

TRANSPORTATION SCIENCE

Modeling Crew Itineraries and Delays in the National Air Transportation System

Journal:	<i>Transportation Science</i>
Manuscript ID	Draft
Manuscript Type:	Original Manuscript
Date Submitted by the Author:	n/a
Complete List of Authors:	Wei, Keji; Dartmouth College, thayer school of engineering Vaze, Vikrant; Dartmouth College, Department of Engineering Sciences
Keywords/Area of Expertise:	robust crew pairing, delay and disruption propagation, parameter estimation and validation

SCHOLARONE™
Manuscripts

Only

Modeling Crew Itineraries and Delays in the National Air Transportation System

Keji Wei, Vikrant Vaze

Thayer School of Engineering, Dartmouth College, Hanover, New Hampshire 03755

keji.wei.th@dartmouth.edu, vikrant.s.vaze@dartmouth.edu

Abstract

We propose, optimize and validate a methodological framework for estimating the extent of crew-propagated delays and disruptions. We identify the factors that influence the extent of crew-propagated delays and disruptions, and incorporate them into a robust crew scheduling model. We develop a heuristic solution approach for solving the inverse of the robust crew scheduling problem to generate crew schedules that are similar to real-world airline crew scheduling samples within a reasonable computational time. We develop a sequence of exact and heuristic techniques to quickly solve the forward problem within a provably small optimality gap for network sizes that are among the largest in robust crew scheduling literature. Computational results using four large real-world airline networks demonstrate that the crew schedules produced by our approach generate propagation patterns similar to those observed in real world. Extensive out-of-sample validation tests indicate that parameters calibrated for one network perform reasonably well for other networks. We provide new insights into the perceived tradeoff between planned costs and delays costs as reflected by actual airline crew schedules. Finally, we present a general approach to estimate crew-propagated delays and disruptions for any given network using our methodological framework under a variety of data availability scenarios.

Keywords: robust crew pairing; delay and disruption propagation; parameter estimation and validation.

1 Motivation

Flight delays and disruptions cost tens of billions of dollars annually to the world economy. The total cost of flight delays in the U.S. in 2007 was especially large, estimated to be approximately \$31.2 billion (Ball et al. 2010). Even though the data from the Bureau of Transportation Statistics (BTS) shows that the level of flight delays has somewhat reduced since its 2007 peak, delay costs still represent considerable amount of resource wastage for the U.S. National Air Transportation System, resulting in profit reductions or losses to the airlines, and additional discomfort and inconvenience to the passengers. Over the last five completed calendar years, i.e., from January 1st 2011 to December 31st 2015, only 79.21% of the domestic flights in the U.S. had a delay of 15 minutes or less (BTS, 2016). During the same period, 1.68% of the U.S. domestic flights were cancelled.

The U.S. Department of Transportation (DOT) classifies flight delays into five main categories. They are Air Carrier Delays, Late Arriving Aircraft Delays, National Aviation System Delays, Extreme Weather Delays, and Security Delays. In particular, during the aforementioned five-year period (Jan 1st 2011 – Dec 31st 2015), Air Carrier Delays (defined as those within the control of the airline, such as those due to maintenance or crew problems, aircraft cleaning, baggage loading, fueling, etc.) accounted for nearly 32% of all the flight delays, while another 34% were attributed to Late Arriving Aircraft Delays (those due to late arrival of the previous flights using the same aircraft) and about 31% were classified as the National Aviation System Delays (defined as those due to non-extreme weather conditions, airport operations, heavy traffic volume, air traffic control, etc.) (BTS, 2016). However, this public data lacks a dedicated category for flight delays due to the propagation of upstream crew delays and disruptions. Delayed or disrupted flights may generate delays and disruptions to subsequent flights because the crew for those flights is delayed, out of position, or unable to operate the scheduled flights without violating government regulations or collective bargaining agreements (CBAs). Presently, these delays (henceforth

called as the *Crew-propagated Delays and Disruptions* in this paper) are considered to be a subset of the, rather broad, category of Air Carrier Delays. Accurate estimation of crew-propagated delays and disruptions is critical not only as a step toward fully understanding the aviation system performance, but also for informing government policy and air carrier decisions related to airline crew scheduling. In this paper, we develop an approach to estimate crew-propagated delays and disruptions, similar to aircraft-related delays and disruptions reported by DOT's Late Arriving Aircraft Delays category.

There is yet another motivation for conducting this research. Public data sources lack information about not only the crew-propagated delays and disruptions, but also the crew itineraries themselves. Many past studies related to flight delay propagation assume knowledge of aircraft connections, crew connections, and passenger connections. There is a large amount of literature focusing on airline recovery optimization (see Petersen et al. (2012), for a recent example, and Barnhart and Vaze (2015a) for a detailed review) which uses aircraft, crew and passenger schedules as inputs to their computational case studies. There is also a growing body of literature on various strategies for mitigation of airport and airspace congestion and delays (e.g., Jacquillat and Vaze, 2016; Jacquillat and Odoni, 2015; Vaze and Barnhart, 2012), which assumes the knowledge of aircraft, crew and passenger connections. However, information on only the aircraft connections is available publicly through the BTS (Bureau of Transportation Statistics) website. In terms of the passenger connections, a prior study by Barnhart, Fearing and Vaze (2014) used a statistical estimation approach to come up with estimates of passenger itineraries, passenger delays and disruptions. However, similar estimates of crew itineraries or even a validated methodology to come up with such estimates is not available. Apart from aiding in future research studies such as those mentioned above, such estimates of crew itineraries would also be beneficial for assessing the full impact of any delay mitigation strategy being considered by airlines and/or government. In this paper, similar to the work on passenger itinerary and delays estimation by Barnhart, Fearing and Vaze (2014), we develop a crew itinerary estimation methodology, and to

1
2
3 generate a database of estimated crew itineraries that will enable accurate estimation of crew-
4
5 propagated delays and disruptions consistent with their real-world values.
6
7

8
9 Note that we *do not* attempt to develop an approach to generate crew itineraries that are identical to
10
11 the real-world airline crew itineraries in *every* possible aspect. Such objective would not only be
12
13 extremely difficult to attain, but also likely cause overfitting issues based on our limited confidential
14
15 data sample of crew itineraries which we use for the estimation purposes. Instead, we develop a robust
16
17 process to generate crew itineraries that are similar to the real-world airline crew itineraries in their
18
19 potential for crew-propagated delays and disruptions. There are possibly other, non-delay related,
20
21 aspects of crew itineraries that could be of relevance for other purposes. However, our focus in this
22
23 paper is on ensuring that our process is accurate and stable in terms of the estimation of crew-
24
25 propagated delays and disruptions.
26
27
28
29

30 Crew-propagated delays and disruptions for an airline during a particular time period are strongly
31
32 dependent on the chosen set of crew pairings. A crew pairing is defined as a sequence of flight legs
33
34 covered by a crew member that follows a number of rules and regulations, and can be considered to be
35
36 the smallest self-contained unit comprising crew itineraries. Thus, estimating crew pairings is at the core
37
38 of the challenge of estimating crew-propagated delays and disruptions. A crew pairing problem is the
39
40 one of generating a set of crew pairings that covers all available flights legs. Crew pairing decisions have
41
42 a significant effect on airlines' planned and operational costs. Operational costs include delay and
43
44 disruption costs due to irregular operations and they are often amplified by the propagation of delays
45
46 and disruptions through crew connections. Thus, understanding the extent, causes and impacts of
47
48 propagation of delays and disruptions is essential for developing methods to reduce them, and to
49
50 improve the overall aviation system performance. Previous research studies on crew pairing generation
51
52 have focused on minimizing the planned crew costs, and sometimes a subset of the various components
53
54
55
56
57
58
59
60

of operational costs, but none of them have focused on quantifying and understanding the crew-propagated delays and disruptions in the real-world airline networks. Therefore, there is a need to develop approaches and models that produce crew schedules that are similar to real-world crew schedules in terms of the crew-propagated delays and disruptions.

Finally, we note that the present research project was originally motivated and funded by an aircraft delay modeling limitation faced by the U.S. Federal Aviation Administration’s (FAA) Office of Performance Analysis. The FAA analyzes and forecasts, on a monthly basis, aircraft delays at the nation’s major airports, with the objective of identifying airports with significant potential for delays months in advance, so that appropriate actions may be taken to prevent or mitigate such delays. A discrete events simulation platform developed by the FAA for this purpose models aircraft-based delay propagation using public data on aircraft itineraries (BTS, 2016), but does not account for crew-based propagation effects due to lack of crew itinerary data. This is believed to contribute to an underestimation of the propagated delays, and served as the motivation for conducting the research work presented in this paper.

1.1 Crew Pairing Optimization Problem

A crew pairing consists of a sequence of duties, where a duty is defined as the set of tasks to be performed by a crew during a given day. Duties are connected by rest periods. Each duty is made up of a set of consecutive flights with some gaps between them. These gaps are called sit times. A pairing should begin and end at the *crew base* which is usually the domicile of a crew member. Both pairings and duties are subject to various regulations and contractual restrictions. Typically, these include the following.

- The total flying time within a duty cannot exceed an upper bound. There is also an upper bound on the total elapsed time within a duty.

- There is a lower bound on the sit time which guarantees that the crews have enough time to connect between two consecutive flight legs within a duty.
- The rest time between duties should be greater than or equal to a minimum rest time which ensures that the crew is sufficiently rested between duties.
- There is typically an upper limit on the number of duties within a pairing.

In addition to these various feasibility rules, even if we ignore the operational cost considerations, crew pairings also have a highly non-linear pay structure. For a typical North American airline, the planned cost of a pairing p is the maximum of two terms: sum of the costs c_d of all its duties and a fixed fraction (ζ) of the total time away from base ($TAFB_p$). Thus the planned cost of a pairing p is given by:

$$c_p = \max \left\{ \left(\sum_d c_d \right), \zeta * TAFB_p \right\} \quad (1)$$

For each duty d in a pairing, the planned cost (c_d) is defined as the maximum of three terms, a minimum guaranteed pay (δ) per duty, flying time (fly_d) of the duty, and a fixed fraction (ε) of the duty elapsed time ($elapsed_d$). Parameters $\delta, \varepsilon, \zeta$ may vary across different carriers. Thus the planned cost of duty d can be written as

$$c_d = \max \{ \delta, fly_d, \varepsilon * elapsed_d \} \quad (2)$$

The objective of the deterministic crew pairing problem is to minimize the planned crew cost and is usually modeled as a set partitioning problem (Barnhart and Vaze, 2015b). We denote the set of flights by F and the set of pairings by P . a_{ip} is 1 if pairing p contains flight leg i , and it is 0 otherwise. x_p is a binary decision variable which equals 1 if pairing p is chosen in the crew pairing solution, and it is 0 otherwise. Then the crew pairing problem can be formulated as

$$\text{Min } \sum_{p \in P} c_p x_p$$

Subject to

$$\sum_{p \in P} a_{ip} x_p = 1, \quad \forall i \in F, \quad (3)$$

$$x_p \in \{0,1\}, \quad \forall p \in P, \quad (4)$$

1.2 Literature Review

As mentioned at the beginning of this section, from an application perspective, our work is motivated by the work of Barnhart, Fearing and Vaze (2014). Using one year of flight delays data from BTS (BTS, 2016) and a one-quarter sample of confidential passenger booking data from an airline, they estimated passenger itinerary flows and developed insights into factors that affect the performance of the U.S. National Air Transportation System from a passenger perspective. They developed a methodology to model historical travel and delays for passengers. From a methodological perspective, our work, however, is fundamentally different from that of Barnhart, Fearing and Vaze (2014). While Barnhart, Fearing and Vaze (2014) used a statistical approach to estimate passenger itineraries, estimation of crew itineraries is considerably more complicated because of the much more complex rules governing what constitutes a legal itinerary for the crew, thus making a statistical estimation approach unsuitable for our task. Also, while the number of possible passenger itineraries per day can be in thousands for a large airline, the number of legal crew itineraries is usually larger by several orders of magnitude, often making it very difficult, or impossible, to even enumerate all of them exhaustively.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Estimating delay propagation through crew connections is also much more complex than estimating the same through aircraft connections due to the more complex nature of crew work regulations than aircraft maintenance regulations (Barnhart and Vaze, 2015b). Lan, Clarke and Barnhart (2006) came up with a method for quantifying aircraft-based delay propagation. However, there is little literature on how to quantify propagation through crew connections. Several past studies on crew pairing optimization have tried to identify and capture one or more dimensions of a crew schedule that affect the extent of propagation. Broadly, these past studies can be divided into three main categories. First category of studies aims to incorporate one or more features that affect the ease of recovering the crew schedules after a disruption. Here, crew schedule recovery refers to the set of reactive measures available to an airline to bring its crew schedule back on track after a disruption and it typically includes alternatives such as delayed flight departures, crew swaps, reserve crews, flight cancellations, etc. Second category of studies aims to generate crew pairings that are difficult to get disrupted and/or have a low disruption cost. Studies in the first and second categories deal exclusively with crew schedules without capturing the relationship between crew-based and aircraft-based propagation. The third category attempts to capture this interdependence.

Studies in the first category usually focus only on one or two specific factors that improve recovery potential. For example, Shebalov and Klabjan (2006) use the objective function of maximize the number of move-up crews to achieve robustness in crew scheduling, wherein a move-up crew for a flight is a crew that is not actually assigned to that flight but can be feasibly and legally assigned to it. On the other hand, Gao, Johnson, and Smith (2009) extend the fleet purity idea proposed by Smith and Johnson (2006) to both fleet purity and crew base purity. The crew base purity idea is related to the notion of restricting the number of crew bases allowed to serve each airport in order to increase the opportunities to find a move-up crew in crew recovery. Shebalov and Klabjan (2006) as well as Gao, Smith and Johnson

(2009) capture the potential for crew swaps, which is an important dimension of crew recovery process, but do not explicitly capture the extent of delay propagation through crew connections.

