# Theoretical and computational analysis of cell migration in complex tissue environments

Thesis by
Zitong (Jerry) Wang

In Partial Fulfillment of the Requirements for the
Degree of
Doctor of Philosophy

**Caltech**

CALIFORNIA INSTITUTE OF TECHNOLOGY
Pasadena, California

2025
Defended May 31, 2024

# ACKNOWLEDGEMENTS

The past five years of pursuing my PhD have been a rollercoaster of ups and downs, but I'm glad to have spent it in a place where I truly felt at home. One of the most unique aspects of this journey at Caltech was the opportunity to think about things purely for the joy of thinking about them, even if they might seem inconsequential in the grand scheme of things. More importantly, I got to explore these random curiosities with people who are wonderfully passionate and unbelievably kind. It saddens me to realize that I may never have an experience quite like this again, but at the same time, I'm incredibly grateful for having had this privilege that many others do not get.

To my advisor, Matt, I want to express my deepest gratitude for being an incredible source of support throughout my entire PhD journey. Your patience and respect for your students' ability to figure things out on their own, coupled with your knack for providing just the right amount of guidance, have been instrumental in building my confidence as an independent researcher. Your consideration for the challenges I faced outside of work and your willingness to help in any way possible made you not just an exceptional advisor but also a true mentor. I firmly believe that without your guidance, I wouldn't have found such joy in my PhD journey or had the courage to continue my path in academic research. If more advisors in academia were like you, I'm sure many of the issues faced by graduate students would be significantly reduced.

To my amazing wife, thank you for being by my side since day one of my PhD. We've faced our fair share of challenges, especially during the two and a half years of separation due to the pandemic, but your unwavering support for my work never faltered. Your decision to leave your job in China to join me in the United States is a testament to your love and commitment to our relationship. I'm so grateful for the sacrifices you've made and for being my rock throughout this incredible journey.

To my Canadian brothers, Abdullah and Han, thank you both for being awesome friends and colleagues I had the privilege of meeting during grad school. Thank you for all the life chats, SGV food runs, science discussions, and for all the support you have given me. The two of you never hesitate to say the difficult things, calling me out every time I am being unreasonable or immature, of which there were many. Your presence during the major milestones in my life has been invaluable, and I

will work hard to make sure we continue to be part of each other's lives for years to come.

To my amazing labmates, thank you all for being apart of my life for the last five years, for sharing with me countless meals, laughs, gossips, and advice. Especially Guru, who sat next to me for four years and helped me navigate the early stages of my PhD. Guru, your guidance in effectively communicating with Matt and your dedication to your own work were a constant source of inspiration, pushing me to work harder and aim higher in my own research. Even though I will probably never achieve your level of self-discipline, simply knowing that such discipline is humanly possible is an invaluable gift in itself.

To everyone who has been a part of my thesis and candidacy committee, Lior Pachter, Long Cai, Frederick Eberhardt, Akil Merchant, Michael Elowitz, Erik Winfree, and Markus Meister. I want to express my sincere gratitude for your invaluable contributions to my academic journey. Your insights, feedback, and guidance have been instrumental in shaping my research and helping me grow as a scientist. I appreciate the time and effort you've invested in reviewing my work, attending my presentations, and providing constructive criticism. I feel privileged to have had the opportunity to learn from such accomplished and knowledgeable individuals.

Mom and Dad, I want to thank you for your unconditional love and support, which have been the foundation of my success. Even when I chose paths that differed from your expectations, you never hesitated to encourage me and believe in my abilities. Pursuing a PhD was infinitely easier knowing that I had your backing, no matter what.

As I wrap up this chapter of my life and embark on new adventures, I'll always carry with me the lessons, memories, and relationships that have shaped me during my time at Caltech. To everyone who has been a part of this incredible journey, I offer my heartfelt gratitude. Your impact on my life is immeasurable, and I'm forever thankful for your support, guidance, and love.

I wish to dedicate this thesis to my little brother Jeremy for making fun of me and my paper when it got desk rejected, now your name is forever attached to my work.

# ABSTRACT

Cells sense and respond in spatially structured environments, including soils and tissue. My Ph.D. projects centered on developing new theoretical models and computational methods to understand how cells migrate in complex environments.

The first project is more theoretical in nature, leveraging information theory to study how the spatial organization of cell signaling pathways are adapted to the cell's natural environment. In tissue and soil, cells must localize to their targets by navigating distributions of extracellular ligands that are spatially discontinuous, consisting of local concentration peaks, due to binding a non-uniform network of ECM fibers. It is unclear how cells navigate patchy environments while not getting trapped in local concentration peaks. To answer this question, we framed navigation as a problem of maximizing mutual information in space and developed a computational algorithm for computing signaling pathway architectures that maximize mutual information in simulated natural environments. We found that for cells in tissues and soils, dynamic localization of membrane receptors dramatically boosts sensing precision and enables cells to navigate to chemical sources 30 times faster, but this receptor localization strategy is relatively inconsequential for cells in purely diffusive environments. Further, we found that anisotropic receptor dynamics previously observed in immune cells and growth cones are nearly optimal as predicted by our model.

The second project is more computational in nature, leveraging multiplexed tissue imaging to understand T-cell migration in tumor microenvironments. Immunotherapies can halt or slow down cancer progression by activating either endogenous or engineered T-cells to detect and kill cancer cells. T-cells must infiltrate the tumor core for immunotherapies to be effective. However, many solid tumors resist T-cell infiltration, challenging the efficacy of current therapies. In collaboration with clinician scientists at Cedars-Sinai Medical Center, we developed an integrated deep learning framework, Morpheus, that takes large-scale spatial omics profiles of patient tumors, and combines a formulation of T-cell infiltration prediction as a self-supervised machine learning problem with a counterfactual optimization strategy to generate minimal tumor perturbations predicted to boost T-cell infiltration. We applied Morpheus to 368 metastatic melanoma and colorectal cancer samples assayed using 40-plex imaging mass cytometry, discovering cohort-dependent, combinatorial perturbations, involving CXCL9, CXCL10, CCL22 and CCL18 for melanoma and CXCR4, PD-1, PD-L1 and CYR61 for colorectal cancer, predicted to support

T-cell infiltration across large patient cohorts. Using only raw image data, Morpheus also identified distinct therapeutic strategies for different patient strata such as cancer stage or fatty liver presence. Our work presents a paradigm for counterfactual-based prediction and design of cancer therapeutics using spatial omics data.

# PUBLISHED CONTENT AND CONTRIBUTIONS

Wang, Zitong Jerry, Alexander M Xu, et al. (2023). "Generating counterfactual explanations of tumor spatial proteomes to discover effective strategies for enhancing immune infiltration". In: *bioRxiv*. DOI: 10.1101/2023.10.12.562107.
Z.J.W. conceived the project, developed the software package, carried out the analysis, and wrote the paper. (**First and co-corresponding author**).

Wang, Zitong Jerry, and Matt Thomson (2022). "Localization of signaling receptors maximizes cellular information acquisition in spatially structured natural environments". In: *Cell Systems* 13.7, pp. 530–546. DOI: 10.1016/j.cels.2022.05.004.
Z.J.W. conceived the project, carried out the analysis, developed the simulation codes, derived the theoretical results, and wrote the paper. (**First and co-corresponding author**).

# TABLE OF CONTENTS

# LIST OF ILLUSTRATIONS

# LIST OF TABLES

*C h a p t e r   1*

# INTRODUCTION

## 1.1   Understanding cell migration in naturalistic environments

Traditional approaches to studying cell navigation often use highly simplified environmental models, where signals are either uniform or monotonic, neglecting the complex spatial structure in natural cell environments (Berg, Howard C, and Purcell, 1977; Hu et al., 2010; Mugler, Levchenko, and Nemenman, 2016; Endres and Wingreen, 2008). Classic work, beginning with the seminal paper by Berg and Purcell (1977), studied cell sensing in homogeneous environments (Berg, Howard C, and Purcell, 1977). This and subsequent works were extended to study the detection of spatially varying concentrations, where monotonic gradients remain the canonical environmental model (Hu et al., 2010; Mugler, Levchenko, and Nemenman, 2016; Endres and Wingreen, 2008). Recent work has started to address spatial complexity (Chou et al., 2011), but much work remains to understand how cell navigation strategies are affected by natural signal distributions, particularly spatially correlated fluctuations. Such complexity can pose challenges to cell engineering applications, such as CAR-T cell responses to tumor microenvironments (Martinez and Moon, 2019). Fundamentally, it is not clear what sense and response strategies are well adapted to operate in environments where signals take on complex spatial structures.

Modern signal processing theory shows that sensing strategies must adapt to the statistics of the input signals, suggesting that spatial sensing in cells should be adapted to the spatial structure of signaling molecules in the cells' native environments (Candès and Wakin, 2008). For example, when designing electronic sensor networks sensing spatial phenomena, adapting sensor placement to the spatial statistic of the signal can significantly improve information acquisition (Krause, Singh, and Guestrin, 2008). Furthermore, spatial navigation where sensing plays a key role may also benefit from sensor placement adaptation, as suggested by work from both robot and insect navigation (Iida and Nurzaman, 2016; Huston et al., 2015). For example, when navigating turbulent plumes, locusts actively move their antennae to odorant locations in order to acquire more information on source location (Huston et al., 2015). In the context of cell navigation, interstitial gradients can potentially

trap cells in local concentration peaks (Weber et al., 2013). Cells that can adapt sensing to patchy structure of the gradient may overcome local traps.

## 1.2 Experimental tools for profiling tissue signaling environments

A major challenging in studying how cells migration in complex environments, including tissue, is the difficulty associated with profiling the chemical cues present in vivo. Recent advancements in spatial omics technologies enables us to capture the spatial organization of cells and molecular signals in intact human tumors with unprecedented molecular detail, revealing the relationship between localization of different cell types and tens to thousands of molecular signals (Moffitt, Lundberg, and Heyn, 2022). Migration of cells such as T-cells and tumor cells is modulated by a rich array of signals within the tumor microenvironment (TME) such as chemokines, adhesion molecules, tumor antigens, immune checkpoints, and their cognate receptors (Lanitis et al., 2017). Recent advances in *in situ* molecular profiling techniques, including spatial transcriptomic (Rodriques et al., 2019; Eng et al., 2019) and proteomic (Giesen et al., 2014; Goltsev et al., 2018) methods, simultaneously capture the spatial relationship of tens to thousands of molecular signals and cell localization in intact human tumors with micron-scale resolution. Imaging mass cytometry (IMC) is one such technology that uses metal-labeled antibodies to enable simultaneous detection of up to 40 antigens and transcripts in intact tissue (Giesen et al., 2014).

## 1.3 Biomedical challenges associated with in vivo cell migration

The immune composition of the tumor microenvironment (TME) plays a crucial role in determining patient prognosis and response to cancer immunotherapies (Fridman et al., 2017; Binnewies, Mikhail, et al., 2018; Bruni, Angell, and Galon, 2020). Immunotherapies that alter the immune composition using transplanted or engineered immune cells (chimeric antigen receptor T-cell therapy) or remove immunosuppressive signaling (checkpoint inhibitors) have shown exciting results in relapsed and refractory tumors in hematological cancers and some solid tumors. However, effective therapeutic strategies for most solid tumors remain limited (Hegde and D. S. Chen, 2020; Choe, Williams, and Lim, 2020; Pitt et al., 2016). The TME is a complex mixture of immune cells, including T-cells, B cells, natural killer cells, and macrophages, as well as stromal cells and tumor cells (Fridman et al., 2017). The interactions between these cells can either promote or suppress tumor growth and progression, and ultimately impact patient outcomes. For example, high levels

of tumor-infiltrating lymphocytes (TILs) in the TME are associated with improved prognosis and response to immunotherapy across multiple cancer types (Haslam and Prasad, 2019; Lee and Ruppin, 2019). Conversely, an immunosuppressive TME characterized by low levels of TILs is associated with poor prognosis and reduced response to immunotherapy (Pittet, Michielin, and Migliorini, 2022). Durable, long-term clinical response of T-cell-based immunotherapies are often constrained by a lack of T-cell infiltration into the tumor, as seen in classically "cold" tumors such as triple-negative breast cancer or pancreatic cancer, which have seen little benefit from immunotherapy (Bonaventura, Paola, et al., 2019; Savas et al., 2016; Tsaur et al., 2021). The precise cellular and molecular factors that limit T-cell infiltration into tumors is an open question.

