

Learning Probabilistic Trajectory Models of Aircraft in Terminal Airspace from Position Data

Shane T. Barratt*, Mykel J. Kochenderfer†, and Stephen P. Boyd*

Abstract—Models for predicting aircraft motion are an important component of modern aeronautical systems. These models help aircraft plan collision avoidance maneuvers and help conduct offline performance and safety analyses. In this paper we develop a method for learning a probabilistic generative model of aircraft motion in terminal airspace, the controlled airspace surrounding a given airport. The method fits the model based on a historical dataset of radar-based position measurements of aircraft landings and takeoffs at that airport. We find that the model can generate realistic trajectories, can provide accurate predictions, and captures the statistical properties of aircraft trajectories. Furthermore, the model trains quickly, is compact to store, and allows for efficient real-time inference.

Index Terms—Air traffic control, predictive models, machine learning, unsupervised learning, Gaussian mixture model, clustering methods.

I. INTRODUCTION

THE United States's National Airspace System is undergoing a multibillion dollar modernization to increase safety, efficiency, and predictability of air transport, also known as NextGen [1]. A major focus of the modernization is on the portion of flight that occurs in the terminal airspace surrounding an airport, which represents the greatest safety risk. In fact, 56% of worldwide fatal air-related accidents occur in the takeoff, initial climb, final approach, and landing phases of the flight, even though these four phases constitute only 6% of flight time [2].

Key to the development and assessment of air traffic control technologies are probabilistic trajectory models of aircraft in the terminal region. These trajectory models can be used to predict the future trajectories of aircraft, detect conflicts, and recommend avoidance maneuvers. In addition, they can be used to assess the safety and operational performance of new technologies and procedures in simulation before deployment. It is critical, however, that the trajectory models be accurate representations of how aircraft actually behave in the airspace. This paper focuses on the construction of probabilistic trajectory models of aircraft in terminal airspaces from radar-based surveillance data.

Related work. We begin by summarizing proposed methods for modeling aircraft trajectories in unstructured airspace,

Manuscript received TODO. This material is based upon work supported by the National Science Foundation Graduate Research Fellowship under Grant No. DGE-1656518.

All authors are affiliated with Stanford University, Stanford, CA 94305:

* Department of Electrical Engineering; † Department of Aeronautics and Astronautics. {sbarratt, mykel, boyd}@stanford.edu

Digital Object Identifier.

²Top is generated, bottom is held-out test samples.

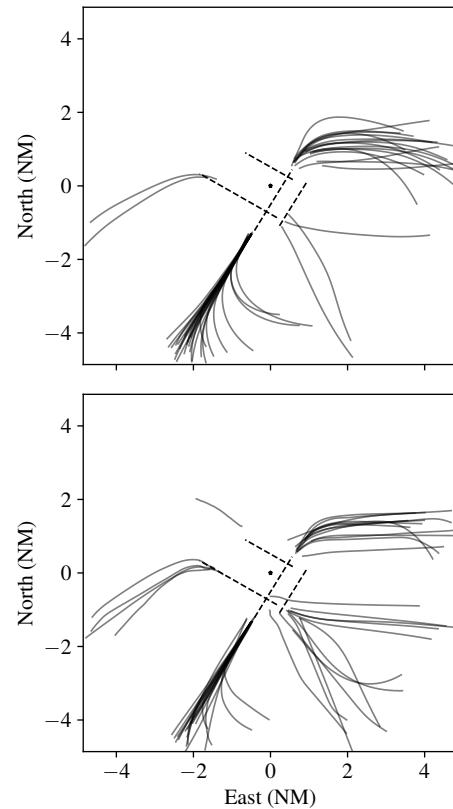


Fig. 1. Lateral view of takeoffs from KJFK. One figure is of held-out test trajectories and one is of trajectories generated by the model. Try to guess which one is real and which one is generated before looking at their true identities.²

i.e., outside of the terminal region. One line of work involves predicting future aircraft motion by estimating the aircraft's state and propagating that estimate through physical equations of motion [5], [6]. Other methods include synthesizing trajectories by combining individual segments defined by different modes of operation [7] and estimating vertical paths based on the weight of the aircraft [8]. Another method proposed by Kochenderfer *et al.* [9], [10] is to learn Bayesian network statistical representations of dynamic variables from historical radar data. Given trajectory predictions and their covariances for two aircraft, Paielli and Erzberger [11] derive a method to solve for the *conflict probability*, defined as the probability that the aircraft become too close. Other methods have been proposed to efficiently calculate conflict probabilities, including probability flow [12], [13] and importance sam-

pling [14]. Many of these trajectory prediction models are used in production aviation systems, including the Center TRACON Automation System (CTAS) [15], [16], and for evaluating the Traffic Collision Avoidance System (TCAS) [17] and ACAS X, the next generation of TCAS [18]. Realistic airport simulations are also used to validate NextGen [19].

There are also hybrid dynamics models of aircraft in unstructured airspace, which explicitly take into account both continuous and discrete dynamic behavior. Prandini *et al.* [20] propagate aircraft state through a stochastic hybrid dynamical system model to perform conflict detection. Combining the stochastic hybrid dynamical system model with a method for inferring the aircraft's navigational intent leads to a commonly used set of modeling tools for aircraft trajectory prediction [21], [22], [23]. A recently proposed method incorporates the hybrid dynamics model with a Bayesian intent model [24].

These methods, for the most part, focus on learning probabilistic trajectory models for aircraft in unstructured airspace. Modeling aircraft in terminal regions is a fundamentally different task that cannot be solved with models for unstructured operation. It is difficult to apply them because the structure of the arrival and departure procedures at airports leads to different behavior than in unstructured airspace.

One way to learn models for the terminal region is by using the airport's published instrument procedures. However, doing this leads to a model that has no knowledge of the statistical variation in how the procedures are carried out. One effective way to capture the statistical variation is to learn from historical surveillance data.

Given the large amount of surveillance data available, several papers have proposed using machine learning techniques to learn models of terminal airspace. One method involves clustering *turning points*, defined as spatial positions where a substantial change of heading occurs. Gariel *et al.* [25] cluster turning points across a dataset of trajectories, treat trajectories as a sequence of these clusters, then cluster those discrete sequences using the least common subsequence algorithm. Mahboubi and Kochenderfer [26] find that the turning point model does not perform well on real, noisy radar data, and extend the work of Gariel *et al.* by representing the transitions between turning points at small airports using a Gaussian hidden semi-markov model.

Another promising direction is to cluster trajectories at the level of position measurements, which is also the principle behind air traffic flow modeling of en route traffic [27], [28]. The main challenge in clustering this way is that the sequences are of different lengths. Li *et al.* [29], [30] use the density-based spatial clustering of applications with noise (DBSCAN) algorithm to cluster time-aligned trajectories and use the learned clusters for anomaly detection around airports. They also experiment with clustering the trajectories after performing principal component analysis. Hong and Lee [31] deal with the varying length issue by clustering with the dynamic time warping (DTW) algorithm to discover traffic patterns. They regress an aircraft's arrival time based on its similarity to the learned clusters and discuss how their method could augment terminal radar approach controllers. Conde *et al.* [32] deal with the varying length issue by resampling

trajectories to make them the same length. They then cluster the resampled trajectories using the DBSCAN algorithm and construct classifiers to assign new flights to clusters. Their main application area is resource management. These three lines of research are the closest to ours.

Proposed method. In this paper, we propose a method for learning probabilistic trajectory models of aircraft in terminal airspaces directly from position measurements. The model training procedure does not rely on other (likely unavailable) information, *e.g.*, velocity measurements, flight plans, or air traffic control instructions. Through an unsupervised clustering algorithm, the method is able to discover the departure and approach procedures for an airport, without ever accessing the airport's instrument procedures. The method then fits a generative model based on the intra-cluster covariance matrices.