Studies in the second category apply a variety of robust planning approaches to airline crew scheduling. Yen and Birge (2006) develop a two-stage stochastic programming-based crew scheduling model that implements a simplified recovery model for the second stage. Schaefer et al. (2005) adjust the cost of each crew pairing to include a combination of planned costs and a linear approximation of delay costs, wherein the delay cost approximation function is fine-tuned based on a discrete events simulation software named SimAir (Rosenberger et al., 2002). The delay cost is assumed to be a function of four attributes, namely, 1) sit time between consecutive flights within a duty, 2) rest time between consecutive duties, 3) total flying time in a duty, and 4) total elapsed time in a duty. The rationale behind these choices is that the potential for the propagation of delays and disruptions is greater when crew's sit and rest times are too short and when the per-duty flying and elapsed times are too long. Yen and Birge (2006) as well as Schaefer et al. (2005) account for the differences in the delay propagation potential of different crew pairings, but neither captures recovery actions such as crew swaps.

Because delays can propagate due to both late arriving or unavailable aircraft and late arriving or unavailable crew, there are many interdependencies between the effects of aircraft schedules and crew schedules on the propagation of delays and disruptions through the overall flight network. Aircraft scheduling and crew scheduling stages of the airline schedule planning process are conventionally solved in a sequential manner. However, recognizing the interdependencies between the two stages, in terms of both planned costs and delay propagation potential, some recent studies (such as Dunbar, Froyland and Wu (2012), and Weide, Ryan and Ehrgott (2009)) have developed integrated robust optimization models. Weide, Ryan and Ehrgott (2009) attempt to increase the buffer in crew connection times when crew changes aircraft. Dunbar, Froyland and Wu (2012) focus explicitly on minimizing delay

costs while ignoring planned costs. Cacchiani and Salazar-González (2016) solve an integrated fleet assignment, aircraft routing and crew pairing problem with a weighted average objective function that incorporates robustness solely as measured by the number of aircraft changes between successive flights in a crew itinerary. Mercier, Cordeau and Soumis (2005) also incorporate aircraft changes by the crew as a measure of robustness in their integrated aircraft routing and crew scheduling model. Ehrgott and Ryan (2002), and Tam et al. (2011) describe and evaluate a bi-criteria optimization approach to balance the planned crew costs and a single robustness measure which penalizes crew connections between two successive flights covered by the same crew but different aircraft wherein the connection time is small relative to the expected flight delays. Studies in this category usually emphasize on aircraft changes and crew sit times but do not focus much on crew recovery potential, crew rest times, per-duty flying times, or per-duty elapsed times.

In summary, past research studies in airline crew scheduling have identified various features of airline crew schedules that affect the crew-propagated delays and disruptions. However, no prior study has combined these different features into a single optimization model. Additionally, while some past studies, such as Yen and Birge (2006), have attempted to incorporate the actual delay costs into the crew pairing optimization models, these models have been highly simplified due to computational tractability issues associated with more detailed models. Finally, and most importantly, all the aforementioned studies have focused, implicitly or explicitly, on finding an “optimal” crew schedule with respect to a known optimization formulation. The problem that we solve in this paper can be thought of as the inverse of this problem. Given an actual crew pairing sample, our goal is to reverse engineer the process used, and the problem solved, by the airlines to generate the crew pairings that the airlines actually used. This will enable us to generate crew pairings that are similar to the crew pairings used by the airlines for other airlines, and/or other aircraft families, and/or other time periods than the ones for which an actual crew pairing data sample is available. It is common knowledge that the major airlines

typically use advanced optimization solvers to generate crew pairings. Furthermore, in addition to minimizing planned costs, most airlines are known to directly or indirectly pay attention to the secondary goal of reducing of delay and disruption costs as well. However, the exact models and algorithms used by a particular airline for crew pairing optimization are proprietary. Therefore, in this paper we reverse engineer airlines’ crew pairing generation process with the objective of being able to generate crew pairings that are similar to the actual airline-generated crew pairings in terms of their potential for crew-propagated delays and disruptions.

1.3 Contributions and Outline

The research presented in this paper makes five main contributions. First, we propose a comprehensive crew-pairing optimization formulation that minimizes the combination of planned costs and various features of crew schedules that make the crew schedules vulnerable to propagation of delays and disruptions. We combine six such features, some of which have been proposed individually in some prior studies, while others are new. Second, we solve this model provably to near-optimality by combining known ideas such as branch-and-bound and delayed column generation as well as a sequence of new heuristic ideas developed by us. The sizes of the networks in the problems solved by us far exceed those solved in past literature studies on robust or recoverable crew pairing optimization. Third, we embed this crew pairing generation problem in an upper-level calibration framework wherein a parameterized crew pairing optimization problem is solved repeatedly by varying the parameters until the resulting crew pairings are *similar* to those used by the airlines. This upper-level calibration problem represents the inverse crew pairing generation problem mentioned in Section 1.2. Fourth, we develop a local-search heuristic for solving the upper-level calibration problem. Our algorithm is motivated by that of Schaefer et al. (2005) and borrows some features from theirs. Ours, however, is the first study to formulate and solve this inverse crew pairing generation problem. Finally, we generate and validate

crew pairing solutions that are similar to those used by the airlines in the real-world in terms of their potential for crew-propagated delays and disruptions. The out-of-sample testing results demonstrate the accuracy and stability of our modeling framework and algorithms. One important conclusion is that the ratio between the planned crew cost and approximate delay costs is found to be stable across airlines and aircraft types.

The rest of the paper is structured as follows. Section 2 describes our overall modeling approach and problem formulation. Section 3 details the solution approach including the various exact algorithms and heuristic ideas developed by us to solve this challenging problem. Section 4 describes our computational case studies in terms of data and pre-processing, and presents evidence of the computational tractability of our approach. Section 5 describes the calibration and validation results obtained from our series of computational experiments. Section 6 further validates our results in terms of the crew-propagated delays and disruptions, and also describes how to use our results for estimating crew-propagated delays and disruptions for any given network. Finally, Section 6 discusses the main conclusions and the directions for future research.

2 Modeling Framework

Our objective is to generate crew itineraries that are similar to the real-world crew itineraries as measured by the extent of crew-propagated delays and disruptions. Therefore, we first need to develop an appropriate similarity metric for comparing two crew pairing solutions with each other, for any given flight network. Defining similarity directly based on the actual costs of propagated delays and disruptions is problematic for multiple reasons. Propagated delays and disruptions depend on not only the crew schedules but also the underlying root (i.e., non-propagated) delays and disruptions, as well as the operational recovery actions used by the airline. Exact recovery actions used by the airlines are typically not public knowledge. Hence these are difficult to model accurately. Moreover, while planning

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

the crew schedule, the airline itself is unaware of the exact set of root delays and disruptions that it will face on a given day of operations. For these reasons, accurate calculation of the costs of crew-propagated delays and disruptions is impossible. Instead, we measure the similarity between crew pairing solutions in terms of their *potential* for crew-propagated delays and disruptions. As explained in Section 1.2, the potential for crew-propagated delays and disruptions is a function of various features of a crew schedule. In Section 2.1, we classify these features into four categories and select six specific representative features for inclusion in our model. Then, in Section 2.2, we provide the mathematical formulation for the crew pairing optimization problem that we use as the basis of our calibration framework. Finally, in Section 2.3, we give a clear mathematical formulation for our calibration problem of minimizing the distance (i.e., maximizing the similarity) between the estimated and actual crew pairing solutions in terms of their potential for crew-propagated delays and disruptions as quantified by these six features.

2.1 Representative Features

In the absence of sufficient schedule buffers and recovery opportunities, delays and disruptions propagate to downstream flights leading to additional operating costs. Therefore, besides planned crew costs (i.e., crew salaries) that we introduced in Section 1.1, airlines often consider some of these buffers and/or recovery opportunities during crew scheduling to reduce these extra operational costs. There are a variety of mechanisms through which delays and disruptions affect downstream flights, the simplest is when the sit time buffer (defined as scheduled sit time minus minimum required sit time) or the rest time buffer (defined as scheduled rest time minus minimum required rest time) between two consecutive flights in a crew pairing is less than the arrival delay of the first flight. This leads to a propagated delay for the second flight unless some recovery action, such as a crew swap, is able to prevent it. Thus the sit time buffers, rest time buffers, and crew recovery potential clearly affect crew-

propagated delays and disruptions. However, if these two flights are scheduled to be operated by the same aircraft, then this delay to the second flight would be unavoidable due to aircraft-based propagation, irrespective of whether the crew is on-time. Note that, as per the DOT classification, delay propagation in such situations is classified as aircraft-based propagation causing the crew-propagated delays and disruptions to be counted as zero to avoid double-counting. Thus, whether or not the crew travels with the same aircraft affects what is considered as the crew-propagated delays and disruptions. Finally, if flight delays result in violation of any of the crew duty regulations and/or CBA rules, such as the total flying time in a duty or total elapsed time in a duty, for a later flight then the later flight becomes inoperable by its scheduled crew, resulting in either a flight cancelation or some crew recovery action. Thus, the available buffers (defined as maximum allowable value minus scheduled value) in total flying time or total elapsed time in a duty also affect crew-propagated delays and disruptions. This discussion motivates our classification of features affecting crew-propagated delays and disruptions as well as our choice of the representative features.

We divide the features affecting crew-propagated delays and disruptions into four categories to ensure that no important category of factors is missing from our model. We name these four categories as Aircraft Change, Push-Back, Crew Legality, and Crew Swaps. This categorization serves two main purposes: (1) it highlights the variety of ways in which delays and disruptions can propagate through crew connections; and (2) it facilitates any future revisions or extensions of the feature-set based on the methodologies we have developed.

2.1.1 Aircraft Change

We first motivate this category with an example. Consider two flights, Flight 1 and Flight 2, scheduled to be operated consecutively by the same crew within the same duty. If Flight 1 and Flight 2 are scheduled to be operated by the same aircraft as well, then irrespective of whether Flight 1's is delayed or not, by

the time the aircraft is ready to operate Flight 2, the crew will typically be ready as well. Thus, depending on the relative values of the arrival delay of Flight 1 and the buffer in the aircraft turn-around time between Flight 1 and Flight 2, there will be either no delay propagation or there will be some delay propagation attributed to the late arriving aircraft. However, no *crew-propagated* delay or disruption will occur. On the other hand, if Flight 1 and Flight 2 are scheduled to be operated by different aircraft, then to avoid delay propagation from Flight 1 to Flight 2, the crew on Flight 1 will need to exit that aircraft, reach the aircraft scheduled to operate Flight 2 and get ready to start operating it before the scheduled departure time of Flight 2. In this scenario, depending on the arrival delay to Flight 1 and the buffer in crew connection time between Flight 1 and Flight 2, there could be delay propagation through crew connection. Therefore, if the crew needs to change aircraft between its consecutive flights within a duty, there is a potential for crew-propagated delays and disruptions. Hence whether or not the crew stays with the aircraft during a connection between two flights in the same duty is an important factor affecting crew-propagated delays and disruptions (U.S. G.A.O, 2008). Specifically, the number of times a crew switches aircraft within a duty in the pairing is included as one of the representative features of the potential for crew-propagated delays and disruptions in our model.

2.1.2 Push-Back

When a flight’s arrival is delayed, and the same crew within the same duty is scheduled to operate a subsequent flight, which is not scheduled to be operated by the same aircraft, a simple policy is to delay the subsequent flight until its scheduled crew is ready to operate it, regardless how severe the delay is. We call this as the push-back strategy (Rosenberger et al. 2002). Similarly, when the arrival of the last flight in a crew duty (which is not the last duty in the crew pairing) is delayed, push-back strategy may be used to delay the departure of the first flight in the crew’s next duty regardless of how severe the delay is and irrespective of whether or not the same aircraft is scheduled to operate the two flights. Note that,

under the push-back strategy, delay propagates through the crew connection when the buffer in crew connection time or crew rest time exceeds the arrival delay of the first flight. Thus, crew sit time buffer (between flights not scheduled to be operated by the same aircraft) and crew rest time buffer are important factors affecting crew-propagated delays and disruptions, and hence both of these are used in our model as features representative of the potential for crew-propagated delays and disruptions.

2.1.3 Crew Legality

When developing crew schedules, airlines must adhere to FAA crew safety regulations and CBAs regarding the maximum flying time in a duty and maximum elapsed time in a duty. For example, if FAA regulations limit a pilot to 8 hours of flying time during a duty, then no airline is permitted to schedule a pilot to fly for more than 8 hours during a single duty. If the scheduled flying time is exactly 8 hours or just under 8 hours, even a small delay to one of the earlier flights in the pilot's duty could cause the actual flying time to exceed 8 hours, thus disallowing the pilot to operate the last flight in his/her originally scheduled duty until the completion of a rest period, either leading to a flight schedule disruption such as cancellation or large delay, or triggering a crew recovery action such as a crew swap or use of reserve crew. Note that, under this scenario, a crew-propagated delay or disruption happens even when the crew connection time buffer is large and/or the crew is not scheduled to change aircraft between flights. A similar argument holds when the scheduled elapsed time in a duty is equal to or just under the maximum allowable duty elapsed time. Thus, the buffer in flying time in a crew duty and the buffer in elapsed time in a crew duty are important factors affecting crew-propagated delays and disruptions, and hence both of these are used in our model as representative features.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

2.1.4 Crew Swaps

As mentioned in Section 1.2, crew schedule recovery refers to the set of reactive actions available to an airline to bring its crew schedule back on track after a disruption and it typically includes alternatives such as delayed flight departures, crew swaps, reserve crews, flight cancellations, etc. While delaying flight departures is the default alternative, under significant disruption events, it can result in very large and expensive delays. A flight cancellation cannot be done in isolation and typically it leads to cancellation of one or more other flights scheduled to be operated by the same aircraft and requires extensive amounts of passenger rebooking. Use of reserve crews is constrained by the availability of the reserve crews and is typically an expensive strategy as well. A crew swap involves assigning a late arriving crew to operate a flight with a later departure time than its originally scheduled flight and instead using a different crew to operate the earlier flight. For being able to swap crews, the two swapped pairings must be from the same crew base, either crew must be qualified to operate the subsequent flights in either pairing, and finish the respective pairings on the same day (Shebalov and Klabjan 2006).

