

# Conformal Automation for Air Traffic Control using Convolutional Neural Networks

S. J. van Rooijen, J. Ellerbroek, C. Borst, E. van Kampen  
Control and Simulation, Faculty of Aerospace Engineering  
Delft University of Technology (TU Delft)  
Delft, The Netherlands

**Abstract**—Lack of trust has been identified as an obstacle in the introduction of workload-alleviating automation in air traffic control. The work presented in this paper describes a concept to generate individual-sensitive resolution advisories for air traffic conflicts, with the aim of increasing acceptance by adapting advisories to different controller strategies. These personalized advisories are achieved using a tailored convolutional neural network model that is trained on individual controller data. In this study, a human-in-the-loop experiment was performed to generate datasets of conflict geometries and controller resolutions, with a velocity obstacle representation as a learning feature. Results show that the trained models can reasonably predict command type, direction and magnitude. Furthermore, a correlation is found between controller consistency and achieved prediction performance. A comparison between individual-sensitive and general models showed a benefit of individually trained models, confirming the strategy heterogeneity of the population, which is a critical assumption for personalized automation.

**Keywords**—Strategic conformance, machine learning, solution space diagram, velocity obstacles, consistency, decision-support

## I. INTRODUCTION

The introduction of automated decision support tools for Air Traffic Control (ATC) is widely seen as unavoidable to keep up with air traffic growth [1], [2]. However, in the past, the introduction of support tools for separation management has shown that such tools are frequently left unused [3]–[5]. Westin argues that automation conformance can be one of the reasons for low acceptance of automated advisories [6]. For instance in the separation task, often multiple resolution options are available for a given conflict, and resolutions proposed by automation do not always coincide with human strategy [7].

Recent work has argued that making advisories *strategically conformal* can be an effective way to increase trust and acceptance of automation [6], [8]. If a decision support tool has the ability to adapt to individual strategies, acceptance should increase. A promising method to create automation that is able to predict human strategies is machine learning [9]. A decision support tool for ATC based on machine learning would be able to adapt to different controller preferences without full knowledge of the underlying decision-making dynamics.

A proof-of-concept of this by Regtuit et al. has shown that machine learning techniques are able to identify and replicate conflict resolution strategies from simple, synthetic traffic situations [10]. While the results from this study are promising, its method is not easily extended to capture realistic, multi-aircraft situations. The work presented in this paper proposes and tests a machine learning model that should be able to predict Air Traffic Controller (ATCo) commands based on a realistic traffic situation.

To capture relevant aspects of realistic traffic situations in a usable form, the Solution Space Diagram (SSD) [11], also

known as velocity obstacle diagram [12], will be used as visual input to the learning algorithm. The SSD is a visual representation of a velocity space of conflicting and conflict-free velocities. Numerous previous studies have shown that an SSD image contains sufficient information concerning an air traffic situation to make an informed decision regarding conflict resolution [11], [13]–[15]. Learning from an image is advantageous because it eliminates the need for manually designing features based on heuristics and assumptions. The proposed model will use *convolutional neural networks*, a machine learning technique that is especially well-suited for visual learning [16].

To provide training data, and to test controller consistency and the SSD as a machine learning feature, a human-in-the-loop experiment is carried out. The experimentally generated datasets are used to train individual-sensitive predictive models using a supervised learning algorithm. Individual model performance is compared to controller consistency and an inter-participant comparison and a comparison with general group-based models are performed to test the effect of individual-sensitive automation.

The remainder of this paper is structured as follows: The proposed concept is introduced in section II. Experiment design is discussed in section III, and results are provided in section IV, followed by a discussion in section V and conclusions & recommendations in section VI.

## II. CONCEPT

Achieving conformal automation has similarities to a machine learning field called Imitation Learning, where machine learning is used to learn to perform tasks based on expert demonstrations. The aim is to generalize these observations to unseen situations, usually using supervised learning [17].

Artificial neural networks with a deep architecture are a powerful method for supervised learning as they automatically devise learning features, without reliance on pre-engineered parameters [18]. This research follows a similar approach, where the input images are provided in the form of rasterized SSDs. In particular, Convolutional Neural Networks (CNNs) [19] have shown their potential of training on image data, for example in learning to play computer games [20], self-driving cars [21], [22] and competing with world-class board game champions [23].

Similar to a regular neural network, a CNN consists of multiple stacked neurons. In the case of a CNN, every input is connected to a pixel-value of the original image, as shown in Figure 1. On each layer of the network, individual neurons are connected to the weighted combination of the outputs of the previous layer, which is passed through a convolution filter. These filters slide over the entire image to create a

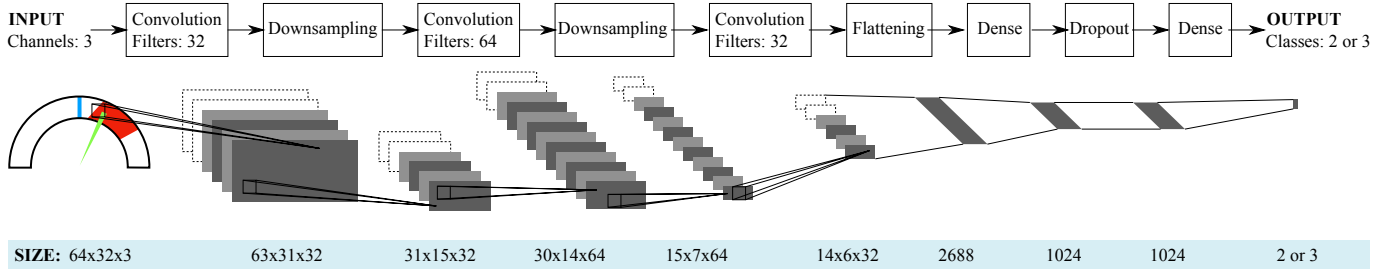


Fig. 1: The Convolutional Neural Network structure used in this research. All weights in one plane are identical. Size values are given in pixels. Adapted from LeCun (1998) [16].

new output map. By using different filters, specific visual features can be extracted from the input image, such as small corners or edges. In this study, filter size is always  $2 \times 2$  pixels, and filters progress over the image with a stride of 1 pixel. The visual features are subsequently combined to form compositions while progressing through the net, using successive layers of filters. To avoid overfitting and reduce computational cost, the data is also downsampled in between filter steps [24]. In the final steps, the data is flattened, and reduced in steps to the expected number of output classes (e.g., when a conflict resolution maneuver can be either using heading, speed, or a direct-to, output classes are *heading*, *speed*, and *direct-to*, resulting in an output size of three). To reduce overfitting, a dropout layer is added here that sets the weights of a fraction of the neurons to zero at each epoch during training [25].