The main limitations of the three most similar lines of work are: 1) they use ad-hoc methods for dealing with the varying length issue; 2) they do not use the learned clusters to construct a full probabilistic model of motion; and 3) the DTW and DBSCAN algorithms have trouble scaling to large datasets (the papers use datasets on the order of 500 trajectories whereas we use datasets with more than 60 000 trajectories).

Our method deals with the varying length issue, inspired by Li *et al.* [29], by time-aligning trajectories around the time they are closest to the runways and extrapolating trajectories to be the same length. This alignment and extrapolation allows us to compare trajectories directly using the Euclidean norm. Our method constructs a probabilistic model from the clusters, which leads to accurate inference and realistic generation of trajectories. Finally, our method uses the k-means clustering algorithm and a Gaussian mixture model, making it cheap to store, efficient to sample from, efficient to perform inference with, and most importantly easy to understand and ultimately implement.

We find that our method performs well on real, noisy data, and is able to capture relevant statistical properties of the terminal airspace environment. We find that trajectories sampled from the model are realistic. In Fig. 1, we show samples generated from an instance of our model trained on takeoffs at the KJFK airport next to held-out test trajectories. We evaluate our method throughout the paper on the KJFK airport, for illustration purposes, though applying the method to other airports is trivial, as we will see. In fact, another major advantage of our method is that it is relatively automatic; it can produce a generative model from radar data with little to no supervision. Our method can create generative models for thousands of airports, automatically, from (poor, noisy, and intermittent) radar position data.

Limitations. Our method has several limitations that must be addressed. It is unable to incorporate the intent of the air traffic controller, is unable to model the complex spatio-temporal interaction between aircraft, and can only fit models based on historical data (and thus cannot be used to test new instrument procedures).

Summary of contributions:

- A linear time trajectory reconstruction procedure given noisy, irregular position measurements.

- A procedure for automatic unsupervised learning of a Gaussian mixture model for aircraft motion that captures correlation between position dimensions over time.
- Publicly available codes for constructing the model, running the experiments, and pre-trained models for the KJFK airport.³

II. RADAR DATA

Our method expects a set of trajectories of vehicles (in this case, aircraft in a terminal airspace) denoted $\tau_1, \dots, \tau_{N_{\text{traj}}}$. Each trajectory is a set of time and position measurement pairs $\tau = \{(t_i, \hat{p}_i)\}_{i=1}^{N_{\text{meas}}}$ where the time is $t_i \in \mathbf{Z}_+$ and the measurements are $\hat{p}_i \in \mathbf{R}^3$.

A. FAA radar dataset

The radar dataset we use in this paper is derived from the Federal Aviation Administration (FAA) multi-sensor fusion tracker [33], which fuses radar detections of a single aircraft by multiple sensors into a single track. The tracker uses primary radars, mode A/C SSRs, mode S SSRs, Wide Area Multilateration systems, and ADS-B receivers. Each of these surveillance sensor systems has different accuracies, updating rates (1–12 Hz), and susceptibility to anomalies. The data was preprocessed and fused into timestamped entries of target address (used to group entries into flights), position (measured in WGS84 latitude/longitude), pressure altitude (meters) and optionally geometric altitude, lateral velocity, and vertical velocity. We ignore lateral/vertical velocities in our analysis, as it is unclear whether they are from the position measurements or are actual velocity measurements. The data contains six months of flights starting March 2012 from three general locations: Central Florida, New York City, and Southern California.

We describe the procedure for filtering data corresponding to the terminal airspace of an airport, defined by the center of its runways $p_{\text{ref}} \in \mathbf{R}^3$ in earth-centered, earth-fixed (ECEF) coordinates. We convert each position in the dataset to east-north-up (ENU) coordinates determined by the origin p_{ref} . We keep only position measurements that are less than 5 NM in the “east” or “north” dimensions and less than 3000 feet in the “up” dimension, *i.e.*, in the terminal airspace of the airport. Then we iterate through measurements associated with each target address. Since target addresses are reused, we split sequences of measurements that have a gap in measurements of more than 30 s. This filtering results in a set of trajectories for the airport in the format expected by our method, where each trajectory is a set of time and aircraft position measurement pairs $\tau = \{(t_i, \hat{p}_i)\}_{i=1}^{N_{\text{meas}}}$. We shift $t_i = t_i - \min_i t_i$ so that the first measurement in each trajectory is always at $t = 0$.

B. Landing and takeoff separation

We are interested in modeling landings and takeoffs, and since these are different modes of operation, we work with them separately. There are many ways to separate a dataset into landings and takeoffs; we used the following approach.

³www.github.com/sisl/terminal-airspace-models

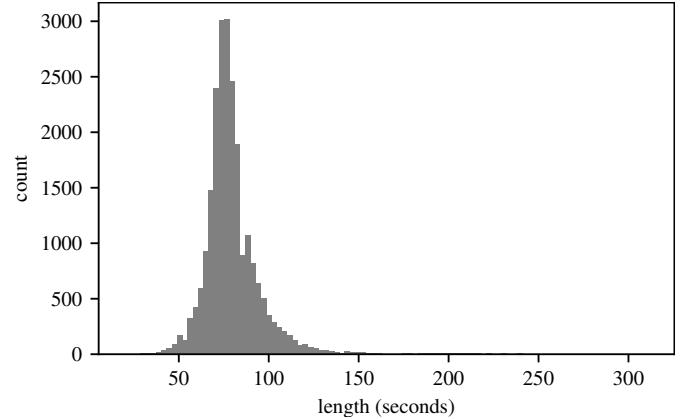


Fig. 2. Histogram of the length of takeoff trajectories at KJFK in seconds. The median is 70 seconds.

For each trajectory, let $T = \max_i t_i$ be the time of the last measurement in that trajectory, $c = \operatorname{argmin}_i \|p_i\|_2$ be the index of the measurement closest to the origin (*i.e.*, the center of the runways), \dot{z}_{avg} be the average “up” velocity in ft/min, and R be approximately the length of the longest runway (2000 m for KJFK). The trajectory is then classified as a

- landing if $(\|p_c\|_2 < R) \wedge (\|p_T\|_2 > R) \wedge (\dot{z}_{\text{avg}} < -200 \text{ ft/min} \wedge (t_c/T > 0.95))$
- takeoff if $(\|p_c\|_2 < R) \wedge (\|p_0\|_2 > R) \wedge (\dot{z}_{\text{avg}} > 200 \text{ ft/min} \wedge (t_c/T < 0.05))$

We used this rule because it seemed to work well with the FAA radar dataset.

For landings, we ignore the position measurements after the index c , as we let the measurement closest to the origin be the last measurement. Similarly, for takeoffs, we ignore the position measurements before the index c . For landings, we reverse time in the trajectory by setting $t_i = T - t_i$ making landings and takeoffs both start near the origin. Reversing time for landings imposes the constraint that for both landings and takeoffs $t_i = 0$ is the point that the trajectory is closest to the runways. From now on, we assume that landings and takeoffs are worked with separately.

C. Preprocessing

The only preprocessing we do is to scale each dimension in each measurement by a separate positive scalar. Pre-scaling allows the dimensions to be compared equally when using ℓ_2 norms, and is also basic standardization practice in machine learning. For aircraft trajectories, we scale the “up” dimension by 10. This is a heuristic, guided by the fact that glide slopes are approximately 3° to 5° ($1/\sin 5^\circ \approx 11$) and by the ratio of the dimensions of the terminal airspace ($5 \text{ NM}/3000 \text{ feet} \approx 10$). We reverse the scaling after we learn the model.