Compared to other crew recovery actions such as cancellations or reserve crews, swaps are typically a less expensive strategy economically, and therefore airlines find it beneficial to increase the crew-swapping opportunities. In order to increase crew-swapping opportunities, Gao, Johnson and Smith (2009) introduce a new concept, called crew base purity, which involves restricting the number of crew bases serving each airport. They found that improving the crew base purity can significantly increase crew-swapping opportunities and thus reduce the cost of crew recovery. They describe the idea of using adjacency graph to quantify the extent of crew swapping potential. In an adjacency graph, airports are represented by nodes and the existence of an arc implies that there is at least one flight connecting the two nodes. For a specific airline’s network, distance between two airports in an adjacency graph is

defined as the minimum number of arcs that need to be traversed to go from one airport to the other. Crews serving airports that are more distant from the crew base, lead to fewer crew swapping opportunities and thus lower recovery potential. In our model, the number of times a crew visits an airport which is at a distance of 2 or more from its base is used as a feature representative of the crew swapping potential and hence representative of the potential for crew-propagated delays and disruptions. Other features indicating the crew recovery potential, such as, the number of reserve crew members available at various airports, could also be potentially included as representative features. However, we did not include them because of lack of data on the availability of reserve crews in our dataset.

2.2 Robust Crew Pairing Formulation

The six representative features identified in Section 2.1 were integrated into a mathematical model that generates crew pairings that are similar to the real-world airline crew pairings. Our model formulation is motivated by the work of Schaefer et al. (2005). Schaefer et al. (2005) used a penalty method for quantifying and maximizing the robustness of a crew schedule. They optimize the total expected operational cost of a crew pairing solution, which is defined as the sum of the planned cost and a linear function of four specific attributes of each crew pairing serving as proxy measures of its robustness. They assume that the aircraft are always available and hence no delay propagates through the aircraft connections. Also, the recovery method is assumed to be push-back only. Finally, they assume that the operational cost of a crew pairing solution is the sum of the operational costs of individual chosen pairings, and that interaction between crew pairings does not have an effect on the operational costs. We retain this last assumption, but partially relax the first and second assumption as follows. Similar to the four attributes chosen by Schaefer et al. (2005), we also include, in our representative features set, the scheduled sit time when crew changes planes, the scheduled rest time between duties, flying time in

a duty, and elapsed time in a duty. Additionally, we also include as one of our representative features, the number of times a crew changes aircraft between successive flights within a duty. This provides a partial proxy for the additional delay that may result from the late arriving aircraft. Similarly, we also include the crew base purity, as measured by the number of times the crew arrives at an airport whose distance from the base is 2 or greater in the adjacency graph. We define these as the instances of violation of the crew base purity. Crew base purity provides a proxy for the crew recovery potential through crew swaps, as described in Section 2.1. Thus, we used the following six features.

Feature 1: Scheduled sit time when crew changes aircraft.

Feature2: Scheduled rest time between duties.

Feature 3: Flying time in a duty.

Feature 4: Elapsed time in a duty.

Feature 5: Number of crew base purity violations.

Feature 6: Number of aircraft changes by the crew within a duty.

Our method of incorporating these features into the crew pairing optimization model is an extension of the penalty method developed by Schaefer et al. (2005). For any pairing p , let c_p be its planned cost, and f_p be the penalty cost as a function of feature i . Then the total cost (\bar{c}_p) of pairing p is defined as

$$\bar{c}_p = c_p + \sum_{i=1}^6 f_p(i) \tag{5}$$

For Features 1 and 2, as the value approaches the largest acceptable value, the potential for crew-propagated delays and disruptions increases. For instance, the FAA requires that a crew must receive

rest if the crew has already flown for 8 hours in a duty. As the scheduled flying time in a duty increases, the potential for this pairing becoming illegal during operation goes up because of increased likelihood of violation of this rule. Similarly, for Features 3 and 4, as the value approaches the smallest acceptable value, the potential for crew-propagated delays and disruptions goes up. For feature i , let δ_i for feature i denote the relevant bound, that is, upper bound for Features 1 and 2 and lower bound for Features 3 and 4. For example, for Feature 4, δ_4 is the minimum rest time as allowed by the FAA regulations and the CBAs. For δ_4 of 10 hours, rest periods shorter than 10 hours in length are not permitted. Let $Count(i, p)$ be the number of times that feature i occurs in pairing p , and let $V_{i,p}^j$ be the value of the j^{th} occurrence of feature i in pairing p . For instance, if pairing p has three duties with elapsed times of lengths 10, 12, 5 hours respectively, then $Count(2, p) = 3$, $V_{2,p}^1 = 10$, $V_{2,p}^2 = 12$, and $V_{2,p}^3 = 5$. We use parameters α_i to represent the maximum penalty, and β_i to represent the slope in feature i 's penalty function. So, for the first four features, the function $f_p(i)$ is defined as:

$$f_p(i) = \sum_{j=1}^{Count(i,p)} \max(\alpha_i - \beta_i |V_{i,p}^j - \delta_i|, 0), \forall i \in \{1, 2, 3, 4\} \quad (6)$$

The form of function $f_p(i)$ described by Equation (6) is similar to that used by Schaefer et al. (2005). It assumes that $f_p(i)$ is additive across the effects of all occurrences of feature i in pairing p . Also, it assumes that, within a range, the effect of the value of the feature in each occurrence is linear and increases as the value of the feature gets increasingly closer to the relevant bound δ_i . At the bound, the effect has the maximum value α_i , because this leaves zero buffer in case of any prior delays or disruptions, and hence creates the maximum potential for crew-propagated delays and disruptions. Farthest away from the bound (i.e., at a distance of $\frac{\alpha_i}{\beta_i}$), the effect is zero. This is because large enough buffers almost fully eliminate any potential for crew-propagated delays and disruptions.

For defining the penalty function for Feature 5, we observe that if most airports are directly connected to the crew base, the airline has a greater potential for an efficient recovery from disruptions by finding a move-up crew. Also, as for Feature 6, we observe that if most crews are to stay with the aircraft, then most of the delay propagation would be attributed to late arriving aircraft, rather than being counted as part of the crew-propagated delays and disruptions. So we penalize the number of occurrences of crew changing aircraft and the number of occurrences of crew base purity violations. Let parameters γ_i , $i \in \{5,6\}$, denote the penalty weights for Features 5 and 6. With $Count(i,p)$ defined the same way as that for Features 1 through 4, the function $f_p(i)$ for Features 5 and 6 is defined as:

$$f_p(i) = \gamma_i * Count(i,p), \forall i \in \{5,6\} \tag{7}$$

Note that this expression is simpler than the one for the effects of Features 1 through 4 because the number of aircraft changes and the number of crew base purity violations directly have an effect on the potential for crew-propagated delays and disruptions, as against the effects of Features 1 through 4 which depend on the difference between the feature value and a relevant bound. This results in a crew pairing optimization model given by

$$Min \sum_{p \in P} \left(c_p + \sum_{i=1}^6 f_p(i) \right) x_p \tag{8}$$

Subject to

$$\sum_{p \in P} a_{ip} x_p = 1, \quad \forall i \in F \tag{9}$$

$$x_p \in \{0,1\}, \quad \forall p \in P \tag{10}$$

2.3 Calibration Framework

There are several ways of conceptualizing our calibration problem. Given the optimization model (8-10), we could consider the calibration problem as one of estimating the parameters $\alpha_i, i \in \{1,2,3,4\}, \beta_i, i \in \{1,2,3,4\}$ and $\gamma_i, i \in \{5,6\}$. Thus it is what is sometimes called as an inverse optimization problem. While an inverse linear optimization problem has been shown to be another linear optimization problem (Ahuja and Orlin, 2001), and hence is easy to solve, similar results do not exist for an inverse integer optimization problem (IIOP). Recently, Lamperski and Schaefer (2015) developed an approach to formulate the IIOP as an integer optimization problem with exponentially larger size. Others have proposed heuristic approaches for solving variants of the IIOP (see Duan and Wang, 2011; and Wang 2013; for recent examples). However, these are computationally intensive and deal with only small-sized problems. A crew pairing optimization problem, on the other hand, typically consists of millions of (or more) variables, and is typically solved using complex, resource intensive algorithms such as branch-and-price (Barnhart et al., 1998). Therefore, solving an inverse version of such a problem is extremely challenging for realistic problem sizes and no existing study has addressed this challenge successfully.

Alternatively, the calibration problem could also be considered to be a type of supervised machine learning problem where the goal is to generate crew pairing solutions similar to those in the labeled training data by learning the parameters $\alpha_i, i \in \{1,2,3,4\}, \beta_i, i \in \{1,2,3,4\}$ and $\gamma_i, i \in \{5,6\}$. This labeled training data, represented by a set of crew pairings, is in the form of a set of sequences of flights where each sequence is operated by the same crew. This is in a non-standard structure for supervised machine-learning, and the mechanism through which the parameters affect the labels is also very complicated. Thus, none of the typical supervised learning approaches, such as support vector machines or neural networks, to name a few, are directly applicable.

This discussion suggests that our calibration problem has several unique attributes, and is computationally much more expensive compared with what existing methods have been shown to be able to solve. Therefore, we propose a new mathematical framework and a solution heuristic for solving this calibration problem. First, in this section we describe the framework and the relevant mathematical notation. Then, in Section 3, we describe the solution heuristic.

Let us denote the set of parameters by $PARAMS$. Thus, $PARAMS = \{\alpha_1, \alpha_2, \alpha_3, \alpha_4, \beta_1, \beta_2, \beta_3, \beta_4, \gamma_5, \gamma_6\}$. Let $\hat{x}(PARAMS)$ be the crew pairing solution generated by solving the optimization model (8-10) for a given set of $PARAMS$ values, and let x be the real-world airline's scheduled crew pairing solution in our sample data. So, $\hat{x}(PARAMS) = \operatorname{argmin}_{\sum_{p \in P} a_{ip} x_p = 1 \forall i \in F; x_p \in \{0,1\} \forall p \in P} \sum_{p \in P} [(c_p + \sum_{i=1}^6 f_p(i)) \cdot x_p]$.

Also, let $F^{\hat{x}}(i) = \sum_{p \in P} f_p(i) \hat{x}_p$ and $F^x(i) = \sum_{p \in P} f_p(i) x_p$ be the values of the i^{th} component of the penalty function corresponding to the crew pairing solutions \hat{x} and x respectively. Then the calibration problem is formulated as follows:

$$\min \sum_{i=1}^6 |F^{\hat{x}}(i) - F^x(i)| \tag{11}$$

Subject to

$$\hat{x} = \operatorname{argmin}_{\sum_{p \in P} a_{ip} x_p = 1 \forall i \in F; x_p \in \{0,1\} \forall p \in P} \sum_{p \in P} \left(c_p + \sum_{i=1}^6 f_p(i) \right) x_p \tag{12}$$

$$f_p(i) = \sum_{j=1}^{Count(i,p)} \max \left(\alpha_i - \beta_i |V_{i,p}^j - \delta_i|, 0 \right), i \in \{1,2,3,4\} \tag{13}$$

$$f_p(i) = \gamma_i * Count(i,p), i \in \{5,6\} \tag{14}$$

Note that this formulation minimizes the L1 norm of the difference between $F^{\hat{x}}(i)$ and $F^x(i)$. Alternatively, we could consider minimizing other norms (such as L2 norm) as well. Our computational experiments with L1 and L2 norms showed that these two alternative formulations did not lead to any significant changes in our results.

3 Solution Approach

In order to generate crew pairings that are similar to those scheduled by the airline, we need to solve the calibration optimization problem represented by (11-14). The similarity or closeness between the two crew pairing solutions provides a measure of success of the calibration process. However, in order to truly assess the stability of this approach, we need to perform out-of-sample testing. As described in Section 4, we use one set of sample data to calibrate the parameters and then use the same calibrated values with another set of sample data (from a different airline, and/or different aircraft family, and/or different time period) to assess the stability of our approach. But before that, we need to develop a heuristic to solve this very difficult problem represented by (11-14). Note that the right hand side of constraint (12), in and by itself, is a very challenging problem for large network sizes. It is a type of robust crew pairing optimization problem. No prior study in the literature has solved robust crew pairing problems of size as large as those of the networks used in this paper. Therefore, we develop and implement new heuristic approaches to solve the inverse of this already very difficult problem. In this section, we describe the solution approach. Then, in Section 4, we present our computational results.

We begin this section by describing, in Section 3.1, our overall heuristic for solving the calibration problem. This involves repeatedly solving instances of the robust crew pairing optimization problem (8-10). Section 3.2 summarizes the overall approach used to solve this robust crew pairing optimization problem, which itself includes repeatedly solving instances of the LP (linear programming) relaxation of this integer optimization problem. The solution to the LP relaxation of the crew pairing optimization

problem involves repeatedly solving instances of a sub-problem called the pricing problem. The process for solving this pricing problem is described in Section 3.3.

3.1 Local Search Heuristic for the Calibration Problem

We use a local search method for solving the optimization problem given by (11-14). The algorithm is as follows:

```
INITIALIZE:

Set all parameters to zero, i.e.,  $\alpha_i = \beta_i = 0, i \in \{1,2,3,4\}$  and  $\gamma_i = 0, i \in \{5,6\}$ .

LOOP:

    FOR i=1:4

        Perform local grid-search by varying  $(\alpha_i, \beta_i)$  values to minimize  $\sum_{i=1}^6 |F^{\hat{x}}(i) - F^x(i)|$ . Update  $(\alpha_i, \beta_i)$  values.

    END FOR

    FOR i=5:6

        Perform local line-search by varying  $\gamma_i$  values to minimize  $\sum_{i=1}^6 |F^{\hat{x}}(i) - F^x(i)|$ . Update  $\gamma_i$  values.

    END FOR

IF no parameter values got updated in the last iteration of the outer LOOP then EXIT.

END LOOP
```

Note that the local line-searches and local grid-searches mentioned in the algorithm require us to examine various combinations of PARAMS values to calculate the $\sum_{i=1}^6 |F^{\hat{x}}(i) - F^x(i)|$ value. Examining each combination of PARAMS values requires solving the robust crew pairing optimization problem given by (8-10). Next, we discuss the process for solving this problem.

