

Measuring and Modeling En Route Flight Efficiency: the US Experience

Yulin Liu¹; Mark Hansen¹; Wenzhe(Emma) Ding¹;

Michael Ball²; Cara Chuang²; David Lovell²

1. University of California, Berkeley

2. University of Maryland, College Park

July 21st, 2016

Outline

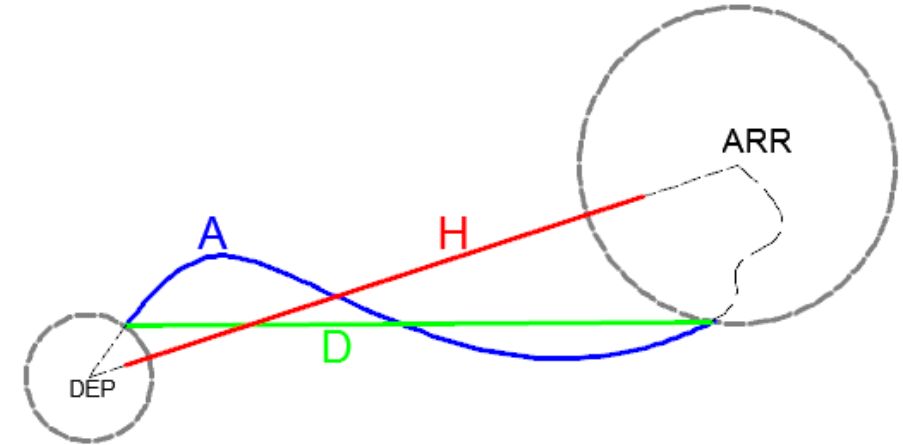
- Introduction
- Data Sources and Preliminary Statistical Analysis
- Macroscopic Model
- Microscopic Model
- Route Selection Model
- Conclusions

Background

- FAA and EuroControl published metrics to evaluate flight en route inefficiencies;
- Understanding the causal factors behind the inefficiency is of great importance.

Defining En route Inefficiency

$$\frac{A - H}{H}$$



- A: Actual flown distance;
- D: Great circle distance between local entry and exit point;
- H: Achieved distance (related to great circle distances from exit/entry points to arcs surrounding arrival/departure airports).

Sources: <https://www.eurocontrol.int/sites/default/files/content/documents/single-sky/pru/news-related/2013-05-08-slides-workshop-achieved-distance.pdf>

Project Goals

- “Evaluation of En Route performance Measures”
- Support FAA in developing en route efficiency performance metrics
- For selected metrics, identify reasons for inefficiency
 - NAS route structure
 - Traffic management initiatives
 - Winds
 - Convective weather
- Eventually allow comparison with other ANSPs such as Eurocontrol

Overview

- Develop models to evaluate how flight en route inefficiency varies based on geographic and seasonal factors;
- Apply clustering algorithms to raw trajectory data for selected OD pairs
- Include cluster membership as explanatory variable in multi-variate model of en route inefficiency
- Analyze within-cluster and between-cluster contributions to en route inefficiency

Outline

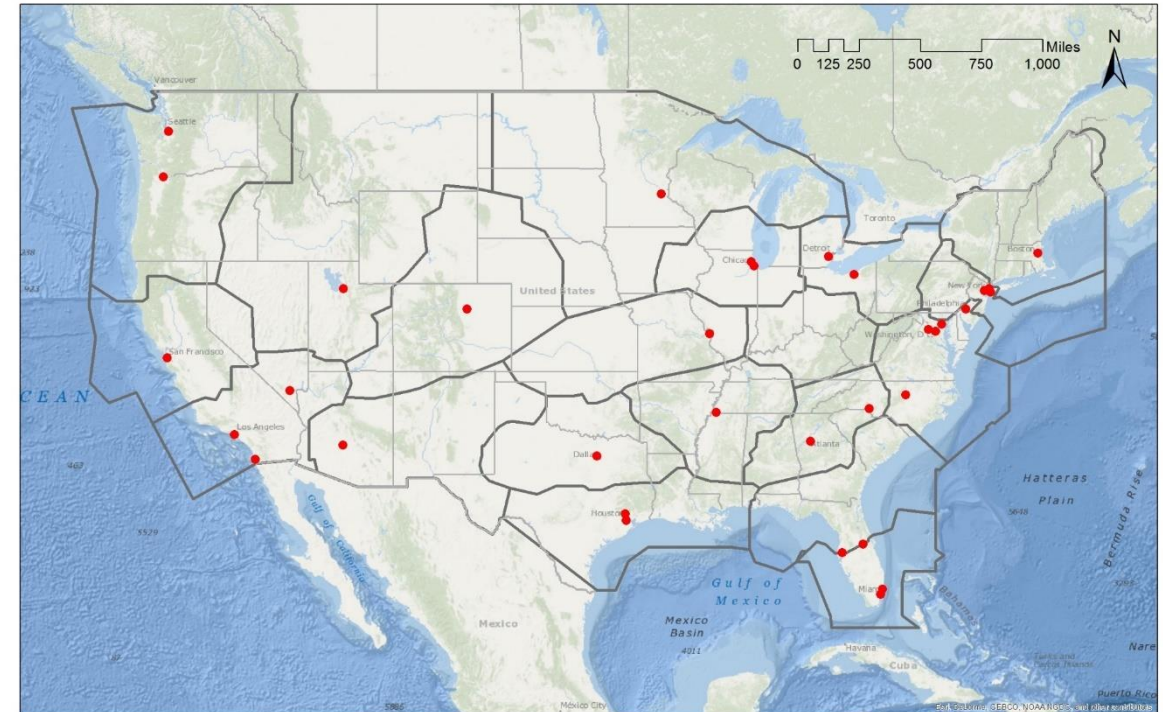
- Introduction
- **Data Sources and Preliminary Statistical Analysis**
- Macroscopic Model
- Microscopic Model
- Route Selection Model
- Conclusions

Data Sources

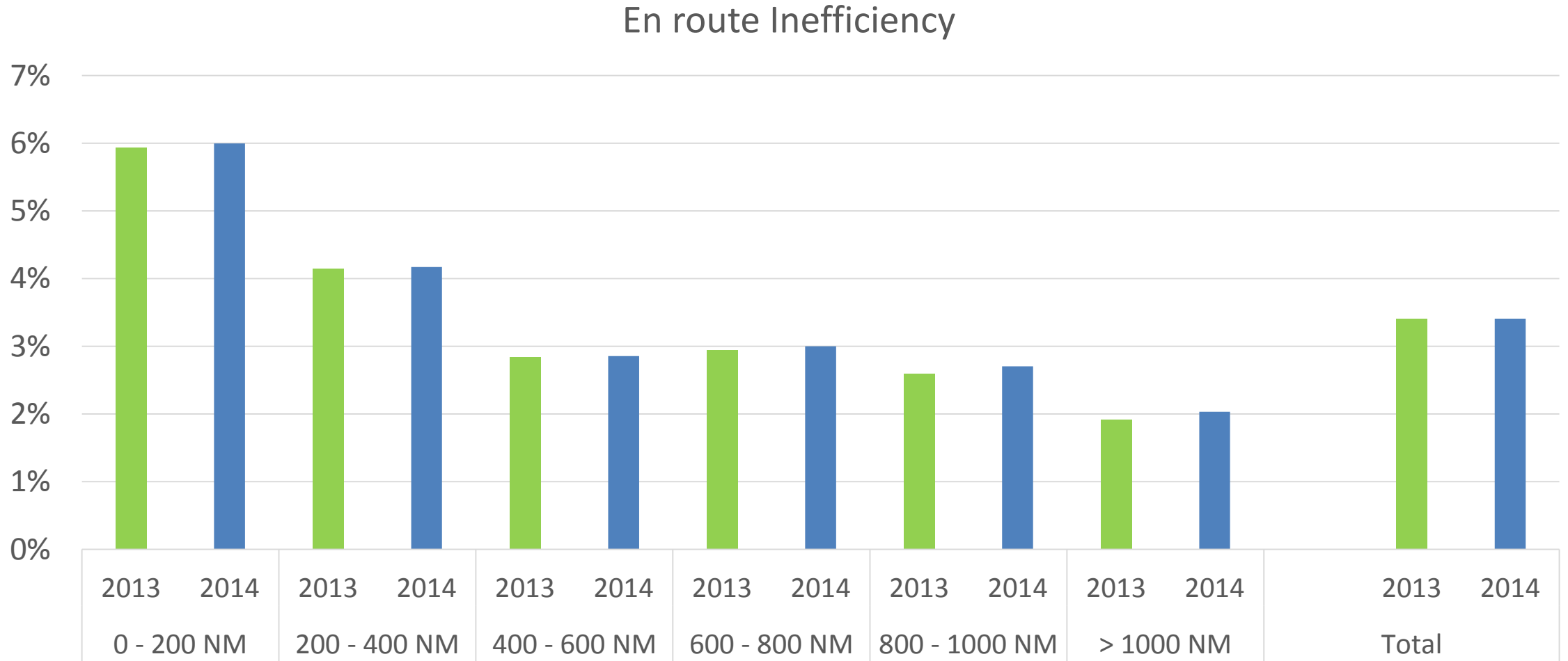
- Flight Event Data
 - From FAA Traffic Flow Management System (TFMS);
 - Flight level records: arrival/ departure airports, D40A100 Actual/Great Circle/ Achieved distances and etc.
- Flight Track Data
 - From FAA Traffic Flow Management System (TFMS);
 - Radar position data (4D trajectory): typically one-minute position updates from individual flights.

Summary Statistics

- Only focusing on 34 core airports scattering around US;
- Around 3 million flights per year in/out of core 34 airports, accounting for about 50% of total flights in/out of the US;