III. MODEL LEARNING

This section outlines the steps one performs to construct the model. We reconstruct the trajectories, cluster them, and then fit a generative model to the clusters.

A. Trajectory reconstruction

There are several complications with the trajectories in their original format. First, some trajectories do not have measurements at every time between their first and last measurement, and some have more than one at a given timestep (which may happen when multiple sensors provide a reading or because of the time-discretization of measurements). Second, the measurements themselves are noisy. Third, we need to extrapolate some trajectories to make them the same duration, which is a requirement for our clustering algorithm. We make the trajectories a common length T_{com} that we choose to be near the median of the trajectory durations across the data. (One can adjust the common length for their specific application; we used median length in this paper.) Fig. 2 shows the histogram of trajectory durations for takeoffs at KJFK. We ignore trajectories that are significantly shorter than T_{com} seconds, as our uncertainty increases the longer we extrapolate. Trajectories shorter than T_{com} are extrapolated to be length T_{com} and trajectories longer than T_{com} are smoothed then truncated to length T_{com} .

We accomplish interpolation, smoothing, and extrapolation by solving the quadratic optimization problem (1), which includes three objective terms. The optimization problem is

$$\underset{P}{\text{minimize}} \quad \|AP - \hat{P}\|_F^2 + \lambda_1 \|D_2 P\|_F^2 + \lambda_2 \|D_3 P\|_F^2 \quad (1)$$

where the optimization variable $P \in \mathbf{R}^{N \times 3}$ is the reconstructed trajectory and its i th row is the position of the reconstructed trajectory at time i , the length $N = \max\{T, T_{\text{com}}\}$, the diagonal matrix $A \in \mathbf{R}^{N \times N}$ has A_{ii} equal to one if there are one or more measurements at time i and zero otherwise, each row in the measurement matrix $\hat{P} \in \mathbf{R}^{N \times 3}$ is the average of the measurements at that time and zero otherwise⁴, $D_2 \in \mathbf{R}^{N-2 \times N}$ is the second-order difference matrix representing the acceleration operator, $D_3 \in \mathbf{R}^{N-4 \times N}$ is the third-order difference matrix representing the jerk operator, $\lambda_1, \lambda_2 \in \mathbf{R}_+$ are regularization hyper-parameters, and $\|\cdot\|_F$ is the Frobenius norm of a matrix. We include the form of the (banded) matrices D_2 and D_3 for the sake of completeness

$$\begin{aligned} (D_2)_i &= e_i - 2e_{i+1} + e_{i+2} \\ (D_3)_i &= -e_{i-1} + 2e_i - 2e_{i+2} + e_{i+3} \end{aligned} \quad (2)$$

where the expression $(A)_i$ denotes the i th row of the matrix A .

The first term of (1) encourages the reconstructed trajectory to be close to the measured points and the second and third terms encourage it to have low acceleration and jerk on average. The regularization parameters λ_1 and λ_2 are selected individually for each trajectory through out-of-sample validation. We perform out-of-sample validation by randomly holding out measurements from that trajectory, fitting trajectories with varying λ_1 and λ_2 on the measurements that were not held out, and selecting the parameters that have the lowest loss (with $\lambda_1 = \lambda_2 = 0$) on the held-out measurements. Fig. 3 shows the reconstructed trajectories on the Pareto-optimal surface

⁴This formulation is equivalent to having an objective function where the ℓ_2 loss is averaged over each measurement per time step.

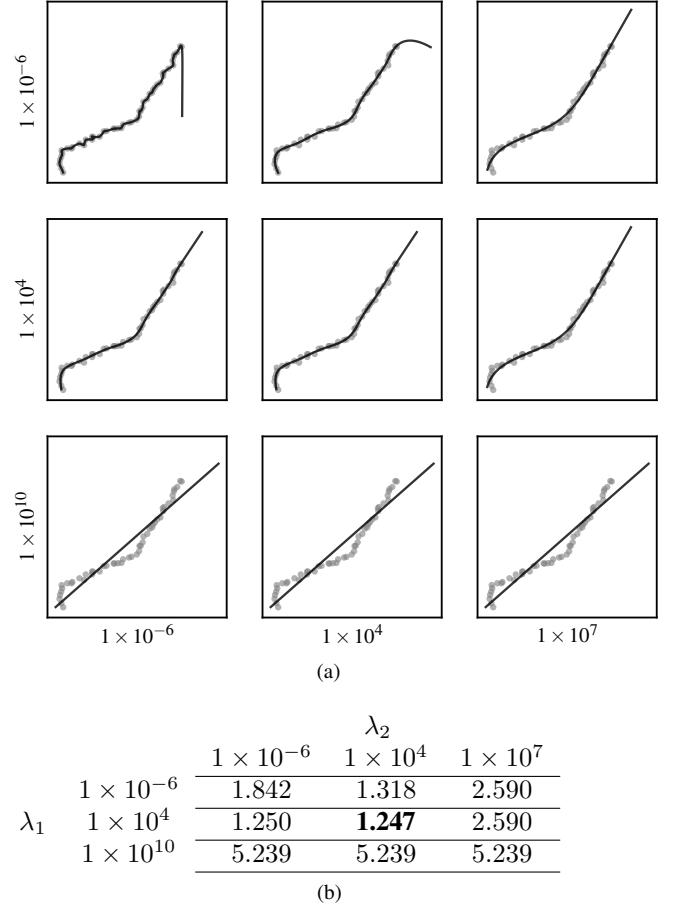


Fig. 3. Trajectory reconstruction procedure. (a) Lateral view of measurements and reconstructed trajectory (including a 20s extrapolation) for various settings of the regularization parameters. A low λ_1 and λ_2 leads to a jagged trajectory (top left) whereas high λ_1 or λ_2 leads to a linear trajectory (bottom right). (b) Out-of-sample validation losses for various regularization parameters. The optimal out-of-sample validation loss is bolded and achieved by $\lambda_1 = 1 \times 10^4$ and $\lambda_2 = 1 \times 10^4$.

that were generated by varying λ_1 and λ_2 , as well as the respective out-of-sample losses, for a simulated aircraft trajectory.

The formulation could potentially be improved by switching to a different loss, *e.g.*, Huber or ℓ_1 , and by adding constraints on the reconstructed trajectory. We stick to the form in (1) mainly for the sake of simplicity of implementation; solving (1) reduces to solving a 9-banded linear system, which has an asymptotic complexity of $O(N)$. We believe that the trajectory reconstruction procedure described here could be useful for trajectory smoothing for a wide variety of vehicles.

B. Clustering

After the trajectories are reconstructed, they all have the same length. Vectorizing the trajectories (stacking the columns of the matrix P) makes each trajectory a column vector $p^{(i)} \in \mathbf{R}^{3T_{\text{com}}}$, $i = 1, \dots, N_{\text{traj}}$. The first, second, and third partitions of the vector are the east, north, and up positions over time.