3.2 Crew Pairing Solution Approach

The deterministic crew pairing optimization problem presents two main challenges. First, the number of feasible pairing variables is extremely large for major airlines' networks, making it difficult to solve even the LP relaxation of the problem. Second, the existence of integrality constraints adds additional complexity. Typically, the crew pairing optimization problem is solved by techniques such as branch-and-price (Barnhart et al., 1998). Branch-and-price techniques combine ideas from the branch-and-bound algorithm for solving integer optimization problems with delayed column generation ideas for solving large-scale linear optimization problems. Given a feasible solution, delayed column generation technique requires finding columns with negative reduced cost. This is sometimes called as the pricing problem. The reader is referred to Kasirzadeh et al. (2015) for a detailed review of the state-of-the-art in solving the deterministic crew pairing optimization problems.

Note that unlike previous researchers that studied the crew pairing optimization problem, our goal is not just to be able to solve the robust crew pairing problem once. Instead, its solution constitutes a sub-problem within our overall calibration optimization process described in Section 3.1, and the overall calibration algorithm requires solving hundreds of these individual crew pairing optimization problems. Therefore, our computational performance requirements for solving the individual robust crew pairing problems are far more stringent than most prior studies in the literature. Unlike prior studies, we cannot afford to wait for several hours to solve the crew pairing optimization problem. Instead of using column generation at each node of the branch-and-bound tree, which is very time consuming, we use a heuristic

strategy to solve this problem. As explained in Section 4, this strategy helps us obtain solutions that are provably within a small optimality gap. This strategy can be summarized as follows, and it refers to two other algorithms, *Algorithm A* and *Algorithm B*, which are described in Section 3.3.

Heuristic Solution Strategy for the Robust Crew Pairing Optimization Sub-problem

Step 1: Form the Restricted Master Problem (RMP) by including only a small subset of columns and relaxing the integrality constraints.

Step 2: Solve the RMP to find a set of dual variable values.

Step 3: Using the dual variables from Step 2, solve the pricing problem with Algorithm B to identify if one or more variables have negative reduced cost. If so, add *all* variables with negative reduced costs to RMP's column pool and go back to Step 2; else go to Step 4.

Step 4: Using the dual variables from Step 2, solve the pricing problem with Algorithm A to identify if one or more variables have negative reduced cost. If so, add *all* variables with negative reduced costs to RMP's column pool and go back to Step 2; if not go to Step 5.

Step 5: Fix the largest fractional variable to 1 and solve RMP again. Check if an integer solution is obtained. If not, go back to Step 2; else stop.

This algorithm was developed after experimenting with various alternative heuristic ideas, and each step was chosen carefully based on the computational performance with and without it. Our computational experiments revealed that Step 5 helps improve the computational performance substantially while increasing the optimality gap by very little or nothing. Also, we found that decomposing the pricing problem's solution process into two steps, i.e. using Algorithm B in Step 3 and Algorithm A in Step 4, was a vital part of the computational speedup that we achieved. Without this, we would not have been able

to finish all our experiments in reasonable amounts of time to accomplish this research project. More details about this two-step approach are provided in Section 3.3.

3.3 Solution to the Pricing Problem

Researchers have proposed and implemented a variety of methods for solving the pricing problem. Garfinkel and Nemhauser (1970) used enumeration-based approaches, that is, approaches in which all paths are enumerated in the sub-problem. However, the success of this approach is highly dependent on the size of problem, and is difficult to solve for large network sizes. AhmadBeygi, Cohn, and Weir (2009) proposed an integer programming approach for solving the pricing problem. This approach is easy to implement using commercial solvers, but has poor performance for large networks such as the ones used in our study. *Multi-Label Shortest Path* (MLSP) is a commonly used algorithm for solving the pricing problem to generate crew pairings (Desaulniers et al. 2005; Vance et al. 1997). Unlike the deterministic crew pairing optimization problem, our robust crew pairing problem involves a more complicated objective function that includes planned costs as well as six different types of penalty costs. Furthermore, as will be explained in Section 4, our network size is the largest among all existing research studies addressing any variety of the robust crew pairing problem. Therefore, we cannot directly use any of the existing methods to solve the problem to near optimality in a limited time. Therefore, we develop a new two-step approach to solve the pricing problem to optimality.

For the pricing problem, the objective function is to minimize the reduced cost of the chosen pairing which is equal to the objective function coefficient of the chosen pairing minus the sum of the dual variables corresponding to all flights included in this pairing. Irnich and Desaulniers (2005) frame this problem as a Shortest Path Problem with Resource Constraints (SPPRC) and utilize a dynamic programming approach to solve it. The core idea is to build paths in a flight network by extending them into all feasible directions, and to identify those paths which represent feasible crew pairings with

negative reduced costs. The efficiency of this approach depends on being able to identify and eliminate those paths such that these paths themselves and all their extensions are guaranteed to be suboptimal. These non-useful paths are discarded by using a *dominance sub-algorithm* based on a set of dominance rules. However, the standard dominance algorithm is too slow for our computational requirements when applied to our large-scale flight networks. In order to accelerate the dominance algorithm, we tested different variations of dominance rules. We found that if we remove the rule that requires that both the dominant and dominated paths to have the same crew base (Algorithm B), the speed of the dominance step significantly increases. However, this simplification risks eliminating some paths that could have negative reduced cost and hence Algorithm B is not guaranteed to identify all paths with negative reduced costs. Algorithm B, therefore, serves as an intermediate step, which helps us identify some negative reduced cost pairings in a very short computational time. However, if it is unable to find any such pairing, then we revert to the full implementation of the dominance algorithm (named Algorithm A) to perform a comprehensive search for negative reduced cost pairings. In our computational experiments we find that, in most of the iterations, Algorithm B is able to identify enough variables with negative reduced costs and add them to the RMP's column pool, thus requiring us to implement Algorithm A much more sparingly and hence speeding up the pricing problem solution process dramatically. The two dominance algorithms are presented below.

Algorithm A: This is a dominance algorithm with an exact implementation, similar to that described by Irnich and Desaulniers (2005), wherein only a path starting with the same crew base can dominate another path. The reader is referred to Irnich and Desaulniers (2005) for more details of this algorithm.

Algorithm B: This is a dominance algorithm with an exact implementation similar to that described by Irnich and Desaulniers (2005) except that a path starting with the same or different crew base can dominate another path.

The exact set of labels used by Algorithm A in our robust crew pairing implementation is as follows. Note that Algorithm B uses all but the last of the labels listed below.

1. Number of duties covered so far by the path.
2. Total flying time so far in the current duty of the path.
3. Total elapsed time so far in the current duty of the path.
4. A constant multiple (ζ) of the total elapsed time so far of the path minus the sum of the dual contributions of all flights included so far in the path.
5. Total flying time so far in the current duty plus the sum of the costs of previous duties in the path minus the sum of the dual contributions of all flights included so far in the path.
6. Minimum guaranteed pay of the current duty plus the sum of costs of previous duties in the path minus the sum of the dual contributions of all flights included so far in the path.
7. A constant multiple of (ε) of the total elapsed time so far in the current duty plus the sum of costs of previous duties in the path minus the sum of the dual contributions of all flights included so far in the path.
8. Crew base (starting point) of the path.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

4 Case Study

In this section, we apply the models presented in Section 2 and the solution methods presented in Section 3 to four networks from two airlines across multiple time periods. We use confidential airline data containing crew scheduling samples acquired from these two airlines to calibrate and validate our parameterized crew pairing models. The data sources and data preprocessing steps are described in Section 4.1, while the computational performance of our models is highlighted in Section 4.2. Section 5 presents the detailed calibration and validation results which demonstrate the accuracy and stability of our approach.

4.1 Data Source and Data Preprocessing

We acquired crew schedules data sample from one major regional carrier (RC) and one major network legacy carrier (NLC) in the United States. The regional carrier has a homogenous fleet consisting of only one fleet family and the data available to us spanned two full months, namely, March and April 2014. The network legacy carrier's operations consisted of several different fleet families. However, for our computational experiments we chose only the three largest networks, namely, those operated by A320, B737 and B757 aircraft types, because the remaining aircraft types represent much smaller sized networks. The network legacy carrier's crew scheduling data sample available to us spanned one full year, from August 2013 to July 2014.

Aside from the confidential data in the form of the crew scheduling samples, we also used the Airline On-Time Performance (AOTP) database from the BTS website (BTS, 2016) which contains on-time arrival data for domestic flights by all major U.S. carriers. Most importantly for our purposes, AOTP provides tail number for each flight, which is a key piece of information useful to track aircraft rotations in real-world airline schedules.

1
2
3 Since our data is obtained from two separate sources, some data preprocessing steps, including data
4 cleaning and merging, need to be performed before using the data for model calibration and validation.
5
6 Preprocessing consisted of two major steps. The first step is to get aircraft tail numbers for all flights in
7
8 the crew scheduling samples by matching each flight in the sample with exactly one flight in the AOTP
9
10 database. This is performed by matching departure airport, arrival airport, departure time, scheduled
11
12 arrival time and the airline code. We are able to match around 95% of all the flights in the airline crew
13
14 scheduling samples.
15
16
17
18
19

20 The second preprocessing step was data filtering to account for the limitations of the AOTP database.
21
22 Because tail number information is missing for some flights in the AOTP, we use, as input to our models,
23
24 only those flights for which the tail number is present in the AOTP database. Since only domestic flights
25
26 information is provided in the AOTP, we removed all international flights from our crew scheduling
27
28 sample as well. Typically, cockpit crews are assigned to operate aircraft belonging to only one fleet
29
30 family in a given pairing. Indeed, almost all the crew pairings in our confidential crew scheduling data
31
32 contained flights operated by a single fleet family. We eliminated the few crew pairings (and their
33
34 corresponding flights) which cut across multiple fleet families in our crew schedule data. As a result, the
35
36 crew pairing problem can be considered separately for each fleet type. For the regional carrier, we use
37
38 first week of March as our calibration dataset and first week of April as validation dataset. For network
39
40 legacy carrier, we use first week of January as our calibration data and first week of one month in each
41
42 quarter, namely, February, April, July and October, as our validation dataset. We also eliminated all crew
43
44 pairings (and all flights in those crew pairings) when at least one flight in that crew pairing had to be
45
46 removed for any of the reasons mentioned above. Overall, this resulted in removal of approximately 15-
47
48 20% of all crew pairings and approximately 10-15% of all flights in our network legacy carrier crew
49
50 schedule sample across different time periods in the sample. Also it resulted in approximately 10-15% of
51
52
53
54
55
56
57
58
59
60

all our crew pairings and approximately 10-15% of all the flights in our regional carrier crew schedule sample across different time periods in the sample.

We were able to get the true values of the planned crew cost parameters, such as δ , ε , and ζ , the values of lower limits on crew sit times and crew rest times, and the upper limits on the maximum duty flying time and maximum duty elapsed time from both the airlines represented in our data samples. Finally, note that all our data is related to cockpit (and not cabin) crew schedules, which are more stringent in their regulation and hence are expected to be responsible majority of the crew-propagated delays and disruptions. Hence all our analysis is restricted to cockpit crew schedules only, and hence deals with a large part, but not all, of the crew-propagated delays and disruption. Note that this is a limitation of the available data and not of our methodology. Our methodology is valid if we were to perform a similar analysis with cabin crew scheduling data.

4.2 Computational Experiments

In this section, we demonstrate the computational performance of our robust crew pairing solution heuristic. In order to solve the inverse problem of robust crew pairing parameter calibration, we need to use a computationally efficient method for solving the robust crew pairing problem itself in the first place. This is so because we will need to solve this problem several (hundreds of) times in order solve one instance of the inverse problem. Furthermore, as demonstrated in Section 5, the size of each crew pairing problem solved by us is bigger than the network size used by any of the previous studies in robust crew pairing literature. CPLEX 12.5 solver with its default settings is used to solve all the linear and integer optimization problems. An 8-thread / 4-core Intel® i7-X5600 CPU with 8GB RAM and Windows 7 Professional as the operating system was used for all computational experiments.

Table 1. Computational Performance of Our Heuristic

Network Size (Flights)	Pricing Approach	Root LP Bound	Our Integer Solution	Gap	Solution Time (hours)
102	Algorithm A Only	398.36	398.37	0.025%	<0.1
	Algorithm B + Algorithm A	398.36	398.37	0.025%	<0.1
3300	Algorithm A Only	12760.43	12832.12	0.56%	10
	Algorithm B + Algorithm A	12760.43	12832.12	0.56%	2

For demonstrating the computational performance of our crew pairing optimization approach, we consider two networks in our crew scheduling sample data: one with 102 flights and the other with 3300 flights. We first solve the root node LP relaxation to optimality to get a lower bound on the optimal objective function value of the integer optimization problem. This is listed in the third column of Table 1. Then using the method described in Section 3.2, we obtain a feasible, but not necessarily optimal, solution of the integer optimization problem. Its objective function value is listed in the fourth column. Fifth column gives the gap between the values in the third and fourth columns by dividing the difference between the two by the value in the third column. The last column gives the total runtime for obtaining the solution in the fourth column. Note that the gap listed in the fifth column gives an upper bound on the true optimality gap of our heuristic solution.

Second column lists the solution approach. We first list the performance of our overall heuristic using exact SPPRC method, i.e., without using Algorithm B. We also list performance when the pricing problem is solved using the heuristic that involves use of both Algorithm A and Algorithm B. Across all cases listed in Table 1, the gap was at most 0.56%. Furthermore, for the small network, Algorithm B does not help in speeding up because Algorithm A alone is sufficient to solve this small problem within a few minutes of overall computational time. However, in case of the large network with 3300 flights, the combined use of Algorithm A and Algorithm B, as described in Section 3.2, significantly reduces the overall computational time from 10 hours to 2 hours. Similar improvements were observed in all of our large

network instances. This demonstrates the value of using our modified two-step pricing problem solution process.

5 Calibration and Validation Results

5.1 Calibration Results

In this section, we display our calibration results (1) to identify the relative importance of different robustness-related features of the crew schedules for different airline types and aircraft families that affect crew-propagated delays and disruptions, and (2) to assess their stability and identify trends across airline types and aircraft families. Table 2 lists the estimated parameters resulting from our calibration process as described in Section 4.1, while Table 3 lists the penalty function values corresponding to each feature.