The neurons in each layer are specified by an activation function. In this study, all layers but the final layer use a Rectified Linear Unit (ReLU) activation function, because of its computational efficiency [26]. The final layer uses a softmax function [27], which transforms the layer of real values into a vector of probabilities per output class  $\sigma(\mathbf{z})$ . The softmax function is defined by Equation 1, where all entries of  $\sigma$  are real, within  $[0, 1]$ , and add up to 1.  $K$  is the dimension of input vector  $\mathbf{z}$ .

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \text{ for } j = 1, \dots, K \quad (1)$$

A cross-entropy function is used to take these probabilities into account in the loss function [28]. The cross entropy  $H$  is calculated for  $M$  output classes by comparing the probability vector  $\sigma$  resulting from the softmax function to the one-hot encoded target vector  $y_i$ :

$$H(y, \sigma) = - \sum_i^M y_i \log \sigma_i \quad (2)$$

The calculated losses, averaged over a minibatch of samples are used to update the network parameters  $\theta$ . In this study, this update is done using a first-order gradient-based optimization algorithm called Adam. Empirical results show that it outperforms previously popular optimizers such as AdaGrad, RMSProp and SGDNesterov [29].

The performance of CNNs is naturally determined by their

inputs, as they should capture all relevant information for the model to make a prediction. This research utilizes the Solution Space Diagram (SSD) as input to the neural network. The SSD was originally designed as a decision support tool that integrates various critical parameters of the Conflict Detection and Resolution problem [30]–[32]. This was later extended to complexity and workload analysis [11], [33]–[35], and automated conflict resolution [15], [36]. Based on the findings of these studies, this research hypothesizes that an SSD image as learning feature contains sufficient information concerning air traffic conflicts to make an informed decision.

Figure 2 illustrates how the SSD is constructed. It consists of an area of reachable velocities, bounded by concentric limits of minimum and maximum operating speeds. This reachable space is reduced by triangular velocity obstacles that correspond to the set of velocity vectors that would lead to a loss of separation with a nearby aircraft. The color of the velocity obstacles is determined by the time to closest point of approach (CPA), divided into three ranges: red ( $t_{cpa} < 60s$ ), orange ( $60 < t_{cpa} < 120s$ ) and gray ( $t_{cpa} > 120s$ ). The remaining free space corresponds to the reachable conflict-free speed/heading combinations. In the SSD, the current heading and speed are indicated by the green vector, see Figure 2. As an additional feature, the target heading to the exit waypoint of the respective aircraft is shown with a blue line. This information is part of the normal tasks of an air traffic controller, and may influence conflict resolution decisions.

In the current study, the above method is used to learn individual controller strategies for observed traffic conflicts. Following Westin’s approach [37], this study defines strategy in terms of three control variables: resolution *type* (heading, speed, or direct-to), resolution *direction* (left or right, speed increase or decrease), and resolution *magnitude* ( $[0, 10]$  deg,  $[10, 45]$  deg and  $> 45$  for heading resolutions, and  $[200 - 250]$  kts, and  $[250 - 290]$  kts for speed resolutions) With the SSDs of all detected conflicts as input, three independent CNNs (as described in Figure 1) are trained to match the SSD of the controlled aircraft to the observed controller strategy: one CNN for resolution type, one for resolution direction, and one for resolution magnitude.

### III. EXPERIMENT

An experiment was performed to provide training and testing data for the model described in the previous section. In this experiment, participants were asked to control a sector with

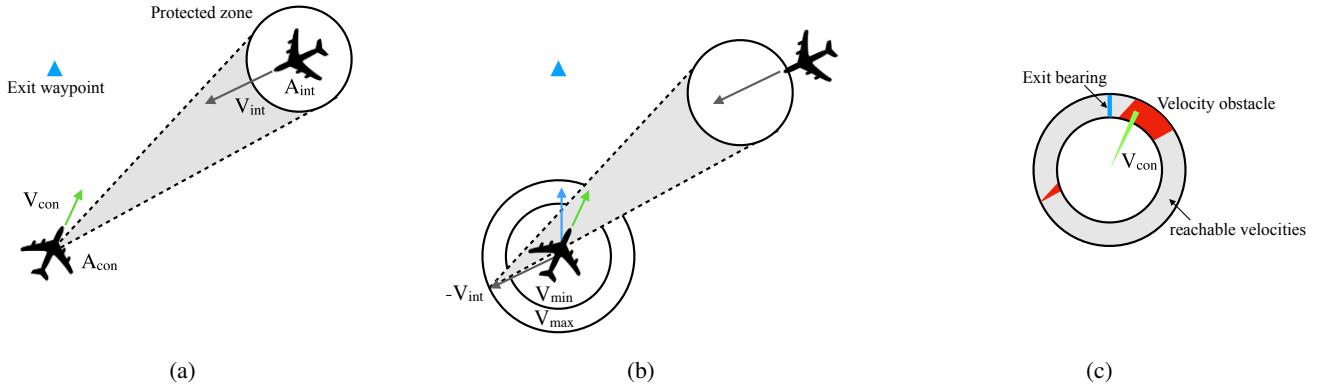


Fig. 2: Construction of the SSD: (a) The grey area indicates the set of conflicting relative velocities. (b) The area is displaced by the velocity vector of  $A_{intruder}$  to generate the velocity obstacle of  $A_{intruder}$  for  $A_{controlled}$ , (c) The SSD is created by limiting the solution area with the minimum ( $V_{min}$ ) and maximum ( $V_{max}$ ) velocity of  $A_{controlled}$ . Adapted from Mercado et al. (2010) [11].

multiple incoming aircraft, where they guided the aircraft to their exit waypoint as efficiently as possible, while avoiding losses of separation. The shape of the sector – inspired by Amsterdam Sector South 1 – is shown in Figure 3. Maneuvering was restricted to the horizontal plane, i.e., all aircraft flew at the same flight level and only heading, speed, and direct-to commands were allowed by means of a command interface (Figure 4). These restrictions decreased the solution space and degrees of freedom, which enabled better comparison between controllers in terms of consistency and strategy.

#### A. Conditions

Two traffic scenarios (S1 and S2) were used in the experiment with identical sector geometry but different traffic flows. Each scenario consisted of 10 conflict pairs. The main aircraft flow was always directed towards the north, and was crossed by traffic on several headings from the east. All conflicts were crossing conflicts, with conflict angles between 45 and 135 degrees. Every conflict angle in the set  $\{45, 55, \dots, 135\}$  was visited at least four times during the entire experiment. The conflicts were chronologically spaced in a way to minimize interference between conflicts.

Both scenarios were performed two times, resulting in a total of four 20-minute scenarios, generating 80 minutes of data per participant. The order of the conditions was randomized between participants to avoid learning biases.

#### B. Participants

The population consisted of 12 participants. All participants were novices, with varying experience in performing ATC control tasks. Half of the participants had knowledge of ATC concepts, but no working experience in controlling a sector. The other half of the participants previously participated in an ATC introductory course at the Dutch Aerospace Research Laboratory NLR, and therefore had some experience in controlling air traffic.