Since aircraft in terminal airspace, for the most part, follow pre-defined procedures and only slightly deviate from them, it is reasonable to believe that the trajectories partition well

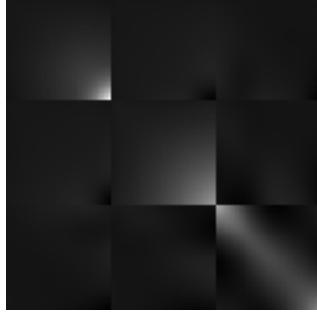


Fig. 4. Sample covariance matrix for the most probable cluster for takeoffs at KJFK. Higher values of covariance are a lighter shade. The first third of the matrix corresponds to the east dimension, the second third corresponds to north dimension over time, and the final third corresponds to up dimension over time. Uncertainty increases with time, as aircraft are more likely to deviate from their nominal path later in the takeoff maneuver.

into clusters. Each cluster should have trajectories that follow a similar path. We can measure the similarity between two trajectories $p^{(1)}$ and $p^{(2)}$ using their Euclidean distance $\|p^{(1)} - p^{(2)}\|_2$. Given a number of groups, denoted K , we want to choose K centers μ_j and assignments $c_i \in \{1, \dots, K\}$ to each trajectory to minimize the total distance between each point and its assigned cluster center. We compute an approximate solution to this problem using the K-means algorithm with K-means++ initialization [34], [35].

The K-means procedure results in K archetypal trajectories μ_j and N_{traj} assignments c_i . For each cluster, we calculate the intra-cluster sample covariance matrices $Q_j \in \mathbf{S}_{++}^{3T_{\text{com}}}$. Fig. 4 provides a visualization of one of the intra-cluster covariance matrices. We let the vector π denote cluster frequencies, where π_j is the fraction of trajectories in cluster j .

C. Gaussian mixture model

Our intra-cluster sample covariance matrices are noisy, as they are estimated with fewer samples than there are entries in the matrix. We would like to remove as much of this noise as possible, so we perform a low-rank approximation to the covariance matrix Q_j . The solution to the low-rank approximation problem comes from the eigendecomposition (or equivalently the singular value decomposition) of Q_j . The solution is $\tilde{Q}_j = U_j \Sigma_j U_j^T$, where U_j is a matrix with the first r singular vectors, and Σ_j is a diagonal matrix with the first r singular values. The columns of U_j can be interpreted as the *principal deviations* for motion in cluster j . (Note that if we save just the principal deviations instead of Q_j it significantly reduces the overall size of the trained model and speeds up sampling and inference.) As a final step, we perform the inverse of the preprocessing in Section II, which involves dividing the “up” dimension by 10 and reversing time for landings.

We use $r = 5$ principal deviations for aircraft at KJFK because it captures most of the intra-cluster variance, though in practice r should be chosen using the procedure in Sec-

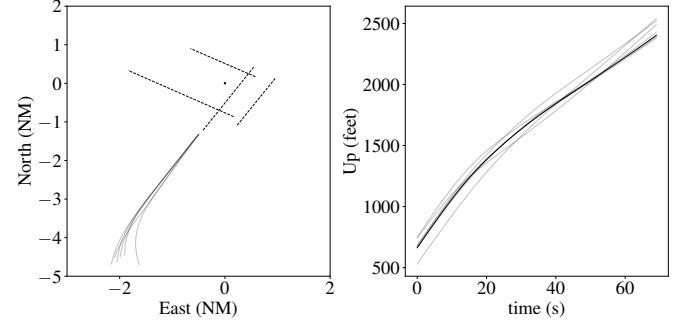


Fig. 5. Archetypal trajectory (black) and top-5 principal deviations (gray) for the most probable cluster of the model for takeoffs at KJFK with $K = 20$.

tion III-D. The 5 principal deviations for the most probable cluster for takeoffs at KJFK are in Fig. 5.

We now construct a generative probabilistic model. Our model for aircraft trajectories in terminal airspaces is the following Gaussian mixture model (GMM)

$$\begin{aligned} j &\sim \text{Categorical}(\pi) \\ p &\sim \mathcal{N}(\mu_j, \tilde{Q}_j). \end{aligned} \quad (3)$$

It is a well known fact that conditioning a Gaussian distribution on measurements gives another Gaussian distribution. This makes GMMs easy to sample from and perform inference with.

This fact leads to a convenient interpretation of the assumed trajectory distribution. Because the marginal distribution of a Gaussian is another Gaussian, the first measurement of the trajectory comes from a Gaussian distribution. Then the position at time t conditioned on the previously sampled points similarly follows a Gaussian distribution. Drawing a full trajectory from the model is equivalent to successively drawing from Gaussian distributions dependent on the past.

1) *Positive Definiteness of Covariance Matrices:* Because of the low-rank approximation step, our covariance matrices are now positive semidefinite and no longer positive definite, which is a requirement to have a valid multivariate density. We can deal with this inconsistency by restricting the density to a $\text{rank}(Q_j)$ subspace where the Gaussian distribution is supported, which results in the following density:

$$f(x) = (\det^*(2\pi Q_j))^{-\frac{1}{2}} \exp(-\frac{1}{2}(x - \mu_j)^T Q_j^\dagger (x - \mu_j)), \quad (4)$$

where Q_j^\dagger is the Moore-Penrose inverse of Q_j and \det^* is the pseudo-determinant, or the product of the non-zero singular values [36].

Similarly, to sample from cluster j we sample $z \sim \mathcal{N}(0, I_r)$ and then emit $\mu_j + U_j \Sigma_j^{1/2} z$. The sample is just an archetypal trajectory plus a random linear combination of the principal deviations weighted by the square root of the singular values, hence the name principal deviations. Fig. 6 shows a cluster center and samples from that cluster’s distribution for the most probable cluster at KJFK.

D. Automatic model generation

We now discuss the claim, made in the Introduction, that the method can go from radar data to generative model

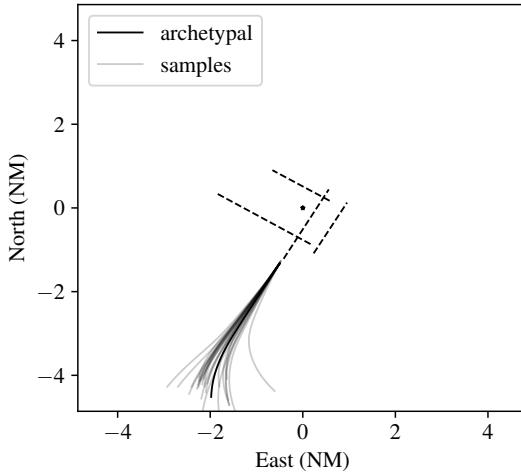


Fig. 6. Samples from the most probable cluster takeoff model at KJFK with $K = 20$. The cluster covariance naturally captures plausible takeoffs at KJFK.

with little to no supervision. The trajectory reconstruction procedure (Section III-A) can be performed independently on each trajectory and requires no parameters. The clustering procedure (Section III-B) requires only one parameter, K . The generative model fitting (Section III-C) also requires only one parameter, r . We can select K and r , automatically, using an out-of-sample validation procedure that depends on the downstream task. First, we split the dataset into a training and held-out set. We fit generative models with varying values of K and r , calculate the performance measure (dependent on the downstream task) of the generative model on the held-out data, and keep the model that has the highest performance on the held-out set. We demonstrate this procedure for two downstream tasks, generation and inference, in Section IV.

IV. EXPERIMENTS

We evaluate our method on takeoffs and landings at the KJFK airport. Details on the layouts of KJFK are given in the Appendixes. The plots in the paper are for the KJFK airport with $K = 50$ clusters for takeoffs and $K = 70$ clusters for landings, as those seemed to be reasonable choices for the number of clusters for the purpose of illustration. In practice, K should be chosen using the procedure in Section III-D.

