# Multi-Objective Optimisation of Safety Related Systems

Richard M. Everson and Jonathan E. Fieldsend *Member, IEEE*

*Abstract*— Many safety related and critical systems warn of potentially dangerous events; for example, the Short Term Conflict Alert (STCA) system warns of airspace infractions between aircraft. Although installed with current technology such critical systems may become out of date due to changes in the circumstances in which they function, operational procedures and the regulatory environment. Current practice is to 'tune' by hand the many parameters governing the system in order to optimise the operating point in terms of the true positive and false positive rates, which are frequently associated with highly imbalanced costs.

In this paper we cast the tuning of critical systems as a multi-objective optimisation problem. We show how a region of the optimal receiver operating characteristic (ROC) curve may be obtained, permitting the system operators to select the operating point. We apply this methodology to the STCA system, using a multi-objective $(1 + 1)$-evolution strategy, showing that we can improve upon the current hand-tuned operating point, as well as providing the salient ROC curve describing the true-positive versus false-positive trade-off. We also provide results for three-objective optimisation of the alert response time in addition to the true and false positive rates. Additionally, we illustrate the use of bootstrapping for representing evaluation uncertainty on estimated Pareto fronts, where the evaluation of a system is based upon a finite set of representative data.

*Index Terms*— Evolutionary computation, multiple objectives, safety related systems.

## I. INTRODUCTION

**M**ANY safety related systems can be regarded as two-class classifiers: they classify a particular set of inputs or features into classes that might be labelled *dangerous* and *benign*. Classifications into the *dangerous* class raise an alarm and generally require some sort of human intervention. The specific example with which this paper is concerned is the Short Term Conflict Alert (STCA) system in operation in the United Kingdom and elsewhere. STCA monitors aircraft locations from ground radar and provides advisory alerts to air traffic controllers if a pair of aircraft are likely to become dangerously close. The STCA system is designed to raise a warning to air traffic controllers if there is a developing conflict between aircraft, giving them time to redirect the aircraft.

Taking its input from ground radar, the STCA system is independent of the aircraft, and cannot know the intentions of the pilots or air traffic controllers who may be aware of a potential conflict and already taking measures to avoid it.

For this reason, and because STCA must make conservative predictions, there are necessarily *nuisance alerts* as well as *genuine alerts*. There is clearly a trade-off between genuine and nuisance alerts and it is desirable to minimise the number of nuisance alerts in order to maintain the air traffic controllers' confidence in STCA.

The Receiver Operating Characteristic (ROC, see [1] for a recent review) is useful for displaying and assessing the performance of two-class classifiers. The ROC curve displays the false positive rate versus true positive rate for a particular classifier as the classification threshold or parameters of the classifier are varied. This visual representation of the operating possible operating points for the classifier permits the system designer to select the optimal parameters with a knowledge of how true and false positive rates will vary as the parameters are altered. Regarding STCA as a two-class classifier, which partitions pairs of radar tracks into *dangerous* or *serious* and *benign* classes, allows ROC analysis to be applied in which *genuine alerts* are true positives, while *nuisance alerts* are false positives.

The STCA system became operational for part of UK airspace in 1988 [2] and versions capable of coping with complex terminal control airspaces have been in operation since 1994. Since its introduction there have been incremental changes to the software and it is now used across the UK and elsewhere. Importantly, however, there have been changes in the volume and nature of air traffic together with changes to the management of the airspace monitored by STCA. Bringing new software into service involves a lengthy period of testing and scrutiny, even for *advisory* systems such as STCA; consequently, staff at the National Air Traffic Services (NATS, the principal civil air traffic control service for the United Kingdom) undertake parameter reviews in which they adjust (*tune*) the operating parameters of the STCA system in order to reduce the number of nuisance alerts, while maintaining the genuine alerts. This tuning is performed on the basis of a large (170 000) database of track pairs containing historical and recent encounters. The great number of parameters (at least 1500) determining the behaviour of STCA make tuning a highly skilled and laborious business. However, despite a recent step towards automation [2], the optimal receiver operating characteristics of the STCA system have not been known.

In this paper we introduce an approach to resolving these optimisation problems using multi-objective optimisation techniques based on evolutionary algorithms [3]–[5]. We cast the true positive and false positive rates obtained by STCA

as two opposing objectives to be maximised and minimised respectively. This allows us to obtain the optimal ROC curve from which the operating point can be chosen with a full knowledge of the trade-off between genuine versus nuisance alert rates.

In section II we describe the STCA system used in the UK; and in section III we describe the current optimisation process of STCA within the UK air traffic service, together with previous attempts at the automation of its optimisation. In section IV we discuss the relation of ROC analysis to the more general theory of Pareto optimality; based on this, in section V we describe the multi-objective optimisation technique approach to discovering the ROC curve for the system, and provide results in section VI. The paper concludes with a discussion in section VII. A preliminary report on this work appeared in [6].

## II. THE SHORT TERM CONFLICT ALERT SYSTEM

Here we focus on the Short Term Conflict Alert system (STCA) which is used widely within Europe by civil aviation authorities, in order to alert air traffic controllers to potential airspace infringements by aircraft pairs (i.e., two aircraft which may become too close). STCA is not strictly a safety critical system—*a system containing computer, electronic or electromechanical components whose failure may cause threat to life and limb or severe damage to property*[1]—but rather a component of the NATS 'safety net', providing *advisory* alerts to air traffic controllers of potential airspace proximity violations. Nonetheless, it exhibits many of the characteristics of a safety critical system: it must be highly reliable, transparent and verifiable. Its importance is highlighted by the fact that it is thought that one of the factors contributing to the midair collision over the border between Germany and Switzerland in July 2002 was that the STCA system in the relevant Swiss control station was switched off for maintenance [8].

### A. Overview

Figures 1 and 2 give an overview of the operation of the STCA system, which incorporates a highly complex and proprietary algorithm. Ground radars track the aircraft in a given airspace and those adjoining, and every four seconds (a STCA cycle) create *track pairs* of all possible combinations of aircraft. A *coarse* filter (Figure 1) is used first to remove all those pairs which are simply too far away from each other to be of concern. Potential conflict pairs are then processed in the core of STCA by three *fine* filters: a mixture of three models; a *linear prediction* filter; a *current proximity* filter; and a *manoeuvre hazard* filter (Figure 2). The boolean outputs of these fine filters are combined by the *alert confirmation module*, and aircraft pairs which are in danger of becoming too close are highlighted and alerted on the air traffic controllers' screens. The STCA is concerned with detecting airspace conflicts that may occur in the near future (around two minutes), so that air traffic controllers may be warned and the situation rectified in sufficient time.