#### C. Procedure

The experiment was performed on a personal computer with a 30" screen showing the sector under control, and a touch control device window similar to the touch control devices used at the Dutch air navigation service provider LVNL.

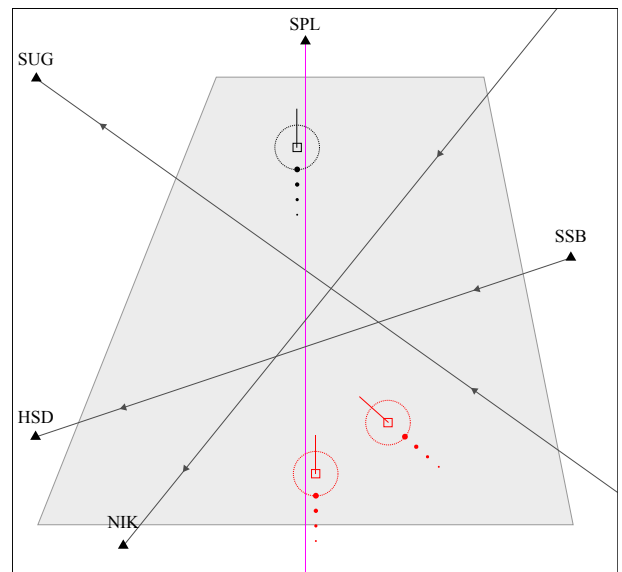


Fig. 3: The 50nm  $\times$  60nm sector as displayed in ATC simulator SectorX. The magenta lines depict the main traffic flows north and west-bound. Three aircraft are visible of which two are in conflict. The circles surrounding the aircraft indicate the protected zones ( $D = 5\text{nm}$ ).

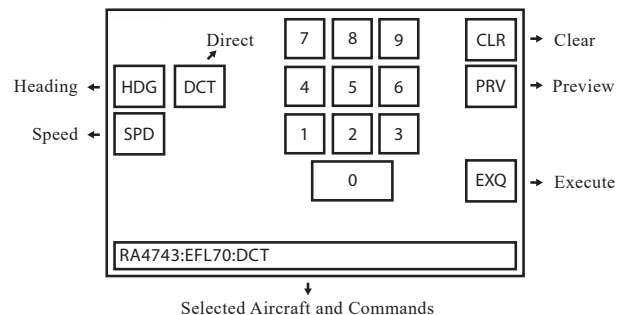


Fig. 4: The command interface that participants used in controlling the aircraft.

Participants could control aircraft through this interface using a computer mouse. Before measurements, each participant performed three training runs of increasing difficulty to get acquainted with the ATC simulator. Training runs lasted 90s, but could be prolonged when required.

Each time a command was given, the application saved the command together with all aircraft states at that moment plus an image of the controlled aircraft's SSD. Commands were defined in terms of *type* (heading, speed, or direct-to), *direction* (left or right, speed increase or decrease) and *magnitude* (magnitude of the heading or speed change). During supervised model training, the SSD image functions as input while the given command functions as target.

#### D. Data preprocessing

The SSD dataset consists of 128x128 pixel RGB samples which were preprocessed for computing time improvements. First, the SSD images were rotated track-up, so that the velocity vector of the controlled aircraft always points upward. Second, because the most relevant information in the SSD is located in the upper half of the image (aircraft are not likely to make turns larger than 90 deg), the lower half of the SSD was cropped away to decrease training times.

#### E. Dependent measures

To avoid misleading results from an imbalanced dataset, model performance is evaluated using the *Matthews Correlation Coefficient (MCC)* [38] for training and validation. For example, if a controller chooses to use heading instead of speed commands in 95% of the situations, simply always guessing 'heading' would result in a 0.95 accuracy score. This would give a false sense of model performance, because 0% of the speed commands are predicted. In such cases, the MCC metric gives less biased results than accuracy [39]. The definitions of accuracy and the two-class MCC are shown in Equations (3) and (4) respectively, where  $TP$  = true positive,  $TN$  = true negative,  $FP$  = false positive,  $FN$  = false negative.

$$\text{accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (3)$$

with possible values [0, 1].

$$\text{MCC} = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (4)$$

with possible values [-1, 1]. As MCC is a more critical metric, MCC values are often lower compared to accuracy values for identical models.

In this study, consistency is considered as the relative occurrence of the most often-used command. For resolution *type* and resolution *direction* this is determined by summing over all commands per participant:

$$\text{consistency} = \max\left( \frac{\sum_{\text{class 1}}}{\sum_{\text{class 1} + \dots + \text{class n}}}, \dots, \frac{\sum_{\text{class n}}}{\sum_{\text{class 1} + \dots + \text{class n}}} \right) \quad (5)$$

where for example class 1, 2, and 3 are respectively HDG, SPD, and DIRECT-TO for command *type* consistency. Using this definition, consistency equals 1 when a single command type

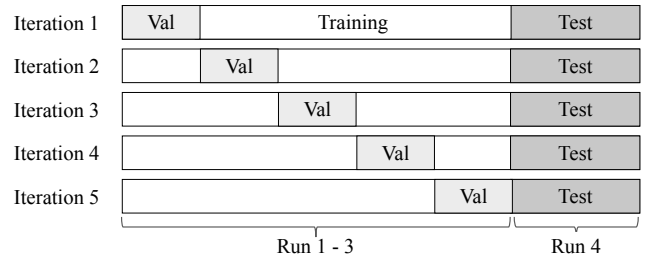


Fig. 5: Stratified 5-fold cross-validation with a separate test set.

is used and 0.5 when they are evenly balanced. For resolution *magnitude*, a different equation is used:

$$\text{consistency} = \frac{\sum \text{unique values possible}}{\sum \text{unique values used}} \quad (6)$$

where the consistency of a participant decreases when a broader range of magnitude values is used.

#### F. Model training

The data from the experiment is used in two ways: First, the dataset from each participant is used to train individual models. Then, five general models are trained on random samples drawn from the combined participant data.

The individual models are trained following the procedure illustrated in Figure 6. After preprocessing, the test data is separated into a training set (experiment runs 1-3) and a test set (experiment run 4). The training data is used to train three personalized models per participant, one for each resolution dimension (type, direction, and magnitude), see Table I for a summary of the selected training procedure settings. During training, K-fold validation is used to select a model that obtains the highest prediction performance. With K-fold validation,  $1/k^{th}$  of the data is iteratively reserved for validation, as illustrated in Figure 5. Using this validation data, the best model is then selected as a final model. This combination of model training and validation has proven to have lower variance and bias of performance measures compared to other methods [40]. In the current experiment, due to the limited quantity of data available, five folds are applied during training.

The five general models are trained on random samples of all (between-subject) data, in the same way that the individual models are trained. An equal number of random samples is used as is available in the individual model training. The mean performance of these five models will be taken, and will be compared to the individual models as a (non-individually sensitive) baseline.

