

# NIST IAD DSE Pilot Evaluation Plan

Version 1.1, updated on Monday, April 4th 2016

## CONTENTS

<b>I</b>	<b>Introduction</b>	1
<b>II</b>	<b>Data</b>	1
II-A	Core sets . . . . .	1
II-B	Data access . . . . .	2
<b>III</b>	<b>Cleaning Task (finding and eliminating errors in dirty traffic detector data)</b>	2
III-A	Description . . . . .	2
III-B	Training data . . . . .	3
III-C	Test data . . . . .	3
III-D	Performance metrics . . . . .	4
III-E	Submissions . . . . .	4
<b>IV</b>	<b>Alignment Task (relating different representations of the same traffic event)</b>	5
IV-A	Description . . . . .	5
IV-B	Training data . . . . .	6
IV-C	Test data . . . . .	6
IV-D	Performance metrics . . . . .	6
IV-E	Submissions . . . . .	6
<b>V</b>	<b>Prediction Task (applying techniques to guess traffic-related events)</b>	7
V-A	Description . . . . .	7
V-B	Training data . . . . .	7
V-C	Test data . . . . .	7
V-D	Performance metrics . . . . .	7
V-E	Submissions . . . . .	7
<b>VI</b>	<b>Forecasting Task (applying predictive analytics to forecast traffic flow)</b>	7
VI-A	Description . . . . .	8
VI-B	Training data . . . . .	8
VI-C	Test data . . . . .	8
VI-D	Performance metrics . . . . .	8
VI-E	Submissions . . . . .	8
<b>VII</b>	<b>System Descriptions</b>	8
<b>VIII</b>	<b>Schedule</b>	8
<b>IX</b>	<b>Rules</b>	9
	<b>References</b>	9
	<b>Appendix A: Summary data table</b>	10
	<b>Appendix B: Data organization structure and example data files</b>	10

## I. INTRODUCTION

This document describes the plan for the National Institute of Standards and Technology (NIST) Information Access Division (IAD) Data Science Evaluation (DSE) Series Pilot Evaluation. This pilot evaluation is a continuation and extension of the Data Science Pre-Pilot Evaluation that was run in 2015, and precedes the first full evaluation of the DSE series. The primary goals of the pilot are to

- further develop and exercise the evaluation process in the context of data science prior to running larger scale evaluations in this context
- serve as an archetype for the development of future evaluation tasks, datasets, and metrics.

The pilot will consist of data and tasks set in the automotive traffic domain, which was picked for its reliability and because large amounts of public data are available.

The pilot is not meant to solve any particular problem in the traffic domain. Instead, the objective is for the developed measurement methods and techniques to apply to a broad range of use cases, regardless of the data type and structure.

The tasks included in the pilot are illustrated in Figure 1 and consist of:

- 1) **Cleaning:** finding and eliminating errors in dirty data.
- 2) **Alignment:** relating different representations of the same object in different data sources.
- 3) **Prediction:** determining possible values for an unknown variable based on known variables.
- 4) **Forecasting:** determining future values for a variable based on past values.

Each task can be completed independently; however, the tasks were designed so that the output of the cleaning task is an input to the remaining tasks. Hence, this pipelining of tasks (where the output of one task serves as the input to another task) forms a *workflow*. This workflow, will have two phases: in the first phase the original errorful data, along with the output of the data cleaning task, will be available as input for tasks downstream in the workflow; in the second phase, the ground truth for the data cleaning task will be available as input for tasks downstream in the workflow. A high-level summary of the pilot workflow is depicted in Figure 2.

## II. DATA

### A. Core sets

There are six datasets that will be made available as part of this evaluation. These datasets form a common core of data that are available for training and development for all tasks (with the exception of the cleaning task, which limits the data available; see Section III for details). The common core of data consists of the following sets, which are further described

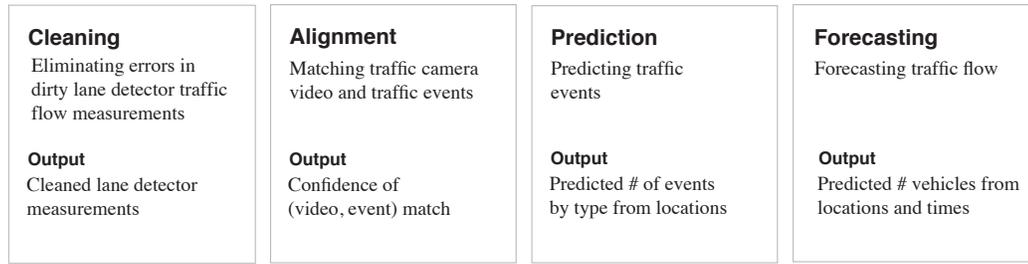


Figure 1. The pilot evaluation is comprised of four primary tasks: Cleaning, Alignment, Prediction, and Forecasting.

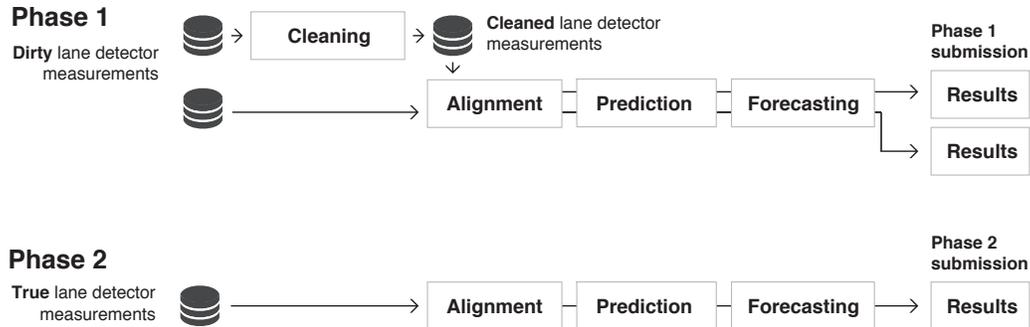


Figure 2. The pilot organizes four primary tasks into two phases, one in which lane detector measurements are “dirty,” and one in which “true” lane detector measurements are provided.



Figure 3. Six types of data are available in the pilot evaluation: Lane detector measurements, Maps, Traffic camera video, Traffic events, Weather data, and U.S. Census data.

in Appendix A Table III, with more details about the exact content of the sets and size estimates.

