

Attack and Defense in Cellular Decision-Making: Lessons from Machine LearningThomas J. Rademaker,¹ Emmanuel Bengio,² and Paul François¹¹*Department of Physics, McGill University, Montreal, Quebec H3A 2T8, Canada*²*School of Computer Science, McGill University, Montreal, Quebec H3A 2A7, Canada* (Received 12 February 2019; revised manuscript received 8 May 2019; published 26 July 2019)

Machine-learning algorithms can be fooled by small well-designed adversarial perturbations. This is reminiscent of cellular decision-making where ligands (called antagonists) prevent correct signaling, like in early immune recognition. We draw a formal analogy between neural networks used in machine learning and models of cellular decision-making (adaptive proofreading). We apply attacks from machine learning to simple decision-making models and show explicitly the correspondence to antagonism by weakly bound ligands. Such antagonism is absent in more nonlinear models, which inspires us to implement a biomimetic defense in neural networks filtering out adversarial perturbations. We then apply a gradient-descent approach from machine learning to different cellular decision-making models, and we reveal the existence of two regimes characterized by the presence or absence of a critical point for the gradient. This critical point causes the strongest antagonists to lie close to the decision boundary. This is validated in the loss landscapes of robust neural networks and cellular decision-making models, and observed experimentally for immune cells. For both regimes, we explain how associated defense mechanisms shape the geometry of the loss landscape and why different adversarial attacks are effective in different regimes. Our work connects evolved cellular decision-making to machine learning and motivates the design of a general theory of adversarial perturbations, both for *in vivo* and *in silico* systems.

DOI: [10.1103/PhysRevX.9.031012](https://doi.org/10.1103/PhysRevX.9.031012)Subject Areas: Biological Physics, Complex Systems,
Interdisciplinary Physics**I. INTRODUCTION**

Machine learning is becoming increasingly popular with major advances coming from deep neural networks [1]. Deep learning has improved the state of the art in automated tasks like image processing [2], speech recognition [3], and machine translation [4], and has already seen a wide range of applications in research and industry. Despite their success, neural networks suffer from blind spots: Small perturbations added to unambiguous samples may lead to misclassification [5]. Such adversarial examples are most obvious in image recognition; for example, a panda is misclassified as a gibbon or a handwritten 3 as a 7 [6]. Real-world scenarios exist, like adversarial road signs fooling computer vision algorithms [Fig. 1(a)] [7] or adversarial perturbations on medical images triggering incorrect diagnosis [8]. Worse, adversarial examples are often transferable across algorithms (see Ref. [9] for a recent review), and certain universal perturbations fool any algorithm [10].

Categorization and inference are also tasks found in cellular decision-making [11]. For instance, T cells have to discriminate between foreign and self-ligands, which is challenging since foreign ligands might not be very different biochemically from self-ligands [12,13]. Decision-making in an immune context is equally prone to detrimental perturbations in a phenomenon called ligand antagonism [14]. Antagonism appears to be a general feature of cellular decision-makers: It has been observed in T cells [15], mast cells [16], and other recognition processes like olfactory sensing [17,18].

There is a natural analogy to draw between decision-making in machine learning and in biology. In machine-learning terms, cellular decision-making is similar to a classifier. Furthermore, in both artificial and cellular decision-making, targeted perturbations lead to faulty decisions even in the presence of a clear ground-truth signal. As a consequence, arms races are observed in both systems. Mutating agents might systematically explore ways to fool the immune cells via antagonism, as has been proposed in the HIV case [19–21]. Recent examples might include neoantigens in cancer [22,23], which are implicated in tumor immunoediting and escape from the immune system. These medical examples are reminiscent of how adversaries could generate black-box attacks aimed

Published by the American Physical Society under the terms of the Creative Commons Attribution 4.0 International license. Further distribution of this work must maintain attribution to the author(s) and the published article's title, journal citation, and DOI.

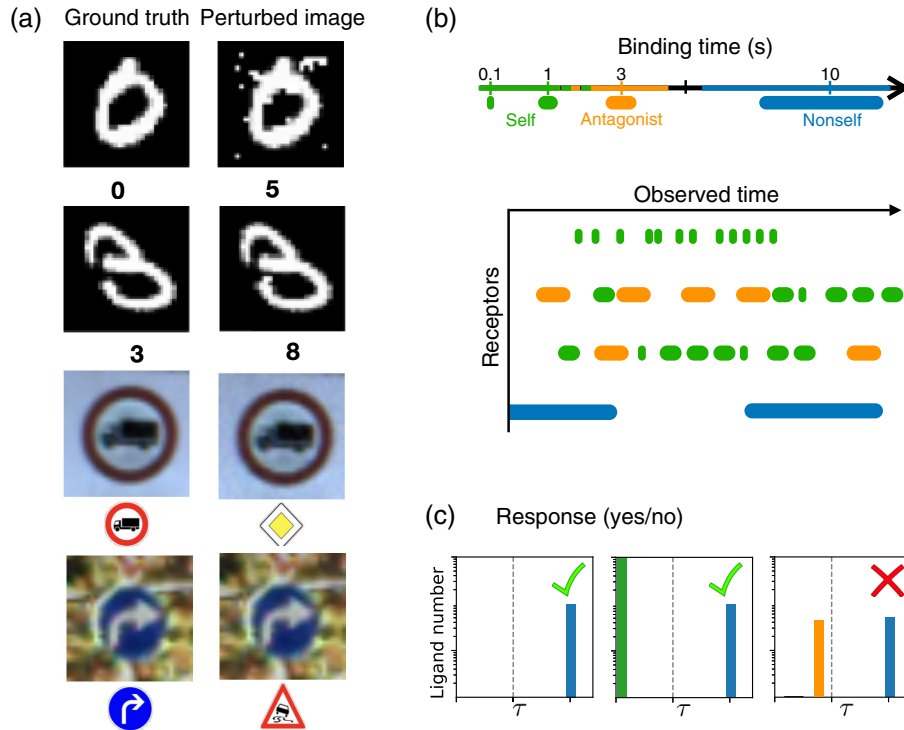


FIG. 1. Ligand discrimination and digit recognition tasks. (a) Adversarial examples on digits and road signs. Reproduced from Ref. [7], Courtesy of Nicolas Papernot. Left column displays original images with categories recognized by machine-learning algorithms, right column displays images containing targeted perturbations leading to misclassification. (b) Schematics of ligand-binding events showing typical receptor occupancy through some observed time during cellular decision-making using T-cell terminology (self vs nonself). The colored bars correspond to self (green), antagonist (orange), and nonself (blue) ligands binding to receptors. Their lengths are indicative of the binding time τ_i , whereas their rate of binding measures the on rate k_i^{on} . (c) Different ligand distributions give different response. The vertical dotted line indicates quality τ_d . Decision should be to activate if one observes ligands with $\tau > \tau_d$, so on the right of the dotted line. In an immune context, T cells respond to ligand distributions of agonists alone and agonists in the presence of nonagonists (with very small binding times τ), while the T cell fails to respond if there are too many ligands just below threshold τ_d .

to fool neural networks [7]. Strategies for provable defenses or robust detection of adversarial examples [24,25] are currently developed in machine learning, but we are still far from a general solution.

In the following, we draw a formal correspondence between biophysical models of cellular decision-making displaying antagonism on the one hand, and adversarial examples in machine learning on the other hand. We show how simple attacks in machine learning mathematically correspond to antagonism by many weakly bound ligands in cellular decision-making. Inspired by kinetic proofreading in cellular decision-making, we implement a biomimetic defense for digit classifiers, and we demonstrate how these robust classifiers exhibit similar behavior to the nonlinear adaptive proofreading models. Finally, we explore the geometry of the decision boundary for adaptive proofreading and observe how a critical point in the gradient dynamics emerges in networks robust to adversarial perturbations. Recent findings in machine learning [26] confirm the existence of two regimes, which are separated by a large nonlinearity in the activation

function. This inspires us to define two categories of attack (high dimensional, small amplitude, and low dimensional, large amplitude) both for models of cellular decision-making and neural networks. Our work suggests the existence of a unified theory of adversarial perturbations for both evolved and artificial decision-makers.

A. Adaptive proofreading for cellular decision-making

Cellular decision-making in our context refers to classification of biological ligands in two categories, e.g., “self vs nonself” in immunology or “agonist vs nonagonist” in physiology [13,27,28]. For most of those cases, qualitative distinctions rely on differences in a continuously varying property (typically a biochemical parameter). Thus, it is convenient to rank different ligands based on a parameter (notation τ) that we call quality. Mathematically, a cell needs to decide if it is exposed to ligands with quality $\tau > \tau_d$, where τ_d is the quality at the decision threshold. Such ligands’ triggering responses are called agonists. A general

problem then is to consider cellular decision-making based on ligand quality irrespective of ligand quantity (notation L). An example can be found in immune recognition with the lifetime dogma [13], where it is assumed that a T cell discriminates ligands based on their characteristic binding time τ to T-cell receptors (this is of course an approximation, and other parameters might also play a role in defining quality; see Refs. [29–31]). Ligand discrimination is a nontrivial problem for the cell, which does not measure single-binding events but has access only to global quantities such as the total number of bound receptors [Fig. 1(b)]. The challenge is to ignore many subthreshold ligands ($\tau < \tau_d$) while responding to few agonist ligands with $\tau > \tau_d$ [13,15,32]. In particular, it is known experimentally in many different contexts that the addition of antagonistic subthreshold ligands can impair proper decision-making [Fig. 1(c)] [15–17].

To model cellular decision-making, we use the general class of “adaptive sorting” or “adaptive proofreading” models, which account for many aspects of immune recognition [14,33] and can be shown to capture all relevant features of such cellular decision-making close to a decision threshold [34]. An example of such a model is displayed in Fig. 2(a). Importantly, we have shown previously that many other biochemical models present similar properties for the steady-state response as a function of the input ligand distribution [35]. In the following, we summarize the most important mathematical properties of such models. An analysis of the detailed biochemical kinetics of the model of Fig. 2(a) is presented in the Appendix A.

We assume an idealized situation where a given receptor i upon ligand binding (on rate k_i^{on} , binding time τ_i) can exist in N biochemical states (corresponding to phosphorylation stages of the receptor tails in the immune context [36,37]). Those states allow the receptor to effectively compute different quantities such as $c_n^i = k_i^{\text{on}} \tau_i^n$, $0 \leq n \leq N$, which can be done with kinetic proofreading [36,38,39]. In particular, ligands with larger τ give a relatively larger value of c_N^i due to the geometric amplification associated with proofreading steps. We assume that receptors are identical, so that any downstream receptor processing by the cell must be done on the sum(s) $C_n = \sum_i c_n^i = \sum_i k_i^{\text{on}} \tau_i^n$. We also consider a quenched situation in which only one ligand is locally available for binding to every receptor. In reality, there is a constant motion of ligands, such that k_i^{on} and τ_i are functions of time and stochastic treatments are required [11,40,41], but on the timescale of primary decision-making, it is reasonable to assume that the ligand distribution does not change much [15].

Adaptive proofreading models rely on an incoherent feed-forward loop, where an output is at the same time activated and repressed by bound ligands via two different branches in a biochemical network [Fig. 2(a)]. An explicit biochemical example is shown in the right panel of Fig. 2(a). Here, activation occurs through a kinetic

proofreading cascade (green arrow or box) and repression through the inactivation of a kinase by the same cascade (red arrow or box). The branches engage in a tug of war, which we describe below.

For simplicity, let us first assume that only one type of ligand with binding time τ and on rate k_{on} is presented. We call L the quantity of ligands. Then, in the absence of saturation, the total number of n th complex C_n of the proofreading cascade along the activation branch will be proportional to $k_{\text{on}} L \tau^n$. This branch is the activation part of the network where the response is activated.

We now assume that the m th complex of the cascades is inactivating a kinase K specific to C_m , so that $K \propto (k_{\text{on}} L \tau^m)^{-1}$ for L big enough. This branch is the repression part of the network. K is assumed to diffuse freely and rapidly between receptors so that it effectively integrates information all over the cell (recent work quantified how this cross talk can indeed improve detection [42]). m is an important parameter that we vary to compare different models. K then catalyzes the phosphorylation of the final complex of the cascade so that we have for the total number C_N ,

$$\dot{C}_N = K C_{N-1} - \tau^{-1} C_N, \quad (1)$$

and at steady state,

$$C_N \propto \frac{k_{\text{on}} L \tau^N}{k_{\text{on}} L \tau^m} = \tau^{N-m}. \quad (2)$$

The L dependence cancels, and C_N is a function of τ alone. From this, it is clear that ligand classification can be done purely based on C_N , the total number of complexes, which is a measure of ligand quality. In this situation, it is easy to define a threshold τ_d^{N-m} that governs cell activation ($C_N > \tau_d^{N-m}$) or quiescence ($C_N < \tau_d^{N-m}$). Biochemically, this can be done via the digital activation of another kinase shared among all receptors [15,33].

This model can be easily generalized to a mixture of ligands with different qualities. To do so, in the previous derivations all quantities accounting for the total complex C_n of the form $k_{\text{on}} L \tau^n$ can be replaced by $\sum_i k_i^{\text{on}} L_i \tau_i^n$, calling L_i the quantity of ligands with identical k_i^{on} , τ_i . We then define the generalized output of the biochemical network as

$$T_{N,m} = \frac{\sum_i k_i^{\text{on}} L_i \tau_i^N}{\sum_i k_i^{\text{on}} L_i \tau_i^m}. \quad (3)$$

Similar equations for an output $T_{N,m}$ can be derived for many types of networks, as described in Ref. [35]. For this reason we focus in the following on the properties of $T_{N,m}$, forgetting about the internal biochemistry giving rise to this behavior. Notice here that by construction $N > m > 1$, but other cases are possible with different biochemistry; for