#### G. Model testing

After training, each individual model is tested with the participant's corresponding test dataset (i.e. the 4<sup>th</sup> run), to test the individual prediction accuracy of the trained models, as shown in Figure 7a. Additionally, each individual model is tested with all *other* participants' test datasets, to evaluate how individually-sensitive each model is (Figure 7b). The general models are also tested with all participants' test datasets (Figure 7c).

## IV. RESULTS

Using the data from twelve participants in the human-in-the-loop experiment, twelve individual models and five general

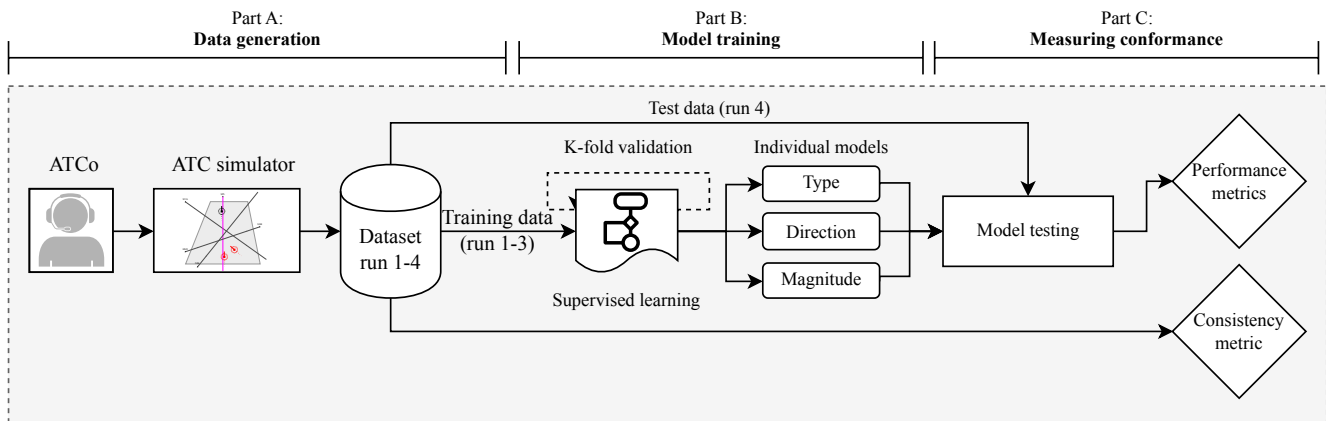


Fig. 6: Data generation and training & testing of the individual models for one participant. The dataset consists of input (SSD images) and target (commands) data. The models are used to predict a command for a given SSD image. Model performance is based on prediction accuracy.

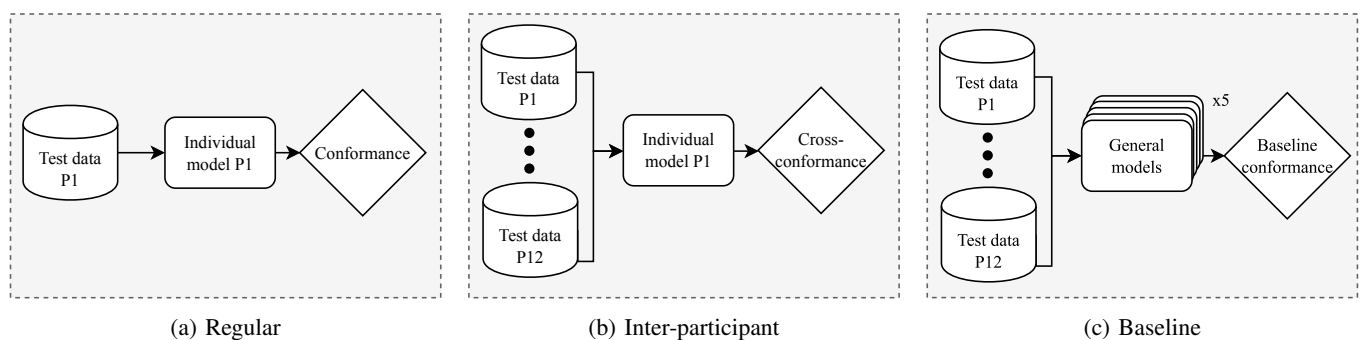


Fig. 7: Three validation steps for participant 1 (P1).

TABLE I: (Hyper)parameters during training.

Parameters	Value
Optimization algorithm	Adam
Output activation	Softmax classifier
Loss function	Categorical entropy
Train/val/test ratio	60%/15%/25%
K-folds	5
Mini batch-size	32 samples
Steps-per-epoch	$2 \times \text{training samples} / \text{batch-size}$
Epochs	30
Learning rate	0.01
Dropout rate	20%
Input image dimensions	128x128 px

models were trained. In this section, the training phase is illustrated with an example of convergence of performance in the training phase. Subsequently, this section presents the individual model results, individual model performance as a function of participant consistency, an inter-participant test of model performance, and a comparison of individual models to the average general model performance. Here, performance is measured using the *MCC* (see section III-E), which ranges between  $-1$  and  $1$ . Because negative correlation never occurred, all *MCC* result figures are clipped to a range of  $[0, 1]$ .

#### A. Training convergence

In the training phase, data from the first three experiment runs is used to train several candidate models. Using the K-fold method illustrated in Figure 5, five candidate models are trained, of which the performance is validated using five different subsets of the data. Figure 8 shows the training progress in terms of these validation results for the individual model of Participant 1, with training epoch on the x-axis, and the resulting *MCC* score on the y-axis. Here, the spread around each line depicts the range between the least and best performing folds per control variable during training, which lasts 25 epochs. It can be seen that with successive epochs, *MCC* values increase, which indicates that the neural network successfully ‘learns’ from the data samples. In most cases, the models reach *MCC* scores  $> 0.95$  during *training*, a performance level that is not achieved in the validation steps, as can be seen in Figure 8. This difference between training and validation performance indicates overfitting on the training data. The spread shows that validation *MCC* can differ more than 0.2 per fold, which is a relatively large amount compared to the mean value.

#### B. Model performance on individual test data

After training (Figure 8), the individual models are applied to the test datasets of each participant (Run 4). Figure 9 shows the achieved *MCC* scores per control variable. In this figure,

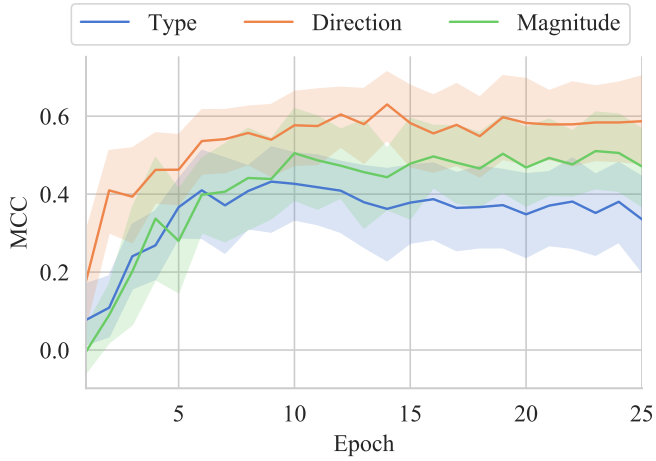


Fig. 8: Validation performance during training of P1’s individual model. The spread indicates the maximum and minimum performance for each fold per control variable.