- Lane Detector Measurements<sup>1</sup>
- Traffic Events
- Traffic Camera Video
- U.S. Census and American Community Survey (ACS)\*
- OpenStreetMap (OSM) Maps\*
- Weather Data\*

The common core datasets are briefly described in Figure 3 and presented in more detail in Appendix A (Table III).

### B. Data access

Each above dataset marked with a \* will be made available through its respective organization’s url. The other datasets will be supplied through an Amazon S3 Storage bucket

<sup>1</sup>The Lane Detector Measurement dataset consists of output from multiple traffic speed and automobile count detectors at different locations along the road. Each individual detector is referred to as a **lane detector**. Lane detectors at the same location are aggregated into **traffic zones**. Figure 4 illustrates this notion of a traffic zone.

available to registered pilot participants.<sup>2</sup> Figure 5 gives the organizational structure of the data that will be supplied on Amazon S3.

### III. CLEANING TASK (FINDING AND ELIMINATING ERRORS IN DIRTY TRAFFIC DETECTOR DATA)

#### Summary 1: Cleaning task

This task involves cleaning traffic lane detector data flow values where a portion of its measurements were manipulated to be erroneous.

**Input.** traffic lane detector measurements with erroneous traffic flow values and speeds.

**Output.** cleaned traffic lane detector measurements with cleaned traffic flow values.

#### A. Description

In this task, participants are asked to clean traffic lane detector measurements containing incorrect flow values, pro-

<sup>2</sup>See <https://aws.amazon.com/s3/> for information about Amazon S3 Storage.

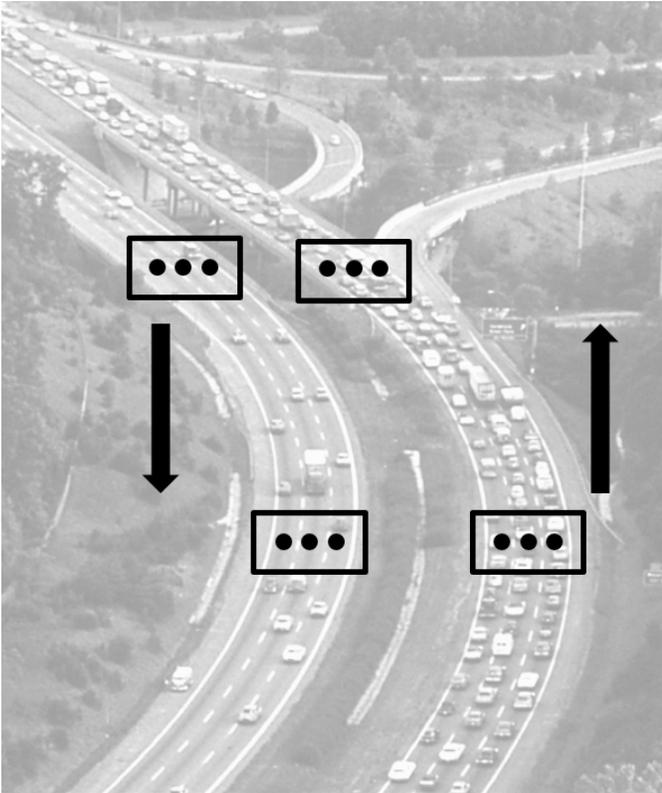


Figure 4. Traffic zones are an aggregation of traffic lanes. The dots are traffic lane detectors and the boxes are the traffic zones. In this figure, there are 12 traffic lane detectors: 4 traffic zones with 3 detectors in each zone. There is one traffic zone detector per road segment in each direction. In the data, there are different numbers of lane detectors in each traffic zone. This photo was obtained from National Archives [1].

viding correct traffic flow values for the erroneous traffic flow measurements. Detecting which values are erroneous is implicit to this task. The traffic flow measurement is the number of vehicles that have passed within a predetermined number of seconds. For each lane detector, that predetermined number of seconds is specified in the interval field and is often 60 seconds, i.e., the traffic flow value is often the number of vehicles to have passed through the detector in the previous minute.

It is important to note that the traffic speed measurements (given in miles per hour) of the sensors also contain errors, although the output of the task is restricted to cleaned traffic flow measurements.

### B. Training data

No cleaned traffic flow data will be provided as part of this task, making it “unsupervised.” However, other data without introduced errors will be available for this task, namely:

- 1) The traffic lane detector inventory (providing detector information, including each detector’s lane id and zone id), and
- 2) The OpenStreetMap data.

Unlike other tasks in the pilot evaluation, only a subset of the provided data is permitted for use in completing the cleaning task. See Rule 1.

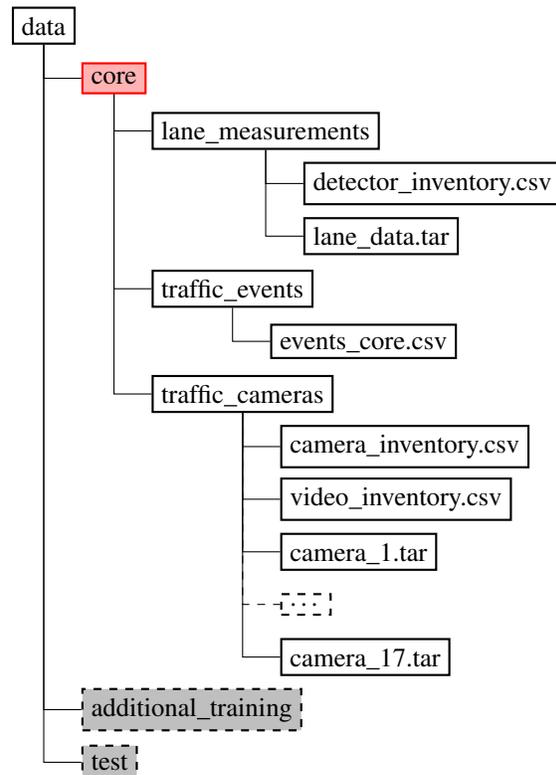


Figure 5. Organizational hierarchy for the pilot core data on Amazon S3. See other file hierarchy figures in the Appendix for the grey-ed out sub-hierarchies.

### Rule 1: Cleaning task rule: allowable data

Participants may only use data from the **traffic lane detector inventory**, **traffic lane detector measurements**, and the **OpenStreetMap** data sets to develop their systems.