*Chapter 2*

# RECEPTOR LOCALIZATION MAXIMIZES INFORMATION ACQUISITION IN NATURAL CELL ENVIRONMENTS

## 2.1 Abstract

Cells in natural environments like tissue or soil sense and respond to extracellular ligands with intricately structured and non-monotonic spatial distributions that are sculpted by processes such as fluid flow and substrate adhesion. In this work, we show that spatial sensing and navigation can be optimized by adapting the spatial organization of signaling pathways to the spatial structure of the environment. We develop an information-theoretic framework for computing the optimal spatial organization of a sensing system for a given signaling environment. We find that receptor localization maximizes information acquisition in simulated natural contexts, including tissue and soil. Receptor localization extends naturally to produce a dynamic protocol for continuously redistributing signaling receptors, which when implemented using simple feedback, boosts cell navigation efficiency by 30-fold. Broadly, our framework readily adapts to studying how the spatial organization of signaling components other than receptors can be modulated to improve cellular information processing.

## 2.2 Introduction

Cells sense and respond in spatially structured environments, where signal distributions are determined by a range of chemical and physical processes from substrate binding to fluid flow (Fowell and Kim, 2021). In tissue and soil, distributions of extracellular ligands can be spatially discontinuous, consisting of local ligand patches (B.-G. Yang et al., 2007; M. Weber et al., 2013; Russo et al., 2016; Milde, Bergdorf, and Koumoutsakos, 2008; Kicheva et al., 2007; Sarris et al., 2012; Kennedy et al., 2006; Raynaud and Nunan, 2014; Hodge, 2006; De Anna et al., 2021; Nunan et al., 2001; Lim et al., 2015; Philipsborn et al., 2006). In tissue, diffusive signal-

ing molecules are transported by interstitial fluid through a porous medium. These molecules are then captured by cells and a non-uniform network of extracellular matrix (ECM) fibers, taking on a distribution that is stable and highly reticulated (B.-G. Yang et al., 2007; Kennedy et al., 2006; M. Weber et al., 2013; Russo et al., 2016; Kicheva et al., 2007; Sarris et al., 2012). For example, ECM-bound chemokine (CCL21) gradients extending from lymphatic vessels take on stable spatial structures, characterized by regions of high ligand concentration separated by spatial discontinuities (M. Weber et al., 2013). Similar observations have been made for the distribution of other chemokines, axon guidance cues and morphogens in tissues (Kicheva et al., 2007; Sarris et al., 2012; Kennedy et al., 2006; Lim et al., 2015). In soil, a heterogeneous pore network influences the spatial distribution of nutrients by dictating both the locations of nutrient sources as well as where nutrients likely accumulate (Raynaud and Nunan, 2014; Hodge, 2006; De Anna et al., 2021; Nunan et al., 2001). Free-living cells detect chemical cues released by patchy distributions of microorganisms, where molecules are moved via fluid flow and diffusion (Raynaud and Nunan, 2014; Hodge, 2006). Cells in these and other natural environments experience surface ligand profiles with varying concentration peaks, non-continuity and large dynamic range (Kennedy et al., 2006; Dlamini, Kennedy, and Juncker, 2020), differing strongly from smoothly varying, purely diffusive environments.

Modern signal processing theory shows that sensing strategies must adapt to the statistics of the input signals, suggesting that spatial sensing in cells should be adapted to the spatial structure of signaling molecules in the cells' native environments (Candès and Wakin, 2008). For example, when designing electronic sensor networks sensing spatial phenomena, adapting sensor placement to the spatial statistic of the signal can significantly improve information acquisition (Krause, Singh, and Guestrin, 2008). Furthermore, spatial navigation where sensing plays a key role may also benefit from sensor placement adaptation, as suggested by work from both robot and insect navigation (Iida and Nurzaman, 2016; Huston et al., 2015). For example, when navigating turbulent plumes, locusts actively move their antennae to odorant locations in order to acquire more information on source location (Huston et al., 2015). In the context of cell navigation, interstitial gradients can potentially trap cells in local concentration peaks (M. Weber et al., 2013). Cells that can adapt sensing to patchy structure of the gradient may overcome local traps.

Traditional approaches to studying cell sensing often use highly simplified environmental models, where signals are either uniform or monotonic, neglecting the

complex spatial structure in natural cell environments (Berg, Howard C, and Purcell, 1977; Hu et al., 2010; Mugler, Levchenko, and Nemenman, 2016; Endres and Wingreen, 2008). Classic work, beginning with the seminal paper by Berg and Purcell (1977), studied cell sensing in homogeneous environments (Berg, Howard C, and Purcell, 1977). This and subsequent works were extended to study the detection of spatially varying concentrations, where monotonic gradients remain the canonical environmental model (Hu et al., 2010; Mugler, Levchenko, and Nemenman, 2016; Endres and Wingreen, 2008). Recent work has started to address spatial complexity (Chou et al., 2011), but much work remains to understand how cell sensing strategies are affected by natural signal distributions, particularly spatially correlated fluctuations. Such complexity can pose challenges to cell engineering applications, such as CAR-T cell responses to tumor microenvironments (Martinez and E. K. Moon, 2019). Fundamentally, it is not clear what sense and response strategies are well adapted to operate in environments where signals take on complex spatial structures.

Interestingly, empirical observations suggest that cells might modulate the placement of their surface receptors to exploit the spatial structure of ligand distribution in its environment (Pignata et al., 2019; Bouzigues et al., 2007; Nieto et al., 1997; Buul et al., 2003; Shimonaka et al., 2003; Yokosuka et al., 2005; Mossman et al., 2005). For example, some axon guidance receptors can dynamically rearrange on the surface of growth cones (Pignata et al., 2019; Bouzigues et al., 2007). In such cases, receptors rearrange constantly, adjusting local surface densities in response to changes in ligand distribution across the cell surface. Some chemokine receptors in lymphocytes exhibit similar spatial dynamics (Nieto et al., 1997; Buul et al., 2003; Shimonaka et al., 2003). However, other chemokine receptors remain uniform even when their ligands are distributed non-uniformly (Vicente-Manzanares and Sánchez-Madrid, 2004). In addition, during antigen recognition, T-cell receptors (TCRs) take on different placements, ranging from uniform to highly polarized, depending on the density of antigen molecules on the surface of the opposing cell (Majzner et al., 2020). Thus, across a diverse range of cell surface receptors, we see different, even contradictory rearrangement behavior in response to changes in environmental structure. It remains unclear whether dynamic receptor rearrangement has an overarching biological function across disparate biological contexts.

Inspired by previous works that applied information maximization principles to understand the design of biological systems for signal processing (Tkačik, Walczak, and Bialek, 2009; Sokolowski and Tkačik, 2015; Monti, Lubensky, and Ten Wolde,

2018; Tkačik, Callan, and Bialek, 2008; Dubuis et al., 2013; Petkova et al., 2019; Tkačik and Gregor, 2021; Cheong et al., 2011), we formulated an information-theoretic framework showing spatial localization of cell surface receptors is an effective spatial sensing strategy in natural cell environments, but relatively inconsequential in purely diffusive environments. Our framework allows us to solve for receptor placements that maximizes information acquisition in natural environments, while generating such environments using existing computational models of tissue and soil microenvironments. We find that anisotropic receptor dynamics previously observed in cells are nearly optimal. Specifically, information acquisition is maximized when receptors form localized patch at regions of maximal ligand concentration. Optimizing receptor placement offers significant gain in information acquisition over uniformly distributed receptors, but only in native cell habitats, leading to an average of ∼ 1 bit of information gain in tissues and soils but only ∼ 0.01 bits in purely diffusive gradients. The optimal strategy maximizes information by taking advantage of patchy ligand distribution in natural environments, reallocating sensing resources to a small but high signal region on the cell surface, while explicitly "ignoring" ligand information at low signal regions.

Our framework extends naturally to produce a dynamic protocol for continuously relocalizing receptors in response to a dynamic environment. We show that a simple feedback scheme implements this protocol within a cell, and significantly improves cell navigation. Compared to cells with uniform receptor placement, cells using this scheme achieve more than 30-fold improvement in their ability to localize to the peak of simulated interstitial gradients. Furthermore, our model accurately predicts spatial distribution of membrane receptors observed experimentally (Pignata et al., 2019; Bouzigues et al., 2007; Nieto et al., 1997; Buul et al., 2003; Shimonaka et al., 2003). Importantly, our framework easily extends to study how spatial organization of many different cellular components, beyond receptor placement, affects information processing (see Discussion). Taken together, our model serves as a useful conceptual framework for understanding the role of spatial organization of signal transduction pathways in cell sensing, and provides a sensing strategy that is both effective in natural cell environments and amenable to cell engineering.

## 2.3 Results

### An optimal coding framework allows the computation of optimal receptor placement given spatial signal statistics



Figure 2.1: **Adapting receptor placements to signal (input) statistic of natural cell environments** (A) (Left) tuning sensor placement can boost the performance of electronic sensor network. (Right) cell surface receptors also function as a sensor network, taking as inputs ligand profiles $C$ across the cell surface and producing as outputs a profile of receptor activity $A$ across the cell membrane. The optimal receptor placement strategy $\phi^* : c \rightarrow r$ maps each ligand profile to a receptor placement, such that the mutual information $I(C; A)$ is maximized. (B) The problem of optimal receptor placement formulated as a resource allocation problem over parallel, noisy communication channels. The $i$-th channel represents the $i$-th region of the cell membrane, with input $C_i$, output $A_i$ and receptor number $r_i$. The input statistic $p(c)$ depends on the environment, and the measurement kernel $p(a_i|c_i, r_i)$ is modeled as a Poisson counting process. The general formulation of the optimal strategy $\phi^*$ allocates $N$ receptors to $m$ channels for each ligand profile $c$, such that $I(C, A)$ is maximized (Equation 2.2). The local formulation selected the receptor placement $\phi^*(c)$ that maximizes $I(\hat{c}, \hat{a})$, where $\hat{c}$ is a Poisson random vector with mean equal to $c$ (Equation 2.4).

Figure 2.1: **(continued)** (C) i. Approximating input statistic by simulating natural environments and sampling ligand profiles $\{c\}$ by tiling cells uniformly across the environment; ii. modeling ligand distribution in tissue microenvironment by incorporating diffusion, advection, ECM binding, degradation, and cell uptakes. iii. modeling ligand distribution in soil microenvironment by generating bacteria distributed in spatial patches, releasing diffusive ligands.

We are interested in optimal strategies for a task we refer to as spatial sensing. Spatial sensing is an inference task where a cell infers external profiles of varying ligand level across its surface from an internal profile of varying receptor activity across its membrane. This is a useful model task since optimizing performance on this task should improve the cell's ability to infer diverse environmental features.

We developed a theoretical framework to study whether manipulating the placement of cell surface receptors can improve spatial sensing performance. Optimizing spatial sensing by tuning receptor placement is analogous to optimizing distributed electronic sensor network by adjusting the location of sensors, which has been extensively studied in signal processing (Krause, Singh, and Guestrin, 2008). In the optimization of distributed sensor networks monitoring spatial phenomena (Figure 2.1A), it is well known that adjusting the placement of a limited number of sensors can significantly boost sensing performance, where the optimal placement strategy is dictated by the statistics of the input signals (Krause, Singh, and Guestrin, 2008; Caselton and Zidek, 1984). The collection of a limited number of receptors on the cell surface also functions as a distributed sensor network, sensing a spatial profile of varying ligand concentration (Figure 2.1A). Therefore, we hypothesized that receptor placement can be tuned to improve spatial sensing, and that the optimal strategy depends on the statistics of ligand profiles that cells typically encounter. Unlike traditional works in sensor optimization which focuses on finding a single "best" placement (Krause, Singh, and Guestrin, 2008), cells can rearrange their receptors within a matter of minutes (Bouzigues et al., 2007), leading to a potentially much richer class of strategies.