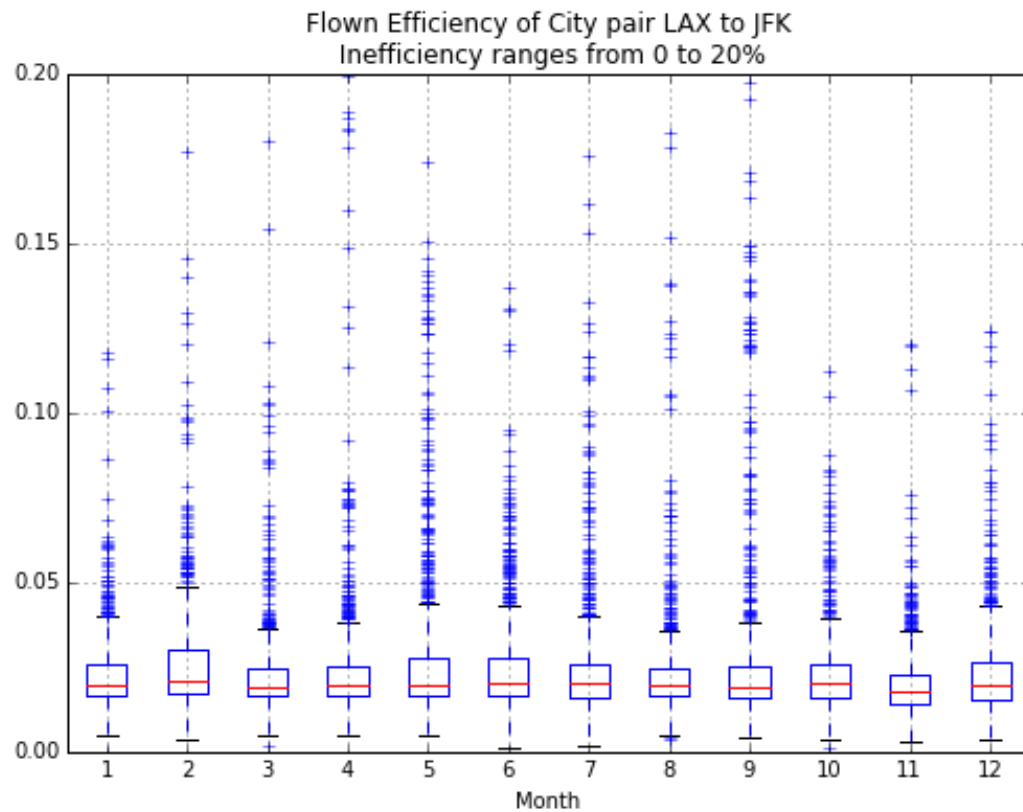


En Route Inefficiency vs Great Circle Distance

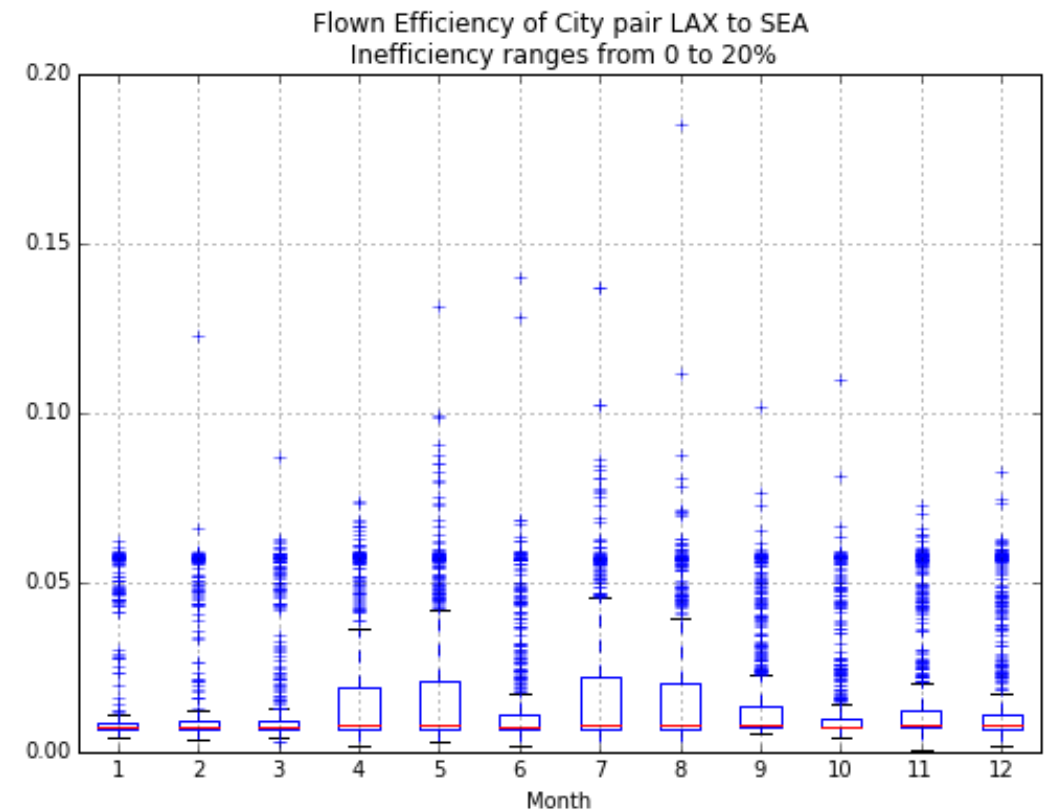


Inefficiency of Typical Airport Pairs (2013)

LAX to JFK (2.45%)

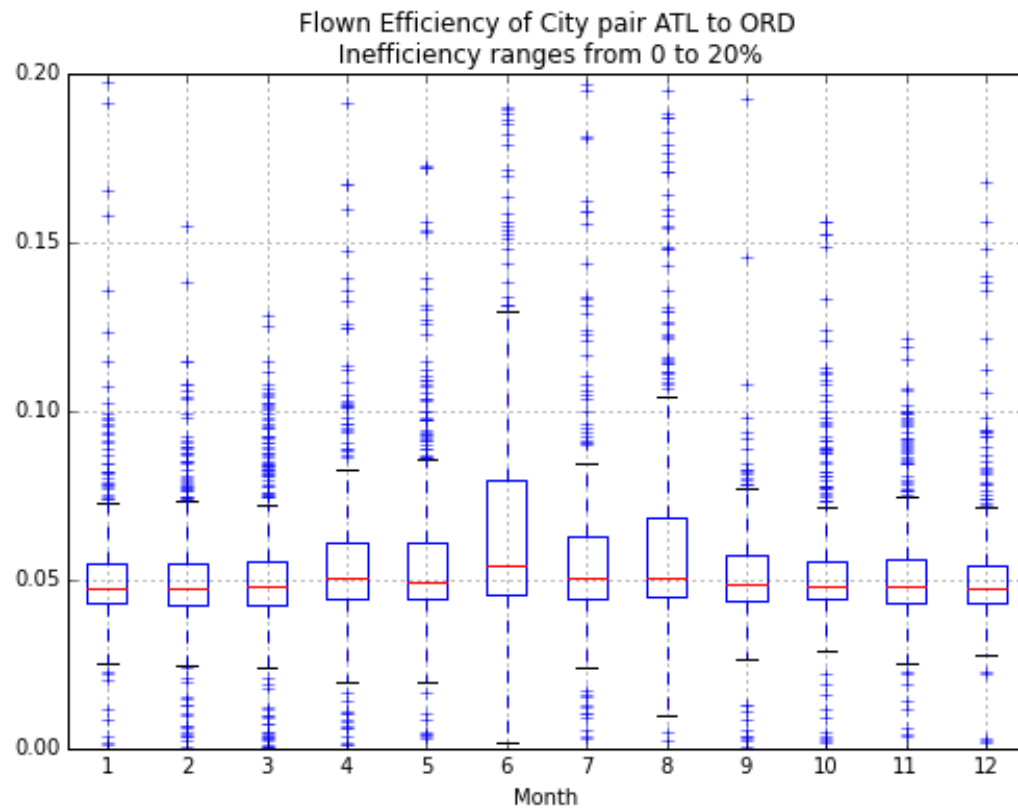


LAX to SEA (1.61%)

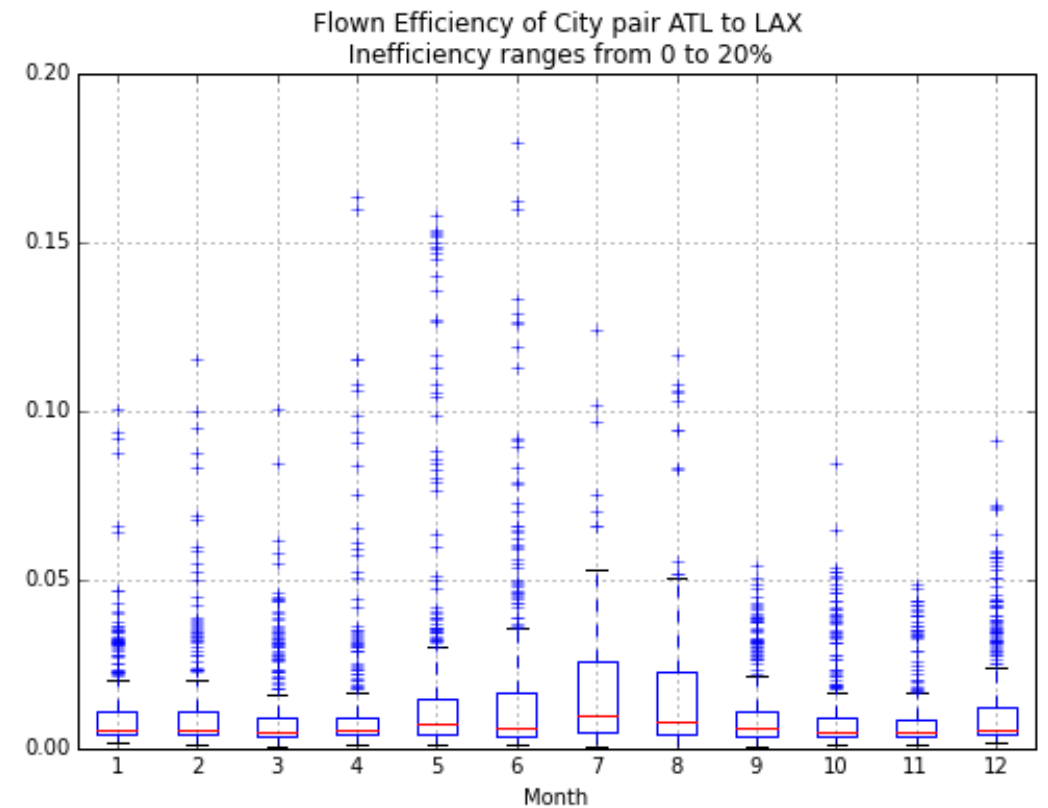


Inefficiency of Typical Airport Pairs (2013)

ATL to ORD (6.86%)



ATL to LAX (1.28%)



Outline

- Introduction
- Data Sources and Preliminary Statistical Analysis
- **Macroscopic Model**
- Microscopic Model
- Route Selection Model
- Conclusions

Method

- Use Regression analysis to explore the fixed effects of departure/ arrival airports, month and great circle distance;
- First to build a model that no interaction terms are included. (Model I, 6M observations, 82 Variables)

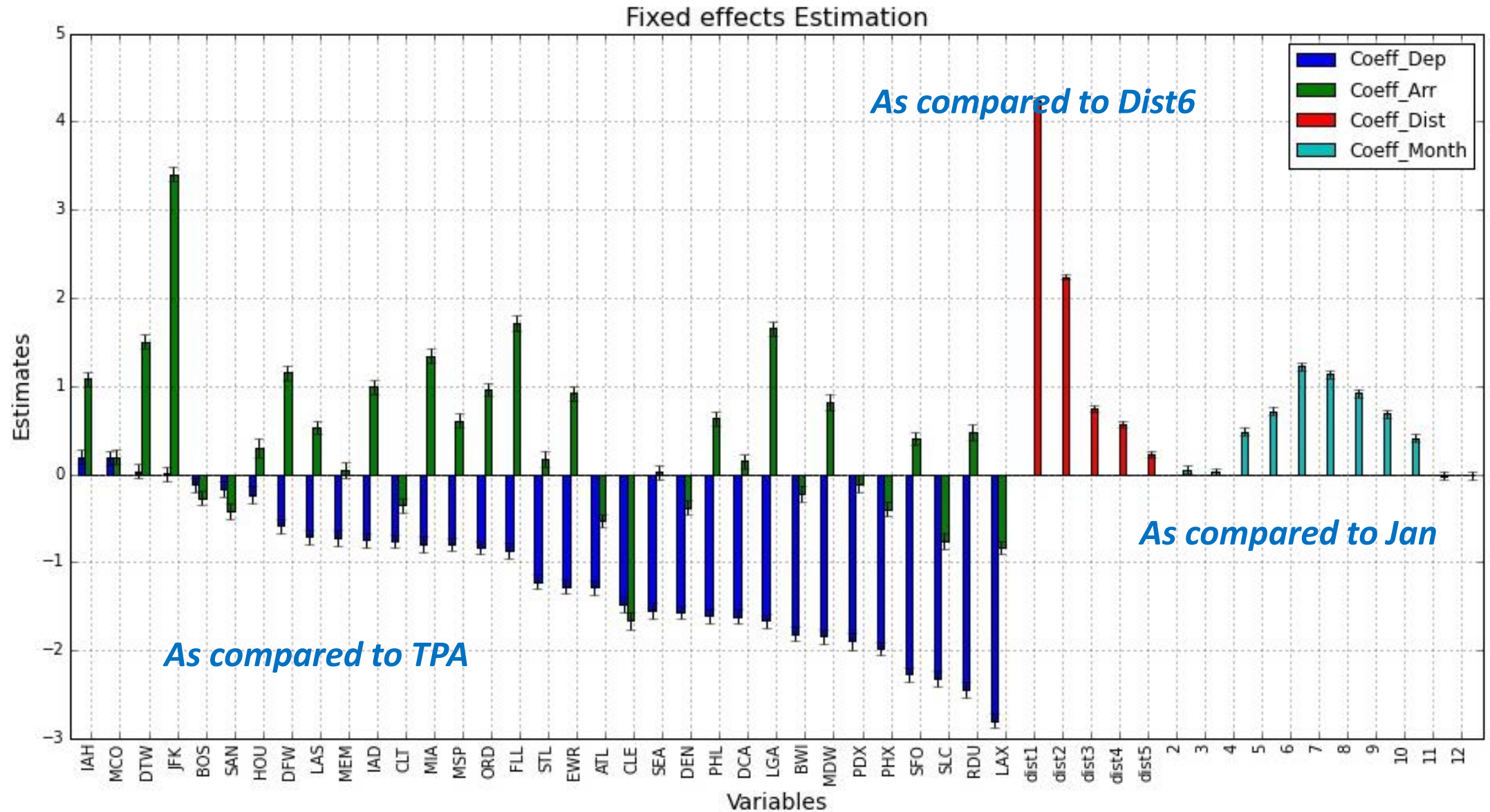
$$ineffi = \beta_{dep} \cdot X_{dep} + \beta_{arr} \cdot X_{arr} + \beta_{mon} \cdot X_{mon} + \sum_i^5 \beta_i \cdot Dist_i$$

- Then build a model that considers the interaction between airports and month. (Model II, 6M observations, 808 Variables)

$$ineffi = \beta'_{int1} \cdot X_{dep} \cdot X_{mon} + \beta'_{int2} \cdot X_{arr} \cdot X_{mon} + \sum_i^5 \beta_i \cdot Dist_i$$