### C. Test data

For this task, the test data will be all measurements of a selected subset of lane detectors. The detectors included in the test set will be selected by their zone id, and each detector will be specified by its lane id.

Errors were introduced into the test set measurements. Two of the attributes, `flow` and `speed` were changed. How the data is changed is a function of the lane detector and the measurement time.

#### 1) Flow errors.

- For each measurement from August 2015 or earlier (all lane detectors), with probability  $0 < p_n < 1$ ,  $p_n = 0.03$  uniformly at random, each value will have its flow altered.
- For each measurement September 2015 or later, depending on the measurement time, year, and the lane detector id, time  $t$ , lane detector  $d$ , with probability  $0 < p_{td} < 1$  uniformly at random, each value will have its flow altered.

#### 2) Speed errors.

- For each measurement from August 2015 or earlier (all lane detectors), with probability  $0 < p_n < 1$ ,  $p_n = 0.03$  uniformly at random, each value will have its speed altered. The altering of the speed values is independent of the altering of the flow values.
- For each measurement September 2015 or later, depending on the time and the lane detector id, for time  $t$  and lane detector  $d$ , with probability  $0 < p_{td} < 1$  uniformly at random, each value will have its speed altered. The altering of the speed values is independent of the altering of the flow values.

The test data will be available in multiple files, each containing the sensor outputs for a single month and named `cleaning_test_yy_mm.csv`, where `yy_mm` refers to the year and month during which the measurements were taken. For ease of submission, the measurements will be sorted first by the measurement timestamp and then by the lane id. Each test measurement will be labeled with a unique `trial_id` field.

For the cleaning task, and only for the cleaning task, participants will be permitted to interact with the test data (See Rule 2). Participants may interact with the traffic lane detector data for all tasks because the traffic lane detector data is training data for the other tasks.

<b>Rule 2: Cleaning task rule: interacting with test data</b>
For this task, and only for this task, participants <b>will be permitted</b> to interact with the test data.

#### D. Performance metrics

For this task, system output will be measured using as the cost the mean absolute value of the error. Formally, let  $n$  be the total number of measurements in the test set, which  $n$  is the number of trials. For each measurement  $i$ , let  $\widehat{fl}_i$  be the estimated traffic flow for measurement  $i$ , and let  $fl_i$  be the correct traffic flow for measurement  $i$ . Then the mean absolute value error is:

$$cost = \frac{\sum_{i=1}^n |\widehat{fl}_i - fl_i|}{n} \quad (1)$$

Furthermore, for additional analysis, an implicit detection decision will be extracted from the submissions. Any measurement that was changed from the test data will be treated as a correction of a detected erroneous measurement. Therefore, participants are encouraged to check that changes to measurements are intentional and not due to round-off or formatting errors.

The *cost* is a generalization of the *percentage error* metric for detection tasks<sup>3</sup>. In addition to this metric, we will be piloting a new cleaning metric as a second alternative metric.

<sup>3</sup>Like its detection counterpart, this metric is sensitive to the proportion of positives (dirty data measurements for this task), and given the (small) amount of positives, this metric may favor more conservative systems.

For this metric,  $x_i$  is the value of the provided flow measurement for measurement  $i$ . This alternative metric, the alternative cleaning cost (*cost<sub>alt</sub>*), can be defined by the following cost function:

$$cost_{alt} = \frac{\sum_{i=1}^n \left( 1 - c_d * \min \left( 1, \frac{|\widehat{fl}_i - x_i|}{c_{flmax}} \right) \right) |\widehat{fl}_i - fl_i|}{n} \quad (2)$$

In this metric, rather than just taking the arithmetic mean of the absolute errors in flow, each absolute error in flow is multiplied by a weight (or a percentage) that takes into account how much the estimated flow  $\widehat{fl}_i$  differs from the provided flow  $x_i$ . The greater the absolute difference  $|\widehat{fl}_i - x_i|$ , the lower the absolute error in flow is weighted.

The multiplied weight is  $\left( 1 - c_d * \min \left( 1, \frac{|\widehat{fl}_i - x_i|}{c_{flmax}} \right) \right)$ . There are two constants that parameterize this alternative cost function. The first,  $c_{flmax}$  specifies the maximum amount of change that results in a decreased weight to the absolute error. For this evaluation,  $c_{flmax} = 20$ , meaning that changing the flow value by more than  $c_{flmax}$  does not decrease the weight to the error. The second constant  $c_d$ , specifies the weight at which the error is reduced. For this evaluation,  $c_d = 0.4$ .

For example, first suppose that  $\widehat{fl}_i = 10$ ,  $fl_i = 20$ ,  $c_{flmax} = 20$  and  $c_d = 0.4$ .

- 1) If  $x_i = 10$ ,  $cost_{alt} = (1 - 0.4 * 0) * 10 = 10$ .
- 2) If  $x_i = 30$ ,  $cost_{alt} = (1 - 0.4 * 1) * 10 = 6$ .
- 3) If  $x_i = 15$ ,  $cost_{alt} = (1 - 0.4 * 0.25) * 10 = 9$ .

In summary, the new metric *cost<sub>alt</sub>* is the score with a discounted penalty for corrections to the measurements. Meaning, that as the system corrects more aggressively, the total error in flow is discounted more, up to a change  $c_{flmax}$  vehicles per interval. Although the total flow error is discounted, changing the flow value to err more may increase the cost because the error in flow increases the cost more than is decreased by the discount for correcting.

By default, the scoring metric will be the cleaning cost. Submitting two versions of a cleaning system, one for each metric is encouraged. Although the main metric is the *cost*, both metrics will be computed on all submissions for analysis purposes.

#### E. Submissions

System output results must be provided in multiple files, one per test data file, using standard ASCII encoding.

For every test file `cleaning_test_yy_mm.csv`, the corresponding submission file shall be named `cleaning_subm_yy_mm_{team}_{submission}.tsv`. Every submission file will have the same number of lines in the same order as the corresponding test file, and consist of a two values per line, each separated by a tab character. The first value is the `trial_id`, and the second is the *cleaned\_flow*, i.e., the corrected traffic flow for the corresponding measurement. Headers in the submission file are optional.