Before presenting the general optimization problem, we set up the mathematical framework through the lens of information theory. Consider a two-dimensional (2D) cell with a 1D membrane surface. By discretizing the membrane into $m$ equally sized regions, we modeled the membrane-receptor system as $m$ parallel communication channels (Figure 2.1B). The input to these $m$ channels is $C = (C_1, ..., C_m)$, where $C_i$ is a random variable representing the amount of ligands at the $i$-th membrane region.

The receptor profile $r = (r_1, .., r_m)$ denotes the amount of receptors allocated to each membrane region. The output $A = (A_1, ..., A_m)$ is the amount of active receptors across the membrane, which depends on $c$ and $r$ through $p(A = a|c, r)$, the measurement kernel. Consider a placement strategy $\phi : c \to r$, mapping a ligand profile to a receptor placement (Figure 2.1B). For a fixed number of receptors $N$, we are interested in the choice of $\phi$ that maximizes the mutual information $I(C; A)$ between the channels' inputs $C$ and outputs $A$, defined as,

$$I(C; A) = \sum_{c \in C} \sum_{a \in A} p(c, a) \log \frac{p(c, a)}{p(c)p(a)}. \tag{2.1}$$

The mutual information, in units of bits, quantifies the "amount of information" obtained about $C$ by observing $A$. It is minimized when $C$ and $A$ are independent, and maximized when one is a deterministic function of the other. Importantly, note the choice of m (membrane bins) sets an upper bound on the mutual information, hence sets the scale for all information value reported in this paper (see Supplement, Section 2.10 for derivations of this relation). Mathematically, the optimal strategy $\phi^*$ can be written as

$$\phi^*_{p(c)} = \underset{\substack{\forall c \ \phi(c) \geq 0 \\ \sum_i \phi_i(c) = N}}{\operatorname{argmax}} \ I(C; A \mid \phi, p(c)), \tag{2.2}$$

where $N$ is the total number of receptors which is taken to be a constant. Note the mutual information converges towards its upper bound as $N$ increases (Figure 2.8A). The mutual information is agnostic to the decoding process in that it does not assume any details about downstream signaling, nor the exact environmental features a cell may try to decode, expanding the scope of our results.

To solve for $\phi^*$, we needed to specify both a measurement kernel $p(a|c, r)$ and an input statistic $p(c)$. We modeled $p(a|c, r)$ assuming that each receptor binds ligands locally and activates independently. Furthermore, each local sensing process is modeled as a Poisson counting process. These assumptions yield the following measurement kernel,

$$p(A = a \mid C = c, \ r) = \prod_{i=1}^{m} \frac{\mu_i{}^{a_i}}{a_i!} e^{-\mu_i}, \tag{2.3}$$

where $\mu_i = r_i(\frac{c_i}{c_i + K_d} + \alpha \frac{K_d}{c_i + K_d})$ is the average number of active receptors at the $i$-th membrane region. $K_d$ is the equilibrium dissociation constant and $\alpha$ accounts for constitutive activity of receptors observed in cells, including many GPCRs, which

we take to be small ($\alpha \ll 1$) (Seifert and Wenzel-Seifert, 2002; Slack and Hall, 2012). The bracket term represents the probability of receptor activation, and the fractional term $\frac{K_d}{c_i+K_d}$ ensures it is always less than 1 (Buchwald, 2019).

Next, we specify the input statistic $p(c)$ for three classes of environments: soil, tissue, and monotonic gradient. For each class of environment, we constructed $p(c)$ empirically, by computationally generating a ligand concentration field as the steady-state solution of a partial-differential equation (PDE), and sampling ligand profiles ($\{c\}$) from them by evaluating the PDE solution around cells placed at different spatial locations (Figure 2.1C-i) (for details see Supplement, Section 2.10). Putting the empirical measure on the samples $\{c\}$ approximates the true distribution of $C$. For soil, we follow mathematical models from (Melke et al., 2010) and (Raynaud and Nunan, 2014), modeling diffusive ligands released from a group of soil bacteria whose spatial distribution agrees with the statistical properties of real soil colonies (Figure 2.1C-iii, Figure 2.2A). For tissue, we adopted models from (Milde, Bergdorf, and Koumoutsakos, 2008) and (Rejniak et al., 2013), where they modeled diffusive ligands released from a localized source, perturbed by in vivo processes such as interstitial fluid flow and heterogeneous ECM binding, leading to an immobilized interstitial gradient (Figure 2.1C-ii, Figure 2.2B). We also considered a monotonic gradient (Figure 2.2B) as an exponential fit to the simulated interstitial gradient. Fitting ensures any difference between the two environments are due to differences in local structures, not global features such as gradient decay length or average concentration. It is important to note that the overall framework can accommodate any choice of $p(c)$ and $p(a|c)$ beyond what we have considered.

We are interested in the functional relationship between ligand profiles $\{c\}$ and their optimal receptor placements $\{\phi^*(c)\}$. To this end, we computed the optimal receptor placement for each sampled profile $c$ individually, reducing the general formulation to a local formulation. Given ligand profile $c$, random vector $\hat{c}$ represents local fluctuations of $c$ due to stochasticity of reaction-diffusion events. In the case of unimolecular reaction-diffusion processes, it can be shown that $\hat{c}$ is a Poisson vector with mean equal to $c$, solution of the PDE. Therefore, we can solve for $\phi^*(c)$ locally by maximizing the mutual information between $\hat{c}$ and the resulting output $\hat{a}$,

$$\phi^*(c) = \underset{\substack{r \geq 0 \\ \sum_i r_i = N}}{\operatorname{argmax}} \; I(\hat{c}, \hat{a} \mid r), \qquad (2.4)$$

where $p(\hat{a}) = \sum_c p(\hat{a}|\hat{c} = c)p(\hat{c} = c)$. The main difference between the general formulation of (2.2) and local formulation of (2.4) is their dependence on the input

statistic $p(c)$. In the general formulation, the strategy $\phi^*_{p(c)}$ is explicitly parametrized by $p(c)$. In the local formulation, $\phi^*$ is independent of the choice of $p(c)$. However, differences in $p(c)$ between environments will still crucially affect the set of optimal receptor profiles that cells will actually adopt. This is because changing $p(c)$ changes the region of the domain of $\phi^*$ that is most relevant, thus changing the optimal receptor profiles that are actually used in different environments. For example, suppose environment A and B have input statistic $p_A$ and $p_B$, and any ligand profile observed in A is not observed in B, and vice versa. Although $\phi^*$ is the same between A and B, this function is being evaluated on entirely different ligand profiles in A compared to B, so that receptor profiles observed in the two environment will likely be very different, in ways dictated by differences between their input statistic $p_A$ and $p_B$. As a result, the statistical structure over the space of ligand profiles plays an important role in determining which receptor placement is effective, even when the placements are computed locally for each ligand profile.

# Receptor localization yields optimal spatial sensing in natural environments



Figure 2.2: **Receptor localization optimizes information acquisition in natural environments.** (A), computationally generated ligand concentration fields using PDE models of soil (left), tissue (interstitium) (middle), and simple exponential gradient (right, fitted to tissue with correlation index $R^2 = 0.98$). (B), i) Example of optimal receptor profile $\phi^*(c)$ (colored) and the corresponding ligand profile $c$ (gray); ii) entropy for each optimal receptor placements in $\{\phi^*(c)\}$ colored by environment, colored triangles indicate the entropy of three receptor placements shown in i). (C), optimal efficacy $\eta$ colored by environments.

Figure 2.2: **(continued)** (D), i) efficacy for soils of varying values of $\sigma^2_{\text{bacteria}}$, ii) efficacy for tissue with varying values of $k_{\text{ECM}}$, and for exponential gradients fitted to each tissue (gradient). Stars correspond to parameter values used to generate Panel A–C and E. (E), scatterplot where each dot corresponds to a single pair of $c$ and $\phi^*(c)$, where $c$ is sampled from environments as illustrated in Figure 2.1C-i; $\eta_c$ is defined in Equation 2.8. Across all panels, $N = 1000$, $K_d = 40nM$, $\alpha = 0.1$, $m = 100$.

Optimal strategies of receptor placement are similar for soil and tissue environment, where receptors are highly localized within membrane positions experiencing high ligand concentrations. Figure 2.2B-i shows three examples of optimal receptor placements $\phi^*(c)$ (colored) with the corresponding ligand profile $c$, one from each class of environments shown in Figure 2.2A. In all three cases, the peak of each optimal receptor profile is oriented towards the position of highest ligand concentration. Compared to monotonic gradient, receptor profiles optimized for the ligand profiles sampled from tissue and soil are highly localized, with around 80% of receptors found within 10% of the membrane. In general, the optimal strategy consistently allocates more receptors to regions of higher ligand concentration, but in a highly nonlinear manner. Figure 2.2B-iii shows, across all sampled ligand profiles $\{c\}$, the peak of receptor profiles always align with the peak of ligand profiles. But instead of allocating receptors proportional to ligand level, receptors tend to be highly localized to a few membrane positions with the highest ligand concentrations.

Indeed, Figure 2.2B-ii shows that optimal receptor profiles tend to have low entropy. The entropy of receptor profile $r$, defined as $H(r) = -\sum_i \frac{r_i}{N} \log\left(\frac{r_i}{N}\right)$, can be used as a measure of localization. Note the maximal value of this entropy measure is limited by the number of membrane bins m. Low entropy corresponds to receptor profiles where most receptors are concentrated to a few membrane positions, forming localized patches. Such high degree of localization is partly explained by low receptor numbers. When receptors are limited, information gain per receptor within each membrane channel is approximately independent of receptor number (for details see Supplement, Section 2.10). Thus, the optimal solution allocates all receptors to the channel with the highest information content (see Figure 2.9). In addition, receptors are more localized for sensing in soil and tissue because locally, they exhibit greater spatial variations in ligand concentration compared to simple gradients (Figure 2.2A) (for details see Supplement, Section 2.10). Absolute ligand concentration also influences the optimal strategy, which we take to be dilute in agreement with empirical measurements (Xiangdan Wang et al., 2008; Clark et al.,

2015). In saturating environments, the optimal solution completely switches, allocating most receptors to regions of lowest ligand concentrations (for details see Supplement, Section 2.10, Figure 2.10). In summary, the optimal placement strategy $\phi^*$ in the environments studied can be approximated by a simple scheme, where receptors localize to form patches at positions of high ligand concentration.

Optimally placed receptors significantly improve information acquisition relative to uniform receptors, especially in soil and tissue environments. To make this statement precise, we quantified the efficacy of a receptor placement strategy $\phi : c \to r$ with respect to a set of ligand profiles $\{c\}$. First, we denote by $I_\phi$ the average information cells acquire by adapting a the placement strategy $\phi$,

$$I_\phi = \langle I(\hat{c}; \hat{a} \mid \phi) \rangle_c \tag{2.5}$$

where $\langle \cdot \rangle_c$ denotes averaging across the set of sample ligand profiles $\{c\}$, and recall $\hat{c}$ is a Poisson-distributed random vector with mean $c$. The efficacy of $\phi$ is the absolute increase in average information cells acquire by adapting the strategy $\phi$ compared to a uniform strategy $\phi^u$, where receptors are always uniformly distributed,

$$\eta(\phi) = I_\phi - I_{\phi^u}, \tag{2.6}$$

We are specifically interested in the optimal efficacy $\eta(\phi^*)$, and refer to it as $\eta$ when the dependency is clear from context. Since we will often compare the optimal and uniform strategy, we will denote $I_{\phi^*}$ as $I_{\text{opt}}$ and $I_{\phi^u}$ as $I_{\text{unif}}$. For a particular $\eta(\phi^*)$, the set of ligand profiles $\{c\}$ referred to in its definition is always the same set that $\phi^*$ is optimized for. The larger $\eta$ is, the more beneficial it is for cells to place receptors optimally rather than uniformly. We found that $\eta$ is an order of magnitude larger for soil and tissue environment compared to a simple gradient (Figure 2.2C). This difference persists across cells of different size and across a wide range of receptor parameter values (Figure 2.2C, Figure 2.7B, Figure 2.11). In other words, placing receptors optimally rather than uniformly benefits cells in complex, natural environments significantly more than cells in simple, monotonic gradients. Note that differences between tissue and monotonic gradient are due to differences in local spatial structure, not global features such as gradient decay length or global average concentration, as both parameters were made to be identical between the two environments. Lastly, although the mutual information is an exponential measure, so an improvement by one bit has different meaning depending on the baseline $I_{\text{unif}}$, this fact does not hinder interpretation of $\eta$ as $I_{\text{unif}}$ is similar between the three environments (Figure 2.11A).

