

Unifying speed-accuracy trade-off and cost-benefit trade-off: a motor control model

Olivier Sigaud^{1,2,*}, Kevin Monfray^{1,2},

1 Institut des Systèmes Intelligents et de Robotique, UMR 7222, UPMC Univ Paris 06, Paris, France

2 Institut des Systèmes Intelligents et de Robotique, UMR 7222, CNRS, Paris, France

* E-mail: Corresponding Olivier.Sigaud@upmc.fr

Abstract

Two basic phenomena interact in the way the speed of our reaching movements is determined. First, we tend to reach faster a target that looks more rewarding, despite the additional muscular cost of a faster movement. Second, when we need to be more precise, our movement takes more time. So far, these two phenomena have been studied in isolation despite their obvious interdependency. In particular, two recent computational models of motor control address the first phenomenon. They explain the emergence of the time of movement as resulting from a cost-benefit trade-off arising from the summation of a temporally discounted reward and a cost that increases for faster movements. However, these models do not account for the second phenomenon, i.e. the dependency between movement time and precision requirements, resulting in a speed-accuracy trade-off and formally expressed by Fitts' law. Another model addresses the role of this speed-accuracy trade-off in determining movement time, but does not take the cost of movement into account.

In this paper, we propose a framework that unifies the cost-benefit trade-off and the speed-accuracy trade-off to explain movement properties related to time. With respect to the cost-benefit trade-off models, precision constraints are incorporated through the derivation of a new optimization criterion that considers probabilistic reaching of a rewarding target that may be missed if the motion is too fast.

Using this computational model, we investigate the more global trade-off arising from the interactions between movement time, cost and accuracy. We show that this model accounts for Fitts' law and for other well-established results in the motor control literature.

1 Introduction

There has been a recent progress in motor control research on understanding how the time of a reaching movement is chosen. In particular, two recent models from Shadmehr et al. [1] and Rigoux&Guigon [2] proposed an optimization criterion that involves a trade-off between the muscular effort and the subjective value of getting the reward, hence a cost-benefit trade-off (CBT). On one hand, reaching a target faster requires a larger muscular effort (refs?). On the other hand, the subjective value of reaching a target decreases as the time needed to reach the target is increased (refs?). As a result, the net expected return consisting of the subjective value minus the muscular effort is optimal for a certain time, as illustrated in Fig. 1(A).

However, these models do not account directly for basic facts about the relation between movement difficulty and movement duration as captured more than fifty years ago by Fitts' law [3]. According to this law, the smaller a target, the slower the reaching movement. This is well explained by the so-called *speed-accuracy trade-off* (SAT) stating that, the faster a movement, the less accurate it is, hence the higher the probability to miss the target.

In contrast with the models of [1] and [2], the model of Dean [4] takes the SAT into account. The key difference with respect to [1] and [2] is that, instead of maximizing a reward, this model maximizes a *reward expectation*, i.e. the reward times the probability to get it.

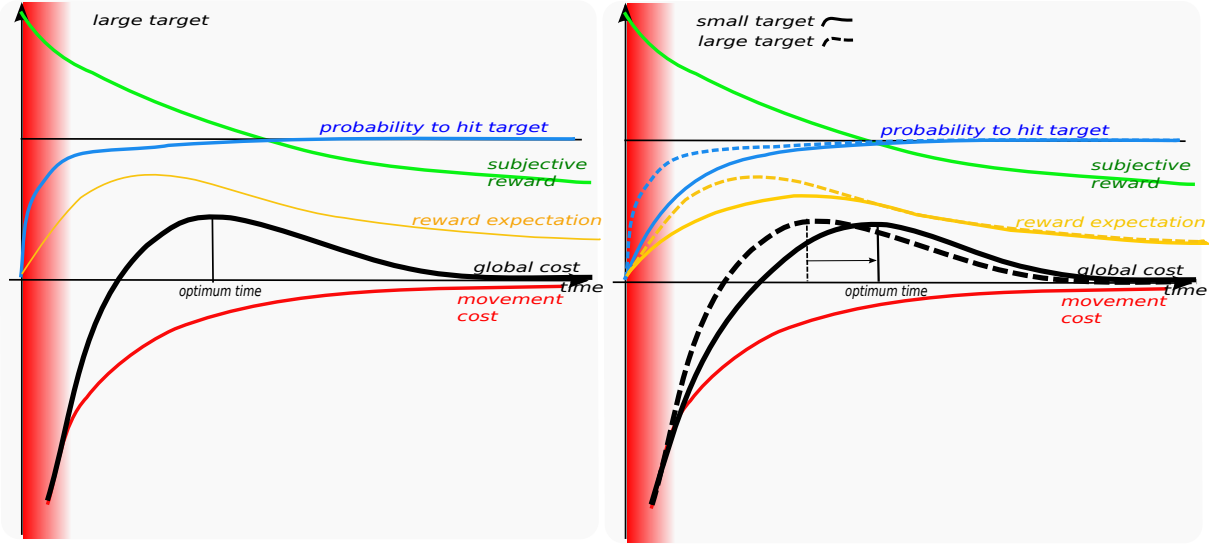


Figure 1. Influence of movement time on cost related quantities. Green: subjective utility of hitting the target; red: muscular energy cost; black: global cost versus reward trade-off. The red area denotes infeasible short times; blue: probability to hit the target; orange: reward expectation (subjective reward times probability). A: Sketch of the models in [1] and [2]. The subjective utility of hitting the reward decreases over time as one is less interested in gains that will occur in a distant future than at the present time. Hitting is less and less costly in terms of efforts as the movement is performed more slowly. The expected outcome, resulting from the sum of the subjective reward and the (negative) cost reaches a maximum for a certain time. When the outcome is negative (outside the useful interval), one should not move. B: Sketch of the presented model. In the case of a larger target, the hitting probability is higher for faster movements (solid lines) than for a smaller target (dashed lines). As a result, the maximum of the reward expectation is shifted towards longer time for smaller targets, and the optimum movement time is also longer for smaller targets.

However, the model proposed in [4] is an abstract model of movement time selection that looks for an optimal trade-off between an externally decayed reward and a SAT that relates the probability of missing to movement time. As such, it does not account for movement execution, neither for the choice of a motor trajectory and its impact on the cost of movement. The model does not explain Fitts' law, it rather incorporates its consequences into an abstract model of the SAT that is fitted to experimental data. The mathematical design of the model is based on several simplifying assumptions and it predicts optimal movement times that are systematically shorter than those observed with subjects. The authors of [4] discuss that this may result from the fact that the model does not take the cost of movement into account.

In this paper, we show that the models of Shadmehr [1] and Rigoux [2] as well as the model of Dean [4] can be unified into a model that solves the difficulties faced by these previous models.

This unification is simply implemented by including sensory and motor noise into the optimal control model proposed in [2], shifting from a deterministic account of the movement to a stochastic one, in line with the models of [5–9].

As a matter of fact, in the models of [1] and [2], the target is given as a single point and the movement is considered as always reaching it, irrespective of the size of the target. In order to fully account for Fitts' law, one must consider the intrinsic dispersion of reaching movements towards a target and the effect of sensory and muscular noise on this dispersion (e.g. [5], see [10] for a review), which is not the

case of the models of [1] and [2].

Considering expectation is a way to account for the fact that, in case of a miss, one would not get the reward, so the global outcome of the movement would only consist of its incurred cost.

The mathematical way to capture this intuition is presented in Section 4.1 and illustrated in Fig. 1(B). Technically, the reward and muscular activation terms in the optimization criterion proposed in [2] are simply replaced by reward and cost expectation terms.