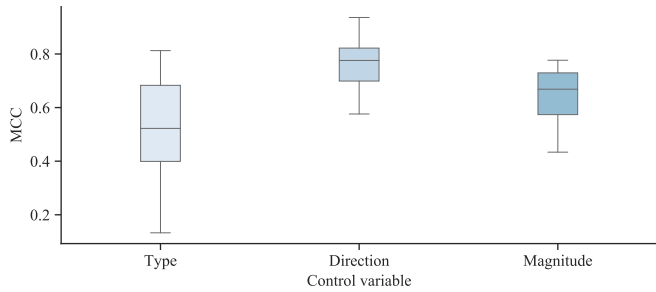


Fig. 9: Model test-performance per control variable.

the large variability in performance (particularly for the *type* control variable) indicates that the personalized predictions are not equally effective across the entire population of participants. The *direction* prediction shows the highest MCC score (mean = 0.76, SD = 0.11), while *type* (mean = 0.52, SD = 0.21) and *magnitude* (mean = 0.64, SD = 0.12) predictions achieve lower performances.

A potential reason for poor performance of the trained model is low participant consistency: in some cases, the participant data on which the model is trained does not show sufficiently consistent behaviour across different runs, and between conflicts that do appear comparable in the SSD. Figure 10 shows the normalized consistency (as defined in Section III-E) per participant and control variable. Here, it can be seen that while some participants are relatively consistent (participants 5, 7 and 10), other participants (particularly 8 and 11) show more erratic decision-making. Figure 10 also shows that participant consistency varies per control variable. For instance, participants can be very consistent in the *type* of resolution they choose, but are less consistent in the *direction* they choose for their resolutions.

The effect of participant consistency on the performance of the trained model can be evaluated by observing the correlation between consistency and model performance. To illustrate this, Figure 11 shows the mean model performance (the mean over all folds and abstraction levels), against the mean consistency per participant. When a Pearson Correlation Coefficient test is

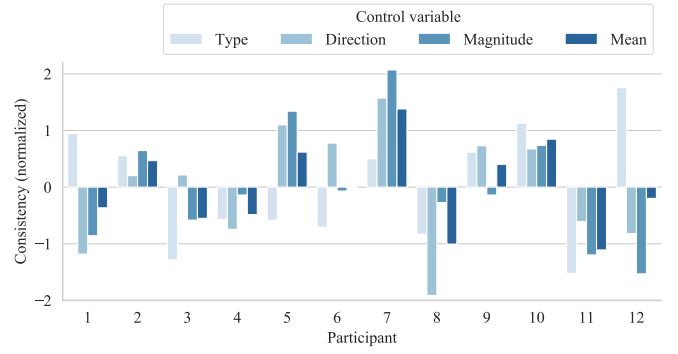


Fig. 10: Consistency scores per participant split per control variable.

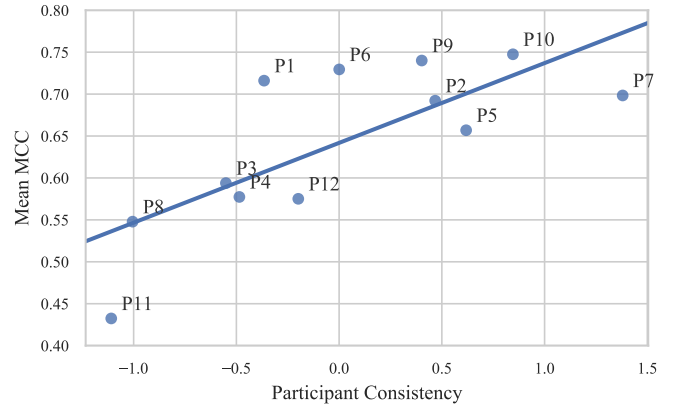


Fig. 11: Participant consistency vs individual model performance.  $R^2 = 0.56$ .

applied to this data, a positive correlation ( $r = 0.75$ ,  $p = .005$ ) can be found between participant consistency and individual model MCC. This supports the assumption that the personal models of more consistent participants perform better than the models of their less consistent counterparts.

### C. Model performance on inter-participant data

A way to evaluate whether the personalized models are indeed individual-sensitive, is to test the models against all other participant test datasets. Figure 12 shows the results of using the models of each participant on the test data of all participants. In these spider-plots, the model performance (MCC value) in terms of *type* (blue), *direction* (orange), and *magnitude* (green) is shown for each participant’s test data, along twelve radials of each chart. For the individual models of participants 1, 6, 7, 9, and 10 it can be seen that overall performance is highest when the model is applied to the test data of the corresponding participant. For instance, for the individual model of participant 1, a mean performance of MCC = 0.72 is achieved when the model is applied on the test data of participant 1, compared to an average MCC of 0.37 when testing with other participants’ data. This difference indicates that participant 1 makes different decisions in similar situations compared to the rest of the population. Other participants’ models show more uniform MCC scores, regardless of which test set is used.

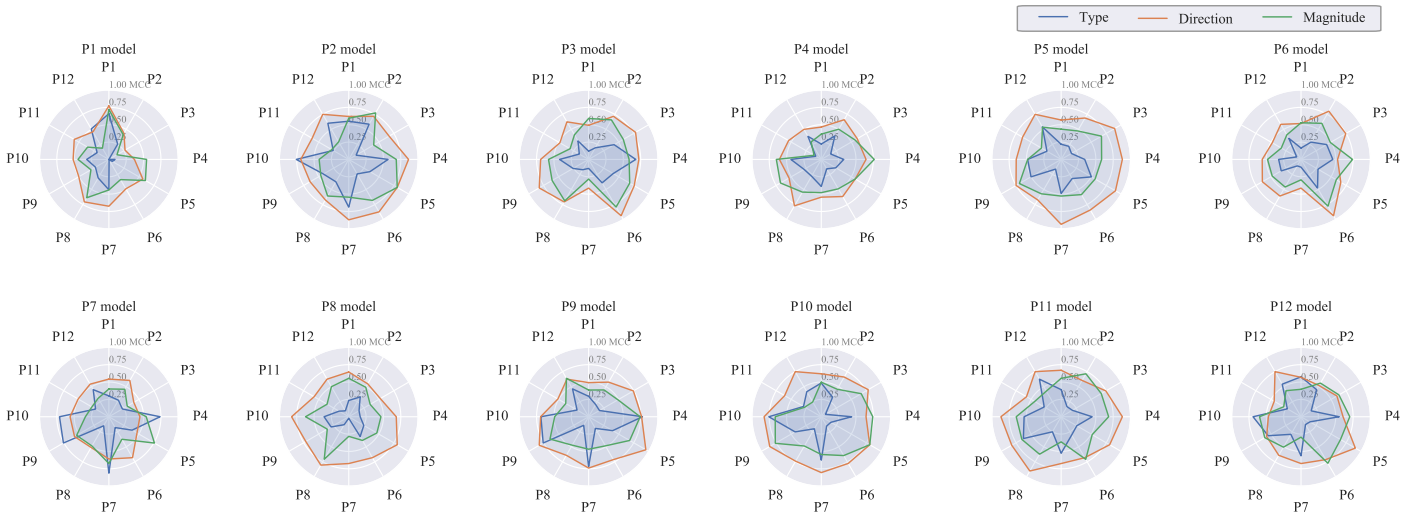


Fig. 12: Performance (in MCC) of individual models tested on the test datasets of all other participants.

