

1 **Clustering Surrogate Safety Indicators to**
2 **Understand Collision Processes**

3 Nicolas Saunier (corresponding author)

4 Assistant Professor, Department of Civil, Geological and Mining Engineering

5 Polytechnique Montréal

6 nicolas.saunier@polymtl.ca

7 Mohamed Gomaa Mohamed

8 PhD Candidate, Department of Civil, Geological and Mining Engineering

9 Polytechnique Montréal

10 mohamed.gomaa@polymtl.ca

11 5333 words + 8 figures + 1 tables

12 August 1, 2013

1 **ABSTRACT**

2 As time series are collected through more and more pervasive devices carried by users and vehicles,
3 new tools are necessary to understand and mine the large amounts of transportation data being thus
4 generated. This work proposes a new similarity measure for time series that is applied to surrogate
5 measures of safety and other indicators characterizing road user interactions. The new similarity
6 measure based on the aligned longest common sub-sequence is paired with a custom clustering
7 algorithm that does not require to set the number of expected clusters and remains interpretable
8 through the use of prototype indicator profiles as cluster representatives. The method is applied
9 to five indicators, including time to collision and probability of collision, for a large real world
10 dataset of traffic videos of collisions and conflicts. The results confirm the general assumption of
11 surrogate methods for safety analysis that some interactions without a collision have very similar
12 processes to collisions. It also highlights the danger of using a significant proportion of candidate
13 interactions without a collision that seem to share little similarities with collisions.

1 INTRODUCTION

2 As Moore's law makes sensors and computers ubiquitous in vehicles, road environments and on
3 users, more and more data is collected continuously on all vehicles and road users. Examples are
4 location data from vehicle and personal GPS sensors, road user trajectories extracted from traffic
5 cameras, vehicle kinematic data and engine operational data from on-board diagnostics (OBD) de-
6 vices. Storing this data results in large datasets of temporal measurements characterizing different
7 elements of the road system. This data has the potential to be useful for several transportation
8 applications, e.g. activity patterns, vehicle-based and site-based safety diagnosis, calibration and
9 validation of macroscopic and microscopic models, behaviour observations at various space and
10 time scales. However the promise of this "big data" can only be fulfilled if new methods are
11 developed to deal with it and mine the large datasets that can be accumulated. Though aggregat-
12 ing spatio-temporal data over time and space wastes the potential of data collected at much finer
13 resolutions, analysis often rely on reduced data for lack of expertise and tools, and for practical
14 reasons.

15 Of particular interest is the development of methods for the surrogate analysis of safety.
16 Traditional collision-based diagnosis methods have several shortcomings that have been repeatedly
17 covered in previous work, e.g. in (1, 2). There is therefore a search for proactive methods that do
18 not require to wait for accidents to occur. These surrogate methods rely on the observation of
19 all interactions and the measure of their "severity" or proximity to a potential collision through
20 continuous safety indicators such as the time to collision (TTC). These observations are more and
21 more commonly obtained automatically through vehicle-based sensors such as data loggers (3) or
22 dedicated devices installed for example for naturalistic driving studies (4) and site-based sensors
23 such as video cameras with video analysis software (2, 5, 6). However, most analyses of this data
24 and, to the authors' knowledge, all analyses of surrogate safety indicators rely on the aggregation
25 of the temporal indicators into a single value. The most commonly used in traffic conflict analyses
26 is the minimum TTC, or a severity level based on the TTC at a specific instant coupled with the
27 road user speed at the same instant in the case of the Swedish traffic conflict technique (7). This is
28 a terrible loss of information that could partially explain the mixed results to validate and transfer
29 surrogate measures of safety. As stated in (8), "the problem with values taken at a certain time is
30 that they do not incorporate any information before or after the chosen moment, creating a risk that
31 even very different encounters might be classified in the same category".

32 This paper is a follow up on (9) that relied on contextual information and aggregated mea-
33 sures of road users' individual speeds and speed differential to cluster interactions with and without
34 a collision. The purpose is to better understand collision processes and the similarities between all
35 interactions. This will help determine whether all interactions without a collision can be used as
36 surrogates for collisions. There is preliminary evidence that this is not the case and that some cate-
37 gories of interactions or interactions of different severity levels may not be associated with safety,
38 i.e. that more interactions would translate into more collisions over the long run, and even be in-
39 dicators of a good level of safety through the promotion of driver awareness and learning through
40 interactions with other road users (10).

41 This paper presents ongoing work on the development of a method to compare and cluster
42 time series or profiles of interaction indicators, including surrogate measures of safety. A new
43 similarity measure based on the longest common sub-sequence (11) is proposed to better measure
44 indicator profile similarity by taking into account the rate of change. A custom clustering algorithm
45 is developed that does not require to set the number of expected clusters and remains interpretable

1 through the use of prototype indicator profiles as cluster representatives. The method is demon-
2 strated on a large video dataset of interactions with and without a collision. The final contribution
3 of this paper is the development of an open source library that implements the proposed methods
4 and the release of the exact code used to produce the presented analyses in order to encourage re-
5 producibility and wider adoption of the methods. The background is presented in the next section.
6 It is followed by a description of the method, which is then demonstrated on a real world dataset.
7 Finally the paper is concluded and future work is discussed.

8 **BACKGROUND**

9 **Surrogate Safety Analysis**

10 There is a growing body of literature on surrogate methods for safety analysis and readers are
11 referred to these PhD theses (*12, 13, 14*) and the TRB white paper (*15*) for an introduction to
12 the field and a coverage of the early techniques known as traffic conflict techniques (TCT). The
13 defining characteristic of relevant traffic events for safety is the collision course, which is the
14 situation in which two road users would collide if their movements remain unchanged (taken from
15 the definition of a traffic conflict as “an observable situation in which two or more road users
16 approach each other in time and space to such an extent that there is a risk of collision if their
17 movements remain unchanged” (*16*)). Identifying a collision course at a given instant therefore
18 requires to predict road users’ future positions from their current and past positions. The default
19 motion prediction method is to assume that the road users will move at constant velocity. The
20 choice is rarely justified, does not yield robust measurements and does not take the context (the
21 road, e.g. in a curve, and traffic) into account which results in unrealistic motion prediction (e.g.
22 going off the road or into a wall). New prediction methods have been proposed in (*17, 18*) with
23 open source implementations (*6*).

24 For surrogate safety analysis to be objective, a number of quantitative safety indicators have
25 been proposed in the literature to measure the proximity to a potential collision, or probability of
26 collision, and the severity of the potential collision. TTC is the best known of these indicators. It
27 is defined for a given motion prediction method as the time required for two road users to collide
28 following the predicted trajectories. If several predicted trajectories are available, with correspond-
29 ing probabilities, the expected TTC can be computed (*2*). Many other safety indicators, including
30 post-encroachment time (PET), deceleration to safety time, etc., have been presented over the years
31 (see (*12, 13, 14*) and their references for more details).

32 **Interpreting Interactions and Safety Indicators**

33 In most TCTs, a specific value of a continuous safety indicator is used and compared to a threshold
34 to distinguish, usually for diagnosis purpose, the most severe conflicts from “safer” interactions,
35 defined as a situation in which two road users are within some distance. For example, Hydén
36 (*19*) used the TTC just before one of the road users attempts an evasive action called the time to
37 accident with a threshold of 1.5 s to define severe conflicts. The Federal Highway Administration
38 (FHWA) designed the piece of software Surrogate Safety Assessment Model (SSAM) to perform
39 the analysis of trajectory data extracted from microscopic simulation software (*20*). SSAM uses a
40 predefined threshold for different safety indicators to identify the most severe conflicts among all
41 road user interactions (e.g. the default threshold on minimum TTC is 1.5 s). The most severe value
42 of safety indicators is typically used to summarize them, for example minimum values for spatio-
43 temporal indicators (e.g. distance, TTC and predicted PET) or maximum values for probability

1 of collision or deceleration to safety time. However, as argued in the introduction and in (8),
2 narrowing down the whole interaction to a single value leads to losing a lot of information. Even
3 the work of Minderhoud and Bovy (21) highlighted in (8) still condenses the whole indicator profile
4 into a single measure through integration.