In addition to the large difference in information gain ($\eta$) between natural environments and simple gradients, Figure 2.11 shows similar differences exist when comparing other metrics assessing the benefit of optimizing receptor placement, such as the relative information gain (($I_{\text{opt}} - I_{\text{unif}})/I_{\text{unif}}$) and the absolute increase in the number of distinguishable input signals ($2^{I_{\text{opt}}} - 2^{I_{\text{unif}}}$). For example, optimizing receptor placement increases the number of distinguishable input states by 40 in tissue, while optimizing the same receptors in the fitted gradient leads to an increase of 1 (Figure 2.11A). Note that in the limit of strong constitutive receptor activity, all placement strategies become equivalent to uniformly distributed receptors. Since receptor activation in the absence of ligands reduces statistical dependence between ligand level and receptor activity, the average information acquisition $I_\phi$ for any strategy $\phi$ converges to zero, driving information gain $\eta(\phi)$ compared to the uniform strategy to zero (Figure 2.11A).

For both soil and tissue environment, the optimal efficacy $\eta$ depends on a key parameter in their respective PDE model. We illustrate this dependence by adjusting the value of each respective parameter, sampling new ligand profiles $\{c\}$, solving for optimal placements $\{\phi^*(c)\}$, and computing $\eta$. Figure 2.2D shows how $\eta$ changes as we adjust environmental parameters. In soil, $\eta$ drops substantially as ligand sources (bacteria) become more aggregated (Figure 2.2D-i), corresponding to an increase in the parameter $\sigma^2_{\text{bacteria}}$ of the random process used to model bacterial distribution (star corresponds to empirical value from (Raynaud and Nunan, 2014)). This result is intuitive since increasing the extent of aggregation of sources makes the environment appear more like a simple gradient generated from a single ligand source. In tissue, optimal efficacy dropped when most ligands were found in solution, instead of bound to the ECM (Figure 2.2D-ii), corresponding to low ECM binding rate ($k_{\text{ECM}}$). For reference, star indicates the empirical value of $k_{\text{ECM}}$ for the chemokine CXCL13 (B.-G. Yang et al., 2007). Compared to its fitted monotonic gradient, $\eta$ in the interstitial gradient remain significantly higher for all ECM binding rates (Figure 2.2D-ii). In tissue, gradients made up of ECM-bound ligands are ubiquitous, suggesting the optimization of receptor placement is highly relevant.

Optimal efficacy ($\eta$) is larger in soil and tissue because ligand profiles that cells encounter in such environments tend to be more patchy, having most of the ligands concentrated in a small subset of membrane regions. We make this statement precise by quantifying patchiness of a ligand profile $c$ using a measure of sparsity,

$$\text{sparsity}(c) = 1 - \frac{\overline{c}}{c_{\text{RMS}}}, \tag{2.7}$$

where the root-mean-square $c_{\text{RMS}} = \sqrt{\frac{1}{m} \sum_i c_i^2}$ and $\overline{c} = \frac{1}{m} \sum_i c_i$ is the average concentration of $c$ across the membrane. A ligand profile with a sparsity of one has all ligands contained in a single membrane region, whereas a uniform distribution of ligands has a sparsity of zero. Next, we defined an efficacy measure $\eta_c$ for each ligand profile $c$,

$$\eta_c = I(\hat{c}; \hat{a} \mid \phi^*) - I(\hat{c}; \hat{a} \mid \phi^u) = I_{\text{opt},c} - I_{\text{unif},c}, \tag{2.8}$$

where again $\phi^u$ denotes uniform receptor distribution. Unlike $\eta$ as defined in Equation 2.6, $\eta_c$ does not involve the averaging across the entire set $\{c\}$ through $\langle \cdot \rangle_c$, it measures improvement in information gain for only a single ligand profile $c$. The larger $\eta_c$ is, the more useful the optimal placement is for sensing $c$ compared to a uniform profile. Each dot in Figure 2.2E corresponds to a ligand profile sampled from an environment, as illustrated in Figure 2.1C-i. Figure 2.2E shows that 1) across a wide range of concentrations, sparser ligand profiles tend to induce higher efficacy $\eta_c$, and 2) ligand profiles sampled from soil and tissue tend to be sparser compared to profiles from the corresponding monotonic gradient. Taken together, since signals cells encounter in natural environments tend to have sparse concentration profiles, cells can improve their spatial sensing performance by localizing receptors to regions of high ligand concentration.

In summary, the value of optimizing receptor placement as a sensing strategy depends strongly on the environmental structure. Patchy ligand distribution found in tissues and soils makes optimizing receptor placement a highly effective sensing strategy. Our result demonstrates that uncovering effective cell sensing strategies requires a careful consideration of the spatial structure of the cells' natural habitat.

**Spatial sensing via the optimal strategy is robust to imprecise placements caused by biological constraints**



Figure 2.3: Optimal efficacy $\eta(\phi^*)$ is robust to minor deviations in receptor placement away from the optimal form. (A), the effect of different degrees of shifting and flattening applied to a receptor profile (black curve). (B), colors of heat map represent ratio of perturbed efficacy $\eta(\phi^p)$ to optimal efficacy $\eta(\phi^*)$ for different combinations of shifting and flattening, computed for ligand profiles $\{c\}$ sampled from either soil or tissue; call-out boxes corresponds to different sets of perturbations, showing the average of the optimal $\{\phi^*(c)\}$ (gray) and perturbed $\{\phi^p(c)\}$ (red) receptor placements, after all ligand profile peaks were centered; red number indicates the value on heat map; cell radius = 10 µm.

Despite the optimal strategy $\phi^*$ being highly localized and precisely oriented, we found that neither features are necessary to achieve high efficacy. Given the stochastic nature of biochemical processes in cells, this robustness is crucial as it makes the strategy feasible in cells. Fortunately, receptors do not need to adopt $\phi^*$ precisely in order to obtain substantial information gain. To illustrate, we perturb the optimal placements and show that sensing efficacy persists when receptors partially align with ligand peak and localize weakly. For soil and tissue, we circularly shift and flatten (by applying a moving average) all optimal receptor profiles $\{\phi^*(c)\}$ computed from sample ligand profiles to obtain $\{\phi^p(c)\}$, the corresponding set of perturbed profiles. Different degrees of shifting and flattening represents different degrees of misalignment and weakened localization, respectively. Figure 2.3A shows results of different perturbations (colored) applied to a receptor profile (black). To assess the effect of these perturbations on sensing, we compute the efficacy $\eta(\phi^p)$ of the perturbed profiles, and compare it to the optimal efficacy $\eta(\phi^*)$. The heatmap in

Figure 2.3B shows the ratio of perturbed to optimal efficacy for various combinations of perturbations, across soil and tissue. Figure 2.3B-i shows two examples of perturbations (red dots) that drastically alter the receptor profile while still achieving high efficacy. The red and gray curve in the call-out box represents what the "average" perturbed and optimal profiles look like, respectively. They are obtained by circularly shifting each profile in $\{\phi^p(c)\}$ and $\{\phi^*(c)\}$ so the peak of $c$ is center, followed by averaging across the set of shifted profiles element-wise. Clearly, highly localized receptors (> 80% of receptors found within 10% of membrane) are not necessary for effective sensing. In fact, compared to uniformly distributed receptors, a modest enrichment of receptors oriented towards the ligand peak (4 folds relative to uniform) already provides significant information gain (Figure 2.3B) — a behavior of membrane receptors that has been observed in cells (McClure et al., 2015). Such robustness holds across different cell sizes and efficacy metric (Figure 2.12). In tissue, the heatmap of Figure 2.3B-i also shows that weakly localized receptors (large flatten factor) are more robust to misalignment (large shift factor). Altogether, this robustness (Figure 2.3) suggest that biochemical implementations of receptor localization could improve sensing in natural or engineered cells even in the presence of stochastic fluctuations that induce imperfect localization. Moreover, receptor localization previously observed in cells is sufficient to obtain significant information gain.

**Optimization framework extends naturally to produce a dynamic protocol for sensing time-varying ligand profiles**



Figure 2.4: **A dynamic receptor placement protocol based on maximizing rate of information gain.** (A), schematic showing a cell moving along a path (gray curve) sensing a sequence of ligand profiles $\{c_t\}$ at points (crosses) along the path, using receptor placements $\{r_t^*\}$ generated by the dynamic protocol. (B), accounting for transport cost, the optimal placement strategy is modified to localize receptors to an intermediate position between subsequent ligand peaks or form multiple receptor peaks.

Our framework extends naturally to produce a dynamic protocol for rearranging receptors in response to dynamically changing ligand profiles. So far, we have viewed ligand profiles as static snapshots and considered instantaneous protocols for receptor placement. In reality, cells sense while actively exploring their environment, so that the ligand profile it experiences is changing in time, both due to intrinsic changes in the environment state as well as due to the motion of the cell. As the ligand profile $c_t$ changes over time, we want the receptor profile $r_t$ to change in an "efficient" manner to improve information acquisition (Figure 2.4A). Specifically, we obtain a dynamic protocol by extending our framework to account for both information acquisition and a "cost" for changing receptor location. We quantify this cost using the Wasserstein-1 distance $W_1(r_A, r_B)$, which is the minimum distance receptors must move across the cell surface to redistribute from profile $r_A$ to $r_B$ (for details see Supplement, Section 2.10). For a cell sensing a sequence of ligand profiles $\{c_t\}_{t=1}^T$ over time, the optimal receptor placement $r_t^*$ for $c_t$ now depends additionally on $r_{t-1}^*$, the optimal placement for the previous ligand profile,

$$r_t^* = \underset{\substack{r \geq 0 \\ \sum_i r_i = N}}{\mathrm{argmax}}\ I\left(\hat{c}_t; \hat{a} \mid r\right) - \gamma W_1\left(r_{t-1}^*, r\right), \tag{2.9}$$

where $p(\hat{a}) = \sum_c p(\hat{a}|\hat{c}_t)p(\hat{c}_t)$, and $\gamma \geq 0$ represents the cost of moving one receptor per unit distance. The cost $\gamma$ implicitly encodes a time scale for receptor redistribution. Smaller $\gamma$ means less "cost" is associated with redistributing receptors, hence the receptor profile becomes more dynamic. The exact relationship between $\gamma$ and the speed of receptor redistribution depends on both receptor properties and the environment, see Figure 2.13A–D for an example of how receptor speed scales with $\gamma$. For $\gamma = 0$, this formulation reduces to the original formulation of Equation 2.4. This dynamic formulation admits a natural interpretation as maximizing information rate (information per receptor-distance moved) instead of absolute information. For $t = 1$, we define $r_t^*$ according to the original formulation. Hence, we refer to the dynamic protocol of Equation 2.9 as the general optimal strategy since it encompasses $\phi^*$. Figure 2.4B illustrates two salient features of this dynamic protocol. Firstly (left), when the peak of the previous receptor profile $r_{t-1}^*$ is near the peak of the current ligand profile $c_t$, $r_t^*$ is obtained by shifting receptors towards the current ligand peak but not aligning fully. Secondly (right), when the peak of the previous receptor profile is far from the current ligand peak, some receptors are moved to form an additional patch at the current ligand peak refer to Figure 2.13F to see how changing $\gamma$ affects the receptor behavior in Figure 2.4B). Receptor properties such as the strength of constitutive receptor activity ($\alpha$) also affect receptor redistribution dynamics. Decreasing $\alpha$ increases mutual information, making the cost of redistribution less relevant, leading receptors to localize more readily to align with new ligand peaks (Figure 2.13E). Although the formulation of Equation 2.9 is quite complex, this general optimal strategy can be achieved by a simple receptor feedback scheme.