[1]The working definition adopted by an ACM Special Interest Group on Computer-Human Interaction (SIGCHI) workshop [7] and typical of definitions of safety critical systems.
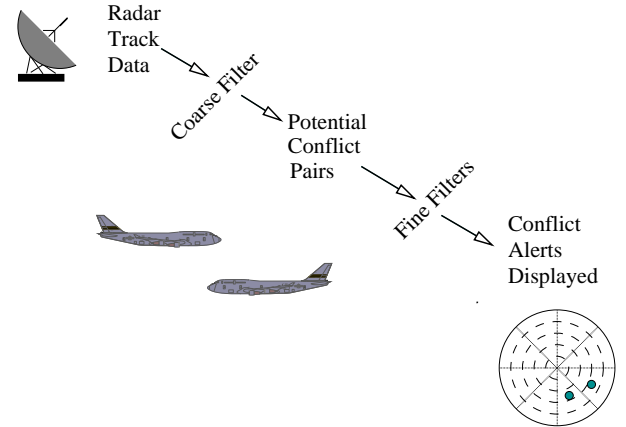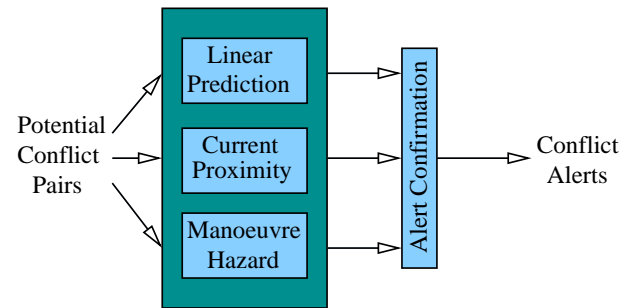


Fig. 1.   High level view of the STCA model.



Fig. 2.   STCA fine filters.

The minimum separation that is counted as an air proximity conflict depends on a number of criteria (for example, the airspace location and available radar cover). Generally in the UK in controlled airspace it ranges between 3, 5 or 10 nautical miles horizontally and 1000ft vertically. The linear prediction filter checks for loss of horizontal or vertical separation assuming that the aircraft continue on in a straight line at their current headings and speeds. The current proximity filter merely checks for a current loss of separation and the manoeuvre hazard filter classifies potential conflicts when either or both of the aircraft are turning. The combination of the binary classifications from the three fine filters by the alert confirmation module (Figure 2) is relatively sophisticated. During the confirmation process alerts from a track pair are checked within a moving time window, and if they are in conflict for a number of successive radar cycles (typically two or three), then an alert is passed onto the controller, although alerts from the current proximity filter are relayed more rapidly.

### B. Parameterisation

Each portion of the UK airspace is marked as distinct *region types*. For instance *en route* describes the airspace between airports, while regions where aircraft circle until permission is given to land are designated as *stack*. Since aircraft in different region types tend to have different types of flight behaviour, separate parameter sets are used for each one of the

region types. The particular parameter set used for classifying a track pair therefore depends upon the region types of the two aircraft; additional rules are used to determine the relevant parameter set if the aircraft have different region types.

The busy airspace above London, with which this study is concerned, is divided into 16 of these different region types. This multiplicity of parameter sets leads to a great number of parameters that can be adjusted to affect the performance of the STCA system. There are 96 parameters pertaining to the three fine filters, which means that the system uses approximately 1550 parameters (the coarse filter using fewer than 20). This includes both floating point and integer values over many varying ranges. Note that it is not feasible to adjust the parameters for the filters of each region type independently of the other region types, because track pairs involving pairs of regions lead to significant interactions between the parameters of different region types. On the other hand, as we describe below, only approximately two thirds of the available parameters are routinely adjusted.

The three central components of STCA are readily understood and their operation is capable of verification by practitioners, which is a common feature to the majority of critical systems in use. Regulatory authorities are very uneasy about using black-box techniques, such as artificial neural networks, in which function mappings are not easily described or understood. As we have described, the filter components of the STCA system themselves do, however, possess a large number of user determined parameters, which affect the operation of the system and therefore whether or not the system alerts pairs as being in potential conflict. The STCA program may be thought of as a decision tree, particular branches of which are followed depending upon the aircraft track pair being processed and the thresholds which are determined by the operational parameters of STCA. Note that the operational parameters affect the classification produced by altering the thresholds and model parameters; changes in parameter values do not affect the logical routes that may be taken through the decision tree. The logical structure of the program is incrementally altered by NATS Operational Analysis & Support group as new versions of the software are introduced. However, routine tuning of the system does not affect the logical structure.

## III. OPTIMISATION OF STCA

The STCA system is in operation in the four UK air traffic control centres and at other air traffic control centres in Europe, so appropriate parameter setting must be chosen for each particular locale. Moreover, changes in the volume of air traffic, changes in local air traffic operational procedures and changes in the regulatory environment mean that the STCA operational parameters must be reviewed and updated in order to prevent the system becoming out of date. In the UK all serious near-miss encounters are reviewed under the auspices of the Airprox Board (see for example [9]). In addition NATS regularly assesses the efficacy of the STCA system by running an off-line version with a database comprised of recent general traffic encounters together with historical serious encounters.

The two samples permit the nuisance alert rate for general traffic to be monitored together with the warning time provided for genuine alerts.

### A. Manual Optimisation

As shown in Table I, each encounter is categorised by NATS staff into one of five categories of diminishing severity; category 4 encounters are semi-automatically categorised, but all others are manually annotated. Note that without knowledge of a pilot's intentions or the instructions a pilot has received, it is very difficult to predict whether an ascending or descending aircraft will level off at a specified height or 'bust' through the level potentially leading to a conflict. Errors in predicting level off clearly lead to nuisance alerts and as such we ignore category 3 encounters (as recommended by NATS).

STCA performance on the database is assessed using the Conflict Alert Management Performance Analysis Package (CAMPAP), which runs the STCA system on the database and analyses the performance for each category in each region [10], [11]. Using CAMPAP, the Operational Analysis & Support group within NATS has over the last 10 years, through manual adjustment of the parameters, tuned STCA to achieve the best balance between genuine and nuisance alerts. In essence this has been achieved by skilled staff running different parameter settings through the CAMPAP simulation, by changing one or more of the values in current use, and assessing the performance on the collated data.