Table 2. Parameter Results

Feature Type	Parameter	RC	NLC-A320	NLC-B737	NLC-B757
Type 1	α_1	1	0.3	0.3	0.8
	β_1	1.5	0.65	1	3
Type 2	α_2	0.5	0	0	1.5
	β_2	0.15	0	0	1.1
Type 3	α_3	0.4	1.1	3	1
	β_3	1.4	0.4	1.25	1
Type 4	α_4	2	1.65	2	3.8
	β_4	1.5	0.5	0.7	1.3
Type 5	γ_5	0	0.025	0.08	0
Type 6	γ_6	0.05	0.07	0.4	0

Note that, while we had access to crew schedules data on the networks of all six aircraft families used by the network legacy carrier (NLC), we only present results using the three biggest networks (corresponding to A320, B737 and B757's) because the numbers of flights in the networks of the other three aircraft families are too small to be interesting for our analysis. Additionally, we present results using the network of the regional carrier (RC) which was operated by a homogeneous fleet consisting of only one aircraft family. So Tables 2 and 3 present results using four distinct networks, namely, regional carrier's complete network (*RC*) excluding the flights that were filtered out in pre-processing, and the network legacy carrier's networks using A320 fleet family (*NLC-A320*), B737 fleet family (*NLC-B737*), and B757 fleet family (*NLC-B757*). The numbers of flights in the RC, NLC-A320, NLC-B737 and NLC-B757 networks are 2432, 1200, 1840, and 147, respectively. All experiments are conducted over a seven day time horizon and the maximum number of duties allowed in a single crew pairing is 4 in all cases.

Table 3. Penalty Function Values

Cost Type		Planned	Type 1	Type 2	Type 3	Type 4	Type 5	Type 6	%age
RC	Airline Sample	4590.40	105.87	42.05	0	3.15	0	33.90	3.87%
	Without Calibration	4143.38	293.72	46.94	4.08	93.60	0	43.34	10.41%
	With Calibration	4244.99	91.23	40.66	0.41	8.44	0	34.80	3.97%
NLC-320	Airline Sample	4800.45	5.26	0	15.91	19.39	3.42	28.07	1.48%
	Without Calibration	4515.51	13.05	0	62.24	148.29	5.30	37.45	5.57%
	With Calibration	4539.24	5.27	0	9.48	12.53	3.38	28.00	1.28%
NLC-737	Airline Sample	7448.00	0.77	0	53.74	32.75	9.44	182.4	3.61%
	Without Calibration	6696.80	9.35	0	300.36	280.52	12.56	280.00	11.65%
	With Calibration	6773.94	2.05	0	15.77	12.21	10.88	175.20	3.09%
NLC-757	Airline Sample	976.18	0.44	1.08	0.47	3.52	0	0	0.56%
	Without Calibration	925.08	0.41	2.17	2.28	22.25	0	0	2.85%
	With Calibration	927.09	0	0	0.77	6.26	0	0	0.75%

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

In addition to the penalty function values corresponding to each of the six features, Table 3 lists the planned cost values for comparison purposes. In the last column of Table 3, the total penalty cost as a percentage of the planned cost is listed. For each of the four networks, there are three rows. First row provides the components of the objective function evaluated using the actual airline-provided crew schedules, and the calibrated parameter values listed in Table 2. The second row provides the components of the objective function evaluated using the calibrated parameter values (listed in Table 2), and the crew pairings obtained by solving the crew pairing optimization problem by setting all parameters to 0. Finally, row column provides the components of the objective function evaluated using the calibrated parameter values (listed in Table 2), and the crew pairings obtained by solving the crew pairing optimization problem by setting all parameters to their calibrated values (listed in Table 2). Note that the calibration algorithm does not explicitly attempt to match the planned cost values, because our aim is to match the potential for crew-propagated delays and disruptions. Yet, for all four networks, the planned costs of the crew pairing solution generated by our approach are found to be closer to the actual airline-provided crew schedules with calibration than without calibration. Across all networks and all cost types, the with-calibration cost values were found to be closer (in most cases significantly closer), than the without-calibration cost values, to the airline-provided crew pairing solutions in 22 out of the 23 network-cost type combinations in Table 3. Note that we have excluded the network-cost type combinations where all three values are zeros, which happens in 5 instances. In Section 5.2, we provide a metric for easy comparison of this degree of closeness in the form of a percentage error measure.

Looking at Tables 2 and 3, we can already notice several differences between the four networks. Some of these differences reflect the different crew pay and crew legality rules across the four networks. For the NLC-B737 and NLC-A320 networks, crew pay does not depend on the time away from base (i.e., parameter ζ in (1) equals 0). As a result, there is no tradeoff associated with the length of the rest period. For networks with nonzero ζ values, having short rest periods can cause delay propagation while

having long rest periods can add to the planned crew cost. Absent this tradeoff, for the NLC-A320 and NLC-B737 networks, the optimization simply sets the rest period lengths such that the Type 2 penalty function value is zero irrespective of the values of parameters α_2 and β_2 . For the RC network, we find that irrespective of the values of Type 5 parameters α_5 and β_5 , Type 5 penalty function value zero. Recall that Type 5 penalty function penalizes the number of occurrences of crew-base purity violations. Because of the simple hub-and-spoke structure of the regional carrier's network, most crew travel from a hub to a spoke and back, there isn't much opportunity for changing the Type 5 penalty cost by varying the γ_5 parameter. Therefore, for the RC network, γ_5 value remains 0 even after calibration. Finally the NLC-B757 network is the smallest among the four, due to which crews usually don't have too many alternatives other than staying with the aircraft and the crews do not end up going more than a distance of 1 unit away from the crew base in the adjacency graph. This simplified structure of the network explains why both Type 5 and Type 6 parameters (namely, γ_5 and γ_6) and the corresponding penalty function values are set to 0 for the NLC-B757 network in all cases.

Although these four flight networks vary in size, and although the absolute level of crew-propagated delays and disruptions cannot be directly compared across the four networks, the last column of Table 3 lies in the range from 0.75% to 4% across all four networks, for the airline-provided crew schedules and also for the solution generated by our calibrated model. These numbers are much higher for the crew-schedules generated using the uncalibrated model. These results demonstrate that our approach generates crew schedules whose balance between planned and operational costs is similar to that of the actual crew-scheduling solution used by the airlines. Previous studies involving robust crew pairing optimization, such as Yen and Birge (2006), have emphasized the importance of right tradeoff between planned and operational costs. They test effects of different penalty parameters to control this tradeoff, but do not provide explicit insights into the right tradeoff. Our results, for the first time, allow us to get a measure of the perceived balance between the planned costs and penalty costs as reflected by the

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

airlines’ actual crew scheduling practices. Table 3 suggests that the right balance between planned and penalty costs across the four different networks is in a relatively narrow range of 0.75% to 4% and is thus quite stable across airlines and aircraft families.

Unlike previous studies in literature, such as Schaefer et al. (2005), Shebalov and Klabjan (2006), Gao, Johnson, and Smith (2009), which focus on minimizing a subset of the potential factors affecting crew-propagated delays and disruptions, we use a more comprehensive approach by including wider variety of factors. Unlike previous robust crew pairing research by Yen and Birge (2006) whose work considers total expected cost of future actions due to disruptions as a component of the crew pairing objective function, our approach can give a separate ratio between the penalty cost corresponding to each feature of operational cost and the planned cost. By allowing penalty costs of each component to be assessed separately, we get a clearer understanding of the relative importance of each component as perceived by the airlines.

Table 3 provides some preliminary evidence of the effectiveness and accuracy of our calibration framework. However, there are several shortcomings of using the in-sample penalty costs to assess the similarity of our solutions to airline-provided crew schedules. First, this in-sample comparison has an inherent bias because we are using the same data samples to calibrate and test the accuracy. We address this concern in Section 5.2 by presenting results of computational experiments where the parameter calibration is performed using one dataset and then other datasets corresponding to different time periods and/or different aircraft families and/or different airline types, are used to assess the out-of-sample accuracy of our approach. Second, we are measuring the closeness of the two solutions using penalty functions, which themselves depend on the calibrated parameter values. To make our comparisons more meaningful, we need to be able to compare the distributions of the actual feature values, which we do in Section 5.3. Finally, an argument could be made that all the methods

used by us for evaluating the accuracy of our approach depend on the features that we deem to be proxies for crew-propagated delays and disruptions. While many of these were chosen and are well-supported by previous research, we cannot claim them to be precise measures of crew-propagated delays and disruptions. Therefore, a true test of the performance of our approach can only be conducted by comparing the actual crew-propagated delays and disruptions. This concern is addressed in Section 6.1.

5.2 Out-of-Sample Validation Results

This section demonstrates the accuracy and stability of our results through out-of-sample validation. First, in Tables 4 through 7, we present results where the calibration and validation datasets belong to two different time periods for the same airline and for the same fleet family. For the regional carrier network, we choose March 2014 as the calibration set and April 2014 as the validation set. For the three networks belonging to the network legacy carrier, we select the first week of one month from each quarter to represent flight schedule through a full year. Specifically, we use January 2014 data for calibration and perform validation using datasets from April 2014, July 2014 and October 2013. Additionally, February 2014 dataset is also used to perform validation for a scenario where the calibration and validation datasets are not too far apart in time from each other. The intent of this validation is to test the validity of using parameters calibrated using one time period to predict crew schedules for another time period for the same airline and same aircraft family. If the results are found to be stable across time periods, then this allows us to use crew scheduling data samples from one period to estimate crew schedules for other periods and thus reduces our data requirements if we were to estimate crew-propagated delays and disruptions across long periods of time.

Let C_i be the Type i , $i \in \{1, \dots, 6\}$, penalty cost associated with the crew schedule generated by our approach, and $C_i^{Airline}$ be the Type i , $i \in \{1, \dots, 6\}$, penalty cost associated with the corresponding

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

airline-provided crew schedule. Then we define the Absolute Percentage Error (APE) as $\frac{|C_i - C_i^{Airline}|}{\sum_i^N C_i^{Airline}}$, where N is the total number of components of the penalty cost function, i.e. the total number of robustness features. Note that this is not a commonly used method of error representation, but it is chosen because it offers certain advantages for our problem setting. First, the choice of denominator ($\sum_i^N C_i^{Airline}$) in the APE expression guarantees that we do not have issues related to division by zero. Contrast this choice of denominator with a more standard $C_i^{Airline}$ term as the denominator, which would have led to division-by-zero issues in many cases, such as Type 3 error for RC network, as presented in Table 3. Second, since we use the same denominator for each feature, it is easy and meaningful to compare the percentage errors across different features.

Tables 4 through 7 list the APEs for each feature and also the average and maximum values across features. The columns titled “Before” and “After” list the errors for crew pairing solutions generated using uncalibrated and calibrated parameters respectively. Note that as described in Section 5.1, in all cases the penalty function evaluation is performed using the calibrated parameters. Looking at the results presented in Tables 4 through 7, several observations can be made. Comparing to before calibration errors, errors are substantially lower in most cases after calibration. The improvement is especially clear when looking at the average or maximum values of the APEs across types. Average and maximum APEs are reduced substantially by the calibration process and a reduction is observed across all calibration and validation datasets. In many cases the reduction is by one or more orders of magnitudes. While there is slight worsening of the APEs in the out-of-sample validation datasets compared with the in-sample calibration dataset, especially for the RC network, the out-of-sample APEs are consistently reduced by the calibration process demonstrating the stability and effectiveness of our approach. Moreover, seasonality is not found to play a significant role in terms of the errors. The out-of-sample validation errors did not worsen and stayed stable as the time between the calibration and

validation datasets increased from 1 month (for February validation dataset) to 6 months (for the July validation dataset). The consistently lower error values for the crew pairings generated by using calibrated compared with the uncalibrated parameters, as measured individually, using averages, or using maximum values, indicate that our approach produces crew pairings that are stable across time periods of up to several months.

Table 4. Cross-Time Period Validation Results: RC

Data	Calibration (Mar)		Validation (Apr)	
	Before	After	Before	After
Type 1	101.6%	7.9%	58.9%	28.8%
Type 2	2.6%	0.8%	1.9%	0.3%
Type 3	2.2%	0.2%	1.6%	0.1%
Type 4	48.9%	2.9%	26.2%	0.2%
Type 5	0	0	0	0
Type 6	5.1%	0.5%	2.9%	0.4%
Average	26.7%	2.0%	15.3%	5.0%
Maximum	101.6%	7.9%	58.9%	28.8%

Table 5. Cross-Time Period Validation Results: NLC-A320

Data	Calibration (Jan)		Validation (Feb)		Validation (Apr)		Validation (Jul)		Validation (Oct)	
	Before	After	Before	After	Before	After	Before	After	Before	After
Type 1	10.8%	0.0%	12.0%	3.08%	9.3%	2.5%	11.3%	2.6%	21.1%	5.5%
Type 2	0	0	0	0	0	0	0	0	0	0
Type 3	64.3%	8.9%	59.8%	4.01%	39.2 %	12.9%	40.8%	8.6%	70.9%	11.5%
Type 4	178.9%	9.5%	228.9%	8.18%	140.3%	31.0%	134.6%	23.1%	299.2%	0.3%
Type 5	2.6%	0.1%	1.9%	0.74%	3.0%	1.6%	3.2%	2.6%	4.4%	3.3%
Type 6	13.0%	0.1%	12.0%	1.52%	9.4%	1.9%	8.2%	3.7%	18.9%	5.5%
Average	44.93%	3.10%	52.43%	2.92%	33.53%	8.32%	33.02%	6.77%	69.08%	4.35%
Maximum	178.9%	9.5%	228.9%	8.18%	140.3%	31.0%	134.6%	23.1%	299.2%	11.5%