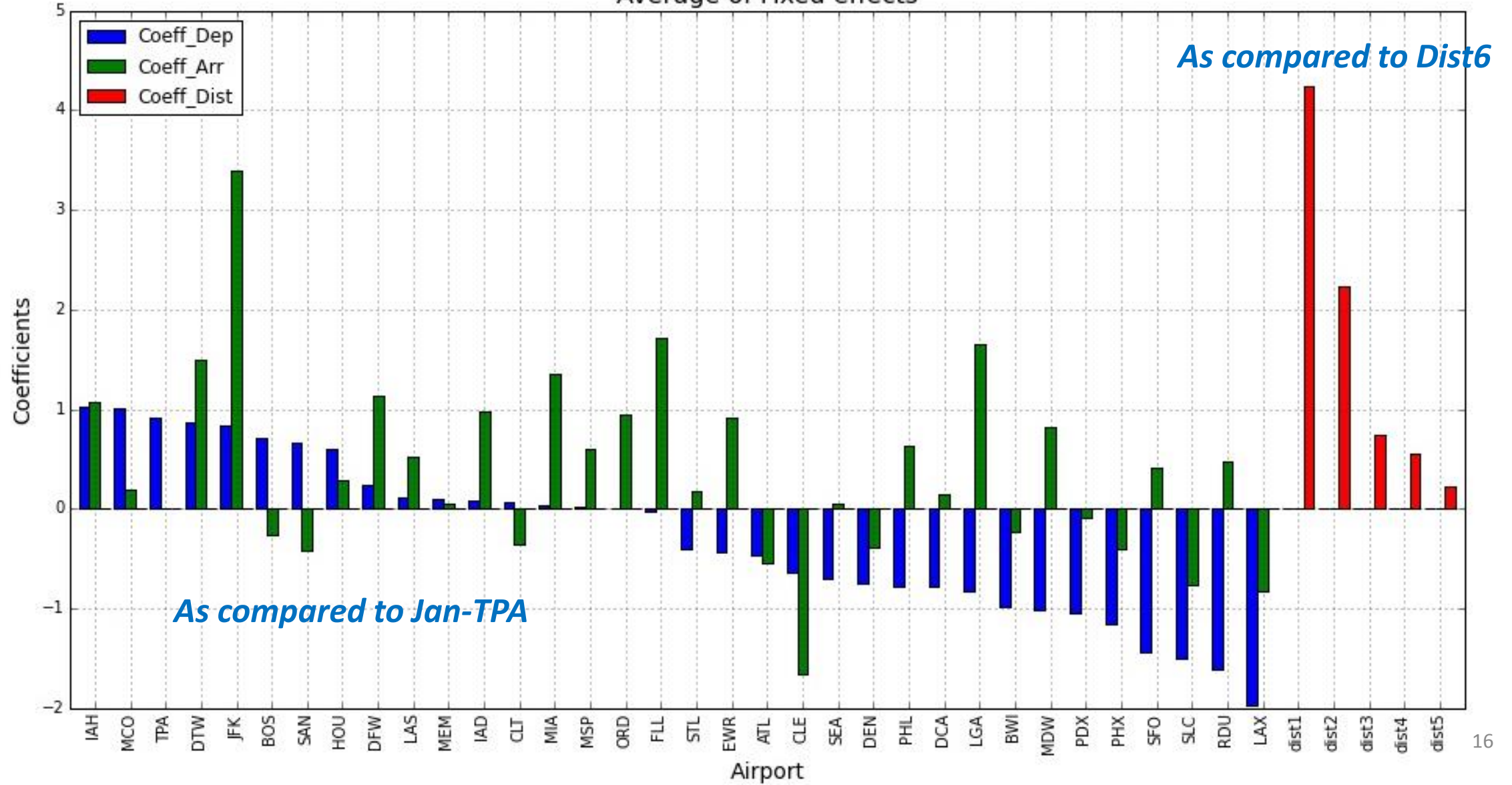
- $Dist_1: 0 - 200 \text{ NM}; Dist_2: 200 - 400 \text{ NM}; Dist_3: 400 - 600 \text{ NM};$
- $Dist_4: 600 - 800 \text{ NM}; Dist_5: 800 - 1000 \text{ NM}; Dist_6: > 1000 \text{ NM}$

Model I - Estimation

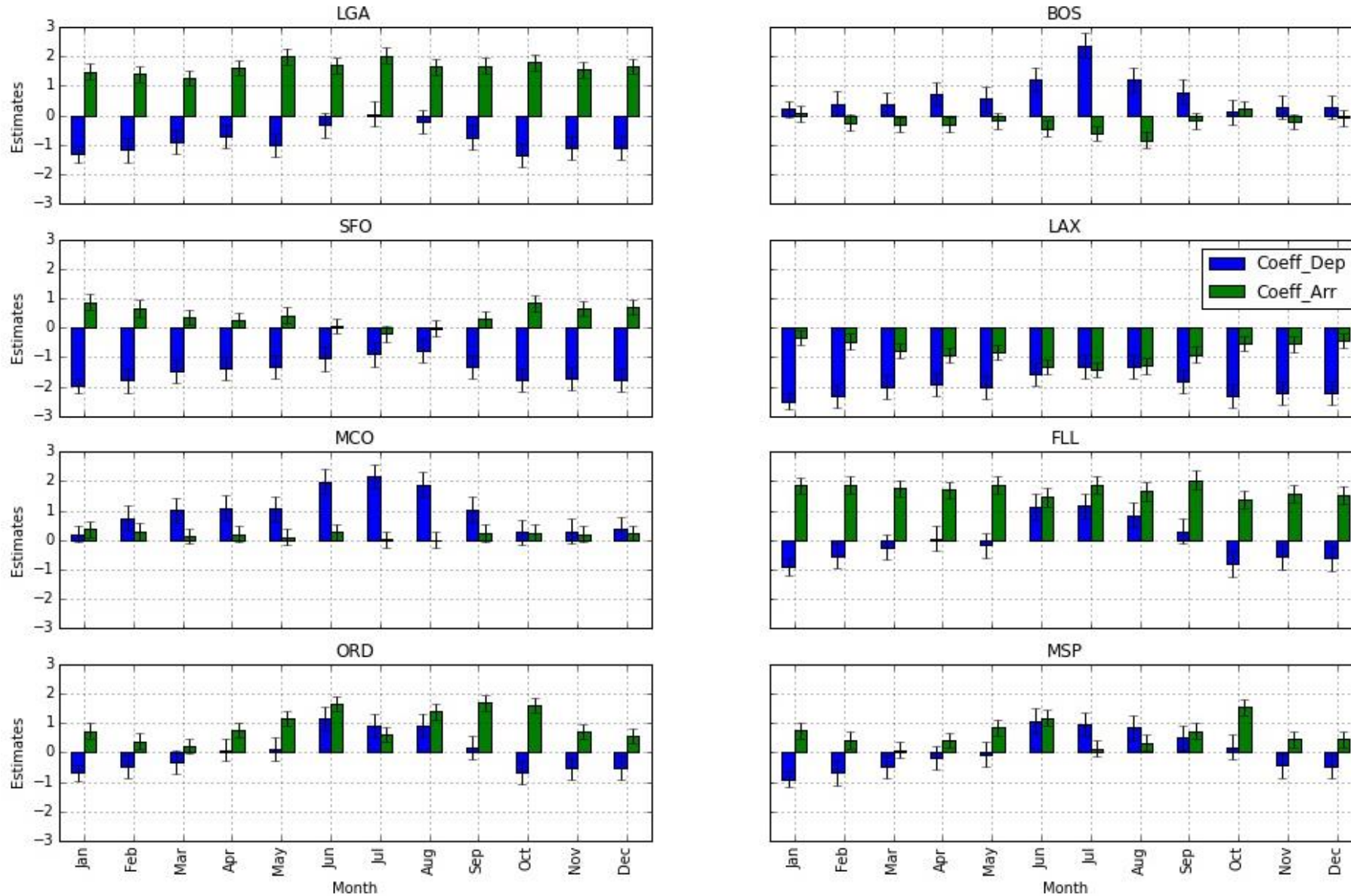


Model II - Estimation

Average of Fixed effects



Model II – comparison



Outline

- Introduction
- Data Sources and Preliminary Statistical Analysis
- Macroscopic Model
- **Microscopic Model**
- Route Selection Model
- Conclusions

Trajectory Clustering

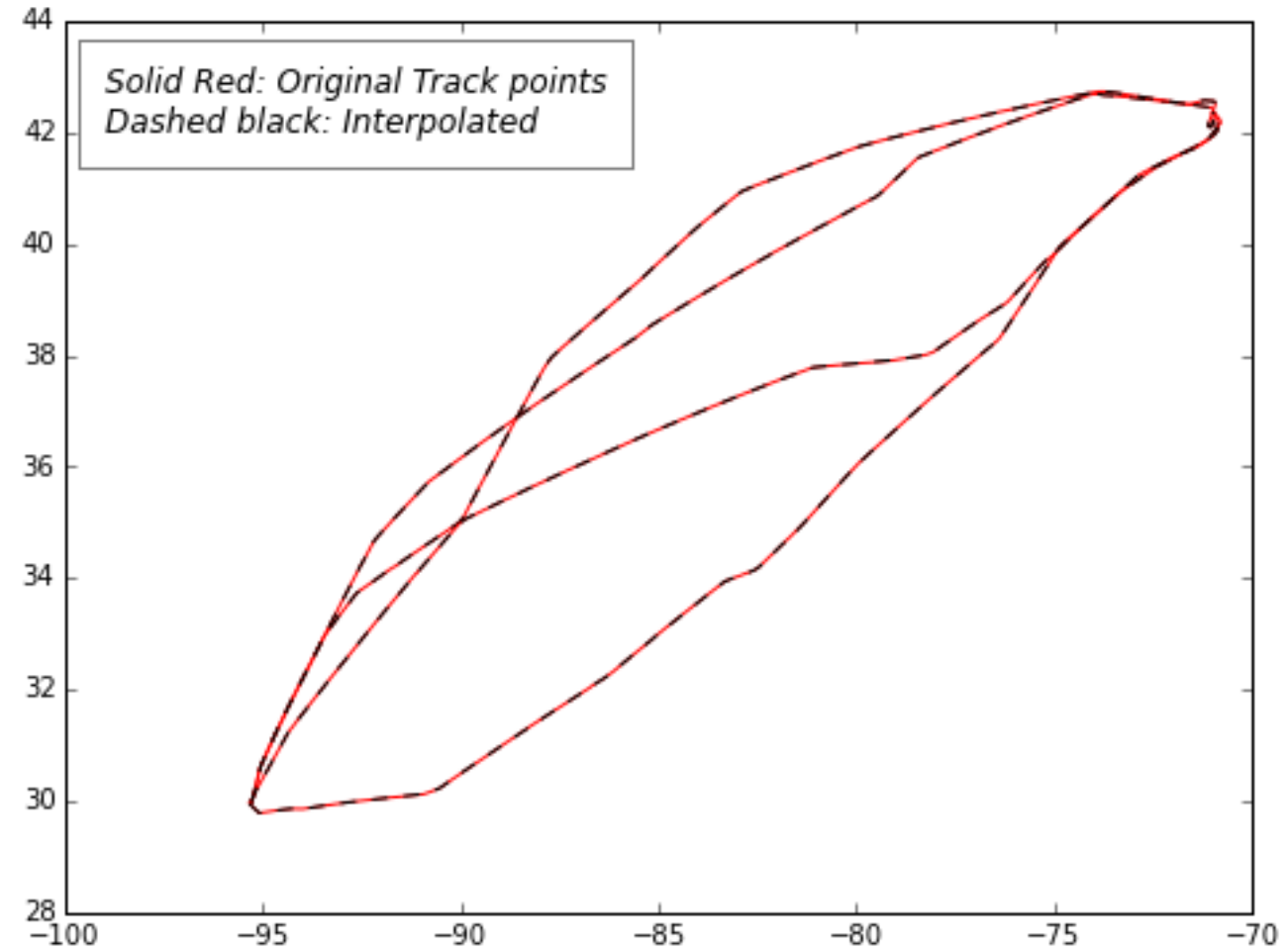
- A Trajectory is defined as the set of horizontal tracking points and times for a specific flight
- A flight typically has $\sim 10^2$ tracking points
- No two flights have exactly the same trajectory
- We want to find sets of flights with similar trajectories
- Standard clustering techniques need to be adapted because
 - Different trajectories have different numbers of tracking points
 - There is no inherent correspondence between the points of different trajectories
 - The number of tracking points is very large

Clustering Algorithms

- Step 0: Trajectory Cleaning
 - Exclude geometric-discontinuity and time-discontinuity trajectories;
 - Exclude trajectories starting/ending outside terminal areas.
- Step 1: Trajectory resampling
 - Get trajectories with equal numbers of points;
 - Linear Interpolation (with respect to tracking distance).
- Step 2: Principal Component Analysis (PCA)
 - Dimension reduction & Trajectory smoothing;
 - First five modes can capture more than 97% of variations.
- Step 3: Clustering
 - Trajectory classifications;
 - DBSCAN algorithm is applied to the PCA mode matrix to get representative clusters.