As iterative evolution of the STCA system has occurred, and the airspace in the UK is partitioned into ever more disparate region types, this task clearly becomes more arduous. As an indication of the increasing complexity it may be noted that since the work of Beasley et al. [2] in 2002 the increase in the number of fine filter parameters and regions has led to an increase of roughly 500 in the number of STCA parameters.

### B. Weighted Objective Optimisation

Beasley et al. [2] recognised that the current approach of tweaking the system variables by hand may be suboptimal, and so applied the tabu search heuristic in an attempt to automate the process. In this work a single objective was maximised. The objective was a weighted sum of the number of genuine alerts gained and lost in comparison with a base parameter set; the number of nuisance alerts gained and lost in comparison with the base parameter set; and a measure of the difference in warning times for alerts, again in comparison with the base parameter set.

The problem when optimising a weighted sum of objectives is knowing the appropriate weights a priori to operate at a point on a Pareto front whose location is not known in advance. Indeed, slightly different shaped fronts can lead weighted sum optimisers to return drastically different operating points [12].

The tabu search optimiser [2] was also found to be susceptible to trapping in local minima and required manual analysis of the parameter space to re-start the search. Perhaps in the light of these considerations, the original iterative person-based adjustment is still in use by NATS.

TABLE I
ENCOUNTER CATEGORIES USED BY NATS.

| C | Alert | Description |
|---|---|---|
| 1 | Necessary | Serious or potentially serious encounter with a significant collision risk for which alerts and additional warning time are considered highly desirable. |
| 2 | Desirable | Serious encounters, which involved an actual or potential loss of separation, but little risk of collision, where alerts and additional warning time are considered desirable. |
| 3 | Unnecessary | Level off with risk encounters where a standard level off prevented a conflict. The desirability of alert for these encounters is dependent on where (and to some extent when) they occur. In busier airspace, such as stacks, they may be seen as an unnecessary distraction. Whereas in some less busy areas of airspace they may be seen as a valuable safety net (some controllers may reaffirm level off instructions when STCA indicates that a level bust would lead to conflict). |
| 4 | Undesirable | No actual or potential conflict. An alert would be considered a nuisance. |
| 5 | Bad data | Bad data for which alerts are generally considered a nuisance but are commonly deemed beyond the remit of STCA and therefore not usually taken into account during a parameter review. |

## IV. ROC ANALYSIS & PARETO OPTIMALITY

If we wish to satisfy the two opposing objectives of true positive maximisation and false positive minimisation, when the classes are skewed and the costs imbalanced it does not make sense to try and optimise a single objective function as illustrated in the previous section. If the costs of an incorrect classification were known the expected cost for any parameter set could be calculated [13] and used as a single objective function [14]. However, this procedure requires accurate specification of the misclassification costs which are seldom accurately known; indeed it is often desirable to present the user with a ROC curve from which the best operating point can be selected. A common method is to employ the Neyman Pearson criterion: a maximum false-positive rate is specified, which then determines the true-positive rate.

Alternatively, some other summary measure of the ROC curve, such as the area under the ROC curve (AUROC) could be used as a measure of the quality of a set of parameters [15], [16]; this overall measure could then be used as an objective to be optimised with respect to the system parameters.

Of course, all these measures based upon the ROC curve require knowledge of the ROC curve, which hitherto has been unavailable for the STCA system. In this section we show how multi-objective evolutionary algorithms (MOEAs) may be used to derive the ROC curve for the STCA system. However, we take the view that summarising the ROC curve neglects the true value of the curve, namely providing the user with an analysis of the trade-offs inherent in choosing an operating point. In this manner we can entirely circumvent the problematic *a priori* setting of objective weights encountered in [2].

### A. The ROC curve and Pareto optimality

In general we consider a classifier $g(\mathbf{x}; \boldsymbol{\theta})$ which gives an estimate of the probability that a feature vector $\mathbf{x}$ belongs to one of two classes. We assume that the classifier depends upon a vector of adjustable parameters $\boldsymbol{\theta}$, and we denote by $T(\boldsymbol{\theta})$ the classifier's true positive classification rate (measured on a particular dataset of interest), while the false positive rate is denoted by $F(\boldsymbol{\theta})$.

A ROC curve is frequently obtained by varying the probability threshold separating the two classes. As the threshold is varied from zero to one a non-decreasing ROC curve in the $(F, T)$ plane is obtained for any particular fixed set of parameters, and different ROC curves are obtained for different parameters. In this work, we consider the classification threshold to be subsumed in the parameter vector and seek to discover the set of parameters (including threshold) that simultaneously minimise $F(\boldsymbol{\theta})$ and maximise $T(\boldsymbol{\theta})$. In fact, the STCA classifier is a hard classifier, yielding only a binary classification rather than an estimate, however imprecise, of the probability of class membership. Nonetheless, we may still seek the set parameter values that yield the optimal true-positive versus false-positive trade-offs. (See, for example, [1] for extensive discussions of ROC curves for hard and soft classifiers.)

A general multi-objective optimisation problem seeks to simultaneously extremise $D$ objectives:

$$y_i = f_i(\boldsymbol{\theta}), \qquad i = 1, \ldots, D \tag{1}$$

where each objective depends upon a vector $\boldsymbol{\theta}$ of $P$ parameters or decision variables. It is convenient to assume that all the objectives are to be minimised, so for the STCA system we minimise the pair of objectives $(-T(\boldsymbol{\theta}), F(\boldsymbol{\theta}))$. The parameters may also be subject to the $J$ constraints:

$$e_j(\boldsymbol{\theta}) \geq 0, \qquad j = 1, \ldots J \qquad (2)$$

so that the multi-objective optimisation problem may be expressed as:

$$\text{minimise} \quad \mathbf{y} = \mathbf{f}(\boldsymbol{\theta}) = (f_1(\boldsymbol{\theta}), \ldots, f_D(\boldsymbol{\theta})) \qquad (3)$$

$$\text{subject to} \quad \mathbf{e}(\boldsymbol{\theta}) = (e_1(\boldsymbol{\theta}), \ldots, e_J(\boldsymbol{\theta})) \geq 0 \qquad (4)$$

where $\boldsymbol{\theta} = (\theta_1, \ldots, \theta_P)$ and $\mathbf{y} = (y_1, \ldots, y_D)$.

When faced with only a single objective an optimal solution is one which minimises the objective given the model constraints. However, when there is more than one objective to be minimised solutions may exist for which performance on one objective cannot be improved without sacrificing performance on at least one other. Such solutions are said to be *Pareto optimal* [3]–[5] and the set of all Pareto optimal solutions is said to form the Pareto front.