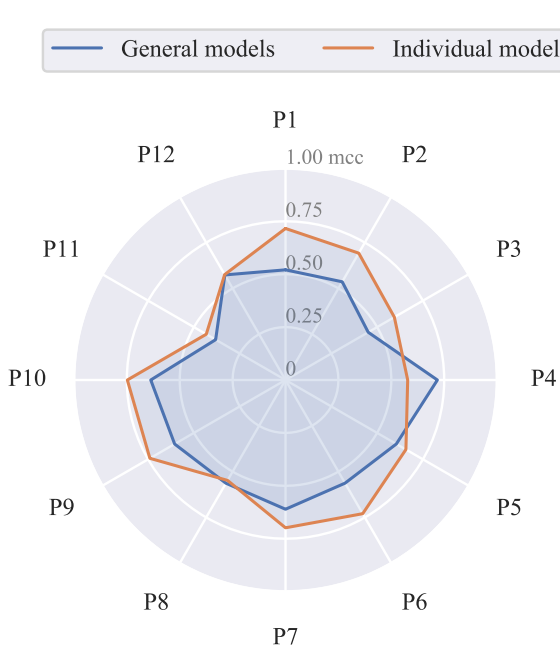


Fig. 13: Performance (in MCC) of each participant's individual model compared to the mean performance of five general models, averaged over all abstraction levels.

#### D. Comparison between individual and general models

A second way to test whether the trained models are individual-sensitive is to compare individual model performance to the performance of the general models, when applied to the test data of each respective participant. Figure 13 shows the average individual model performance per participant, compared to the average general model performance per participant. The chart shows that most individual models outperform the mean of the general models, but some cases show near equal or even worse (P4 and P8) performance, possibly caused by a strategy change in the final run.

A paired t-test shows that the individual models perform

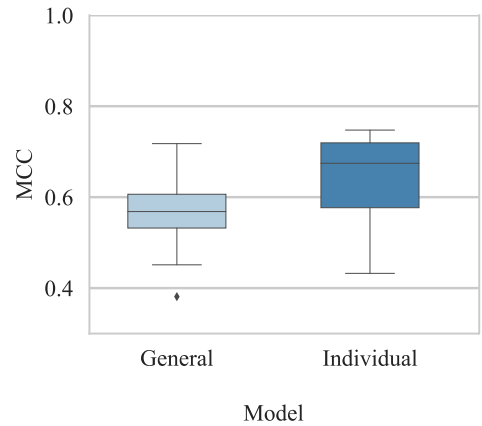


Fig. 14: Comparison of model performance between general and individual models.

significantly better ( $t(11) = 2.9, p = 0.02$ ) than the general models in terms of MCC, see Figure 14. The individual models provide a mean 0.08 (SD = 0.10) MCC improvement over the general models. The personalized approach is most effective for participant 1, whose individual models score 0.20 MCC higher than the baseline.

## V. DISCUSSION

The aim of this study was to create individual-sensitive models of controller strategy by training a set of convolutional neural networks on a visual representation of traffic conflicts. A human-in-the-loop experiment was performed to generate training data for the model creation.

It is a common problem in machine learning that such model training requires a large amount of data. To mitigate this problem, the performed experiment considered only a subset of the types of conflict that controllers can encounter in their sector. Throughout the experiment, similar conflicts were presented to each participant multiple times, by only introducing conflicting traffic from the east, with a limited number of crossing angles. In addition, altitude differences were not taken into account, nor were altitude changes accepted as

a conflict resolution. Increasing the breadth of encountered conflicts could increase the validity and applicability of the trained models, but would require more training data. Results show that model performance per validation fold can fluctuate up to 0.2 MCC, which shows sensitivity to random state-action pair sampling from the training set. This indicates that even for this limited subset of conflicts, performance improvements can be gained with exposure to more training data.

Evidence that convolutional neural networks are able to interpret the SSD is shown by the accuracy increase during training. This might, however, also indicate that the network overfits on the training data at pixel-level without generalization to new samples. For the current dataset, the ability to generalize is demonstrated by the fact that the validation performance during training follows an upward trend similar to the training performance. Test results using the separate test set further confirm this. On the other hand, overfitting on the training samples does occur to a certain extent, indicated by the performance difference between the training and test sets.

It is expected that performance improvements can be achieved by a formal grid search to optimize the hyperparameters of the model. However, since every training iteration generates 360 models (12 participants  $\times$  3 control variables  $\times$  2 repetitions  $\times$  5 cross-validation folds), iterating over parameters is computationally expensive and meticulous. The network architecture used in this research is kept constant for all control variables, i.e., it can predict command *type*, *direction* and *magnitude* by only altering the last fully connected layer. Designing separate network architectures that are tailored to each control variable could consequently improve performance.

Resolution *magnitude* predictions are obtained using classification (i.e., dividing possible magnitudes over a limited number of bins) to more easily compare results to *type* and *direction* predictions. As resolution magnitude values are, in reality, more finely grained, this caused the accuracy of *magnitude* predictions to be reliant on the classification, i.e., the granularity of the bins. Therefore, more precise predictions are expected to be achieved using regression instead of classification.

Data quantity is often a limiting factor in machine learning experiments with human data. A method to increase effectiveness of the available data is *reinforcement learning*. One option is to start with supervised learning – incorporating human strategy – and to improve the models with more experience using reinforcement learning [23]. Another approach could be to use *inverse* reinforcement learning, which could learn a personalized reward function that is subsequently used to train a model through interaction with a simulated environment.

Nonetheless, the achieved MCC scores using the current methods indicate a considerably better than random prediction, even for the general models (up to 81% accuracy, averaged over all three control variable models). This illustrates that the SSD indeed captures sufficient information concerning air traffic conflicts to base predictions on – given the simplified scenarios used in this research – confirming the hypothesis. However, apart from the conflict geometry information included in the SSD, research on control heuristics indicates that traffic flows, sector geometry and other controller goals are also play significant drivers of ATCo strategy [41], [42]. Including these sector characteristics in the visual network input could further increase prediction conformance. Moreover,

model predictions could be taken to a higher-level decision stage such as aircraft selection and resolution geometry (e.g. ‘behind’ or ‘in front’), as proposed in Westin et al.’s consistency framework [37]. Finally, it should be noted that humans base their decisions on a dynamic situation, while the SSD only provides a static representation of each conflict situation. Using multiple consecutive SSD frames could be a way to incorporate the dynamics of a conflict situation into the model.