Resampling Example

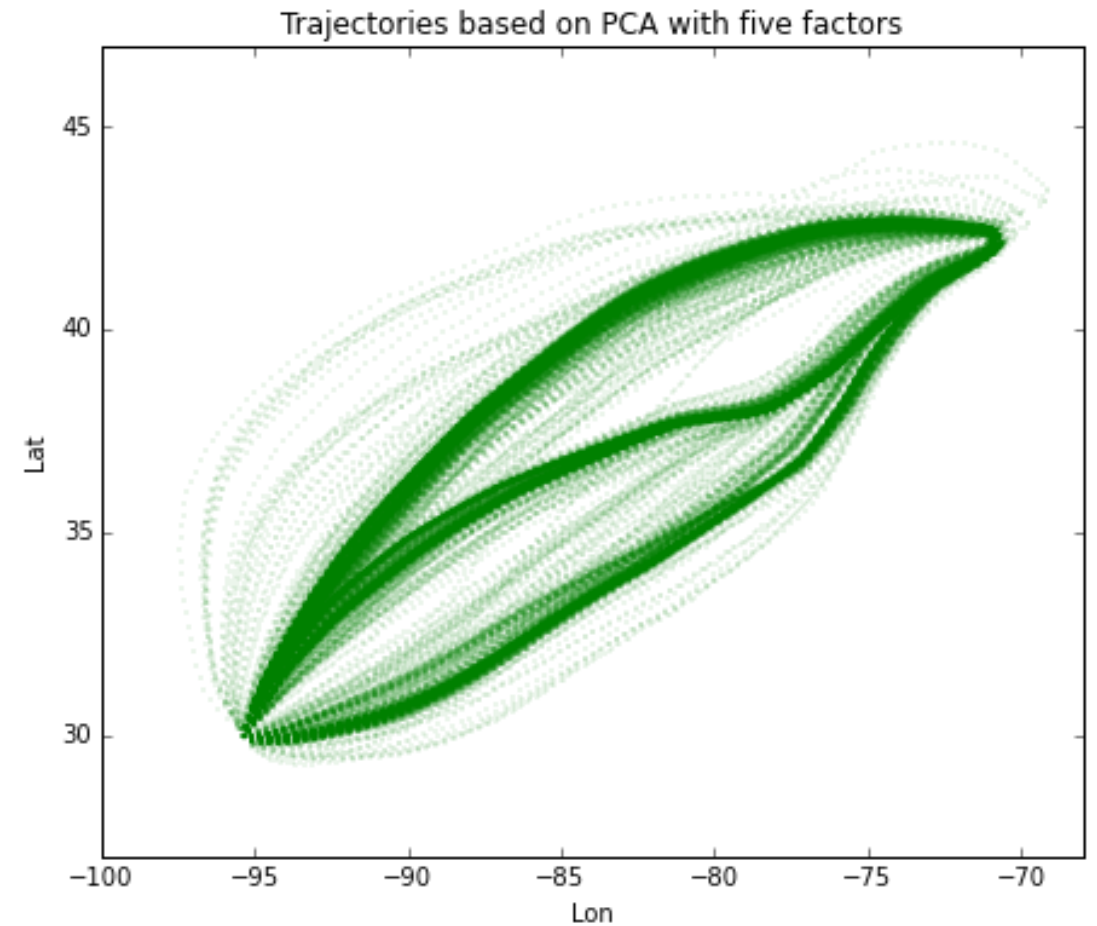
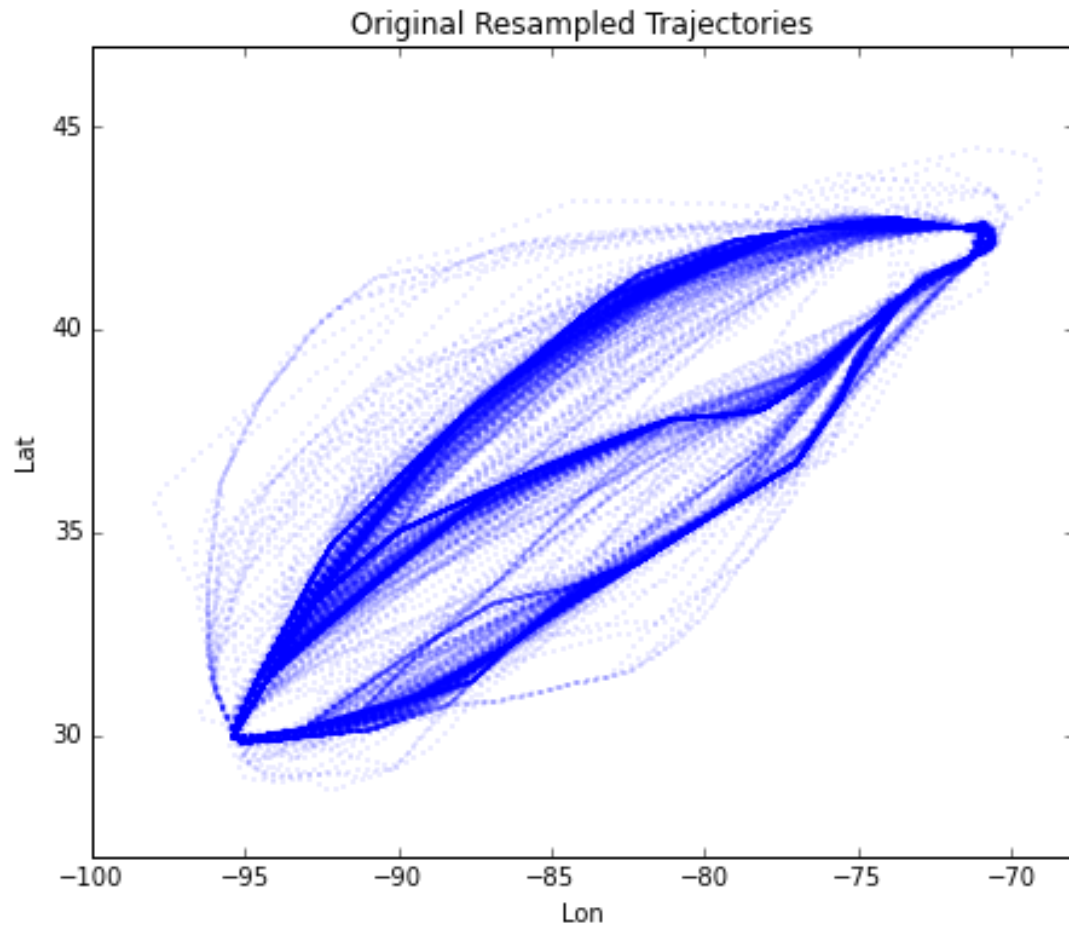
- Linear interpolation between the start and end tracking location for each route
- 100 pseudo points are predicted locations at:
 - Initial location (d_0)
 - $d_0 + \text{trajectory distance}/99$ (d_1)
 - $d_1 + \text{trajectory distance}/99$ (d_2)
 - ...
 - Final trajectory location (d_{100})



Dimension Reduction

- Reduce the dimension of trajectories – save computational time
- Improve the quality of clustering – Principal Component Analysis (PCA) can help to filter off noise and smooth the data
- Using PCA, we found that five factors can capture almost all the variation in the original 900 variables e.g.
 - 99% for IAH → BOS
 - 96% for FLL → JFK
 - 94% for ORD → DCA

Example of Dimension Reduction (IAH-BOS)



Clustering

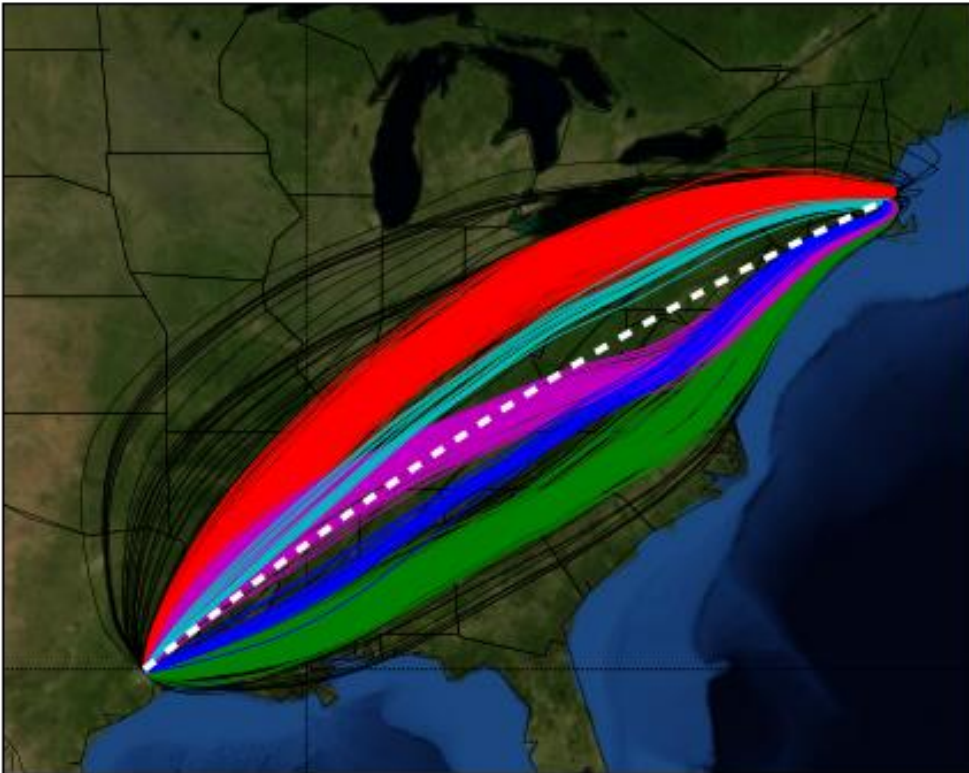
- Use trajectory factor scores to find sets of trajectories that are similar to each other
- Measure similarity between any two trajectories using Euclidean distance between their respective factor scores
- Tried two clustering methods
 - K means
 - DBSCAN
- Prefer DBSCAN because
 - Do not need to pre-determine number of clusters
 - Allows trajectories to be identified as outliers
 - Can limit variation within each cluster
 - However DBSCAN does require two parameters to be specified

DBSCAN Results for Six Airport Pairs

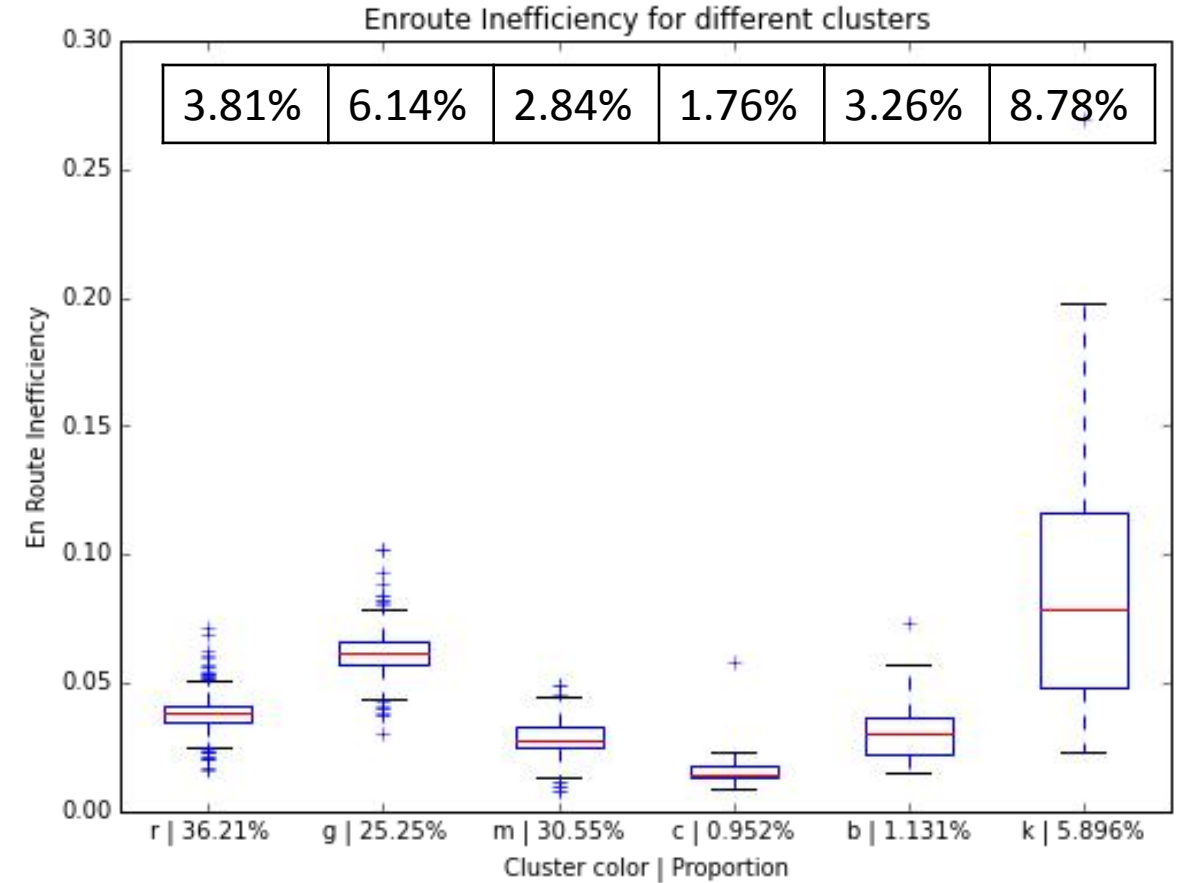
- Input parameters (distance threshold and minimum points) determined using trial-and-error for each pair