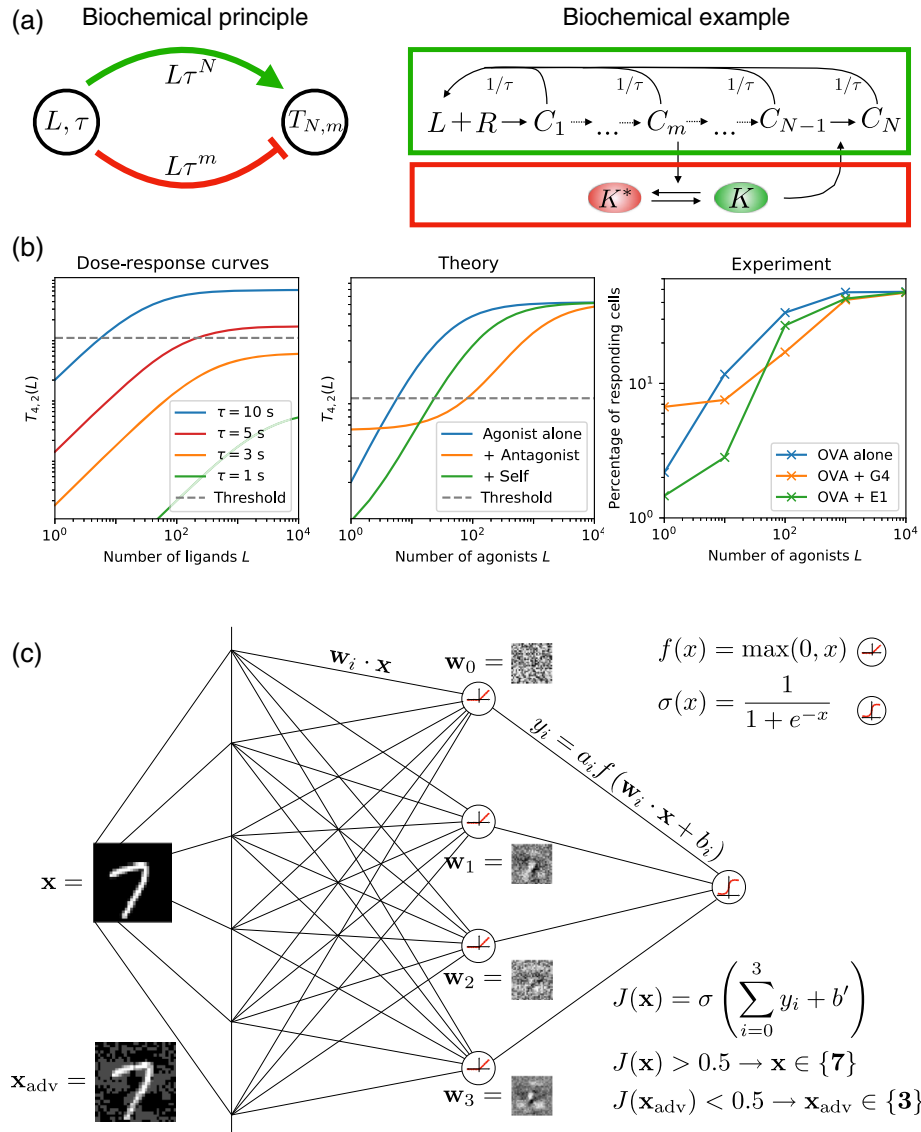


FIG. 2. Adaptive proofreading and neural network. (a) Left: Adaptive proofreading networks have an activating and repressing branch with different weights on τ . Right: Detailed adaptive proofreading network adapted from Ref. [34]. Ligand L binds to receptor R to form unphosphorylated complex C_0 . The receptor chain is iteratively phosphorylated until reaching state C_N along the activating branch (green). At every stage C_i , the ligand can unbind from the receptor with ligand-specific rate τ^{-1} . At C_m , the repressing branch (red) splits by inhibiting the kinase K , which mediates the feed-forward mechanism. (b) Dose-response curves for pure ligand types and mixtures in both adaptive proofreading models and experiments on T cells (redrawn from Ref. [32]). Details on the models and parameters used are given in Appendix B. For experiments, OVA are agonist ligands, G4 and E1 are ligands known to be below threshold but showing clear antagonistic properties. (c) Schematic of the neural network used for digit recognition. We explicitly show the four weight vectors W_i learned in one instance of the training, the activation function J , and an adversarially perturbed sample \mathbf{x}_{adv} .

instance, examples in olfaction correspond to the case $N = 1, m = 0$ [17] (see also another example in Ref. [34]). Also notice that if kinetic parameters of the ligands are not identical, the dependence on L_i does not cancel out, which will be the origin of most of the key phenomena that we describe below.

Figure 2(b) shows theoretical and experimental curves of a realistic adaptive proofreading model (including minimum concentration for repression of kinase K , etc.; see Appendix B for the full model and parameter values).

We choose $(N, m) = (4, 2)$ so that the qualitative features of the theoretical curves match the experimental curves best. Adaptive proofreading models give dose-response curves plateauing at different values as a function of parameter τ , allowing to perform sensitive and specific measurement of this parameter. For small τ (e.g., $\tau = 3$ s), one never reaches the detection threshold [dotted line on Fig. 2(b), left panel] even for many ligands. For slightly bigger $\tau = 10$ s $> \tau_d$, the curve is shifted up so that detection is made even for a small concentration of agonists.

Nontrivial effects appear if we consider mixtures of ligands with different qualities. Then, the respective computation made by the activation and repression branch of the network depends in different ways on the distribution of the presented ligand binding times. For instance, if we now add L_a antagonists with lower binding time $\tau_a < \tau$ and equal on rate k^{on} , we have $T_{N,m} = [(L\tau^N + L_a\tau_a^N)/(L\tau^m + L_a\tau_a^m)]$, which is smaller than the response τ^{N-m} for a single type of ligand, corresponding to ligand antagonism [Fig. 2(b), middle panel] [14,15,43,44]. In the presence of many ligands below the threshold of detection, the dose-response curve is simultaneously moved to the right but with a higher starting point (compared to the reference curve for “agonist alone”), as observed experimentally [Fig. 2(b), right panel, data redrawn from Ref. [32]]. Different models have different antagonistic properties based on the strength of the activation branch (N) relative to the repression branch (m). More mathematical details on these models can be found in Refs. [14,33,34].

B. Neural networks for artificial decision-making

We compare cellular decision-making to decision-making in machine-learning algorithms. We constrain our analysis to binary decision-making (which is of practical relevance, for instance, in medical applications [8]), using as a case study image classification from two types of digits. These images are taken from MNIST [45], a standard database with 70 000 pictures of handwritten digits. Even for such a simple task, designing a good classifier is not trivial, since it should be able to classify irrespective of subtle changes in shapes, intensity, and writing style (i.e., with or without a central bar for a 7).

A simple machine-learning algorithm is logistic regression. Here, the inner product of the input and a learned weight vector determines the class of the input. Another class of machine-learning algorithm is feed-forward neural networks: interconnected groups of nodes processing information layerwise. We choose to work with neural networks for several reasons. First, logistic regression is a limiting case of a neural network without hidden layers. Second, a neural network with one hidden layer more closely imitates information processing in cellular networks, i.e., in the summation over multiple phosphorylation states of the receptor-ligand complex (nodes) in a biochemical network. Third, such an architecture reproduces classical results on adversarial perturbations such as the ones described in Ref. [6]. Figure 2(c) introduces the iterative matrix multiplication inside a neural network. Each neuron i computes $\mathbf{w}_i \cdot \mathbf{x}$, $i \in [0, 3]$, adds bias b_i , and transforms the result with an activation function $f(x)$. We choose to use a rectified linear unit, which returns 0 when its input is negative and the input itself otherwise. The resulting $f(\mathbf{w}_i \cdot \mathbf{x} + b_i)$ is multiplied by another weight vector with elements a_i summed up with a bias defining a

scalar quantity $x = \sum_i a_i f(\mathbf{w}_i \cdot \mathbf{x} + b_i) + b'$. Finally, we obtain the score $J(\mathbf{x})$ (a probability between 0 and 1 for the input to belong to a class) by transforming x with the logistic function $\sigma(x)$. Parameters of such networks are optimized using classical stochastic gradient descent within a scikit implementation [46]; see Appendix B. As an example, in Fig. 2(c), a 7 is correctly classified by the neural network [$J(\mathbf{x}) > 0.5$], while the adversarial 7 is classified as a 3 [$J(\mathbf{x}_{\text{adv}}) < 0.5$].

II. RESULTS

We first summarize the general approach followed to draw the parallel between machine learning and cellular decision-making. We limit ourselves to simple classifications where a single decision is made, such as “agonist present vs no agonist present” in biology or “3 vs 7” in digit recognition. As input samples, we consider pictures in machine learning and ligand distributions in biology. We define a ligand distribution as the set of concentrations with which the ligands with unique binding times are present. This ligand distribution corresponds to a picture that is presented as a histogram of pixel values; the spatial correlation between pixels is lost, but their magnitude remains preserved. Decision-making on a sample is then done via a scoring function (or score). This score is computed either directly by the machine-learning algorithm (score J) or by the biochemical network via the concentration of a given species (score $T_{N,m}$). For simple classifications, the decision is then based on the relative value of the score above or below some threshold (typically, 0.5 for neural networks where the decision is based on sigmoidal functions, or some fixed value related to the decision time τ_d for biochemical networks).

The overall performance of a given classifier depends on the behavior of the score in the space of possible samples (i.e., the space of all possible pictures or the space of all possible ligand distributions). Both spaces have high dimensions: For instance, the dimension in the MNIST picture corresponds to number of pixels $28 \times 28 = 784$, while in immunology, ligands can bind to roughly 30 000 receptors [15]. The score can thus be thought of as a nonlinear projection of this high-dimensional space in one dimension. We study how the score behaves in relevant directions in the sample space and how to change the corresponding geometry and position of decision boundaries (defined as the samples where the score is equal to the classification threshold). We show that similar properties are observed both close to typical samples and to the decision boundary. It is important to notice at this stage that the above considerations are completely generic on the biology side and are not necessary limited to, say, immune recognition. However, we show that adaptive proofreading presents many features reminiscent of what is observed in machine learning.

A. Fast gradient sign method recovers antagonism by weakly binding ligands

In this framework, from a given sample, an adversarial perturbation is a small perturbation in sample space giving a change in score reaching (or crossing) the decision boundary. We start by mathematically connecting the simplest class of adversarial examples in machine learning to antagonism in adaptive proofreading models. We follow the original fast gradient sign method (FGSM) proposed by Ref. [6]. The FGSM computes the local maximum adversarial perturbation $\eta = \epsilon \operatorname{sgn}(\nabla_x J)$ (where sign is taken elementwise). $\nabla_x J$ represents the gradient of the scoring function categorizing images in two different categories (such as 3 and 7 in Ref. [6]). Its elementwise sign defines an image that is added to the initial batch of images with small weight ϵ . Examples of such perturbations are shown in Fig. 2(c) (bottom left) and Fig. 6(a) for the 3- vs 7-digit classification problem. While to the human observer, the perturbation is weak and changes only the background, naive machine-learning algorithms are completely fooled by the perturbation and systematically misclassify the digit.

Coming back to adaptive proofreading models, we apply FGSM for the computation of a maximally antagonistic perturbation. To do so, we need to specify the equivalent of pixels in adaptive proofreading models. A natural choice is to consider parameters associated with each pair (index i) of receptor or ligands, namely, k_i^{on} (corresponding to the rate at which ligands bind to receptors, also called on rate [47]) and τ_i (corresponding to quality). If a receptor i is unoccupied, we set its k_i and τ_i to 0 [48]. We then compute gradients with respect to these parameters.

As a simple example, we start with the case $(N, m) = (1, 0)$, which also corresponds to a recently proposed model for antagonism in olfaction [17], with the role of k^{on} played by inverse affinity κ^{-1} , the role of τ played by efficiency η , and the spiking rate of the olfactory receptor neurons is $J(T_{N,m})$, which can be interpreted as a scoring function in the machine-learning sense. In this case, $T_{1,0}$ simply computes the average quality τ_{av} of ligands presented weighted by k_i^{on} (models with $N > m > 0$ give less intuitive results as we show in the following). It should be noted that while this computation is formally simple, biochemically it requires elaborated internal interactions because a cell cannot easily disentangle influence of individual receptors; see Refs. [14,17] for explicit examples.

Starting from the computation of $\nabla_x J$ with respect to parameters k_i^{on} and τ_i , the FGSM perturbation is

$$\eta = \epsilon \operatorname{sgn} \left(\frac{\partial_{\tau_i} J}{\partial_{k_i^{\text{on}}} J} \right) = \epsilon \operatorname{sgn}(A) \operatorname{sgn} \left(\frac{k_i^{\text{on}}}{\tau_i - T_{1,0}} \right), \quad (4)$$

where $A = \{[J'(T_{1,0})]/(\sum k_i^{\text{on}})\} > 0$. Notice in the above expression that since derivatives act on different

parameters, an ϵ -sized perturbation of a given parameter is expressed in its corresponding unit. For simplicity, we do not explicitly write the conversion factor between units (this is for mathematical convenience and does not impact our results). From the above expression, we find that an equivalent maximum adversarial perturbation is given by three simple rules [Fig. 3(a)]:

- (i) decrease all τ_i by ϵ .
- (ii) decrease k_i^{on} by ϵ for ligands with $\tau_i > T_{1,0}$.
- (iii) increase k_i^{on} by ϵ for ligands with $\tau_i < T_{1,0}$.

The key relation to adversarial examples from Ref. [6] comes from considering what happens to the unbound receptors for which both k_i^{on} and τ_i are initially 0. Let us consider a situation with L identical bound ligands with ($k^{\text{on}} = 1$, binding time τ) giving response $T_{1,0}^{\text{before}} = \tau$, where τ itself is of order 1 [i.e., much bigger than the ϵ -sized perturbation on the binding time considered in Eq. (4)]. The three rules above imply that we are to decrease the binding time by ϵ and that all R previously unbound receptors are now to be bound by ligands with $k^{\text{on}} = \epsilon$, with small binding time ϵ . We compute the new response to be

$$T_{1,0}^{\text{after}} = \frac{L(\tau - \epsilon) + \epsilon R \epsilon}{L + \epsilon R} = \frac{\tau - \epsilon + \frac{\epsilon R}{L} \epsilon}{1 + \frac{\epsilon R}{L}}. \quad (5)$$

If there are many receptors compared to initial ligands, and assuming $\epsilon \ll \tau$, the relative change

$$\frac{T_{1,0}^{\text{after}} - T_{1,0}^{\text{before}}}{T_{1,0}^{\text{before}}} \simeq -\frac{\frac{\epsilon R}{L}}{1 + \frac{\epsilon R}{L}} \quad (6)$$

is of order 1 when $\epsilon R \sim L$, giving a decrease comparable to the original response instead of being of order ϵ as we naturally expect from small perturbations to all parameters. Thus, if a detection process is based on thresholding variable $T_{1,0}$, a significant decrease can happen with such a perturbation, potentially shutting down response. Biologically, the limit where ϵR is big corresponds to a strong antagonistic effect of many weakly bound ligands. Examples can be found in mast cell receptors for immunoglobulin: Weakly binding ligands have been suggested to impinge a critical kinase, thus, preventing high-affinity ligands to trigger response [16], a so-called ‘‘dog in the manger’’ effect. Another example is likely found in detection by Natural Killer cells [27]. A similar effect called ‘‘competitive antagonism’’ is also observed in olfaction where ligands with strong inverse affinity can impinge action of other ligands [17]. One difference in olfaction is that for competitive antagonism, the concentration C is of order 1 while the affinity κ^{-1} is big; conversely, here the concentration R is big while k^{on} is low. Since we consider the product of both terms, both situations lead to similar effects, but our focus on a small change of k^{on} makes the comparison with machine learning more direct.