A. Archetypal trajectories

We first sanity check the model by visualizing the learned centers. Fig. 7 shows the learned centers for landings at KJFK and Fig. 8 shows the centers for takeoffs at KJFK. The centers seem to capture the plausible ways of approaching and departing KJFK. Overall, the glide slope for landings seems to be of constant slope and the climb rates for takeoffs seem to start high and then level off. (The different phases for takeoffs likely correspond to the form of the takeoff procedures at KJFK.) Several of the archetypal trajectories look super-imposed, which seems redundant. But in fact, different archetypes can capture not only the procedures, but also approaches at different glide slopes, airspeeds, and wind

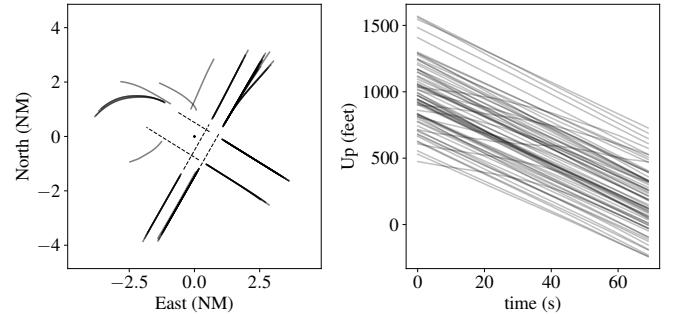


Fig. 7. Model for landings at KJFK with $K = 70$. Left, lateral view of archetypal trajectories. Right, archetypal trajectories vertical profile. Many of the archetypal trajectories overlap in lateral view but not in vertical profile. Each center corresponds to a slightly different glide slope.

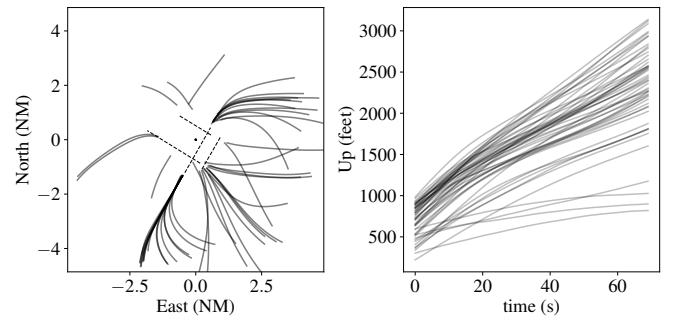


Fig. 8. Model for takeoffs at KJFK with $K = 50$. Left, archetypal trajectories lateral view. Right, archetypal trajectories vertical profile. As with landings, each cluster corresponds to different climb rates.

conditions. For example, a Challenger 300 lands slower than a much larger Boeing 747, so we might end up with clusters for smaller and larger aircraft. These nuances are not coded into the procedure, but rather learned directly from the data.

One observation worth discussing is that there seem to be clusters that represent invalid approaches/departures, *e.g.*, a runway does not exist there or clusters have a negative/constant vertical profile. These invalid clusters correspond to a small π_j , and thus low probability, under the model. That there are clusters representing invalid maneuvers likely means that there are anomalous trajectories in the data. These anomalous trajectories could be filtered using outlier detection, but we leave that for future work.

B. Statistical properties

To evaluate the statistical properties of the model, we randomly split the trajectories into a held-out set and training set. Then, we trained the model on the training set and evaluated its performance on the held-out set. To evaluate whether or not the model is able to capture the overall statistics of vehicle motion in the environment, we devised several visual and quantitative experiments and tested the model over different values of K .

We have a set of held-out test samples and independently sampled trajectories from the model. We first plot the overall histogram of positions, with a logarithmic color scale, for generated landings at KJFK (Fig. 9), held-out landings at

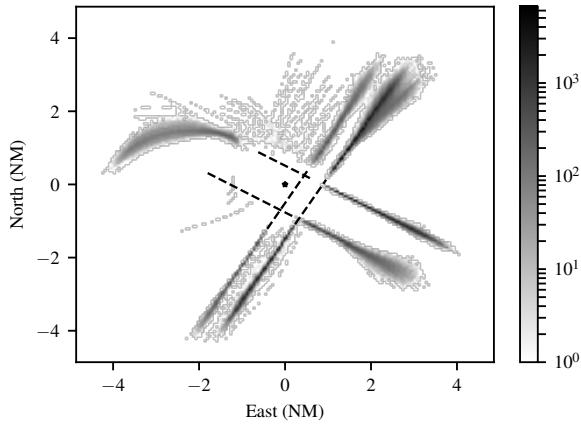


Fig. 9. Lateral view of log-histogram of **generated** landing flight tracks. Compare with Fig. 10.

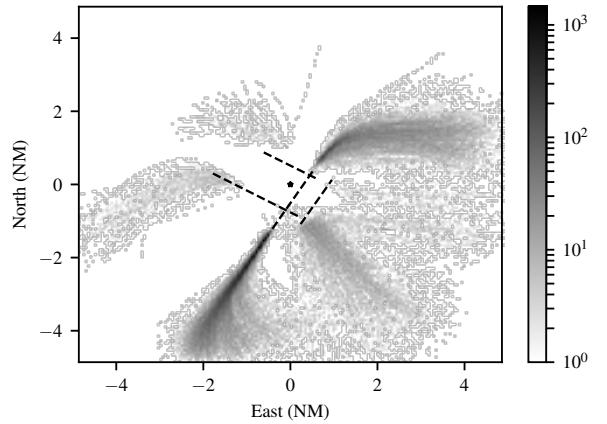


Fig. 11. Lateral view of log-histogram of **generated** takeoff flight tracks. Compare with Fig. 12.

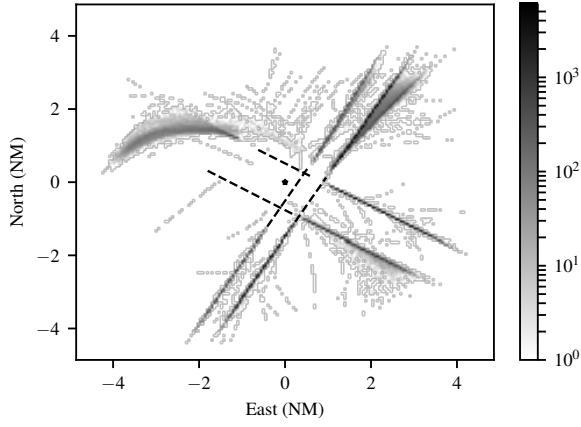


Fig. 10. Lateral view of log-histogram of **real held-out** landing flight tracks. Compare with Fig. 9.

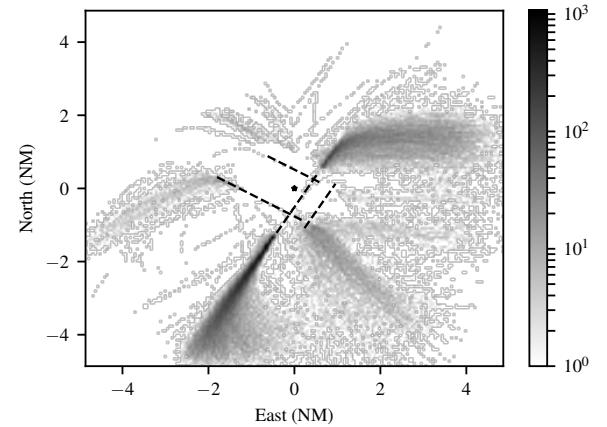


Fig. 12. Lateral view of log-histogram of **real held-out** takeoff flight tracks. Compare with Fig. 11.

KJFK (Fig. 10), generated takeoffs at KJFK (Fig. 11) and held-out takeoffs at KJFK (Fig. 12). The $5 \times 5 \text{ NM}^2$ grid is split into 400×400 bins making each bin $23.15 \times 23.15 \text{ m}^2$.