IAH \rightarrow BOS (1679 of original 1817)

DBSCAN applied to PCA mode matrix

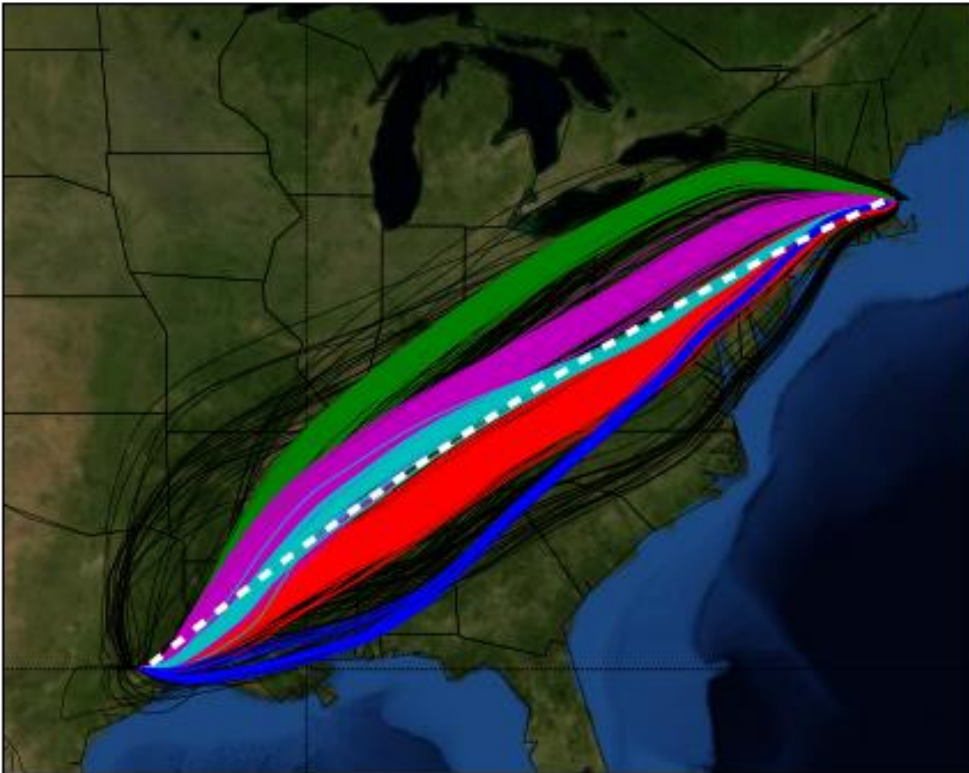


Black lines are classified as outliers

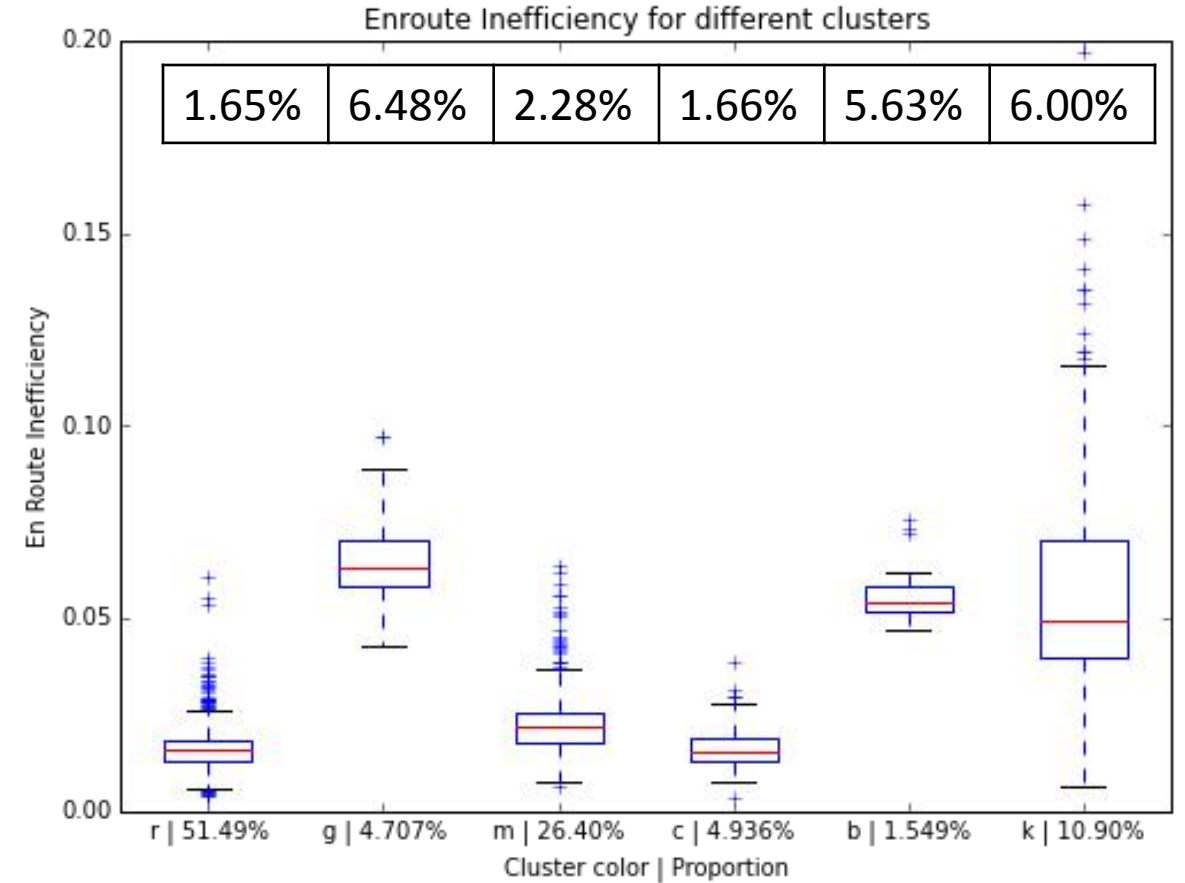


BOS → IAH (1742 of original 1883)

DBSCAN applied to PCA mode matrix

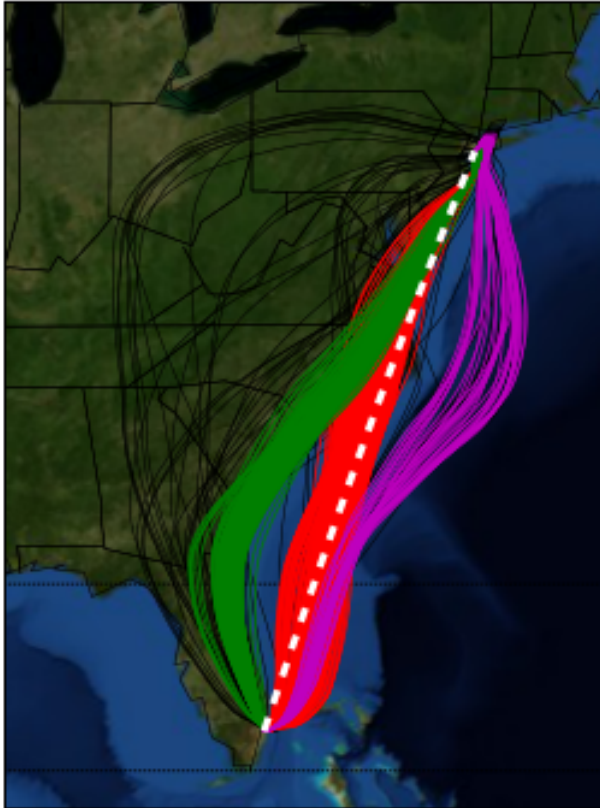


Black lines are classified as outliers

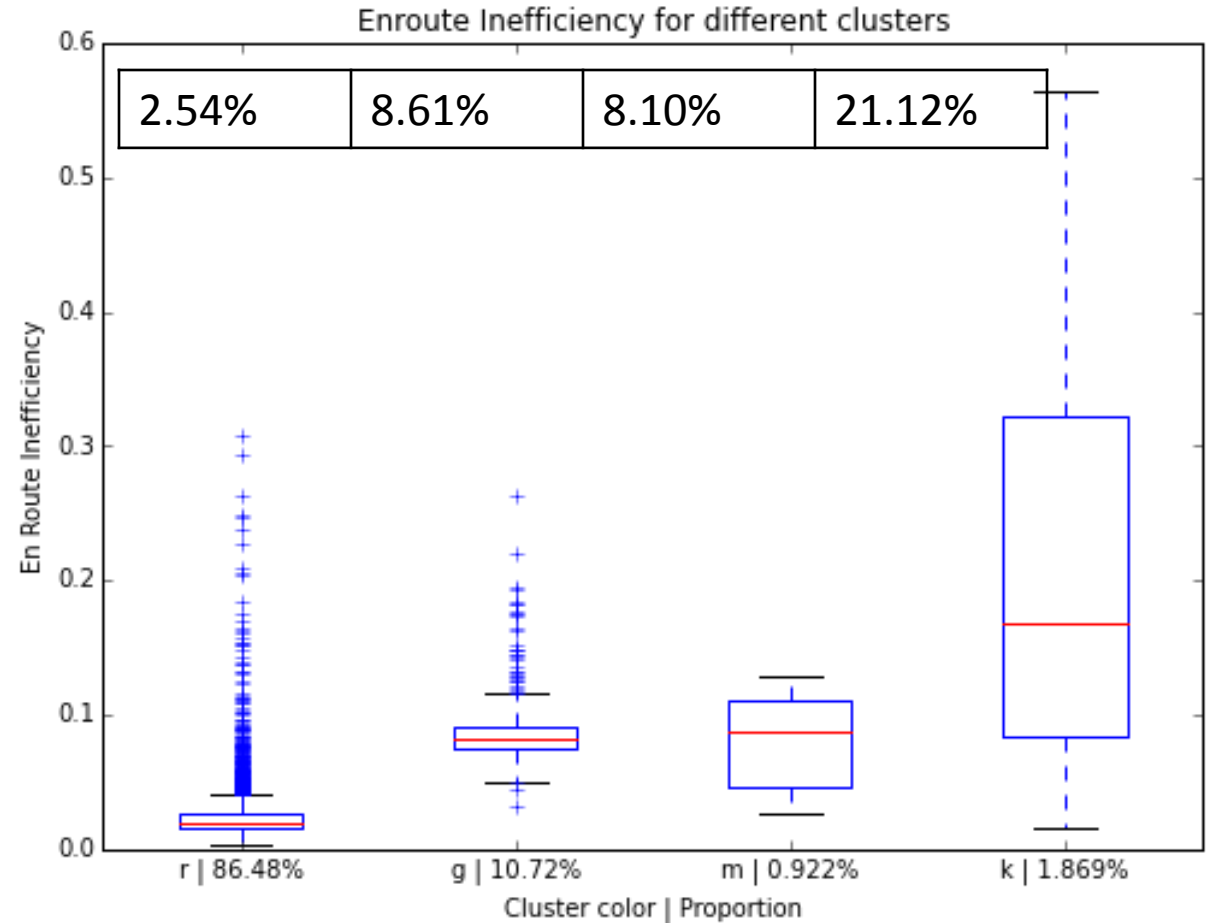


FLL → JFK (4011 of original 4267)