With respect to Fig. 1(A), the new model includes an additional term that expresses the dependency between the size of the target and the probability to hit it, as expressed by the SAT. As can be seen in Fig. 1(B), if the target is smaller, then the probability to get it is smaller for a given time, thus the reward expectation should itself be smaller. As a result, the optimum time resulting from the optimal combination of this reward expectation with the cost of movement should shift to longer times, which is qualitatively consistent with Fitts' law.

Beyond a theoretical motor control model, this paper proposes a computational model that is tested against a set of experimental phenomena.

2 Results

The goal of the computational studies hereafter was to investigate the behavioral properties of the proposed model in order to highlight the differences with respect to [1], [2] and [4].

In a first step we illustrate some basic motor control properties of the model at the level of trajectories and motor cost. In a second step, we examine the complex relationship between movement time, final dispersion and the expected gain arising from this model. In a third step, we show that the model globally accounts for Fitts' law. Finally, we show that it generates an asymmetric velocity profile, where the asymmetry increases with the amplitude of signal-dependent motor noise.

2.1 Movement cost is not a symmetric function of target location

In order to investigate basic properties of the motor control model, it is demonstrated on a simulated two-joint planar arm with six muscles taken from the literature and illustrated in Fig. 9. This simulated arm model is described in more details in Section 4.3.

The cost obtained from the different initial configurations described in Section 4.3.4 is shown in Figure 2.

As expected, one can observe that the smaller the distance to the target, the lower the cost. Furthermore, starting from the left hand-side of the goal point results in a lower cost than starting from the right hand-side. This is explained by the fact that the optimal muscular strategy for performing these movements differs depending on the side. Actually, movements starting from the right are performed by moving simultaneously the elbow and the shoulder whereas when starting from the left, only the shoulder is involved, leading to a lower cost. The muscular activations corresponding to these two situations are shown in Figure 3.

One can see that the pattern of activation varies depending on the initial configuration of the movement, corresponding to two different optimal strategies: when starting from a short distance to the left hand-side, only the shoulder is moved, whereas when starting from a long distance to the right hand-side, both the elbow and shoulder are moved, resulting in a more complex muscular activation strategy. One can also see that co-contraction is avoided, consistently with the minimum intervention principle [6].

2.2 Relations between movement time, expected gain and final dispersion

The model presented in this paper is designed to investigate the complex relationships between movement duration, expected gain and final dispersion. In order to perform this investigation, we chose a set of five

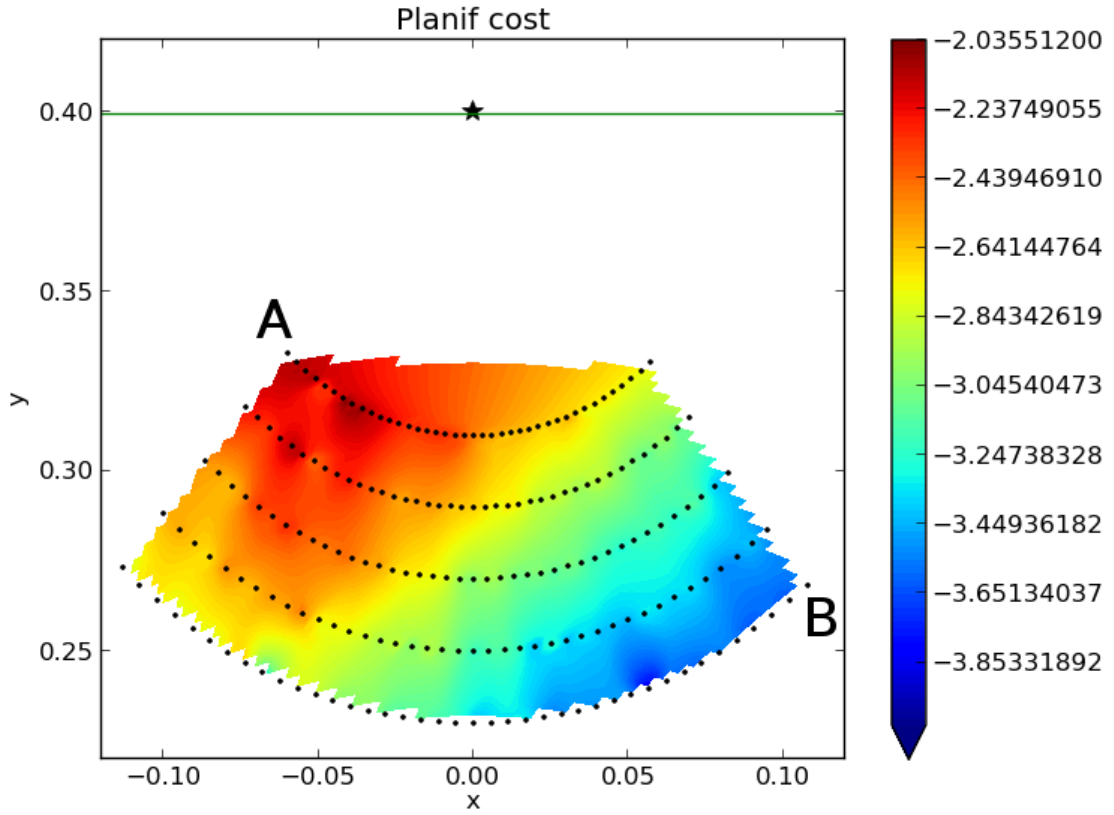


Figure 2. Cost of reaching movements towards the star at $(x = 0, y = 0.4)$. The movement is stopped when the end-effector crosses the green line at $y = 0.4$ (see Methods for details). The color of a point in the reachable space illustrates the cost of a reaching movement from that point. The color-cost correspondence is given by the scale on the right hand-side. As expected, the smaller the distance to the target, the lower the cost. Furthermore, starting from the left hand-side of the goal point results in a lower cost than starting from the right hand-side.

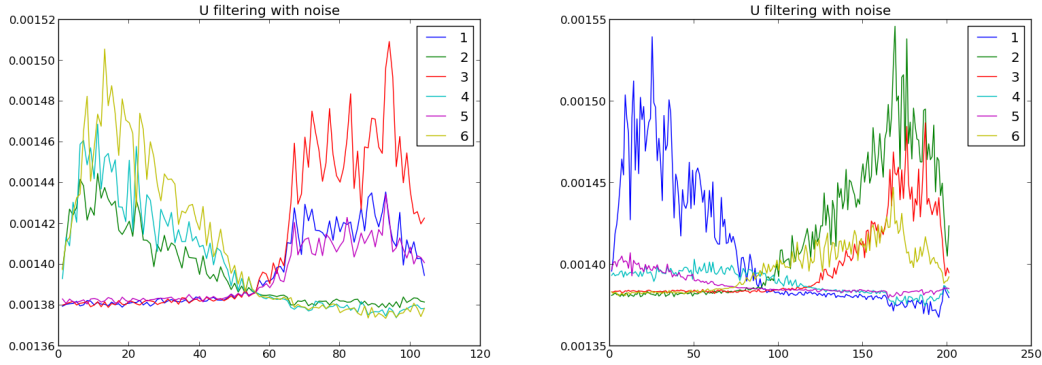


Figure 3. Muscular activations when performing a reaching movement, either starting from point A in Figure 2 (left) or from point B in the same Figure (right). The numbers in the legend correspond to the muscles numbers in Figure 9(b). One can see that the pattern of activation is very different depending on the initial configuration of the movement, corresponding to two different optimal strategies: when starting from A, only the shoulder is moved, whereas when starting from B, both the elbow and shoulder are moved, resulting in a more complex muscular activation strategy. One can also see that co-contraction is avoided, consistently with the minimum intervention principle.

target sizes ($\{1\text{mm}, 2\text{mm}, 6\text{mm}, 10\text{mm}, 20\text{mm}\}$). For each target size, we optimized a specific controller (see Methods). For each of these controllers corresponding to each target size, we recorded movement time, final dispersion and performance. The corresponding data is shown in Fig. 4.

