

1 Mouse tracking as a window into decision making

2 Mora Maldonado, Ewan Dunbar & Emmanuel Chemla

3 January 14, 2018

4 1 Introduction

5 In the past ten years, mouse-tracking has become a popular method to target
6 the processes underlying decision making in different domains, ranging from
7 phonetic competition (Spivey, Grosjean, & Knoblich, 2005), and syntactic and
8 pragmatic processing (Farmer, Anderson, & Spivey, 2007; Tomlinson, Bailey, &
9 Bott, 2013) to social cognition (Freeman & Ambady, 2010; Freeman, Dale, &
10 Farmer, 2011). All these studies have worked on the assumption that motor
11 responses are prepared in parallel to cognitive processing and performed in a
12 cascade manner (i.e., as fast as they can be executed) (Song & Nakayama, 2006;
13 Freeman & Ambady, 2010, Spivey & Dale, 2006). As a result, features in mouse
14 trajectories could be indicators of specific decision processes, revealing their dy-
15 namics with fine-grained temporal resolution (Freeman et al., 2011; Freeman
16 & Ambady, 2010; Hehman, Stoller, & Freeman, 2014)¹. Mouse-tracking stud-
17 ies typically present participants with a *two-alternative forced choice*, where
18 they have to make a choice using the options appearing in the top left or right
19 corner of the screen. Whenever a decision involves two independent processes
20 –such as a change of mind–, mouse trajectories are expected to be displayed as
21 two movements, whereas a single smooth and graded movement would reflect a
22 commitment with an initial choice (see ??, Wojnowicz et al., 2009). Of course,
23 one could still imagine other types of decision, such as a single but late com-
24 mitment, or a decision made after uncertainty or doubt. These might have a
25 different reflection on mouse trajectories.

26 The existence of intuitions about how decision processes should be mapped
27 into mouse paths (i.e. linking hypotheses) has allowed researches to draw con-
28 clusions about the cognitive processes underlying their experimental manipula-
29 tions. Dale and Duran’s (2011) approach to negation processing is an example
30 of this. Negation has been traditionally understood as an operator that reverses
31 the truth conditions of the sentence, inducing in “extra-step, or mental oper-
32 ation” in online processing (Wason & Johnson-Laird, 1972, REF). To test the
33 dynamics of negation integration, Dale and Duran tracked mouse trajectories

¹In this sense, mouse trajectories are equivalent to eye movements. However, the main advantage of mouse tracking over eye tracking is the simplicity of the set up, which can be even tested online.

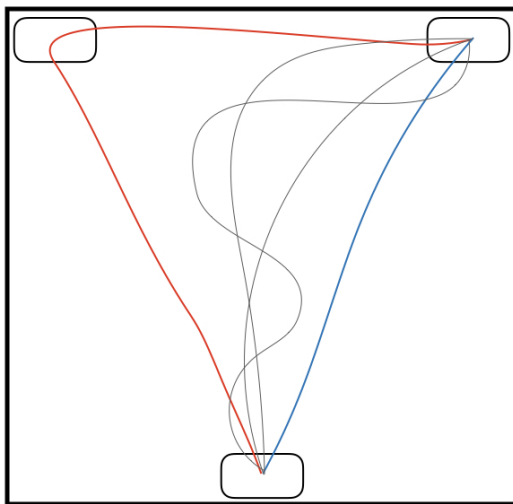


Figure 1: Shape of trajectories underlying distinct decision processes

34 as participants performed a Truth-Value Judgment Task (TVJT), where they
 35 had to verify the truth of general statements such as *Cars have (no) wings*.
 36 The authors found that mouse trajectories presented more shifts towards the
 37 alternative response (i.e. they were less straightforward) when evaluating nega-
 38 tive than affirmative true sentences. These results were interpreted as evidence
 39 for a ‘two-step’ processing of negation, where truth conditions for the positive
 40 content are first derived and negated only afterwards, as a second step².

41 As observed, the impact of experimental manipulations in the shape of tra-
 42 jectories has been taken as evidence for underlying decision patterns. This
 43 association, however, has never been explicitly tested. While no one can deny
 44 that mouse trajectories are sensitive to experimental manipulations –such as
 45 negation–, it is unclear whether this can be directly interpreted as reflecting
 46 decision making.

47 Our main goal here is to test the connection between cognition (decision mak-

²Several studies have suggested that the positive argument might play an important role in negation processing (REF). This pattern of results, however, seem to depend on the amount of contextual support given for the sentence: ‘two-step’ negation processing seems to occur specifically for sentences presented out of the blue, whereas no difference between negative and positive sentences arises when the right contextual support is provided. How to explain this pattern of results has been at the center of the debate in the negation processing literature (see Y. Tian and R. Breheny XX for review), but we will not explore it here.

ing) and action (mouse trajectories): Are decision processes always reflected in mouse trajectories? If yes, how are they reflected? One could easily imagine a situation where two different trajectory shapes do not correspond to two different decision mechanisms, but to a single one (due to uncertainty, noise, etc.). Differences in trajectories can therefore underlie something other than decision shift.

In this paper, we address these questions by identifying the features in mouse-trajectories corresponding to two different types of decisions: *straightforward* decisions (i.e. single commitment) and *shifted* decisions (i.e. a change of mind). First, we present a *validation* experiment where, instead of taking mouse trajectories as indicators of cognitive processes, we *manipulate* whether or not our stimuli trigger a flip in decision (Section 1). The data from this validation experiment (i.e. two groups of ‘quasi-decisions’) is then used to feed a *linear discriminant analysis* (henceforth, LDA), trained to classify trajectories depending on the underlying decision (Section 2). After comparing the performance of LDA classifier to other traditionally used mouse-tracking measures (Section 3), the LDA classifier will be further tested with new data, obtained from a replication of Dale and Duran’s (2011) experiment on negation processing. If there is a change of decision triggered by negation, trajectories corresponding to negative trials should be classified together with (as) trajectories underlying decision change in the validation experiment.

2 Manipulating decision making: Validation Experiment

We developed an experiment where participants had to perform a *two-alternatives forced task*: at each trial, they were presented with a coloured frame surrounding the screen and they had to determinate whether the frame was blue or red. Responses were made by clicking on the “blue” or “red” buttons, allowing the recording of mouse-movements during each trial. Importantly, responses were considered accurate if they described the color at the moment of the click. In order to mirror decision processing, we manipulated whether the color of the frame remained stable or changed at some point during the trial. While in the latter case the initial choice will be the accurate response (*straightforward* trials), in the former, participants were forced to swap their answer (*switched* trials), mimicking a change of decision. Note that, since we are only *mirroring* decision making, we refer to these decision processes as ‘*quasi-decisions*’. An illustration of the procedure is provided in [Figure 2](#).

Participants Fifty four participants (F=27) were recruited using Amazon Mechanical Turk. Two subjects were excluded from the analyses because they did not use a mouse to perform the experiment. All of them were compensated with 0.5 USD for their participation, which required approximately 5 minutes.