In the  $\{submission\}$  part of the name, indicate which metric is being optimized for by appending a “M1” at the end optimized for the cleaning cost ( $cost$ ), and a “M2” at the end if the system is optimized for the alternate cleaning cost ( $cost_{alt}$ ). If no metric is specified in the submission name, it will be assumed that the system is optimized for the original cleaning cost, as if the submission were appended with a “M1.”

Flow values for all trials must be provided. Participants are encouraged to submit the trials in the order of the trials in each test file, but at a minimum, participants must include the `trial_id` in each row. Files should be submitted without header rows.

#### IV. ALIGNMENT TASK (RELATING DIFFERENT REPRESENTATIONS OF THE SAME TRAFFIC EVENT)

##### Summary 2: Alignment task

This task involves matching traffic events with the traffic video segments containing those events.

**Input.** video segments from traffic cameras and traffic event data around that camera.

**Output.** a confidence value for each ( $v =$  video segment,  $e =$  traffic event) pair, where a greater confidence value indicates a greater belief that  $v$  and  $e$  refer to the same event.

##### A. Description

The goal of this task is to analyze video from camera feeds to detect an event and match it to a separate inventory of traffic events. This task may be divided into two steps:

- From video, detect the occurrence of one or more traffic events.
- Match the detected events to events in the event instances list.

In this section, a *recorded* event is a traffic event that is present in the video and a *reported* event is a traffic event present in the traffic event inventory. In this task, systems must detect events recorded in the video and match them to reported traffic events.

Table I lists and summarizes all of the different traffic event types in the traffic event data. The types of traffic events that are to be detected are a subset of those described in Table I, namely:

- Accidents and Incidents,
- Obstructions, and
- Device Status.

Note that some recorded events may not be reported, e.g., if the situation did not last long, did not disturb traffic, or resolved itself on its own. For example, a disabled vehicle may not be reported because the driver restarted the car after a few minutes. Systems should not link recorded events with events that are not reported in the traffic event inventory.

Each video segment has accurate time information, and the location of the source camera of each video segment is given in latitude and longitude. Some video feeds may have the

Table I  
TRAFFIC EVENT TYPES AND SUBTYPES

Traffic Event Type	Traffic Event Subtype
Delay Status Cancellation	Delay and disruption
Accidents And Incidents	Abandoned vehicle, accident, accident involving a semi trailer, accident involving a truck, disabled vehicle, hazardous material spill, incident, injury accident, minor accident, multi vehicle accident, numerous accidents, serious accident, vehicle on fire
Device Status	Sign down, traffic lights not working
Disasters	Brush fire, major flood, wildfire
Disturbances	Bomb alert, security alert, security incident
Incident Response Equipment	Other
Obstruction	Animal struck, debris on roadway, downed cables, drawbridge open, fallen trees, obstruction on roadway, subsidence
Pavement Conditions	Surface water hazard
Precipitation	Snow
Roadwork	Bridge construction, bridge maintenance operations, construction work, emergency maintenance, overgrown grass, overgrown trees, paving operations, road construction, road maintenance operations, road marking operations, road widening, storm drain, water main work, work in the median, work on underground services
Special Events	Concert, fair, major event, parade
Sporting Events	Sports event
System Information	Test message
Traffic Conditions	Traffic congestion
Visibility And Air Quality	Visibility reduced
Warning Advice	Alert, police at scene
Winds	Hurricane, strong, winds, tornado

direction the source camera is facing water-marked (“E” for East, for example).

The Traffic Events Instances file lists all *reported* traffic events in a comma-delimited format with geolocation information, a description of the location (intersection, for example), and the category of the traffic event. *In this additional dataset, all original timestamps will have been removed.*

The output of this task is a confidence value for each video segment and traffic event pair. All pairs must be evaluated, and must be evaluated independently from other pairs, i.e., a  $(v_x, e)$  pair may not be used to compute the value for a  $(v_y, e)$  pair.

When a recorded event overlaps multiple video segments, every video segment that contains some part of the recorded event will be considered a match with the corresponding reported event.

<b>Rule 3: Alignment Task Rule: Independent Trials</b>
All ( $v$ = video segment, $e$ = traffic event) pairs must be <b>evaluated independently</b> .

### B. Training data

All the common core data described in Section II will be available for training and development purposes.

An additional “Traffic Event” dataset will also be provided for training. This training set involves a subset of traffic cameras, denoted as training cameras, and for each of those cameras, the list of traffic events with timestamps will be provided. In more detail, this additional training data consists of:

- **Video:** 15-minute video segments from the training cameras, available from the core data, at `core_data/traffic_videos/camera_name`
- **Reported events:** All traffic events reported as having taken place at a distance less than  $d = 500$  meters from each training camera’s location. These traffic events will include timestamps. These events will be supplied in the alignment subfolder of the training data, see Appendix B (Figure 6).

### C. Test data

For the test data set, participants will be tested on different traffic cameras. Those cameras are denoted as the test cameras.

The test data consists of video segments and traffic events, similar to the training data described in Section IV-B, but the events provided will not have timestamps. In order to consider all traffic events reported in the vicinity of a camera, the test data will consist of:

- **Video:** 15-minute video segments from the source test cameras. These test video segments will be supplied in the `videos/camera_name` subfolder from the `test/alignment` folder.
- **Reported events:** All events reported as having taken place at a distance less than  $d = 500$  meters from each test camera’s location. *All original timestamp data will have been removed.*

See Appendix B (Figure 7) for the organizational structure of the test data and where to find it on Amazon S3.

Note that the cameras present in the training set may not be the same as those in the test set.

In addition to the video lists and the lists of events, there will be one test file per video camera which contains all of the possible (video, event) pairs indexed by a `trial_id`. These files will be named `alignment_{camera_name}_test.csv`. These files will contain trials, where each trial is three comma-separated values, which are the `trial_id`, `video_id`, and `event_id`.

### D. Performance metrics

A (video segment, reported event) pair is referred to as a *trial*. When the video segment contains a recorded event that

corresponds to the given reported event, this is referred to as a *target trial*. If a trial is not a target trial, it is considered a *non-target trial*. When a system outputs a non-match for a target trial, the resulting error is called a “miss,” and when a system outputs a match for a non-target trial, the resulting error is called a “false alarm.”

Each trial will be treated as a match or non-match by comparing the system output to a certain threshold; trials greater than or equal to the threshold will be considered matches, and all others will be considered non-matches. By using the sorted system outputs as thresholds, the system’s misses and false-alarms can be calculated at all possible *a posteriori* thresholds.