**Simple feedback scheme rearranges receptors to achieve near-optimal information acquisition**



Figure 2.5: **Positive feedback scheme redistributes receptors to achieve near-optimal sensing efficacy for both static and dynamic signals.** (A), the cell is modeled as a one-dimensional membrane lattice with a well-mixed cytosol. Receptors are subject to three redistribution mechanisms: endocytosis ($k_{off}$), activity-dependent incorporation into membrane ($hA_iR_{cyto}$), membrane diffusion ($d_m$); Value of $h$ sets the feedback strength between receptor activity and the rate with which receptors incorporate into the membrane; $h = 4 \times 10^{-3}$ s$^{-1}$, $d_m = 1 \times 10^{-2}$ µm$^2$ s$^{-1}$, $k_{off} = 1 \times 10^{-1}$ s$^{-1}$ (see Supplement, Table 2.2). (B), receptor profiles (yellow) generated by simulating the feedback scheme for an initially uniform set of receptors, against a static ligand profile from tissue and soil. (C), ratio of scheme efficacy $\eta(\phi^s)$ to optimal efficacy $\eta(\phi^*)$ for static signals $\{c\}$ sampled from soil and tissue, stars indicate parameter values used for simulation in Panel B. (D), (top) kymograph showing the entire temporal sequence of receptor profiles of a moving cell; (bottom) position of ligand peak aligned in time with position of receptor peak as generated by the feedback scheme.

Figure 2.5: **(continued)** (E), snapshots of receptor profiles taken at select time points. (F), ratio of scheme efficacy $\eta(\phi^s)$ to optimal efficacy $\eta(\phi^*)$ for a sequence of signals $\{c_t\}$ sampled by translating a cell through soil and tissue environment, stars indicate parameter values used for simulation in Panel D-E; cell radius = $10\mu m$ (see Figure 2.15B-C for results with cell radius = $5\mu m$). (G), histogram showing the distribution of ligand peak (gray) and receptor peak (yellow) position on the membrane of the cell from Panel D, dashed black line indicates the direction of the global gradient with respect to membrane positions. See Table 2.2 for feedback scheme simulation parameters.

A positive feedback scheme implements the general optimal strategy (Equation 2.9), organizing receptors into localized pole(s) to achieve near-optimal information acquisition. Asymmetric protein localization is a fundamental building block of many complex spatial behavior in cells, involved in sensing, movement, growth, and division (Macara and Mili, 2008). Many natural localization circuits are well-characterized down to molecular details (Hegemann et al., 2015; Zhu et al., 2020). In fact, even synthetic networks have been experimentally constructed in yeast, capable of reliably organizing membrane-bound proteins into one or more localized poles (Chau et al., 2012). Such works demonstrate the feasibility of engineering new spatial organization systems in cells.

Using a PDE model of a receptor redistribution scheme, we show that simple, local interactions can redistribute receptors to achieve near-optimal information acquisition, for both static and dynamic signals. Figure 2.5A illustrates the three redistribution processes (arrows) in our feedback scheme that affects receptor distribution ($r$), which can be expressed mathematically as

$$\frac{\partial r(x,t)}{\partial t} = D\nabla^2_{\text{memb}}r - k_{\text{off}}\,r + hAR_{\text{cyto}}, \tag{2.10}$$

where $x$ denotes membrane position and $t$ denotes time. The first term represents lateral diffusion of receptors on the membrane with uniform diffusivity $D$. The second term represents endocytosis of receptors with rate $k_{\text{off}}$. The last term represents recruitment/incorporation of receptors to membrane position $i$ from a homogeneous cytoplasmic pool ($R_{\text{cyto}}$) with rates $hA_i$, where $h$ a proportionality constant and $A_i$ is the local receptor activity (see Supplement, Section 2.10 for simulation detail, including how parameter values were derived from literature). This activity-dependent receptor recruitment provides the necessary feedback that enables ligand-dependent receptor redistribution. Recent works suggest activity-dependent receptor recruitment can be achieved through biased docking and fusion

of secretory vesicles carrying the receptors to regions of high receptor activity (Xin Wang et al., 2019; Hegemann et al., 2015; Kinoshita-Kawada et al., 2019). Budding yeasts Ste2 receptors achieve this feedback using an interacting loop with intracellular polarity factor Cdc42 (Hegemann et al., 2015). Note that our feedback scheme is only meant to illustrate one possible implementation of the dynamic rearrangement protocol. Feasible alternatives such as activity-dependent endocytosis or microtubule-dependent receptor redistribution have also been proposed, providing a range of biochemical strategies for implementation (Bouzigues et al., 2007; Suchkov et al., 2010).

Given a fixed ligand profile $c$, Figure 2.5B shows our feedback scheme can, within minutes, localize receptors (yellow) towards the position of maximum ligand concentration. This localization dynamic is robust to changes in $k_{\text{off}}$ and $h$ across at least an order of magnitude (Figure 2.15A). We denote the steady-state receptor profile generated by our scheme in response to ligand profile $c$ as $\phi^s(c)$. As Figure 2.5B shows, scheme-generated profiles are far less localized than their optimal counterpart $\phi^*(c)$. Despite this, Figure 2.5C shows scheme efficacy $\eta(\phi^s)$ are close to that of the optimal value $\eta(\phi^*)$. Recall $\eta(\phi^*)$ measures the absolute increase in average information acquired using optimally placed instead of uniform receptors. Therefore, the scheme efficacy $\eta(\phi^s)$ makes a similar comparison between scheme-driven and uniform receptors. In Figure 2.5C, we see scheme efficacy is robust to variations in both endocytosis ($k_{\text{off}}$) and average membrane incorporation rate ($\langle hA_i \rangle_i$), with other parameters fixed to empirical values (Marco et al., 2007). Stars represent parameters used to simulate profiles in Figure 2.5B.

Our feedback scheme (Equation 2.10) can continuously rearrange receptors in response to changes in ligand profile, exhibiting dynamics similar to the optimal dynamic protocol (Equation 2.9). Figure 2.5D-E shows a time-varying receptor profile, generated by the feedback scheme in a cell translating across the tissue environment. In this dynamic setting, the scheme can still induce asymmetric redistribution of receptors. Figure 2.5D (top) shows this dynamic asymmetry through a kymograph of a sequence of receptor profiles $\{\phi^s(c_t)\}$. As desired, snapshots along this sequence show receptors localize towards regions of high ligand concentration (Figure 2.5E). Receptor placements generated by our scheme exhibit features of the dynamic protocol shown in Figure 2.4B. First, as the ligand peak changes position slightly, the receptor peak gets shifted in the same direction after a delay. Figure 2.5D (bottom) illustrates this phenomena by aligning the time trace of both

peak positions . Here, a shift in the ligand peak (gray) is often followed by a corresponding shift in receptor peak (yellow) after an appreciable delay, hence there is only partial peak-to-peak alignment. Second, if the ligand peak changes position abruptly, a second receptor peak forms, oriented towards with the new ligand peak. Figure 2.5E-iii illustrates this clearly by showing a new receptor peak forming precisely after a large shift in ligand peak position (Figure 2.5D). We assess the performance of our scheme by comparing scheme-generated placements $\{\phi^s(c)\}$ and optimal placements $\{\phi^*(c)\}$ corresponding to the same sequence of ligand profiles $\{c_t\}$. Figure 2.5F shows that for cells moving in soil and tissue, scheme efficacy $\eta(\phi^s)$ (star) is not far from the optimal value $\eta(\phi^*)$. Furthermore, scheme efficacy is robust to variations in endocytosis ($k_{\text{off}}$) and average incorporation rate ($\langle ha_i \rangle_i$). Taken together, our feedback scheme organizes receptors to achieve near-optimal information acquisition, in both static and dynamic environments.

Our feedback scheme can align receptors with the global gradient direction, suggesting that this scheme may allow cells to escape local ligand concentration peaks within interstitial gradients. On the one hand, Figure 2.5G shows that the peak of ligand profiles (gray), as experienced by cells, do not always agree with the direction of the global gradient (dashed line) — a known feature of interstitial gradients (M. Weber et al., 2013). On the other hand, receptors organized by the feedback scheme (yellow) align very well with the global gradient direction. This effect of the feedback scheme comes from its ability to localize receptors and account for past receptor profiles. The latter allows the current receptor profile to carry memory of past ligand profiles that the cell has encountered, enabling a form of spatial averaging over ligand peaks. This alignment of receptors to the global gradient should provide significant boost to cell navigation performance, especially in non-monotonic, interstitial gradients.

**Feedback scheme enables cells to search quickly and localize precisely in simulated interstitial gradients**



Figure 2.6: **In simulated interstitial gradient, cells localize to source quickly and precisely when receptors are redistributed by the feedback scheme instead of uniformly distributed.** (A), (left) interstitial CCL21 gradient, (right) white curves represent haptotactic trajectories of dendritic cells (M. Weber et al., 2013). (B), (top) schematic of a navigation task where a cell (green flag) in a region of an interstitial gradient move towards the source (red flag) by sensing spatially distributed ligands by decoding source direction locally; the ligand field shown is a region of the tissue environment in Figure 2.2A obtained through PDE simulation (see Supplement, Section 2.10); (bottom) red curve shows the tissue ligand field averaged over the y-direction, and black curve is the fitted exponential gradient, scale bar: 10 µm. (C), sample trajectories of repeated simulations of cells navigating with uniform receptors (blue) and with scheme-driven receptors (orange), all scale bars: 10 µm.

Figure 2.6: **(continued)** (D), (left) histogram of time taken to reach source across 600 cells at different starting positions of equal distance from source. note the rightmost bar includes all cells that did not reach the source after 8 hours; (right) bar plot showing percentage of runs completed in 1 hrs (success rate), see also Figure 2.16 for success rate across different simulation parameters. (E), same type of data as in Panel D for cells navigating in an exponential gradient (fitted to the interstitial gradient used to generate Panel D). (F), red stripes (left) represent growth cones moving within specific lamina along a Slit gradient (right schematic), an ellipse-shaped cell used for this simulation to mimic navigating growth cone, scale bar: 40 µm (Xiao and Baier, 2007). (G), (top) schematic of a navigation task where a cell (green flag) senses its environment in order to remain close to source. solid white line represents cell trajectory, dotted white line demarcates a distance of 5 µm from ligand source. (see Table 2.3 for tissue simulation parameters), (bottom) red curve shows the tissue ligand field averaged over the y-direction, and black curve is the fitted exponential gradient, scale bar: 2 µm. (H), sample trajectories of repeated simulations of cells performing task with either uniform or scheme-driven receptors, all scale bars: 2 µm. (I), (left) histogram of time spent by cell at various distance from the ligand source. (measured from source to farthest point on cell, perpendicular to source edge) aggregated across 600 cells starting at different positions, moving at 2 hours near the ligand source; (right) bar plot shows percentage of time spent more than 5 µm from source (error rate), see also Figure 2.16 for error rate across different simulation parameters. (J), same data type as Panel I for cells navigating in an exponential gradient (fit to interstitial gradient of Panel I).

Cells using our feedback scheme effectively localizes to the ligand source of simulated interstitial gradients, while cells with uniform receptors become trapped away from the source by local concentration peaks. Immune cells can navigate towards the source of an interstitial gradient in a directed, efficient manner (Figure 2.6A) (M. Weber et al., 2013). Efficient navigation can be difficult in complex tissue environments, partly due to the existence of local maxima away from the ligand source, potentially trapping cells on their way to the source (Figure 2.6B). By simulating cell navigation using standard models of directional decoding (for details see Supplement, Section 2.10), we found that cells with uniform receptors can indeed become trapped during navigation. Figure 2.6C demonstrates this behavior through the trajectories of individual cells with uniform receptors (blue), as they consistently become stuck within specific locations of the environment. On the other hand, using the same method of directional decoding, cells with scheme-driven receptors (orange) reliably reach the source in an efficient manner. Figure 2.6D illustrates this difference through a histogram of the time it took for a cell to reach the source, created by simulating cells starting at uniformly sampled locations 40 µm from the

source, moving at a constant speed of $2\mu m/$min. Remarkably, for the circuit pa-
rameter values chosen, only 2% of cells (13/600) with uniform receptors reached
the source within 1 hour, compared to 73% of cells (436/600) using the feedback
scheme, boosting success rate by more than 30-folds. In fact, Figure 2.6D shows
that > 97% of cells with uniform receptors fail to reach the source even after 6
hours, as expected due to being trapped. This improvement in success rate persists
across a wide range of scheme parameters (orders-of-magnitude) and directional
decoding schemes (Figure 2.16). We emphasize that the poor performance of cells
with uniform receptors is only partially due to inaccuracy associated with decoding
local gradients. Indeed, cells that only follow local gradients have trouble finding
the global peak (ligand source) in simulated interstitial gradients. We demonstrate
this by simulating cells moving precisely along local gradient directions (direction
of maximal increase in ligand concentration across the cell's surface), such cells
become trapped at local ligand peaks on their way to the source (Figure 2.16C). As
expected, Figure 2.6E shows that the difference in performance between uniform and
scheme-driven receptors is relatively less pronounced in the simple gradient (black
curve Figure 2.6B bottom) — a twofold difference in success rate. We discuss
the analogy between our feedback scheme and the infotaxis algorithm (Vergassola,
Villermaux, and Shraiman, 2007) in the Discussion section .