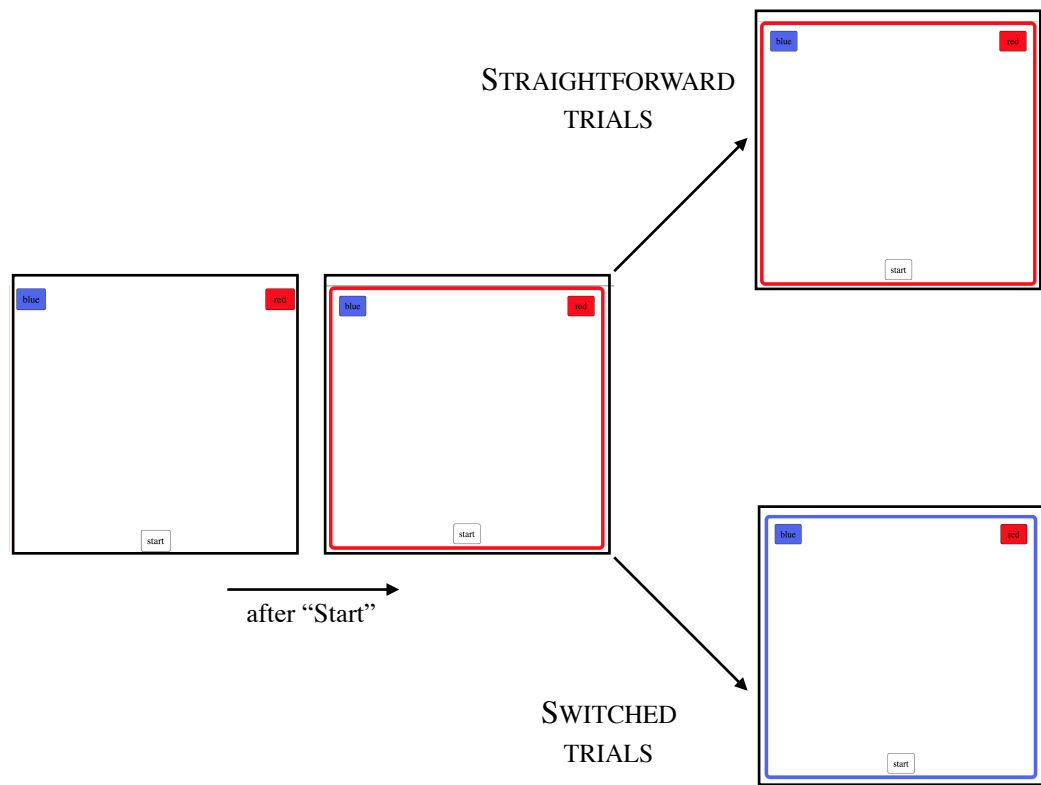


Figure 2: Procedure in Validation Experiment

88 **Design** Each trial instantiated one of two possible DECISION PATTERNS. In
89 *straightforward* trials, the frame color remained stable, and the decision made
90 at the beginning of the trial did not need to be revised. Conversely, in *switched*
91 trials, the color swapped during the trial, forcing a revision of the initial choice.
92 The POINT OF CHANGE in *switched* trials was determined by the position in the
93 y axis, and it could be early, middle or late. The FRAME COLOR at the response
94 point was also controlled: it could be red (right button) or blue (left button).
95 A summary of the design is given in Table 1.

DECISION TYPE	FRAME COLOR		DECISION POINT
Straightforward	Blue		—
	Red		—
Switched	Blue	Red	early ($y=.4$), middle ($y=.7$), late ($y=.9$)
	Red	Blue	early ($y=.4$), middle ($y=.7$), late ($y=.9$)

Table 1: Design in Calibration Experiment

96 To prevent participants from developing a strategy based on staying on the
97 middle of the screen, the proportion of trials was adjusted so that straightfor-
98 ward trials were the majority (32 repetitions per frame color), whereas switched
99 trials had 4 repetitions per frame color and change point. The total number of
100 trials was 88.

101 **Interface** The interface was programmed using JavaScript. Mouse move-
102 ments triggered the extraction of x, y -pixel coordinates (i.e., no constant sample
103 rate). The software was adapted proportionally to the window of the partici-
104 pant’s browser, forming a rectangle: the height was covered at 100 percent and
105 the width was 120 percent of the height. Three buttons were displayed during
106 the experiment (‘start’ and response buttons). Their size was also determined
107 by the browser window (i.e., approximately 20 percent of the total width). The
108 ‘start’ button was placed at the bottom center of the screen. The two re-
109 sponse boxes were located at the top left (‘blue’) and top right (‘red’) corners
110 of window. This location was constant across participants, and handedness was
111 controlled. In each trial, mouse movements were recorded between start-clicks
112 and response-clicks. The x, y -pixel trajectory was saved together with its raw
113 time. Afterwards, the positions were normalised according to participants’ win-
114 dow size, to allow comparisons between subjects. The normalization was done
115 by considering the start button at the $[0,0]$ point, the “blue” button corner at
116 $[-1,1]$ and the “red” button at $[1,1]$.

117 **Data treatment** Mouse-tracking data are particularly variable trial to trial.
118 On one hand, variations in response times imply different quantity of x, y posi-
119 tions per trial, making difficult the comparisons between items. On the other
120 hand, in our design, positions are extracted based on mouse movements, and de-
121 vices with more or less sensibility could influence the number of samples taken
122 during the trial. In order to compare mouse trajectories, we normalized the

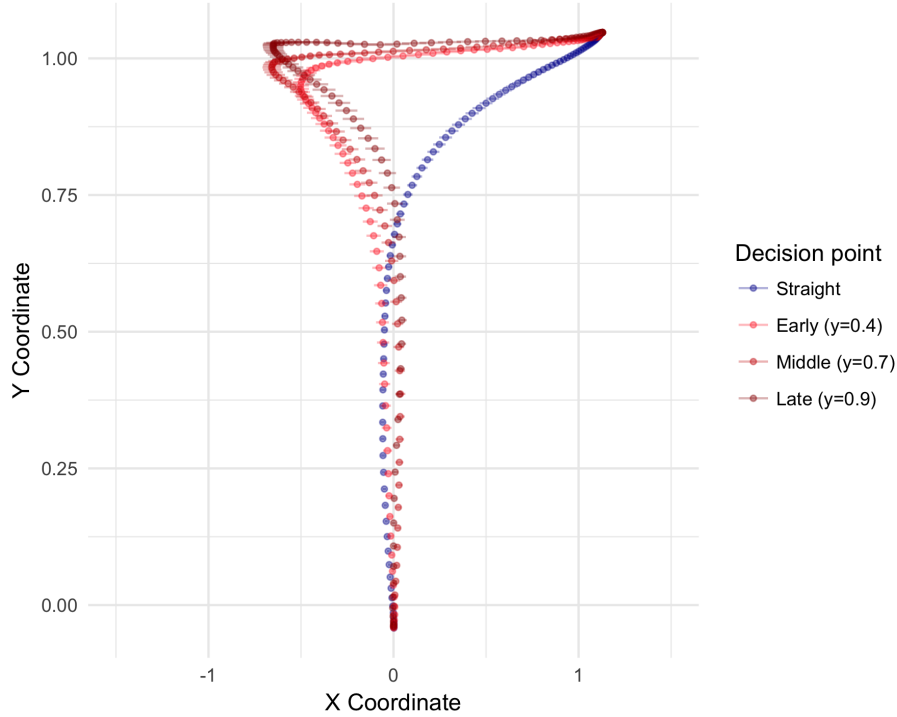


Figure 3: Mean trajectories per class

time course into 101 proportional times steps (percentage of trial duration). This normalization, as all the other calculations, was performed in the Spyder environment using Python 2.7.