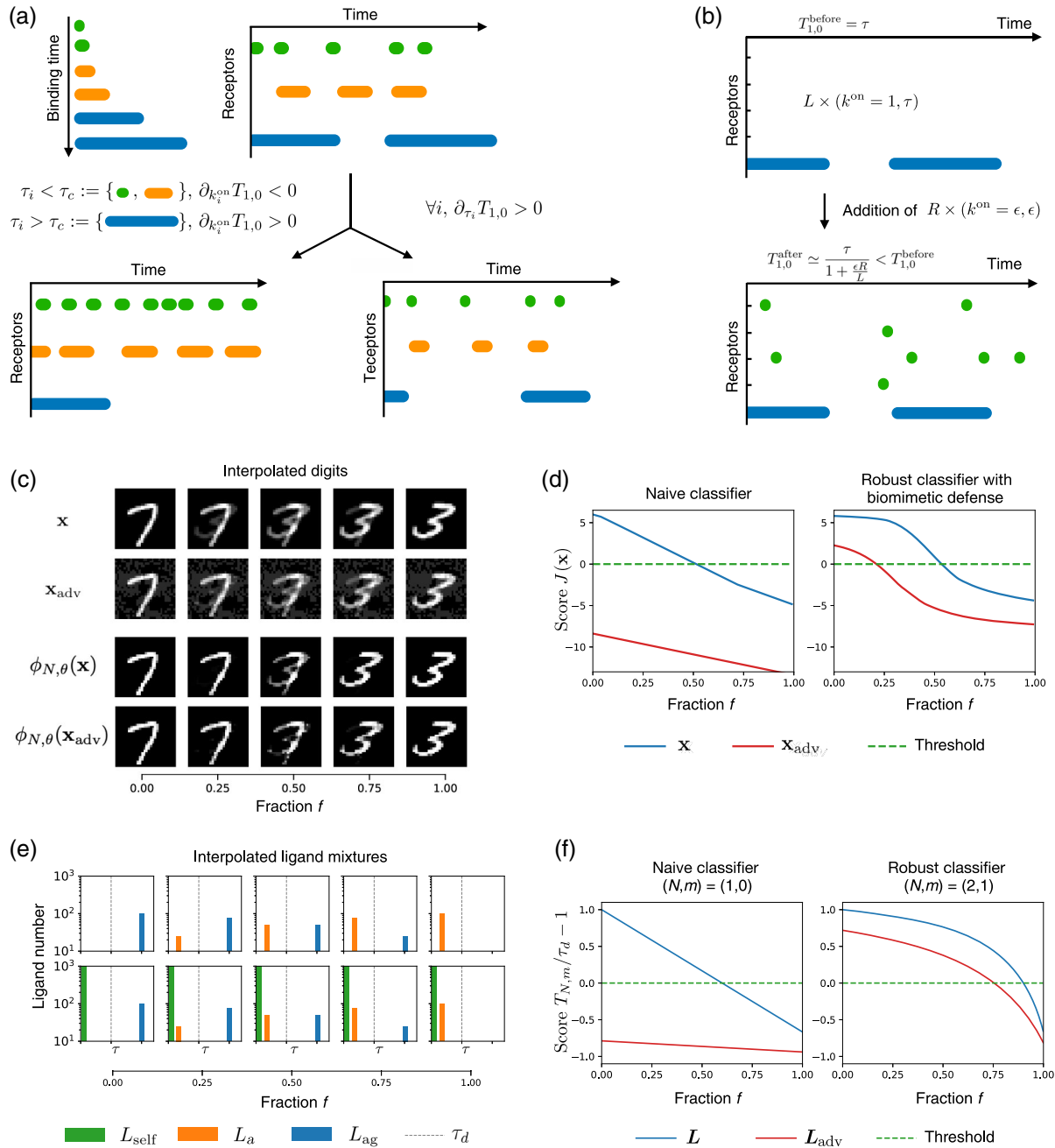


FIG. 3. Schematics of FGSM applied to immune recognition. (a) We compute how to lower the response for the receptor occupancy through a given period of time by changing k_i^{on} and τ_i . Bottom left: Increasing k_i^{on} for ligands with $\tau_i < \tau_d$ and decreasing k_i^{on} for ligands with $\tau_i > \tau_d$ reduces the weighted average $T_{1,0}$ (change in frequency of the colored bars). Bottom right: Decreasing τ_i for all ligands decreases $T_{1,0}$ (change in length of the colored bars). (b) Response to nonself-ligands is lowered from $T_{1,0}^{\text{before}}$ to $T_{1,0}^{\text{after}}$ upon addition of R ligands with small binding time ϵ . (c) Interpolated digits with and without adversarial perturbation along the interpolation axis between $\vec{7}$ ($f = 0$) and $\vec{3}$ ($f = 1$). Adversarial perturbations are computed via the FGSM with $\epsilon = 0.2$. For the biomimetic defense $\phi(N, \theta)$, we choose $N = 5$ and $\theta = 0.5$. (d) Scoring function $J(\mathbf{x})$ on pictures of (c) without (left) and with (right) the biomimetic defense. The classification threshold is indicated by the dashed green line at $J = 0$. Samples with $J > 0$ are classified as 7, otherwise 3. (e) Interpolated ligand mixtures with and without self-ligands along the interpolation axis between agonist ($f = 0$) and antagonist ($f = 1$). Here, $(L_{\text{ag}}, \tau_{\text{ag}}) = (100, 6)$, $(L_a, \tau_a) = (100, 1)$, $(L_{\text{self}}, \tau_{\text{self}}) = (1000, 0.1)$. (f) Scoring function on ligand mixtures of (e) for a naive immune classifier $(N, m) = (1, 0)$ (left) and a robust immune classifier $(N, m) = (2, 1)$ (right). The threshold is indicated by a dashed green line at $T_{N,m}/\tau_d - 1 = 0$. $T_{N,m}/\tau_d - 1 > 0$ corresponds to the detection of agonists and below corresponds to no detection. In both digit recognition and ligand discrimination, the naive networks interpolate the score linearly and are sensitive to adversarial perturbations, while the score for robust networks is flatter, closer to the initial samples for longer, thus, more resistant to perturbation.

B. Behavior across boundaries in sample space and adversarial perturbations

To further illustrate the correspondence, we compare the behavior of a trained neural network classifying 3' s and 7' s with the adaptive proofreading model $(N, m) = (1, 0)$ for more general samples. We build linear interpolations between two samples on either side of the decision boundary for both cases [Figs. 3(c)–3(f), linear interpolation factor f varying between 0 and 1]. This interpolation is the most direct way in sample spaces to connect objects in two different categories. The neural network classifies linearly interpolated digits, while the adaptive proofreading model classifies gradually changing ligand distributions.

We plot the output of the neural network x just before taking the sigmoid function σ defined in Fig. 2(c), and similarly, we plot $T_{N,m}/\tau_d - 1$ for adaptive proofreading models. In both cases, the decision is thus based on the sign of the considered quantity. In the absence of adversarial or antagonistic perturbations, for both cases, we see that the score of the system almost linearly interpolates between values on either side of the classification boundary [top panels of Figs. 3(d) and 3(f), blue curves]. However, in the presence of adversarial or antagonistic perturbations, the entire response is shifted way below the decision boundary [top panels of Figs. 3(d) and 3(f), red curves], so that, in particular, the initial samples at $f = 0$ (image of 7 or ligand distribution above threshold) are strongly misclassified.

Goodfellow *et al.* [6] proposed the linearity hypothesis as an explanation for this adversarial effect: Adding $\eta = \epsilon \operatorname{sgn}(\nabla_x J)$ to the image leads to a significant perturbation on the scoring function J of order ϵd , with d the usually high dimensionality of the input space. Thus, many weakly lit-up background pixels in the initial image can conspire to fool the classifier, explaining the significant shift in the scoring function in Fig. 3(d) top panel. The linearity hypothesis is consistent with the linearity we observe on the interpolation line, even without adversarial perturbations. A more quantitative explanation based on averaging is given in Ref. [49] on a toy model that we reproduce below to further articulate the analogy: After defining a label $y \in \{-1, +1\}$, a fixed probability p , and a constant η , one can create a $(d + 1)$ -dimensional feature vector x ,

$$y \in \{-1, +1\}, \quad x_1 \sim \begin{cases} +y, & \text{w.p. } p, \\ -y, & \text{w.p. } 1 - p, \end{cases} \quad (7)$$

$$x_2, \dots, x_{d+1} \in \mathcal{N}(\eta y, 1).$$

From this, Tsipras *et al.* build a 100% accurate classifier in the limit of $d \rightarrow \infty$ by averaging out the weakly correlated features x_2, \dots, x_d , which gives the score $f_{\text{av}} = \mathcal{N}[\eta y, (1/d)]$. Taking the sign of f_{av} will coincide

with the label y with 99% confidence for $\eta \geq 3/\sqrt{d}$. But such a classification can be easily fooled by adding a small perturbation $\epsilon = -2\eta y$ to every component of the features, since it will shift the average by the same quantity $-2\eta y$, which can still be small if we take $\eta = O(1/\sqrt{d})$ [49].

We observe a very similar effect in the simplest adaptive proofreading model. The strong shift of the average $T_{1,0}$ in Eq. (5) is due to weakly bound receptors ϵR , which play the same role as the weak features (components x_2, \dots, x_{d+1} above), hiding the ground truth given by ligands of binding time τ (equivalent to x_1 above) to fool the classifier. We also see a similar linearity on the interpolation in Fig. 3(f) top panel. There is thus a direct intuitive correspondence between adversarial examples in machine learning and many weakly bound ligands. In both cases, the change of scoring function (and corresponding misclassification) can be large despite the small amplitude ϵ of the perturbation. Once this perturbation is added, the system in Fig. 3 still interpolates between the two scores in a linear way but with a strong shift due to the added perturbation.

C. Biomimetic defense for digit classification inspired by adaptive sorting

Kinetic proofreading, famously known as the error-correcting mechanism in DNA replication [38,39], has been proposed as a mechanism for ligand discrimination [36]. In the adaptive proofreading models we study here, kinetic proofreading allows the encoding of distinct τ dependences in the activation or repression branches [33]. The primary effect of kinetic proofreading is to nonlinearly decrease the relative weight of weakly bound ligands with small binding times, thus, ensuring defense against antagonism by weakly bound ligands. Inspired by this idea, we implement a simple defense for digit classification. Before feeding a picture to the neural network, we transform individual pixel values x_i of image \mathbf{x} with a Hill function as

$$x_i \leftarrow \phi_{N,\theta}(x_i) = \frac{x_i^N}{x_i^N + \theta^N}, \quad (8)$$

where N (coefficient inspired by kinetic proofreading) and $\theta \in [0, 1]$ are parameters we choose. Similar to the defense of adaptive proofreading where ligands with small τ are filtered out, this transformation squashes grayish pixels with values below threshold θ to black pixels; see Fig. 3(c) bottom panels.

In Fig. 3(d), bottom panel, we show the improved robustness of the neural network armed with this defense. Here, the adversarial perturbation is filtered out efficiently. Strikingly, with or without adversarial perturbation, the score now behaves nonlinearly along the interpolation line in sample space: It stays flatter over a broad range of f until suddenly crossing the boundary when the digit switches identity (even for a human observer) at $f = 0.5$. Similarly,

for adaptive sorting with $(N, m) = (2, 1)$, antagonism is removed, and the score exhibits the same behavior of flatness followed by a sudden decrease on the interpolation line. Thus, similar defense displays similar robust behavior of the score in sample space.

D. Gradient dynamics identify two different regimes

The dynamics of the score along a trajectory in sample space can thus vary a lot as a function of the model considered. This motivates a more general study of a worst-case scenario, i.e., gradient descent towards the decision boundary for different models. Krotov and Hopfield studied a similar problem for a MNIST digit classifier encoded with generalized rectified polynomials of variable degrees n [50] (reminiscent of the iterative FGSM introduced in Ref. [51]). The general idea is to find out how to most efficiently reach the decision boundary and how this depends on the architecture of the decision algorithm. Krotov and Hopfield identified a qualitative change with increasing n , accompanied by a better resistance to adversarial perturbations [26,50].

We consider the same problem for adaptive proofreading models and study the potential-derived dynamics of binding times for a ligand mixture with identical k_{on} when following the gradient of $T_{N,m}$ (akin to a potential in physics). The adversarial goal is to fool the classifier with a minimal change in a given example (or in biological terms, how to best antagonize it). We iteratively change the binding time of nonagonist ligands $\tau < \tau_d$ to

$$\tau \leftarrow \tau - \epsilon \frac{\partial T_{N,m}}{\partial \tau} \quad (9)$$

while keeping the distribution of agonist ligands with $\tau > \tau_d$ constant. In the immune context, these dynamics can be thought of as a foreign agent selected by evolution to antagonize the immune system. Some biological constraints will force ligands to stay above threshold, so the only possible evolutionary strategy is to mutate and generate antagonists ligands to mask its nonself part. Such antagonistic phenomena have been proposed as a mechanism for HIV escape [19,20] and associated vaccine failure [21]. Similar mechanisms might also be implicated in the process of tumor immunoediting [23].

From a given ligand mixture with few ligands above threshold and many ligands below threshold, we follow the dynamics of Eq. (9) and display the ligand distribution at the decision boundary for different values of N , m as well as the number of steps to reach the decision boundary in the descent defined by Eq. (9) (Fig. 4; see also Fig. 5 for another example with a visual interpretation). We observe two qualitatively different dynamics. For $m < 2$, we observe strong adversarial effects, as the boundary is almost immediately reached and the ligand distribution

barely changes. As m increases, in Fig. 4(a) the ligands in the distribution concentrate around one peak. For $m = 2$, a qualitative change occurs: The ligands suddenly spread over a broad range of binding times, and the number of iterations in the gradient dynamics to reach the boundary drastically increases. For $m > 2$, the ligand distribution becomes bimodal, and the ligands close to $\tau = 0$ barely change, while a subpopulation of ligands peaks closer to the boundary. Consistent with this, the number of ϵ -sized steps to reach the boundary is 3 to 4 orders of magnitude higher for $m > 2$ as it is for $m < 2$.

E. Qualitative change in dynamics is due to a critical point for the gradient

The qualitative change of behavior observed at $m = 2$ can be understood by studying the contribution to the potential $T_{N,m}$ of ligands with very small binding times $\tau_\epsilon \sim 0$. Assuming without loss of generality that only two types of ligands are present (agonists $\tau_{\text{ag}} > \tau_d$ and spurious $\tau_{\text{spurious}} = \tau_\epsilon$), an expansion in τ_ϵ gives, up to a constant, $T_{N,m} \propto -\tau_\epsilon^m$ for small τ_ϵ [see Fig. 4(b) for a representation of this potential and Appendix C for this calculation]. In particular, for $0 < m < 1$, $[(\partial T_{N,m})/(\partial \tau_\epsilon)] \propto -\tau_\epsilon^{m-1}$ diverges as $\tau_\epsilon \rightarrow 0$. This corresponds to a steep gradient of $T_{N,m}$ so that the system quickly reaches the boundary in this direction. The ligands close to $\tau_\epsilon \sim 0$ then quickly localize close to the minimum of this potential [unimodal distribution of ligand for small m in Figs. 4(a) and 4(b)].

The potential close to $\tau_\epsilon \sim 0$ flattens for $1 < m < 2$, but it is only at $m = 2$ that a critical point for the gradient [i.e., characterized by $\partial^2 T_{N,m}/(\partial \tau_\epsilon)^2 = 0$] appears at $\tau_\epsilon = 0$. The critical point qualitatively modifies the dynamics defined by Eq. (9). For $m \geq 2$, due to the new local flatness of this gradient, ligands at $\tau = 0$, the dynamical critical point of Eq. (9), are pinned by the dynamics. By continuity, the dynamics of the ligands slightly above $\tau_\epsilon = 0$ are critically slowed down, making it much more difficult for them to reach the boundary. This explains both the sudden broadening of the ligand distribution and the associated increase in the number of steps to reach the decision boundary. Conversely, an inflexion point (square) appears in between the minimum (circle) and $\tau_\epsilon = 0$ [Fig. 4(b)]. Ligands close to the inflexion point separate and move more quickly towards the minimum of potential, explaining the bimodality at the boundary (if we continue the dynamics past the boundary, all ligands with nonzero binding times will collapse to the minimum of the potential). For both larger N and larger m , we obtain flatter potentials and a larger number of iterations. In Appendix D, we further describe the consequence of adding proofreading steps on the position of the boundary itself, using another concept of machine learning called “boundary tilting” [52] (Fig. 6 and Table I).