The notion of *dominance* may be used to make Pareto optimality clearer. A decision vector $\boldsymbol{\theta}$ is said to *strictly dominate* another $\boldsymbol{\phi}$ (denoted $\boldsymbol{\theta} \prec \boldsymbol{\phi}$) iff

$$\begin{aligned} f_i(\boldsymbol{\theta}) \leq f_i(\boldsymbol{\phi}) \quad \forall i = 1, \ldots, D \quad \text{and} \\ f_i(\boldsymbol{\theta}) < f_i(\boldsymbol{\phi}) \quad \text{for some } i. \end{aligned} \qquad (5)$$

Less stringently, $\boldsymbol{\theta}$ *weakly dominates* $\boldsymbol{\phi}$ (denoted $\boldsymbol{\theta} \preceq \boldsymbol{\phi}$) iff

$$f_i(\boldsymbol{\theta}) \leq f_i(\boldsymbol{\phi}) \quad \forall i = 1, \ldots, D. \qquad (6)$$

A set of $M$ decision vectors $\{\boldsymbol{\theta}_i\}$ is said to be a *non-dominated set* if no member of the set is dominated by any other member:

$$\boldsymbol{\theta}_i \not\prec \boldsymbol{\theta}_j \quad \forall i, j = 1, \ldots, M. \qquad (7)$$

A solution to the minimisation problem (3) is thus *Pareto optimal* if it is not dominated by any other feasible solution, and the non-dominated set of all Pareto optimal solutions is the Pareto front. Recent years have seen the development of a number of evolutionary techniques based on dominance measures for locating the Pareto front; see [3], [5], [17] for recent reviews.

## V. OPTIMISATION USING MULTI-OBJECTIVE EVOLUTIONARY ALGORITHMS

Anastasio, Kupinski & Nishikawa introduced the use of multi-objective evolutionary algorithms (MOEAs) to optimise ROC curves, illustrating the method on a synthetic data [18] and for medical imaging problems [19]. Here we used a similar methodology, albeit with improved convergence properties.

The multi-objective evolutionary algorithm used in this study is a stochastic search algorithm, based on a simple $(1 + 1)$-evolution strategy (ES), similar to that introduced in [20]. In outline, the procedure for locating the Pareto front/ROC curve, operates by maintaining an archive, $A$, of mutually non-dominating solutions, $\boldsymbol{\theta}$, which is the current approximation to the Pareto front/ROC curve. At each stage of the algorithm

---

**Algorithm 1** A MO $(1 + 1)$-ES for STCA optimisation.

Inputs:
$N$     Number of ES generations

```
1:     A := initialise()
2:     n := 0
3:     while n < N :
4:         θ := select(A)
5:         θ' := perturb(θ)
6:         (T(θ'), F(θ')) := STCA(θ')
7:         if θ' ⋡ φ ∀φ ∈ A:
8:             A := {φ ∈ A | φ ⊀ θ'}
9:             A := A ∪ θ'
10:        end
11:        n := n + 1
12:    end
```

---

some solutions in $A$ are copied and perturbed. Those perturbed solutions that are dominated by members of $A$ are discarded, while the others are added to $A$ and any dominated solutions in $A$ are removed. In this way the estimated Pareto front $A$ can only advance towards the true Pareto front. This algorithm, unlike earlier versions [20], maintains an archive which is unconstrained in size, permitting better convergence properties [21].

Algorithm 1 describes in more detail the algorithm as applied to the optimisation of the STCA system. Following the current operating practice of NATS and [2], we choose to optimise only 912 of the $> 1500$ parameters affecting the STCA system; these parameters are those parameters which have different values in different regions after tuning by NATS. Furthermore we restrict these parameters to the ranges over which they are adjusted by NATS.

The archive or frontal set $A$ is initialised by drawing parameters for the STCA system uniformly from their feasible ranges; in addition the current 'best' parameter set from manual tuning $\boldsymbol{\theta}^\star$ is added to $A$. Of course many of these randomly selected parameter vectors are dominated by other parameter vectors and these dominated parameters are deleted from $A$ so that $A$ is a non-dominated set (7). In fact, in the work reported here, we found that of 100 randomly initialised parameters only $\boldsymbol{\theta}^\star$ and one other parameter vector remained in $A$ after dominated parameters were removed.

Following initialisation, the loop on lines 4–11 of Algorithm 1 is repeated for the desired number of iterations. At each iteration a single parameter vector $\boldsymbol{\theta}$ is selected from $A$; selection may be uniformly random, but partitioned quasi-random selection [21] was used here to promote exploration of the front. The selected parent vector is perturbed to generate a single *child* (line 5). Each individual parameter in the parent vector is perturbed with equal probability (0.2 here); the perturbations themselves are made by adding a random number to the parent parameter value. Yao *et al.* [22] have shown that perturbations drawn from heavy-tailed distributions
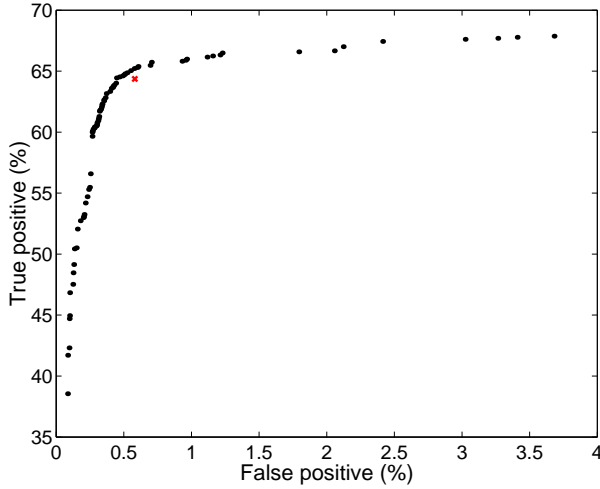
Fig. 3. Dots show estimates of the Pareto optimal ROC curve for STCA obtained after 6000 evaluations of the $(1+1)$-ES multi-objective optimiser. The cross indicates the manually tuned operating point $\boldsymbol{\theta}^{\star}$.
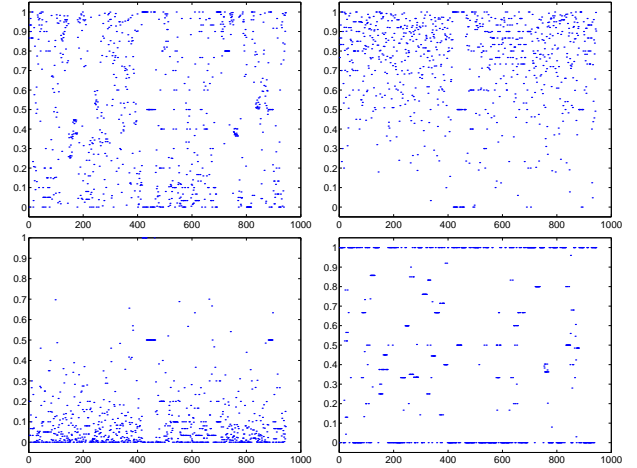