Overall performance Inaccurate responses, corresponding to 4% of the data, were removed from the analyses. Mean trajectories for each DECISION TYPE and DECISION POINT are illustrated in Figure 3. These trajectories suggest that participants made a decision as soon as they were presented with the colour frame, and revised this decision if needed. When they were forced to change their choice, this switch was reflected in mouse trajectories. **MM: Do we want to include other graphs?**

3 Classifying decision processes with LDA

Different ‘quasi-decisions’ (i.e. DECISION TYPES) have a different impact on mouse trajectories, as observed in Figure 3. To identify the features characteristic of each class (*switched* vs. *straightforward*), we use a Linear Discriminant

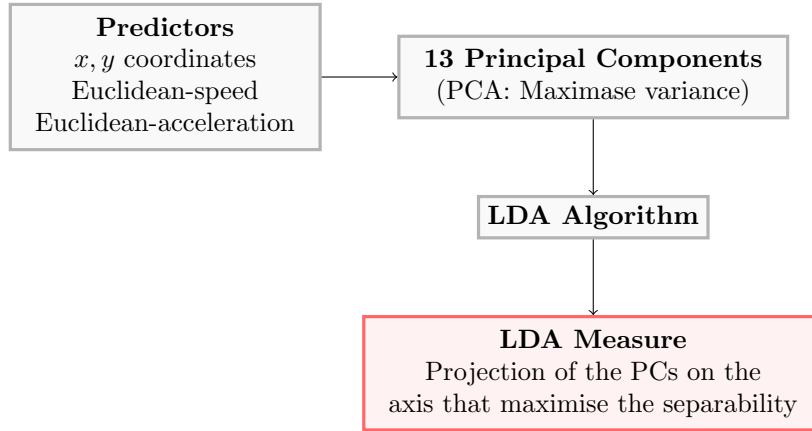


Figure 4: Diagram

137 Analysis method for classification.

138 **Description of LDA classifier** The LDA is an optimal solution to classify
 139 continuous data –such as trajectories– into two or more classes –such as decision
 140 patterns (straightforward vs. switched quasi-decisions). In a nutshell, the LDA
 141 algorithm assumes that different classes have a common covariation matrix,
 142 and finds the linear combination of predictors that gives maximum separation
 143 between the classes. This linear combination of predictors is obtained as a linear
 144 coefficient and it can be used to form a decision rule for the classification.

145 The predictors used by the classification algorithm were: (a) the x, y co-
 146 ordinates, (b) Euclidean-based velocity, and (c) Euclidean-based acceleration.
 147 While mouse coordinates give us absolute spatio-temporal information about
 148 where you are when, velocity and acceleration data provides information about
 149 how did we arrive there³. To avoid collinearity, we applied a Principle Compo-
 150 nent Analysis (PCA) to identify the 13 principal components on these predictors,
 151 and performed the LDA on these principal components. We obtained a *LDA*
 152 *measure* for each trial, which is basically the placement of the trajectory on
 153 the axis that maximise separability. A scheme of the procedure is provided in
 154 [Figure 4](#)

155 **Performance of the LDA classifier** The result of applying the procedure
 156 described in [Figure 4](#) to the trajectories in the validation experiment is il-
 157 lustrated in [Figure 5](#). To evaluate the overall performance of the classifier,
 158 we calculated the area under the ROC curve. The ROC curve diagnoses the
 159 classification ability as a function of the degree of sensitivity (percentage of

³It's worth noticing that none of these predictors can be obtained as a linear combina-
 tion of the others (although velocity and acceleration data can indeed be obtained by *some*
 combination of coordinates).

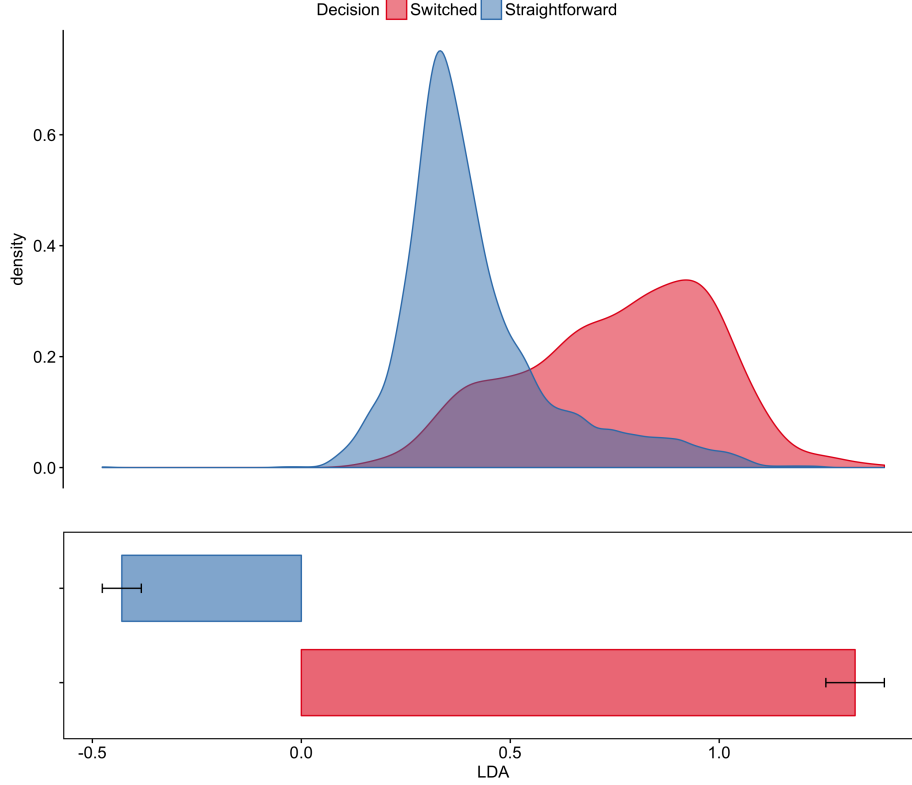


Figure 5: Classifier performance

change-of-mind trajectories correctly identified) and specificity (percentage of rapid-decision trajectories correctly identified) of the classifier. Given that an ideal ROC curve will hug the top left corner, the larger the AUC the better the classifier performance.

