete values, as all shown values are actual observations without smoothing.

Our plots compare the MWM implementation of the five pairing systems with the BBP implementation of Dutch.

Simulation Results

All simulations use the following parameters, unless noted otherwise:

- number of players n : 32
- number of rounds: 7
- strength range: between 1400 and 2200
- maximum allowed color difference β : 2
- sample size: 100 000

These values were chosen to be as realistic as possible, based on parameters of more than 320 000 real-world tournaments uploaded to the website chess-results.com.¹ The experiments were run on a compute server using version

¹The data was kindly provided by Heinz Herzog, author of the FIDE-endorsed tournament manager [Swiss-Manager](https://chess-results.com) (Herzog 2020b) and chess-results.com (Herzog 2020a).

20.04.1 of the Ubuntu operating system. It is powered by 48 Intel Xeon Gold 5118 CPUs running at 2.3 GHz and 62.4 GiB of RAM. We emphasize that the standard real-life challenge at a tournament, that is, computing a single pairing via a maximum weight matching for a tournament round can be solved in a fraction of a second on a standard laptop.

Ranking Quality

The pairing system of a Swiss-system tournament has a major impact on the obtained ranking quality, as Figure 2 shows. Burstein and Random2 achieve the best ranking quality, followed by Dutch and Dutch BBP. Random has a worse ranking quality and Monrad performs by far the worst. For other strength ranges, Figure 3 shows consistent results.

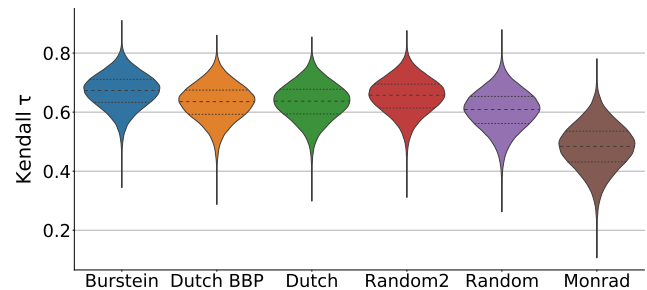


Figure 2: Ranking quality measured by normalized Kendall τ . A higher value means a better ranking quality.

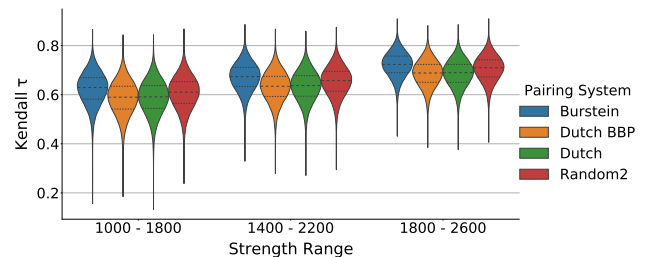


Figure 3: Ranking quality measured by normalized Kendall τ for different strength ranges.

Comparing Dutch to Dutch BBP shows that they behave very similarly, with slight advantage for Dutch. This is remarkable, since Dutch BBP is based on complex and rigid declarative criteria that are time-tested, while Dutch is the output of our easy-to-understand, purely matching-based approach. Together with the performance of Burstein and Random2 this shows that more transparent pairing systems can outperform the state-of-the-art Dutch BBP in terms of ranking quality.

We provide additional experimental results on the ranking quality in the appendix. There we present consistent results also for fewer or more players, for other strength range sizes, and for different player strength distributions. Additionally, we elaborate on how our flexible maximum weight matching model enabled us to detect the exact reason why certain pairing systems produce better rankings, which might help designing better pairing systems in the future.

Fairness

The highly complex pairing criteria of the FIDE were designed with a focus on two fairness goals phrased as quality criteria, (Q1): minimizing the number of float pairs and (Q2): minimizing the absolute color difference.

Criterion (Q1) is at the heart of Swiss-system tournaments as pairing players of equal score ensures well-balanced matches. As Figure 4 shows, Burstein, Dutch, and Random2 beat Dutch BBP in terms of the number of float pairs. In the appendix we show consistent results for other simulation parameters.

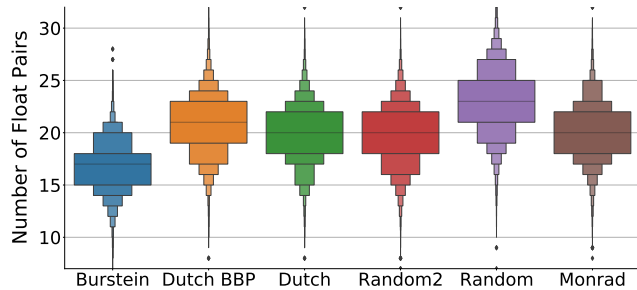


Figure 4: Number of float pairs out of the $7 \cdot 16 = 112$ matches of the tournament. Recall that floating is often unavoidable due to the size of the score group. A lower number indicates a better implementation of criterion (Q1).