Fig. 4. Normalised parameter values for operating points on the Pareto optimal ROC curve shown in Figure 3. Panels correspond to parameters for: low $F$, low $T$ (bottom-left); medium $F$, medium $T$ (top-left); high $F$, high $T$ (top-right); and the manually tuned operating point $\boldsymbol{\theta}^{\star}$ (bottom-right).

facilitate convergence by promoting exploration and we draw perturbations from a Laplacian density, $p(x) \propto e^{-|x/w|}$, whose width is set equal to one tenth the feasible range of the parameter being perturbed; perturbations that lie outside the feasible range are resampled.

The true $T(\boldsymbol{\theta}')$ and false $F(\boldsymbol{\theta}')$ positive rates for the perturbed vector are evaluated by running the STCA/CAMPAP system with parameters $\boldsymbol{\theta}'$ on the test database of track pairs (Table I). Following NATS practise, we consider category 1 and 2 alerts to be true positives, while category 4 alerts are treated as false positives. The relatively small number of category 3 and 5 alerts are ignored. If the child $\boldsymbol{\theta}'$ is not dominated by any of the parameter vectors in $A$, any parameter vectors in $A$ that $\boldsymbol{\theta}'$ dominates are deleted from the archive (line 8) and $\boldsymbol{\theta}'$ is added to $A$ (line 9). These two steps ensure that $A$ is always a non-dominated set whose members dominate any other solution encountered thus far in the search.

In a $(\mu + \lambda)$−ES, $\mu$ parameter vectors are perturbed to generate $\lambda$ new vectors. That is, $\mu$ parameter vectors are selected (whose performances have already been evaluated); these *parents* are copied and have their parameter values perturbed in order to generate $\lambda$ *children*. Optimisation schemes with $\lambda > 1$ are attractive because the evaluation of the children may be performed in parallel. The computational cost of evaluating a single set of STCA parameters within CAMPAP is fairly high, at approximately 5 minutes. However, the system is written in a proprietary variant of PASCAL, which necessitates it be run on a Compaq Alpha machine. Since only a single Alpha was available to us, we used a $(1+1)$-ES, which has been shown to perform well compared to $(\mu + \lambda)$ MOEA implementations [23].

## VI. RESULTS

In this paper we present a conservative application of the MOEA method to STCA optimisation. It is conservative in that the ranges of parameters to be varied are limited by the current ranges of that parameter across the 16 region types within the currently applied STCA parameterisation of NATS. This means effectively we are only concerned with adjusting $2/3$ of the model parameters (still a significant number!), and the parameters are confined to regions of decision space with which personnel at NATS have considerable experience.

### A. True and false positive optimisation

Initially we optimised the true and false positive rates for a database comprised of manually and semi-automatically categorised encounters. The database included historical track pairs leading to serious or potentially serious encounters together with general traffic track pairs from two weeks in 2001.

Even this conservative optimisation approach produces some striking results. Figure 3 shows the estimates of the Pareto optimal ROC curve obtained using the multi-objective optimiser after $N = 6000$ evaluations (approximately 12 days computation). The current NATS operating point is also plotted as a cross. The optimisation has located an ROC curve consisting of 76 points ranging from $38.5\%$ to $67.9\%$ true positive and $0.1\%$ to $3.7\%$ false positive. In addition the manually tuned STCA operating point $\boldsymbol{\theta}^{\star}$ lies *behind* (is dominated by) several operating points on the estimated ROC curve. Although the improvement over $\boldsymbol{\theta}^{\star}$ is relatively small in percentage terms, the quantity of track pairs processed by the STCA system means that a significant reduction in the *number* of false alerts could be achieved while maintaining the current genuine alert rate. We regard as more important, however, the production of the ROC curve itself, because it reveals the true-positive versus false-positive trade-off, permitting the operating point to be chosen. In fact it may be observed that the current operating point $\boldsymbol{\theta}^{\star}$ is close to the corner of the Pareto optimal curve. Choosing an operating point to the left of the corner would result in a rapidly diminishing genuine alert rate for little gain in the nuisance alert rate; whereas operating points to the right of the corner provide small increases in the
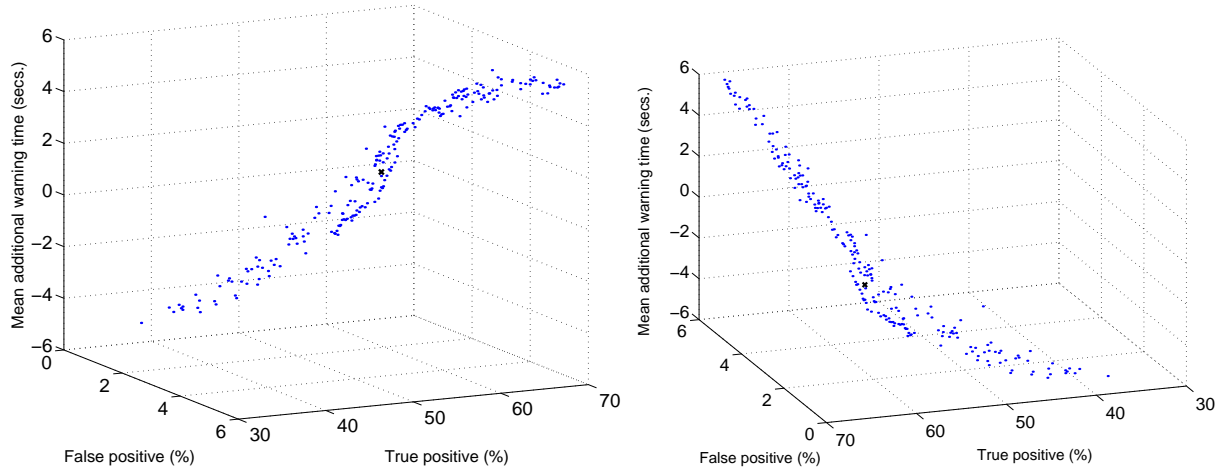
Fig. 5. Three objective estimated Pareto front for STCA. The cross indicates the current operating point $\theta^\star$.

true positive rate at the expense of relatively large increases in the false positive rate.