The cross-validation was performed as follows: Calibration data were partitioned into 10 bins that kept the proportion of *straightforward* and *switched* trajectories constant (75/25 proportion). For each bin, we took the complementary set of data (the rest 90%) to train the classifier. The data contained in the bin were used as test set to diagnose the classifier performance. In other words, we obtained a ROC curve and its AUC for each test bin (bins=10). The performance of the LDA classifier was compared to *baseline*, equivalent the worst possible outcome, and a *topline*, which was what we would expect from the classifier under the best possible conditions. For the baseline, we used a random classifier that assigned labels by sampling from a beta distribution centred at the probability of straightforward trials; the topline was computed by testing and training the original LDA classifier with the same set of data. The mean AUCs values for the LDA, the baseline and the topline in each bin are given in

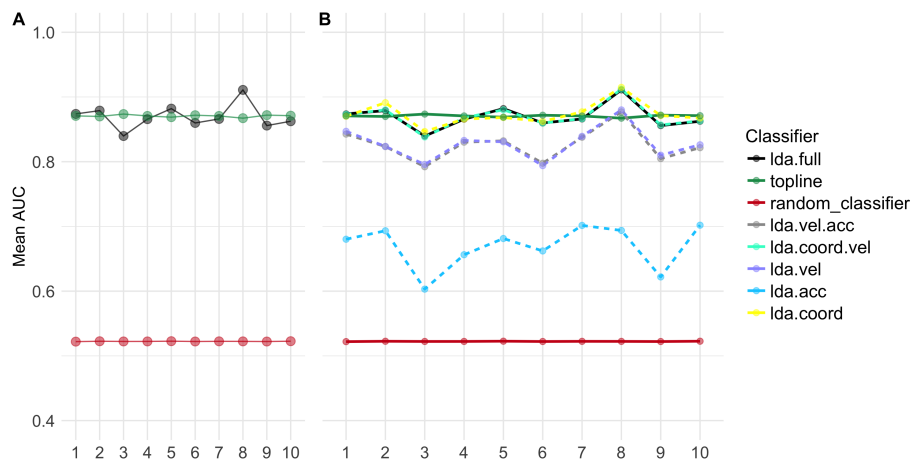


Figure 6: Mean AUC values

Figure 6a.

To assess whether the performance of the LDA classifier was statistically different from baseline and topline performances, we tested how likely was to obtain the attested differences if the null hypothesis was the case; namely, if the LDA classifier was as good as both baseline and topline. The difference in AUC means between each pair of classifiers was tested against a distribution of the possible mean-differences observable when the outcome is independent from the classifier. The sampling distribution under the null hypothesis was computed by a resampling the observed data ('permutation distribution').

In Table 2a, we report the results of performing a one-tail test on these mean differences (observed vs. under null hypothesis). As expected, our original LDA is significantly better than a random classifier at categorising trajectories into the two classes. Conversely, there is no significant difference between the performance of our LDA and the topline (i.e., the LDA is not significantly different from the best possible classification).

Meaningful features and optimal predictors Our original LDA classifier takes as predictors both absolute and relative spatio-temporal features (coordinates, speed and acceleration). While x, y coordinates provide information about the absolute position at each specific time step (*where* you are *when*), speed and acceleration contribute to knowing *how* one has arrived to a given position. Some of these features, however, might not be relevant for the classification. By comparing classifiers trained with different predictors, we can gather information about which features of mouse trajectories are more relevant for the distinction between decision processes (i.e. for the classification).

We trained five additional LDA classifiers obtained by subsetting the three original LDA predictors. If both absolute and relative features are required to

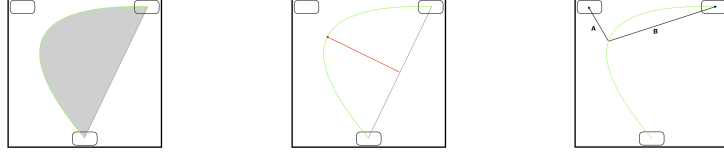
		(a)		(b)				
	ORIGINAL LDA (coords, speed, acc)	BASELINE	TOPLINE	LDA WITH DIFFERENT PREDICTORS				
				coords, vel	vel, acc	coords	vel	acc
AUC (mean)	.87	.52	.87	.87	.83	.87	.82	.67
Mean Difference	-	.35	-.002	-.0004	.04	-.006	.04	.2
<i>p</i> value	-	<.001	0.58	.5	<.001	.68	<.001	<.001

Table 2: Cross-validation results for the LDA classifier. The performance of the LDA was compared to the one of (a) Random and Topline classifiers and (b) LDA classifiers with different predictors.

predict the decision type, we would expect our “full” original LDA classifier to be better than any other classifier that takes only a subset of these original predictors. The performance of these additional classifiers was diagnosed in the same way as before, by computing the AUC for each of the 10 test bins. Figure 6b illustrates the mean AUCs values for each of these classifiers, together with the original LDA, the baseline and topline. Pair-wise comparisons with the original LDA were done by testing whether the observed mean differences would be expected under the null hypothesis (i.e. no difference in performance between classifiers). Table 2b summarises the comparisons between each of these classifiers and our original LDA.

The original LDA does not significantly differ from classifiers that contain the coordinates among their predictors, suggesting that the distinction between *straightforward* and *switched* ‘quasi-decisions’ might be solely explained by the information contained in the x, y coordinates. In contrast, the original LDA is significantly better than classifiers that use only speed and acceleration as predictors. These comparisons therefore reveal that, for classifying our validation data, absolute spatio-temporal features (x, y coordinates) are generally better predictors than relative features (speed and acceleration). That is to say, it seems to be more relevant to know where are you when than how you got there.

Note that this does not mean that changes in decision across the board do not have an impact on speed and acceleration. Indeed, it has been suggested that the speed and acceleration components can capture the level of commitment towards the response, such that a change of decision (*switched* trajectories) might have associated a specific speed/acceleration pattern (Hehman et al 2014). The fact that our data is not well captured by these components might be a specific property of the ‘quasi-decisions’ we are simulating, or generated by the fact that our data are too noisy to be able to see differences in speed and acceleration (relation with the fact that we have online data).



(a) Area Under the Curve (b) Maximal Deviation (c) Maximal LogRatio

$$\sum H[(x_t - x_{t-1})(x_{t-1} - x_{t-2})] \quad (\sum H[(a_t - a_{t-1})(a_{t-1} - a_{t-2})]) - 1$$

Figure 7: Traditional mouse-tracking measures.

4 LDA *versus* traditional mouse-tracking analyses

The LDA classifier is *a priori* the optimal solution to the type of discrimination problem examined here. However, when addressing the question of whether a change of decision has occurred, previous studies have used alternative techniques to analyse mouse trajectories. In what follows, we will compare the performance of our LDA classifier to the one of other measures commonly used in mouse tracking studies. We focus on measures that mainly assess the spatial disorder in trajectories, which is typically taken to be indicative of unpredictability and complexity in response dynamics (Hehlman et al 2014).