Figure 5 focuses on criterion (Q2) and shows that an absolute color difference very similar to the one guaranteed by Dutch BBP can be achieved via our MWM engine. The pairing system Random even outperforms Dutch BBP in this regard. In the appendix, we provide additional experiments with different numbers of rounds and numbers of players that lead to consistent results. Also, we report there on experiments in which an even stronger color difference constraint is enforced, and observe the impact on the obtained ranking quality and the number of float pairs. Interestingly, the obtained ranking quality is almost the same but this comes at a cost of a slightly increased number of float pairs.

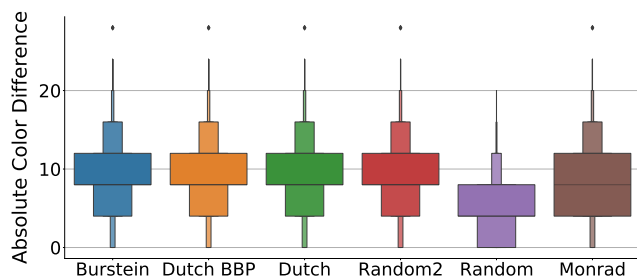


Figure 5: Absolute color difference after 6 rounds. A lower *acd* means a better color distribution. Recall that a $acd \geq n$ for each odd round, while $acd = 0$ is possible after each even round.

Hence, our maximum weight matching approach with edge weights that prioritize matches within score groups and sec-

ondly optimize for color balance is on a par with the sophisticated official FIDE criteria for criterion (Q2) and it even outperforms them for criterion (Q1). Thus, our more transparent approach ensures the same color balance quality but achieves even fewer float pairs. Moreover, our approach also allows for a different trade-off between criteria (Q1) and (Q2) that does not affect the obtained ranking quality.

Conclusion

The experimental results of our MWM engine with Burstein or Random2 demonstrate that it is possible to outperform the state-of-the-art FIDE pairing criteria in terms of both ranking quality and fairness, i.e., criteria (Q1) and (Q2), with a novel efficient mechanism that is more transparent and intelligible to all participants. The direct comparison of our MWM Dutch engine versus Dutch BBP shows that even if the same pairing system is used, MWM achieves the same ranking quality but is more powerful since it yields an improved fairness. We believe that the key to this is the direct formulation of the most important criteria as a maximum weight matching problem.

The only scenario for which we might advise against using our mechanism is when the arbiter has no access to a computing device. In order to manually produce pairings in our framework, the arbiter would need to calculate the edge weights and then execute Edmonds' blossom algorithm. Instead, the FIDE (FIDE 2020, Chapter C.04.3.D) provides manually executable rules. However, these rules include exhaustive search routines that can make the execution very slow, i.e., highly exponential in the number of players (Biró, Fleiner, and Palincza 2017). Therefore, the ill-fated arbiter has to choose between learning to execute Edmonds' blossom algorithm and following a cumbersome exponential-time pairing routine.

A clear advantage of our mechanism is that it is easily extendable: as Random and Random2 already demonstrate, a new pairing system can be implemented simply by specifying how edge weights are calculated. Similarly, as we have also demonstrated, the color balance can be adjusted by simply changing the parameter β . By thinning out the edge set in our graph, we can also reach an alternating black-white sequence for each player instead of just minimizing the color difference in each round. Also, the flexibility of the maximum weight matching approach proved to be essential for uncovering the driving force behind the achieved high ranking quality: the normalized strength difference. Hence, our approach was not only valuable for computing better pairings but also in the analysis of the obtained ranking quality.

Last but not least, the flexibility of the MWM engine likely allows to incorporate additional quality criteria like measuring fairness via the average opponent ratings. Also quality criteria of other games and sports tournaments organized in the Swiss system can be integrated into the model.