5 There are few examples of the use or interpretation of continuous traffic event indica-
6 tors over certain time intervals. Some studies on driver behaviour have relied on speed profiles.
7 Parkhurst (22) examined the shape of speed profiles to understand the driver behaviour at urban
8 and rural non-signalized intersections. Lareshyn et al. (23) classified the speed profiles, extracted
9 using automated video analysis, of vehicles making left turns at a signalized intersection and inter-
10 acting with oncoming traffic and crossing pedestrians. Among the three types of pattern recogni-
11 tion techniques tested, cluster analysis (k-means), supervised learning (k-nearest neighbours), and
12 dimension reduction, k nearest neighbours was found to perform well with respect to the human
13 observer annotations.

14 The goal of using whole safety indicators time series is to better understand collision pro-
15 cesses and how interactions with and without a collision compare. Indeed, the work of Davis et al.
16 (24) on a small set of traffic events suggests that the evasive actions undertaken by road users in-
17 volved in conflicts may be of a different nature than the ones attempted in collisions. The work
18 of Svensson and Hydén (10) provides some evidence that interactions with fairly high severities
19 could be associated with improved safety because they are frequent and severe enough to create
20 and maintain awareness among road users.

21 **Time-Series Clustering**

22 The objective of clustering is to classify the data into groups (clusters) with similar characteristics.
23 Because the groups are not known, clustering is also called unsupervised classification. Many al-
24 gorithms have been proposed in the machine learning literature (25), e.g. hierarchical, based on
25 density, centroids, statistical distributions, etc. A time series is a data type that represents a se-
26 quence of observation vectors $X(t) = [x_1(t), \dots, x_n(t)]$ as a function of time t , usually at discrete
27 instants. Time series can be univariate (one variable per observation, $n = 1$) or multivariate (many
28 variables per observation, $n \geq 2$). The readers are referred to (26, 27) for surveys of clustering
29 methods for time series data. Among the different algorithms developed in various domains, most
30 attempt to reduce the dimensionality of the data to enhance the clustering performance. For exam-
31 ple, Vlachos et al. (28) used k-means clustering incrementally at different levels (resolution) based
32 on discrete wavelet transformation (DWT) decomposition.

33 The key component of many clustering methods is the measure of the similarity or dis-
34 tance between pairs of elements in the series. The Euclidean distance is popular, but it requires
35 that both time series have the same length and it is sensitive to distortions (e.g. shifting along the
36 time axis) and noise. The development of elastic distance and similarity measures, such as dy-
37 namic time warping (DTW) and longest common sub-sequence similarity (LCSS), overcome the
38 previous drawbacks. Both DTW and LCSS are implemented using dynamic programming. DTW
39 attempts to find the best alignment between two time series by minimizing the distance between
40 them. Conversely, LCSS finds the length of the longest matching sub-sequence by comparing ev-
41 ery point of the two time series using a given matching method. Morris and Trivedi (29) evaluated
42 different similarity measures (HU, PCA (Principle Component Analysis), DTW, LCSS, PF (Picia-
43 relli and Foresti (30)), Modified Hausdorff) and clustering methods for trajectories as a first step
44 to understand road user behaviour. After tests on six different datasets, the authors concluded that

1 LCSS was consistently the top performer.

2 A relevant example for transportation of multivariate time series clustering is trajectory
 3 clustering which is done in surrogate safety analysis (17) and robotics applications. The objective
 4 is to cluster a dataset of observed trajectories into the main motion patterns. Bennewitz et al.
 5 (31) learnt the motion patterns of people in a scene using the Expectation Maximization (EM)
 6 algorithm, which enabled a robot to update its behaviour accordingly. Hu et al. (32) modelled road
 7 user activities with a fuzzy self-organizing neural network. One of the main applications is future
 8 motion prediction. Recently Morris and Trivedi (33) proposed a 3-stages hierarchical learning
 9 framework to analyze object activities and to predict future activities, as well as to detect abnormal
 10 events. The authors used LCSS as a similarity measure and spectral clustering algorithm (34) for
 11 the trajectory clustering.

12 PROPOSED APPROACH

13 Real time series from transportation will have varied lengths. This is the case for the analysis
 14 of safety indicators investigated in this paper. The choice is made to avoid pre-processing the
 15 data that would introduce distortions and may lead to loss of information for example through
 16 re-sampling. Methods that can deal with the data as it is, without pre-processing, are therefore
 17 preferred. Among the various such methods, the LCSS is favoured as it is flexible and can be
 18 adapted to specific purposes.

19 The Aligned Longest Common Sub-sequence

20 Let $X = [X(t_1), \dots, X(t_n)]$ and $Y = [Y(t_1), \dots, Y(t_m)]$ be two time series of respective length n and m of
 21 safety indicators characterizing two interactions (the series may be multivariate, e.g. if concatenat-
 22 ing several indicator measurements at each instant). Let $Head(X)$ be the series $[X(t_1), \dots, X(t_{n-1})]$.
 23 Given a real number $\delta > 0$ and a matching function $match$ for the elements of the series (e.g. for
 24 univariate series and a given real number $\varepsilon > 0$, $d_\varepsilon(a, b)$ is *true* if $|a - b| \leq \varepsilon$, *false* otherwise),
 25 the length $LCS_{\delta, match}(X, Y)$ of the longest common sub-sequence is computed as

- 26 • 0 if $m = 0$ or $n = 0$,
- 27 • $1 + LCS_{\delta, match}(Head(X), Head(Y))$ if $match(X(t_n), Y(t_m))$ is *true* and $|n - m| \leq \delta$,
- 28 • $max(LCS_{\delta, match}(Head(X), Y), LCS_{\delta, match}(X, Head(Y)))$ otherwise.

29 This is typically computed in a matrix S using dynamic programming where $S_{i,j}$ is the *LCS*
 30 for the respective sub-sequences of X and Y $[X(t_1), \dots, X(t_i)]$ and $[Y(t_1), \dots, Y(t_j)]$. The matrix is of
 31 size $(n + 1, m + 1)$ and initialized to zero. The $S_{i,j}$ is then iteratively computed using the following
 32 algorithm:

- 33 • for $i \in [1, \dots, n]$
- 34 – for $j \in [max(1, i - \delta), \dots, min(m, i + \delta)]$
- 35 * if $match(X(t_i), Y(t_j))$
- 36 · $S_{i,j} = S_{i-1,j-1} + 1$
- 37 * else
- 38 · $S_{i,j} = max(S_{i-1,j}, S_{i,j-1})$