Table 6. Cross-Time Period Validation Results: NLC-B737

Data	Calibration (Jan)		Validation (Feb)		Validation (Apr)		Validation (Jul)		Validation (Oct)	
	Before	After	Before	After	Before	After	Before	After	Before	After
Type 1	3.1%	0.5%	3.9%	0.9%	2.0%	0.3%	2.6%	0.9%	3.1%	0.7%
Type 2	0	0	0	0	0	0	0	0	0	0
Type 3	88.4%	13.6%	76.3%	7.8%	48.9%	13.4%	68.1%	12.8%	93.8%	6.9%
Type 4	88.8%	7.4%	86.9%	5.2%	60.7%	10.3%	70.9%	10.5%	100.2%	8.2%
Type 5	1.1%	0.5%	1.1%	1.0%	2.0%	0.6%	2.3%	0.7%	0.4%	1.1%
Type 6	35.0%	2.6%	30.2%	14.1%	26.3%	7.4%	28.0%	6.8%	35.3%	18.0%
Average	36.07%	4.10%	33.07%	4.83%	23.32%	5.33%	28.65%	5.28%	38.80%	5.82%
Maximum	88.8%	13.6%	86.9%	14.1%	60.7%	13.4%	70.9%	12.8%	100.2%	18.0%

Table 7. Cross-Time Period Validation Results: NLC-B757

Data	Calibration (Jan)		Validation (Feb)		Validation (Apr)		Validation (Jul)		Validation (Oct)	
	Before	After	Before	After	Before	After	Before	After	Before	After
Type 1	0.5%	8.0%	0.0%	2.9%	0.0%	0.0%	3.7%	0.8%	3.8%	3.8%
Type 2	19.8%	19.6%	0.0%	0.0%	161.6%	0.0%	16.0%	0.8%	105.5%	10.2%
Type 3	32.8%	5.4%	14.1%	3.9%	54.7%	0.0%	17.4%	19.8%	20.3%	9.6%
Type 4	339.9%	49.7%	185.1%	92.0%	1080.8%	25.1%	309.9%	79.1%	321.7%	39.4%
Type 5	0	0	0	0	0	0	0	0	0	0
Type 6	0	0	0	0	0	0	0	0	0	0
Average	65.50%	13.78%	33.20%	16.47%	216.18%	4.18%	57.83%	16.75%	75.22%	10.50%
Maximum	339.9%	49.7%	185.1%	92.0%	1080.8%	25.1%	209.9%	79.1%	321.7%	39.4%

As the next step in our out-of-sample validation process, Tables 8 through 11 present cross-validation results where the validation is performed on a dataset which belongs to a different airline, a different fleet family, and in some cases, a different time period compared with the calibration dataset. This constitutes an important test of our calibration approach because realistically we cannot expect crew scheduling samples to be available for all combinations of airlines, fleet types and time periods. Instead, if we are able to access a small crew scheduling sample from one airline, for one fleet family, for one

time period, it is desirable to use that sample to calibrate parameters of a model that can then be used to generate crew schedules for other airline types, other fleet types and/or other time periods. Tables 8 through 11 present these cross-validation results where the four chosen combinations of networks and time periods are RC network for March 2014, NLC-B737 network for January 2014, NLC-B757 network for January 2014, and NLC-A320 network for January 2014. For each network in each table, we use the parameter sets obtained by calibration over the networks listed in the top row.

Each table represents results of validation using a single network and time period combination specified in the table caption. The intent of this validation is to test whether our parameters calibrated for one combination of airline, fleet family and time period still perform well for other combinations of airlines, fleet families and time periods.

Table 8. Validation across Airline, Fleet Family and Time Period for RC Network for March 2014

Data	RC (Calibration)		NLC-A320		NLC-B737		NLC-B757	
	Before	After	Before	After	Before	After	Before	After
Type 1	101.6%	7.9%	21.4%	6.7%	6.1%	0.6%	43.2%	1.1%
Type 2	2.6%	0.8%	22.0%	16.5%	0	0	0.8%	7.5%
Type 3	2.2%	0.2%	28.8%	11.6%	31.1%	16.8%	18.5%	1.8%
Type 4	48.9%	2.9%	50.8%	21.1%	26.4%	7.4%	149.5%	51.2%
Type 5	0	0	0.3%	1.1%	0.4%	2.3%	0	0
Type 6	5.1%	0.5%	6.5%	1.4%	15.7%	7.7%	0	0
Average	26.73%	2.05%	21.63%	9.73%	13.28%	5.80%	35.33%	10.27%
Maximum	101.6%	7.9%	50.8%	21.1%	31.1%	16.8%	149.5%	51.2%

Table 9. Validation across Airline, Fleet Family and Time Period for NLC-A320 Network for January 2014

Data	NLC-A320	RC	NLC-B737	NLC-B757
	(Calibration)			

	Before	After	Before	After	Before	After	Before	After
Type 1	10.8%	0.0%	36.9%	38.2%	2.4%	0.2%	34.1%	4.8%
Type 2	0	0	22.0%	16.5%	0	0	6.6%	6.3%
Type 3	64.3%	8.9%	1.1%	0.0%	52.6%	5.3%	29.7%	0.5%
Type 4	178.9%	9.5%	85.9%	4.0%	73.6%	3.0%	895.2%	52.4%
Type 5	2.6%	0.1%	0	0	2.9%	1.7%	0	0
Type 6	13.0%	0.1%	8.2%	0.6%	25.8%	16.6%	0	0
Average	44.93%	3.10%	25.68%	9.88%	26.22%	4.47%	160.93%	10.67%
Maximum	178.9%	9.5%	85.9%	38.2%	73.6%	16.6%	895.2%	52.4%

Table 10. Validation across Airline, Fleet Family, and Time Period for NLC-B737 for January 2014

Data	NLC-B737 (Calibration)		RC		NLC-A320		NLC-B757	
	Before	After	Before	After	Before	After	Before	After
Type 1	3.1%	0.5%	60.2%	28.9%	17.1%	7.4%	27.6%	2.9%
Type 2	0	0	7.2%	22.6%	0	0	38.6%	0.3%
Type 3	88.4%	13.6%	2.6%	0.3%	120.1%	13.4%	37.0%	3.3%
Type 4	88.8%	7.4%	118.7%	0.2%	242.9%	20.6%	685.7%	80.6%
Type 5	1.1%	0.5%	0	0	3.5%	3.6%	0	0
Type 6	35.0%	2.6%	12.1%	5.8%	44.0%	33.8%	0	0
Average	36.07%	4.10%	33.47%	9.63%	71.27%	13.13%	131.48%	14.52%
Maximum	88.8%	13.6%	118.7%	28.9%	242.9%	33.8%	685.7%	80.6%

Table 11. Validation across Airline, Fleet Family, and Time Period for NLC-B757 Network for January 2014

Data	NLC-B757 (Calibration)		RC		NLC-A320		NLC-B737	
	Before	After	Before	After	Before	After	Before	After
Type 1	0.5%	8.0%	21.4%	25.0%	2.5%	2.6%	0.0%	0.0%
Type 2	19.8%	19.6%	50.4%	19.6%	0	0	0	0
Type 3	32.8%	5.4%	3.1%	1.4%	69.4%	6.1%	70.4%	6.4%

Type 4	339.9%	49.7%	87.0%	33.5%	101.5%	16.5%	48.3%	0.1%
Type 5	0	0	0	0	2.3%	1.8%	3.2%	5.2%
Type 6	0	0	6.3%	0.9%	5.7%	2.5%	13.9%	4.0%
Average	65.50%	13.78%	28.03%	13.40%	30.23%	4.92%	22.63%	2.62%
Maximum	339.9%	49.7%	87.0%	33.5%	101.5%	16.5%	70.4%	6.4%

Tables 8 through 11 show that the average and maximum errors (APEs) when using parameters after calibration are much smaller, when compared to those using parameters before calibration, for all combinations of calibration and validation datasets. However, when compared with the calibration error, the validation errors are typically larger when using a different airline type and/or a different fleet family's calibration parameters than using their own calibrated parameters. This is especially obvious in Table 8 where the calibration is performed using the RC network for March 2014 and the validation is performed using the three NLC networks for January 2014. This seems to suggest that the three NLC networks are more "similar" to each other in terms of their calibrated parameters values than the similarity between NLC and RC networks. This is not surprising given that the underlying set of flights representing the RC network exhibits many differences in the network structures, schedules and flight durations when compared with the three NLC networks. Moreover, Tables 9 and 10 together suggest that the parameters for the NLC-A320 and NLC-B737 networks are especially similar to each other as reflected by their low cross-validation errors with respect to each other. This phenomenon can also be explained by the fact that A320 and B737 aircraft families are similar to each other in that they are both single aisle, twin-engine aircraft with similar seating capacity and range capabilities causing their flight networks to also look similar to each other.

Thus Tables 8 through 11 provide several interesting insights. First, they demonstrate that the out-of-sample validation errors are considerably lower using the calibrated than the uncalibrated parameters even when the calibration was performed using a crew scheduling sample from, airline type and/or fleet family. However, we also note that the error reduction by using the calibrated rather than uncalibrated

parameters is greater when the calibration and validation datasets are more similar, in terms of airline type and fleet family. This suggests that, on the one hand, when estimating crew schedules for a given flight network, it is advisable to use a parameter set that has been calibrated using a flight network that shares as many of its attributes as possible. On the other hand, though, using *any* set of calibrated parameters is still likely to be considerably better than using uncalibrated parameters. Even if the calibrated parameters are from a different time period, different airline and/or different fleet family, they improve the accuracy considerably compared with the uncalibrated parameters, i.e., compared with solving the deterministic crew scheduling problem. Thus, while it is advisable and beneficial to have a wide variety of airline crew schedule samples, our calibration approach enhances the degree of similarity of the generated crew pairing solution with the actual pairing solution used by the airline even when crew sampling data is relatively scarce.

5.3 Validating Crew Pairing Distributions

All the calibration and validation results in Section 5.2 have measured the similarity between two crew pairing solutions in terms of the differences in the values of penalty function components in the objective function. However, these penalty function components themselves are functions of the calibrated parameters. To avoid this dependence and to make our comparisons fairer, we perform additional validation of our results by directly comparing the distributions of the features that affect crew-propagated delays and disruptions for our results against the distributions of those features for the airline-provided crew pairing solutions. Note that, we do not expect our solution to result in precisely the same crew pairings as those scheduled by the airline. Instead, our goal is to ensure that the distributions of the features affecting crew-propagated delays and disruptions are similar between our solution and the airline-provided solution so that the two crew pairing solutions possess similar potential for crew-propagated delays and disruptions. This validation approach of ours is similar in spirit

to that used by Barnhart, Fearing and Vaze (2014) to compare the distributions of features of passenger itinerary flows.

We consider the following distributions for validation purposes.

1. Distribution of flying time in a duty.
2. Distribution of elapsed time in a duty.
3. Distribution of scheduled sit times.
4. Distribution of scheduled rest times.

Note that these correspond to Features 1 through 4 described in Section 2.2. The last two features are not included in this type of validation separately to avoid redundancy because they are simply counts of occurrences of crew base purity violations and aircraft changes within a crew duty, and hence are fully represented by the penalty function comparisons in Section 5.2.

Chi-square statistic (Lewis and Burke, 1969) and the Kolmogorov-Smirnov statistic (Bradley, 1968) are two commonly used metrics for comparing two distributions to each other. The lower the values of these statistics, the more similar are the two distributions. Table 12 compares the distributions of these four features for the crew pairing solutions generated by our model (using parameters both before and after calibration) with those of the airline-provided crew pairing solutions. We present calibration and validation results for each of the four networks. For the RC network, the calibration is performed using March 2014 dataset and validation on April 2014 dataset while for the three NLC networks, the calibration is performed using January 2014 dataset and validation is performed using February 2014 dataset. Note that we do not present the rest time distributions for the NLC-A320 and NLC-B737 networks for the reasons mentioned in Section 5.1. These results presented in Table 12 further reinforce the conclusion that the calibrated models generate crew-pairing solutions that are very similar to those provided by the airline in terms of the distributions of the potential for crew-propagated delays and

disruptions, when tested on both calibration as well as out-of-sample validation datasets. In almost all cases, the calibrated parameters yield a better fit to real-world distributions compared with the uncalibrated ones and in many cases the improvement is large. We found similar results for the several other cross-validation experiments conducted by us, but we don't present them here due to space considerations.

Table 12. Validating Distributions of Crew Pairing Solution Features

Dataset	Feature	Chi-Square Statistic		Kolmogorov-Smirnov Statistic	
		Before	After	Before	After
RC, Calibration (March 2014)	Flying Time	48.86	0.15	0.67	0.33
	Elapsed Time	94.19	36.22	0.40	0.20
	Sit Time	89.67	106.84	0.50	0.50
	Rest Time	6.95	1.27	0.33	0.33
RC, Validation (April 2014)	Flying Time	40.55	7.26	0.67	0.33
	Elapsed Time	83.15	3.58	0.40	0.20
	Sit Time	157.16	7.64	0.75	0.50
	Rest Time	6.63	17.74	0.33	0.33
NLC-A320, Calibration (January 2014)	Flying Time	73.47	1.25	0.75	0.25
	Elapsed Time	171.22	1.32	0.75	0.25
	Sit Time	42.13	13.3	0.40	0.20
NLC-A320, Validation (February 2014)	Flying Time	60.36	8.95	0.75	0.25
	Elapsed Time	160.34	10.75	0.50	0.25
	Sit Time	55.16	31.90	0.40	0.20
NLC-B737, Calibration (January 2014)	Flying Time	281.96	49.77	0.67	0.33
	Elapsed Time	284.30	21.63	0.60	0.40
	Sit Time	81.07	27.98	0.75	0.50
NLC-B737, Validation (February 2014)	Flying Time	113.51	20.09	0.50	0.17
	Elapsed Time	262.42	11.15	0.60	0.20
	Sit Time	54.66	33.15	0.25	0.25
NLC-B757, Calibration	Flying Time	3.45	0	0.25	0

(January 2014)	Elapsed Time	4.85	3.40	0.40	0.20
	Sit Time	3.67	0.37	0.50	0.25
	Rest Time	3.92	9.23	0.33	0.33
NLC-B757, Validation (February 2014)	Flying Time	2.38	2.31	0.25	0.25
	Elapsed Time	3.32	2.23	0.20	0.20
	Sit Time	0.22	2.25	0.25	0.25
	Rest Time	7.10	2.98	0.33	0.33

In Figures 1, 2 and 3 we present the histograms of sit times, duty flying times and duty elapsed times for a sample network (NLC-A320) for two sample validation months (October 2013 and April 2014), using calibration parameters from January 2014, in order to provide a visual comparison of the feature distributions. Note that the histograms for other networks and/or other time periods provide similar insights and hence are skipped due to space considerations. We divided the data into unequal sized bins to create narrower bins in the feature ranges that particularly affect crew-propagated delays and disruptions. This enables us to focus the histogram comparisons especially on the shorted sit times, which are closer to the minimum sit time limits, and longer duty flying times and longer duty elapsed times, which are closer to the corresponding maximum limits. The sit time histograms in Figure 1 indicate that in eight of the 10 bins across the two validation months, the after calibration values are closer, than the before calibration values, to those of the actual airline samples. The duty flying time histograms in Figure 2 indicate a particularly strong improvement in the similarity to actual airline samples after calibration and values in all eight bins are improve significantly after calibration. Finally, duty elapsed time histograms in Figure 3 also indicate considerably improved similarity after calibration in most bin values. These out-of-sample validation histograms further confirm our finding that the calibration process improves the similarity between the crew pairings generated by our approach and those used by the airline, when evaluated based on the actual distributions of the features that affect crew-propagated delays and disruptions.