References

- Appleton, D. R. 1995. May the best man win? *Journal of the Royal Statistical Society: Series D (The Statistician)*, 44(4): 529–538.
- Bierema, J. 2017. BBP Pairings, a Swiss-system chess tournament engine. <https://github.com/BieremaBoyzProgramming/bbpPairings>. Accessed: 2021-01-06.
- Biró, P.; Fleiner, T.; and Palincza, R. 2017. Designing chess pairing mechanisms. In *10th Japanese-Hungarian Symposium on Discrete Mathematics and Its Applications May 22-25, 2017, Budapest, Hungary*, 77.
- Chatterjee, K.; Ibsen-Jensen, R.; and Tkadlec, J. 2016. Robust draws in balanced knockout tournaments. In *Proceedings of the 25th International Joint Conference on Artificial Intelligence*, 172–179.
- Dezső, B.; Jüttner, A.; and Kovács, P. 2011. LEMON—an open source C++ graph template library. *Electronic Notes in Theoretical Computer Science*, 264(5): 23–45.
- Dirac, G. A. 1952. Some theorems on abstract graphs. *Proceedings of the London Mathematical Society*, 3(1): 69–81.
- Edmonds, J. 1965. Paths, Trees, and Flowers. *Canadian Journal of Mathematics*, 17: 449–467.
- Elmenreich, W.; Ibounig, T.; and Fehérvári, I. 2009. Robustness versus performance in sorting and tournament algorithms. *Acta Polytechnica Hungarica*, 6(5): 7–18.
- Elo, A. E. 1978. *The rating of chessplayers, past and present*. Arco Pub.
- FIDE. 2020. FIDE Handbook. <https://handbook.fide.com/>. Accessed: 2021-01-07.
- FIDE SPP Commission. 2020. Probability for the outcome of a chess game based on rating. <https://spp.fide.com/2020/10/23/probability-for-the-outcome-of-a-chess-game-based-on-rating/>.
- Friendly, M.; and Denis, D. 2005. The early origins and development of the scatterplot. *Journal of the History of the Behavioral Sciences*, 41(2): 103–130.
- Glickman, M. E.; and Jensen, S. T. 2005. Adaptive paired comparison design. *Journal of Statistical Planning and Inference*, 127(1-2): 279–293.
- Gupta, S.; Roy, S.; Saurabh, S.; and Zehavi, M. 2018. When rigging a tournament, let greediness blind you. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, 275–281.
- Herzog, H. 2020a. Chess-Results.com, the international Chess-Tournaments-Results-Server. <https://chess-results.com/>. Accessed: 2020-12-07.
- Herzog, H. 2020b. Swiss-Manager. <http://www.swiss-manager.at/>. Accessed: 2020-12-07.
- Hintze, J. L.; and Nelson, R. D. 1998. Violin plots: a box plot-density trace synergism. *The American Statistician*, 52(2): 181–184.
- Hofmann, H.; Wickham, H.; and Kafadar, K. 2017. Value plots: Boxplots for large data. *Journal of Computational and Graphical Statistics*, 26(3): 469–477.
- Hoshino, R. 2018. A Recursive Algorithm to Generate Balanced Weekend Tournaments. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32.
- Irving, R. 1985. An efficient algorithm for the “stable roommates” problem. *Journal of Algorithms*, 6(4): 577–595.
- Kendall, M. G. 1945. The treatment of ties in ranking problems. *Biometrika*, 239–251.
- Kim, M. P.; and Williams, V. V. 2015. Fixing tournaments for kings, chokers, and more. In *Proceedings of the 24th International Joint Conference on Artificial Intelligence*.
- Kolmogorov, V. 2009. Blossom V: a new implementation of a minimum cost perfect matching algorithm. *Mathematical Programming Computation*, 1(1): 43–67.
- Kujansuu, E.; Lindberg, T.; and Mäkinen, E. 1999. The stable roommates problem and chess tournament pairings. *Divulgaciones Matemáticas*, 7(1): 19–28.
- Larson, J.; Johansson, M.; and Carlsson, M. 2014. An Integrated Constraint Programming Approach to Scheduling Sports Leagues with Divisional and Round-Robin Tournaments. In Simonis, H., ed., *Integration of AI and OR Techniques in Constraint Programming*, 144–158. Cham: Springer International Publishing. ISBN 978-3-319-07046-9.
- Milvang, O. 2016. Probability for the outcome of a chess game based on rating. <https://pairings.fide.com/images/stories/downloads/2016-probability-of-the-outcome.pdf>.
- Ólafsson, S. 1990. Weighted matching in chess tournaments. *Journal of the Operational Research Society*, 41(1): 17–24.
- Scarf, P.; Yusof, M. M.; and Bilbao, M. 2009. A numerical study of designs for sporting contests. *European Journal of Operational Research*, 198(1): 190–198.
- Spearman, C. 1904. The Proof and Measurement of Association between Two Things. *The American Journal of Psychology*, 15(1): 72–101.

Appendix

Ranking Quality

In the following we discuss additional simulation experiments that measure the obtained ranking quality for various parameter settings.

Different Tournament Sizes

We start with experimental results demonstrating that our findings on the ranking quality remain valid for tournaments of different sizes in terms of number of players and number of rounds.

Usually it is expected that a player who wins all matches also wins the tournament, without being tied for the first place. This can be ensured by playing at least $\lceil \log_2 n \rceil$ rounds: four rounds for 16 players, five rounds for 32 players and six rounds for 64 players. Most tournaments are five or seven rounds long, according to data from chess-results.com (Herzog 2020a).

In general, more rounds lead to higher ranking quality, although with diminishing effect, as Figure 6 shows. In terms of the achieved ranking quality, the MWM engine with Burstein outperforms Dutch BBP in all cases, except for the unrealistic case of a tournament with only two rounds.

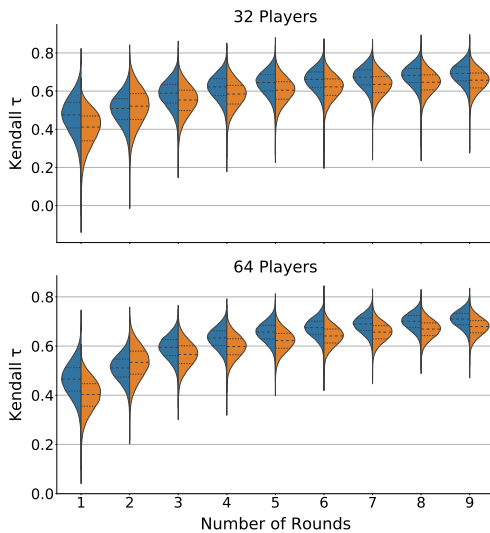


Figure 6: Ranking quality after 1-9 rounds, 32 or 64 players with strength range 1400-2200. Results for Burstein are shown in blue, Dutch BBP results are shown in orange.