1 The maximum value of the matrix is the *LCS*. To be comparable, independently of the
 2 indicator respective lengths, the associated similarity measure $LCSS_{\delta,match}(X,Y)$ and distance
 3 $DLCS_{\delta,match}(X,Y)$ are typically derived as (11)

$$LCSS_{\delta,match}(X,Y) = \frac{LCS_{\delta,match}(X,Y)}{\min(n,m)} \quad (1)$$

$$DLCS_{\delta,match}(X,Y) = 1 - LCSS_{\delta,match}(X,Y) \quad (2)$$

$$(3)$$

4 The parameter δ was introduced in (11) to control how far in time elements of the two series
 5 can be matched. This is not suited for series that have different lengths, which is not tested in (11).
 6 As an example (plotted in FIGURE 1), $LCSS_{4,d_{0.1}}([0, 1, \dots, 19], [10, 11, \dots, 19]) = 0$ (no similarity)
 7 while $[10, 11, \dots, 19]$ is an exact sub-sequence of $[0, 1, \dots, 19]$. A solution is to use a simpler version of
 8 LCS without δ (which is equivalent to choosing $\delta = +\infty$): $LCSS_{+\infty,d_{0.1}}([0, 1, \dots, 19], [10, 11, \dots, 19]) =$
 9 1 (maximum similarity). This causes other issues as it allows any value to match any other value
 10 irrespective of the rate of change in the series (however, the order in the series is always re-
 11 spected). Take for example the series $X = [0, 1, \dots, 19]$ and $Y = [0, 2, \dots, 18]$ (plotted in FIGURE 1):
 12 Y increases at twice the rate of X (with a step of 2 instead of 1), but is still a sub-sequence of
 13 X . If for a given application series evolving at different rates of change are considered dissim-
 14 ilar, computing *LCSS* without δ is inappropriate as this example shows: $LCSS_{+\infty,d_{0.1}}(X,Y) = 1,$
 15 while $LCSS_{1,d_{0.1}}(X,Y) = 0.2$. Two other examples for two safety indicators, distance and TTC,
 16 are showed in FIGURE 3 and illustrate how similarity is over-estimated by the traditional *LCS*
 17 computation, while a finite δ does not allow to compute the similarity of the series because they
 18 are most similar parts must be aligned.

19 It follows that the existing formulations of the longest common sub-sequence, with or with-
 20 out δ , are insufficient to measure the similarity of series if the series are simply shifted with respect
 21 to each other or if series with different rates of change should be considered different. That is why
 22 a new similarity measure is introduced that finds the best alignment of two series while taking into
 23 account a finite δ , allowing to take into account the rates of change. The length *ALCS* of the aligned
 24 longest common sub-sequence is computed by simply shifting the two series with respect to each
 25 other, i.e. by adding an integer parameter *shift* to the *LCS* computation (replacing the condition
 26 $|n - m| \leq \delta$ by $|n - shift - m| \leq \delta$) and taking the maximum *LCS* for all possible *shift* values. The
 27 corresponding aligned similarity measure *ALCSS* and distance *DALCS* are defined accordingly.

28 Another benefit is to use the longest common sub-sequence itself. The indices correspond-
 29 ing to the elements of the series that are matched to obtain the longest common sub-sequence are
 30 obtained by “decoding” the process of the computation of the *LCS*. For example, the longest com-
 31 mon sub-sequence of series $X = [1, 3, 5, 6, 7]$ and $Y = [1, 2, 3, 4, 6, 7, 8]$, using $d_{0.1}$ and finite δ , are
 32 respectively $[0, 1, 3, 4]$ and $[0, 2, 4, 5]$ meaning that the element in position 0 of X matches element 0
 33 of Y , element 1 in X matches element 2 of Y , etc. From these indices can be computed the average
 34 difference of the corresponding indices which corresponds to the “optimal” alignment of one series
 35 with respect to the other. The alignment corresponding to the *ALCS* is obtained by applying the
 36 *shift* corresponding to the maximum *LCS* to the optimal alignment of the longest common sub-
 37 sequence indices. This is very useful to visualize the data and validate the similarities: FIGURE 2
 38 shows the alignment obtained for two TTC indicators considered completely similar if aligned and

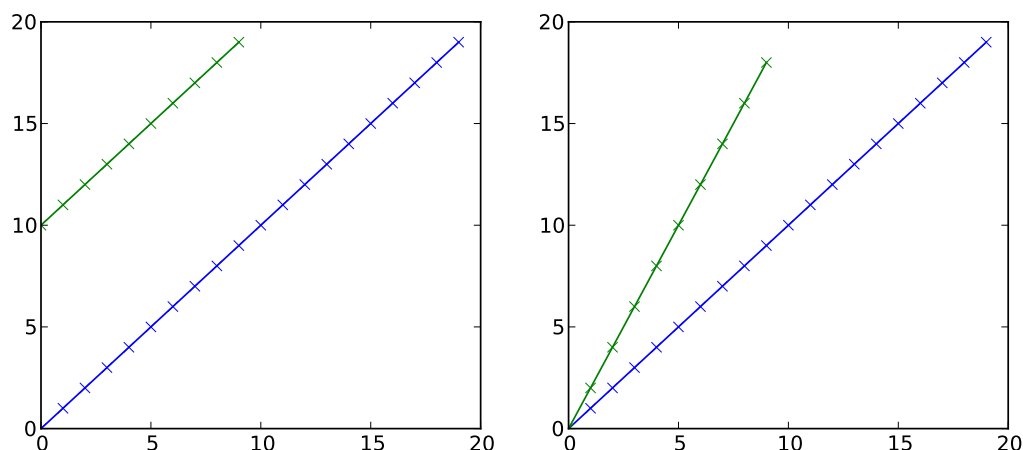


FIGURE 1 Examples of simple series that illustrate the advantages of using a finite δ and aligned longest common sub-sequence. The series in each plot have maximum similarity if using $\delta = +\infty$. This is desired in the plot on the left since it is an exact sub-sequence, but not on the right if the rate of change is taken into account.

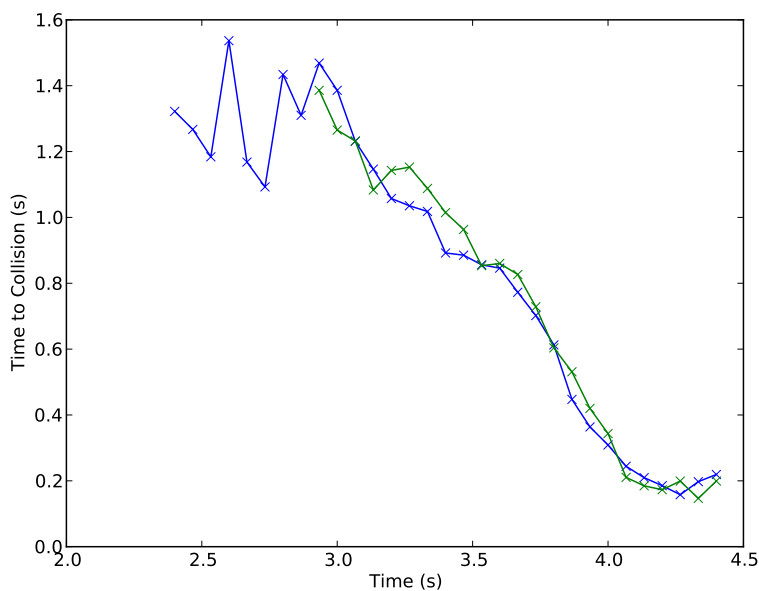
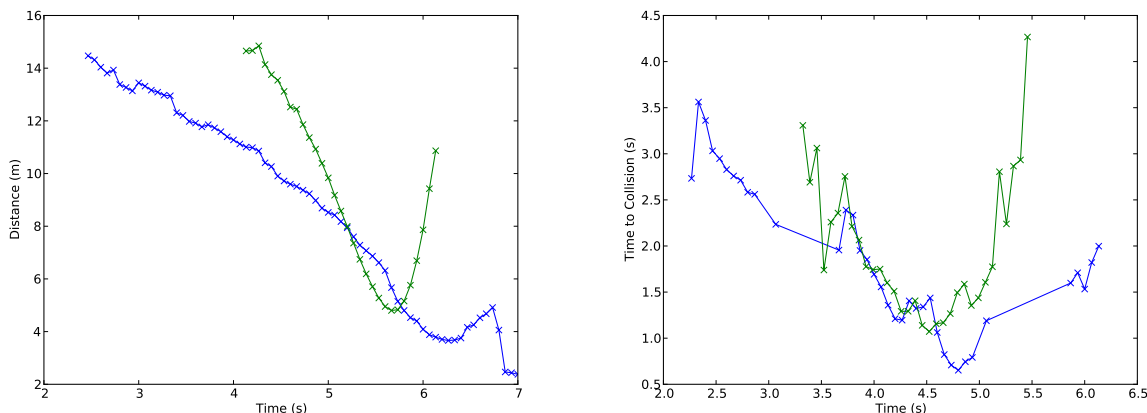