DBSCAN applied to PCA mode matrix

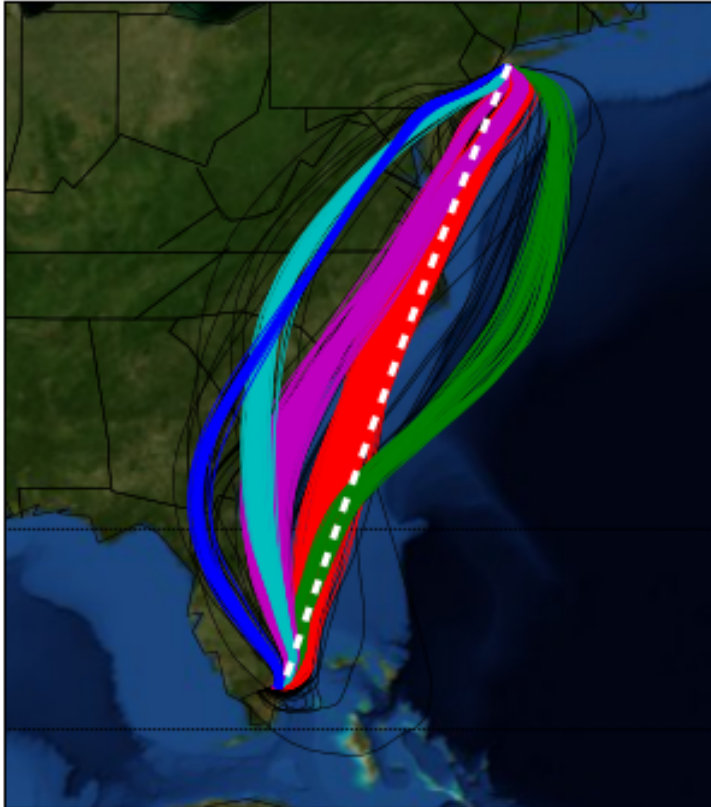


Black lines are classified as outliers

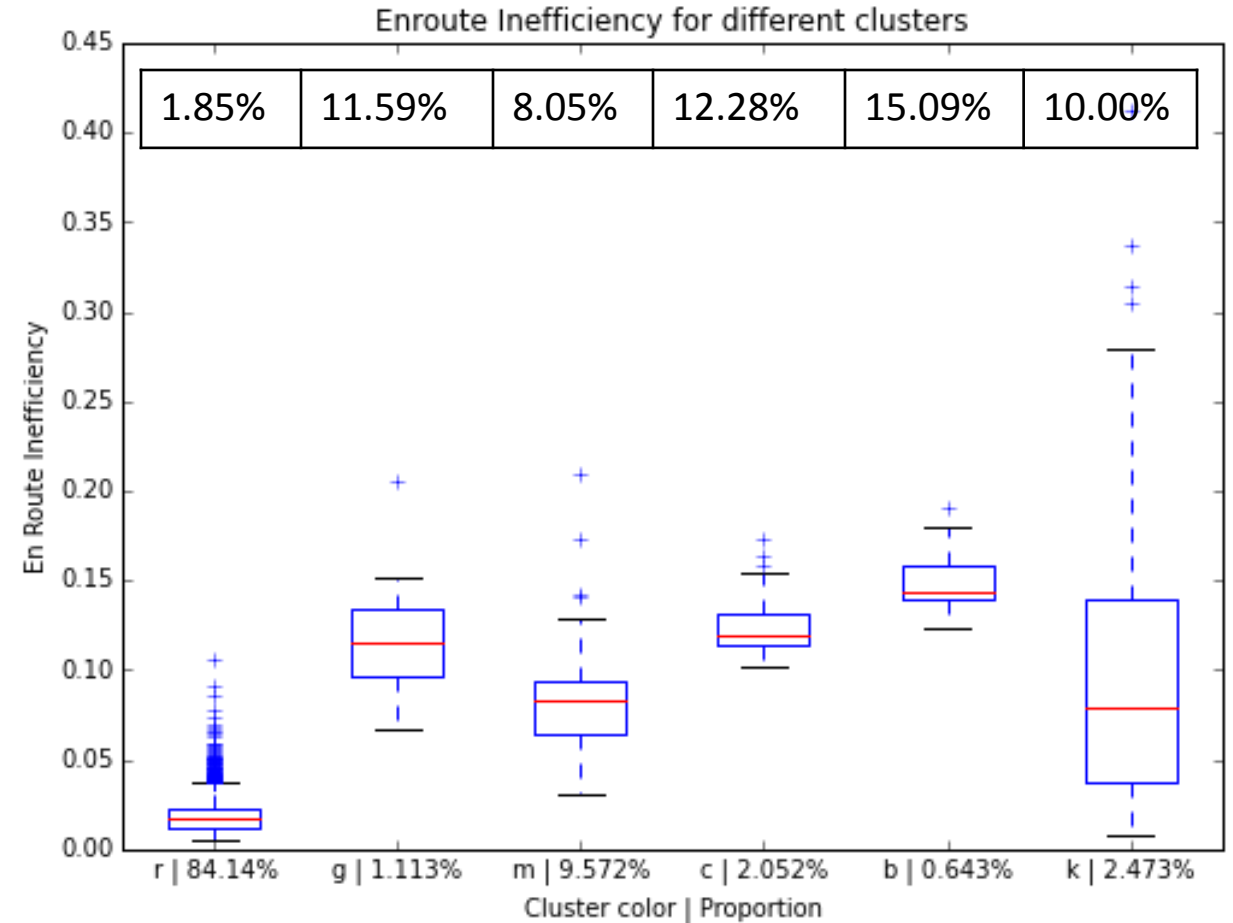


JFK → FLL (4043 of original 4273)

DBSCAN applied to PCA mode matrix

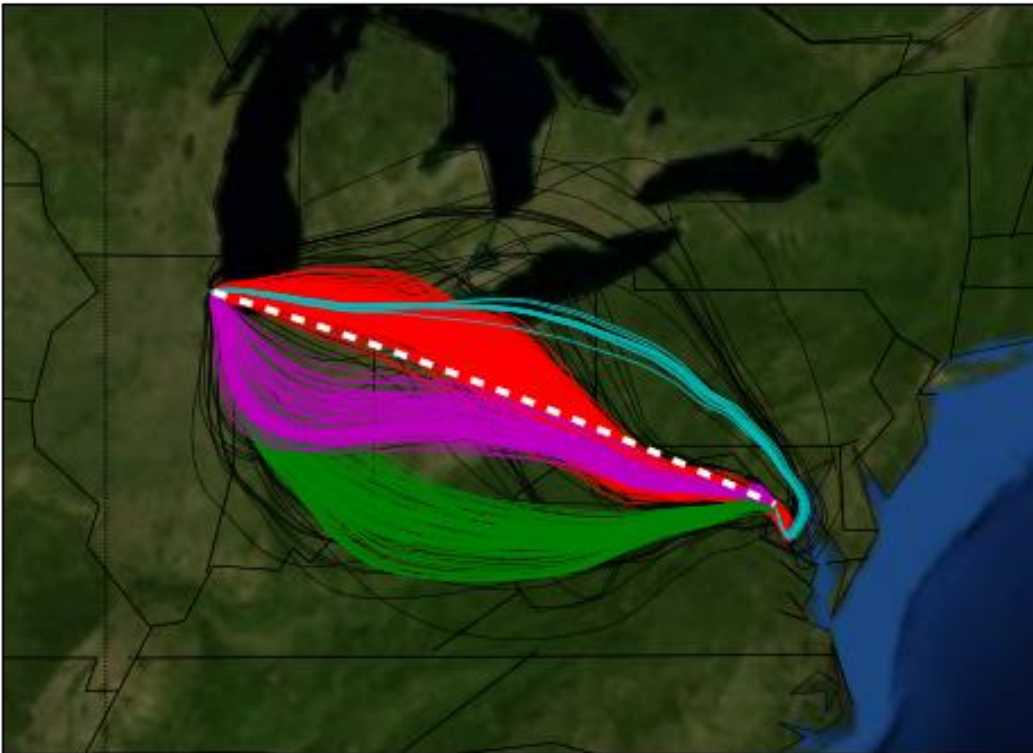


Black lines are classified as outliers

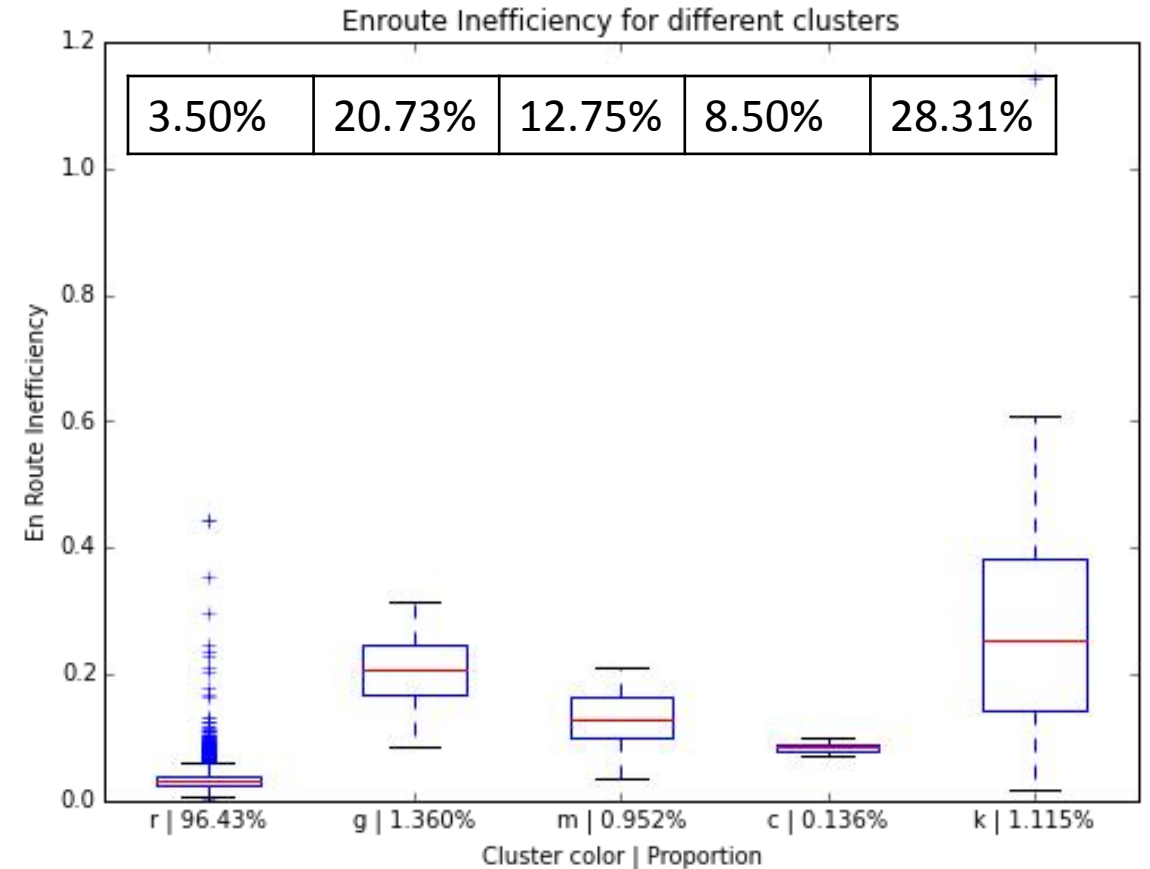


ORD \rightarrow DCA (7349 of original 7574)