We also would like to make sure that the model captures realistic physical properties of the trajectories outside of just positions, so we calculate the sample longitudinal velocity, vertical velocity, and turn rate for generated samples and compare that to those in held-out data. Fig. 15 in the Appendix shows the generated/held-out turn-rate distributions for takeoffs at KJFK. (The figures for vertical velocity and turn rate were omitted for the sake of space.)

We now perform the out-of-sample validation procedure described in Section III-D for two tasks: generation and prediction. For generation, we calculate the Kullback-Leibler divergence (KL-divergence) to measure the similarity between the (smoothed) empirical distributions of position, longitudinal velocity, vertical velocity, and turn rate,

$$D_{\text{KL}}(P\|Q) = - \sum_i P(i) \log \frac{Q(i)}{P(i)}. \quad (5)$$

We then average the KL-divergence across the four distributions to get an overall measure of distributional similarity.

For prediction, we give the model the first 10 measurements of every held-out trajectory, calculate the closest cluster, and calculate the posterior mean conditioned on those 10 measurements. We then calculate the root-mean-square (RMS) error between the posterior mean (the prediction) and the actual rest of the trajectory. Fig. 13 shows these two metrics, generation and inference, for varying values of K at KJFK. (The parameter r was set to 5 in all experiments.)

For KJFK, from the out-of-sample validation procedure, we can conclude that the optimal number of clusters K for prediction for takeoffs is 10 (RMS=322.22) and for landings is 20 (RMS=106.86). The inflated RMS value for takeoffs is likely because it is challenging to predict where an aircraft will go given just the direction it took off in, but for landings the model is able to identify the procedure from only 10 measurements. Increasing K , in the case of prediction, leads to *overfitting*, where the generative model fits the training data too well and fails to generalize to held-out data. However, for generation, increasing K does not lead to overfitting. This result is likely because different clusters may represent, e.g., different aircraft types and wind conditions, which are not coded into the airport procedures. These clusters are not

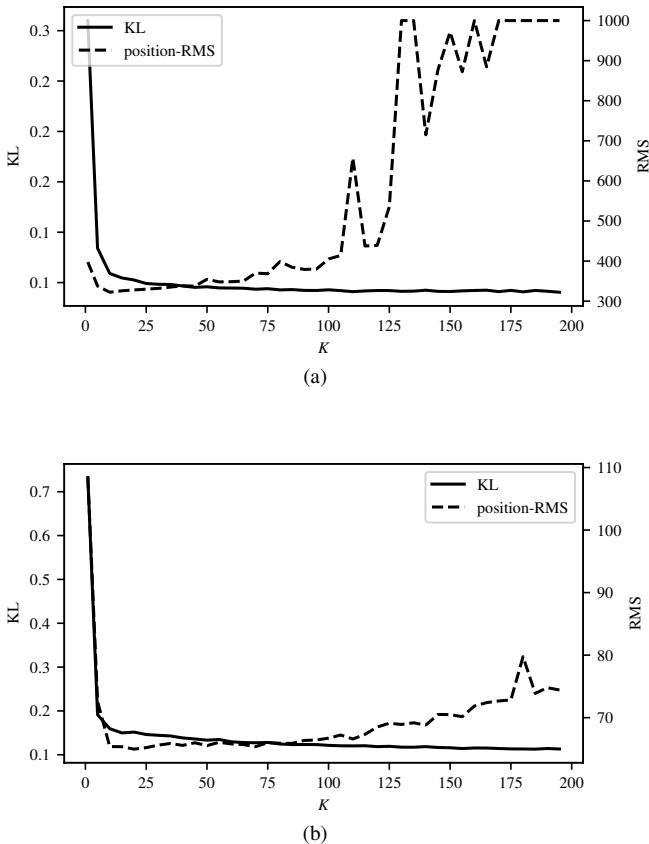


Fig. 13. Model for KJFK takeoffs. Log-plot of average KL-divergence and prediction RMS error vs. number of clusters (lower is better in both cases). The RMS error was clipped at 1000 m for visualization purposes. (a) takeoffs, optimal KL is 0.002 and the optimal RMS is 322.22 m ($K = 10$); (b) landings, optimal KL is 0.005 and the optimal RMS is 106.86 m ($K = 20$).

important for RMS-based prediction, but are important for accurately capturing the statistical properties of the terminal airspace environment. We now investigate the visual fidelity of samples generated from the model.

C. Pilot turing test

In Fig. 1 we showed samples from the model and held-out test samples side-by-side and inquired the reader to guess which one was generated before seeing their true identities. Throughout the experiments, we found the generated samples to be realistic. To evaluate how realistic the samples are scientifically, we devised what we call a “Pilot Turing Test” (PTT), akin to the Turing Test in Artificial Intelligence [37].

To perform a PTT, we first take 25 random samples generated from the model and 25 random held-out samples. We added a small amount of random noise to the real and fake samples to obscure radar anomalies. Then, we give the pilot 25 pairs of generated and real trajectories in random order, where each plot has a lateral view of the trajectory and its vertical profile, as in Fig. 8. We ask the pilot to label one trajectory in each pair as real and one as generated. Then we calculate the accuracy, defined as the fraction of trajectories the pilot labeled as real that were actually real. We want the accuracy

TABLE I
PILOT TURING TEST ACCURACIES

Pilot 1		Pilot 2	
Landings	Takeoffs	Landings	Takeoffs
52 %	52 %	44 %	68 %

to be around 50 %, which indicates that the generated and real trajectories are indistinguishable to an expert.

We conducted a PTT with two licensed pilots for landings/takeoffs at KJFK. The first pilot (Pilot 1) is a professional Flight Test Captain for a U.S. airline. The second pilot (Pilot 2) is a professional Boeing 747 pilot with over 13 000 hours of total flying experience. The resulting accuracies for both pilots can be found in Table I. The files used to conduct the PTT can be found in the code. On average, we achieved an accuracy of 54 %, which means the pilots performed slightly better than random guessing. We can conclude that our model is able to generate trajectories that are practically (visually) indistinguishable from real trajectories.

V. APPLICATIONS

This section describes how to generate trajectories and how to perform inference. Our model can be used for other applications, including anomaly detection, using the clusters to adjust procedures, and reactionary motion planning.

A. Generation

As described in Section III, we can draw samples from the assumed probability distribution over aircraft motion. We draw samples by first sampling from the categorical distribution over clusters, sampling the vector $z \sim \mathcal{N}(0, I_r)$, and then emitting the trajectory $\mu_j + U_j \Sigma_j^{1/2} z$. If we further condition the categorical distribution over clusters on external factors such as time of day or previously sampled clusters, we can construct 3D simulations of aircraft operating in particular terminal airspaces. We could even integrate it into an airport queuing model [38] to construct even more realistic simulations.

The ability to generate samples could be particularly useful for safety and performance evaluation of decision support systems. For example, the model could be used for verification of AutoATC, an automated air traffic control system for nontowered airports. The model would provide 3D simulations of the terminal airspace, feed the resulting position traces into a simulation of AutoATC, and identify failure cases where AutoATC failed to provide the correct guidance to avoid a collision. It could also be used to perform ablation studies, e.g., by increasing the “pressure” on the airport by increasing the number of clusters sampled per time period or scaling the covariance matrices in the clusters to emit less realistic trajectories.