Two of the most commonly used methods of mouse tracking **spatial analysis** are the *Area Under the Curve* and the *Maximal deviation* (henceforth, AUC and MD respectively) (see Freeman & Ambady, 2010). The AUC is the geometric area between the observed mouse-trajectory and an idealised straight-line trajectory drawn from the start to the end points, whereas the MD is the point that maximises the perpendicular distance between this ideal trajectory and the observed path (Figure 7; REF). For both measures, higher values are associated with higher deviation peaks towards the alternative; values closer to zero (or below) suggest trajectory close to ideal. Another frequently used measure to estimate the complexity of the trajectory is based on quantifying the number of times the trajectories goes back and forth along the x-axis (horizontal flips, Dale and Duran 2011), as illustrated in Figure 7.

While all these measures serve to evaluate the degree of complexity of the path, they might fail to distinguish between ‘two-step’ and ‘uncertain’ decision processes –i.e., trajectories with a true deviation to the alternative or centred on the middle of the screen⁴. In order to assess more directly whether mouse trajec-

⁴For instance, a late medium-size deviation towards the alternative could underly a two-step decision whereas an early but big-size deviation towards the alternative might very well be considered just noise. However, measures such as the AUC might not be able to make a

257 stories have a meaningful deviation towards the alternative, the distance to both
 258 target and alternative responses should be taken into account. For instance, the
 259 *ratio of the target distance to the alternative distance* can be calculated for each
 260 x, y position. While ratio values closer to 1 suggest a position near the middle,
 261 higher values indicate a deviation towards the alternative response.

262 AUC, MD, x-coordinates flips and point that maximises the log-ratio (Max-
 263 imal LogRatio, henceforth) were calculated for the validation data. Following
 264 Dale and Duran (2011, and other studies on error corrections), we also anal-
 265 ysed the *acceleration component* (AC) as a function of the number of changes
 266 in acceleration (NB: This is not the same as computing the number of times
 267 the acceleration changes direction, going from positive to negative acceleration,
 268 as D&D claim). Since stronger competition between alternative responses is
 269 typically translated into steeper acceleration peaks, changes in acceleration can
 270 be interpreted as decision points (Helman et al 2014). Figure 8 illustrates the
 271 distribution and mean values for each ‘quasi-decision’.

272 The same cross-validation procedure described in the previous section was
 273 used to diagnose the performance of each of these measures. The mean AUC
 274 values for each of these measures are illustrated in Figure 9. Table 3 summarises
 275 the result of comparing the LDA performance to the one of each alternative
 276 measure.

	ORIGINAL LDA	AUC	MD	MAXIMAL LOGRATIO	X-COORD. FLIPS	AC
AUC (mean)	.87	.62	.81	.81	.73	.53
Mean Difference	-	.24	.06	.06	.14	.34
p value	-	<.001	<.001	<.001	<.001	<.001

Table 3: Cross-validation results for the LDA classifier. The performance of the LDA was compared to the one of five commonly used measures in mouse-tracking studies.

277 Overall, these comparisons reveal that the LDA is significantly better at clas-
 278 sifying validation data than other commonly used measures. While the difference
 279 with the classifier is in all the cases significant, mean AUCs values suggest that
 280 the MD and the Maximal LogRatio are better at distinguishing decision pro-
 281 cesses than other measures such as the AUC, the number of X-coordinates flips
 282 and the Acceleration Component. These two measures are the only ones cal-
 283 culated based on coordinates themselves and therefore give more importance
 284 to spatio-temporal information. In other words, both the MD and the max-
 285 imal log-ratio give different weight to positions depending the moment when
 286 they occurred, and therefore are more sensitive to the moment at which devia-
 287 tion occurred. This information seems to be essential for the classification, as
 288 observed in Section 3.

significant distinction between them.

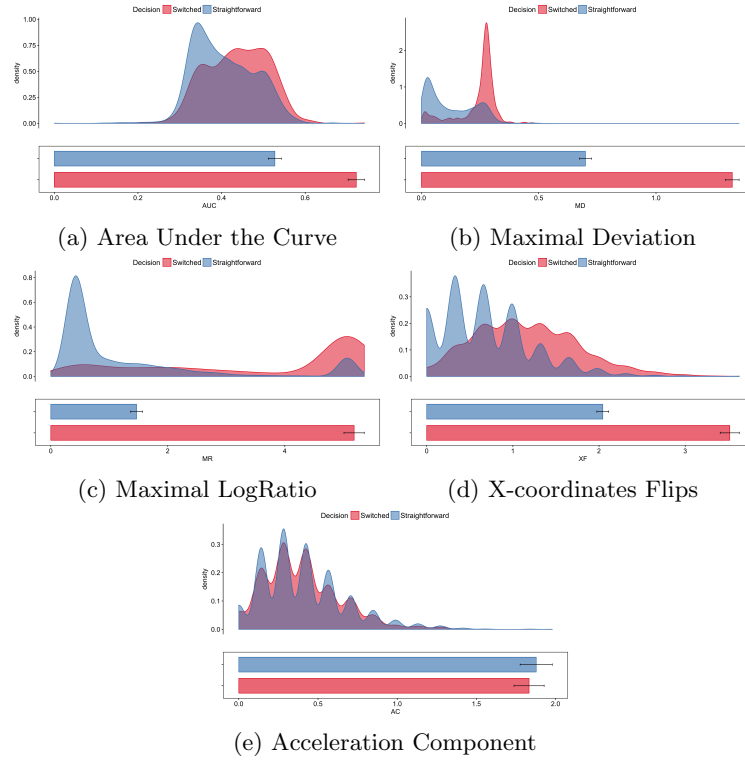


Figure 8: Different measures applied to Dale & Duran replication.

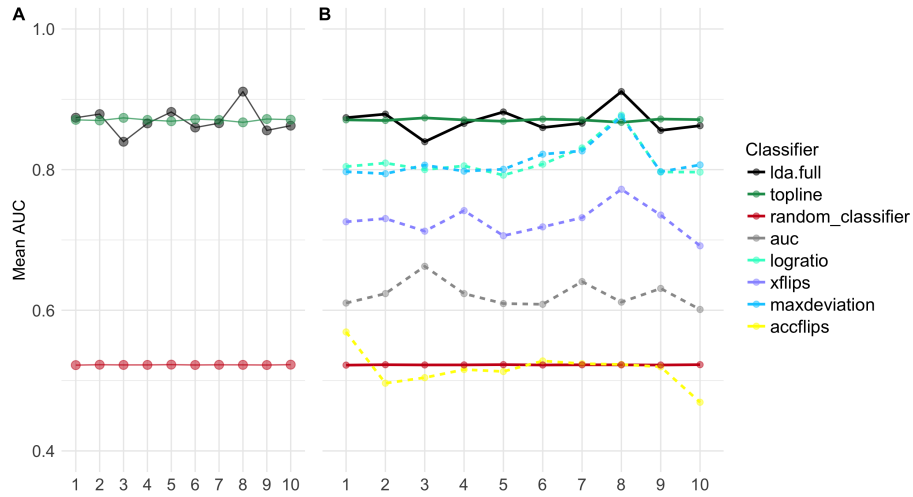


Figure 9: Mean AUC values

289 Finally, we had noticed that velocity and acceleration were not good predic-
 290 tors for the LDA classifier. Indeed, the performance of the Acceleration Com-
 291 ponent overlaps here with the one of the baseline (i.e. random classification),
 292 suggesting that this type of information is not helpful for the distinction.