The performance measure will be based on a decision cost function (DCF) representing a linear combination of the miss and false alarm error rates at a threshold  $\tau$ . The cost function for this task will be:

$$DCF(\tau) = \frac{|\text{misses}(\tau)|}{|\text{target trials}|} + 100 \times \frac{|\text{false alarms}(\tau)|}{|\text{non-target trials}|} \quad (3)$$

The overall performance measure will be the minimum DCF value obtained considering all  $\tau$ . Hence the cost for a system on the alignment task is:

$$\text{cost} = \min_{\tau}(DCF) \quad (4)$$

For this task, it is possible that only a subset of the submitted values will be scored.

### E. Submissions

System output results must be provided in a single file per test camera, using standard ASCII encoding. The file will be named `align_subm_{camera_name}_{team}_{submission}.tsv`.

Each test camera’s submission file will list results for all (video, event) pairs, and will index each pair by trial id. The results must be listed one per line as two tab-separated fields:

- `trial_id`, from the test file
- `confidence value`, the computed confidence between the event and the video segment. Confidence scores are real numbers between 0 and 1, inclusive ( $0 \leq \text{confidence value} \leq 1$ ).

All trials, all possible combinations of event id (from the events list supplied for that camera) and video ids (from the video inventory file for that camera), must appear in the submission file, and every confidence value must be computed independently. Participants are encouraged to submit the trials in the order of the trials in each test file, but at a minimum, participants must include the `trial_id` in each row. Files should be submitted without header rows.

See Appendix B (Figure 8) for an example submission file.

## V. PREDICTION TASK (APPLYING TECHNIQUES TO GUESS TRAFFIC-RELATED EVENTS)

### Summary 3: Prediction task

This task involves predicting the number and types of traffic incidents in a region over a time period.

**Input.** geographical bounding boxes and time intervals.

**Output.** predicted counts for each specified type of traffic event.

#### A. Description

For this task, participants will develop a system that can predict the number and types of traffic events by type for a given (geographical bounding, interval of time) pair. This task will consider only a subset of the available event types, given in Table I. The traffic event types considered for this task are denoted  $\mathcal{E}$  and consist of:

- Accidents and Incidents.
- Roadwork.
- Precipitation.
- Device Status.
- Obstruction.
- Traffic Conditions.

The task output will be a list of counts of traffic events for each event type listed above.

#### B. Training data

All the common core data described in Section II will be available for training and development purposes.

#### C. Test data

The test data consists of a series of trials, where each trial is a (location, time interval) pair. Each trial is specified as follows:

- 1) a *trial\_id*.
- 2) The decimal latitude and longitude coordinates of the trial geographical bounding box.<sup>4</sup>
- 3) Start and end of the time window.<sup>5</sup>

For the evaluation of this task the participants will be provided with an ASCII tab separated file containing a list of input trials. For instance, a trial for the greater DC Metro area for a month of March of 2015 would be a single line that would consist of the following tab separated fields:

```
1 38.78 -77.25 39.0 -76.75
2015-03-01T00:00:00 2015-03-31T23:59:59
```

The following statements are true:

- The area of each bounding box is between  $0.25 \text{ km}^2$  and  $8 \text{ km}^2$ .

<sup>4</sup>Each bounding box is defined with the (latitude, longitude) pair of the North-West corner of the bounding box and the (latitude, longitude) pair of the South-East corner.

<sup>5</sup>Each timestamp follows the ISO 8601 format of combined date and time in UTC, i.e. in 'YYYY-MM-DDThh:mm:ss' format. The timezone is Eastern Standard time.

- Time intervals are at least 3 hours long and at most 31 days long.

Each (location, time interval) pair is selected with the number of traffic events in mind.

#### D. Performance metrics

Task performance will be scored with the average of each trial's Root Mean Squared Error (RMSE).

The RMSE is computed for each trial  $j$ . In each trial  $j$ , there are  $|\mathcal{E}|$  types of events for which to predict the number of events of type  $e$ , denoted  $\hat{e}_j$ . For each event type  $e$  and trial  $j$ ,  $e_j$  denotes the true count of events of that type.

$$RMSE(j) = \sqrt{\frac{\sum_{e \in \mathcal{E}} (\hat{e}_j - e_j)^2}{|\mathcal{E}|}} \quad (5)$$

An event is considered to be in the test window if any part of the event overlaps with the test timeframe and location, i.e., if any of the traffic event timestamps fall within the test window.

The RMSE for all the trials is then averaged to get a cost, where  $n$  is the number of trials:

$$cost = \frac{1}{n} \sum_{i=1}^n RMSE(i) \quad (6)$$

#### E. Submissions

System output must be provided in a single file using standard ASCII encoding and named `prediction_submissions_{team}_{submission}.tsv`.

The output for each trial, i.e., (location, time interval) pair, must consist of a single line with each line containing the *trial\_id* of the trial. Each line must consist of a list of seven tab-separated fields. The first field is the *trial\_id* and the remaining six fields consist of real values of the estimated number of events per type, ordered as follows:

- 1) Accidents and Incidents.
- 2) Roadwork.
- 3) Precipitation.
- 4) Device Status.
- 5) Obstruction.
- 6) Traffic Conditions.

Results for all trials must be submitted. Participants are encouraged to submit the trials in the order of the trials in each test file, but at a minimum, participants must include the *trial\_id* in each row. Files should be submitted without header rows.

## VI. FORECASTING TASK (APPLYING PREDICTIVE ANALYTICS TO FORECAST TRAFFIC FLOW)

### Summary 4: Forecasting task

This task involves providing future traffic flow values.

**Input.** a list of locations and times to forecast traffic flow.

**Output.** for each forecast location and time, a list of

forecasted values for every specified time interval.

### A. Description

In this task, participants will leverage past traffic information and current conditions (weather, maps) to forecast vehicle flows<sup>6</sup> on major roads. Vehicle flow is defined as the number of vehicles to pass a fixed mile marker on a road (for *all* lanes if there is more than one) for a period of time. All flow information stems from the “Traffic Lane Detector Measurements” set which provides flow, average speed, and lane occupancy information on one or more lanes at specific mile markers. These locations are often composed of multiple lanes, and therefore covered by multiple lane detectors, which report the flow, average speed, and lane occupancy information.