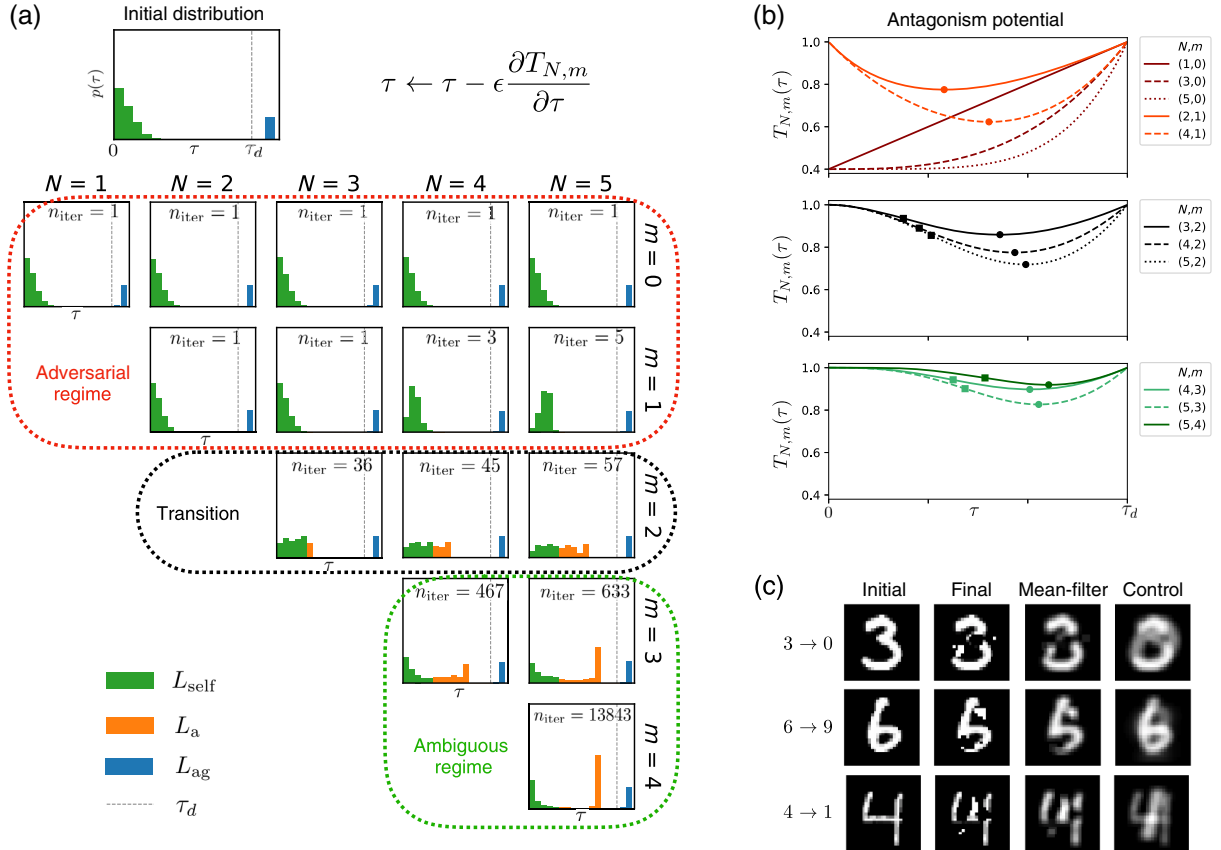


FIG. 4. Characterization of the decision boundary following gradient-descent dynamics. (a) Ligand distribution at the decision boundary by applying iterative gradient descent (top right of the panel) to an initial distribution (top left). For various cases (N, m) , we change the binding time of self-ligands along the steepest gradient until reaching the decision boundary. n_{iter} indicates the number of iterations needed to reach the decision boundary. We identify the adversarial regime (red), the ambiguous regime (green), and a transition (black) depending on m . (b) $T_{N,m}$ for mixtures of ligands at τ_d and ligands at τ , as a function of τ for various (N, m) . Antagonism strength is maximal when $T_{N,m}$ is minimal. Minima and inflexion points are indicated with a circle and square. (c) Few-pixel attack as a way of circumventing proofreading or local contrast defense, while creating ambiguous digits. We add a 3×3 mean filter to demonstrate the ambiguity of digits at the decision boundary. The control image is the mean-filtered initial digit combined with the locally contrasted average target digit. Note that also the control is lacking a clear ground truth.

F. Categorization of attacks

The transition at $m = 2$ is strongly reminiscent of the transition observed by Krotov and Hopfield in their study of gradient dynamics similar to Eq. (9) [50]. In both our works, we see that there are (at least) two kinds of attacks that can bring samples to the decision boundary. The FGSM corresponds to small perturbations to the input in terms of L_∞ norm leading to modifications of many background pixels in Ref. [50] or many weakly bound ligands for the adaptive proofreading case, also similar to the meaningless changes in x_2, \dots, x_d described above in Eq. (7) [49].

Defense against the FGSM perturbation is implemented through a higher degree n of the rectified polynomials in Ref. [50], while in adaptive proofreading, this is done through critical slowing down of the dynamics of Eq. (9) for $m > 2$. The latter models are nevertheless sensitive to another kind of attack with many fewer perturbations of the

inputs but with bigger magnitude. This attack corresponds to digits at the boundary where few well-chosen pixels are turned on in Ref. [50]. For adaptive proofreading models, this attack leads to the ligand distribution becoming bimodal at the decision boundary. Three important features are noteworthy. First, the latter perturbations are difficult to find through gradient descent [as illustrated by the many steps to reach the boundary in Fig. 4(a)]. Second, the perturbations appear to be meaningful: They correspond to interpretable features and interfere with the original sample. These perturbations make it difficult or even impossible to recover the ground truth by inspecting the sample at the decision boundary. Digits at the boundary for Ref. [50] appear indeed ambiguous to a human observer, and ligand distribution peaking just below threshold is potentially misinterpreted biologically due to inherent noise. This ambiguity has actually been observed experimentally in T cells, where strong antagonists are also weak agonists

[15,32], meaning that T cells do not take reliable decisions in this regime. Lastly, in Ref. [26] it has been observed in machine learning that memory capacity considerably increases for high n due to the local flattening of the landscape close to memories (ensuring that random fluctuations do not change memory recovery). A similar effect in our case is observed: The antagonism potential is flattened out with increasing N , m so that any spurious antagonism becomes at the same time less important and lies closer to the decision boundary.

G. Biomimetic defenses against few-pixel attacks

It is then worth testing the sensitivity to localized stronger attacks of digit classifiers, helped again with biomimetic defenses. The natural analogy is to implement attacks based on strong modification of a few pixels [53].

For this problem, we choose to implement a two-tier biomimetic defense: We implement first the transformation defined in Eq. (8) that will remove influence of the FGSM types of perturbations by flattening the local landscape as in Fig. 3(d). In addition, we choose to add a second layer of defense where we simply average out pixel values locally. This can be interpreted biologically as a process of receptor clustering or time averaging. Time averaging has been shown to be necessary in a stochastic version of adaptive proofreading [32,33], where temporal intrinsic noise would otherwise make the system cross the boundary back and forth endlessly. In the machine-learning context, local averaging has been recently proposed as a way to defend against few-pixel attacks [54], which thus can be considered as analogous to defending against biochemical noise.

We then train multiple classifiers between different pairs of handwritten digits. Following the approach of the “one-pixel” attack [53], we consider digits classified in the presence of this two-tier defense, then we sequentially fully turn pixels on or off ranked by their impact on the scoring function until we reach the decision boundary. Details on the procedure are described in Appendix E. A good defense would manifest itself similarly to the Krotov-Hopfield case [50], where no recognizable (or ambiguous) digits are observed at the boundary.

Representative results of such few-pixel attacks with biomimetic defenses are illustrated in Fig. 4(c). The “final” column shows the misclassified digits after the attack and the “mean filter” column shows the local average of the final digits for further comparison, with other examples shown in Fig. 7 and details on the behavior of scoring functions in Fig. 8. Clearly, the attacked samples at the boundaries hide the ground truth of the initial digit, and as such cannot be considered as typical adversarial perturbations. Samples at the boundary are out of distribution but preserve structure comparable to written characters (e.g., attacks from 0 to 1 typically look like a greek ϕ ; see Fig. 7). This makes them impossible to classify as arabic digits even for a human observer. This is consistent with the

ambiguous digits observed for big n by Krotov and Hopfield [50]. In other cases, samples at the boundary between two digits actually look like a third digit: For instance, we see that the sample at the boundary between a 6 and a 9 looks like a 5 (or a Japanese ζ). This observation is consistent with previous work attempting to interpolate in latent space between digits [55], where at the boundary a third digit corresponding to another category may appear. We also compare in Fig. 4(c) the sample seen by the classifier at the boundary after the biomimetic defenses with a “control” corresponding to the average between the initial digit and the target of the attack [corresponding to the interpolation factor $f = 0.5$ in Figs. 3(c) and 3(d)]. It is then quite clear that the sample generated by the attack is rather close to this control boundary image. This, combined with the fact that samples at the boundary still look like printed characters without clear ground truth indicate that the few-pixel attacks implemented here actually select for meaningful features. The existence of meaningful features in the direction of the gradient have been identified as a characteristic of networks robust to adversarial perturbation [49] similar to the results of Ref. [50] and our observation for adaptive proofreading models above.

III. DISCUSSION

Complex systems (*in vivo* or *in silico*) integrate sophisticated decision-making processes. Our work illustrates common features between neural networks and a general class of adaptive proofreading models, especially with regard to mechanisms of defense against targeted attacks. Parallels can be drawn between these past approaches since the models of adaptive proofreading presented here were first generated with *in silico* evolution aimed to design immune classifiers [33]. Strong antagonism naturally appeared in the simplest simulations and required modification of objective functions very similar to adversarial training [6].

Through our analogy with adaptive proofreading, we are able to identify the presence of a critical point in the gradient of response as the crucial mediator of robust adversarial defense. This critical point emerges due to kinetic proofreading for a cellular decision network and essentially removes the spurious adversarial directions. Another layer of defense can be added with local averaging. This is in line with current research on adversarial robustness in machine learning, showing that robust networks exhibit a flat loss landscape near each training sample [56]. Other current explorations include new biomimetic learning algorithms, giving rise to prototypelike classification [57]. Adversarial defense strategies including nonlocal computation and nonlinearities in the neural network are also currently under study [54]. The mathematical origin of the effectiveness of those defenses is not yet entirely clear, and identification of critical points in the gradient might provide theoretical insights into it.

More precisely, an interesting by-product of local flatness, where both the gradient and second derivative of the score are equal to zero, is the appearance of an inflexion point in the score and thus a region of maximal gradient. The effect of an inflexion point is visible in Figs. 3(d) and 3(f): While the score of nonrobust classifiers is linear when moving towards the decision boundary, the scoring function of classifiers resistant to adversarial perturbations is flat at $f = 0$ and significantly changes only when the input becomes ambiguous near the inflexion point. The reason why this effect is important in general is that a combination of local flatness and an inflexion point is bound to strongly influence any gradient-descent dynamics. For instance, for adaptive proofreading models, the ligand distribution following the dynamics of Eq. (9) changes from unimodal to bimodal at the boundary, creating ambiguous samples. For a robust classifier, such samples are thus expected to appear close to the decision boundary since they coincide with the larger gradients of the scoring function. As such, they could correspond to meaningful features (contrasting the adversarial perturbations), as we show in Fig. 4(c) with our digit classifier with biomimetic defense. Examples in image classification might include the meaningful adversarial transformations between samples found in Ref. [49] or the perturbed animal pictures fooling humans [58] with chimeric images that combine different animal parts (such as spider and snake), leading to ambiguous classifications. Similar properties have been observed experimentally for ambiguous samples in immune recognition: Maximally antagonizing ligands have a binding time just below the decision threshold [15]. We interpret this property as a consequence of the flat landscape far from the decision threshold leading to a steeper gradient close to it [32,34].

We use machine-learning classification and implement biomimetic defense by relying on a single direction, since that is what emerges in the most simple version of adaptive proofreading models that we consider here. In general, however, the space of inputs in machine learning is much more complex, and there are more than two categories, even in digit classification. One possible solution is to break down multilabel classification into a set of binary classification problems, but this might not always be appropriate. Instead, the algorithm effectively has to learn representations, such as pixel statistics and spatial correlations in images [2]. With a nonlinear transformation to a low-dimensional manifold description, one could still combine information on a global level in ways similar to parameter τ . The theory we present here could then apply once the mapping of the data from the full-dimensional space to such a latent space is discovered.

Case in point, Tsipras *et al.* proposed a distinction in machine learning between a robust but probabilistic feature [x_1 in Eq. (7)] and weakly correlated features [x_2, \dots, x_d in Eq. (7)] [49], both defining a single direction in latent space. They then observed a robustness-accuracy trade-off

due to the fact that an extremely accurate classifier would mostly use a distribution of many weakly correlated features (instead of the robust but randomized feature) to improve accuracy. The weight to put in the decision on either feature (robust or weak) would depend on the training. Our work shows the natural connection between weak features in this theory and weak ligands in the biological models [see discussion below Eq. (7)]. In the biological context, the standard situation is that all ligands are treated equally. Then, one can show mathematically that for such networks performing quality sensing irrespective of quantity, antagonism necessarily ensues [34], as we further identify here using the FGSM transformation. This latter result can be reformulated in terms of machine learning [49] in the following compact way: Perfectly robust classification (i.e., with no antagonism) is impossible in biology if all receptors are equivalent. But biology also provides evidence that robustness can nevertheless be improved by applying local nonlinear transformation such as the biomimetic defense of Eq. (8). Elaborating on the distinction between robust and weak features proposed in Ref. [49], nonlinear transformations should specifically target weak correlated features. Explorations of generalized nonlinear transformations in image feature space [26,50] might lead to further insights into the possible nonlinear transformations defending against adversarial perturbations. We learn in particular from biology that the major effect of nonlinearity is to change the position of maximally adversarial perturbations in sample space. Perfect robustness might be impossible in general, yet similar to cellular decision-making the most effective perturbations may shift from a pile of apparently unstructured features for naive classifiers to a combination of meaningful features for robust classifiers, giving ambiguous patterns at the decision boundary (allowing us to further distinguish between ambiguous and adversarial perturbations).

From the biology standpoint, new insights may come from the general study of computational systems built via machine learning. In particular, systematic search and application of adversarial perturbations in both theoretical models and experiments might reveal new biology. For instance, our study of Fig. 4, inspired by gradient descent in machine learning [50], establishes that cellular decision-makers exist in two qualitatively distinct regimes. The difference between these regimes is geometric by nature through the presence or absence of a dynamical critical point in the gradient. The case $m < 2$ with a steep gradient could be more relevant in signaling contexts to separate mixtures of inputs, so that every weak perturbation *should* be detected [42]. For olfaction, it has been suggested that strong antagonism allows for a rescaling of the distribution of typical odor molecules, ensuring a broad range of detection irrespective of the quantity of molecules presented [17]. The case $m \geq 2$ is much more resistant to adversarial perturbations and could be most relevant in

an immune context where T cells filter out antagonistic perturbations. This might be relevant for the pathology of HIV infections [19–21] or, more generally, it could provide explanations of the diversity of altered peptide ligands [59]. We also expect similar classification problems to occur at the population level, e.g., when T cells interact with each other to refine individual immune decision-making [60,61]. Interestingly, there might be a trade-off between resistance to such perturbations (in particular, to self antagonism, pushing towards higher m in our model) and the process of thymic selection which relies on the fact that there should be sensitivity to some self-ligands [62] (pushing towards lower m in our model).