Different Strength Range Sizes

Here we vary the used strength range size, i.e., we sample the player strengths from different intervals. A smaller strength range size corresponds to a tournament among players with similar strength and larger strength range sizes model tournaments with more heterogeneous players. The results depicted in Figure 7 show that also for different strength range sizes the MWM engine with Burstein or Random2 outperforms Dutch BBP in terms of ranking quality and that Dutch is on a par with Dutch BBP.

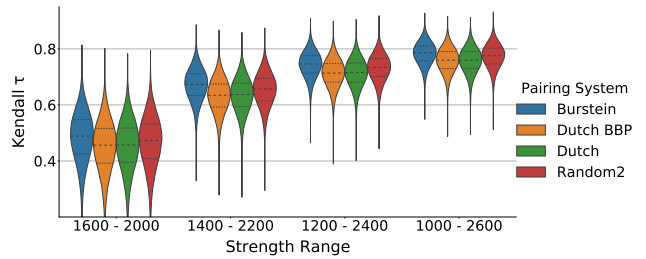


Figure 7: Ranking quality measured by normalized Kendall τ for different strength range sizes.

A higher strength range size results in higher ranking quality and less variance. The increasing ranking quality can be explained by a higher mean strength difference, which results from a larger strength range size. Variance decreases, because match results become more predictable.

The difference in ranking quality between Burstein and Dutch BBP is much higher for a strength range size of 400 compared to 800 and 1200. For small strength range sizes in all Dutch BBP paired matches it is more likely that a weaker player wins against a stronger opponent, while for Burstein at least some matches are still predictable.

Different Player Strength Distributions

We provide additional experimental results that indicate that our findings hold independently of the employed player strength distributions, i.e., we get the same behavior also for non-uniform distributions. Since no data is available that let's us estimate how realistic player strength distributions look like, we focus on several natural candidates that deviate strongly from uniform distributions.

First, we considered player strength distributions that are derived from exponential distributions. For this, we consider in Figure 8 a case with many strong players and only a few weak players and in Figure 9 a case with many weak players and only a few strong players within the given strength range size. We also considered player strength distributions de-

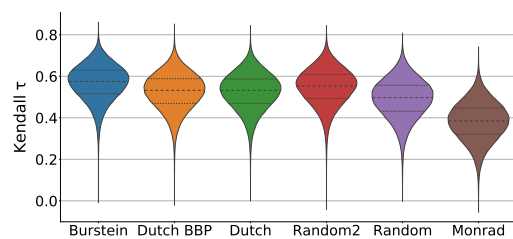


Figure 8: Ranking quality measured by normalized Kendall τ for 32 players with an exponential player strength distribution in the range $[1400, 2200]$ with mean at 2000.

derived from a normal distribution with a mean exactly in the middle of the strength range size and a standard deviation of a fourth of the strength range size. See Figure 10 for the corresponding results.

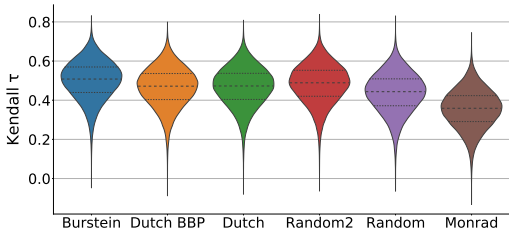


Figure 9: Ranking quality measured by normalized Kendall τ for 32 players with an exponential player strength distribution in the range [1400, 2200] with mean at 1600.

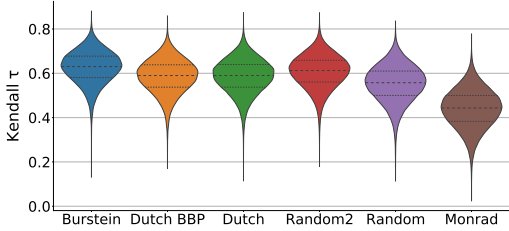


Figure 10: Ranking quality measured by normalized Kendall τ for 32 players with a normally distributed player strength distribution in the range [1400, 2200] with mean at 1800 and standard deviation of 200.

Finally, we investigated a player strength distribution that is derived from uniformly sampling player strengths from the real-world distribution of Elo scores of all 363 275 players listed by FIDE², restricted to the desired strength range. Figure 11 shows also very similar results for this case.

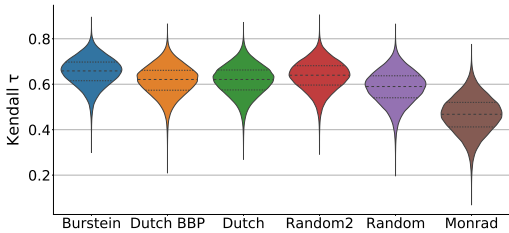


Figure 11: Ranking quality measured by normalized Kendall τ for 32 players uniformly sampled from the real-world distribution of Elo scores restricted to the range [1400, 2200].

Ranking Quality via Spearman ρ and NDCG

For comparison reasons, we provide an evaluation of the achieved ranking quality via the Spearman ρ and the normalized discounted cumulative gain (NDCG) measures.