Figure 4 gives an indication of how the parameters which could be altered during the optimisation vary as the Pareto front is traversed. Each of the four panels in figure 4 shows the 912 variable parameters, each normalised to the interval $[0, 1]$, so that 0 represents the minimum value it was permitted to assume during optimisation and 1 represents the maximum. The bottom-right panel shows the parameters at the manually tuned operating point $\theta^\star$; many of the parameters are at their extreme values because we choose the allowable ranges to be defined by the extremal values located by NATS manual optimisation. There is a resemblance between these parameters and the parameters corresponding to the middle of the Pareto front ($F = 0.52\%$, $T = 64.87\%$) shown in the top-left panel. The bottom-left and top-right panels show $\theta$ corresponding to the extreme ends of the front. These appear to have a qualitatively different character. We observe that there is a discernible bias toward the minimum allowable values in the parameters at the bottom-left end of the front ($F = 0.08\%$, $T = 38.55\%$) and trend toward the maximum allowable parameter values in the parameters at the top-right end ($F = 3.68\%$, $T = 67.86\%$). This may indicate that further optimal solutions can be found by permitting the optimisation to range over parameter values beyond those currently employed by NATS.

### B. Warning time optimisation

In addition to the trade-off between correct alerts and incorrect alerts, it is desirable to increase the warning time of genuine alerts given to air traffic controllers. Current practise is to compare a new parameter set with the current operating point by calculating the mean increase or decrease in warning times over the coincident genuine warnings of the two parameter sets. Using the same method we can compare all our frontal operating points with the current operating point. Furthermore we can use this extra objective to create a three-objective optimisation problem in which we seek to maximise
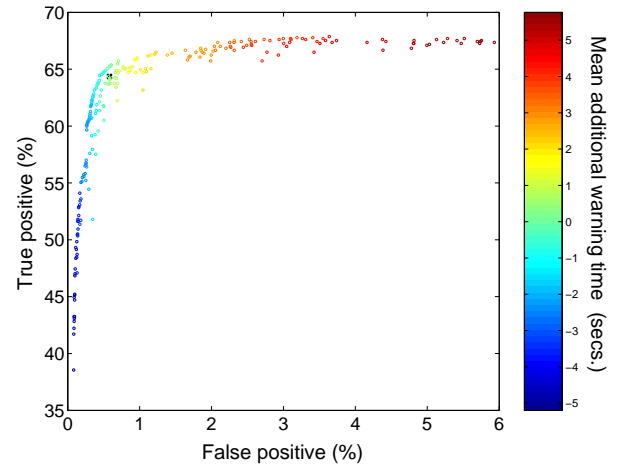


Fig. 6. Three objective estimated Pareto front for STCA, the third objective (mean additional warning time) represented in colour. The cross indicates the manually tuned operating point $\theta^\star$.

the mean warning time and true positive rate, while minimising the false positive rate.

Again we use a $(1 + 1)$-ES, with the same parameters as the previous experiment. We initialise the algorithm with the frontal points discovered in the previous optimisation (which by definition also form an estimated Pareto front in the 3 objective case). The front located after 5000 generations looks like a twisted ribbon, as shown in Figure 5. As before the current operating point $\theta^\star$ lies behind the discovered front.

It is interesting to observe that as the number of correct warnings increases the mean additional warning time is also seen to increase. This is shown clearly in Figure 6 where front is plotted as an ROC curve in two dimensions with the warning time in colour. We also point out that Figure 6 shows that the increases in genuine and nuisance alert rates close to the corner of the Pareto ROC curve are obtained without any significant change in the warning time.

The three-objective front contains almost four times as many points as the initial two-dimensional front. However, as Figure
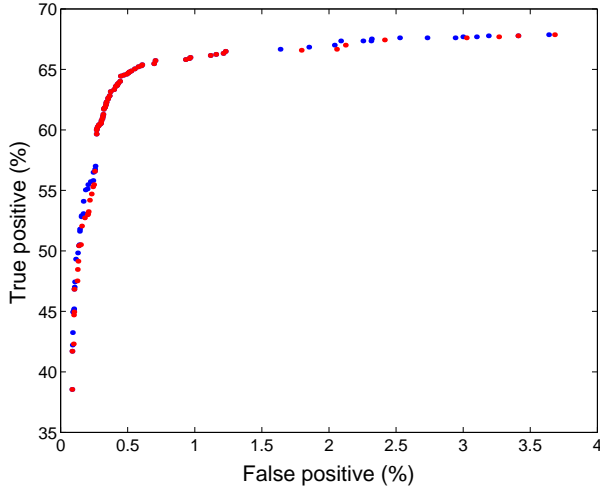
Fig. 7. Solutions from the 3-D optimisation that are are not dominated in the $F$-$T$ plane by the 2-D front. Red dots indicate the 2-D front (Figure 3) and blue dots indicate solutions from the 3-D front.
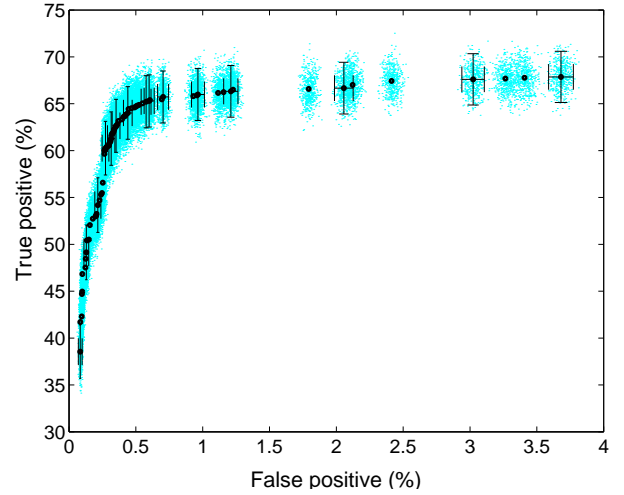


Fig. 8. Uncertainty of points on the estimated Pareto optimal ROC curve evaluated using bootstrapping. Error bars indicate two-standard deviation intervals calculated using equation (9) for a few representative points.

7 shows, if the three dimensional points are projected into the $F$–$T$ plane a few of them dominate or are mutually non-dominating with solutions from the initial two-dimensional optimisation. As the front has shifted only marginally forward at the edges, we may infer that we are providing a fairly good approximation of the true Pareto optimal ROC curve for the problem subject to the constraints on the parameters.