293 We have shown that (i) a rough manipulation of decision making processes has
 294 a direct impact on mouse trajectories; (ii) a linear discriminant analysis (LDA)
 295 using absolute-temporal information is enough to accurately distinguish these
 296 quasi-decisions; and (iii) this LDA does a better classification than other tradi-
 297 tional mouse-tracking measures. But, can our LDA can classify more complex
 298 decision processes, such as the ones involved in sentence verification tasks?

299 5 Extension to linguistic data

300 How well does our LDA, trained on “quasi-decisions”, classify new trajectories,
 301 which underly cognitive processes that might or not correspond to different
 302 decision patterns? To address this question, we further apply our classifier to
 303 data obtained from a replication of Dale and Duran (2011) experiment.

304 Dale and Duran (2011) found differences in the processing of true positive
 305 and negative sentences when people performed a simple truth-value judgement
 306 task. These results were interpreted as indicating that negation underlies an
 307 abrupt shift in cognitive dynamics. If this is indeed the case, we would expect
 308 mouse trajectories corresponding to verifying negative sentences to pattern with
 309 *switched* trajectories from the validation experiment. This pattern of results
 310 would provide additional support to the hypothesis that, at least in this context,
 311 processing negation does involve a two-step derivation (i.e. an unconscious
 312 change of decision)⁵.

313 5.1 Experiment

314 Participants were asked to perform a truth-judgment task, where they had to
 315 decide whether a sentence (e.g. Cars have wheels) is true or false according to
 316 common knowledge. Each of the sentences could be either a true or a false state-
 317 ment in its negated or non-negated form. Unlike Dale and Duran’s experiment,
 318 the complete statement was presented in the middle of the screen after partici-
 319 pants pressed “start” (i.e. no self-paced reading). The “true” and “false” boxes
 320 appear at the top-left or top-right corners of the screen, in the same way as in
 321 our validation experiment. An illustration of the sentences used as examples is
 322 provided in [Table 4](#).

⁵Note that the data used to train the LDA correspond to “quasi-decisions”, namely it is only an approximation of what should happen during an unconscious change of decision, such as the one expected for negation processing. As a result, we expect some additional variability on the negation results (since there are some aspects of the decision process that the LDA won’t be able to capture).

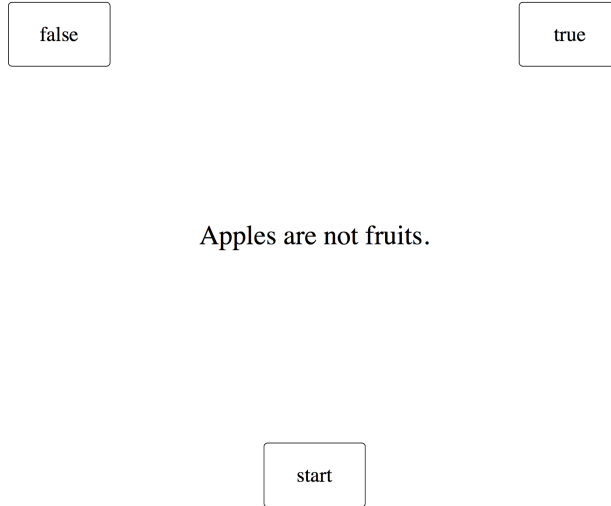


Figure 10: Trial Example Dale & Duran Replication

Participants 53 English native speakers were tested using Amazon Mechanical Turk. They were rewarded for their participation. The experiment lasted approximately 10 minutes.

Design The experimental design consisted in two fully crossed factors: TRUTH VALUE (true, false) and SENTENCE POLARITY (negative, positive). We had a total of 4 conditions, and each participant saw 4 instances of each condition (16 sentences).

Truth value	Polarity	Example
True	Positive	Cars have wheels.
	Negative	Cars have no wings.
False	Positive	Cars have wings.
	Negative	Cars have no wheels.

Table 4: Design

Interface and data treatment The interface and data treatment were the same as the ones used for the calibration experiment. Mouse trajectories' time course was normalised into 101 time steps.

5.2 Results and discussion

Replicating Dale and Duran (2011) All participants responded correctly more than 75% of the time. No participant was discarded based on accuracy.

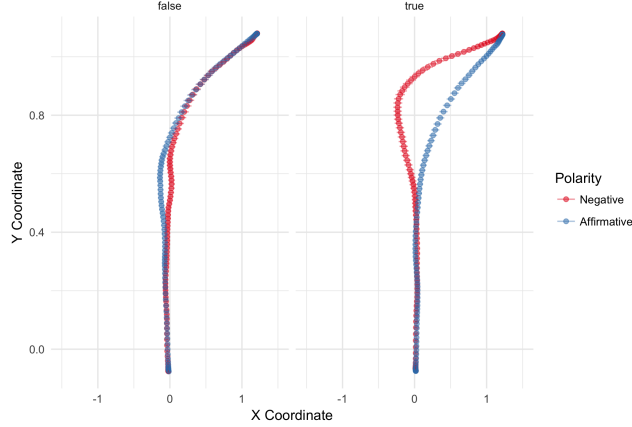


Figure 11: Mean trajectories for accurate trials

Only accurate trials were taken into account for the analysis. Figure 11 illustrates mean trajectories for the four possible trial conditions.

To assess whether we replicate Dale and Duran’s results, we calculated the x -coordinate flips (see above) and analysed them with a linear mixed-effects model (Baayen, Davidson, and Bates, 2008), taking TRUTH, POLARITY and their interaction as predictors. We included random intercepts per subject and a random slope with the interaction of both factors. P -values were obtained by comparing the omnibus model to a reduced version of itself, where the relevant factor was removed. This pipeline mirrors the model performed by Dale & Duran.

Unlike Dale and Duran, we did not perform analyses based on the acceleration component (i.e. acceleration flips). In our validation experiment, this quantitative measure was unable to distinguish mouse trajectories underlying different “quasi-decisions”. The origin of this inadequacy is hard to determine: it could be a property of the kind of decisions we are manipulating, or just a consequence of noisy data. We reasoned that if the different decision processes involved in a rather simple task were not captured by the acceleration component, this measure might also be unable to classify more complex processes, such as the ones at play in a sentence verification task.

The model revealed a main effect of POLARITY, such that negation significantly increases x -coordinate flips by 0.76 ($\chi^2 = 10.11; p = .0014$), and a significant interaction TRUTH \times POLARITY ($\chi^2 = 22.7; p < .001$), such that the difference between negative and positive sentences is bigger for the true than for the false statements. There was no significant effect of TRUTH ($\chi^2 < 1; p = .5$). Table 5 summarises the pattern of means and estimates for both ours and Dale and Duran’s results.