Our feedback scheme can also help cells remain within a highly precise region along
a chemical gradient. During certain developmental programs, cells must restrict
their movements within a region along a gradient in order to form stable anatomical
structures. Growth cones demonstrate an extraordinary ability in accomplishing
this task. Axon projections of retinal ganglion cells can remain within a band of
tissue (lamina) of only $3 - 7\mu m$ wide, at a specific point along a chemical gradient
(Figure 2.6F) (Xiao and Baier, 2007; Xiao, Staub, et al., 2011). Figure 2.6G
illustrates how we assess our scheme's ability to achieve this level of precision.
We initiate a cell at a gradient source and track the proportion of time the cell
was more than $5\,\mu m$ away from the source. As the cell moves along the gradient,
uneven ligand distribution in the environment can lead the cell to move erroneously
away from the source. Figure 2.6H shows that cells with uniform receptors (blue)
can indeed make excursions away from the source. But cells with the feedback
scheme (orange) reliably stay close to the source for an extended period of time.
We quantify this difference by pooling from 600 trajectories of cells starting at
different positions along the source, decoding source direction and navigating for
2 hours. Figure 2.6I shows the number of time steps the cells collectively spent at

specific distances from the source. For the circuit parameter values chosen, cells with uniform receptors are found more than 5 µm away from the source 15% of the time (22204/144000 steps). On the other hand, cells with the feedback scheme do so only 2% of the time (3287/144000 steps), a 7-fold reduction in error rate. This difference in error rate persists for a wide range of scheme parameters and directional decoding schemes (Figure 2.16D,F). Similar improvement in performance is found for cells navigating in fitted exponential gradients (black curve Figure 2.6G bottom). Figure 2.6J shows the error rate is reduced by 10-fold from cells with uniform to scheme-driven receptors (10% vs. 1%). This result is intuitive as the gradients used for this task has extremely short decay length (5 µm) to mimic in vivo gradients that growth cones encounter. As a result, the fitted exponential becomes very similar to the simulated interstitial gradient. Taken together, our feedback scheme is functionally effective in simulated patchy gradients found in tissue, enabling cells to solve common navigation tasks with significantly improved accuracy and precision.

**Optimal efficacy accurately predicts experimental observations of membrane-receptor distribution**



Figure 2.7: **Optimal efficacy $\eta$ predicts observed distributions of cell surface receptors using their surface expression level and binding affinity.** (A) observed membrane distributions of receptors in heterogeneous environments, i. white arrowheads indicate Slit receptor Robo1 of commissural growth cones navigating in an interstitial Slit gradient (Pignata et al., 2019), ii. chemokine receptor CCR5 of human T lymphocytes subject to a CCL5 gradient (Nieto et al., 1997), iii. (left) transmission image of growth cone, white arrowhead indicates direction of GABA gradient, (right) bright dots represent GABA$_A$R redistributing in response to a GABA gradient (Bouzigues et al., 2007), iv. C5aR-GFP remains uniformly distributed in response to a point source of a C5aR agonist, delivered by micropipette (white dot), open arrowheads point to leading edges of cells (Servant et al., 1999). Scale bars i-iii: $5\mu m$, iv: $10\mu m$. (B) optimal efficacy $\eta$ for different values of $K_d$ and $N$; values computed using the tissue environment, where the ratio between average ligand concentration and $K_d$ is fixed, $\alpha = 0.1$; red dots correspond to receptors that polarize in heterogeneous environments (CCR2, CXCR4, CCR5, GABA$_A$R, Robo1), white dots represent receptors that are constantly uniform (IL-2R, TNFR1, TGF$\beta$R2, CR3, C5aR), roman numerals correspond to receptors in Panel A, see Table 2.4 for receptor data.

In addition to generating optimal sensing strategies for simulated environments, our framework can be used to predict receptor distribution of natural cell surface receptors (Figure 2.7A), using both the environmental structure in which the receptors function and their biological properties. In addition to environmental structure, receptor properties such as cell surface expression level ($N$) and binding affinity ($K_d$) also play a role in determining the optimal strategy by affecting the measurement

kernel (Equation 2.3). For a simulated tissue environment, Figure 7B shows that despite offering greater than twofold gain in information ($\eta > 100$) when $N$ is small, optimizing receptor placement offers nearly zero gain in information ($\eta \ll 1$) when $N/K_d$ is large. High $N$ and low $K_d$ improve information acquisition by allowing the receptor activities to be more sensitive to changes in input level, and since the total amount of information available to the cell is fixed, the amount of additional gain that can be made by optimizing receptor placement is reduced.

Figure 2.7B suggests that for real cell surface receptors, we may be able to predict their membrane distribution by specifying both their environment and biological parameters ($N$, $K_d$). Specifically for receptors functioning in tissue, we predict those with parameters that fall within the high $\eta$ regime (Figure 2.7B) are more likely to adapt the optimal localized distribution. Although data are limited, empirical observations of real receptors agree with this prediction. Comparing data across cell surface receptors from multiple cell types found in human tissue, Figure 2.7B show that receptors (red dots) with parameters corresponding to large $\eta$ have been observed to localize in non-uniform environments (Figure 2.7A-i-iii). Importantly, the localized receptors concentrate at the region of the membrane with the highest ligand concentration, consistent with the theoretically optimal strategy. Such localization is clearly illustrated in Figure 2.7A-iii, where GABA receptors localize precisely to the membrane region experiencing the highest ligand concentration, as indicated by the white arrow. Receptors (white dots) with parameters corresponding to small $\eta$, however, are always uniformly distributed (Figure 2.7A-iv), even when the environment is non-uniform. Furthermore, although Figure 2.7B is based on a fixed $\alpha$ (constitutive receptor activity), the striking relationship between receptor organization and optimal efficacy $\eta$ holds for values of $\alpha$ spanning at least two orders-of-magnitude (Figure 2.17). More detailed comparisons between the experimental receptor distributions and the theoretical optimum is unfortunately not possible, because detailed, quantitative descriptions of the ligand profile that the cells were sensing are not available. This agreement between theory and observations is not meant to imply that evolution optimizes receptor placement. Indeed, there are key caveats such as variations in receptor expression over time and differences between the environments of different receptors. Our theory does, however, provide a framework for studying natural variations in the spatial organization of receptors, such as differences observed between chemotactic receptors in the same T-cell (Nieto et al., 1997).

## 2.4 Discussion

A rich collection of works, spanning diverse areas including developmental biology, systems biology, and neuroscience, put forth the idea of optimizing mutual information to predict the design of biological systems that process information (Tkačik, Walczak, and Bialek, 2009; Sokolowski and Tkačik, 2015; Monti, Lubensky, and Ten Wolde, 2018; Tkačik, Callan, and Bialek, 2008; Dubuis et al., 2013; Petkova et al., 2019; Tkačik and Gregor, 2021; Cheong et al., 2011). For example, information maximization principles have been applied to derive fundamental limits on the fidelity of information transfer in biochemical networks (Mehta et al., 2009; Cheong et al., 2011). Inspired by these works, we formulated an information-theoretic framework that enables us to compute effective cell sensing strategies across different environments. We applied the framework to different signaling microenvironments, including tissues and soils, to discover a receptor localization strategy that significantly improves both cell sensing and navigation. More broadly, our work has a series of conceptual and practical implications. Our theory suggests a functional role for spatial organization in cellular information processing, conceptually showing how spatially organized intracellular components can be used by cells to more accurately infer the state of its external environment, here through sensing and chemoreceptors. Furthermore, our theory conceptually shows how spatial organization of a cell's sensing apparatus can actually reflect spatial structure of its environment. Similar results are found in neuroscience, but it is interesting to see how such an efficient coding perspective can help understand spatial organization within a cell. Lastly, our theory has practical consequences for cell engineering. Currently, most synthetic circuits function without spatial modulation and are studied in well-mixed compartments. Our work shows how spatial control over synthetic sense and response architectures can provide new strategies for engineering circuits that function in natural environments.

### Connection between information acquisition and navigation

We showed that a receptor placement strategy aimed at maximizing information rate can boost cell navigation performance. Since information content increases towards the ligand source, receptors are more likely to move towards the side of the membrane closer to the source rather than away, enforcing movement up gradients. Furthermore, the trade-off between information acquisition and receptor redistribution in Equation 2.9 can be viewed as combining exploitative and exploratory tendencies, where larger redistribution "cost" favors exploitation. This strategy is similar in

principle to the infotaxis algorithm (Vergassola, Villermaux, and Shraiman, 2007), where one can view receptors as "navigating agents", whose movements guide the cell towards the target. Although the idea is quite intuitive, the exact relationship between navigation and information acquisition requires further investigation. On the one hand, the feedback scheme is most effective in the case of limited sampling of inputs (Figure 2.16B,E), which suggests maximizing information content indeed helps with navigation. On the other hand, moving receptors to maximize information rate is significantly more effective as a navigation strategy compared to only maximizing absolute information (Figure 2.14).

## 2.5 Data availability

All analysis, simulation and plotting scripts are openly available at: `https://github.com/neonine2/receptor-code`. All data generated in this work is openly available at: `http://dx.doi.org/10.22002/D1.2149`.

## 2.6 Material availability

This paper did not generate new reagents.

## 2.7 Acknowledgements

We thank Michael Elowitz and Erik Winfree for scientific discussions and Dominik Schildknecht, Han Kim, Guruprasad Raghavan, Pranav Bhamidipati, Abdullah Farooq, for feedback on the manuscript, Inna-Marie Strazhnik for illustrations, and Angela Anderson for editorial advice. We also would like to thank Eugenio Marco and Katarzyna Rejniak for technical advice with receptor feedback and tissue simulations, respectively. The authors would like to acknowledge the Heritage Medical Research Institute and Packard Foundation for funding and intellectual support.

## 2.8 Author contributions

Conceptualization, Z.W. and M.T.; Methodology, Z.W. and M.T.; Manuscript writing, Z.W. and M.T.; Supervision, M.T.; Funding acquisition, M.T.

## 2.9 Declaration of interests

The authors declare no competing interests.

## 2.10    Supplemental methods

**Formulation of optimization problem**

In this paper, we developed a theoretical framework to study whether manipulating the placement of cell surface receptors can improve the spatial sensing performance. Optimizing spatial sensing by tuning receptor placement is analogous to optimizing distributed electronic sensor network by adjusting the location of sensors, which has been extensively studied in signal processing (Krause, Singh, and Guestrin, 2008). Before presenting the general optimization problem, we set up the mathematical framework through the lens of information theory. Consider a two-dimensional (2D) cell with a 1D membrane surface. By discretizing the membrane into $m$ equally sized regions, we modeled the membrane-receptor system as $m$ parallel communication channels (Figure 2.1B). The $i$-th channel takes as input $C_i \in \mathbb{N}_0$, a random variable denoting ligand count at the $i$-th region of the membrane surface. Given $r_i \in \mathbb{N}_0$ receptors, this channel produces as output $A_i \in \mathbb{N}_0$, a number of active receptors that is random due to stochastic nature of receptor activation and randomness in $C_i$. Given $m$ channels representing the entire cell membrane, our model comprised four key mathematical objects: ligand profile $C = (C_1, ..., C_m)$, receptor placement $r = (r_1, .., r_m)$, active receptor profile $A = (A_1, ..., A_m)$, and measurement kernel $P(A = a \mid C = c, r)$. The input $C \sim p(c)$ is now the entire ligand profile across the cell surface. Each realization $c$ of $C$ has probability $p(c)$ of being observed. We explain below how $p(c)$ can be constructed to represent statistics of ligand profiles cells naturally encounter (see **Input statistic**). The receptor profile $r$ denotes the number of receptor allocated to each membrane region. The output $A \sim p(a)$ is the number of active receptors across the membrane, which depends on $c$ and $r$ through $p(a|c, r)$, the measurement kernel. We explain below how this kernel can be modeled (see **Measurement kernel**).