As illustrated in Fig. 1(B), the probability to reach the target depends on the size of the target and the time of the movement (or its velocity). More precisely, if the target is smaller, fast movements should fail more often. Thus, as a result of including the accuracy constraint in the model, the optimal movement time resulting from the model described in Fig. 1(B) should be always longer than the optimal movement time resulting from the model of [2].

This is what is observed in Fig. 4(A). One can see that movement time increases when the target is smaller, and also increases with the movement distance, consistently with Fitt's law.

In Fig. 4(B), one can see that the net expected return is smaller for a larger distance, because the muscular effort for performing a larger movement is larger.

Most importantly, the net expected return increases with the size of the target. There are two explanations for this fact. First, it means that the benefit in terms of subjective value from reaching the target faster is higher than the increase in cost resulting from a faster movement. Second, if the target is larger, less precision constraints on the movement can result in a better trajectory in terms of muscular activations to reach the target.

Finally, Figure 5 shows an example of the obtained dispersion corresponding to the five targets from a distance of 18cm, using 100 trajectories from the same starting point at $x = 0$ to the target.

One can observe that, when the target is smaller, dispersion is reduced to increase the probability of reaching the target. In order to reduce dispersion, the motion is performed slower, as illustrated in Fig. 4. However, for very small targets, the probability to miss the target is not null. This means that, with our parameter settings, it is more optimal to pay the price of a few failed movements than to move slow enough to succeed at all times. There are even target sizes for which reaching may fail whatever the movement velocity.

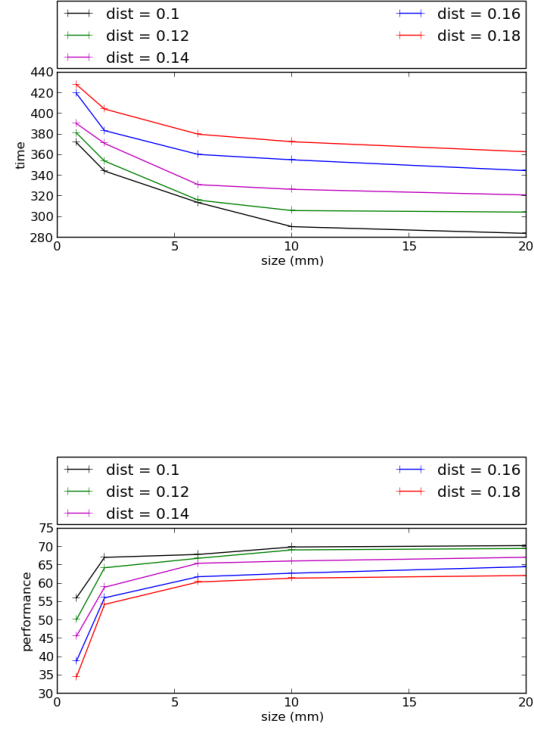


Figure 4. A: Time of movement for various target sizes and distances. The larger the target, the faster the movement. Additionnally, the further the target, the longer the movement. B: Expected movement gain for various target sizes and distances. The larger the target, the higher the gain. Additionnally, the further the target, the lower the gain.

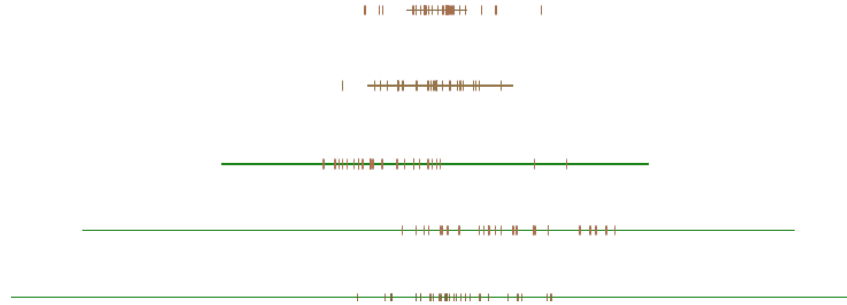


Figure 5. Dispersion resulting from the CEPS controllers optimized for five different target sizes, a movement distance of 18cm and a noise amplitude of 0.4. One can observe that the controller often misses small targets whereas it does not use all the potential dispersion for large targets.

2.3 Reproduction of Fitts' law

Fitts' law states that movement time (MT) is linear in its difficulty index (DI), this index being bigger for longer movements and smaller targets. Fitts' law is written:

$$MT = a + \underbrace{b \cdot \log_2 \left(\frac{D}{W} \right)}_{DI} \quad (1)$$

where D is the distance of the movement (denoted with A for amplitude in other papers), W is the width of the target and a and b are linear coefficients. This law was initially studied for one dimensional movements, and then extended for many other contexts [12–17].

From the data presented in Section 2.2 and using (1), we compute DI values for different distances D and target widths W . Figure 6 shows the resulting movement time MT over DI .

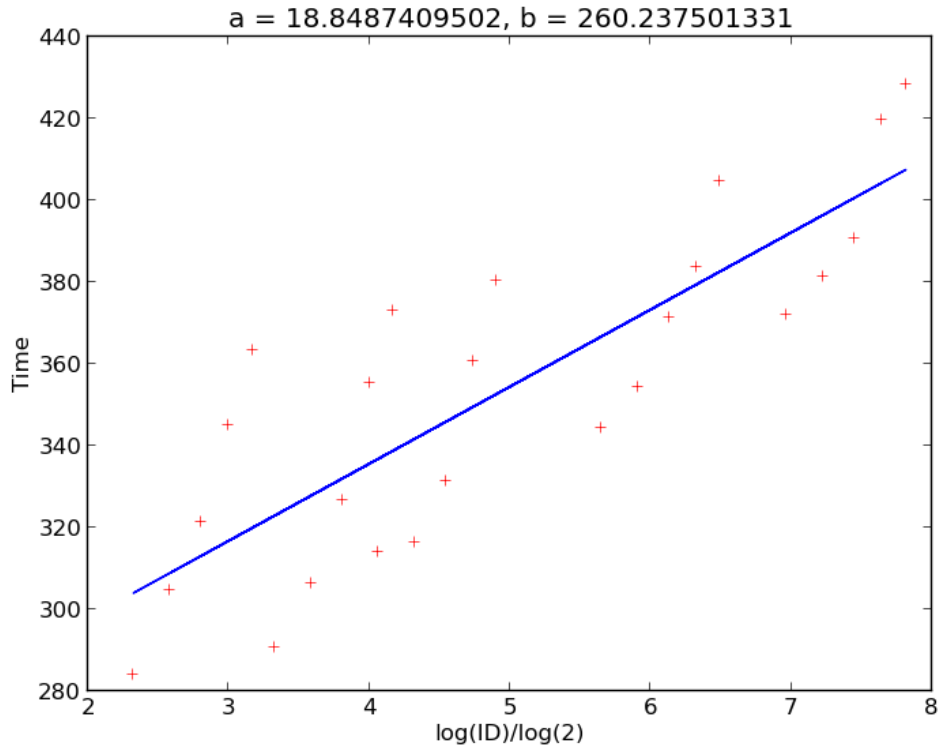


Figure 6. Reproduction of Fitts law based on the results of Section 2.2.