These results seem to replicate Dale and Duran’s findings: Verifying true negated sentences produces less straightforward trajectories than true positive

Condition	x -flips	x -flips in D&D
T/no negation	2.22	1.13
T/negation	3.67	1.71
F/no negation	2.82	1.24
F/negation	2.9	1.34
Estimate Polarity	.76	0.35
Estimate Truth	.07	0.13
Estimate Truth \times Polarity	1.35	0.47

Table 5: Mean and effect estimates

sentences (i.e. negation gives rise to more ‘curvy’ trajectories). The values obtained in the two experiments, however, are slightly different; namely, our results present higher range of values (see Table 5). Note that, in our experiments, the mouse-position was not sampled at a fixed rate (see above), creating additional noise which could be responsible for the range difference^{maybe footnote?}.

More generally, our findings pattern with a broader set of psycholinguistic studies, which, using different techniques, have shown that verifying negative sentences involve computing the positive content at an early processing stage (CITE).

Classifier performance Two different LDA classifiers, trained with validation data, were applied to the new experimental data. The first classifier was our original LDA, which had as predictors x, y coordinates as well as distance-based velocity and acceleration. The second LDA had only x, y coordinates as predictors. Validation results (see above) suggest that the simpler model, which only relies on absolute information, might be sufficient to classify the two basic kinds of decision processes. That is to say, the simple model might fit the data as well as a more complex model, and be interpreted more straightforwardly.

The relevant difference in processing between affirmative and negative sentences is expected to arise specifically for *true* statements – there is an interaction with truth values. Consequently, we analysed the performance of both classifiers when applied to true trials. Figure 12 illustrates the distribution and means of the resulting *lda measure*.

To assess how well these classifiers dissociate between positive and negative trials, we bootstrapped new samples from the original set of data (iterations=1000) and calculated the area under the ROC curve for each of these samples. In order to estimate the power of the classification, we measured the performance after reducing the sample size. Figure 13A shows the mean AUC values obtained after applying the same procedure to different sample sizes. Note that these values are generally lower than the ones obtained in the validation experiment. This is not surprising given that the classifier is being trained and tested with different sets of data.

Could the observed performance be expected if negative and positive trials

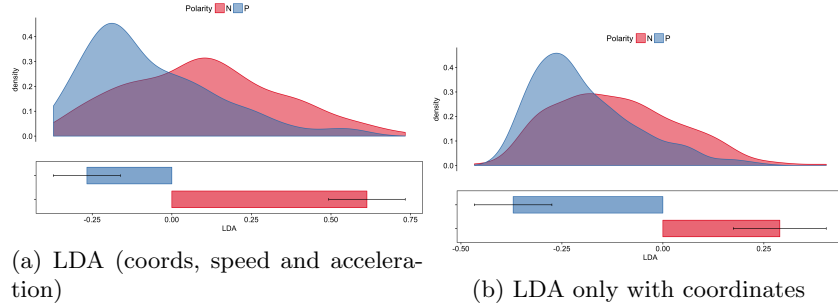


Figure 12: Two LDA classifiers applied to *true* trials.

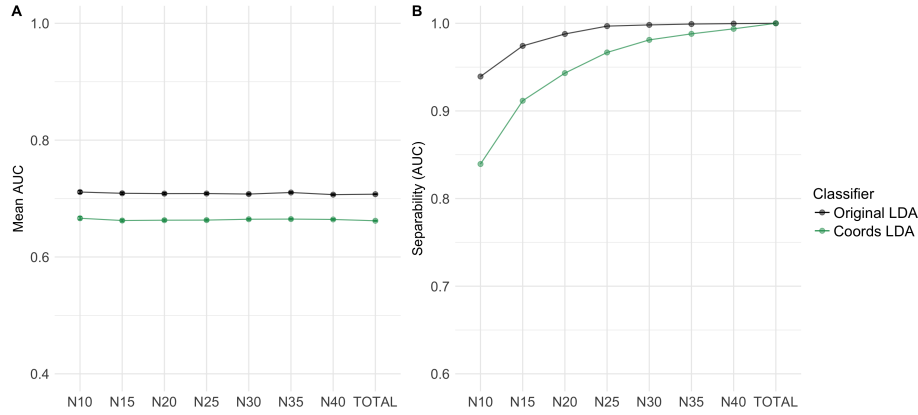


Figure 13: **Performance of LDA classifiers.** A. Mean AUC values over bootstrapped data (iterations=1000) for different sample sizes; B. Difference of classifier performance when applied to scrambled vs. original set of data.

396 were actually not different from each other? Are these AUC values significantly
 397 different from the ones one could have obtained from applying the LDA to a set
 398 of data where there is no difference between experimental conditions (i.e. *null*
 399 *hypothesis*)? We calculated the AUC values for a set of data where the exper-
 400 imental labels (positive, negative) were scrambled. This distribution of AUC
 401 values under the *null hypothesis* was compared to the performance observed
 402 on the original set of data. Figure 13B illustrates the separability of the two
 403 classifications for each sample size.

404 The LDA classifier trained with “quasi-decisions” can indeed make a rele-
 405 vant distinction between experimental conditions. That is to say, the contrast
 406 between negative and positive trials has something in common with the contrast
 407 between straightforward and switched trials obtained from a simulated decision
 408 process (validation experiment). The fact that negation has similar properties
 409 as pseudo “switched-decisions” suggests that verifying negative sentences might

underlie a change of decision, as proposed by Dale & Duran (2011), among others. However, while mouse trajectories corresponding to negative and ‘switched’ trials do share basic properties (e.g. shape), they seem to differ on how they are placed in the “change of decision” spectrum; namely, they occupy different parts of the quasi-decision-based LDA continuum (compare Figure 5 and Figure 12). This is not surprising given that we are dealing with different cognitive processes (very simple ‘quasi-decisions’ vs. sentence-verification based decision).

Lastly, while the classifiers’ comparison in Figure 6 indicated that relative spatio-temporal features, such as distance-based acceleration and speed, were not essential for the classification of ‘quasi-decisions’, these features do seem to play a role in the classification of sentence-verification data. When comparing the performance of both LDAs in Figure 13, we observe that the *full* classifier –which takes all features as predictors– makes a better distinction than the simplified one.

Other mouse-tracking measures In the validation experiment, the performance of the LDA classifier was shown to be significantly better than the one of other mouse-tracking measures. Does this difference remain when these measures are applied to the new experimental data (and hence to slightly different decision processes)? While it’s true the LDA trained on validation data can make a distinction between negative and positive trials, it might not be the *best* possible strategy for classification.

We address the question of whether different measures differ on their ability to find the observed effect by applying the same procedure as before: we calculated the mean area under the ROC curve for different sample sizes (cf. Figure 14A), and contrasted these values against the null hypothesis (i.e. the values we would have obtained if there had been no difference between the experimental conditions; Figure 14B)⁶.

The results in Figure 14A suggest that most measures perform a worse classification than the one observed for the validation data (compare with Figure 9). Given that a decrease in performance is attested across the board and not only for classifiers trained with validation data, this difference must be driven by properties of the new data set. Specifically, the sentence-verification data might be more variable, such that both negative and positive trials might underlie instances of different decision processes.