Consider a placement strategy $\phi : c \rightarrow r$, that maps a ligand profile to a receptor placement. In our general optimization problem (Figure 2.1B), we are interested in the choice of $\phi$ that maximizes the amount of information the cell can obtain regarding $C$ by observing $A$, for a fixed number of receptors $N$. Formally, we quantify this information using the mutual information,

$$I(C; A) = \sum_{c \in C} \sum_{a \in A} p(c, a) \log \frac{p(c, a)}{p(c)p(a)}. \qquad (2.11)$$

The mutual information is minimized when $C$ and $A$ are independent, and maximized when one is a deterministic function of the other. Since $p(c, a) = p(a|c, r =$

$\phi(c))p(c)$, each summand in the mutual information will be affected by the choice of $\phi$. Taken together, we arrive at our general formulation of the optimal strategy $\phi^*$:

$$\phi^*_{p(c)} = \underset{\substack{\forall c\ \phi(c)\geq 0 \\ \sum_i \phi_i(c)=N}}{\text{argmax}}\ I(C; A \mid \phi, p(c)), \tag{2.12}$$

where $N$ is the total number of receptors. The subscript $p(c)$ is meant to emphasize the dependence of the optimal strategy on the input statistics.

To solve for $\phi^*_{p(c)}$, we needed to specify both a measurement kernel $p(a|c, r)$ and an input statistic $p(c)$. The input statistic $p(c)$ for an environment represents the probability that a cell in that environment will encounter the ligand profile $c$.

**Measurement kernel**

We model $p(a|c, r)$ assuming that each receptor binds ligands locally and activates independently of other receptors. These assumptions allow us to factorize $p(a|c, r)$ as follows,

$$P(A = a \mid C = c,\ r) = \prod_{i=1}^{m} P(A_i = a_i \mid C_i = c_i,\ r_i). \tag{2.13}$$

Each local sensing process involves probabilistic ligand-receptor interaction which can be viewed as a Bernoulli process. In this way, the number of active receptors follows a Binomial distribution, which can be approximated with the Poisson distribution when the probability of successful binding event is low. Indeed, experimental measurements have shown that receptor occupancy is well approximated by the Poisson distribution (Ueda et al., 2001), such that

$$P(A_i = a_i \mid C_i = c_i,\ r_i) = \frac{\mu_i^{a_i}}{a_i!}e^{-\mu_i}, \tag{2.14}$$

where $\mu_i = r_i\left(\frac{c_i}{c_i+K_d} + \alpha\frac{K_d}{c_i+K_d}\right)$. The bracket term represents the probability of activation for a receptor experiencing $c_i$ ligands. $K_d$ is the equilibrium dissociation constant and $\alpha$ represents constitutive receptor activity, which we take to be small ($\alpha \ll 1$). In other words, the number of active receptors $A_i$ given ligand count $c_i$ is a Poisson random variable with mean $\mu_i$. Equation (2.13) and (2.14) together specify the measurement kernel.

**Input statistic**

Next, we specify the input statistic $p(c)$ which will be determined by spatial distribution of ligands, thus differ between different classes of environment. Suppose a

circular cell samples its environment by binding nearby ligands. The cell will encounter certain spatial profiles of ligands more often than others, and such statistics will likely depend on the type of environment the cell lives in. In this work, we studied three classes of environments: soil, tissue, and monotonic gradient. Closed form models do not exist for ligand profile statistics of natural environments. Therefore, we take an empirical approach, generating instances of each environment as the steady-state solution of PDE models and directly sample ligand profiles from them (see Section 2.10 for details on all PDE models). For soil, we adopted mathematical models from (Melke et al., 2010) and (Raynaud and Nunan, 2014), modeling diffusive ligands released from a group of soil bacteria whose spatial distribution agrees with the statistical properties of real soil colonies (Figure 2.1C-iii, Figure 2.2A). For tissue, we adopted models from (Milde, Bergdorf, and Koumoutsakos, 2008) and (Rejniak et al., 2013), where they modeled diffusive ligands released from a localized source, perturbed by in vivo processes such as interstitial fluid flow, non-uniform ECM binding and cell uptake, to represent an interstitial gradient (Figure 2.1C-ii, Figure 2.2B). We also considered a simple (monotonic) gradient (Figure 2.2C) which is an exponential fit to the simulated interstitial gradient (Figure 2.2B). Fitting ensures any difference between the two environments are due to differences in local structures, not global features such as gradient decay length or average concentration. For each environment, we obtain a ligand concentration field $c(x)$ as the steady-state solution of a PDE. Then, we tile it with a cell of fixed size and evaluate the concentration field along each cell membrane to obtain a set of ligand profiles denoted $\{c\}$ (Figure 2.1C-i). Putting the empirical measure on the samples $\{c\}$ approximates the true distribution of $C$. It is important to note that although we modeled $p(c)$ and $p(a|c)$ in these ways, the overall framework can accommodate any alternative choices of model.

For these choices of $p(c)$ and $p(a|c)$, we aimed to study the functional relationship between ligand profiles $\{c\}$ and their optimal receptor placements $\phi^*(c)$. To this end, we optimized receptor profiles for each sampled profile $c$ individually, reducing the general problem to a local formulation. Given ligand profile $c$, the random vector $\hat{c}$ represents local fluctuations of $c$ due to stochasticity of reaction-diffusion events. In the case of unimolecular reaction-diffusion processes, it can be shown that $\hat{c}$ is a Poisson random vector with mean equal to $c$, solution of the PDE. Therefore, we can solve for $\phi^*(c)$ locally by maximizing the mutual information between $\hat{c}$ and

the resulting output $\hat{a}$,

$$\phi^*(c) = \underset{\substack{r \geq 0 \\ \sum_i r_i = N}}{\text{argmax}} \ I(\hat{c}, \hat{a} \mid r), \tag{2.15}$$

where $p(\hat{a}) = \sum_c p(\hat{a}|\hat{c} = c)p(\hat{c} = c)$ and $N$ is the total receptor number. We assume $r$ to be real-valued instead of integer-valued when solving (2.15), this is reasonable as long as $N$ is not too small.

The main difference between the general formulation of (2.12) and local formulation of (2.15) is their dependence on the input statistic $p(c)$. In the general formulation, the strategy $\phi^*_{p(c)}$ is explicitly parametrized by $p(c)$. In the local formulation, $\phi^*$ is independent of the choice of $p(c)$. However, differences in $p(c)$ between environments will still crucially affect the set of optimal receptor profiles that cells will actually adopt. This is because changing $p(c)$ changes the region of the domain of $\phi^*$ that is most relevant, thus changing the optimal receptors profiles that are actually used in different environments. For example, suppose environment A and B have input statistic $p_A$ and $p_B$ with non-overlapping support, meaning that any ligand profile observed in A is not observed in B, and vice versa. Although $\phi^*$ is the same between A and B, this function is being evaluated on entirely different ligand profiles in A compared to B, so that receptor profiles observed in the two environment will likely be very different, in ways dictated by differences between their input statistic $p_A$ and $p_B$. As a result, the statistical structure over the space of ligand profiles plays an important role in determining which receptor placement is effective, even when the placements are computed locally for each ligand profile.

The constrained nonlinear optimization problem of (2.15) was evaluated using the fmincon routine of MATLAB 2021 ("MATLAB Optimization Toolbox" 2021a). The Sequential Quadratic Programming algorithm was used to ensure accurate solutions that may exist near the boundary of the feasible region. Furthermore, the analytical gradient of the objective function, shown in equation Equation 2.28, was supplied to ensure faster convergence.

**Bin number and mutual information**

An important point to emphasize is that the choice of m (number of discrete membrane bins) sets a scale for all information values reported in the paper, because the mutual information ($I(C; A)$) is bounded by the entropy of its input which scales logarithmically with $m$. To illustrate this relationship, we consider a cell sensing a uniform distribution of ligands, where it experiences $\bar{c}$ molecules on average. In

this simplified setting, the number of ligand molecules at each of the $m$ membrane bin are represented by Poisson random variable with mean $\bar{c}/m$. Using the fact that the mutual information is bounded above by its input entropy, we get the following bound,

$$I(C; A) \leq mH(C_i),\qquad(2.16)$$

where $C_i \sim \text{Pois}(\bar{c}/m)$. Substituting $H(C_i)$ for the entropy of a Poisson random variable, we can rewrite the bound above as,

$$I(C; A) \leq m\left(\frac{\bar{c}}{m}[1 - \log(\bar{c}/m)] + e^{-\bar{c}/m}\sum_{k=0}^{\infty}\frac{(\bar{c}/m)^k \log(k!)}{k!}\right).\qquad(2.17)$$

For large m, this expression simplifies to yield an upper bound on the mutual information that scales logarithmically with $m$,

$$I(C; A) \leq \bar{c}(1 - \log(\bar{c})) + \bar{c}\log(m).\qquad(2.18)$$

Figure 2.8A shows this upper bound for $\bar{c} = 1$. As further validation of Equation 2.18, Figure 2.8A shows that as we increase the number of receptors, $I(C; A)$ converges toward the derived upper bound (red). Furthermore, the result that optimizing receptor placement is significantly more beneficial in natural environments compared to simple gradients holds for a wide range of membrane bin numbers, as shown in Figure 2.8B where the absolute information gain ($\eta$, Equation 2.6) is significantly larger in natural environments compared to simple gradients, for a wide range of m values.



Figure 2.8: **Effect of the number of membrane bins (m) on mutual information and optimal efficacy, Related to Figure 1 and 2** (A) The maximum mutual information achievable (red line, Equation 2.18) as a function of m, the number of membrane bins. As the number of receptors per bin (N/m) increases, the mutual information converges to its maximum value. (B) the optimal efficacy for different choices of m, computed across tissue, soil, and simple gradient, $\alpha = 0.01, K_d = 40, N = 1000$.

**Theoretical properties of Poisson channels**

In information theory, the Poisson channel is a canonical model used to study communication of information by random discrete occurrences in time that obey Poisson statistics. We show that we can map our receptor activation model directly onto this canonical model. As a result, we make use of existing results from information theory regarding the Poisson channel to 1) show that the localized receptor placement strategy described in the main text holds across most reasonable biochemical models of receptor activation and 2) provide intuition for how different factors such as ligand concentration can alter the optimal strategy.

**Mapping receptor model to the canonical Poisson channel model** We begin by showing how a single membrane-receptor channel can be mapped to the canonical scalar Poisson model. The same argument applies for mapping multiple parallel membrane-receptor channels to the canonical vector Poisson model introduced in the next section.