FIGURE 2 Example of alignment of two very similar TTC indicators ($LCSS_{2,d_{0.2s}} = 0.2$ and $ALCSS_{2,d_{0.2s}} = 1$).

- 1 barely similar otherwise ($\delta = 2$ and $d_{0.2s}$). The distance and TTC indicators are also aligned in
- 2 FIGURE 3: these examples of safety indicator profiles should not be considered similar, at least

1 not to the degree implied by the *LCSS* with infinite δ .



	Distance	TTC
$LCSS_{+\infty}$	0.87	0.64
$LCSS_2$	0.35	0.12
$ALCSS_2$	0.42	0.42

FIGURE 3 Examples of pairs of profiles for the interaction distance and TTC indicators that are more similar using *LCSS* with infinite δ than using *ALCS* and a finite δ . The matching function used is d_ε with $\varepsilon = 1$ m for the distance indicator and $\varepsilon = 0.2$ s for the TTC indicator. The series are aligned according to the aligned longest common sub-sequence.

2 Clustering Method

3 The primary goal of this work is to compare road user interactions, characterized by a set of
 4 continuous safety indicators. Choosing a data representation and a similarity method rules out
 5 some clustering methods. The choice of keeping the indicators in their original shape with variable
 6 lengths rules out for example classical clustering algorithms such as k-means since the concept of
 7 a centroid is not defined. All the clustering algorithms that operate on a similarity matrix could
 8 be used. Several were investigated: spectral clustering was in particular tested at length and used
 9 in a first version of this work (35). The method is fast and takes as only input the predetermined
 10 number of groups. Finding the number of clusters by trial and error proved to be a challenge, and
 11 the resulting clusters were not always easy to interpret.

12 The algorithm used for the results presented in this paper is a slight variation of the al-
 13 gorithm previously developed to cluster motion patterns (17). This type of algorithm trades the
 14 parameter of the number of clusters for a maximum distance or minimum similarity between in-
 15 stances of the same cluster: when a new instance is too different from the existing clusters, a new
 16 one is created for it. The other idea is to use the original data as representatives, or prototypes, for
 17 each cluster. That provides a visual and more interpretable representation of each cluster. The last
 18 idea is to favour “long” instances, in this case indicators with long time periods of observation.
 19 This is done in two ways: first by sorting the indicators according to their length, to start consid-
 20 ering first the longest indicators, and second by keeping the longer prototype indicator when two
 21 cluster are merged. This solves partially the problem of dependency of the results to the algorithm

1 initialization which is a well-known limitation of many clustering algorithms such as k-means
 2 (initialization of the cluster centroids) and of the one proposed in (17). The algorithm parameter
 3 is therefore the minimum similarity for two indicators to be in the same cluster: when learning
 4 prototypes, an indicator will be added as a new prototype if its maximum similarity to all existing
 5 prototypes is lower than the parameter.

6 EXPERIMENTAL RESULTS

7 The proposed method to cluster interactions and their safety indicators is tested on a large dataset
 8 of 295 traffic videos of collisions and conflicts between motor vehicles collected at an intersection
 9 in Kentucky. This unique dataset has already been used in past studies (2), most notably in the
 10 first paper that compared the characteristics of interactions with and without a collision (9). The
 11 definition of conflicts used by the people who collected and sorted the data is unknown: a visual
 12 review confirms that most match the accepted definition, but will be referred to as interactions
 13 without a collision.

14 As shown in (2), an interaction can be well described by several symmetrical indicators
 15 based on speed and positions independently of the road user absolute positions. The indicators
 16 characterize one road user’s motion with respect to the other, as if it was stationary. These indi-
 17 cators based on positions and speed are: the distance of the road users’ centroids, the minimum
 18 distance separating the road users (from the feature-based tracking algorithm), the speed differen-
 19 tial (the norm of the velocity difference), the angle of the velocities, and the collision course angle
 20 (the angle between the velocity difference and the vector the links the road users’ centroids). To
 21 these indicators are added two safety indicators, TTC and probability of collision, calculated in
 22 (2) using motion prediction methods based on prototype trajectories representing the main motion
 23 patterns.

Indicator	Threshold ε	Minimum Clustering Similarity	Number of Clusters
Distance (Dist)	1 m	0.3	6
Speed differential (SD)	1.5 m/s	0.4	4
Velocity angle (VA)	0.15 rad	0.4	4
Time to collision (TTC)	0.2 s	0.3	4
Probability of Collision (PoC)	0.1	0.5	6

TABLE 1 Thresholds ε for d_ε used in the computation of the aligned normalized similarity *ALCSS* with $\delta = 2$, with the minimum similarity used for clustering and the resulting number of clusters.

24 The choice is made for this study to cluster the interactions based on each indicator sep-
 25 arately, for the following ones: distance, speed differential, velocity angle, TTC and probability
 26 of collision. For each indicator, a threshold is chosen by trial and error to match the profiles us-
 27 ing the aligned normalized similarity *ALCSS* with $\delta = 2$. The matching function is d_ε with the
 28 thresholds ε listed in TABLE 1. An additional criteria is added to remove very short indicators that
 29 do not contain much information (if not favouring longer indicators in the clustering algorithm,
 30 the shortest indicators would tend to be the most similar to the others as they can easily match at
 31 least some sub-sequence of a long indicator). The minimum length is 10 frames, i.e. 0.67 s, and

1 actually applies only to safety indicators since they may not be computed for all instants. The
 2 others can be computed as long as the two road users co-exist in the scene. The software code
 3 used to compute the similarities, the clustering algorithm and the results presented in this paper are
 4 available in the open source Traffic Intelligence project (6) and on the page dedicated to this paper
 5 (<http://nicolas.saunier.confins.net/data/saunier14trb.html>).

6 The choice is also made to not display and analyze clusters with too few indicators. The
 7 minimum number in the following results is 5 instances, including the prototype. Different min-
 8 imum similarities for clustering were tested by trial and error for the different indicators and are
 9 listed in TABLE 1. For all figures from 4 to 8, each cluster prototype is plotted using dots. Interac-
 10 tions with and without a collision are displayed respectively in red and blue. The numbers beside
 11 each cluster number are in order: the percentage of collisions, the number of collisions and the
 12 number of indicators in the cluster.