Our correspondence could also be useful for the theoretical modeling and understanding of cancer immunotherapy [22]. So-called neoantigens corresponding to mutated ligands are produced by tumors. It has been observed that in the presence of low-fitness neoantigens, the blocking of negative signals on T cells (via checkpoint inhibitor blockade) increases the success of therapy [63]. This suggests that those neoantigens are ambiguous ligands: weak agonists acting in the antagonistic regime. Without treatment, negative signals prevent their detection (corresponding to an adversarial attack), but upon checkpoint inhibitor blockade, those ligands are suddenly visible to the immune system, which can now eliminate the tumor. Importantly, differential responses are present depending on the type of cancer, environmental factors, and tumor microenvironment [23]. This corresponds to different background ligand distributions in our framework, and one can envision that cancer cells adapt their corresponding adversarial strategies to escape the immune system. Understanding and categorizing possible adversarial attacks might thus be important in predicting the success of personalized immunotherapy [64].

We connect machine-learning algorithms to models of cellular decision-making, and in particular, their defense strategies against adversarial attacks. More defenses against adversarial examples might be found in the real world, for instance, in biofilm forming in bacteria [65], in size estimation of animals [66], or they might be needed for proper detection of physical 3D objects [67] and road signs [68]. Understanding the whole range of possible antagonistic perturbations may also prove crucial for describing immune defects, including immune escape of cancer cells. It is thus important to further clarify possible scenarios for fooling classification systems in both cell biology and machine learning.

ACKNOWLEDGMENTS

We thank Joelle Pineau and members of the François group for useful discussions. We thank three anonymous reviewers for comments and suggestions. P. F. is supported by a Simons Investigator in Mathematical Modelling of Living Systems grant, an Integrated Quantitative Biology

Initiative grant, Regroupement Québécois sur les Matériaux de Pointe, and a Natural Sciences and Engineering Research Council grant (Discovery Grant). T. J. R. receives funding from the Centre for Applied Mathematics in Bioscience and Medicine (graduate grant), McGill Physics (Schulich grant), and the Fonds de Recherche du Québec—Nature et Technologies (graduate grant). E. B. acknowledges support from the Samsung Advanced Institute of Technology and the Fonds de Recherche du Québec—Nature et Technologies (graduate grant).

APPENDIX A: MATHEMATICAL DETAILS OF THE ADAPTIVE PROOFREADING MODELS

Appendix A 1 contains more details on the derivation of adaptive proofreading models referred to in Sec. I A in the main text. In Appendix B 1, we give the parameters and equations that are used to draw Fig. 2(b) in the main text.

1. Biochemical kinetics

The kinetics for the biochemical network in Fig. 2(b) in the simplest form $[(N, m) = (2, 1)]$ are given by

$$\begin{aligned}\dot{C}_1 &= k^{\text{on}}RL - (\phi K + \tau^{-1})C_1, \\ \dot{C}_2 &= \phi KC_1 - \tau^{-1}C_2, \\ \dot{K} &= \beta(K_T - K) - \alpha C_1 K.\end{aligned}\quad (\text{A1})$$

Here, we assume the T cell has R receptors to which L ligands are bound to form ligand-receptor complexes C_1 and C_2 . The parameters k^{on} and τ^{-1} denote ligand-specific rates, which correspond to an average number of events happening per second (mean of a Poisson-distributed variable). ϕ is the phosphorylation rate for the reaction $C_1 \rightarrow C_2$ (activation branch), which is activated by variable K , and which we call a generic kinase. K itself is inhibited by C_1 (repression branch) with rate α . K_T here is the total number of kinases, and $K_T - K$ the number of inactive kinases. This kinase is shared among all receptors and assumed to diffuse freely and rapidly, so that since K is inactivated by C_1 , (in)activity of K is a measure of the total number of receptors bound. Lastly, β is the activation rate of K . In the steady state, we can solve exactly for C_2 and find

$$C_2 = \phi KC_1 \tau = \frac{L\tau}{\beta/\alpha + L} \simeq \frac{L\tau}{L} = \tau. \quad (\text{A2})$$

Here, $K = [(K_T\beta/\alpha)/(\beta/\alpha + C_1)]$, and as long as $L \gg \beta/\alpha$ the first-order approximation is exact, and the ligand dependence in the nominator and denominator cancels. Without loss of generality, we set $[(\phi K_T\beta)/\alpha] = 1$.

When we consider an environment containing two ligand types with binding times τ_{ag} (agonists) and τ_a (antagonists) at concentrations L_{ag} and L_a , two types of ligand-receptor

complexes can be formed. We call them C_i for agonists and D_i for antagonists. Full equations in the case of $(N, m) = (2, 1)$ are given by

$$\begin{aligned}\dot{C}_1 &= k^{\text{on}}RL_{\text{ag}} - (\phi K + \tau_{\text{ag}}^{-1})C_1, \\ \dot{C}_2 &= \phi KC_1 - \tau_{\text{ag}}^{-1}C_2,\end{aligned}\quad (\text{A3})$$

$$\begin{aligned}\dot{D}_1 &= k^{\text{on}}RL_a - (\phi K + \tau_a^{-1})D_1, \\ \dot{D}_2 &= \phi KD_1 - \tau_a^{-1}D_2, \\ \dot{K} &= \beta(K_T - K) - \alpha(C_1 + D_1)K,\end{aligned}\quad (\text{A4})$$

where we assume that k^{on} is equal for both agonist and antagonist ligands. The main difference here is that variable K integrates global information from both ligand complexes, which results in the steady state in $K = [(K_T\beta/\alpha)/(\beta/\alpha + C_1 + D_1)]$. Moreover, K acts locally on the phosphorylation of both C_1 and D_1 . Finally, the output is given by $T_{2,1} = C_2 + D_2$.

We can generalize this case by assuming that inhibition of the variable K occurs not at the first complex C_1 , but further downstream a kinetic proofreading cascade, namely, at the m th complex $C_m = L_{\text{ag}}\tau_{\text{ag}}^m$ and $D_m = L_a\tau_a^m$. The output variable is then given by $T_{N,m} = C_N + D_N$. Figure 2 (a) shows how information from a single ligand passes through the repression branch (red arrow and box) via K and through the activation branch (green arrow and box) via C_N . The global variable K integrates local information as $K = [(K_T\beta/\alpha)/(\beta/\alpha + C_m + D_m)] \propto (L_{\text{ag}}\tau_{\text{ag}}^m + L_a\tau_a^m)^{-1}$ and catalyzes the phosphorylation of $C_{N-1} = L_{\text{ag}}\tau_{\text{ag}}^{N-1}$ and $D_{N-1} = L_a\tau_a^{N-1}$ to final complex C_N and D_N as

$$\dot{C}_N = KC_{N-1} - \tau_{\text{ag}}^{-1}C_N, \quad (\text{A5})$$

$$\dot{D}_N = KD_{N-1} - \tau_a^{-1}D_N. \quad (\text{A6})$$

In the steady state, the solution for $T_{N,m}$ is then

$$T_{N,m} = C_N + D_N = \frac{L_{\text{ag}}\tau_{\text{ag}}^N + L_a\tau_a^N}{L_{\text{ag}}\tau_{\text{ag}}^m + L_a\tau_a^m}. \quad (\text{A7})$$

This expression for two types of ligands with same k_{on} can be clearly generalized to any type of ligand, giving Eq. (3) in the main text.

APPENDIX B: MATERIALS AND METHODS

In this Appendix, we give the parameters and equations that are used to draw Fig. 2(b) in the main text, and we give the hyperparameters used for training the neural networks classifying 3's and 7's. We refer to the latter in Sec. 1B in the main text.

1. Parameters for Fig. 2(b)

The curves in Fig. 2(b), left panel, come from the model given by

$$T_{4,2}(L) = \frac{1}{\tau_d^2} \frac{L\tau^4}{C_* + L\tau^2}, \quad (\text{B1})$$

with parameter values $C_* = \beta/\alpha = 3000$, $\tau_d = 4s$, and τ as in the legend. The curves in the middle panel of Fig. 2(b) come from

$$T_{4,2}(L) = \frac{1}{\tau_d^2} \frac{L\tau^4 + L_a\tau_a^4}{C_* + L\tau^2 + L_a\tau_a^2}, \quad (\text{B2})$$

with again $C_* = 3000$, $\tau_d = 4s$, and $\tau = 10s$. For blue “agonists alone,” $L_a = 0$, for orange “+ antagonists” $L_a = 10^4$ and $\tau_a = 3s$, and for green “+ self” $L_a = 10^4$ and $\tau_a = 1s$.

2. Hyperparameters for training neural network

We choose our hyperparameters as follows: one hidden layer with four neurons feeding into an output neuron, a random 80/20 training or test split with a 10% validation split. The cross-entropy loss function is minimized via stochastic gradient descent in maximal 300 iterations with a batch size of 200 and an adaptive learning rate initiated at 0.001. The tolerance is 10^{-4} and the regularization rate is 0.1. Most of these parameters are set to their default value, but we find that the training procedure is largely insensitive to the specific choice of hyperparameters.

APPENDIX C: LIGAND DISTRIBUTION AT THE DECISION BOUNDARY

In Appendix C 1, we describe in detail the methods used in the gradient dynamics of changing a ligand distribution to the decision boundary, we provide additional results when adding spatial correlation to the ligand distribution in Appendix C 2, and we calculate the leading order in small binding time τ_e of the gradient $(dT_{N,m})/d\tau_e$ in Appendix C 6. We refer to Appendix C in the main text in Secs. II D and II E, and in Fig. 3(a).

1. Methods

Adaptive proofreading is well suited to characterize the decision boundary between two classes because we can work with an analytical description. We want to know how to most efficiently change the binding time of the spurious binding ligand (with small τ) to cause the model to reach the decision boundary. We take inspiration from Ref. [50] and adapt our approach from the iterative FGSM [51]. At first, we sample the binding times τ_{self} for $L_{\text{self}} = 7000$ self-ligands from a half-normal distribution $|\mathcal{N}(0, \frac{1}{3})|$ and τ_{ag} for $L_{\text{ag}} = 3000$ agonist ligands from a narrowly peaked

normal distribution $|\mathcal{N}(\frac{7}{2}, \frac{1}{10})|$ just above $\tau_d = 3$. We fix the agonist ligand distribution, the signal in the immune picture. Next, we bin ligands in M equally spaced bins with center binding time $\tau_b, b \in 1, \dots, M$, and we compute the gradient for bins for which $\tau_b < \tau_d$,

$$\frac{\partial T_{N,m}}{\partial \tau_b} = \frac{N\tau_b^{N-1}L_b - mT_{N,m}\tau_b^{m-1}L_b}{\sum_{i=1}^M \tau_i^m L_i}, \quad (\text{C1})$$

where L_b is the number of ligands in the b th bin. We subtract this value multiplied by a small number ϵ from the exact binding times, as in Eq. (6) in the main text, and we compute a new output $T_{N,m}$. We repeat this procedure until $T_{N,m}$ dips just below the response threshold τ_d^{N-m} . We then display the ligand distributions. We bin ligands and compute the gradient in batches to prevent the gradient from becoming negligibly small. If we compute the gradient for each ligand with an individual binding time, there will be exactly one ligand with that specific binding time, and because the gradient scales with L , we need to go through many more iterations. Decreasing the bin size and step size ϵ may enhance the resolution, but it is not required. We find good results by considering bins with a bin size of $0.2s$ and $\epsilon = 0.2$.

2. MTL pictures

We can visually recast immune recognition as an image recognition problem by placing pixels on a grid and coloring them based on their binding time with a given scale. We choose to let white pixels correspond to not self ($\tau > \tau_d$), gray pixels to antagonist ligands ($\tau_a < \tau < \tau_d$),

and black pixels to self-ligands $\tau \ll \tau_a$. We are free to introduce any kind of spatial correlation to create “immune pictures” from a ligand distribution. This results in what we term Montreal pictures or “MTL pictures” (Fig. 5). The initial ligand distribution, MTL picture, and scale are given on the left. We perform iterative gradient descent like in the main text and plot the ligand distribution and the corresponding immune pictures at the boundary for various (N, m) . The results are striking. For a T cell operating in the adversarial regime, the signal MTL is unaltered at the decision boundary. At the transition $m = 2$, we see a slight change of color, while in the ambiguous regime, the signal actually changes from MTL to ML, where ML is short for machine learning. As we desire for a robust decision-maker, the response should switch when the signal becomes significantly different. From this, we conclude the only in the robust regime can Montreal turn fully into the city of machine learning.

For the MTL pictures in Fig. 5, we distribute the pixels in the 179×431 frame—equal to R , the number of receptors—as $L_{\text{self}} = 0.60R$, $L_a = 0.12R$, and $L_{\text{ag}} = 0.28R$. We sample τ_{self} from $|\mathcal{N}(0, \frac{1}{3})|$, τ_a from $\tau_d - |\mathcal{N}(0, 13)|$, τ_{ag} from $\tau_d + \mathcal{N}(\frac{7}{2}, \frac{1}{100})$, and we set $\tau_d = 3$. The picture is engineered such that the agonist ligands fill the M and the L, and the antagonists fill the T (which is why the T is slightly darker than the M and L). The self-ligands fill the area around the letters M, T, and L, such that the self ligands with highest binding time surround the T. We choose this example to make the effect of proofreading explicit (and of course because we are based in Montreal and study machine learning). This result is generic, and the

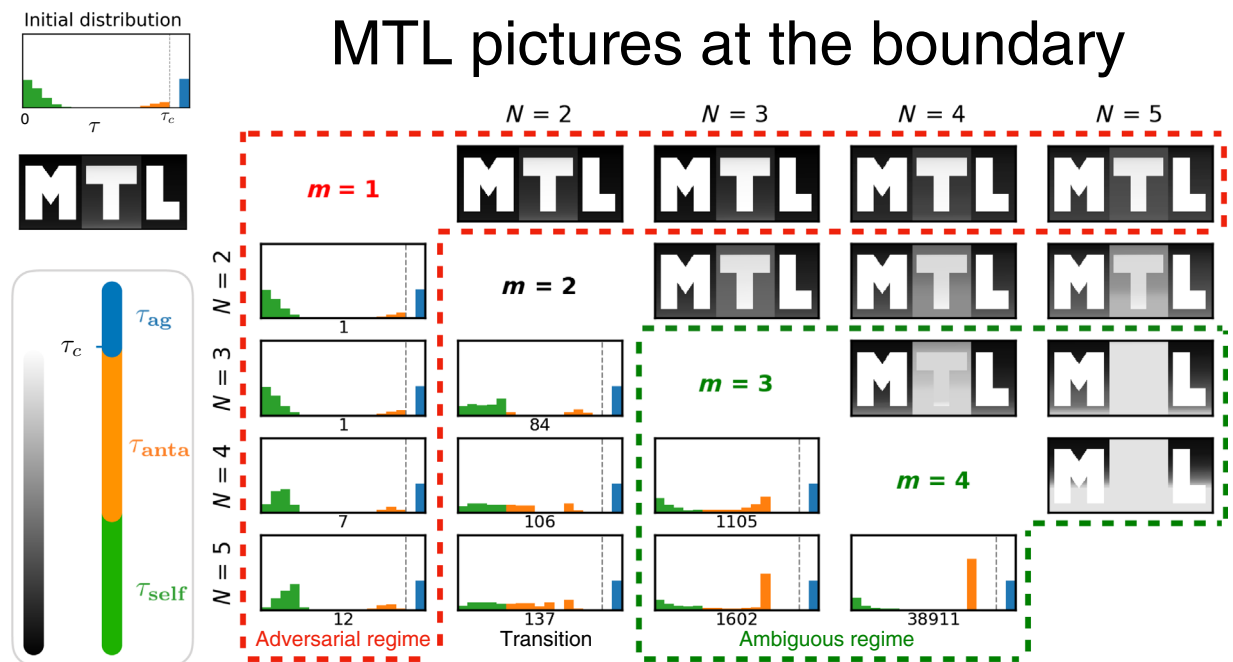


FIG. 5. MTL pictures. Explanation is found in the text.

ambiguity of instances at the decision boundary of a robust model can be visualized with any well-designed image [69].