### C. Robustness of the front

As we described above, the location of the Pareto front is based upon evaluating the STCA system on a representative sample of encounters, and although approximately 170 000 encounters were used, it is important to discover the sensitivity of any putative operating point to the data sample. Indeed, it is especially important not to over-train the system to one particular set of data. Ideally one would optimise the entire STCA system on several independent data sets collected at different times. This, however, is impractical both because of the expense in collecting and annotating the data and of because of the computational expense of multiple optimisations (although this cost might be reduced by initialising new optimisations from fronts obtained in earlier optimisation runs). A further consideration is that serious encounters are (fortunately) rare, so that although independent sets of general traffic may be obtained, the serious encounters would have to be reused. For these reasons we employ a bootstrapping technique [24] in order to estimate the variability in error rates around the front.

The bootstrap method evaluates the error rate on a number of surrogate data sets constructed by sampling the original data set. Suppose that the original data set comprises $N = N_D + N_B$ examples, where $N_D$ is the number of examples in the *dangerous* class and $N_B$ is the number of *benign* examples. A bootstrap sample is constructed by drawing at random with replacement $N$ examples from the original sample. Note that some examples in the original data will be included more than once in a particular bootstrap surrogate, while others will be excluded entirely. The classification rate averaged over

a number of bootstrap replications is just the classification rate evaluated on the original data set, but an estimate of the variability in the classification rate may be obtained from the variation in the classification rates over the bootstrap replications.

Figure 8 shows the true and false positive rates obtained by evaluating the STCA system on 500 bootstrap replications for parameters on the front. While there is considerable spread about each location on the front, these scatter diagrams provide an estimate of the robustness of the parameter set to the data and indicate the range of true and false positive rates that may be expected.

In fact, the variability in rate may be obtained without recourse to numerical sampling. Focusing on the true positive rate $T(\boldsymbol{\theta})$ for a particular parameter set, a bootstrap sample could be constructed as follows: Choose at random, with replacement, $N_D$ examples from the dangerous class and likewise $N_B$ examples from the benign class. Since the true positive rate is $T(\boldsymbol{\theta})$, the probability of obtaining exactly $k$ true positives in the bootstrap sample is given by the binomial distribution:

$$p(k) = \binom{N_D}{k} T^k (1 - T)^{N_D - k} \qquad (8)$$

It is well known (e.g., [25]) that the mean of the binomial density is $N_D T$, so the mean true positive rate over many bootstrap replications is $T$, as expected. Furthermore, the variance of the number of true positive examples in a particular bootstrap replication is $T(1 - T)N_D$, so the variance in the true positive rate is

$$\sigma_T^2 = \frac{T(1 - T)}{N_D} \qquad (9)$$

with a similar expression for the variance of the false positive rate.

In fact, the bootstrap samples were constructed by merely sampling with replacement from the original data set *without*

regard for the number in each class, so although the mean number of dangerous exemplars in each class was $N_D$ it fluctuated from bootstrap sample to bootstrap sample. Nonetheless, it may be shown that when $N_D$ and $N_B$ are even moderately large (greater than about 20) (9) is a very good approximation to the variance in the rate and well describes the scatter around the front shown in Figure 8.

## VII. DISCUSSION

We have presented a straightforward multi-objective optimisation scheme for locating the optimal ROC curve for the Short Term Conflict Alert system employed to give warning of potential breaches in air proximity by aircraft. The results show that parameters yielding a range of genuine and nuisance alert rates are located by the MOEA, thus revealing the genuine versus nuisance alert trade-off and permitting the operating point to be set with explicit knowledge of the trade-off. The idea of dominance is essential to the simultaneous optimisation of both true and false positive alert rates and it is interesting to note that the manually tuned operating point is dominated by several of the solutions found by multi-objective optimisation. In addition we have simultaneously optimised the warning time given for genuine alerts, although we find that significant gains in warning time can only be achieved if the nuisance alert rate is substantially increased.

It should be emphasised that the true and false positive alert rates were evaluated on a database of over 170 000 track pairs, consisting of historical alerts deemed to be serious and two weeks worth of relatively current data, this comprises the same database that is currently used for manual tuning of operational STCA systems for the London sector airspace. It is important current work for skilled staff to inspect the parameter values obtained. However, the bootstrapping of the dataset around the optimised front provides an indication of the robustness of the optimised operating point. While these bootstrap estimates quantify the uncertainty in the optimised front, we remark that it would be beneficial to update a 'probabilistic front' so that new entrants were guaranteed with, say 90%, certainty not to be dominated by other elements of the front. Although it lies outside the scope of this report, we are developing multi-objective optimisers for this purpose and we draw attention to the work of Teich [26] and Hughes [27] who have both discussed optimisation of uncertain objectives.

The optimisations reported here were conservative in that they optimised only the 900 or so parameters that are routinely adapted by NATS, and these parameters were restricted to the ranges used by NATS. Although, as Figure 4 shows, solutions on the front are obtained for parameter values lying between the extremes used by NATS, we look forward to optimising a larger number of parameters and to permitting the parameters to vary over broader ranges.

The Pareto front located by the MOEA is comprised of a discrete set of parameter vectors at which the STCA system could be operated. However, we point out that the work of Scott *et al.* [28] shows that by randomly combining classifiers any operating point on the convex hull of the ROC curve may be obtained. Indeed it is apparent that if the objectives

to be optimised are statistical expectations, then Scott *et al's* work may be readily extended to three or more objectives to obtain an operating point on the convex hull of optimised solutions in many dimensions. It should be noted, however, that although the probabilistic combination of classifiers may lead to provably better average operating points, there are potential legal and ethical ramifications.

The production of the two-dimensional front took approximately twelve days of computer time. However, we emphasise that this was *unattended* computer time, in contrast to the labour-intensive and skilled process by which STCA systems are currently optimised. We anticipate that once an optimised ROC curve has been located for a particular STCA system and database, the subsequent optimisation following incremental incorporation of new cases into the database will be much faster. More rapid optimisation schemes are readily implemented via $(\mu + \lambda)$-ES, which are amenable to coarse parallelisation.

In this paper we have focused on the STCA system as an example safety related system; however, the STCA/CAMPAP system is treated purely as a subroutine of our evolutionary algorithm. Indeed in our implementation, the STCA/CAMPAP programs run on a separate computer. This 'wrapping' of the system to be optimised is important for two reasons. First, it shows that the technique is applicable to any critical system whose operating point is dependent on parameters that must be tuned and whose performance can be automatically evaluated. Second, and more importantly for safety-related systems, the wrapped system has not been modified in any way, thus preserving its integrity and the integrity of any safety case constructed for it.