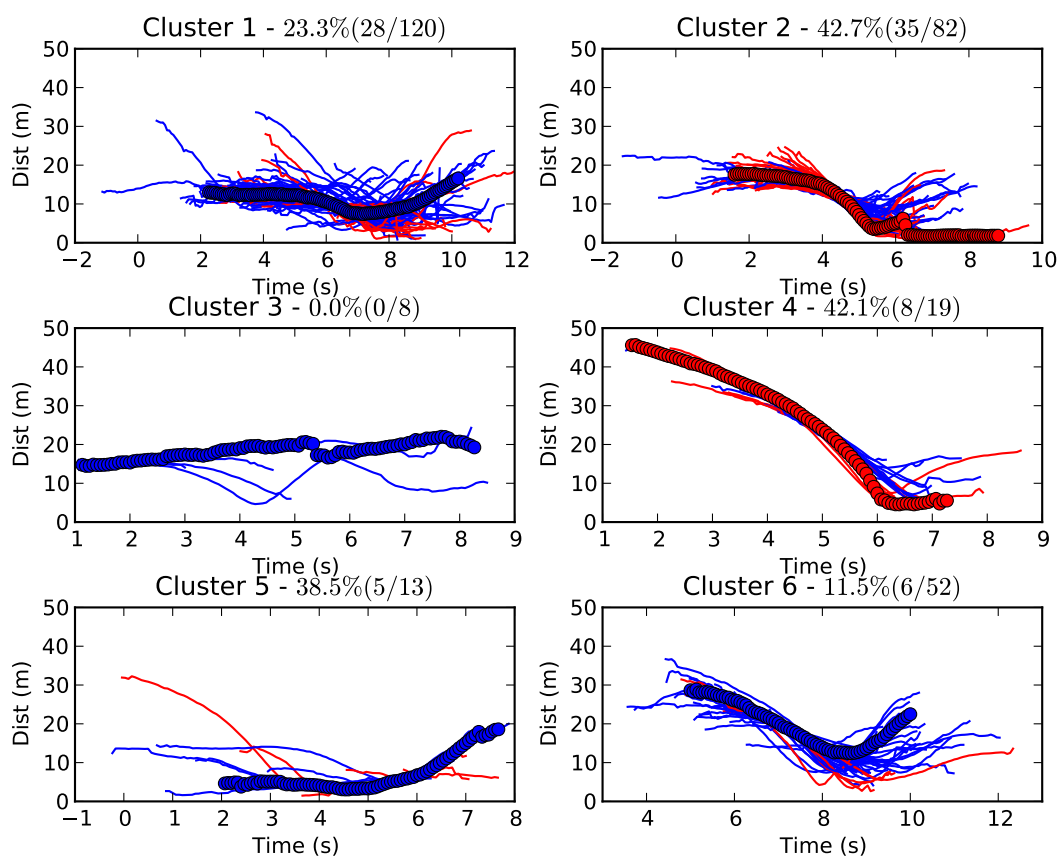


FIGURE 4 Clusters of the distance indicators.

13 The 6 clusters of distance indicators are plotted in FIGURE 4. There are quite different
 14 profiles, from almost flat in cluster 3 to increasing in cluster 5 to decreasing then flat (clusters 2
 15 and 4) or increasing again (clusters 1 and 6). These clusters correspond to varying proportions
 16 of collisions: the clusters 2 and 4 contain the most collisions and have therefore expected shapes
 17 where the distance remains 0 or close to 0 after the collision. It is however remarkable that a

1 majority of the interactions in these clusters do not end up in a collision. Clusters 1 and 6 seem
 2 to correspond to some sort of evasive action since the distance decreases (the road users are on a
 3 collision course), then increases again once the road users start reacting to avoid the collision. The
 4 collisions in cluster 5 may correspond to situations where the road users continue moving after the
 5 shock. The rate of change differs considerably between the clusters.

6 The 4 clusters of SD indicators are plotted in FIGURE 5. The shapes are quite distinctive
 7 and seem relatively homogeneous for each cluster. There is a pattern relating the proportion of
 8 collisions to the highest speed differential: the higher the proportion, the higher the maximum
 9 speed differential. This is related to attempts by road users to avoid the collision, which are stronger
 10 in collisions. The shape of cluster 1 is particularly striking and could be related to rear end or
 11 parallel interactions at similar velocities, followed by a road user turning or changing lane, which
 12 puts the road users on a collision course, followed by a return to the initial conditions or more
 13 evasive actions with higher speed differential.

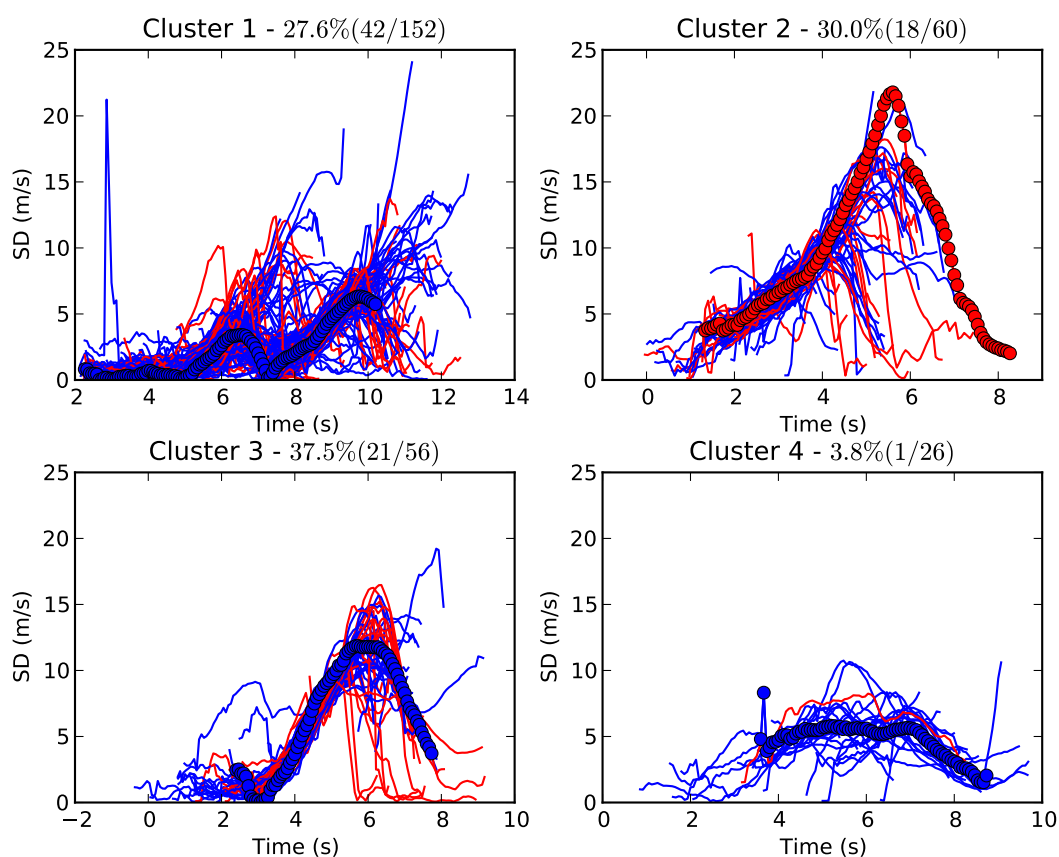


FIGURE 5 Clusters of the SD indicators.

14 The 4 clusters of VA indicators are plotted in FIGURE 6. The VA indicator is very useful
 15 to identify interaction categories, e.g. rear-end, side, head-on, etc., which does change with each
 16 instant and is typically not recorded in collision reports. Interactions in cluster 4 are thus side
 17 interactions, which evolve as the road users try to avoid each other. Cluster 1 contains even more

1 collisions and corresponds to rear-end and parallel interactions where the VA increases (some
 2 attempt at turning or changing lane) then comes back to 0. The clusters 2 and 3 are more difficult
 3 to distinguish. The two must contain situations that start as parallel or rear-end interactions, but
 4 evolve into side interactions, with different angles. One cannot miss in any case some important
 5 differences in profiles, especially at the beginning in clusters 1 and 2. Trying to obtain more
 6 clusters may yield a finer understanding of these clusters.