One can see that we get a clear linear relationships, thus the data is consistent with Fitt's law.

The obtained values of a and b cannot be compared to empirical data from the human motor control literature given the wide variability of these values accross subjects [18, 19].

2.4 Velocity profiles

The final dispersion in reaching trajectories is generated by motor noise. Following the minimum intervention principle from [6], motor noise being proportional to muscular activation, the only way to

decrease motor noise is to decrease muscular activation.

Thus, in order to hit a small target, muscular activations should be small by the end of the movement, which can result in first instance in less co-contraction and then in less velocity. Furthermore, a slower movement provides a better opportunity for state estimation to compensate for delayed feedback about the position of the end effector. Taken together, those two phenomena contribute to the fact that an optimal controller should generate less velocity by the end of the movement for a smaller target. So one way to make sure to hit a small target would be to perform a slow reaching movement.

However, as explained above, a slower movement results in a discounted reward, thus the movement should nevertheless be as fast as possible.

As a consequence, the best option for optimizing reaching accuracy under temporal constraints consists in being very fast in the beginning of the movement and much slower in the end. Thus the velocity profile should be asymmetric. The main drive for this asymmetry being motor noise, the more motor noise, the more asymmetric the movement should be.

This is what we observe in Figure 7, where velocity profiles are generated for different amplitudes of motor noise.

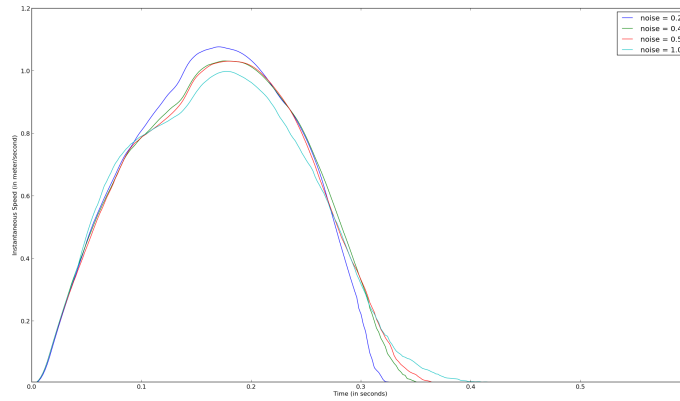


Figure 7. Velocity profiles for different amplitudes of noise.

Incidentally, one can observe on the ascending parts of the profiles that we do not get strictly the shape of a bell curve. This is due to the limited optimization capability of our methods, given the constraint on the number of samples.

3 Discussion

3.1 Positionning

When performing a reaching movement towards a target, three inter-related factors must be determined. The first is the potential outcome of a successful or failed movement, characterized by a discounted reward. Movement time is crucial in this factor because more time means a more discounted reward. The second factor is the cost of the movement, depending on muscular activations and resulting in a velocity profile and a joint-space trajectory. Movement time is crucial in this factor because, for an identical trajectory, shorter time means higher cost. The third factor is the final dispersion, generated by motor noise and imperfectly compensated for by the CNS due to sensory noise and delays. Movement time is crucial in this factor because shorter time means higher dispersion.

The works of [2] and [1] only relate the first two of these factors. By contrast, the model presented in [4] relates the first and the last factors, without consideration for the second.

The model presented in this paper addresses the more global inter-relationship between these three factors and provides an optimality criterion that accounts for the strategy of human subjects in this multi-dimensional choice space.

3.2 Background

The model presented in this paper is consistent with the stochastic optimal control view of motor control [5–9]. It starts from the fact that, for a small target, the faster the movement, the lesser the chance to hit it (**refs**). According to this view, there are two complementary explanations for this fact. First, the motor activation signal descending from the Central Nervous System (CNS) to motoneurons is corrupted with some noise that is proportional to this signal (**refs**). Thus a faster movement means more noise, hence more intrinsic dispersion of the final hit point if the arm was controlled in a purely open-loop way.

Second, this intrinsic tendency to dispersion is compensated for by a feedback control loop which is based on state estimation mechanisms. State estimation itself is based on delayed proprioceptive and perceptive feedback. Thus a faster movement means less time to accurately estimate the state, hence less compensation for dispersion. As a consequence of both mechanisms, a faster movement results on more final dispersion, hence in a lesser probability to hit a small target.

3.3 Contributions

The model proposed here takes the one of [2] as a starting point. This former model reproduces basic characteristics of motor behavior, as expected from the close relationship with previous optimal control models [6, 20–24]. It also explains several phenomena in cost-benefit trade-off tasks [25, 26]. The model presented here is equivalent to the one presented in [2] when the probability to hit the target is set to 1. As a result, it still benefits from the above mentioned properties that are not impacted by this probability. In this section we show that it solves limitations of the models of [4] and [2].

3.3.1 The model accounts for target selection bias

In [27], the authors mention a systematic pointing bias in both the x and y directions for all subjects performing a target hitting experiment. The presence of this bias contrasts with an assumption made in the model of [4] that, spatial errors being symmetric, the optimal choice of x and y should be in the middle of the circular target.

The proposed model shows that this simplifying assumption does not hold. Indeed, as illustrated in Figures 2 and 3, the optimal strategy to reach a target significantly changes depending on the relative position of the starting point and the target point. As a consequence, when aiming at a large target in the center of the sagittal plane in front of himself/herself, an optimal subject would not aim at the center of the target. The optimal aiming point depends on the function relating movement cost to aiming point location, which itself depends on the musculo-skeletal system of the subject. This explains why the bias differs from one subject to the other.

Figures 2 and 3 are obtained in a 2D case with a 1D target whereas Dean’s experiments were performed in 3D for a circular target. Nevertheless, the asymmetry resulting from the proposed model would also be present in a 3D model.

3.3.2 The model directly accounts for Fitts’ law

Motor control results are obtained in [2] in the absence of sensory and motor noise. As such, this model cannot provide a direct account of phenomena relying on the stochasticity of the motor system, such

as Fitts' law. Actually, the model of [2] provides an indirect account of Fitts' law (see [2], Fig. 7A). For obtaining these results, the authors have estimated dispersion as a function of velocity considering a constant velocity over the movement, and they have reconstructed the relationship between Difficulty Index and movement time based on the size of a target that would match this estimated dispersion ([2], personal communication). So Fig. 7A in [2] is based on one target size only.

In [4], an abstract SAT model is directly fitted to human movement data, without directly calling upon a measured movement dispersion.

In contrast, in Section 2.3 we have shown that the proposed model accounts for Fitts' law by using several targets and several starting points. In this model, movement velocity is far from constant and dispersion is measured as an effect of motor noise and imperfect state estimation rather than inferred based on an a priori SAT model.

3.4 Limitations

3.4.1 The model does not account for movement planning

In [27,28], the authors distinguish *movement planning* from *motor planning*. Movement planning consists in choosing where reaching should aim given a set of rewarded and penalized targets and motor variability. By contrast, motor planning consists in specifying movement execution in advance, in terms of muscular activations at each step of the movement, given a chosen target. Movement planning does not take motor costs into account and does not account for movement time. Interestingly, the work of [27,28] is focused on movement planning, thus it does not account for the motor trajectory and the choice of movement time as the models of [1] and [2] do.

The model proposed here might be seen as providing a first stone of the bridge between the work of [27,28] and the one of [1] and [2].

The model presented in this paper cannot explain the capability to immediately combine information about these interacting targets when they are visible, as reported in [29]. Some inference mechanism must be assumed to explain this immediate composition capability. At least part of this inference is probably initiated before movement execution starts. More generally, there is no mechanism in the model proposed here to account for movement preparation (e.g. [30]), though this stage certainly plays a role in the phenomena studied here.