3. Behavior for small binding times

Consider a mixture with L_{ag} ligands at $\tau_{\text{ag}} > \tau_d$ and L ligands with small binding time $\tau_{\text{spurious}} = \tau_e \ll \tau_{\text{ag}}$. To understand the behavior of $T_{N,m}$ as a function of τ_e , we expand $T_{N,m}$ in small variable $\epsilon = [(\tau_e)/(\tau_{\text{ag}})]$ as

$$\begin{aligned} T_{N,m}(\{L_{\text{ag}}, \tau_{\text{ag}}; L, \tau_e\}) &= \frac{\tau_{\text{ag}}^N L_{\text{ag}} + \tau_e^N L}{\tau_{\text{ag}}^m L_{\text{ag}} + \tau_e^m L} \\ &= \frac{1 + \epsilon^N \frac{L}{L_{\text{ag}}}}{1 + \epsilon^m \frac{L}{L_{\text{ag}}}} \tau_{\text{ag}}^{N-m} \\ &\simeq \left(1 + \epsilon^N \frac{L}{L_{\text{ag}}}\right) \left(1 - \epsilon^m \frac{L}{L_{\text{ag}}}\right) \tau_{\text{ag}}^{N-m} \\ &\simeq \tau_{\text{ag}}^{N-m} - \tau_{\text{ag}}^{N-m} \frac{L}{L_{\text{ag}}} \epsilon^m + O(\epsilon^N), \end{aligned}$$

which confirms that up to a constant $T_{N,m} \propto -\epsilon^m \propto -\tau_e^m$ for $m \geq 1$ and $\tau_e \ll \tau_{\text{ag}}$, as well as that

$$\frac{dT_{N,m}}{d\tau_e} \simeq -m \tau_{\text{ag}}^{N-m-1} \frac{L}{L_{\text{ag}}} \epsilon^{m-1} \propto -\tau_e^{m-1}. \quad (\text{C2})$$

APPENDIX D: BOUNDARY TILTING

To further draw the connection between machine learning and adaptive proofreading models, we study a framework to interpret adversarial examples called boundary tilting [52]. We first illustrate this effect on the discrimination of the original MNIST 3 vs 7 problem MNIST from Ref. [6] (Appendix D 1), after which we interpret boundary tilting via proofreading in ligand discrimination (Appendix D 2), and finally, we derive how the addition of a subthreshold ligand at the decision boundary changes the output (Appendix D 3). We refer to these results in the main text at the end of Sec. II E.

1. Digit classification

A typical 3 and 7 (i), the averages $\bar{3}$ and $\bar{7}$ (ii), and the corresponding adversarial examples (iii, iv) are shown in Fig. 6(a). Tanay and Griffin [52] pointed out that the adversarial perturbation generated with the FGSM proposed in Ref. [6] can also be found via $D = \text{sgn}(\bar{3} - \bar{7})$, Fig. 6(a) (v). Note the similarity to the adversarial perturbation from the FGSM $\text{sgn}(w) = \text{sgn}(\nabla_x J)$ [Fig. 6(a) (vi)]. To reveal the linearity of binary digit discrimination, we compute the principal components (PCs) of the traditional training set of 3's and 7's, and project all digits in the test set on PC_1 and PC_2 [Fig. 6(b)]. With a linear support vector classifier (ordinary linear regression) trained on the transformed coordinates PC_1 and PC_2 of the training set, we

achieve over 95% accuracy in the test set. While such an accuracy is far from the state of the art in digit recognition, it is much higher than typical detection accuracy for single cells (e.g., T cells present false negative rates of 10% for strong antagonists [15]). The red and blue stars in Fig. 6(s) denote the average digit $\bar{3}$, $\bar{7}$.

Next, we transform the test set as $3 \rightarrow 3' = 3 - \epsilon_{\text{test}} D$, $7 \rightarrow 7' = 7 + \epsilon_{\text{test}} D$, where $\epsilon_{\text{test}} = 0.4$ is the strength of the adversarial perturbation [Fig. 6(a) (iii)]. $\bar{3}'$ and $\bar{7}'$ move closer in Fig. 6(b), orthogonal to the decision boundary and along the line between the initial averages. This adversarial perturbation moves the digits in what we call an adversarial direction perpendicular to the decision boundary and reduces the accuracy of the linear regression model to a mere 69%.

Goodfellow *et al.* proposed adversarial training as a method to mitigate adversarial effects by FGSM. We implement adversarial training by adding the adversarial perturbation $\epsilon_{\text{train}} D_{\text{train}} = \epsilon_{\text{train}} (\bar{3}_{\text{train}} - \bar{7}_{\text{train}})$ to the images in the training set, computing the new PCs and training the linear regression model. Such adversarial training effectively ‘‘tilts’’ the decision boundary, while preserving 95% accuracy. In the presence of the original adversarial perturbations, we see the effect of the tilted boundary: The perturbation moves digits parallel along the decision boundary, which results in good robust accuracy. This is an illustrative example of the more general phenomenon studied in Ref. [52].

2. Boundary tilting and categorizing perturbations

We consider the change in $T_{N,m}$ for arbitrary N, m upon addition of many spurious ligands. Generalizing Eq. (2) in the main text gives

$$T_{N,m}^{\text{after}} = \frac{L(\tau - \epsilon)^N + \epsilon R \epsilon^N}{L\tau^m + \epsilon R \epsilon^m} = \frac{(\tau - \epsilon)^N + \frac{\epsilon^{N+1} R}{L}}{\tau^m + \frac{\epsilon^{m+1} R}{L}}. \quad (\text{D1})$$

From this expression, we note that $T_{N,m}$ is changing significantly with respect to its initial value upon addition of many weakly bound ligands as soon as $\epsilon^{m+1} R$ is of order L . Thus, the effect described in the main text for weighted averages where $(N, m) = (1, 0)$ also holds for nonlinear computations as long as m is small. It appears that the general strategy to defend against this adversarial perturbation is by increasing m , as previously observed in Ref. [33]. Biochemically, this is done with kinetic proofreading [15,32,36]; i.e., we take an output $T_{N,m}$ with $N > m \geq 1$. Here, the output is no longer sensitive to the addition of many weakly bound self-ligands, yielding an inversion of the antagonistic hierarchy where the strongest antagonizing ligands exist closer to threshold [34]. An extreme case has been proposed for immune recognition where the strongest antagonists are found just below the threshold of activation [15].

We numerically compute how the decision boundary changes when L_{self} ligands at τ_{self} are added to the initial

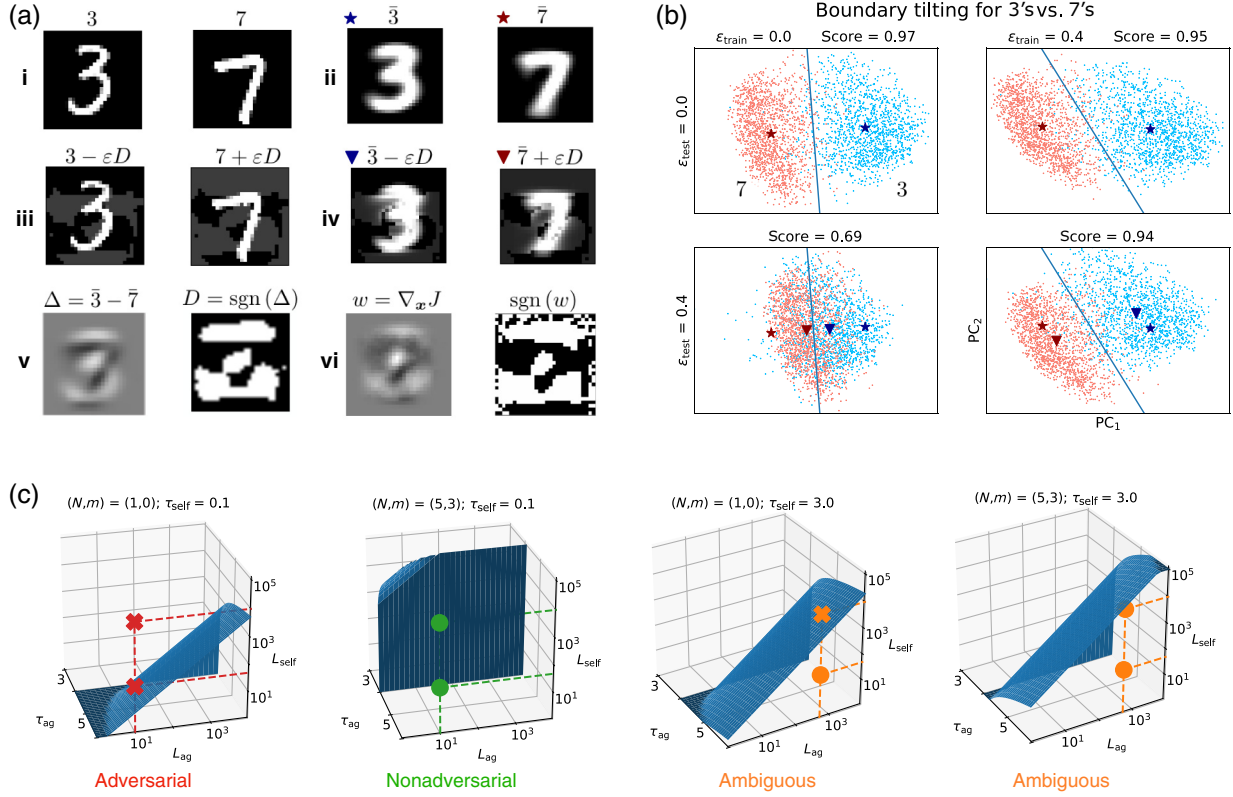


FIG. 6. Boundary tilting in one-dimensional digit classification. (a) (i) Typical 3 and 7 from MNIST. (ii) Average 3, 7 of the traditional test set, (iii, iv) with adversarial perturbation found by (v) subtracting the sign of $\bar{3}$ from $\bar{7}$, which corresponds to (vi), the perturbation found with FGSM. (b) Projection of the digits on the first principal components. The classes are separated by a linear support vector classifier (blue), and the average of the classes with and without adversarial perturbation is shown by the triangle and star. We cycle through permutations of adversarial training and/or adversarial testing. Note how the boundary tilts in the right panels and how the triangle moves parallel to the decision boundary. (c) Decision boundary of the immune model. The region under the surface is the response regime, and the region above is the no-response regime. The classifier with a single proofreading step $(N, m) = (1, 0)$ fails to observe agonists in three of the four marked mixtures, while the robust classifier $(N, m) = (5, 3)$ correctly responds to each indicated mixture.

L_{ag} agonist ligands at τ_{ag} ; i.e., we compute the manifold so that

$$T_{N,m}(\{L_{\text{ag}}, \tau_{\text{ag}}; L_{\text{self}}, \tau_{\text{self}}\}) = \frac{\tau_{\text{ag}}^N L_{\text{ag}} + \tau_{\text{self}}^N L_{\text{self}}}{\tau_{\text{ag}}^m L_{\text{ag}} + \tau_{\text{self}}^m L_{\text{self}}} \quad (\text{D2})$$

is equal to $T_{N,m}(\{L_{\text{ag}}, \tau_d\}) = \tau_d^{N-m}$. We represent this boundary for fixed τ_{self} and variable $L_{\text{ag}}, L_{\text{self}}, \tau_{\text{ag}}$ in Fig. 6(c). Boundary tilting is studied with respect to the reference $L_{\text{self}} = 0$ plane corresponding to the situation of pure L_{ag} ligands at τ_{ag} , where the boundary is the line $\tau_{\text{ag}} = \tau_d$. The case $(N, m) = (1, 0)$ [Fig. 6(c) left panel]

TABLE I. Categories of perturbations.

	Boundary tilting	Gradient when adding one antagonistic ligand
Adversarial	Yes	Steep [$\mathcal{O}(1)$]
Nonadversarial	No	Almost flat [$\mathcal{O}(e^m)$]
Ambiguous	Yes	Weak [$\mathcal{O}(\epsilon)$]

corresponds to a very tilted boundary, close to the plane $L_{\text{self}} = 0$, and a strong antagonistic case. In this situation, assuming $\tau_{\text{ag}} \simeq \tau_d$, each new ligand added with τ_{self} close to 0 gives a reduction of $T_{1,0}$ proportional to τ_d/L_{ag} in the limit of small L_{self} (see next section, Ref. [14]), which is again of the order of the response $T_{1,0} = \tau_{\text{ag}} \simeq \tau_d$ in the plane $L_{\text{self}} = 0$. This response is clearly not infinitesimal, corresponding to a steep gradient of $T_{1,0}$ in the L_{self} direction. We call the perturbation in this case adversarial. This should be contrasted to the case for higher m [Fig. 6(c), middle left) where the boundary is vertical, independent of L_{self} , such that decision-making is based only on the initially present L_{ag} ligands at τ_{ag} . Here, the change of response induced by the addition of each ligand with small binding time τ_{self} is τ_{self}^m due to proofreading a very small number when $\tau_{\text{self}} \simeq 0$ [14]. Contrary to the previous case, the gradient of $T_{N,m}$ with respect to this vertical direction is almost flat and very small compared to the response in the $L_{\text{self}} = 0$ plane. We call the perturbation in this case nonadversarial.

Tilting of the boundary occurs only when τ_{self} gets sufficiently close to the threshold binding time τ_d [Fig. 6(c), right panels]. In this regime, each new ligand added with quality $\tau_{\text{self}} = \tau_d - \epsilon$ contributes an infinitesimal change of $T_{N,m}$ proportional to $[(\tau_d - \tau_{\text{self}})/L_{\text{ag}}] = \epsilon/L_{\text{ag}}$, which gives a weak gradient in the direction L_{self} . But even with such small perturbations one can easily cross the boundary because of the proximity of τ_{self} to τ_d , which explains the tilting. The cases where the boundary is tilted and the gradient is weak are of a different nature compared to the adversarial case of Fig. 6(c), left panel. Here, the boundary is tilted as well, but the gradient is steep, not weak. For this reason, we term the cases in the right panels ambiguous. Similar ambiguity is observed experimentally: It is well known that antagonists (ligands close to thresholds) also weakly agonize an immune response [15]. Our categorization of perturbations is presented in Table I [70].