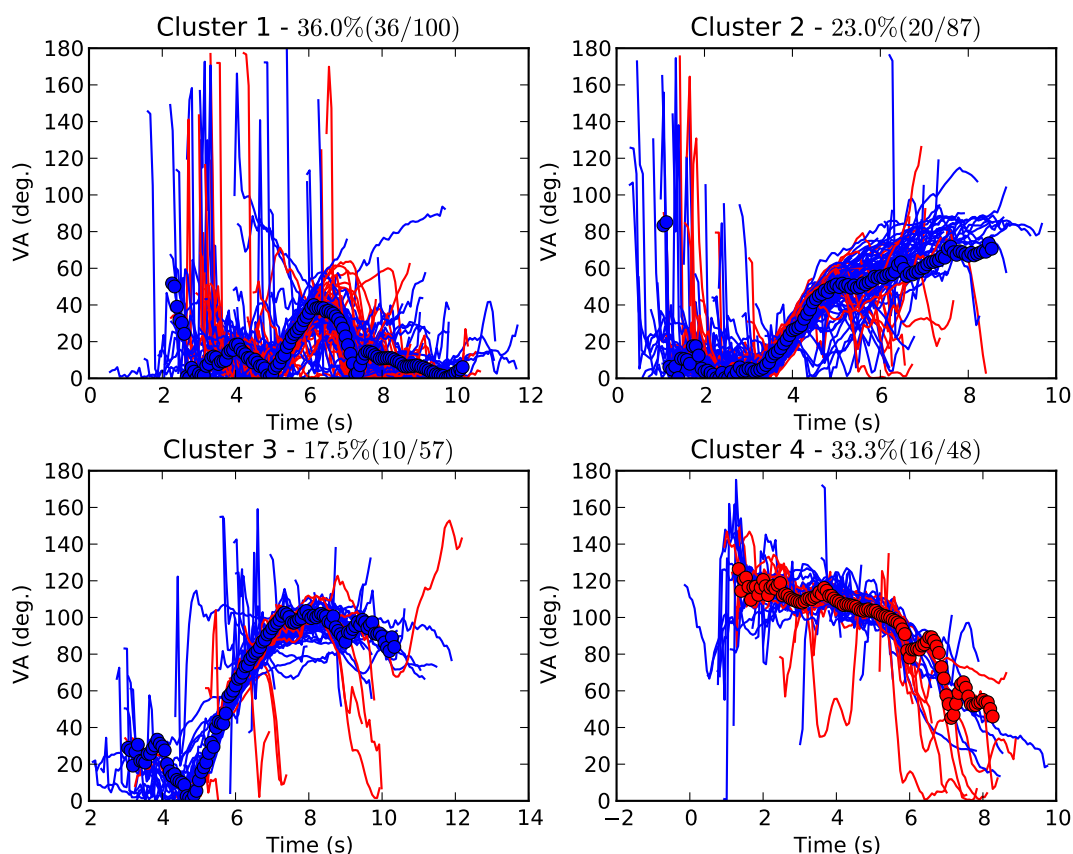


FIGURE 6 Clusters of the VA indicators.

7 The 4 clusters of TTC indicators are plotted in FIGURE 7. It must be noted that there are
 8 only 247 interactions for which TTC can be computed for at least 10 frames. One cannot miss
 9 that the TTC indicators are much noisier than the previous indicators. This is related to the quality
 10 of the data and the more complex process of computing TTC. Although the motion patterns allow
 11 to compute the TTC at more instants, more trajectory prototypes could have made the measures
 12 smoother. Most TTC profiles decrease with time as they are expected for collisions and conflicts.
 13 The clusters 1 and 2 are very interesting because they look similar at first sight. But the proportion
 14 of collisions in cluster 2 is consistent with the profile of its prototype indicator which falls at a
 15 seemingly constant rate as a function of time and reaches almost 0 s. On the contrary, there are
 16 few TTC measures below 0.5 s or even 1 s in cluster 1. There is more variability in the rate of
 17 decrease at the beginning and most profiles increase again after reaching their minimum, which

1 is consistent with a high proportion of interactions without a collision. Cluster 3 contains mostly collisions, with a higher rate of decrease than cluster 2 which explains why they are in different clusters. Finally, cluster 4 contains only one collision and has fairly constant, noisy, TTC values above 1.5 s.

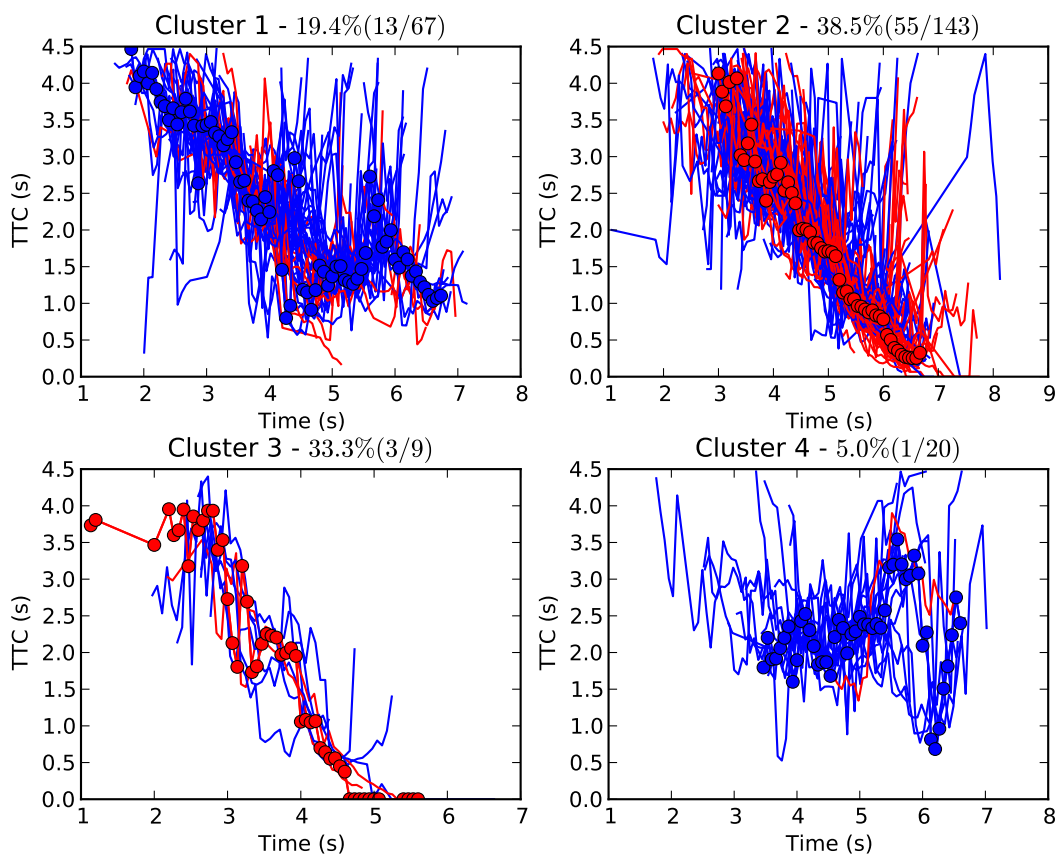


FIGURE 7 Clusters of the TTC indicators.