3.4.2 Muscular effort or activations?

It has been shown (ref) that using $||u^2||$ or muscular effort or... results in very similar movements.

3.4.3 Exponential versus hyperbolic discounting

The proposed model starting from the one in [2], it inherits from this model an exponential discounting of the reward through time. In an alternative model, [1] rather suggests an hyperbolic discounting approach, in line with many other authors (e.g. [31]). At this stage, we consider that the debate between diverse discounting approaches is far from close (see e.g. [32]) and using a different discounting approach would not fundamentally change the results presented in this paper.

3.4.4 Expectation over reward or expected gain

The proposed model computes the expectation over the reward part rather than on the sum of the reward and the movement cost. The intuition behind this choice is that the reward term varies a lot depending on whether the target was hit or not whereas the movement cost is grossly constant over movements from the same point to a same target. If the movement cost was actually constant over movements, it could be left out of the expectation term without harm. To discriminate between both potential models,

one should investigate experimental settings where the cost of movements varies a lot, for instance using force fields. This is left for future work.

3.4.5 Imperfect optimization

Results in Figures 4(C) and 7 show that the incremental optimization process used in this paper (see Section 4) was not given enough iterations to reach a global optimum.

3.5 Predictions

The computational study presented in this paper can be seen as generating a number of predictions that remain to be tested experimentally.

First, in the context of planar reaching movements towards a large target in front of the subject, we predict a tendency to move the distribution of hit points to the left, because movements towards the left are less expensive than movements along the sagittal plane.

Second, the model proposed here progressively optimizes its distribution of hit points based on the gain resulting from previous hits. Actually, the model proposed here addresses a situation that is quite different from the one experimented in the work of Trommershäuser et al. [27, 28, 33]. Their work considers a situation where the subject can see the target and decides where to aim based on this available information. Decision is described as an inference process based on global information. In our model, by contrast, the search for the right hit point dispersion is a local trial-and-error process. Pre-training orients the controller towards an initial distribution of hit points, then the optimization process adapts this controller to a specific target but the controller is not given any prior information about the size or location of this target. It is only informed whether the target was hit or not through the reward feedback.

Thus, experimentally, the model presented in this paper would correspond to a situation where a subject is vaguely informed about the location of the target but has to adapt its reaching movement to maximise the outcome through trial-and-error. To our knowledge, this situation has never been studied experimentally in the literature.

The most closely related situation is the one described in [34], where the target is progressively shown to the subject by plotting more and more random points drawn according to the spatial distribution of the reward. Thus, in a way, the subject discovers the target through time, rather than through trial-and-error.

In the situation corresponding to the proposed model, it would be interesting to determine experimentally the circumstances under which a subject sacrifices accuracy depending on the target location, its size, its rewarding value and timing constraints over the movement. All the corresponding data could be checked against the predictions of the proposed model. In particular, one can anticipate that, if the reward gets null after a short time, subjects should perform the movement very fast at the expense of accuracy, given that hitting a rewarded target only part of the time is still better than receiving no reward at all over all trials.

3.6 From motor control to motor learning

The perspective taken here about our model consisted in considering the proposed methods as a tool to get optimal behaviours with respect to the cost function defined by Eq. (3).

By the way, this method optimizes a parametric controller for a given target size and location by trial-and-error, without knowing these size and location in advance. For a particular context, it empirically optimizes the trade-off between cost and accuracy by tuning the motor input so that velocity generates the optimal dispersion for the given target. In that respect, the model might be considered under a motor learning perspective that would try to explain how we may learn optimal reaching movements from trial-and-error, but this is beyond the scope of this paper.

4 Material and methods

In the first part of this section, we describe the theoretical background of the model. In a second part, we describe how we obtain a computational model that optimizes the cost function described in (3) for different contexts. Finally, the simulated arm and experimental apparatus used to model reaching are described in Section 4.3.

4.1 Mathematical formulation of the model

The cost function $J(\mathbf{u})$ proposed for a control \mathbf{u} in the model of [2] is

$$J(\mathbf{u}) = \int_0^\infty e^{-t/\gamma} [\rho R(\mathbf{s}_t) - \nu L(\mathbf{u}_t)] dt \quad (2)$$

where $R(\mathbf{s}_t)$ is the immediate reward function that equals 1 at the goal point (also called rewarded state) and is null everywhere else. The function $L(\mathbf{u}_t)$ is the movement cost. The authors of [2] take $L(\mathbf{u}_t) = \|\mathbf{u}_t\|^2$, as in many motor control models. The continuous-time discount factor γ accounts for the “greediness” of the controller, i.e. the smaller γ , the more the agent is focused on short term rewards. Finally, ρ is the weight of the reward term and ν the weight of the effort term. In all experiments presented here, based, on the previous work from [2], we took $\gamma = 0.998$, $\rho = 1$ and $\nu = 3000$.

A near optimal deterministic policy to solve this problem is obtained through a computationally expensive variation calculus method (see [35] for details). Given that the policy does not take the presence of noise in the model of the plant into account, the actions must be computed again at each time step depending on the new state reached by the plant which further contributes to the cost of the method. The controller resulting from this model is called the NOPS (for Near-Optimal Planning System) in the rest of this paper.

Now let us consider the integration of accuracy constraints. Instead of a deterministic controller, the new model is based on a stochastic controller where the rewarded state is reached or not. As a result, the outcome of a large set of movements performed with noise is computed as the value of the reward multiplied by the probability to obtain it over the different movements. Mathematically, the value multiplied by the probability is called the expectation.

Taking the probability to reach the target into account as described above, the new optimization criterion is written

$$J(\mathbf{u}) = \int_0^\infty e^{-t/\gamma} \mathbb{E}[\rho R(\mathbf{s}_t) - \nu L(\mathbf{u}_t)] dt \quad (3)$$

where $\mathbb{E}[\cdot]$ stands for the expectation of the cumulated reward, and $R(\mathbf{s}_t)$ equals 1 if the end effector hits the target.

4.2 Incremental stochastic optimization

The optimal control problem arising from a cost function including a reward expectation cannot be solved analytically. The reward expectation itself must be estimated empirically through a set of attempts to hit the target (these attempts are called “roll-outs” hereafter). The more roll-outs, the better the estimate of the reward expectation. In [2], the simpler optimal control problem was solved with a numerical variation calculus method called the “Near-Optimal Planning System” (NOPS) hereafter. This method is computationally expensive, it takes about 10 minutes for generating one reaching trajectory on a standard computer. As a consequence, it cannot be used as such to empirically determine the reward expectation for a given problem configuration.

To circumvent this difficulty, the computational model presented in this paper relies on a two-step approach. First, we approximate the NOPS using a nonlinear function approximation technique named

XCSF. XCSF is a regression algorithm that can approximate a function in a large continuous space [36, 37]. It generates a parametric model of the approximated function as a Gaussian mixture of linear models, i.e. a collection of local linear models bound to Gaussian support functions. A more complete description of XCSF can be found in [38–40]. The result is a parametric controller that approximates the function relating the state \mathbf{s} of the system to the adequate control input \mathbf{u} the NOPS would provide in that state. This approximated function is called the XCSF controller. It is trained by using trajectories generated by the NOPS as input samples, using the cost function described in (2). Its parameters are the weights of all local linear models learned with XCSF.