Besides Kendall τ , Spearman ρ is commonly used for comparing rankings. Here, we use a normalized variant of Spearman ρ , similar to the normalized Kendall τ .

The NDCG measure is not commonly used for comparing rankings. It is used to evaluate search engines, by assigning

a relevance rating to documents and awarding a higher score if highly relevant documents are listed early. Applied to our case, NDCG puts an emphasis on ranking the top players correctly, while ranking the lowest ranked players correctly is basically irrelevant.

As shown in Figure 12 and Figure 13, the results with normalized Spearman ρ and NDCG look almost identical to the results for normalized Kendall τ in Figure 2.

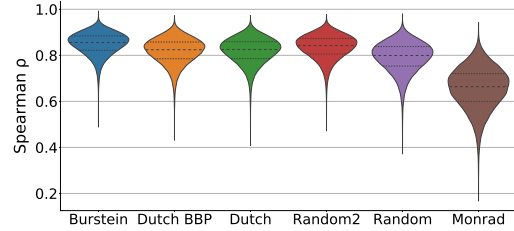


Figure 12: Ranking quality measured by normalized Spearman ρ .

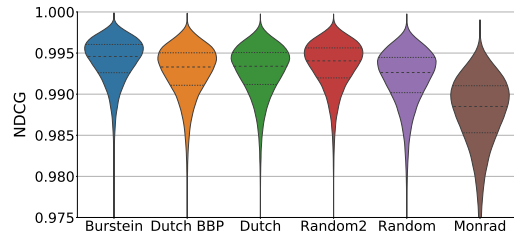


Figure 13: Ranking quality measured by the normalized discounted cumulative gain (NDCG).

Also for different strength ranges or range sizes we get consistent results, as shown in Figures 14, 15, 16 and 17.

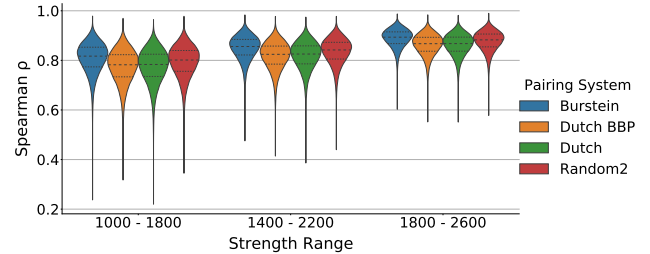


Figure 14: Ranking quality measured by normalized Spearman ρ .

Reasons for High Ranking Quality

Here we provide experiments that shed light on why Burstein, Random2, Dutch, and Dutch BBP reach a better ranking quality than Random and Monrad and why Burstein and Random2 outperform Dutch BBP.

In order to rank players correctly, their relative playing strength must be approximated from match results. We call a match result *paradoxical* if a weaker player wins against a

²See https://ratings.fide.com/download_lists.phtml for details.

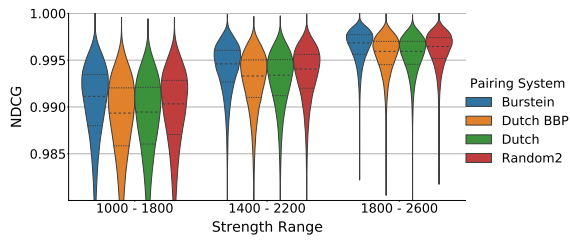


Figure 15: Ranking quality measured by the normalized discounted cumulative gain (NDCG).

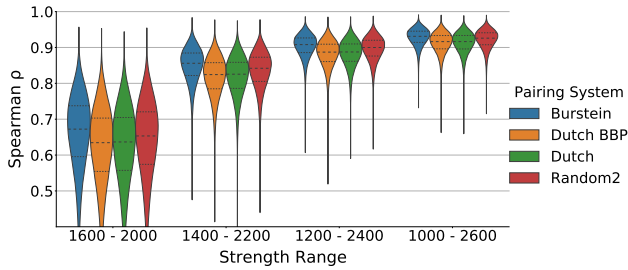


Figure 16: Ranking quality measured by normalized Spearman ρ .

stronger opponent. Paradoxical match results hinder the approximation of both players' strengths, so pairing systems should aim to minimize the number of paradoxical match results. Figure 18 confirms this, by showing a strong negative correlation between the proportion of paradoxical results and ranking quality for Dutch BBP. A similar correlation can be observed for all pairing systems.

The probability of a paradoxical match result increases as the strength difference of paired players decreases because the outcome of those matches is less predictable. In general, a higher mean strength difference in a tournament lowers the number of paradoxical match results, which then leads to better ranking quality. Our results in Section justify the observation that mean strength difference seems to be positively correlated with ranking quality, as mean strength difference is low when using Monrad, medium with Random, and high for Burstein, Random2 and Dutch/Dutch BBP.