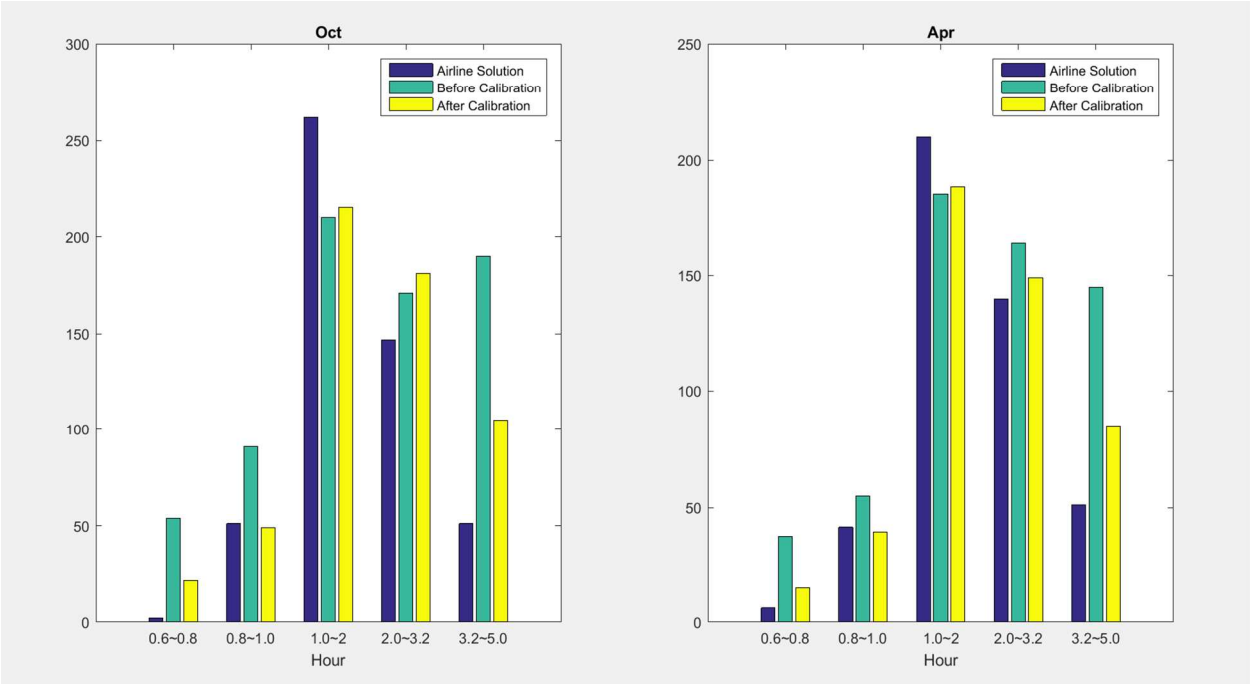


Figure 1: Sit Times Validation Histograms for NLC-A320 Network for October 2013 and April 2014

Months with Parameters Calibrated using January 2014 Data

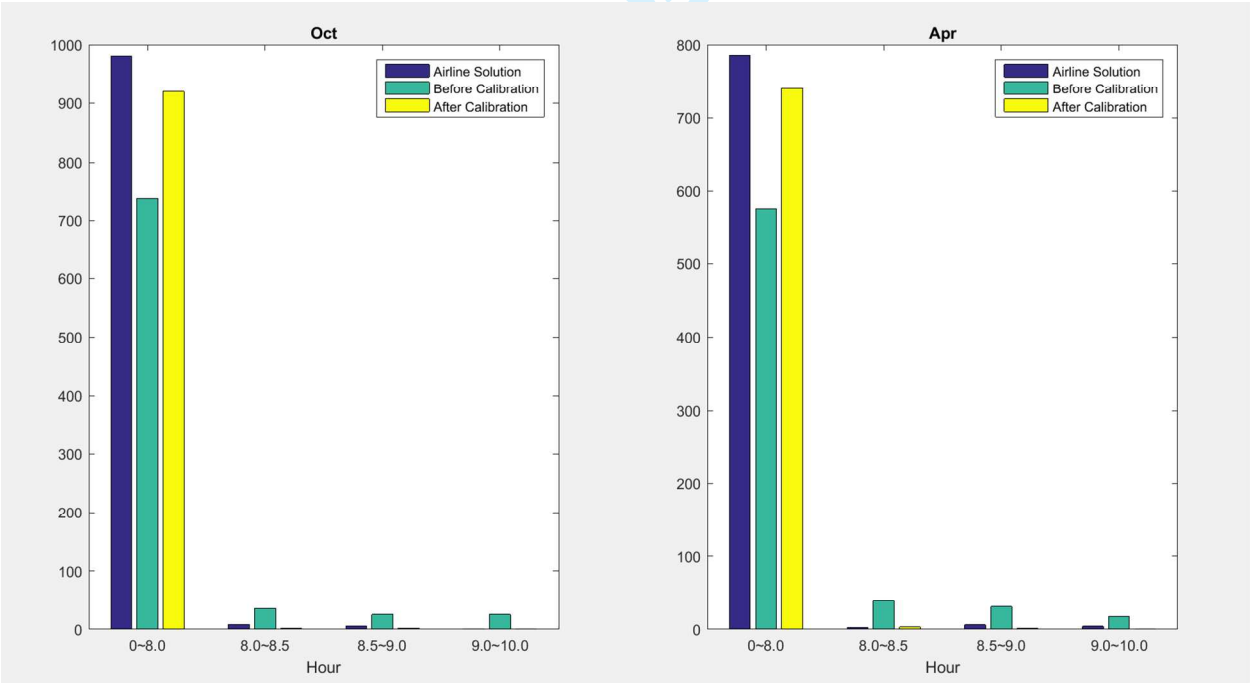


Figure 2: Duty Flying Times Validation Histograms for NLC-A320 Network for October 2013 and April 2014 Months with Parameters Calibrated using January 2014 Data

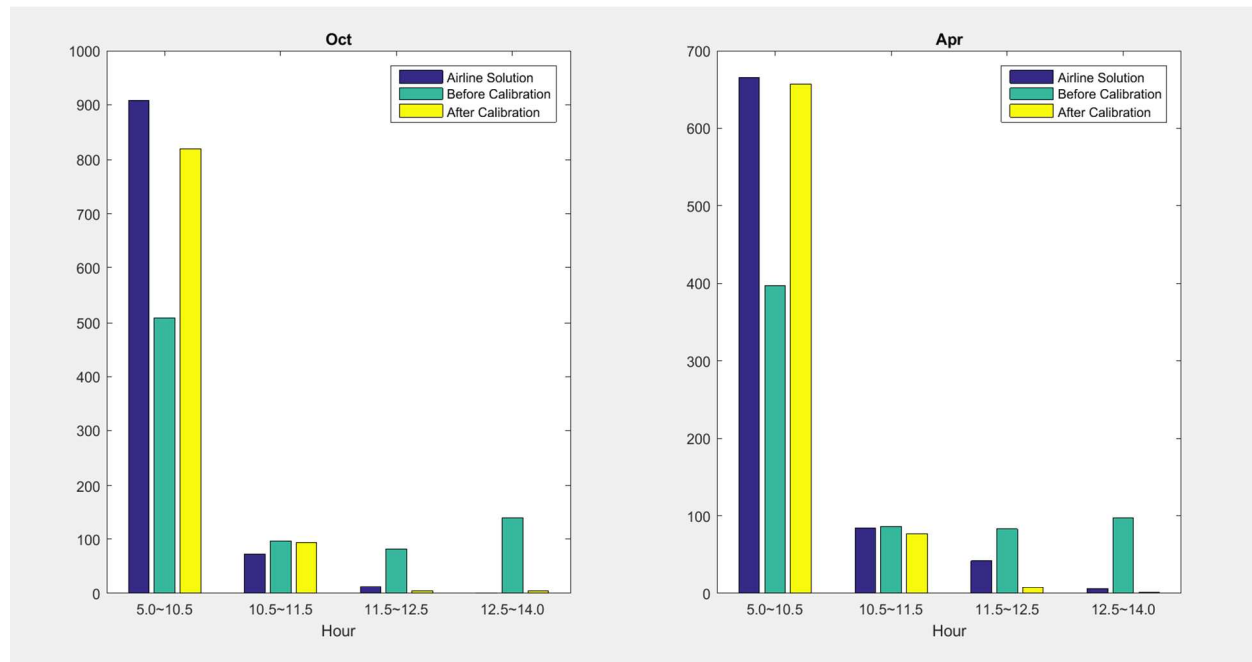


Figure 3: Duty Elapsed Times Validation Histograms for NLC-A320 Network for October 2013 and April 2014 Months with Parameters Calibrated using January 2014 Data

6 Estimation and Validation of Crew-propagated Delays and Disruptions

All the calibration and validation results presented in Section 5 were based on comparisons of different crew pairing solutions in terms of the values of penalty function components and distributions of features that are representative of the crew-propagated delays and disruptions. In this section, we focus directly on the crew-propagated delays and disruptions themselves. In Section 6.1, we provide additional validation of our results in terms of these crew-propagated delays and disruptions. Then, in

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Section 6.2 we present the process for estimating the crew-propagated delays and disruptions for any given airline, fleet family and time period, thus highlighting how the crew pairing solutions generated by our research can be used to estimate crew-propagated delays and disruptions for any given network.

6.1 Validation based on Crew-propagated Delays and Disruptions

As argued in Section 2, crew-propagated delays and disruptions depend not only on crew schedules but also on the root delays and on operational recovery actions used by the airline. As researchers, typically we are not aware of the exact set of recovery strategies and parameters that were used by the airlines because such information often tends to be confidential. Moreover, neither the researchers nor the airlines are aware of the exact levels of future root delays when planning the crew schedules for those future time periods. These limitations make it difficult to accurately calculate and compare the actual values of crew-propagated delays and disruptions. In this section, we use an existing crew-recovery software tool, named *SimAir*, for simulating airline operational recovery and for estimating crew-propagated delays and disruptions. For a detailed description of *SimAir*, the reader is referred to Rosenberger et al. (2002). Some past studies in robust crew scheduling, such as Schaefer et al. (2005), have used *SimAir* to evaluate crew schedule operational performance.

While *SimAir* is one of the most sophisticated existing tools available to researchers for simulating airline operations, its recovery approach does not include crew swaps. Thus, we expect *SimAir* to overestimate the total operational costs to some extent. Moreover, *SimAir* requires distributions of various components of root delays as inputs, but they are not easy to ascertain accurately. *SimAir* divides each flight into a variety of phases such as gate departure, taxi-out, take-off, en-route, touch-down, and taxi-in, and models delays in each phase in detail using separate delay probability distributions. For a given flight, estimating these delay distributions empirically requires making several assumptions and analyzing a lot of historical delay data. As such, this is beyond the scope of our current validation

process. Therefore, in order to provide an example of how to validate our crew pairing solutions based on crew-propagated delays and disruptions, we make assumptions about the delays in different flight phases. For the sake of simplicity, we divide all flights into two categories, namely long-haul and short-haul, and assume the delay distributions to be the same for all flights within a category. To keep the presentation simple, in this section we will focus only on the NLC-B757 network, which is the smallest across the four case study networks described in Section 4.

Table 13 provides a summary of the root delay distributions for all individual flight phases that we used as inputs to SimAir. According to the scheduled block time for each flight, we divide flights into two categories. If the scheduled block time is greater than 2 hours, then we categorize it as a long-haul flight, and otherwise we categorize it as a short-haul flight. Then the actual flying time equals the scheduled block time of the flight minus a random variable with positive expected value. Therefore, the mean of the Gaussian distributed random variable used for modeling the flying time delay is negative, as indicated in the last two rows of Table 13. Values detailed in Table 13 are the same as those suggested as default values in the SimAir User's Manual (2003).

Table 13. Summary of the SimAir Delay Distributions

Event	Distribution	Mean (Standard Deviation)	Unit
Departure Gate Delay	Gaussian Distribution	1 (3)	Minutes
Take-off Service Rate	Gaussian Distribution	20 (2)	Take offs per hour
Touch-Down Service Rate	Gaussian Distribution	20 (2)	Touch downs per hour
Taxi-Out Duration	Constant	5	Minutes
Taxi-In Duration	Constant	5	Minutes
Flying Time Delay (Long-Haul)	Gaussian Distribution	-20 (5.948)	Minutes
Flying Time Delay (Short-Haul)	Gaussian Distribution	-20 (5.665)	Minutes

Using SimAir, we test three crew pairing solutions for the NLC-B757 network, namely, the one actually used by the airline, the one generated by solving the deterministic crew pairing optimization problem (i.e., setting all parameters to 0), and the one generated by our approach using parameters after calibration. For each of these three crew pairing solutions, we conduct 2000 simulation runs. We test the stability of our simulation results by calculating the 95% confidence interval for 15-minute flight on-time performance for the three crew pairing solutions. They are found to be [88.14%, 88.44%], [86.13%, 86.43%], and [87.73%, 88.05%], respectively, for the airline sample, pre-calibration crew pairing solution, and post-calibration crew pairing solution. These narrow intervals not only demonstrate the stability of our results but also show that the on-time performance for the airline’s actual crew pairing solution is much closer to that obtained by our approach after calibration than before calibration. Table 14 provides full details of SimAir result comparisons across the three crew pairing solutions. First row of Table 14 lists the number of simulation runs while the remaining nine rows list various simulation result summary statistics. Interestingly, we find that the airline sample simulation results are closer to after calibration results than the before calibration results in terms of each of these nine criteria, and in many cases the similarity improvement due to calibration is significant. This further validates the effectiveness of our approach.