As described in [41], this controller mimics the NOPS in the limited region where it has been trained, but it reacts several orders of magnitude faster because the result of the optimisation process is “compiled” into the controller parameters. From this much faster controller, it becomes possible to empirically estimate a discounted reward expectation from many roll-outs.

In a second step, the XCSF controller is re-optimized with respect to the cost function (3) for different target sizes and different sets of initial positions. This optimization is performed using a variant of the Cross-Entropy Method (CEM) [42] illustrated in Fig. 8.

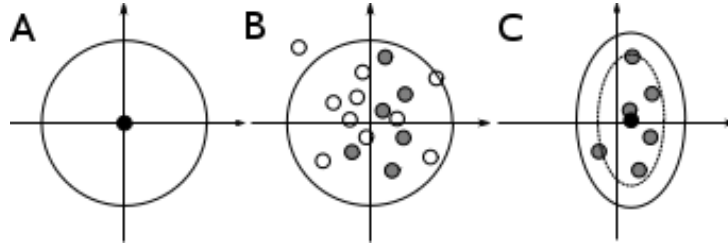


Figure 8. Schematic view of the Cross-Entropy method. A: Start with the normal distribution (μ, σ^2) . B: Draw sample parameters from this distribution, evaluate them and select the best ones (in grey). C: Compute the new μ and σ^2 (adding some noise) and go to A.

Given the initial XCSF controller learned from NOPS demonstrations, the method consists in optimizing the parameters of this controller by a local stochastic search method. New roll-outs are performed with varying parameters for all local linear models around those of the current controllers, and the parameters that give rise to a better performance with respect to the cost function (3) are retained in the new current controller. For more details about the methods, see [41].

We use the JavaXCSF [43] implementation of XCSF, and the main code for the experiments as well as the CEPS algorithm are also implemented in Java. The experiments are run on a Intel Core 2 Duo E8400 @ 3 GHz with 4 GB RAM.

XCSF is tuned as follows. The maximum number of local linear models (population) is set to 200. Learning is stopped after 200,000 iterations. The input are normalized: the target and current positions are bounded by the reachable space and the speed is bounded by $[-100, +100] \text{ rad.s}^{-1}$. The default action $\mathbf{u}_{default}$ is set to a vector of zeros i.e., no muscular activation. After tuning empirically the parameters, the learning rate **beta** is set to 0.1, the accuracy factor **alpha** is set to 1.0 and the deletion threshold **delta** is set to 0.1. Compaction, randomization and multithreading are disabled to improve reproducibility of the results [43].

4.3 Arm model and experimental apparatus

There are several models in the literature that combine a simple two joint planar rigid-body dynamics model with a muscular actuation model. Most of these models [9, 44, 45] are defined in the sagittal plane and ignore gravity effects, an exception being [46] that lies in the vertical plane and takes the gravity force into account.

Apart from this exception, the differences between the models above mostly lie in the muscular actuation component.

Our model is also a two joints planar arm in the sagittal controlled by 6 muscles, illustrated in Fig. 9, where the muscular actuation model is taken from [44] (pp. 356-357) as cited by [45].

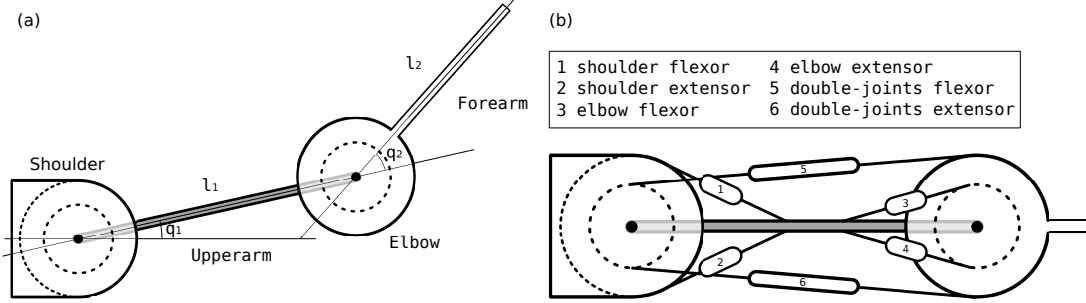


Figure 9. Arm model. (a) Schematic view of the arm mechanics. (b) Schematic view of the muscular actuation of the arm, where each number represents a muscle whose name is in the box.

Figure 9 shows the arm with the positioning of the muscles. Table 1 in Appendix A reminds the nomenclature of all the parameters and variables of the model.

4.3.1 Rigid-body dynamics

The rigid-body dynamics equation of a mechanical system is:

$$\ddot{q} = M(q)^{-1}(\tau - C(q, \dot{q})\dot{q} - g(q) - B\dot{q}) \quad (4)$$

where q is the current articular position, \dot{q} the current articular speed, \ddot{q} the current articular acceleration, M the inertia matrix, C the Coriolis force vector, τ the segments torque, g the gravity force vector and B a damping term that contains all unmodelled effects.

Here, g is ignored since the arm is working in the sagittal plane.

All angles are expressed in radians.

Given the arm parameters $m_1 = 1.4$, $m_2 = 1.1$, $l_1 = .30$, $l_2 = .35$, $s_1 = .11$, $s_2 = .16$, $d_1 = .025$, $d_2 = .045$, where m_i is the mass of segment i , l_i the length of segment i , s_i the inertia of segment i and d_i the distance from the center of segment i to its center of mass, the inertia matrix can be computed as

$$M = \begin{bmatrix} k_1 + 2k_2 \cos(q_2) & k_3 + k_2 \cos(q_2) \\ k_3 + k_2 \cos(q_2) & k_3 \end{bmatrix}, \text{ with } k_1 = d_1 + d_2 + m_2 l_1^2, \quad k_2 = m_2 l_1 s_2, \quad k_3 = d_2.$$

The Coriolis force vector is given by $C = \begin{bmatrix} -\dot{q}_2(2\dot{q}_1 + \dot{q}_2)k_2 \sin(q_2) & \dot{q}_1^2 k_2 \sin(q_2) \end{bmatrix}$.

The damping term is $B = \begin{bmatrix} .05 & .025 \\ .025 & .05 \end{bmatrix} \dot{q}$.

The computation of the torque τ exerted on the system given an input muscular actuation u is explained in the next section.

4.3.2 Muscular actuation

Our muscular actuation model is taken from [44]. It is a simplified version of the one described in [9] in the sense that it uses a constant moment arm matrix A whereas [9] is computing this matrix as a function of the state of the arm.

From [44], we take the maximum force exerted by each muscle as

$$\begin{aligned}
f_{\max} &= [f_{\max 1} \quad f_{\max 2} \quad f_{\max 3} \quad f_{\max 4} \quad f_{\max 5} \quad f_{\max 6}] \\
&= [700 \quad 382 \quad 572 \quad 445 \quad 159 \quad 318]^\top
\end{aligned}$$

and the moment arm matrix is

$$\begin{aligned}
A &= \begin{bmatrix} A_{11} & A_{21} & A_{31} & A_{41} & A_{51} & A_{61} \\ A_{12} & A_{22} & A_{32} & A_{42} & A_{52} & A_{62} \end{bmatrix}^\top \\
&= \begin{bmatrix} .04 & -.04 & 0 & 0 & .028 & -.035 \\ 0 & 0 & .025 & -.025 & .028 & -.035 \end{bmatrix}^\top.
\end{aligned}$$

Finally, given an action u corresponding to a raw muscular activation as output of the controller, the muscular activation is augmented with Gaussian noise using $\tilde{u} = \log(\exp(\kappa \times u_t \times (1 + \mathcal{N}(0, I\sigma_u^2))) + 1)/\kappa$, where \times refers to the element-wise multiplication, I is a 6×6 identity matrix. and $\kappa = 25$ is the Heaviside filter parameter, and the input torque is computed as $\tau = A^\top(f_{\max} \times \tilde{u})$.