DBSCAN applied to PCA mode matrix

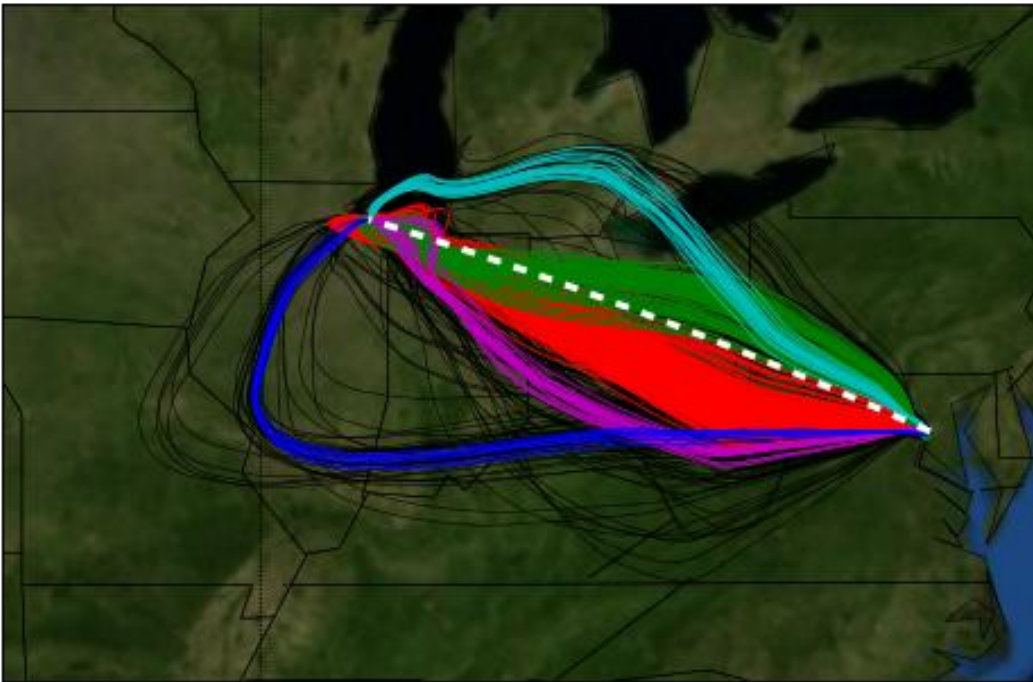


Black lines are classified as outliers

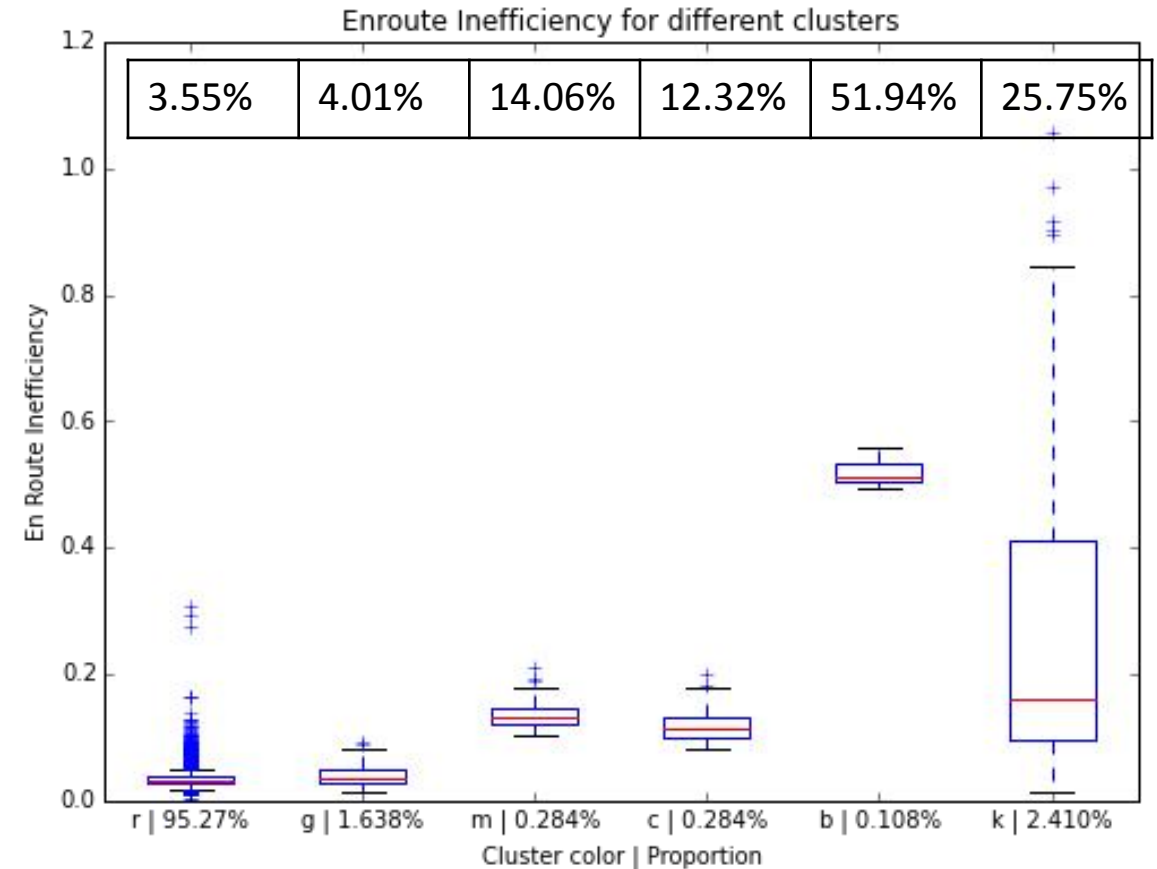


DCA → ORD (7383 of original 7557)

DBSCAN applied to PCA mode matrix



Black lines are classified as outliers



Impact of Cluster Membership on En Route Efficiency

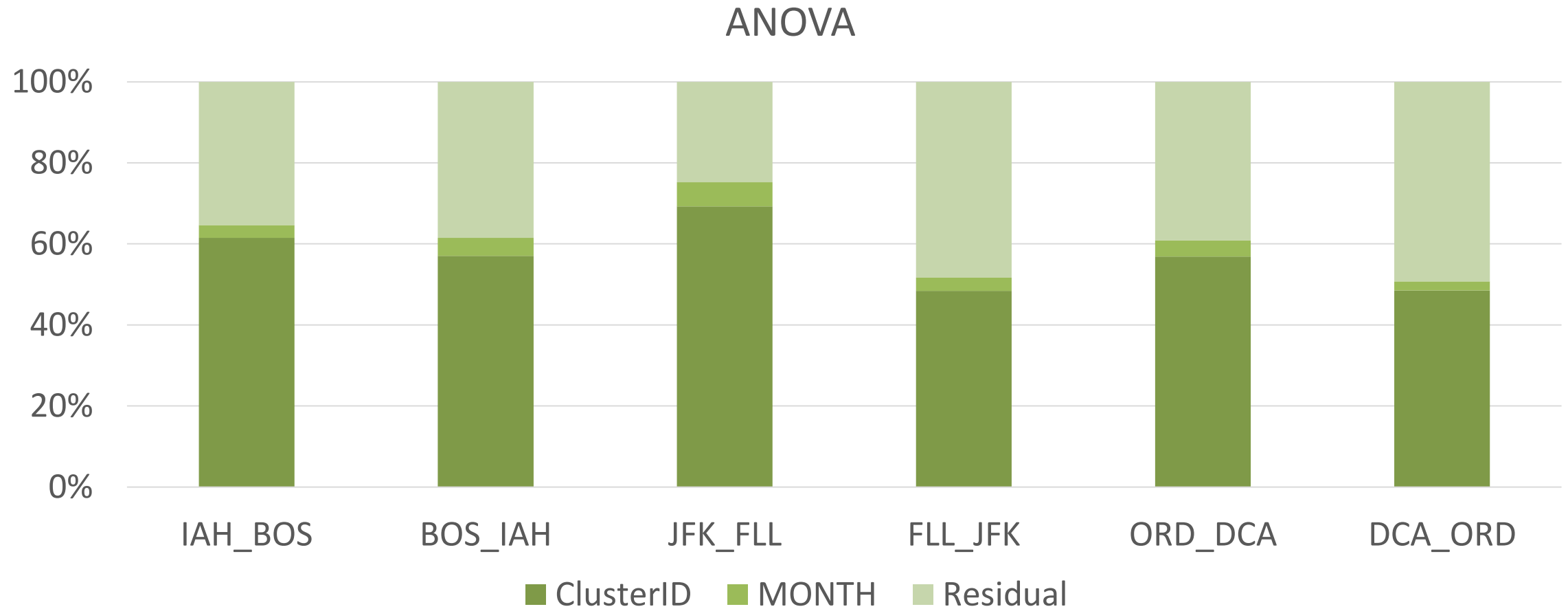
- Build route-specific fixed effect models to capture variations in en route inefficiencies among representative clusters;
- Model specification
 - Separate models for each airport pair
 - $Inefficiency(\%) = \beta_0 + \beta'_1 \cdot X_{month} + \beta'_2 \cdot X_{ClusterID}$;
 - X_{month} and $X_{ClusterID}$ are categorical variables;
 - Cluster ID can be found in previous slides.

Estimation Results



	IAH_BOS	BOS_IAH	JKF_FLL	FLL_JFK	ORD_DCA	DCA_ORD
Intercept	6.019***	8.452***	9.768***	20.924***	27.959***	25.485***
MONTH -2	-0.148	0.023	0.220*	0.042	0.009	-0.096
MONTH -3	-0.266*	-0.191	0.193*	0.317	-0.006	0.165
MONTH -4	0.121	0.312*	0.636***	0.259	0.045	-0.008
MONTH -5	0.184	-0.023	0.112	0.131	0.083	-0.114
MONTH -6	0.105	0.853***	0.367***	0.096	0.688***	0.621***
MONTH -7	0.025	0.215	0.336***	0.393*	0.344**	0.108
MONTH -8	-0.166	0.568***	0.105**	0.374*	0.151	0.964***
MONTH -9	0.093	0.306*	0.083	-0.073	-0.051	0.367*
MONTH -10	-0.090	0.018	0.012	-0.273	-0.109	-0.377*
MONTH -11	-0.074	-0.029	0.223*	-0.020	-0.116	-0.342
MONTH -12	-0.261	-0.071	0.130	-0.228	-0.041	-0.179
Cluster ID – r	-4.328***	-4.831***	-8.108***	-18.470***	-24.538***	-22.026***
Cluster ID – b	0.525***	-2.463***	1.558***	-12.457***	-7.521***	-21.908***
Cluster ID – g	-3.697***	-5.801***	-1.970***	-12.956***	-15.434***	-11.593***
Cluster ID – m	-4.292***	-7.017***	2.240***	-	-19.705***	-13.695***
Cluster ID - c	-0.409	-5.498***	5.058***	-	-	26.114***
R squared	0.6463	0.6147	0.7523	0.5167	0.6083	0.5076

Analysis of Variance



Outline

- Introduction
- Data Sources and Preliminary Statistical Analysis
- Macroscopic Model
- Microscopic Model
- **Route Selection Model**
- Conclusions

Route Selection Model

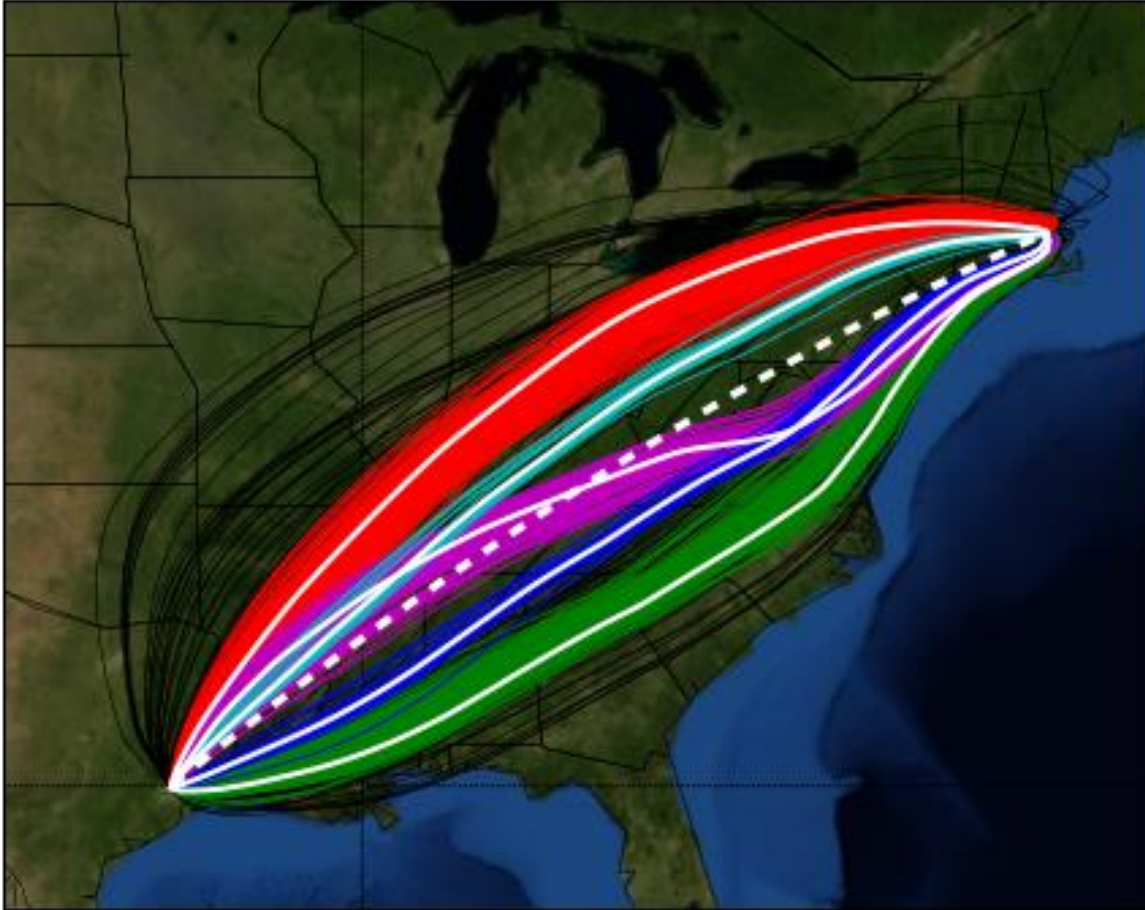
- We try to understand what factors may influence the route selection;
- We are specifically interested in the impacts of convective weather such as thunder storm, shower and etc.

Methodologies

- *Step 0: Trajectory Clustering*
 - Apply trajectory clustering algorithm to classify trajectories into groups;
 - Typically four to five groups can be classified
- *Step 1: “Center” of Clusters*
 - Apply 1-Median algorithm to determine the center for each cluster
- *Step 2: Data Fusion*
 - Obtain (historical) real time weather data for the same time period with trajectory data;
 - Merge the real-time weather data with 4D trajectory data (time series);
 - For each trajectory, determine the aggregated level of weather intensity
- *Step 3: Route Selection Model*
 - Estimate the impact of weather factors and other interesting factors;
 - Multinomial logit model is applied to the final dataset.