Moreover, the LDA classifier seems here to be as powerful as other traditional mouse-tracking measures, such as the Maximal Deviation and the Maximal LogRatio. This finding contrasts with the results of the validation experiment, opening the possibility of using any of these alternative measures to analyse mouse-tracking data from sentence verification tasks. Importantly, the classifier is still a better choice from a conceptual perspective, as long as it does not make any specific assumption about how the change of decision should be reflected

⁶Interestingly, the ranking based on power does not correspond exactly to the one based on the AUC mean values: how good is each measure at detecting the effect is not necessarily equivalent to its absolute performance.

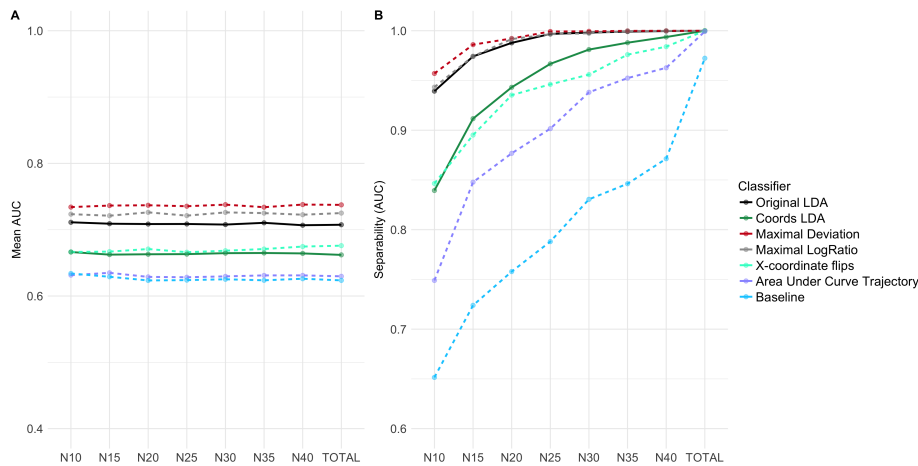


Figure 14: Performance of other measures

by mouse-trajectories.

Baseline Our LDA, trained to classify ‘quasi-decisions’, can separate the two experimental conditions of Dale & Duran’s replication. We have interpreted this result by suggesting that the LDA is distinguishing mouse trajectories that underlie two different decision processes. Alternatively, one could argue that the classification made by the LDA is not based on decision processes, but on some other feature of mouse paths, which happens to be partially shared between conditions in both experiments. For example, the LDA might not be sensitive to decision shift but to cognitive cost, and the contrast between straightforward and switched trials, on one hand, and positive and negative trials, on the other, might have that in common.

To disentangle these possibilities, we asked how our LDA classifies trajectories that might have different shapes but do not underlie two different decision processes (straightforward vs. switched), but a single one. To this end, we constructed a *baseline* set of data, which contained only positive trials from the original data set. These trials were further categorised into two classes depending on whether their response time was above or below the subject mean. We reasoned that shorter response times would correspond to an ‘early commitment’ towards the answer, whereas longer response times would reflect a ‘late commitment’. Importantly, no trial in the *baseline* data set was assumed to underlie a decision shift; thus, the LDA was expected to perform a poor classification.

As illustrated by mean trajectories in Figure 15a, the two classes in the baseline data have slightly different trajectory shapes. The distribution of the LDA measure after testing the classifier on the new data set is shown in Figure 15b. The performance was evaluated following the same procedure applied above (see

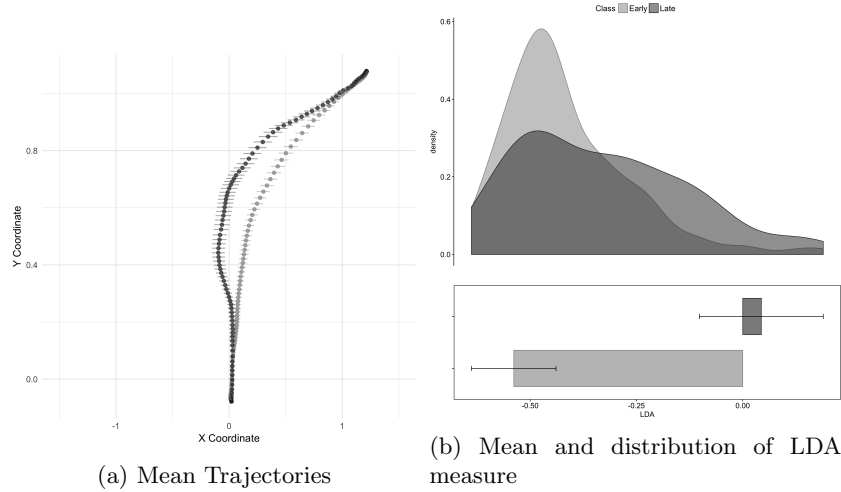


Figure 15: Baseline data set

blue line in Figure 14).

The classification on ‘early’ and ‘late’ categories is less accurate than the one performed in the original data set, to separate negative and positive trials. Differences in trajectories that are not due to the experimental manipulation are poorly captured by the LDA measure: even trajectories that do have some properties in common with *switched* and *negation* trials are not taken to be underlying a change of decision. These findings suggest that our classifier is not just tapping onto trajectory similarity but onto decision processes reflected on mouse trajectories.

6 Conclusion/General discussion

We aimed to investigate the connection between action and cognition by testing one of its specific instances: the mapping of decision making processes into mouse movements. Our findings make three main contributions on this point.

First, by manipulating whether the stimulus triggered or not a change of decision, we have shown –for the first time in a direct way– that mouse trajectories reflect basic decision processing: When participants were forced to change their answer, this switch had a systematic/direct mapping/impact on hand movements (Section X).

Second, we trained a LDA classifier with mouse-trajectories underlying these ‘quasi-decisions’ to determinate whether or not a given trial involved a decision shift. This LDA has been proven to accurately classify not only paths corresponding to other quasi-decisions’, but also mouse-trajectories underlying more complex decision processes, such as sentence verification. While the performance of the classifier –at this stage– might be as good as the one of other

501 commonly used mouse-tracking measures (i.e. Maximal Deviation), it has the
502 unique advantage of not relying on any specific assumption about how trajec-
503 tories should look like. Indeed, we demonstrated that the LDA classifier is not
504 just sensitive to superficial similarity between trajectories, but to the underlying
505 cognitive processes.

506 Lastly, our results also contribute to the research in negation processing.
507 Besides replicating Dale & Duran’ experiment, the classification performed by
508 the LDA suggests that verifying negative sentences involves a decision shift,
509 similar to the one used for the training. We then provide new evidence to
510 the hypothesis that processing negated sentences –at least in out-of-the-blue
511 contexts– involves a two-step derivation, where the positive argument is initially
512 computed.

513 To conclude, we should mention that the differences in the LDA performance
514 across the two experiments can be well understood by noticing that two data
515 sets capture slightly different decision processes. In order to capture more subtle
516 contrasts in decision making, the training data should contain more variation.
517 Further research will explore this possibility.

518 7 Supplementary Materials