Systems will be required to complete a series of *trials*, where each trial is a (location, start time) pair. For each trial, the system must output a sequence of  $m$  real numbers, corresponding to estimated traffic flow for every  $t$  minutes over time interval  $I$ . That is, each value is the predicted number of vehicles to cross the specified location in the specified  $t$  minute interval. For this task,  $I = 60$ ,  $t = 5$  (and therefore  $m = 12$ ).

### B. Training data

All the common core data described in Section II will be available for training and development purposes.

### C. Test data

The test data are specified as a tab-delimited file listing locations and times for which to forecast flow.

- 1) trial\_id (referenced in the submission format, see Section VI-E)
- 2) location latitude (decimal).
- 3) location longitude (decimal).
- 4) short text description of the location, usually an intersection (for example “495/GW Parkway”).
- 5) datetime of start.<sup>7</sup>

### D. Performance metrics

Task performance will be scored with the average of each trial’s Mean Absolute Percentage Error (MAPE).

The MAPE is computed for each trial  $i$ . In each trial  $i$ , there are  $m$  forecasted values, one value for each interval. For each measurement  $j$  of the  $m$  real numbers which are the measurements for trial  $i$ ,  $\widehat{y}_{i,j}$  denotes the number of vehicles forecasted to pass through the trial location during the timespan  $[t * (j - 1), t * j]$  of trial  $i$ , and  $y_{i,j}$  denotes the true number of vehicles that passed through the specified location during that interval of time for trial  $i$ . The *MAPE* for each trial  $i$ , *MAPE*( $i$ ) is computed as:

$$MAPE(i) = \frac{1}{m} \sum_{j=1}^m \frac{|\widehat{y}_{i,j} - y_{i,j}|}{y_{i,j}}. \quad (7)$$

<sup>6</sup>Traffic flow and vehicle flow are used interchangeably.

<sup>7</sup>Following ISO 8601, the datetime format is YYYY-MM-DDThh:mm:ss

The MAPE for all the trials is then averaged to get a cost, where  $n$  is the number of trials:

$$cost = \frac{1}{n} \sum_{i=1}^n MAPE(i) \quad (8)$$

### E. Submissions

System output results must be provided in a single file using standard ASCII encoding. The file must be named forecasting\_{team}\_{submission}.tsv.

For each trial, the submission file will list system outputs by trial\_id: for each (location, time) pair listed in the test file, the  $m = 12$  forecasted values will be reported on a single line, separated by a tab. To identify each trial, the forecasted values for each trial will be prepended by the trial\_id, referencing the trial\_id from the test file.

A resulting line will look like the following:

```
trial_id    value1    value2    ...    value12
```

Participants are encouraged to submit the trials in the order of the trials in each test file, but at a minimum, participants must include the trial\_id in each row. Files should be submitted without header rows. See Appendix B (Figure 11) for an example submission file, and the matching figure in Appendix B (Figure 10) for an example test file.

## VII. SYSTEM DESCRIPTIONS

For each system output submission to a DSE task, a very brief description of each system is requested. Each system description should provide:

- 1) *System Name*. The name of the system.
- 2) *Task*. The task(s) the system(s) was used for: cleaning, alignment, prediction, or forecasting.
- 3) *Team Affiliation*. The names of the submitting team.
- 4) *System Summary*. A brief summary the system in a few sentences (at most a paragraph).
- 5) *Algorithmic Information*. A brief description of additional details about the methods and algorithms used (1-2 paragraphs).
- 6) *Data Sources*. A brief description of the data sources used. This data includes the training data sources as well as additional data fed to the algorithms. Mention which of the provided data sets were used as well as whether there were any external data sets used. Cite any external sources.
- 7) *Development Data*. (optional) a brief description of any development sets that were generated and any brief results.

## VIII. SCHEDULE

The **tentative** key dates for the NIST Pilot evaluation are given below in Table II. In particular, there will be two phases to the pilot, each corresponding to a different workflow, as was presented earlier in Fig. 2:<sup>8</sup>

<sup>8</sup>In describing this workflow, the term “dirty data” is used to refer to the traffic detector data participants are being asked to clean; “cleaned data” to refer to the output traffic flow data from the cleaning task; and “truth data” to refer to the ground-truth data, which is the correct answer to the cleaning task.

- **Phase 1.** Participants are given the common core data along with the dirty lane detector data and will be asked to submit system outputs from the four tasks using this data as input. Additionally, participants will be encouraged to submit the results of the alignment, prediction, and forecasting tasks using the cleaned traffic detector data.
- **Phase 2.** Participants are given the cleaning task truth data and will be asked to run the same systems for the alignment, incident, and flow tasks, using the cleaning task truth data as input.

Table II  
PILOT TENTATIVE SCHEDULE

Register to the evaluation by	July 31st, 2016
Release of training data and cleaning test data	August 1st, 2016
Release of the rest of the test data	October 28th, 2016
First submission deadline (all tasks)	November 28th, 2016
Release of cleaned ground-truth traffic detector data*	November 29th
Second submission deadline	December 6th, 2016
Release of initial results	January 12th, 2016
Workshop (location TBD)	March 2017 (Tentative, to be confirmed)

This schedule has two submission dates:

- **The first submission** (where the bulk of the time is given), for which participants complete tasks using the dirty lane detector data. Note that participants will be submitting results using both the dirty data as well as (optionally) the cleaned data from the cleaning task.
- **The second submission** (after the first submission deadline has passed), where participants use the lane detector ground truth data to complete the alignment, prediction, and forecasting tasks.

The cleaned ground truth may not be released. Furthermore, the cleaned ground truth may only be released to participants that provide sufficient submissions for the tasks in this evaluation. In particular, it is necessary (but perhaps not sufficient) that participants submit a task other than the cleaning that utilizes the traffic lane detector data and provides submissions that use the dirty traffic detector data and the cleaned traffic detector data with it.

## IX. RULES

The evaluation is subject to the following rules and restrictions:

- While participants **will** be allowed to use outside data, these data **must be publicly available**, and participants must include references to the data sources used (and how to obtain them) when submitting evaluation results. No internal or proprietary data are allowed to be used.
- Participants may not interact with the test data in any way, e.g., reading the test csv files or watching the test videos is prohibited. The test data for the cleaning task is the only exception to this rule, since the traffic lane detector data can be used as training data for other tasks.