5 The 6 clusters of PoC indicators are plotted in FIGURE 8. It must be noted that there are
6 only 260 interactions for which PoC can be computed for at least 10 frames. PoC is also a noisy
7 indicator, depending as TTC on motion prediction methods and the existence of potential collision
8 points. There are two main clusters, 1 and 2, and 4 smaller ones. If ranking the cluster from their
9 maximum PoC, it goes from cluster 3, to cluster 5, then 4 and 6, and finally 2: the first 4 have few
10 collisions, while the last one, cluster 2 has the highest proportion and the highest maximum values
11 reaching 0.8, which is consistent. On the other hand, cluster 1 is more difficult to interpret: it is
12 the largest cluster, contains mostly interactions without a collision, and seems to have two peaks.
13 Whether this is related to noisier interactions without a collision or an actual variation of PoC is
14 unclear.

15 To sum up the observations, the methods could produce varying numbers of clusters for
16 each indicator that can be interpreted. There are two main results for all indicators. First, there
17 are clusters with very few collisions, e.g. cluster 4 for the SD indicator and cluster 4 for the TTC

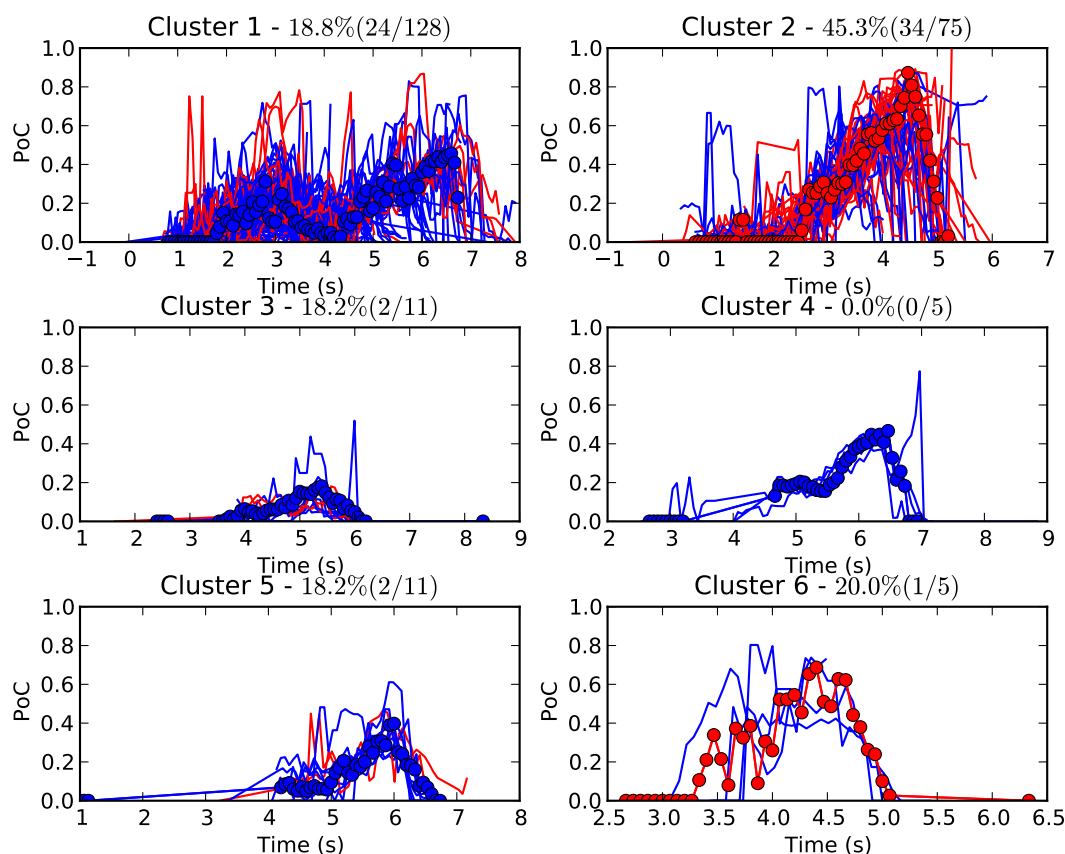


FIGURE 8 Clusters of the PoC indicators.

1 indicator, which, as was noticed in the previous study (9), seems to indicate that some interactions
 2 without a collision are not similar to any collisions. This suggests therefore that the factors that are
 3 associated with these interactions are different from the ones associated with collisions. Second,
 4 it is also clear that even in the clusters with the highest share of collisions (45.3 % for cluster 2 of
 5 the PoC indicator), there is always a majority of interactions without a collision that have similar
 6 processes to the collisions and are therefore good candidate predictors of these collisions. Finally,
 7 there is also a clear trade-off between having few clusters with some degree of variability, which
 8 can be seen in several clusters such as cluster 1 for the SD indicator, and more numerous and more
 9 homogeneous clusters. This choice is up to the analyst and highlights the flexibility of the method
 10 as an exploratory tool.

11 CONCLUSION

12 This paper has introduced a new similarity measure built upon the longest common sub-sequence
 13 that is sensitive to the rate of change of time series and is shown to be adapted to the clustering
 14 of several interaction indicators, including safety indicators such as TTC and the probability of
 15 collision. The number of resulting clusters is relatively small and can be easily interpreted in most
 16 cases. The results yield further credibility to the main hypothesis of surrogate safety analysis that

1 some interactions without a collision have similar processes as collisions and could be used as
2 predictors. It also strengthens the observations made in (9) that not all interactions should be used
3 for surrogate safety analysis, as can be seen for almost each indicator. Another contribution of
4 this work is to release all the necessary code and data samples to allow true reproducibility of the
5 presented work.

6 There is considerable room for further research. The main goal is to cluster interactions
7 with and without a collision based on all, or at least several, indicators simultaneously, i.e. by
8 evaluating the similarity of interactions at a given instants through the similarity of all its indicators
9 at this instant. It is hoped that very strong similarities can thus be identified. Finally, this method
10 can be applied to other time series in transportation, especially in safety, such as the large datasets
11 produced by current naturalistic driving studies.

12 **ACKNOWLEDGEMENT**

13 The authors wish to acknowledge the financial support of the Natural Sciences and Research Coun-
14 cil of Canada (NSERC). They also wish to thank Zu Kim of California PATH and Ann Stansel of
15 the Kentucky Transportation Cabinet for providing the video dataset.

16 **REFERENCES**

- 17 [1] Ismail, K., *Application of computer vision techniques for automated road safety analysis and*
18 *traffic data collection*. Ph.D. thesis, University of British Columbia, 2010.
- 19 [2] Saunier, N., T. Sayed, and K. Ismail, Large Scale Automated Analysis of Vehicle Interactions
20 and Collisions. *Transportation Research Record: Journal of the Transportation Research*
21 *Board*, Vol. 2147, 2010, pp. 42–50, presented at the 2010 Transportation Research Board
22 Annual Meeting.
- 23 [3] Bagdadi, O. and A. Várhelyi, Jerky driving-An indicator of accident proneness? *Accident*
24 *Analysis & Prevention*, Vol. 43, No. 4, 2011, pp. 1359–1363.
- 25 [4] Hallmark, S., D. Mce, K. M. Bauer, J. M. Hutton, G. A. Davis, J. Hourdos, I. Chatterjee,
26 T. Victor, J. Bårgman, M. Dozza, H. Rootzén, J. Lee, C. Ahlström, O. Bagdadi, J. En-
27 gström, D. Zholud, and M. Ljung-Aust, *Initial Analyses from the SHRP 2 Naturalistic Driv-*
28 *ing Study: Addressing Driver Performance and Behavior in Traffic Safety*. Transportation
29 Research Board, 2013.
- 30 [5] Saunier, N. and T. Sayed, A feature-based tracking algorithm for vehicles in intersections. In
31 *Canadian Conference on Computer and Robot Vision*, IEEE, Québec, 2006.
- 32 [6] Saunier, N., *Traffic Intelligence*. [https://bitbucket.org/Nicolas/](https://bitbucket.org/Nicolas/trafficintelligence)
33 [trafficintelligence](https://bitbucket.org/Nicolas/trafficintelligence), 2011-2013, software under the open source MIT License.
- 34 [7] *The Swedish Traffic Conflict Technique*. Sweden, 2005.
- 35 [8] Lareshyn, A., Å. Svensson, and C. Hydén, Evaluation of traffic safety, based on micro-
36 level behavioural data: Theoretical framework and first implementation. *Accident Analysis &*
37 *Prevention*, Vol. 42, No. 6, 2010, pp. 1637–1646.