The premise of individual-sensitive automation is that controllers are sufficiently consistent, and that there is sufficient difference in strategies between different controllers. In this study, the results indeed reveal a correlation between participant consistency and model performance. This is in line with the expectation that more consistent controllers are better suited to base strategic conformal automation on. It should be noted, however, that the validity of the consistency metric is limited to asymmetric crossing scenarios and constrained conflict angles, as used in this research. In reality, different conflict geometries will have different implications on the resolution strategy of a controller. A tailored consistency metric which evaluates all situations in a case-by-case fashion would be better suited to determine the true consistency of a controller.

Both the application of the individual models to inter-participant data, and the comparison of the individual models to the general models illustrate that indeed, different participants show different strategies, and that the individual models are able to capture this. Although the mean accuracy and MCC improvements are not always large in these comparisons, for half of the population, MCC scores increased by more than 0.10 (up to 0.20) when comparing the individual model to the general model. Nevertheless, the difference between the most accurate models (mean = 83% accuracy) and the least accurate models (mean = 61% accuracy) is considerable. For participants 4 and 8, this even resulted in the individual model performing worse than the general model. In these cases, it may be that the participants were still learning during the experiment measurement phase, and may have changed control strategies in between runs. Combined with the results on controller consistency, however, it can be seen that models become more accurate when controllers are more consistent, which is in accordance with previous empirical research where professional ATCos were used [8], [43].

## VI. CONCLUSIONS AND RECOMMENDATIONS

This research evaluated to what extent strategic conformal automation for air traffic control can be achieved using convolutional neural networks through a human-in-the-loop experiment. A 12-participant experiment was devised to generate training data consisting of solution space diagram (SSD) images and conflict resolutions. Achieved model performances show that the SSD contains sufficient information to make accurate predictions on resolution *type*, *direction* and *magnitude* given by controllers in 2D traffic conflict scenarios.

Results show a correlation between the controller consistency metric and achieved model performance, confirming the hypothesis that consistent controllers are more suited for strategic conformal automation. Regardless, the majority of controllers in the population is sufficiently consistent to base the conformal automation on. Personalized models obtain significantly higher prediction accuracies than general models,



indicating that controllers in this experiment exhibit differentiating strategies, i.e., are not homogeneous as a group. This is a critical assumption for strategic conformal automation. However, the performance improvement due to individual modeling substantially differs per controller, ranging from a deterioration to improvements of 0.20 MCC (Matthews Correlation Coefficient) and 12% accuracy. Nonetheless, convolutional neural networks appear to be a feasible method to achieve strategic conformal automation.