Given (4), the simulator uses the Euler method to compute the evolution of the system, with a time step of $\Delta t = 2$ ms.

4.3.3 Motor noise and state estimation

In order to reproduce Fitts' law and to study the structure of movement variability, taking sensory and motor noise into account is necessary. In order to get an adequate dispersion, we considered multiplicative motor noise w_t with $p(w) \sim \mathcal{N}(0, 0.4)$ and additive sensory noise ν_t with $p(\nu) \sim \mathcal{N}(0, 0.0004)$. The feedback delay is 100ms. State estimation was performed with an extended Kalman filter as described in [47].

4.3.4 Experimental apparatus

In most experiments where a subject has to reach a target, either in monkeys (e.g. [48]) or humans [4, 27, 29, 34, 49], the target is displayed on a vertical screen and the subject performs a trajectory in the sagittal plane that the screen intercepts. Our experimental apparatus reproduces such a scene in the sagittal plane.

Figure 10 shows the experimental apparatus for reaching. Cartesian coordinates (in meters) are expressed with respect to the position of the shoulder, taken as origin ($x = 0, y = 0$). The arm is shown together with the screen and the set of starting positions of the tested movements. In all experiments the target is located at the same place, but its size changes. The starting configuration is varied systematically over all movements so as to investigate the effect of the distance travelled through the movement.

From a machine learning and control point of view, the state-space consists of the target articular position q^* , the current articular position q of the arm and its current articular speed \dot{q} . The state $s = (q^*, q, \dot{q})$ has a total of 6 dimensions. The initial state is defined by a position, null speed and the target position. The positions are bounded to represent the reachable space of a standard human arm, with $q_1 \in [-0.6, 2.6]$ and $q_2 \in [-0.2, 3.0]$, as shown in Fig. 10. The action-space consists of an activation signal for each muscle, which also makes a total of 6 dimensions.

To define a reaching movement with the NOPS, a goal point is required. This goal point, shown as a star in Fig. 10, is located at ($x = 0, y = 0.4$), that is 1 mm behind the screen. This position is determined so that, when training the XCSF controller to reproduce the movements generated by the NOPS, training

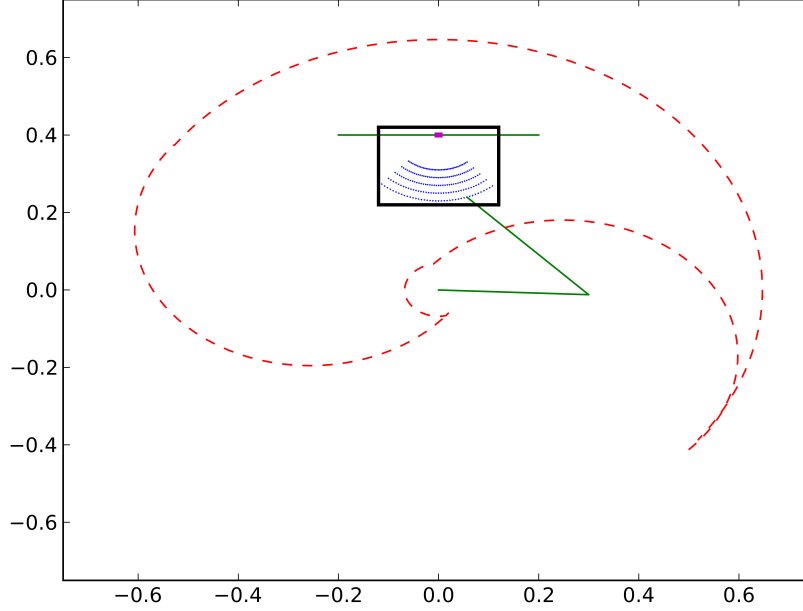


Figure 10. The arm workspace. The reachable space is delimited by a spiral-shaped envelope. The two segments of the arm are represented by two green lines. Initial movement positions are represented with blue dots organized into five sets of different distances to the target. The screen is represented as a dark green line positioned at $y = 0.4$. The origin of the arm is at $x = 0.0, y = 0.0$.

is performed beyond the screen so as to prevent the effect of discontinuities in the approximated control functions.

When using CEPS, reaching the goal point is replaced by hitting the target on the screen. The target is defined as an interval of varying length around $(x = 0, y = 0.399)$. The movement is stopped once the line $y = 0.4$ has been crossed, and the intersect between the trajectory and this line is computed to determine whether the target was hit. The reward for immediately hitting the target without taking the incurred costs into account is set to 40.

In order to train XCSF, we generate trajectories with the NOPS from the initial positions shown in Fig. 10 to the goal point. The starting points have been organized into five groups of different distances with respect to the goal point ($d = 10\text{cm}, 12\text{cm}, 14\text{cm}, 16\text{cm}$ and 18cm respectively). There are 40 initial positions per distance, thus a total of 200 initial positions.

We measure the dispersion over 100 movements towards this target, as well as the average movement time and average movement cost.

We run 200 steps of CEPS on the XCSF controller for each target.

For all the obtained controllers, we measure again the dispersion over 100 movements, the average movement time and average movement cost.

Finally, for all these targets, we record the velocity profile under different noise conditions.

A Nomenclature of arm parameters

Table 1. Parameters of the arm model.

m_i	mass of segment i (kg)
l_i	length of segment i (m)
s_i	inertia of segment i ($kg.m^2$)
d_i	distance from the center of segment i to its center of mass (m)
κ	Heaviside filter parameter
A	moment arm matrix ($\in \mathbb{R}^{6 \times 2}$)
f_{\max}	maximum muscular tension ($\in \mathbb{R}^6$)
M	inertia matrix ($\in \mathbb{R}^{2 \times 2}$)
C	Coriolis force ($N.m \in \mathbb{R}^2$)
τ	segments torque ($N.m \in \mathbb{R}^2$)
B	damping term ($N.m \in \mathbb{R}^2$)
u	raw muscular activation (action) ($\in [0, 1]^6$)
σ_u^2	multiplicative muscular noise ($\in [0, 1]^6$)
\tilde{u}	filtered noisy muscular activation ($\in [0, 1]^6$)
q^*	target articular position ($rad \in [0, 2\pi]^2$)
q	current articular position ($rad \in [0, 2\pi]^2$)
\dot{q}	current articular speed ($rad.s^{-1}$)
\ddot{q}	current articular acceleration ($rad.s^{-2}$)

References

1. Shadmehr R, Orban de Xivry JJ, Xu-Wilson M, Shih TY (2010) Temporal discounting of reward and the cost of time in motor control. *Journal of Neuroscience* 30: 10507–10516.
2. Rigoux L, Guigon E (2012) A model of reward- and effort-based optimal decision making and motor control. *PLoS Computational Biology* .
3. Fitts PM (1954) The information capacity of the human motor system in controlling the amplitude of movement. *Journal of Experimental Psychology* 47: 381–91.
4. Dean M, Wu SW, Maloney LT (2007) Trading off speed and accuracy in rapid, goal-directed movements. *Journal of Vision* 7: 1-12.
5. Harris CM, Wolpert DM (1998) Signal-dependent noise determines motor planning. *Nature* 394: 780-784.
6. Todorov E, Jordan MI (2002) Optimal feedback control as a theory of motor coordination. *Nature Neurosciences* 5: 1226-1235.
7. Todorov E (2004) Optimality principles in sensorimotor control. *Nature Neurosciences* 7: 907-915.
8. Todorov E (2005) Stochastic Optimal Control and Estimation Methods Adapted to the Noise Characteristics of the Sensorimotor System. *Neural Computation* 17: 1084–1108.