However, when looking at results from Dutch BBP only, there is a small negative correlation instead, as Figure 19 shows. This is also true for Dutch, Burstein, and Random2. This seemingly paradoxical correlation can be explained as follows. A better ranking leads to a smaller mean strength difference for these pairing systems. In an optimal ranking, each player is in her correct score group, together with players of similar strength, so the mean strength difference will be low. However, in a suboptimal ranking, some players are in a score group that does not reflect their strength. Therefore, these players are either stronger or weaker compared to the other players in their score group, which results in higher mean strength difference.

Figure 20 shows empirical evidence for this effect: the pairing in round one is always the same, but paradoxical match results due to randomness lead to different rankings,

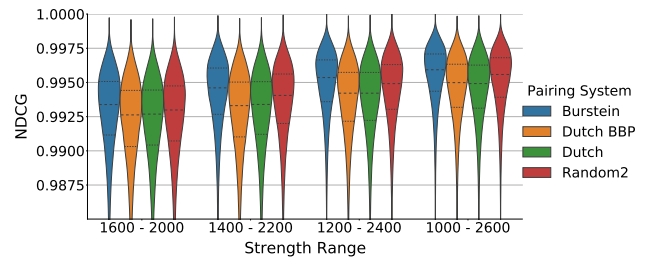


Figure 17: Ranking quality measured by the normalized discounted cumulative gain (NDCG).

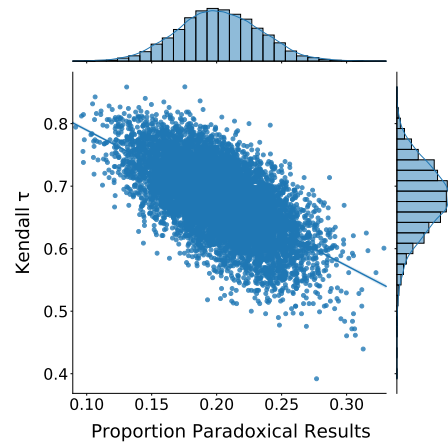


Figure 18: Correlation between paradoxical results and normalized Kendall τ after seven rounds, paired with Dutch BBP.

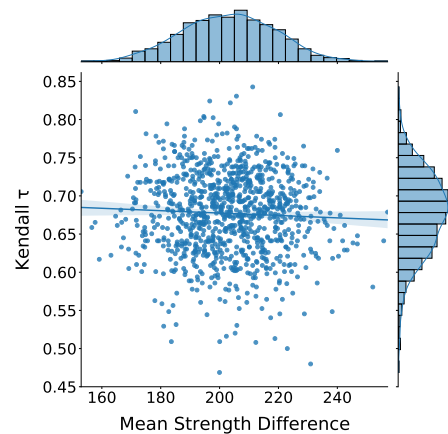


Figure 19: Correlation between mean strength difference and normalized Kendall τ after seven rounds, paired with Dutch BBP.

which then determine the mean strength difference in round two. In our experiment, the same single randomly paired first round was played 10 000 times. Each time the ranking quality after round one and the mean strength difference of the Dutch BBP pairing for round two was recorded. For the

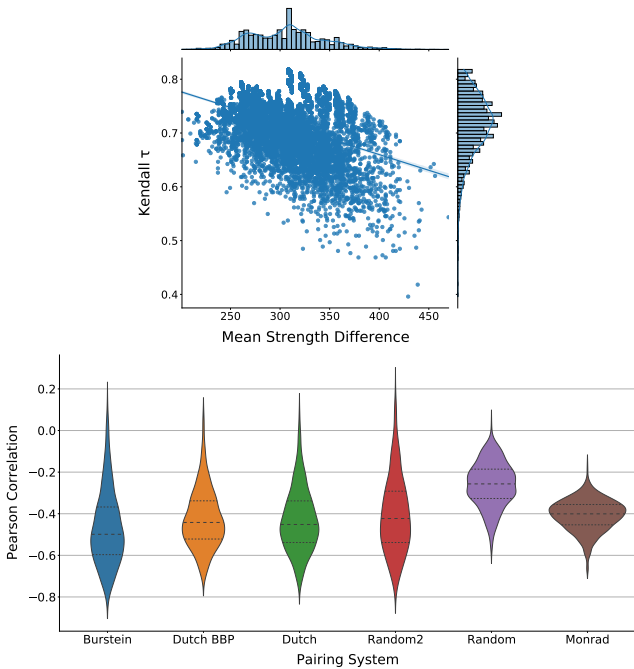


Figure 20: The scatter plot (top) shows the correlation between normalized Kendall τ and mean strength difference. The violin plot (bottom) shows the distribution of Pearson correlation coefficients if that experiment is repeated for 1 000 different tournaments, whose first round was simulated 1 000 times.

analysis, we use the *Pearson correlation coefficient* that is a standard measure for the linear dependence between two variables. In our case, a negative Pearson correlation coefficient indicates that on average, a higher Kendall τ is observed together with a lower mean strength difference.

The problem with mean strength difference is that it does not take into account whether a low mean strength difference was the result of a pairing system’s choice or due to unfavorable score groups. This can be avoided by taking the maximum possible strength difference into account. For this, we define the *normalized strength difference* as the total strength difference divided by the maximum possible total strength difference.