B. Inference

As described in Section III, we can condition each cluster’s Gaussian distribution on partial measurements. By performing full Bayesian inference in the Gaussian mixture model, we can calculate a probability distribution over clusters and time index

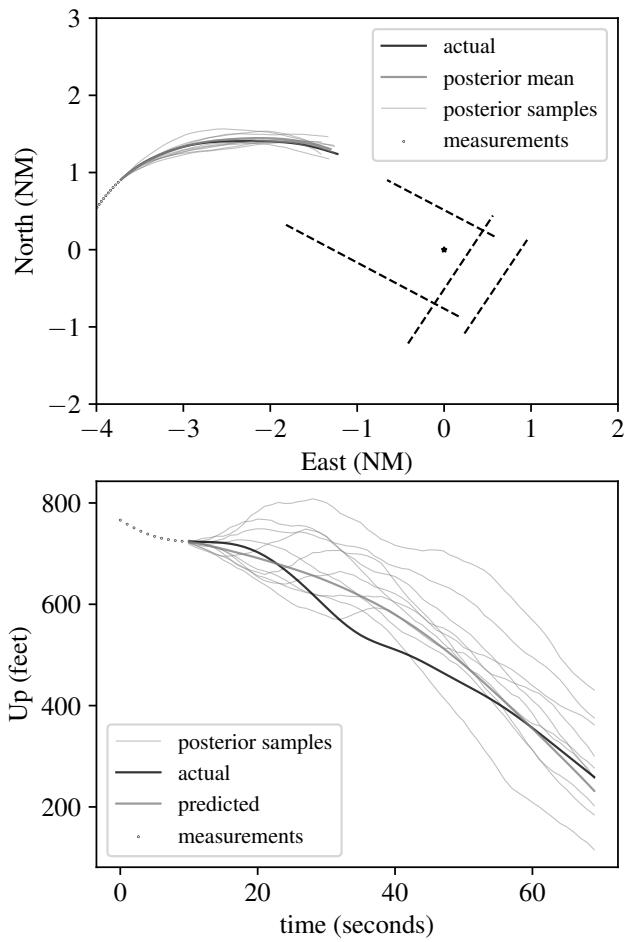


Fig. 14. Demonstration of inference for a KJFK landing model with $K = 20$. Given 10 samples (dots) of a landing trajectory (black), the model produces a prediction of where the aircraft will go (dark gray) and samples from the posterior distribution (light gray). The average RMS error for the prediction is 77.41 m.

in that cluster's distribution. For each of these possibilities, we can also calculate the posterior distribution over future and past measurements by conditioning that cluster's multivariate Gaussian on the partial measurements. We demonstrate the result of this procedure for an approach at KJFK in Fig. 14.

Being able to perform inference directly leads to several applications. If we had, for example, 10 radar measurements of an aircraft approaching an airport, we can immediately provide non-trivial answers to useful questions (we provide the answers for the inference that was performed in Fig. 14):

- 1) Which runway will the aircraft land on? Runway 13L.
- 2) How long until the aircraft lands? ≈ 60 s.
- 3) Which approach procedure is the aircraft performing?
PARKWAY VISUAL RWY 13L/R.

If they were departing, we could answer a similar set of questions:

- 1) How long until the aircraft exits the terminal airspace?
- 2) Which direction will the aircraft exit the terminal airspace?
- 3) Which departure procedure is the aircraft performing?

If we had the posterior distribution for two aircraft, we

could perform a Monte Carlo simulation to calculate a conflict probability by sampling a large number of possible trajectories from both posterior distributions and calculating the proportion that violated some conflict criteria. This posterior distribution could also be used as a model of other aircraft in an on-board planning system for collision avoidance, or as input to an automated air traffic control system.

VI. CONCLUSION

In this paper, we proposed a method for learning probabilistic motion models for aircraft in terminal airspaces directly from position measurements. Our method could be extended to a larger terminal airspace, e.g., < 50 NM from the airport, and be used to aid terminal radar approach controllers as in Hong and Lee [31]. We also believe that a slightly altered version of our method could be applied to vehicles in other contexts. It is likely to work well for highly structured environments, e.g., loading dock for trucks, parking lots, intersections with stoplights/stop signs, or shipping ports. On the other hand, the method not perform as well in unstructured environments, such as interstate freeways and unstructured airspace.

The main limitation and hence direction for future work is to incorporate pairwise interaction. We can assume that aircraft actively avoid each other, which affects the probability distribution over future motion. We believe that modeling pairwise interaction can be achieved with a framework similar to ours, and are currently looking into extending the methodology to incorporate it.

APPENDIX A KJFK AIRPORT

We evaluate our method on the John F Kennedy International Airport (KJFK) airport. KJFK is a towered airport in Queens, New York with four runways: 13R/31L, 4L/22R, 13L/31R, and 4R/22L. As of 2016, KJFK has 1256 aircraft operations per day, 91 % commercial, 7 % air taxi, 2 % transient general aviation and less than 1 % military [39].

APPENDIX B SUPPLEMENTAL FIGURES

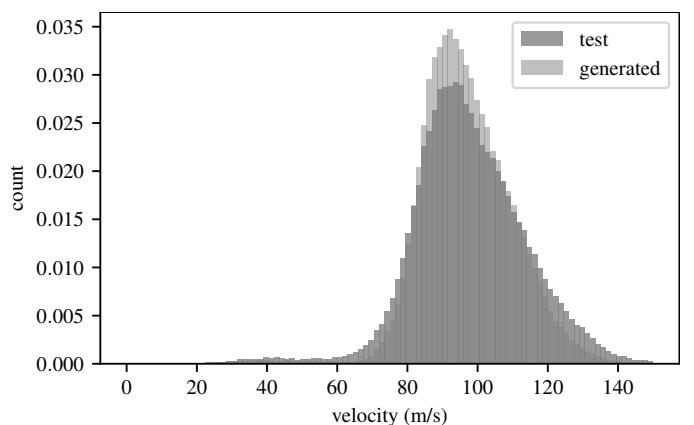


Fig. 15. Model for takeoffs at KJFK with $K = 50$. Histogram of lateral velocities for generated and held-out test data.

ACKNOWLEDGMENT

The authors acknowledge the FAA and the MITRE Corporation for sharing the surveillance data. The authors also thank E. Maki for thoughtful discussion, and R. Abraham and G. Knox for their participation in the Pilot Turing Test.