3. Gradient in the L_2 direction

We recall results from Ref. [34] to show how the addition of subthreshold ligands one at a time changes the output. We first consider $\{L, \tau_d\}$ threshold ligands with output

$$T_{N,m}(L, \tau_d) = \tau_d^{N-m}. \quad (\text{D3})$$

The main result of Ref. [34] is the linear response of $T_{N,m}(L, \tau_d)$ to the addition of $\{L_a, \tau_d - \epsilon\}$ subthreshold ligands,

$$T_{N,m}(\{L, \tau_d; L_a, \tau_d - \epsilon\}) = T(L + L_a, \tau_d) - \epsilon L_a \mathcal{A}(L + L_a, \tau_d) \quad (\text{D4})$$

$$= \tau_d^{N-m} - \epsilon \frac{L_a}{L + L_a} \frac{d}{d\tau} T_{N,m}(L + L_a, \tau)|_{\tau=\tau_d}, \quad (\text{D5})$$

where we use the definition

$$\mathcal{A}(L, \tau_d) = \frac{1}{L} \frac{d}{d\tau} T_{N,m}(L, \tau)|_{\tau=\tau_d} \quad (\text{D6})$$

for the coefficient in a mean-field description. As the derivative $[d/(d\tau)]T_{N,m}(L, \tau)|_{\tau=\tau_d} > 0$ and $\epsilon = \tau_a - \tau_d$, each additional subthreshold ligand at τ_a decreases the output with a value proportional to

$$\frac{\tau_d - \tau_a}{L}. \quad (\text{D7})$$

In the case $(N, m) = (1, 0)$, the mean-field approximation is exact; i.e., the first derivative of $(dT)/(d\tau)$ is the only nonzero derivative given by

$$\mathcal{A}(L, \tau_d) = \frac{1}{L} \frac{d}{d\tau} \tau \Big|_{\tau=\tau_d} = \frac{1}{L}. \quad (\text{D8})$$

With the addition of a single subthreshold ligand $\tau_a \simeq 0$, so that $\epsilon \simeq \tau_d$, the output is maximally reduced by $[\tau_d/(L+1)] \simeq (\tau_d/L)$, a finite quantity, as we describe in the main text. For higher m , the linear approximation holds only for ligands at τ_a close to threshold.

APPENDIX E: FEW-PIXEL ATTACK

In this Appendix, we describe in detail the procedure for the few-pixel attack. We use this to come to our conclusion in Sec. II G and Fig. 4(c) in the main text.

The few-pixel attack connects to ligand antagonism in the sense that few pixels are needed to cause misclassification, corresponding to the addition of few maximally antagonizing ligands to a mixture fooling robust adaptive proofreading models. It is not the most efficient attack against a classifier without biomimetic defense, but it is the most efficient attack against classifiers with biomimetic defense, equivalent to adaptive proofreading models with $m > 1$. For these adaptive proofreading models, there exists a unique maximally antagonistic binding time defined as the binding time that maximally reduces $T_{N,m}$.

With this idea in mind, we decide to make pixels black or white in a controlled manner until the neural network classifies the perturbed initial digit as the target class. In the following, we refer to several stages of the few-pixel attack using Fig. 7. We first compute what we term pixel maps. Pixel maps contain the change of score when making a pixel white or black. In Fig. 7, blue colors correspond to pixels that will lower the score when turned white or black, while red colors are for pixels that will increase the score for the same operation. A gray color means the score is unchanged when whitening or blacking the pixel. The pixel maps are scaled to the maximum change in score. We proceed in merging and sorting the pixel maps from maximum to minimum change in score towards the target class, iteratively following the sorted list to decide which pixels in our digit to turn white or black. We do this until we reach the decision boundary (first iteration in which the digit is misclassified). The final digits in the row above the red rectangle in Fig. 7 are the resulting boundary digits. They already contain perturbations corresponding to real features but have an air of artificiality to them which allows us to fairly easily distill the ground truth. We further apply mean filtering [54], which is a 3×3 convolutional block that computes mean pixel values as

$$y_{i,j} = \frac{1}{9} \sum_{k,l=-1}^1 x_{i+k,j+l}. \quad (\text{E1})$$

Biologically, mean filtering is pure receptor clustering, where a perturbation to a single receptor locally affects other ligands. Such digits are truly ambiguous digits that are tough to classify even as humans. These are the type of digits we expect to find on the decision boundary. Finally,

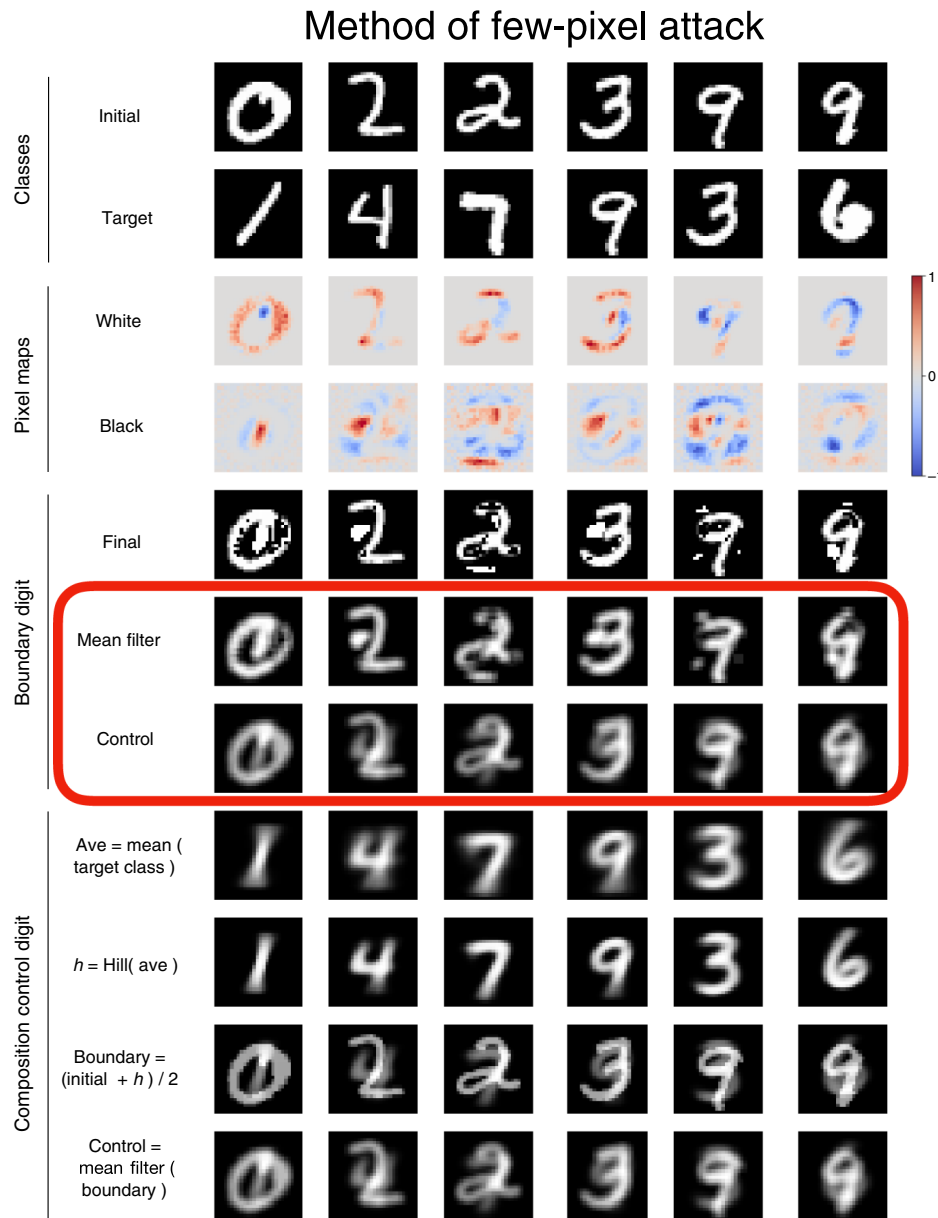


FIG. 7. Method of few-pixel attack. Each column shows how a few-pixel attack causes misclassification of an initial digit to a target class. The important result is the prefiltered boundary digits and the control in the red rectangle. Pixel maps determine which pixels increase (red) or decrease (blue) the score when turning an individual pixel in the initial digit white or black. We merge the pixel maps, sort this list of pixels, and go through it from maximum to minimum change in score until misclassification occurs, resulting in the prefiltered digit. We apply a mean filter to make them look more like real digits, and indeed, these mean-filtered boundary digits closely resemble our control digits at the boundary. The control digits are composed of the mean-filtered initial digit plus locally contrasted [with hill function ($N = 3$; $\theta = 0.5$)] average digit of the target class.

we compare the mean-filtered digit at the decision boundary to the control: The sum of the initial digit and the hill function of Eq. (8) ($N = 3$; $\theta = 0.5$) on the average of all digits in the target class, then mean filter (Fig. 7 for a step-by-step composition). We apply the mean filter to the control to again remove the artificiality of a digit plus an average and make the comparison between boundary digit and control digit fairer. The similarity between the mean-filtered boundary digit and

control digit confirms our intuition that we are actually operating in the space between both classes when misclassification occurs.

We can also apply the mean filter to the initial digit before generating the pixel maps, and during the procedure, check the score on the mean-filtered perturbed image. This gives similar results, as we see by following the trajectory of the score for *boundary null* and *boundary mean*. We show the score explicitly in Fig. 8 for the digits

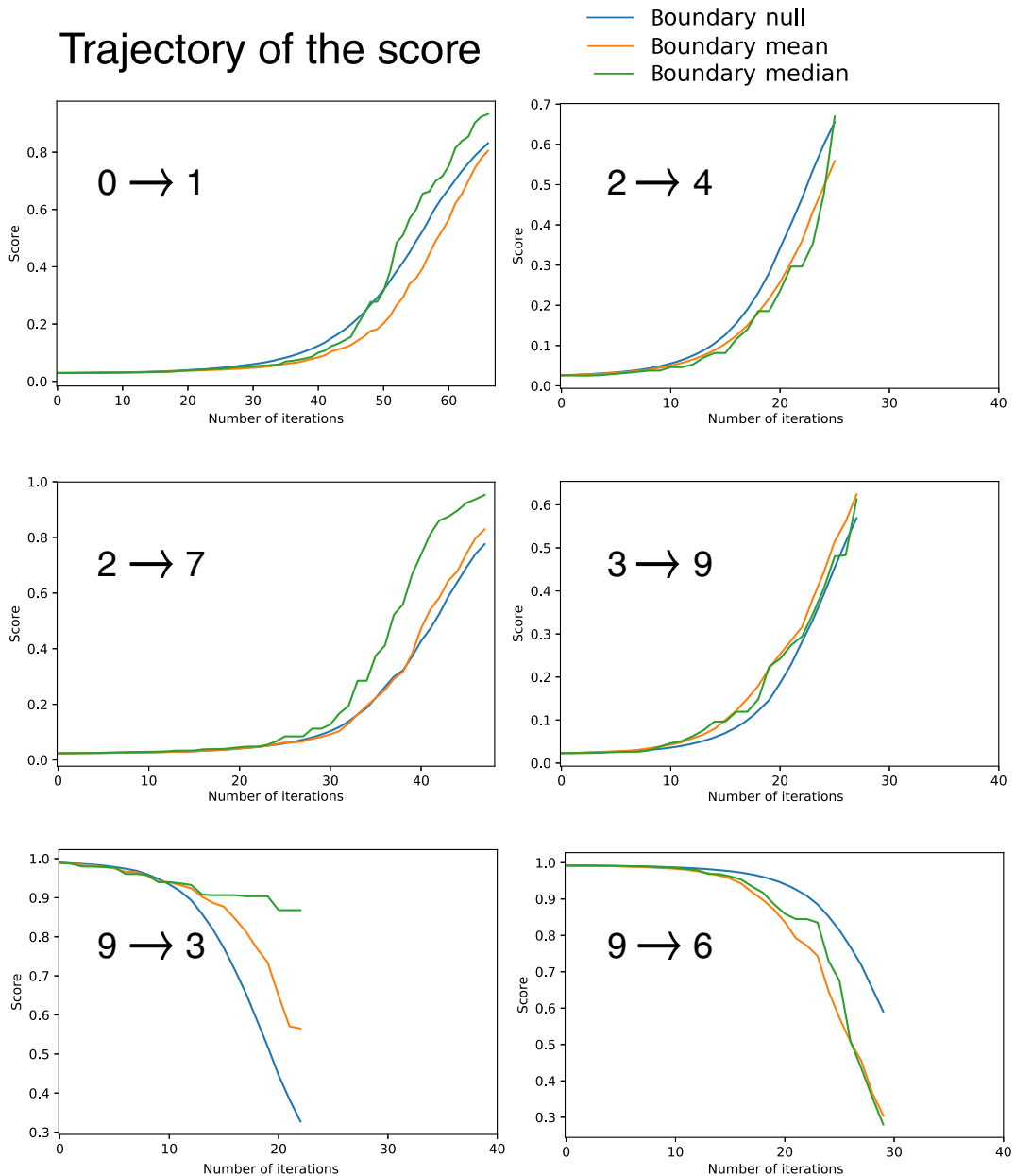


FIG. 8. Trajectory of the scoring functions of the attacks in Fig. 7. The blue, orange, and green lines correspond to various digits (actual digit, mean-filtered digit, median-filtered digit) for which we check the score and terminate when reaching the boundary. The trajectory of the score for the null digit and the mean-filtered digit is generally the same. Moreover, the behavior of the score looks similar to the behavior of $T_{N,m}$ upon addition of maximally antagonizing ligands to a mixture of only agonist ligands in Fig. 3(d) in the main text.

in Fig. 7. The behavior of the score is remarkably similar to the interpolation between ligand mixtures [Fig. 3(f), bottom panel in the main text]. A nonlinear filtering method proposed in Ref. [54] is the median filter, but this one works less well for black and white pixels.

We show examples that are generated when we select for instances where the number of iterations is large enough (20 suffice, we still consider this to be a few-pixel attack, keeping in mind that digits have 784 individual pixels). The authors of Ref. [53] specifically

searched for single-pixel attacks. Examples of single-pixel misclassification exist in our neural networks trained on two types of digits in MNIST too, but these we find noninformative. In cellular decision-making, this case corresponds to adding a single antagonist ligand to a ligand mixture to cause misclassification. This is possible only if the ligand mixture is already very close to the boundary. For such samples, we do not expect ambiguity to appear. Remember that near the boundary, the score landscape is steep, and small additions have a large effect.