For computing the normalized strength difference it is essential to calculate the maximum possible strength difference. For this, we again use our maximum weight matching approach, but this time with a pairing system that maximizes strength difference. In particular, we use a modification of our Burstein edge weights $w(p_i, p_j)$ where we set $\pi(p_i, p_j) := |\text{str}(p_i) - \text{str}(p_j)|$. Remember that $\text{str}(p_i)$ and $\text{str}(p_j)$ are the true strength values of players i and j , respectively. Of course, this new pairing system requires knowledge of all true player strengths, so it cannot be used in realistic settings. We only use it as an analytical tool.

Figure 21 compares the normalized strength difference for Dutch BBP and for our maximum weight matching implementation of Burstein, which clearly beats Dutch BBP in

ranking quality. Firstly, this figure shows a positive corre-

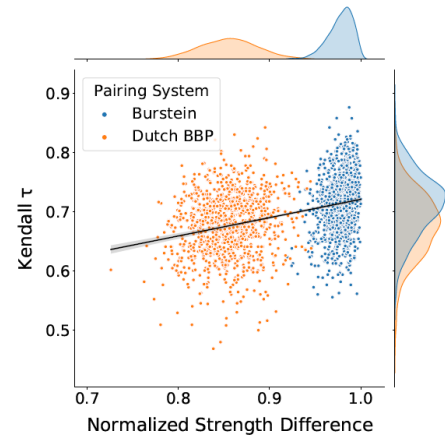


Figure 21: Correlation between ranking quality and normalized strength difference for Burstein and Dutch BBP after seven rounds.

lation between normalized strength difference and normalized Kendall τ for each of Burstein and Dutch BBP after seven rounds. Simulations with each of Dutch, Random2, Random, and Monrad also indicate a similar positive correlation. Secondly, Figure 21 also demonstrates the positive correlation between the normalized strength difference and the ranking quality across pairing systems. In particular, Burstein clearly beats Dutch BBP in normalized strength difference and also in ranking quality. This correlation is true in general: considering all pairing systems, exactly the ones with a high normalized strength difference (Burstein, Random2, Dutch, Dutch BBP) lead to a good ranking quality, while the ones with medium and low normalized strength difference (Random and Monrad) lead to medium and low ranking quality.

To summarize, our flexible maximum weight matching model enabled us to detect the exact reason why certain pairing systems produce better rankings. Our surprising finding is that even though at first sight, a high mean strength difference seems to be the pivotal factor, it is actually a high normalized strength difference that results in a better ranking quality. This discovery might help designing better pairing systems in the future.

Fairness

Here we present additional simulation results that measure the achieved fairness, i.e., results regarding the compliance with the quality criteria (Q1) and (Q2).

Number of Float Pairs

We consider the obtained number of float pairs for different strength ranges and different strength range sizes. Figures 22 and 23 show that we get consistent results for different strength ranges and different strength range sizes. Burstein has by far the lowest number of float pairs, but also Random2 and Dutch perform slightly better than Dutch BBP.

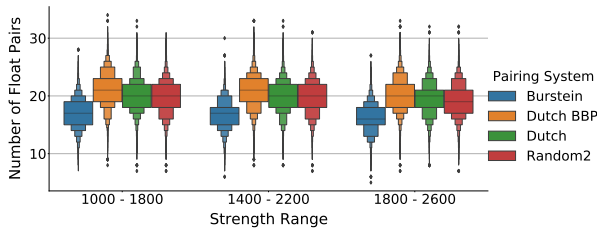


Figure 22: Number of float pairs for different strength ranges.

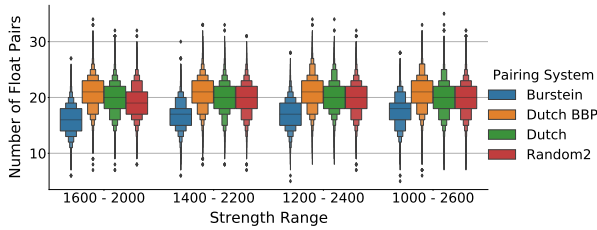


Figure 23: Number of float pairs for different strength range sizes.

Figure 24 shows a direct comparison of the obtained number of float pairs for Burstein and Dutch BBP for different numbers of players and different tournament lengths.

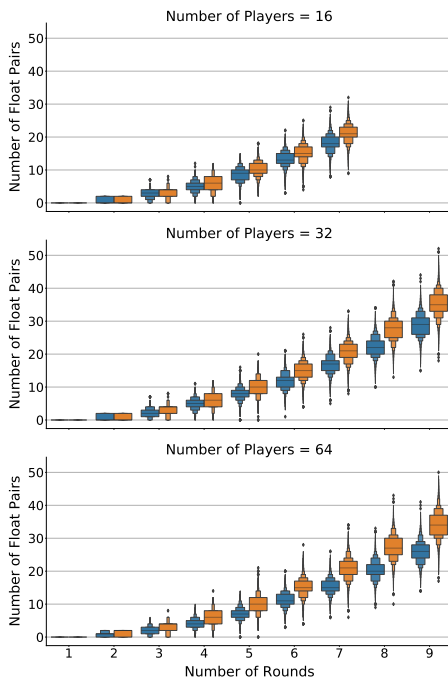


Figure 24: Number of float pairs for different tournament sizes and lengths. The results for Burstein are shown in blue, results for Dutch BBP in orange.