Finally we remark that the majority of the parameters in the STCA filters have direct physical or mechanical interpretation, and that the transparency of the classification process is an important component in assuring the safety case for STCA. However, whether tuned by hand or optimised by a machine algorithm, the operational parameters are inferred from data and we look forward to the construction of safety cases for purely statistical classifiers whose operational parameters are inferred from data and have no ready physical interpretation.

## REFERENCES

[1] T. Fawcett, "ROC Graphs: Notes and Practical Considerations for Researchers," *Machine Learning*, 2004, (Submitted).
[2] J. Beasley, H. Howells, and J. Sonander, "Improving short-term conflict alert via tabu search," *Journal of the Operational Research Society*, vol. 53, pp. 593–602, 2002.
[3] C. Coello, "A Comprehensive Survey of Evolutionary-Based Multiobjective Optimization Techniques," *Knowledge and Information Systems. An International Journal*, vol. 1, no. 3, pp. 269–308, 1999.

[4] C. Fonseca and P. Fleming, "An Overview of Evolutionary Algorithms in Multiobjective Optimization," *Evolutionary Computation*, vol. 3, no. 1, pp. 1–16, 1995.

[5] D. V. Veldhuizen and G. Lamont, "Multiobjective Evolutionary Algorithms: Analyzing the State-of-the-Art," *Evolutionary Computation*, vol. 8, no. 2, pp. 125–147, 2000.

[6] J. Fieldsend and R. Everson, "ROC Optimisation of Safety Related Systems," in *Proceedings of ROCAI 2004, part of the 16th European Conference on Artifi cial Intelligence (ECAI)*, J. Hernández-Orallo, C. Ferri, N. Lachiche, and P. Flach, Eds., Valencia, Spain, 2004, pp. 37–44.

[7] P. Palanque, F. Paternó, and R. Fields, "Designing user interfaces for safety critical systems," in *Bulletin of ACM Speical Interest Group on Computer Human Interaction*. Association of Computing Machinery, 1998, vol. 30, no. 4, available from http://www.acm.org/sigchi/bulletin/1998.4/palanque.html.

[8] "Investigation Report, AX001-1-2/02," Bundesstelle für Flugunfalluntersuchung, Hermann-Blenk-Strasse 16, 38108 Braunschweig, Germany, May 2004, available from http://www.bfu-web.de.

[9] "Analysis of Airprox in UK Airspace (July 2002 to December 2002)," United Kingdom Airprox Board, Hillingdon House, Uxbridge, Middlesex, UB10 0RU, UK, 2003, available from http://www.caa.co.uk/ukab.

[10] G. Beeston, H. Howells, and M. Richards, *PAP User Guide*, 1st ed., National Air Traffi c Services Ltd., Research and Development Group, Department of Technical Research and Development, August 2000.

[11] ——, *Software Description for CAMPAP (Issue 1.03) Module*, 1st ed., National Air Traffi c Services Ltd., Research and Development Group, Department of Technical Research and Development, August 2000.

[12] I. Das and J. Dennis, "A closer look at drawbacks of minimizing weighted sums of objectives for Pareto set generation in multicriteria optimization problems," *Structural Optimization*, vol. 14, no. 1, pp. 63–69, 1997.

[13] R. Duda and P. Hart, *Pattern Classifi cation and Scene Analysis*. New York: Wiley, 1973.

[14] D. Hand, *Construction and assessment of classifi cation rules*. Wiley, 1997.

[15] E. DeLong, D. DeLong, and D. Clarke-Pearson, "Comparing the areas under two or more correlated receiver operating charteristic curves: a non parameteric approach," *Biometrics*, vol. 44, no. 11, pp. 837–845, 1988.

[16] J. Hanley and B. McNeil, "The mean and use of the area under a receiver operating characteristic (ROC) curve," *Radiology*, vol. 82, no. 143, pp. 29–36, 1982.

[17] K. Deb, *Multi-Objective Optimization Using Evolutionary Algorithms*. Chichester: Wiley, 2001.

[18] M. Kupinski and M. Anastasio, "Multiobjective Genetic Optimization of Diagnostic Classifi ers with Implications for Generating Receiver Operating Characterisitic Curves," *IEEE Transactions on Medical Imaging*, vol. 18, no. 8, pp. 675–685, 1999.

[19] M. Anastasio, M. Kupinski, and R. Nishikawa:, "Optimization and FROC analysis of rule-based detection schemes using a multiobjective approach," *EEE Transactions on Medical Imaging*, vol. 17, pp. 1089–1093, 1998.

[20] J. Knowles and D. Corne, "The Pareto Archived Evolution Strategy: A new baseline algorithm for Pareto multiobjective optimisation," in *Proceedings of the 1999 Congress on Evolutionary Computation*. Piscataway, NJ: IEEE Service Center, 1999, pp. 98–105. [Online]. Available: citeseer.nj.nec.com/knowles99pareto.html

[21] J. Fieldsend, R. Everson, and S. Singh, "Using Unconstrained Elite Archives for Multi-Objective Optimisation," *IEEE Transactions on Evolutionary Computation*, vol. 7, no. 3, pp. 305–323, 2003.

[22] X. Yao, Y. Liu, and G. Lin, "Evolutionary Programming Made Faster," *IEEE Transactions on Evolutionary Computation*, vol. 3, no. 2, pp. 82–102, 1999.

[23] J. Knowles and D. Corne, "Approximating the Nondominated Front Using the Pareto Archived Evolution Strategy," *Evolutionary Computation*, vol. 8, no. 2, pp. 149–172, 2000. [Online]. Available: citeseer.nj.nec.com/knowles00approximating.html

[24] B. Efron and R. Tibshirani, *An introduction to the Bootstrap*, ser. Monographs on Statistics and Probability. New York: Chapman & Hall, 1993, no. 57.

[25] A. Papoulis, *Probability, Random Variables and Stochastic Processes*. McGraw-Hill, 1991.

[26] J. Teich, "Pareto-front exploration with uncertain objectives," in *Evolutionary Multi-objective Optimisation, EMO 2001*, ser. LNCS, E. Z. et al, Ed., vol. 1993, 2001, pp. 314–328.

[27] E. Hughes, "Evolutionary multi-objective ranking with uncertainty and noise," in *Evolutionary Multi-objective Optimisation, EMO 2001*, ser. LNCS, E. Z. et al, Ed., vol. 1993, 2001, pp. 329–342.

[28] M. Scott, M. Niranjan, and R. Prager, "Parcel: feature subset selection in variable cost domains," Cambridge University Engineering Department, Tech. Rep. CUED/F-INFENG/TR.323, 1998.