Trajectory Clustering and Center

DBSCAN applied to PCA mode matrix

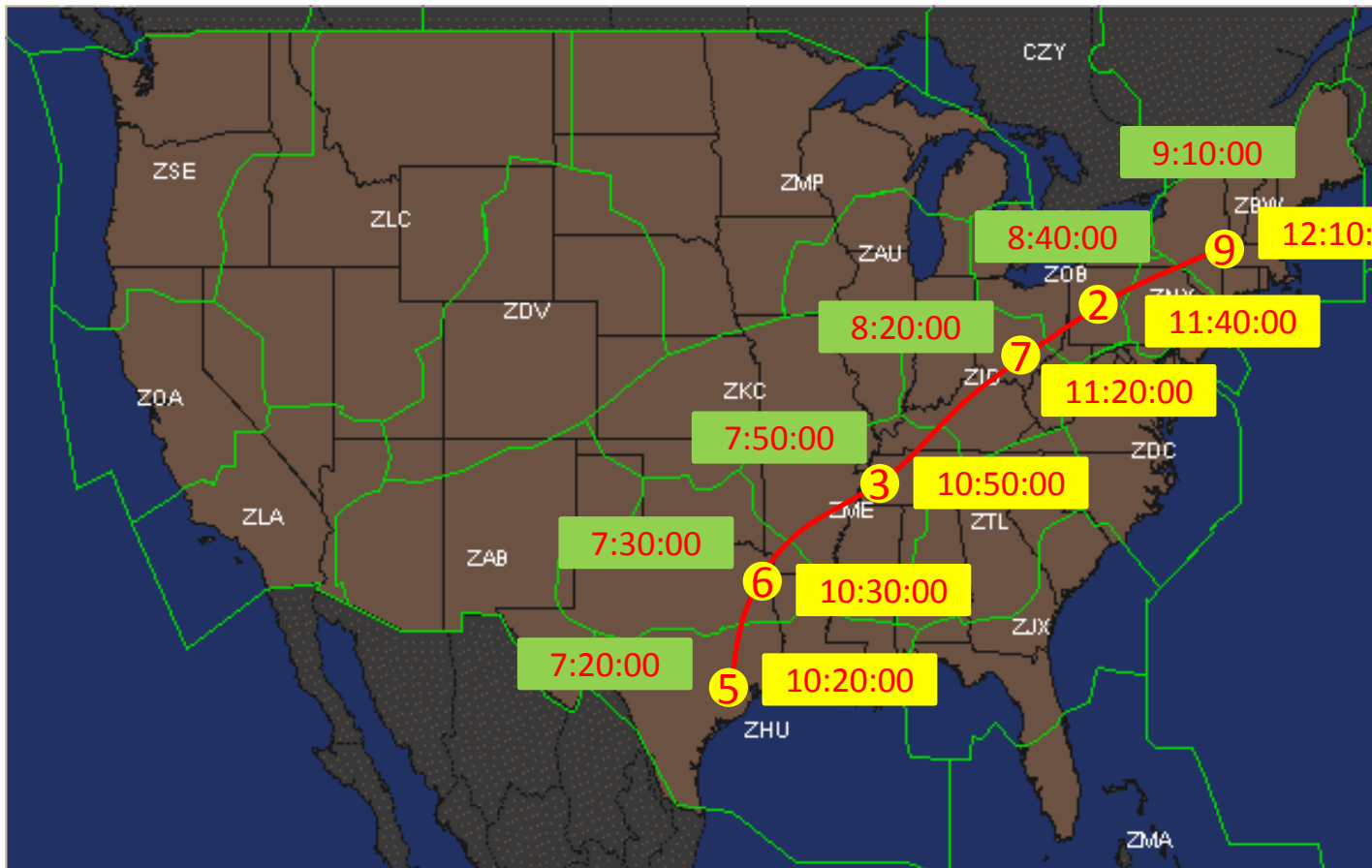


- The case for IAH → BOS in 2013;
- Solid white lines are the centers for cluster groups;
- Five representative clusters are identified (red, blue, green, magenta and cyan), with only 6% of outliers (black).
- 4d information (lat, lon, alt, time) is extracted for the five center trajectories.

Data Fusion

- Hourly weather data
 - Obtained from NOAA;
 - Include hourly (usually less than 30 minute) data indicating whether a station had thunder storms, showers, rains, hails... during a time period;
 - There are more than 2000 weather stations around United States;
- Center-trajectory data
 - “center” is determined by clustering algorithm;
 - 4d information is extracted for each center
- Data Fusion
 - First to exclude trajectories that are classified as outliers, then for each member trajectories, we extract the time stamp of its first track point;
 - Then we replace the departure time of center-trajectories with such time stamp, and calculate the time stamps for the rest of track point for center-trajectories;
 - For each track point of center-trajectories, we calculate the total number of convective weather happened within the ARTCC zone associated with the tracking location in a time interval that is 2 hours prior and after the tracking time stamp;
 - Finally we aggregate the total number of convective weather appeared during the whole trajectory.

Example for Data Fusion



- Red line: one of the center trajectory;
- Yellow dots: track points;
- Yellow boxes: original time stamp for track points along center trajectory;
- Green boxes: time stamp for track points after replacing the departure time;
- Numbers in yellow dots: number of convective weather (e.g. thunder storm) appeared for the ARTCC associated with the track point in a 2-hour time interval;
- Total convective weather intensity:
 - $5+6+3+7+2+9 = 32$

Model Specification

- For the case of IAH \rightarrow BOS, we have five clusters, including 1579 trajectory members;
- Model Specification

$$V_0 = \beta_0 + \beta_1 \cdot ThunderStormLevel_0 + \beta_2 \cdot Shower_0 + \beta_3 \cdot Squall_0 + \beta' \cdot Month + \gamma' \cdot Morning Flight$$

$$V_1 = \beta_1 + \beta_1 \cdot ThunderStormLevel_1 + \beta_2 \cdot Shower_1 + \beta_3 \cdot Squall_1$$

$$V_2 = \beta_2 + \beta_1 \cdot ThunderStormLevel_2 + \beta_2 \cdot Shower_2 + \beta_3 \cdot Squall_2$$

$$V_3 = \beta_3 + \beta_1 \cdot ThunderStormLevel_3 + \beta_2 \cdot Shower_3 + \beta_3 \cdot Squall_3$$

$$V_4 = \beta_1 \cdot ThunderStormLevel_4 + \beta_2 \cdot Shower_4 + \beta_3 \cdot Squall_4$$

- Independent variable
 - Thunder Storm Level: the aggregated thunder storm level for one nominal trajectory;
 - Shower: the aggregated number of total shower appearance for one nominal trajectory;
 - Squall: the aggregated squall ratio for one nominal trajectory;
 - Month: monthly fixed effects, 11 dummy variables;
 - Morning flight: dummy variable, equals to 1 if the local departure time is before 12 pm.

Estimation Results

Multinomial Logit Model Regression Results

```

=====
Dep. Variable:                CHOICE    No. Observations:                1,579
Model:                Multinomial Logit Model    Df Residuals:                1,560
Method:                MLE    Df Model:                19
Date:                Fri, 15 Jul 2016    Pseudo R-squ.:                0.281
Time:                11:53:10    Pseudo R-bar-squ.:                0.273
converged:                True    Log-Likelihood:                -1,827.522
                                LL-Null:                -2,541.302
=====

```

	coef	std err	z	P> z	[95.0% Conf. Int.]	
ASC_R0	3.8766	0.287	13.501	0.000	3.314	4.439
ASC_R1	3.0809	0.235	13.135	0.000	2.621	3.541
ASC_R2	3.2143	0.234	13.729	0.000	2.755	3.673
ASC_R3	-0.2599	0.340	-0.764	0.445	-0.926	0.406
Month 1 - fixed effects	-0.0957	0.267	-0.359	0.720	-0.618	0.427
Month 2 - fixed effects	-0.6622	0.281	-2.356	0.018	-1.213	-0.111
Month 3 - fixed effects	-1.2952	0.282	-4.591	0.000	-1.848	-0.742
Month 4 - fixed effects	-0.4955	0.257	-1.925	0.054	-1.000	0.009
Month 5 - fixed effects	-0.0165	0.248	-0.067	0.947	-0.502	0.469
Month 6 - fixed effects	-0.5894	0.252	-2.338	0.019	-1.083	-0.095
Month 7 - fixed effects	-0.2547	0.241	-1.059	0.290	-0.726	0.217
Month 8 - fixed effects	-0.1762	0.238	-0.742	0.458	-0.642	0.289
Month 9 - fixed effects	-0.6282	0.248	-2.535	0.011	-1.114	-0.143
Month 10 - fixed effects	-0.5649	0.243	-2.323	0.020	-1.042	-0.088
Month 11 - fixed effects	-0.0949	0.242	-0.393	0.695	-0.568	0.379
Morning Flight	-0.3495	0.111	-3.155	0.002	-0.567	-0.132
Thunder Storm	-0.0662	0.015	-4.327	0.000	-0.096	-0.036
Shower	-1.7684	0.807	-2.193	0.028	-3.349	-0.188
Squall	-0.2585	0.133	-1.943	0.052	-0.519	0.002

- Thunder Storm Level, Shower and Squall are continuous variables, all have negative sign and are significant @10% level;
- Comparing to December, it seems that summer seasons have stronger negative correlations with route selection.

Outline

- Introduction
- Data Sources and Preliminary Statistical Analysis
- Macroscopic Model
- Microscopic Model
- Route Selection Model
- **Conclusions**

Conclusions

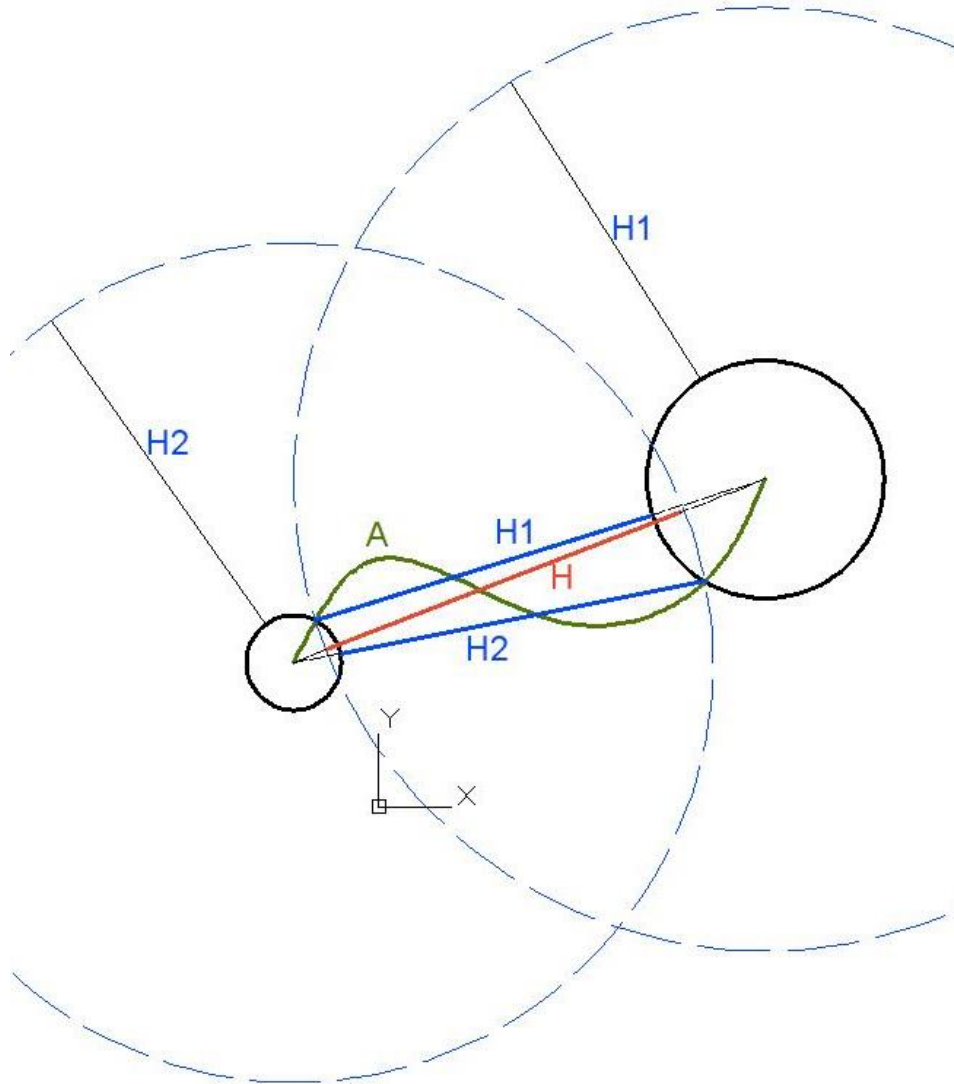
- Flight en route inefficiency is on average 3.5%, but varies significantly with airport pairs and seasons;
- Long-haul flights tend to be more efficient than short-haul flights;
- Flights depart from/ arrive at airports in the east coast tend to be more inefficient than the others;
- Flights in summer seasons are more inefficient;

Conclusions

- For most airport pairs, individual flight trajectories, while unique, can be divided into natural clusters whose members are very similar to one another;
- “Outlier” trajectories not belonging to a cluster account for from 1-10% of the total, depending on the airport pair;
- Cluster membership accounts for about 60% of overall variation in inefficiency
- Flights in summer seasons (May to August) are in general more inefficient than the others, but seasonal variation accounts for only 2-6% of the variation

Backup Slides

Method – on “Achieved distance”



- $H = \frac{H_1 + H_2}{2}$
- Indicate how much closer is the **Entry point** to destination and how much further is the **Exit point** away from origin.

Gap Between Actual and Flight Plan Distance