Also here we consistently get that Burstein achieves much fewer float pairs than Dutch BBP.

Absolute Color Difference

The measured absolute color difference increases slightly with the number of rounds and also with the number of players, as Figure 25 shows.

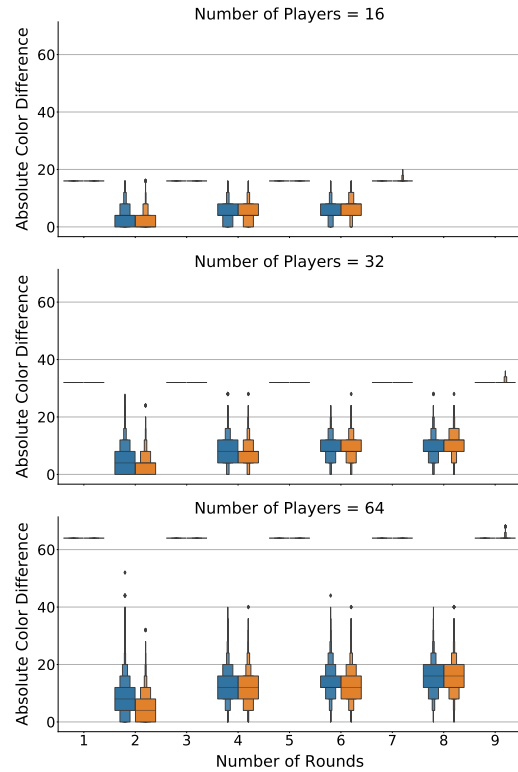


Figure 25: Absolute color difference in rounds 1-9, 16-64 players with strength range 1400-2200. Results for Burstein are shown in blue, Dutch BBP results are shown in orange.

Note that in every odd round, the absolute color difference must be at least n , which can also be seen. All investigated pairing systems almost always meet this lower bound for odd rounds. Interestingly, Dutch BBP seems to perform slightly better in tournaments with at most 4 rounds compared to Burstein, but this tiny advantage vanishes for at least six rounds. We get similar results when comparing with Random2, Dutch, Random, and Monrad.

Lower Maximum Allowed Color Difference

So far, for all our experiments we assumed that the maximum allowed color difference β equals 2, i.e., the difference of the number of matches played with white pieces and the number of matches played with black pieces is at most 2. This is in line with the official FIDE rules. However, due to the flexibility of our maximum weight matching approach, we can easily enforce an even stronger color difference constraint and observe the impact on the obtained ranking quality and the number of float pairs.

Interestingly, as Figure 26 shows, the obtained ranking quality is almost the same even if we look at the extreme case of setting $\beta = 0.1$, which is equivalent to enforcing an

alternating black-white sequence for all players. Notice that setting β to anything in the interval $(0, 0.5]$ implies that the absolute color difference is 0 for all even rounds and n for all odd rounds.

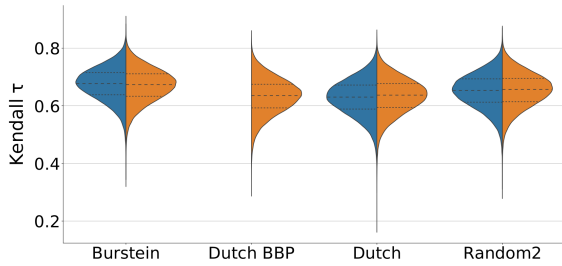


Figure 26: Ranking quality measured by normalized Kendall τ . Results for $\beta = 0.1$ are shown in blue, results for $\beta = 2$ in orange.

Naturally, the high ranking quality for a much more restricted β comes at a cost, which can be measured in the number of float pairs. The obtained number of float pairs is influenced by the maximum allowed color difference β , because for higher β it is easier to fulfill criterion (Q1), i.e., to find suitable matches within the corresponding score group. In our experiments we investigate the increase in the number of float pairs when we assume the extreme case of $\beta = 0.1$. Figure 27 shows that the number of float pairs increases for all pairing systems, compared to the case with $\beta = 2$. However, the increase is only moderate. This result offers a novel trade-off for tournament organizers: when using the MWM engine, they have the choice between keeping the number of floaters down at the cost of a standard color difference, as advised by FIDE, or they opt for slightly more float pairs, but can guarantee an alternating white-black color assignment for each player. The ranking quality is equally high in both variants.

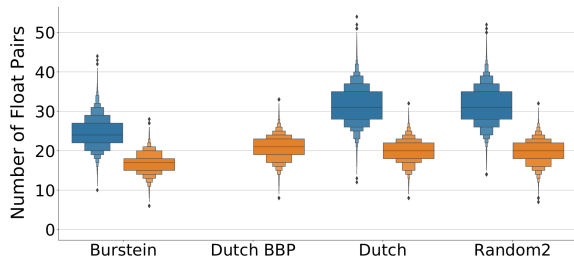


Figure 27: Number of float pairs for 7 rounds. Results for $\beta = 0.1$ are shown in blue, results for $\beta = 2$ in orange.