9. Li W (2006) Optimal control for biological movement systems. Ph.D. thesis, University of California, San Diego.
10. Faisal AA, Selen LPJ, Wolpert DM (2008) Noise in the nervous system. *Nature Reviews Neuroscience* 9: 292–303.
11. Selen LP, Beek PJ, van Dieën JH (2006) Impedance is modulated to meet accuracy demands during goal-directed arm movements. *Experimental Brain Research* 172: 129–138.
12. Soechting JF (1984) Effect of Target Size on Spatial and Temporal Characteristics of a Pointing Movement in Man. *Experimental Brain Research* 54: 121–132.
13. Bootsma RJ, Marteniuk RG, MacKenzie CL, Zaal FT (1994) The speed-accuracy trade-off in manual prehension: effects of movement amplitude, object size and object width on kinematic characteristics. *Experimental Brain Research* 98: 535–41.
14. Laurent M (1994) Fitts' law in two-dimensional task space. *Experimental Brain Research* 100: 144–148.
15. Plamondon R, Alimi AM (1997) Speed/accuracy trade-offs in target-directed movements. *The Behavioral and brain sciences* 20: 279–303; discussion 303–49.
16. Smyrnis N, Evdokimidis I, Constantinidis T, Kastrinakis G (2000) Speed-accuracy trade-off in the performance of pointing movements in different directions in two-dimensional space. *Experimental Brain Research* 134: 21–31.
17. Bootsma RJ, Fernandez L, Mottet D (2004) Behind Fitts law: kinematic patterns in goal-directed movements. *International Journal of Human-Computer Studies* 61: 811–821.
18. Crossman ER, Goodeve PJ (1983) Feedback control of hand-movement and fitts' law. *The Quarterly Journal of Experimental Psychology* 35: 251–278.
19. Scott MacKenzie I (1989) A note on the information-theoretic basis for fitts law. *Journal of Motor Behavior* 21: 323–330.
20. Gordon J, Ghilardi MF, Cooper SE, Ghez C (1994) Accuracy of planar reaching movements. *Experimental Brain Research* 99: 112–130.
21. Shadmehr R, Mussa-Ivaldi FA (1994) Adaptive representation of the dynamics during learning of a motor task. *Journal of Neuroscience* 14: 3208–3324.
22. Guigon E, Baraduc P, Desmurget M (2007) Computational motor control: Redundancy and invariance. *Journal of Neurophysiology* 97: 331–347.
23. Liu D, Todorov E (2007) Evidence for the flexible sensorimotor strategies predicted by optimal feedback control. *The Journal of Neuroscience* 27: 9354–9368.
24. Guigon E, Baraduc P, Desmurget M (2008) Computational motor control: Feedback and accuracy. *European Journal of Neuroscience* 27: 1003–1016.
25. Watanabe K, Lauwereyns J, Hikosaka O (2003) Effects of motivational conflicts on visually elicited saccades in monkeys. *Experimental Brain Research* 152: 361–367.
26. Rudebeck PH, Walton ME, Smyth AN, Bannerman DM, Rushworth MF (2006) Separate neural pathways process different decision costs. *Nature neuroscience* 9: 1161–1168.

27. Trommershäuser J, Maloney LT, Landy MS (2003) Statistical decision theory and the selection of rapid, goal-directed movements. *Journal of the Optical Society of America A, Optics, image science, and vision* 20: 1419-1433.
28. Trommershäuser J, Gepshtein S, Maloney LT, Landy MS, Banks MS (2005) Optimal compensation for changes in task-relevant movement variability. *The Journal of Neuroscience* 25: 7169-78.
29. Trommershäuser J, Maloney LT, Landy MS (2009) The expected utility of movement. In: *Neuroeconomics: Decision Making and the Brain*, Elsevier, 8. pp. 99-111.
30. Cos I, Bélanger N, Cisek P (2011) The influence of predicted arm biomechanics on decision making. *Journal of neurophysiology* 105: 3022-3033.
31. Prévost C, Pessiglione M, Météreau E, Cléry-Melin M, Dreher J (2010) Separate valuation subsystems for delay and effort decision costs. *The Journal of Neuroscience* 30: 14080-14090.
32. Green L, Myerson J (1996) Exponential versus hyperbolic discounting of delayed outcomes: Risk and waiting time. *American Zoologist* 36: 496-505.
33. Trommershäuser J, Maloney LT, Landy MS (2003) Statistical decision theory and trade-offs in the control of motor response. *Spatial vision* 16: 255-275.
34. Battaglia PW, Schrater PR (2007) Humans trade off viewing time and movement duration to improve visuomotor accuracy in a fast reaching task. *Journal of Neuroscience* 27: 6984-6994.
35. Rigoux L (2011) Compromis entre efforts et récompenses : Un modèle unifié de la décision et de la motricité. Thèse de doctorat, Université Pierre et Marie Curie, 4 place Jussieu 75005 Paris.
36. Wilson SW (2001) Function approximation with a classifier system. In: *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO-2001)*. San Francisco, California, USA: Morgan Kaufmann, pp. 974-981.
37. Wilson SW (2002) Classifiers that Approximate Functions. *Natural Computing* 1: 211-234.
38. Butz MV, Kovacs T, Lanzi PL, Wilson SW (2004) Toward a theory of generalization and learning in XCS. *IEEE Transactions on Evolutionary Computation* 8: 28-46.
39. Butz MV, Herbart O (2008) Context-dependent predictions and cognitive arm control with XCSF. In: *Proceedings of the 10th annual conference on Genetic and evolutionary computation*. ACM New York, NY, USA, pp. 1357-1364.
40. Sigaud O, Salaun C, Padois V (2011) On-line regression algorithms for learning mechanical models of robots: a survey. *Robotics and Autonomous Systems* 59: 1115-1129.
41. Marin D, Decock J, Rigoux L, Sigaud O (2011) Learning cost-efficient control policies with XCSF: generalization capabilities and further improvement. In: *Proceedings of the 13th annual Conference on Genetic and Evolutionary Computation*. ACM, pp. 1235-1242.
42. Rubinstein RY (1997) Optimization of computer simulation models with rare events. *European Journal of Operational Research* 99: 89-112.
43. Stalpf PO, Butz MV (2009) Documentation of JavaXCSF. Technical report, COBOSLAB.
44. Katayama M, Kawato M (1993) Virtual trajectory and stiffness ellipse during multijoint arm movement predicted by neural inverse models. *Biological Cybernetics* 69: 353-362.

45. Mitrovic D, Klanke S, Vijayakumar S (2008) Adaptive optimal control for redundantly actuated arms. In: Proceedings of the Tenth International Conference on Simulation of Adaptive Behavior. pp. 93–102.
46. Kambara H, Kim K, Shin D, Sato M, Koike Y (2009) Learning and generation of goal-directed arm reaching from scratch. *Neural networks* 22: 348–61.
47. Guigon E, Baraduc P, Desmurget M (2008) Optimality, stochasticity and variability in motor behavior. *Journal of Computational Neuroscience* 24: 57–68.
48. Kitazawa S, Kimura T, Yin PB (1998) Cerebellar complex spikes encode both destinations and errors in arm movements. *Nature* 392: 494–497.
49. Hudson TE, Maloney LT, Landy MS (2008) Optimal compensation for temporal uncertainty in movement planning. *PLoS Computational Biology* 4: e1000130.