REFERENCES

- [1] “Modernization of U.S. Airspace,” <https://www.faa.gov/nextgen>, 2018.
- [2] Boeing, “Statistical summary of commercial jet airplane accidents,” Boeing, Tech. Rep., 2016.
- [3] M. Enriquez, “Identifying temporally persistent flows in the terminal airspace via spectral clustering,” in *Tenth USA/Europe Air Traffic Management Research and Development Seminar (ATM2013)*, 2013.
- [4] D. Vasquez and T. Fraichard, “Motion prediction for moving objects: a statistical approach,” in *Intl. Conf. on Robotics and Automation*, vol. 4. IEEE, 2004, pp. 3931–3936.
- [5] G. Chatterji, B. Sridhar, and K. Bilimoria, “En-route flight trajectory prediction for conflict avoidance and traffic management,” in *Proc. AIAA Guidance, Navigation, and Control Conf.*, July 1996, p. 3766.
- [6] G. Chatterji, “Short-term trajectory prediction methods,” in *Proc. AIAA Guidance, Navigation, and Control Conf.*, 1999, p. 4233.
- [7] R. Slattery and Y. Zhao, “Trajectory synthesis for air traffic automation,” *Journal of Guidance, Control, and Dynamics*, vol. 20, no. 2, pp. 232–238, 1997.
- [8] A. Warren and Y. Ebrahimi, “Vertical path trajectory prediction for next generation atm,” in *Proc. 17th Digital Avionics Systems Conf.*, vol. 2. IEEE, 1998, pp. F11–1.
- [9] M. Kochenderfer, J. Kuchar, L. Espindle, and J. Griffith, “Uncorrelated encounter model of the national airspace system,” Massachusetts Institute of Technology, Lincoln Laboratory, Project Report ATC-345, 2008.
- [10] M. Kochenderfer, M. Edwards, L. Espindle, J. Kuchar, and D. Griffith, “Airspace encounter models for estimating collision risk,” *Journal of Guidance, Control, and Dynamics*, vol. 33, no. 2, pp. 487–499, 2010.
- [11] R. Paielli and H. Erzberger, “Conflict probability estimation for free flight,” *Journal of Guidance, Control, and Dynamics*, vol. 20, no. 3, pp. 588–596, Oct 1997.
- [12] C. van Daalen and T. Jones, “Fast conflict detection using probability flow,” in *Automatica*, vol. 45, no. 8, Aug 2009, pp. 1903–1909.
- [13] L. Pienaar and T. Jones, “The application of probability flow for conflict detection near airports,” in *Proc. AIAA Guidance, Navigation, and Control Conf.*, Jan 2015, p. 1325.
- [14] J. Chryssanthacopoulos, M. Kochenderfer, and R. Williams, “Improved Monte Carlo sampling for conflict probability estimation,” in *Proc. 51st AIAA/ASME/ASCE/AHS/ASC Structures, Structural Dynamics, and Materials Conf.*, Apr 2010, p. 3012.
- [15] R. Coppenbarger, “Climb trajectory prediction enhancement using airline flight-planning information,” in *Proc. AIAA Guidance, Navigation, and Control Conf.*, 1999, p. 170.
- [16] W. Chan, R. Bach, and J. Walton, “Improving and validating ctas performance models,” in *Proc. AIAA Guidance, Navigation, and Control Conf.*, Aug 2000, p. 4476.
- [17] RTCA, “Minimum Operational Performance Standards for Traffic Alert and Collision Avoidance System II (TCAS II),” RTCA STD DO-185B, Washington, D.C., Tech. Rep., June 2008.
- [18] M. Kochenderfer, J. Holland, and J. Chryssanthacopoulos, “Next generation airborne collision avoidance system,” *Lincoln Laboratory Journal*, vol. 19, no. 1, pp. 17–33, 2012.
- [19] E. Harkleroad, A. Vela, J. Kuchar, B. Barnett, and R. Merchant-Bennett, “Risk-based modeling to support nextgen concept assessment and validation,” *Lincoln Lab Report, ATC-405*, 2013.
- [20] M. Prandini, J. Hu, J. Lygeros, and S. Sastry, “A probabilistic approach to aircraft conflict detection,” *IEEE Trans. Intell. Transp. Syst.*, vol. 1, no. 4, pp. 199–220, Dec 2000.
- [21] J. Yepes, I. Hwang, and M. Rotea, “New algorithms for aircraft intent inference and trajectory prediction,” *Journal of Guidance, Control, and Dynamics*, vol. 30, no. 2, pp. 370–382, 2007.
- [22] I. Hwang and C. Seah, “Intent-Based probabilistic conflict detection for the next generation air transportation system,” in *Proc. of the IEEE*, vol. 96, no. 12, Dec 2008, pp. 2040–2059.
- [23] W. Liu and I. Hwang, “Probabilistic trajectory prediction and conflict detection for air traffic control,” *Journal of Guidance, Control, and Dynamics*, vol. 34, no. 6, pp. 1779–1789, 2011.
- [24] C. Lowe and J. How, “Learning and predicting pilot behavior in uncontrolled airspace,” in *Proc. AIAA Infotech at Aerospace*, 2015, p. 1199.
- [25] M. Gariel, A. Srivastava, and E. Feron, “Trajectory clustering and an application to airspace monitoring,” *IEEE Trans. Intell. Transp. Syst.*, vol. 12, no. 4, pp. 1511–1524, July 2011.
- [26] Z. Mahboubi and M. Kochenderfer, “Learning traffic patterns at small airports from flight tracks,” *IEEE Trans. Intell. Transp. Syst.*, vol. 18, no. 4, pp. 917–926, 2017.
- [27] M. Gariel, *Toward a graceful degradation of air traffic management systems*. Georgia Institute of Technology, 2010.
- [28] A. Marzuoli, M. Gariel, A. Vela, and E. Feron, “Data-based modeling and optimization of en route traffic,” *Journal of Guidance, Control, and Dynamics*, vol. 37, no. 6, pp. 1930–1945, 2014.
- [29] L. Li, M. Gariel, R. Hansman, and R. Palacios, “Anomaly detection in onboard-recorded flight data using cluster analysis,” in *Proc. 2011 IEEE/AIAA 30th Digital Avionics Systems Conf.*, Oct 2011, pp. 4A4–1.
- [30] L. Li, J. Hansman, R. Palacios, and R. Welsch, “Anomaly detection via a gaussian mixture model for flight operation and safety monitoring,” *Transportation Research Part C: Emerging Technologies*, vol. 64, pp. 45–57, 2016.
- [31] S. Hong and K. Lee, “Trajectory prediction for vectored area navigation arrivals,” *Journal of Aerospace Information Systems*, vol. 12, no. 7, pp. 490–502, 2015.
- [32] M. Conde Rocha Murca, R. DeLaura, J. Hansman, R. Jordan, T. Reynolds, and H. Balakrishnan, “Trajectory clustering and classification for characterization of air traffic flows,” in *Proc. 16th AIAA Aviation Technology, Integration, and Operations Conf.*, 2016, p. 3760.
- [33] C. Jagodnik, J. Stella, and D. Varon, “Fusion tracking in air traffic control,” *Journal of Air Traffic Control*, vol. 50, no. 1, May 2008.
- [34] S. Lloyd, “Least squares quantization in PCM,” *IEEE Trans. Inf. Theory*, vol. 28, no. 2, pp. 129–137, Mar 1982.
- [35] D. Arthur and S. Vassilvitskii, “k-means++: The advantages of careful seeding,” in *Proc. ACM-SIAM Symp. on Discrete Algorithms*, Jan 2007, pp. 1027–1035.
- [36] C. Rao, *Linear Statistical Inference and its Applications*. Wiley New York, 1973, vol. 2.
- [37] A. Turing, “Intelligent machinery, a heretical theory,” *Philosophia Mathematica*, vol. 4, pp. 256–260, 1996.
- [38] H. Idris, J. Clarke, R. Bhuvan, and L. Kang, “Queuing model for taxi-out time estimation,” *Air Traffic Control Quarterly*, vol. 10, no. 1, pp. 1–22, 2002.
- [39] “AirNav: KJFK - John F Kennedy Int. Airport,” <http://www.airnav.com/airport/KJFK>, 2018.



Shane T. Barratt is currently working towards the Ph.D. degree in electrical engineering at Stanford University. He received the B.S. degree in electrical engineering and computer science from the University of California, Berkeley, in 2017.



Mykel J. Kochenderfer received the B.S. and M.S. degrees in computer science from Stanford University in 2003 and the Ph.D. degree from The University of Edinburgh in 2006. He was with MIT Lincoln Laboratory, where he worked on airspace modeling and aircraft collision avoidance. Since 2013, he has been an Assistant Professor of aeronautics and astronautics with Stanford University.



Stephen P. Boyd is the Samsung Professor of Engineering, and Professor of Electrical Engineering in the Information Systems Laboratory at Stanford University, with courtesy appointments in Computer Science and Management Science and Engineering. He received the A.B. degree in Mathematics from Harvard University in 1980, and the Ph.D. in Electrical Engineering and Computer Science from the University of California, Berkeley, in 1985, and then joined the faculty at Stanford.