Table 14. SimAir-Based Validation Results

	Airline Sample	Before Calibration	After Calibration
Number of Simulation Runs	2000	2000	2000
15-Minute On-time Performance	88.29%	86.28%	87.89%
Number of Reserve Crew Calls	0	1.99	0.05
Number of Crew Deadheads	0	0.98	0.05
Flight Legs Ferried	3.02%	3.46%	2.95%
Flights Delayed ≤ 0 min	29.41%	28.49%	29.23%
Flights Delayed (0,15] min	58.88%	57.79%	58.66%

Flights Delayed (15,45] min	3.47%	3.37%	3.52%
Flights Delayed > 45 min	0.96%	0.77%	0.83%
Cancelled Flights	7.28%	9.58%	7.76%

6.2 Estimating Crew-propagated Delays and Disruptions

In this paper, we developed an approach to generate crew pairing solutions that are similar to the actual crew pairing solutions used by the airlines in the real world, in terms of their potential for crew-propagated delays and disruptions. As mentioned in Section 1, this work has at least three main types of applications. First, it is the first step toward estimating the extent to which delays and disruptions propagate through crew connections. Second, it allows us to assess and compare the effectiveness of various operational recovery strategies used by the airlines. Finally, it allows us to evaluate and compare the full impact of various candidate strategies for congestion and delay mitigation that are being considered by the airlines, airports, air traffic control system, and the government. In this section, we briefly describe how each of these objectives can be achieved using our results.

Table 2 in Section 5.1 of this paper provides four different sets of parameters representing four different airline networks. The robust crew pairing optimization model (8-10) can be solved for each of these four sets of parameters to come up with an estimated crew schedule for any given airline network of interest. As our results in Section 5 indicate, when picking the right set of parameters, it is advisable to choose a set of parameters that corresponds to a network which is most similar (in terms of airline type, fleet type and time period) to the network of interest. As found in Section 5, however, no matter which parameter set is picked, using calibrated parameters gives a far better fit than solving the deterministic crew pairing model, in all cases. Alternatively, the calibration approach described in Section 3 and 4 could be used to generate more suitable parameters in case a better-matching airline crew schedule sample is available.

Once the crew schedules are estimated, they can be used to estimate the crew-propagated delays and disruptions. Note that, for accurately estimating delays in a historical dataset, some knowledge or assumption regarding the recovery strategies used by the airlines is necessary. For a given set of root delays and for a given operational recovery strategy, our crew schedules can be used to estimate historical crew-propagated delays and disruptions in a relatively straightforward manner. On the other hand, for evaluating and comparing different operational recovery strategies and different congestion or delay mitigation strategies, a root delay simulator such as SimAir can be used in combination with our estimated crew schedules to calculate the full extent of delay reduction achievable by these strategies. Note that, in all cases, total propagated delays and disruptions should be measured by accounting for the propagation through aircraft as well as crew connections. However, the aircraft connections are publicly available and hence are not a bottleneck to this overall process.

7 Conclusion and Future Research

In this paper, for the first time, the inverse of the robust crew pairing generation problem was presented, formulated and solved in order to gain insights into the extent of robustness of the real-world airline crew scheduling processes. The problem was formulated as one of learning the parameters of the robust optimization objective function using real-world airline crew scheduling samples. A heuristic solution approach was developed and implemented. It involved solving the forward problem (the robust crew pairing problem) repeatedly to minimize a similarity measure between the solution of the robust crew pairing problem and the actual airline crew schedule samples by identifying the optimal set of objective function parameters. The forward problem minimizes the sum of planned cost and penalty costs which penalize the crew pairings for six different features that make them vulnerable to propagation of delays and disruptions. In our case studies, this sub-problem itself represented the largest networks in the literature for which a robust crew pairing problem has been solved. A sequence

of exact methods and heuristic ideas were used to solve this robust crew pairing problem to near-optimality. This allowed the overall parameter calibration problem to be solved in reasonable amount of time.

Several new insights were obtained into the airline crew pairing generation process. First, compared with the crew pairings obtained by solving the deterministic crew pairing problem the calibrated parameters led to crew pairings that are considerably closer to the actual airline crew schedules in all our experiments. In most cases, the accuracy improvement was substantial. This suggests that airlines do take into account robustness or the potential for propagation of delays and disruptions when creating their crew schedules. Furthermore, we found that the crew pairings calibrated using four different airline networks perform similar to each other, and much better than the deterministic crew pairing solutions, in terms of their closeness to the actual crew schedules, even when the calibration and evaluation is not conducted on the same network. This suggests that the calibrated parameters are relatively stable, and that even in cases where the data available for model training is for a network somewhat dissimilar to the one of interest, it is better to use the calibrated parameters than the uncalibrated ones. However, for maximizing estimation accuracy, whenever possible, we found it advisable to use parameters calibrated with a network that is as similar to the one of interest as possible, in terms of airline type, fleet type and time period. Finally, this paper presented, for the first time in literature, a measure of the tradeoff as perceived by actual airlines between the crew salary costs and costs of crew-propagated delays and disruptions as reflected by the calibrated robust crew pairing objective functions. Across the four networks, the ratio of the penalty costs (representing the costs of crew-propagated delays and disruptions) and the crew salary costs was found to lie between 0.5% and 4%. Note that this is inferred based on the crew pairings used by the airlines and not based on the actual costs of these delays and disruptions.

In addition to these insights, as described in Section 6.2, this research makes the estimated crew pairing solutions available for further research and analysis. We have made this entire model calibration code as well as the resulting calibrated crew pairing solutions publicly available for future research. These estimated crew pairing solutions are useful to gauge the extent of delays and disruptions that propagate across the airline networks through crew connections. While these crew pairing estimates do give a starting point to estimate the crew-propagated delays and disruptions, the important next step toward accurately estimating historical delay propagation is to develop an understanding of crew recovery strategies used by the airlines in the real world. Once we have access to a historical sample of actual crew recovery actions, then a framework similar to the one developed in this paper could be used to learn the airline crew recovery optimization process as well. This will be the next phase step in our research project.

Acknowledgments

This research is sponsored by the Federal Aviation Administration’s National Center of Excellence for Aviation Operations Research (NEXTOR). We are also thankful to ILOG for providing CPLEX licenses to conduct our computational experiments.

Reference

AhmadBeygi S, Cohn A, Lapp M (2010) Decreasing airline delay propagation by re-allocating scheduled slack. *IIE Trans.* **42**(7): 478–489.

AhmadBeygi S, Cohn A, Weir M (2009) An integer programming approach to generating airline crew pairings. *Comput. Oper. Res.* **36**(4): 1284–1298.

Ahuja RK, Orlin JB (2001) Inverse optimization. *Oper. Res.* **49**(5):771–783.

Barnhart C, Fearing D, Vaze V (2014) Modeling passenger travel and delays in the national air transportation system. *Oper. Res.* **62**(3): 580–601.

- 1
2
3 Barnhart C, Johnson E, Nemhauser G, Savelsbergh M, Vance P (1998) Branch-and-price: Column
4 generation for solving huge integer programs. *Oper. Res.* **46**(3): 316–329.
5
6
7
8 Barnhart C, Vaze V. (2015a) Ch.10. Irregular Operations: Schedule Recovery and Robustness.
9
10 Belobaba P, Odoni A, Barnhart C, eds. *The Global Airline Industry*, 2nd ed (John Wiley & Sons,
11 West Sussex), 263-287.
12
13
14 Barnhart C, Vaze V. (2015b) Ch.8. Airline Schedule Optimization. Belobaba P, Odoni A, Barnhart C, eds.
15 *The Global Airline Industry*, 2nd ed (John Wiley & Sons, West Sussex), 189-222.
16
17
18
19 Bradley J (1968) Distribution-Free Statistical Tests. (Prentice Hall, Englewood Cliffs, NJ).
20
21
22 Bureau of Transportation Statistics (BTS) (2016) TranStats. On-Time Performance:
23 http://www.transtats.bts.gov/Fields.asp?Table_ID=236
24
25
26 Cacchiani V, Salazar-González J-J (2016) Optimal Solutions to a Real-World Integrated Airline Scheduling
27 Problem. *Transportation Sci.* Articles in Advance.
28
29
30
31 Cockpit Work Rule Comparison Fall 2010 (2010).
32
33 Desaulniers G, Desrosiers J, Solomon M (2005) *Column Generation* (Springer, New York).
34
35
36 Duan Z, Wang L (2011) Heuristic algorithms for the inverse mixed integer linear programming problem.
37 *Journal of Global Optimization.* **51**(3):463-471.
38
39
40 Dunbar M, Froyland G, Wu C (2012) Robust airline schedule planning: Minimizing propagated delay in
41 an integrated routing and crewing framework. *Transportation Sci.* **46**(2):204–216.
42
43
44 Engel F (1995) Summary over test runs, Internal report, Carmen Systems AB, Gothenburg, Sweden.
45
46
47 Ehrgott M, Ryan DM (2002) Constructing robust crew schedules with bicriteria optimization. *J. Multi-*
48 *Criteria Decision Anal.* **11**(3):139–150.
49
50
51 Gao C, Johnson E, Smith B (2009) Integrated airline fleet and crew robust planning. *Transportation Sci.*
52 **43**(1): 2–16.
53
54
55
56
57
58
59
60

- Garfinkel R, Nemhauser G (1970) Optimal political districting by implicit enumeration techniques. *Management Sci.* **16**(8):495–508.
- IATA (International Air Transport Association) Annual Review 2015. Available at <https://www.iata.org/about/Documents/iata-annual-review-2015.pdf>.
- Irnich S, Desaulniers G (2005) Shortest path problems with resource constraints. Desaulniers G, Desrosiers J, Solomon MM, eds. *Column Generation* (Springer, New York), 33–65.
- Jacquillat A, Vaze V (2016) Inter-airline Equity in Airport Scheduling Interventions. Under Review.
- Jacquillat A, Odoni A (2015) An Integrated Scheduling and Operations Approach to Airport Congestion Mitigation. *Oper. Res.* **63**(6):1390-1410.
- Klabjan D, Johnson E, Nemhauser G (2001) Solving large airline crew scheduling problems: Random pairing generation and strong branching. *Computational Optim Appl.* **20**(1):73–91.
- Kasirzadeh, A, Saddoune, M, Soumis, F (2015) Airline crew scheduling: Models, algorithms, and data sets. *EURO Journal on Transportation and Logistics*. doi:10.1007/s13676-015-0080-x.
- Lan S (2003) Planning for robust airline operations: Optimizing aircraft routings and flight departure times to achieve minimum passenger disruptions. Ph.D. dissertation, Massachusetts Institute of Technology, Cambridge, MA.
- Lan S, Clarke J-P, Barnhart C (2006) Planning for robust airline operations: Optimizing aircraft routings and flight departure times to minimize passenger disruptions. *Transportation Sci.* **40**(1): 15–28.
- Lewis D, Burke C (1969) The use and misuse of the chi-square test. *Psych. Bull.* **46**: 433–489.
- Mercier A, Cordeau JF, Soumis F (2005) A computational study of Benders decomposition for the integrated aircraft routing and crew scheduling problem. *Comput. Oper. Res.* **32**:1451–1476
- Petersen JD, Sölveling G, Clarke J-P, Johnson EL, Shebalov S (2012) An optimization approach to airline integrated recovery. *Transportation Sci.* **46**(4):482–500.

- Rosenberger J, Schaefer A, Goldsman D, Johnson E, Kleywegt A, Nemhauser G (2002) A stochastic model of airline operations. *Transportation Sci.* **36**:357–377.
- Lamperski J, Schaefer A (2015) A polyhedral characterization of the inverse-feasible region of a mixed-integer program. *Oper. Res. Lett.* **43**(6): 575-578.
- Rosenberger J, Johnson E, Nemhauser G (2003) Rerouting aircraft for airline recovery. *Transportation Sci.* **37**(4): 408–421.
- Shebalov S, Klabjan D (2006) Robust airline crew pairing: Move-up crews. *Transportation Sci.* **40**(3): 300–312.
- Schaefer A, Johnson E, Kleywegt A, Nemhauser G (2005) Airline crew scheduling under uncertainty. *Transportation Sci.* **39**(3):340–348.
- SIMAIR User's Manual (2003) Available at <https://www.ise.nus.edu.sg/project/simair/download/manual.html>.
- Smith B, Johnson E (2006) Robust airline fleet assignment: Imposing station purity using station decomposition. *Transportation Sci.* **40**(4):497–516.
- Tam B, Ehrgott M, Ryan D, Zakeri G (2011) A comparison of stochastic programming and bi-objective optimization approaches to robust airline crew scheduling. *OR Spectrum.* **33**(1):49-75.
- U.S. G.A.O (2008) Commercial Aviation: Impact of Airline Crew Scheduling on Delays and Cancellations of Commercial Flights. GAO-08-1041R. United States.
- Vance P, Atamturk A, Barnhart C, Gelman E, Johnson E, Krishna A, Mahidhara D, Nemhauser G (1997) A heuristic branch-and-price approach for the airline crew pairing problem. Technical Report Technical Report LEC-97-06, Georgia Institute of Technology.
- Vaze V, Barnhart C (2012) Modeling airline frequency competition for airport congestion mitigation. *Transportation Sci.* **46**(4): 512–535.
- Wang L (2013) Branch-and-bound algorithms for the partial inverse mixed integer linear programming

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Problem. *Journal of Global Optimization*. **55**(3):491-506.

Weide O, Ryan D, Ehrgott M (2010) An iterative approach to robust and integrated aircraft routing and crew scheduling. *Comput. Oper. Res.* **37**(5):833–844.

Yen JW, Birge JR (2006) A stochastic programming approach to the airline crew scheduling problem. *Transportation Sci.* **40**(1):3–14.

For Review Only