- Each participating site must send one or more representatives who have working knowledge of the evaluation system to the evaluation workshop. Representatives must give a presentation on their system(s) and participate in discussions of the current state of the technology and future plans. Workshop registration information will be distributed to registered evaluation participants when available. The workshop will be open only to evaluation participants and representatives of interested government and supporting agencies.
- Dissemination:
  - NIST will generate and place on its web site charts of all system results for conditions of interest, but these charts will not contain the site names of the systems involved. Participants may publish or otherwise disseminate these charts, unaltered and with appropriate reference to their source.
  - Participants may not publish or otherwise disseminate comparisons of their performance results with those of other participants without the explicit written permission of each such participant. Furthermore, publicly claiming to “win” or suggest a ranking the evaluation is strictly prohibited. Any misrepresentation of the evaluation or its results is also strictly prohibited.
  - The results reported by NIST are not to be construed, or represented, as endorsement of any participant’s system, or as official findings on the part of NIST or the U.S. Government.

## DISCLAIMER

Certain commercial equipment, instruments, software, or materials are identified in this evaluation plan in order to specify the experimental procedure adequately. Such identification is not intended to imply recommendation or endorsement by NIST, nor is it intended to imply that the equipment, instruments, software or materials are necessarily the best available for the purpose.

## REFERENCES

- [1] Photograph, National Archives Identifier 546711 (Artist Yoichi R. (Robert) Okamoto), “Looking south from beltway bridge over the potomac. the capital beltway circles the virginia-maryland suburbs and provides high speed access to points in the district, 5/1973,” May 1973, DOCUMERICA: The Environmental Protection Agency’s Program to Photographically Document Subjects of Environmental Concern, 1972–1977 Record Group 412: Records of the Environmental Protection Agency, 1944–2006.
- [2] “Maryland chart traffic cameras,” 2015. [Online]. Available: <http://www.chart.state.md.us/travinfo/trafficcams.php>
- [3] “2010 u.s. census,” 2015. [Online]. Available: <http://www.census.gov>
- [4] “Openstreetmap,” 2015. [Online]. Available: <http://www.openstreetmap.org/>
- [5] “NOAA’s integrated surface hourly,” 2015. [Online]. Available: <http://www.ncdc.noaa.gov/isd>
- [6] J. N. Lott, “The quality control of the integrated surface hourly database,” in *84th American Meteorological Society Annual Meeting*, vol. 7.8. Seattle, WA: American Meteorological Society, 2004. [Online]. Available: <http://www1.ncdc.noaa.gov/pub/data/inventories/ish-qc.pdf>
- [7] “NOAA,” 2015. [Online]. Available: <http://www.ncdc.noaa.gov/swdi>

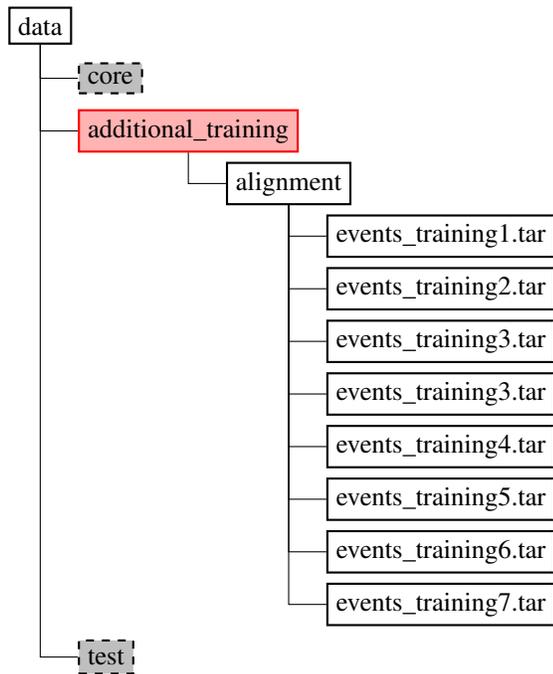


Figure 6. File hierarchy for the additional training data provided for the alignment task. See other file hierarchy figures for the grey-ed out sub-hierarchies.

#### APPENDIX A SUMMARY DATA TABLE

Table III describes each of the data sets, including the specific fields in each case, supplementing the information provided in Section II. The documentation associated with each data source provides a detailed description of all of its corresponding fields.

#### APPENDIX B DATA ORGANIZATION STRUCTURE AND EXAMPLE DATA FILES

This section presents the organizational structure of the data that will be provided to participants. These figures supplement Figure 5, which gives the overall structure of the data in the common core. Additionally, to clarify the descriptions of the test and submission sets throughout the evaluation, select examples for test and submission files are provided. These partial examples contain hypothetical data in that the id numbers in the examples may not correspond to actual id values in the data.

Figure 6 shows the organizational structure of the additional training data provided for the alignment task described in Section IV-B.

Figure 7 shows the organizational structure of the test data for the alignment task described in Section IV-C.

Figure 8 gives a partial example of a submission file for the alignment task described in Section IV-E.

Figure 9 shows the organizational structure of the test data for the forecasting task described in Section VI-C.

Figure 10 gives a partial example for a test set data file for the forecasting task described in Section VI-C.

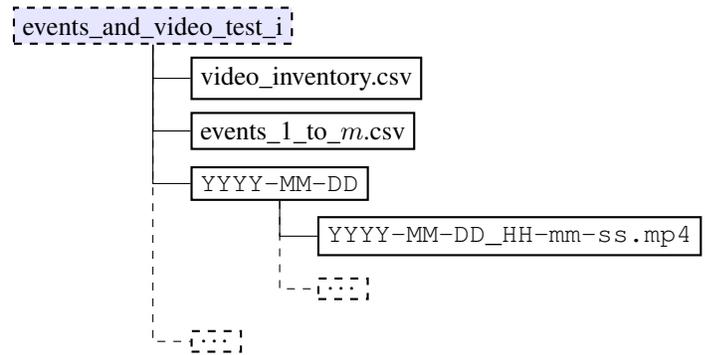


Figure 7. Contents of a *events\_and\_video\_test\_i* folder (test data for the alignment task)

1	0.1	
2	0.01	
		...
m	0.8	

Figure 8. Example submission file for the alignment task. Spaces between fields are tabulation characters. Headers are *not* required and are provided here for clarification. In this example, the test camera has  $m$  video segments and  $n$  events.