- 1 [9] Saunier, N., N. Mourji, and B. Agard, Investigating Collision Factors by Mining Microscopic
2 Data of Vehicle Conflicts and Collisions. *Transportation Research Record: Journal of the*
3 *Transportation Research Board*, Vol. 2237, 2011, pp. 41–50, presented at the 2011 Trans-
4 portation Research Board Annual Meeting.
- 5 [10] Svensson, A. and C. Hydén, Estimating the severity of safety related behaviour. *Accident*
6 *Analysis & Prevention*, Vol. 38, No. 2, 2006, pp. 379–385.
- 7 [11] Vlachos, M., G. Kollios, and D. Gunopulos, Elastic Translation Invariant Matching of Tra-
8 jectories. *Machine Learning*, Vol. 58, No. 2-3, 2005, pp. 301–334.
- 9 [12] Svensson, A., *A Method for Analyzing the Traffic Process in a Safety Perspective*. Ph.D.
10 thesis, University of Lund, 1998, bulletin 166.
- 11 [13] Archer, J., *Methods for the Assessment and Prediction of Traffic Safety at Urban Intersections*
12 *and their Application in Micro-simulation Modelling*. Academic thesis, Royal Institute of
13 Technology, Stockholm, Sweden, 2004.
- 14 [14] Laureshyn, A., *Application of automated video analysis to road user behaviour*. Ph.D. thesis,
15 Lund University, 2010.
- 16 [15] Tarko, A., G. A. Davis, N. Saunier, T. Sayed, and S. Washington, *Surrogate Measures of*
17 *Safety*. White paper, ANB20(3) Subcommittee on Surrogate Measures of Safety, 2009.
- 18 [16] Amundsen, F. and C. Hydén (eds.), *Proceedings of the first workshop on traffic conflicts*,
19 Institute of Transport Economics, Oslo, Norway, 1977.
- 20 [17] Saunier, N., T. Sayed, and C. Lim, Probabilistic Collision Prediction for Vision-Based Au-
21 tomated Road Safety Analysis. In *The 10th International IEEE Conference on Intelligent*
22 *Transportation Systems*, IEEE, Seattle, 2007, pp. 872–878.
- 23 [18] Mohamed, M. G. and N. Saunier, Motion Prediction Methods for Surrogate Safety Analysis.
24 In *Transportation Research Board Annual Meeting Compendium of Papers*, 2013, 13-4647.
25 Accepted for publication in *Transportation Research Record: Journal of the Transportation*
26 *Research Board*.
- 27 [19] Hydén, C., *The development of a method for traffic safety evaluation: The Swedish Traf-*
28 *fic Conflicts Technique*. Ph.D. thesis, Lund University of Technology, Lund, Sweden, 1987,
29 bulletin 70.
- 30 [20] Gettman, D., L. Pu, T. Sayed, and S. Shelby, *Surrogate Safety Assessment Model and Valida-*
31 *tion: Final Report*. FHWA, 2008.
- 32 [21] Minderhoud, M. M. and P. H. Bovy, Extended time-to-collision measures for road traffic
33 safety assessment. *Accident Analysis & Prevention*, Vol. 33, No. 1, 2001, pp. 89–97.
- 34 [22] Parkhurst, D., *Using Digital Video Analysis to Monitor Driver Behavior at Intersections*.
35 Center for Transportation Research and Education (CTRE), Iowa State University, 2006.

- 1 [23] Laureshyn, A., K. Åström, and K. Brundell-Freij, From speed profile data to analysis of
2 behaviour. *IATSS Research*, Vol. 33, No. 2, 2009, pp. 88–98.
- 3 [24] Davis, G. A., J. Hourdos, and H. Xiong, Outline of Causal Theory of Traffic Conflicts and
4 Collisions. In *Transportation Research Board Annual Meeting Compendium of Papers*, 2008,
5 08-2431.
- 6 [25] Alpaydin, E., *Introduction to Machine Learning, second edition*. The MIT Press, second
7 edition ed., 2010.
- 8 [26] Vlachos, M., G. Kollios, and D. Gunopulos, Discovering Similar Multidimensional Trajecto-
9 ries. In *Proc. of 18th International Conference on Data Engineering (ICDE)*, San Jose, CA,
10 2002, pp. 673–684.
- 11 [27] Liao, T. W., Clustering of time series data: a survey. *Pattern Recognition*, Vol. 38, No. 11,
12 2005, pp. 1857–1874.
- 13 [28] Vlachos, M., J. Lin, E. Keogh, and D. Gunopulos, A Wavelet-Based Anytime Algorithm for
14 K-Means Clustering of Time Series. In *Proc. Workshop on Clustering High Dimensionality
15 Data and Its Applications*, 2003, pp. 23–30.
- 16 [29] Morris, B. and M. Trivedi, Learning trajectory patterns by clustering: Experimental studies
17 and comparative evaluation. In *Proceedings of the IEEE International Conference on Com-
18 puter Vision and Pattern Recognition (CVPR)*, 2009, pp. 312–319.
- 19 [30] Piciarelli, C. and G. Foresti, On-line trajectory clustering for anomalous events detection.
20 *Pattern Recognition Letters*, Vol. 27, No. 15, 2006, pp. 1835–1842.
- 21 [31] Bennewitz, M., W. Burgard, G. Cielniak, and S. Thrun, Learning Motion Patterns of People
22 for Compliant Robot Motion. *The International Journal of Robotics Research*, Vol. 24, No. 1,
23 2005, pp. 31–48.
- 24 [32] Hu, W., X. Xiao, D. Xie, T. Tan, and S. Maybank, Traffic Accident Prediction using 3D
25 Model Based Vehicle Tracking. *IEEE Transactions on Vehicular Technology*, Vol. 53, No. 3,
26 2004, pp. 677–694.
- 27 [33] Morris, B. and M. Trivedi, Trajectory Learning for Activity Understanding: Unsupervised,
28 Multilevel, and Long-Term Adaptive Approach. *IEEE Transactions on Pattern Recognition
29 and Machine Intelligence*, Vol. 33, No. 11, 2011, pp. 2287–2301.
- 30 [34] Zelnik-Manor, L. and P. Perona, Self-Tuning Spectral Clustering. In *Advances in Neural
31 Information Processing Systems 17* (L. K. Saul, Y. Weiss, and L. Bottou, eds.), MIT Press,
32 Cambridge, MA, 2005, pp. 1601–1608.
- 33 [35] Mohamed, M. G. and N. Saunier, Classifying Profiles of Surrogate Safety Measures to Under-
34 stand Collision Processes. In *Canadian Multidisciplinary Road Safety Conference*, Montreal,
35 2013.