## REFERENCES

- [1] SESAR Consortium, "The Concept of Operations at a glance," *Single European Sky*, 2007.
- [2] Joint Planning and Development Office (JPDO) and Next Generation Air Transportation System (NextGen), "Concept of operations for the next generation air transportation system," Tech. Rep., 2011.
- [3] M. Bekier, B. R. Molesworth, and A. Williamson, "Tipping point: The narrow path between automation acceptance and rejection in air traffic management," *Safety science*, vol. 50, no. 2, pp. 259–265, 2012.
- [4] T. S. Bolic, "Automation adoption and adaptation in air traffic control," Ph.D. dissertation, University of California, 2006.
- [5] R. Ehrmanntraut, "Full automation of air traffic management in high complexity airspace," *Doctor of Engineering*, Technical University of Dresden. Eyferth, K., Niessen, C., and Spaeth, O.(2003), *A model of air traffic controllers' conflict detection and conflict resolution*. *Aerospace Science and Technology*, vol. 7, no. 6, pp. 409–416, 2010.
- [6] C. Westin, C. Borst, and B. Hilburn, "Strategic Conformance: Overcoming Acceptance Issues of Decision Aiding Automation?" *IEEE Transactions on Human-Machine Systems*, vol. 46, no. 1, pp. 41–52, 2016.
- [7] T. Prevot, J. R. Homola, L. H. Martin, J. S. Mercer, and C. D. Cabral, "Toward automated air traffic control: investigating a fundamental paradigm shift in human/systems interaction," *International Journal of Human-Computer Interaction*, vol. 28, no. 2, pp. 77–98, 2012.
- [8] B. Hilburn, C. Westin, and C. Borst, "Will Controllers Accept a Machine That Thinks Like They Think? The Role of Strategic Conformance in Decision Aiding Automation," *Air Traffic Control Quarterly*, vol. 22, no. 2, pp. 115–136, 2014.
- [9] R. S. Sutton, A. G. Barto, and R. J. Williams, "Reinforcement Learning is Direct Adaptive Optimal Control," *IEEE Control Systems*, vol. 12, no. 2, pp. 19–22, 1992.
- [10] R. M. Regtuit, C. Borst, E.-J. Van Kampen, and M. R. M. van Paassen, "Building Strategic Conformal Automation for Air Traffic Control Using Machine Learning," in *AIAA SciTech Forum*. Kissimmee, Florida: AIAA Information Systems, 2018.
- [11] G. Mercado-Velasco, M. Mulder, and M. Van Paassen, "Analysis of Air Traffic Controller Workload Reduction Based on the Solution Space for the Merging Task," *AIAA Guidance, Navigation, and Control Conference*, vol. AIAA 2010-, no. August, pp. 1–18, 2010.
- [12] P. Fiorini and Z. Shiller, "Motion planning in dynamic environments using velocity obstacles," *The International Journal of Robotics Research*, vol. 17, no. 7, pp. 760–772, 1998.
- [13] S. B. J. Van Dam, A. L. Abeloos, M. Mulder, and M. Van Paassen, "Functional presentation of travel opportunities in flexible use airspace: An EID of an airborne conflict support tool," *Conference Proceedings - IEEE International Conference on Systems, Man and Cybernetics*, vol. 1, no. January, pp. 802–808, 2004.
- [14] J. Ellerbroek, K. Brantegem, M. Van Paassen, N. de Gelder, and M. Mulder, "Experimental evaluation of a coplanar airborne separation display," *IEEE Transactions on Human-Machine Systems*, vol. 43, no. 3, pp. 290–301, 2013.
- [15] Y. I. Jenie, E.-J. van Kampen, C. C. de Visser, J. Ellerbroek, and J. M. Hoekstra, "Selective velocity obstacle method for deconflicting maneuvers applied to unmanned aerial vehicles," *Journal of Guidance, Control, and Dynamics*, vol. 38, no. 6, pp. 1140–1146, 2015.
- [16] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [17] B. Piot, M. Geist, and O. Pietquin, "Bridging the gap between imitation learning and inverse reinforcement learning," *IEEE transactions on neural networks and learning systems*, vol. 28, no. 8, pp. 1814–1826, 2017.
- [18] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. Cambridge, MA: MIT Press, 2018.
- [19] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, "Backpropagation Applied to Handwritten Zip Code Recognition," *Neural computation*, vol. 1, no. 4, pp. 541–551, 1989.
- [20] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fiedelnd, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, p. 529, 2015. [Online]. Available: <http://dx.doi.org/10.1038/nature14236>
- [21] D. a. Pomerleau, "Alvinn: An autonomous land vehicle in a neural network," *Advances in Neural Information Processing Systems 1*, pp. 305–313, 1989.
- [22] Z. Chen and X. Huang, "End-To-end learning for lane keeping of self-driving cars," *IEEE Intelligent Vehicles Symposium, Proceedings*, pp. 1856–1860, 2017.
- [23] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis, "Mastering the game of Go with deep neural networks and tree search," *Nature*, vol. 529, no. 7587, pp. 484–489, 2016. [Online]. Available: <http://dx.doi.org/10.1038/nature16961>
- [24] Y. LeCun, R. Pfeifer, Z. Schreter, F. Fogelman, and L. Steels, "Generalization and network design strategies," Elsevier, Zurich, Switzerland, Tech. Rep., 1989.
- [25] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov, "Dropout: A Simple Way to Prevent Neural Networks from Overfitting," *Journal of Machine Learning Research*, vol. 15, pp. 1929–1958, 2014.
- [26] K. Jarrett, K. Kavukcuoglu, and Y. LeCun, "What is the best multi-stage architecture for object recognition?" in *Computer Vision, 2009 IEEE 12th International Conference on*. IEEE, 2009, pp. 2146–2153.
- [27] C. M. Bishop, *Pattern Recognition and Machine Learning*. New York, USA: Springer, 2006.
- [28] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.
- [29] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *International Conference on Learning Representations (ICLR)*, 2015.
- [30] S. B. J. Van Dam, M. Mulder, and M. van Paassen, "Ecological Interface Design of a Tactical Airborne Separation Assistance Tool," *IEEE Transactions on Systems, Man, and Cybernetics Part A: Systems and Humans*, vol. 38, no. 6, pp. 1221–1233, 2008.
- [31] J. Ellerbroek, M. Visser, S. B. J. van Dam, M. Mulder, and M. M. van Paassen, "Design of an Airborne Three-Dimensional Separation Assistance Display," *IEEE Transactions on Systems, Man, and Cybernetics, part A: Systems and Humans*, vol. 41, no. 6, pp. 863–875, 2011.
- [32] J. Ellerbroek, K. C. R. Brantegem, M. M. van Paassen, and M. Mulder, "Design of a Co-Planar Airborne Separation Display," *IEEE Transactions on Human-Machine Systems*, vol. 43, no. 3, pp. 277–289, 2013.
- [33] P. Hermes, M. Mulder, M. Van Paassen, and H. Huisman, "Solution-Space-Based Analysis of the Difficulty of Aircraft Merging Tasks," *Journal of Aircraft*, vol. 46, no. 6, pp. 1995–2015, 2009.
- [34] S. M. A. Rahman, "Solution Space-based Approach to Assess Sector Complexity in Air Traffic Control," Ph.D. dissertation, 2014.
- [35] J. G. D'Engelbronner, C. Borst, J. Ellerbroek, M. van Paassen, and M. Mulder, "Solution Space-Based Analysis of Dynamic Air Traffic Controller Workload," *Journal of Aircraft*, vol. 52, no. 4, pp. 1146–1160, 2015. [Online]. Available: <http://arc.aiaa.org/doi/10.2514/1.C032847>
- [36] Y. I. Jenie, E.-J. van Kampen, C. C. de Visser, J. Ellerbroek, and J. M. Hoekstra, "Three-Dimensional Velocity Obstacle Method for Uncoordinated Avoidance Maneuvers of Unmanned Aerial Vehicles," *Journal of Guidance, Control, and Dynamics*, vol. 39, no. 10, 2016.
- [37] C. Westin, *Strategic Conformance: Exploring Acceptance of Individual-Sensitive Automation for Air Traffic Control*. PhD Thesis. Delft University of Technology, Netherlands, 2017.
- [38] B. W. Matthews, "Comparison of the predicted and observed secondary structure of T4 phage lysozyme," *Biochimica et Biophysica Acta (BBA)-Protein Structure*, vol. 405, no. 2, pp. 442–451, 1975.
- [39] D. M. Powers, "Evaluation: from precision, recall and F-measure to ROC, informedness, markedness and correlation," *Journal of Machine Learning Technologies*, vol. 2, no. 1, pp. 37–63, 2011.
- [40] R. Kohavi, "A study of cross-validation and bootstrap for accuracy estimation and model selection," in *International Joint Conference on Artificial Intelligence*, vol. 14, no. 2. Montreal, Canada, 1995, pp. 1137–1145.
- [41] B. Kirwan and M. Flynn, "Investigating air traffic controller conflict resolution strategies," EUROCONTROL, Brussels, Belgium, Tech. Rep., 2002.
- [42] S. Fothergill and A. Neal, "Conflict-resolution heuristics for en route air traffic management," *Proceedings of the Human Factors and Ergonomics Society*, pp. 71–75, 2013.
- [43] C. Westin, C. Borst, and B. Hilburn, "An empirical investigation into three underlying factors affecting automation acceptance," *Proceedings of the Fifth SESAR Innovation Days*, pp. 1–9, 2015.

**Sjoerd van Rooijen** received his MSc in Aerospace Engineering from Delft University of Technology in January 2019 (cum laude). The work presented in this paper was part of his thesis.

**Joost Ellerbroek** received the M.Sc. (2007) and Ph.D. (2013) degrees from the Delft University of Technology, The Netherlands. He is currently assistant professor with the CNS/ATM chair at the faculty of Aerospace Engineering of TU Delft. His current work includes, amongst others, urban airspace design for drones, topics of airspace complexity and capacity analysis, data mining applications for ATM, analysis of future ATM concepts such as the extended arrival manager, and new approach procedures, separation algorithms for UAS, and several other ASAS-related studies.

**Clark Borst** received the M.Sc. (2004) and Ph.D. degrees (2009) from the Delft University of Technology, The Netherlands, where he is currently working as an Assistant Professor in aerospace human-machine systems. His work and research interests include developing and evaluating human-centered aviation automation for flight deck and air traffic control applications using Ecological Interface Design principles.

**Erik-Jan van Kampen** obtained his BSc (2004) and MSc (2006) degrees in Aerospace Engineering, and his PhD degree (2010) at Delft University of Technology. The topic of the PhD-research was interval optimization and its application to aerospace problems. His research interests are Intelligent Flight Control, nonlinear adaptive control, and interval optimization.