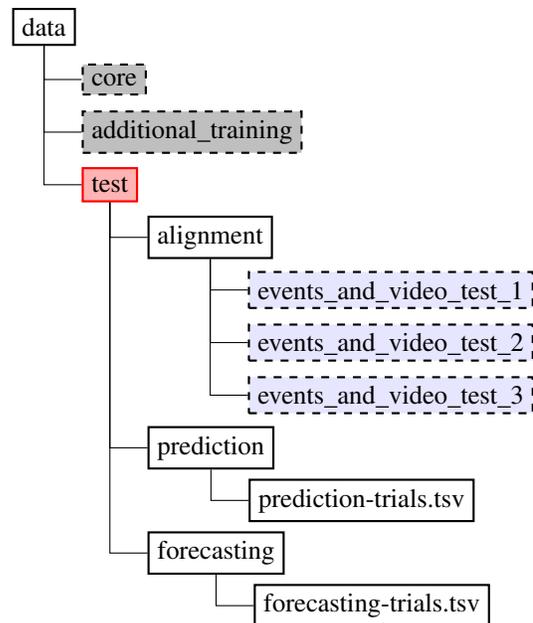


Figure 9. File hierarchy for the evaluation test data. Note that the cleaning task test data is in the “core” section rather than in the “test” section. See other file hierarchy figures for the grey-ed out sub-hierarchies, and Figure 7 for the contents of the folders in blue.

Figure 11 gives a partial example for a submission file for the forecasting task described in Section VI-E.

Table III  
SUMMARY OF AVAILABLE DATASETS.

Data Type	Data Subset	Description
Lane Detector	Lane Detector Inventory	List of all traffic lane detectors as a comma-separated file. Each detector is uniquely identified by its <code>lane_id</code> value, and each detector inventory gives the location of the detector (in decimal latitude and longitude coordinates), the source organization for the measurements of those detectors, the time interval between scheduled measurements, and other relevant information.
	Lane Detector Measurements	Measurements from traffic sensors in locations in the DC Metro area and the Baltimore area. Traffic sensors are placed on both directions of the highways, in each lane. Lane and zone (multiple lanes of the same road going in the same direction) data are provided. The measurements include the following attributes (among others): <ol style="list-style-type: none"> <li>1) <i>Flow</i>: the number of vehicles to have passed through the lane detector since the last scheduled measurement.</li> <li>2) <i>Speed</i>: the average vehicle speed since the last measurement.</li> <li>3) <i>Occupancy</i>: the average percent of time a vehicle was in front of the detector since the last measurement.</li> <li>4) <i>Quality</i>: a data quality field.</li> </ol> <p>This dataset has data from 2006 to 2015, collected from the DC-Maryland-Virginia area, and is approximately 150GB.</p>
Traffic Events	Traffic Event Instances	A traffic event is defined as a situation that involves traffic in the ways described in Table I. Prominent features of an event are injuries, damage to vehicles, hazards to persons, failure of equipment, closure of one or more lanes, debris or roadkill on the road or shoulder, or any obstruction on roadway. Each traffic event listing includes the following fields (among others): <ol style="list-style-type: none"> <li>1) Description.</li> <li>2) Location, both in formatted text (the intersection) and in decimal latitude and longitude.</li> <li>3) Times the event was created, confirmed and closed.</li> <li>4) The type and subtype of the traffic event; the field labels are <code>event_type</code> and <code>event_subtype</code>.</li> </ol> <p>A traffic event is reported in the inventory when it is outlined by or to authorities (911 calls, road kill pickups, etc.). It is confirmed when an authority arrives at the scene (the police arrives, etc.) and deemed over when it was closed by authorities. For an accident, this typically indicates when all lanes have been reopened, damaged vehicles have been removed, and all responders have left the scene. This dataset has data from 2003 to 2015, collected from the DC-Maryland-Virginia area, and is approximately 200MB.</p>
Traffic Camera Video [2]	Camera Inventory	A list of all traffic cameras with their locations, described both in text (the intersection) and in decimal latitude and longitude.
	Camera Video Feeds	Consecutive 15-minute video segments from traffic cameras in Maryland with start times. The traffic cameras may be remotely operated by humans, who can rotate the camera and zoom, which happens when the human operator chooses to look at a traffic situation. Some cameras may have watermarks indicating the direction the camera is facing (E for East, SW for South-West, etc), or the current time. This dataset has data from 2015, collected from Maryland, and is approximately 3TB.
U.S. Census	2010 U.S. Census [3]	Publicly available information including population counts; age, income, and occupation demographics; and household demographics in summary files and PUMS (Public Use Microdata Sample).
	American Community Survey (ACS)	A more frequent survey providing statistics on transportation and commutes, such as the average commute length, the percentage of people who carpool, and the percentage of people who use public transportation. There are 1-year, 3-year, and 5-year surveys as summary Files and PUMS, like the U.S. Census Data.
OpenStreetMap [4]	[No subset]	Map data from from OpenStreetMap, describing the road network in the DC-MD-VA area as well as locations including airports, public transportation stations, and buildings that host large events. These maps also support lookup by latitude and longitude coordinates.
Weather	Integrated Surface (ISD) [5]	A dataset of measurements from weather stations in the DC-MD-VA area with a variable number of measurements. Measurements include station information, temperature, air pressure, weather condition, precipitation, and other elements. The ISD set is quality-controlled. The quality control does not state that it is free of errors or missing data; only that others have looked at it to try to improve the quality of the data. Lott [6] discusses the quality control process that are used in the ISD to check for formatting errors and outliers.
	Severe Weather Data Inventory (SWDI) [7]	A compilation of many types of severe weather, including storms, hail, tornados and lightning strikes.

trial_id	lat	long	description	start-datetime
1	37.00	78.00	"495/GW Parkway North"	2015-09-18T13:00:00
2	37.00	78.00	"495/GW Parkway South"	2015-09-18T14:00:00
			...	
100	36.90	77.10	"495/River Rd North"	2015-09-19T13:00:00

Figure 10. Example test data for the forecasting task. Spaces between fields are tabulation characters. Headers are *not* required and are provided here for clarification.

1	12	13	...	4
2	2	4	...	2
			...	
100	2	5	...	2

Figure 11. Example submission file for the forecasting task. Spaces between fields are tabulation characters. Headers are *not* required and are provided here for clarification.