- [1] Y. LeCun, Y. Bengio, and G. Hinton, *Deep Learning*, *Nature (London)* **521**, 436 (2015).
- [2] A. Krizhevsky, I. Sutskever, and G. E. Hinton, *Imagenet Classification with Deep Convolutional Neural Networks*, in *Proceedings of Advances in Neural Information Processing Systems Conference, Lake Tahoe, Nevada, 2012* (Curran Associates Inc., USA, 2012), pp. 1097–1105.
- [3] G. Hinton, L. Deng, D. Yu, G. E. Dahl, A.-r. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath *et al.*, *Deep Neural Networks for Acoustic Modeling in Speech Recognition: The Shared Views of Four Research Groups*, *IEEE Signal Process. Mag.* **29**, 82 (2012).
- [4] I. Sutskever, O. Vinyals, and Q. V. Le, *Sequence to Sequence Learning with Neural Networks*, in *Proceedings of Advances in Neural Information Processing Systems Conference* (2014), pp. 3104–3112, <http://papers.nips.cc/paper/5346-sequence-to-sequence-learning-with-neural>.
- [5] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, *Intriguing Properties of Neural Networks*, [arXiv:1312.6199](https://arxiv.org/abs/1312.6199).
- [6] I. J. Goodfellow, J. Shlens, and C. Szegedy, *Explaining and Harnessing Adversarial Examples*, [arXiv:1412.6572](https://arxiv.org/abs/1412.6572).
- [7] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, and A. Swami, *Practical Black-Box Attacks against Machine Learning*, in *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security* (ACM, Abu Dhabi, United Arab Emirates, 2017), pp. 506–519.
- [8] S. G. Finlayson, I. S. Kohane, and A. L. Beam, *Adversarial Attacks against Medical Deep Learning Systems*, [arXiv:1804.05296](https://arxiv.org/abs/1804.05296).
- [9] N. Akhtar and A. Mian, *Threat of Adversarial Attacks on Deep Learning in Computer Vision: A Survey*, [arXiv:1801.00553](https://arxiv.org/abs/1801.00553).
- [10] S. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard, *Universal Adversarial Perturbations*, in *Proceedings of the 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (IEEE, New York, 2017), pp. 86–94.
- [11] E. D. Siggia and M. Vergassola, *Decisions on the Fly in Cellular Sensory Systems*, *Proc. Natl. Acad. Sci. U.S.A.* **110**, E3704 (2013).
- [12] N. R. J. Gascoigne, T. Zal, and S. M. Alam, *T-Cell Receptor Binding Kinetics in T-Cell Development and Activation*, *Expert Rev. Mol. Med.* **3**, 1 (2001).
- [13] O. Feinerman, R. N. Germain, and G. Altan-Bonnet, *Quantitative Challenges in Understanding Ligand Discrimination by $\alpha\beta$ T Cells*, *Molecular immunology* **45**, 619 (2008).
- [14] P. François and G. Altan-Bonnet, *The Case for Absolute Ligand Discrimination: Modeling Information Processing and Decision by Immune T Cells*, *J. Stat. Phys.* **162**, 1130 (2016).
- [15] G. Altan-Bonnet and R. N. Germain, *Modeling T Cell Antigen Discrimination Based on Feedback Control of Digital ERK Responses*, *PLOS Biol.* **3**, e356 (2005).
- [16] C. Torigoe, J. K. Inman, and H. Metzger, *An Unusual Mechanism for Ligand Antagonism*, *Science* **281**, 568 (1998).
- [17] G. Reddy, J. D. Zak, M. Vergassola, and V. N. Murthy, *Antagonism in Olfactory Receptor Neurons and Its Implications for the Perception of Odor Mixtures*, *eLife* **7**, e34958 (2018).
- [18] J. Tsitron, A. D. Ault, J. R. Broach, and A. V. Morozov, *Decoding Complex Chemical Mixtures with a Physical Model of a Sensor Array*, *PLoS Comput. Biol.* **7**, e1002224 (2011).
- [19] P. Klenerman, S. Rowland-Jones, S. McAdam, J. Edwards, S. Daenke, D. Laloo, B. Köppe, W. Rosenberg, D. Boyd, A. Edwards *et al.*, *Cytotoxic T-Cell Activity Antagonized by Naturally Occurring HIV-1 Gag Variants*, *Nature (London)* **369**, 403 (1994).
- [20] U.-C. Meier, P. Klenerman, P. Griffin, W. James, B. Köppe, B. Larder, A. McMichael, and R. Phillips, *Cytotoxic T Lymphocyte Lysis Inhibited by Viable HIV Mutants*, *Science* **270**, 1360 (1995).
- [21] S. J. Kent, P. D. Greenberg, M. C. Hoffman, R. E. Akridge, and M. J. McElrath, *Antagonism of Vaccine-Induced HIV-1-Specific CD4+ T Cells by Primary HIV-1 Infection: Potential Mechanism of Vaccine Failure*, *J. Immunol.* **158**, 807 (1997).
- [22] A. Snyder *et al.*, *Genetic Basis for Clinical Response to CTLA-4 Blockade in Melanoma*, *N. Engl. J. Med.* **371**, 2189 (2014).
- [23] T. N. Schumacher and R. D. Schreiber, *Neoantigens in Cancer Immunotherapy*, *Science* **348**, 69 (2015).
- [24] K. Grosse, P. Manoharan, N. Papernot, M. Backes, and P. McDaniel, *On the (Statistical) Detection of Adversarial Examples*, [arXiv:1702.06280](https://arxiv.org/abs/1702.06280).
- [25] E. Wong and Z. Kolter, *Provable Defenses against Adversarial Examples via the Convex Outer Adversarial Polytope*, in *Proceedings of the International Conference on Machine Learning* (2018), pp. 5283–5292, <http://proceedings.mlr.press/v80/wong18a.html>.
- [26] D. Krotov and J. J. Hopfield, *Dense Associative Memory for Pattern Recognition*, in *Proceedings of the Advances in Neural Information Processing Systems Conference* (2016), pp. 1172–1180, <https://papers.nips.cc/paper/6121-dense-associative-memory-for-pattern-recognition>.
- [27] J. Das, *Activation or Tolerance of Natural Killer Cells Is Modulated by Ligand Quality in a Nonmonotonic Manner*, *Biophys. J.* **99**, 2028 (2010).
- [28] F. Lagarde, C. Beausoleil, S. M. Belcher, L. P. Belzunces, C. Emond, M. Guerbet, and C. Rousselle, *Non-Monotonic Dose-Response Relationships and Endocrine Disruptors: A Qualitative Method of Assessment*, *Environ. Health* **14**, 13 (2015).
- [29] C. C. Govern, M. K. Paczosa, A. K. Chakraborty, and E. S. Huseby, *Fast On-Rates Allow Short Dwell Time Ligands to Activate T Cells*, *Proc. Natl. Acad. Sci. U.S.A.* **107**, 8724 (2010).
- [30] A. K. Chakraborty and A. Weiss, *Insights into the Initiation of TCR Signaling*, *Nat. Rev. Immunol.* **15**, 798 (2014).
- [31] M. Lever, H.-S. Lim, P. Kruger, J. Nguyen, N. Trendel, E. Abu-Shah, P. K. Maini, P. A. van der Merwe, and O. Dushek, *Architecture of a Minimal Signaling Pathway Explains the T-Cell Response to a 1 Million-Fold Variation in Antigen Affinity and Dose*, *Proc. Natl. Acad. Sci. U.S.A.* **113**, E6630 (2016).

- [32] P. François, G. Voisinne, E. D. Siggia, G. Altan-Bonnet, and M. Vergassola, *Phenotypic Model for Early T-Cell Activation Displaying Sensitivity, Specificity, and Antagonism*, *Proc. Natl. Acad. Sci. U.S.A.* **110**, E888 (2013).
- [33] J.-B. Lalanne and P. François, *Principles of Adaptive Sorting Revealed by In Silico Evolution*, *Phys. Rev. Lett.* **110**, 218102 (2013).
- [34] P. François, M. Hemery, K. A. Johnson, and L. N. Saunders, *Phenotypic Spandrel: Absolute Discrimination and Ligand Antagonism*, *Phys. Biol.* **13**, 066011 (2016).
- [35] F. Proulx-Giraldeau, T. J. Rademaker, and P. François, *Untangling the Hairball: Fitness-Based Asymptotic Reduction of Biological Networks*, *Biophys. J.* **113**, 1893 (2017).
- [36] T. W. McKeithan, *Kinetic Proofreading in T-Cell Receptor Signal Transduction*, *Proc. Natl. Acad. Sci. U.S.A.* **92**, 5042 (1995).
- [37] G. J. Kersh, E. N. Kersh, D. H. Fremont, and P. M. Allen, *High- and Low-Potency Ligands with Similar Affinities for the TCR: The Importance of Kinetics in TCR Signaling*, *Immunity* **9**, 817 (1998).
- [38] J. J. Hopfield, *Kinetic Proofreading: A New Mechanism for Reducing Errors in Biosynthetic Processes Requiring High Specificity*, *Proc. Natl. Acad. Sci. U.S.A.* **71**, 4135 (1974).
- [39] J. Ninio, *Kinetic Amplification of Enzyme Discrimination*, *Biochimie* **57**, 587 (1975).
- [40] J.-B. Lalanne and P. François, *Chemodetection in Fluctuating Environments: Receptor Coupling, Buffering, and Antagonism*, *Proc. Natl. Acad. Sci. U.S.A.* **112**, 1898 (2015).
- [41] T. Mora, *Physical Limit to Concentration Sensing Amid Spurious Ligands*, *Phys. Rev. Lett.* **115**, 038102 (2015).
- [42] M. Carballo-Pacheco, J. Desponds, T. Gavrilchenko, A. Mayer, R. Prizak, G. Reddy, I. Nemenman, and T. Mora, *Receptor Crosstalk Improves Concentration Sensing of Multiple Ligands*, *Phys. Rev. E* **99**, 022423 (2019).
- [43] R. N. Germain and I. Stefanová, *The Dynamics of T Cell Receptor Signaling: Complex Orchestration and the Key Roles of Tempo and Cooperation*, *Annu. Rev. Immunol.* **17**, 467 (1999).
- [44] B. N. Dittel, IrenaŠtefanova, R. N. Germain, C. A. Janeway, Jr., *Cross-Antagonism of a T Cell Clone Expressing Two Distinct T Cell Receptors*, *Immunity* **11**, 289 (1999).
- [45] Y. LeCun and C. Cortes, *The MNIST Database of Handwritten Digits*, <http://yann.lecun.com/exdb/mnist/>, 1998.
- [46] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, *Scikit-learn: Machine Learning in PYTHON*, *J. Mach. Learn. Res.* **12**, 2825 (2011).
- [47] The on rate is easily confused with the unbinding rate, whose inverse we call the binding time, which indicates the lifetime of the ligand-receptor complex.
- [48] An alternative choice without loss of generality is to consider a situation where for unoccupied receptors, k_i is 0 but τ_i is arbitrary, corresponding to a ligand available for binding.
- [49] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Madry, *Robustness May Be at Odds with Accuracy*, [arXiv:1805.12152](https://arxiv.org/abs/1805.12152).
- [50] D. Krotov and J. J. Hopfield, *Dense Associative Memory is Robust to Adversarial Inputs*, *Neural Comput.* **30**, 3151 (2018).
- [51] A. Kurakin, I. Goodfellow, and S. Bengio, *Adversarial Machine Learning at Scale*, [arXiv:1611.01236](https://arxiv.org/abs/1611.01236).
- [52] T. Tanay and L. Griffin, *A Boundary Tilting Perspective on the Phenomenon of Adversarial Examples*, [arXiv:1608.07690](https://arxiv.org/abs/1608.07690).
- [53] J. Su, D. V. Vargas, and S. Kouichi, *One Pixel Attack for Fooling Deep Neural Networks*, [arXiv:1710.08864](https://arxiv.org/abs/1710.08864).
- [54] C. Xie, Y. Wu, L. van der Maaten, A. Yuille, and K. He, *Feature Denoising for Improving Adversarial Robustness*, [arXiv:1812.03411](https://arxiv.org/abs/1812.03411).
- [55] D. Berthelot, C. Raffel, A. Roy, and I. Goodfellow, *Understanding and Improving Interpolation in Autoencoders via an Adversarial Regularizer*, [arXiv:1807.07543](https://arxiv.org/abs/1807.07543).
- [56] S.-M. Moosavi-Dezfooli, A. Fawzi, J. Uesato, and P. Frossard, *Robustness via Curvature Regularization, and Vice Versa*, [arXiv:1811.09716](https://arxiv.org/abs/1811.09716).
- [57] D. Krotov and J. J. Hopfield, *Unsupervised Learning by Competing Hidden Units*, *Proc. Natl. Acad. Sci. U.S.A.* **116**, 7723 (2019).
- [58] G. Elsayed, S. Shankar, B. Cheung, N. Papernot, A. Kurakin, I. Goodfellow, and J. Sohl-Dickstein, *Adversarial Examples that Fool Both Computer Vision and Time-Limited Humans*, in *Proceedings of Advances in Neural Information Processing Systems Conference* (2018), pp. 3911–3921, <https://papers.nips.cc/paper/7647-adversarial-examples-that-fool-both-computer-vision-and-time-limited-humans>.
- [59] E. R. Unanue, *Altered Peptide Ligands Make Their Entry*, *J. Immunol.* **186**, 7 (2011).
- [60] T. C. Butler, M. Kardar, and A. K. Chakraborty, *Quorum Sensing Allows T Cells to Discriminate between Self and Nonself*, *Proc. Natl. Acad. Sci. U.S.A.* **110**, 11833 (2013).
- [61] G. Voisinne, B. G. Nixon, A. Melbinger, G. Gasteiger, M. Vergassola, and G. Altan-Bonnet, *T Cells Integrate Local and Global Cues to Discriminate between Structurally Similar Antigens*, *Cell Rep.* **11**, 1208 (2015).
- [62] J. N. Mandl, J. P. Monteiro, N. Vriskoop, and R. N. Germain, *T Cell-Positive Selection Uses Self-Ligand Binding Strength to Optimize Repertoire Recognition of Foreign Antigens*, *Immunity* **38**, 263 (2013).
- [63] M. Łuksza, N. Riaz, V. Makarov, V. P. Balachandran, M. D. Hellmann, A. Solovyov, N. A. Rizvi, T. Merghoub, A. J. Levine, T. A. Chan, J. D. Wolchok, and B. D. Greenbaum, *A Neoantigen Fitness Model Predicts Tumour Response to Checkpoint Blockade Immunotherapy*, *Nature (London)* **551**, 517 (2017).
- [64] U. Sahin and Ö. Türeci, *Personalized Vaccines for Cancer Immunotherapy*, *Science* **359**, 1355 (2018).
- [65] J. Yan, M. Deforet, K. E. Boyle, R. Rahman, R. Liang, C. Okegbe, L. E. P. Dietrich, W. Qiu, and J. B. Xavier, *Bow-Tie Signaling in c-di-GMP: Machine Learning in a Simple Biochemical Network*, *PLoS Comput. Biol.* **13**, e1005677 (2017).
- [66] A. Laan and G. de Polavieja, *Sensory Cheating: Adversarial Body Patterns Can Fool a Convolutional Visual System during Signaling*, [bioRxiv, https://doi.org/10.1101/326652](https://doi.org/10.1101/326652) (2018).

- [67] A. Athalye, L. Engstrom, A. Ilyas, and K. Kwok, *Synthesizing Robust Adversarial Examples*, in *Proceedings of the 35th International Conference on Machine Learning* (2018), Vol. 80, pp. 284–293, <http://proceedings.mlr.press/v80/athalye18b.html>.
- [68] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, and D. Song, *Robust Physical-World Attacks on Deep Learning Models*, in *Proceedings of the Conference on Computer Vision and Pattern Recognition* (2018), <https://arxiv.org/abs/1707.08945>.
- [69] Scripts to reproduce Figs. 4(a) and 5 are available at <https://github.com/tjrademaker/advxs-antagonism-figs/>.
- [70] Scripts for boundary tilting in ligand discrimination and digit discrimination are available at <https://github.com/tjrademaker/advxs-antagonism-figs/>.