Recall the receptor model of Equation 2.14 we used to represent the number of active receptors $A$ for a given ligand level $c$, which is motivated by empirical measurements of receptor activity,

$$p(A = a \mid C = c, r) = \frac{\mu^a}{a!} e^{-\mu}, \quad \mu = r\left(\frac{c}{c + K_d} + \alpha \frac{K_d}{c + K_d}\right). \qquad (2.19)$$

Although this model of receptor activation consists of many biochemical details, we can map it directly onto the canonical scalar Poisson model,

$$Y|X \sim \text{Pois}(\alpha X) \qquad (2.20)$$

where $X$ is a scalar input, $Y$ is a scalar output, and $\alpha$ is a scaling variable. Such a model defines a Poisson channel whose output is a Poisson random variable conditioned on the input $X$ with its mean equal to $rX$, where $X$ is an arbitrary input random variable. We map this channel model maps onto our model of receptor activation for a single membrane region, by defining $X, \alpha$ in the following way,

$$\begin{aligned} X := f(C) &= \frac{C}{C + K_d} + \alpha \frac{K_d}{C + K_d}, \\ \alpha &:= r, \end{aligned} \qquad (2.21)$$

where $X$ represents the probability of receptor activation, $r$ denotes the number of receptors. From this set of definitions, it follows that $Y = A$ is the number of

active receptors. Note that $\alpha$ from equation (2.20) is a constant value, rather than a function like the placement strategy $\phi$. Therefore, Equation 2.21 agrees with our local formulation of (2.15), and matches the general formulation of (2.12) if $\phi$ is a constant function. An important consequence of this mapping is that we can now study the quantity $I(X, Y)$ since,

$$I(X; Y) = I(f(C); A) \tag{2.22}$$
$$= I(C; A), \tag{2.23}$$

where the second line follows from the fact that the mutual information is invariant to invertible transformations $f$. Since most physical models of receptor activity ($f$) are strictly increasing functions of ligand count, hence invertible, theoretical properties of $I(X, Y)$ which we discuss here directly applies to many receptor models beyond what is considered in this work, such as models with signal amplification and receptor cooperativity. Note that we are leaving the probability distribution $P(X)$ unspecified, which again makes many of the following results valid for many choices of $f$.

**Relating properties of $I(X, Y)$ to receptor sensing** Having established the relationship $I(X; Y) = I(C; A)$, we now use the scalar Poisson model to illustrate how theoretical properties of the mutual information $I(X, Y)$ agrees with our intuition of ligand sensing via receptor binding. We specialize to the case where $X$ is a non-negative random variable which is sufficient for our problem as $X$ only takes values between 0 and 1. In this setting, Theorem 2 of (Guo, Shamai, and Verdú, 2008) gives the partial information gain for the scalar Poisson channel as,

$$\frac{d}{dr}I(X; Y) = E[X \log X - E[X|Y] \log E[X|Y]]. \tag{2.24}$$

An immediate consequence of Equation 2.24 is that the mutual information $I(X, Y)$ (hence $I(C; A)$) is strictly increasing in the scaling variable (receptor number), which follows from the fact that the right side of Equation 2.24 is non-negative due to Jensen's inequality since $x \log x$ is a convex function. The fact that $I(X, Y)$ is strictly increasing in $r$ agrees with the intuition that increasing the number of receptors should increase the amount of information the cell can acquire about its external environment.

Observe that the right hand side of Equation 2.24 is exactly the minimum mean loss in estimating $X$ based on $Y$ under the loss function $l(x_1, x_2) = x_1 \log(x_1/x_2) - x_1 + x_2$.

Using this fact, one can show that $I(X, Y)$ is a concave function of $r$, which again agrees with the intuition that since the total amount of information available $H(X)$ is fixed, incremental gain in information acquisition must diminish as more receptors are added. Importantly, the fact that $I(X; Y)$ is an increasing, concave function of the scaling variable $r$ holds across all models of receptor activation (Equation 2.21). In particularly, the concavity of $I(X; Y)$ is a general phenomena and not a result of saturation from ligand binding.

**Mapping full membrane-receptor model to vector Poisson channel model** We will now rewrite our local optimization problem of Equation 2.15 using the canonical vector Poisson channel model. By doing so, we will be able to use theoretical properties of the vector Poisson model to provide additional insight into the optimal solution, and expand the result beyond the specific receptor model used in the main text. By a similar argument as in the scalar Poisson case, the full membrane-receptor model considered in our work (main text, Equation 2.3) maps exactly onto the canonical vector Poisson channel model, defined as

$$\boldsymbol{Y}|\boldsymbol{X} \sim \prod_{i=1}^{m} P(Y_i \mid \boldsymbol{X}) = \prod_{i=1}^{m} \text{Pois}(Y_i \mid (\boldsymbol{\Phi X})_i) \tag{2.25}$$

where the random vector $\boldsymbol{X} = (X_1, X_2, ..., X_m)$ maps to the probability of receptor activation across the $m$ discretized membrane regions, the random vector $\boldsymbol{Y} = (Y_1, Y_2, ..., Y_m)$ maps to the random vector of active receptors $\boldsymbol{A} = (A_1, A_2, ..., A_m)$, the channel matrix $\boldsymbol{\Phi} \in \mathbb{R}_+^{m \times m}$ can map onto receptor placement $\boldsymbol{r} = (r_1, r_2, ..., r_m)$ such that $\boldsymbol{\Phi} = \text{diag}(\boldsymbol{r})$. This mapping represents the fact that receptors bind ligands locally and activate independently of other receptors. We introduce $\boldsymbol{\Phi}$ for completeness, showing that this model can accommodate situations where there are crosstalks between channels, leading to non-zero terms in the off-diagonal. Since the equivalent of (2.23) holds for the vector model, we can rewrite the optimization problem in Equation 2.15 as,

$$\boldsymbol{r}^* = \underset{\substack{\boldsymbol{r} \geq 0 \\ \sum_i r_i = N}}{\text{argmax}} \ I(\boldsymbol{X}, \boldsymbol{Y} \mid \boldsymbol{r}). \tag{2.26}$$

By working with Equation 2.26, we derive results that hold for many models of receptor activation, including all models where activity is a monotonic function of ligand level.

**Reformulation of receptor optimization in terms of partial information gain**
We reformulate Equation 2.26 in terms of the partial derivatives $\partial I(\boldsymbol{X}; \boldsymbol{Y})/\partial r_i$,

which provides additional insight into the optimal solution. According to the Karush-Kuhn-Tucker (KKT) conditions, the following must hold at the optimal solution $r^*$ for $1 \leq i \leq m$ ,

$$\frac{d}{dr_i} I(X; Y \mid r^*) = \lambda - \mu_i,$$

$$\mu_i r_i^* = 0,$$

(2.27)

where $\mu_i \geq 0$ and $\lambda$ are the KKT multipliers. Another way to interpret the equations above is that for all channels where the optimal receptor number $r_i^*$ is non-zero, their partial derivatives $\frac{d}{dr_i} I(X; Y) \mid_{r^*}$ must be equal. Put another way, optimal solution occurs when incremental information gain is matched across channels. Since whenever the partial derivatives do not all agree, then one can always move receptors from the channel with a smaller partial derivative to one with higher partial derivative to achieve a higher mutual information.

**Asymptotic of the gradient of mutual information at small** $N$ We show that when total receptor number is low, the partial information gain depends only on the properties of $X$, the probability of receptor activation. This result allows us to directly solve for the optimal solution at the low $N$ regime. Theorem 1 of (L. Wang et al., 2014) gives the gradient of mutual information between input and output of the vector Poisson channel $I(X; Y)$, with respect to the matrix $\Phi$ as,

$$\nabla_\Phi I(X; Y)_{ij} = E[X_j \log(\Phi X)_i] - E[E[X_j|Y] \log E[(\Phi X)_i|Y]]. \quad (2.28)$$

Specializing to the setting where $\Phi$ is a diagonal matrix with $\text{diag}(\Phi) = r$, the derivative of mutual information with respect to receptor number at position $i$ is,

$$\frac{d}{dr_i} I(X; Y) = E[X_i \log X_i] - E[E[X_i|Y] \log E[X_i|Y]]. \quad (2.29)$$

This derivative can be interpreted as the information gain at the $i$-th channel per receptor added.

When total receptor number ($N$) is low, corresponding to all scaling variables ($\{r_i\}$) being small, we can express (2.29) as a function of just the random variable $X$. First, using Lemma 1 from (Dytso, Fauß, and Poor, 2020), we have

$$r_i E[X_i|Y = y] = (y_i + 1) \frac{p_Y(y + 1_i)}{p_Y(y)}$$

$$= r_i \frac{E\left[(X_i)^{y_i+1} e^{-r_i X_i} \prod_{m \neq i} \frac{1}{y_m!} (X_m)^{y_m} e^{-r_m X_m}\right]}{E\left[\prod_m (X_m)^{y_m} e^{-r_m X_m}\right]}.$$

(2.30)

Therefore, for $r_i > 0$, we have

$$E[X_i|\mathbf{Y} = \mathbf{y}] = \frac{E\left[(X_i)^{y_i+1} e^{-r_i X_i} \prod_{m \neq i} \frac{1}{y_m!} (X_m)^{y_m} e^{-r_m X_m}\right]}{E\left[\prod_m (X_m)^{y_m} e^{-r_m X_m}\right]}. \tag{2.31}$$

Now using monotone convergence theorem, we obtain

$$\lim_{r_1,\ldots,r_m \to 0^+} E[X_i \mid \mathbf{Y} = \mathbf{y}] = \frac{E\left[(X_i)^{y_i+1} \prod_{m \neq i} \frac{1}{y_m!} (X_m)^{y_m} e^{-r_m X_m}\right]}{E\left[\prod_m (X_m)^{y_m}\right]} \tag{2.32}$$

The above limit holds for any path, and also holds for all value of $\mathbf{y}$ including zero. Evaluating the above limit at $\mathbf{y} = 0$ we have

$$\lim_{r_1,\ldots,r_m \to 0^+} E[X_i|\mathbf{Y} = 0] = E[X_i]. \tag{2.33}$$

Applying this limit to (2.29) gives the desired result,

$$\lim_{r_1,\ldots,r_m \to 0^+} \frac{\partial}{\partial r_i} I(\mathbf{X};\mathbf{Y}) = E[X_i \log X_i] - E[X_i] \log E[X_i], \tag{2.34}$$

which we denote as

$$\frac{\partial I_0}{\partial r_i} := \lim_{r_1,\ldots,r_m \to 0^+} \frac{\partial}{\partial r_i} I(\mathbf{X};\mathbf{Y}). \tag{2.35}$$

In this limit, the partial derivatives are independent of receptor number. Intuitively, when receptor numbers are low, the effect of diminishing return that comes from having many receptors should be weak. Importantly, this result holds for arbitrary distribution $P(\mathbf{X})$, hence it holds for any environmental statistic and model of receptor activation. As we show in the next section, this fact allows us to solve for the optimal solution $\mathbf{r}^*$ exactly in the low $N$ limit.

**Factors affecting optimal receptor placement**

**Total receptor number**

Optimal receptor placement can be strongly localized when receptors are limited in quantity. When receptor number is small, Equation 2.34 shows that the $\frac{d}{dr_i} I(\mathbf{X};\mathbf{Y})$ becomes independent of receptor number. This result implies that $I(\mathbf{X};\mathbf{Y})$ is maximized when all receptors are allocated to the channel with the largest partial derivative $\frac{d}{dr_i} I(\mathbf{X};\mathbf{Y})$, resulting in strong receptor localization. We can see this by first noting that in the limit of small $N$, Equation 2.34 and Taylor's theorem allows us to write the mutual information as a linear function of $\mathbf{r}$,

$$I(\mathbf{X};\mathbf{Y}) = \sum_{i=1}^{m} \frac{\partial I_0}{\partial r_i} r_i \tag{2.36}$$

Hence, our optimization problem becomes a linear program with the following form,

$$\begin{aligned} \text{maximize} \quad & \boldsymbol{a}^T \boldsymbol{r} \\ \text{subject to} \quad & \mathbf{1}^T \boldsymbol{r} = r_{\text{tot}}, \quad \boldsymbol{r} \geq 0, \end{aligned} \tag{2.37}$$

where $a_i = \partial I_0 / \partial r_i$. Suppose the $a_i$'s are sorted in increasing order (with corresponding $r_i$'s rearranged as well),

$$a_1 \leq a_2 \leq \cdots < a_k = \cdots = a_{m-1} = a_m \tag{2.38}$$

Denote $a_{\max} := a_m$, we have

$$\boldsymbol{a}^T \boldsymbol{r} \leq a_{\max}(\mathbf{1}^T \boldsymbol{r}) = a_{\max} r_{\text{tot}} \tag{2.39}$$

for all feasible $\boldsymbol{r}$, with equality if and only if

$$r_k + \cdots + r_m = r_{\text{tot}}. \tag{2.40}$$

The optimal solution is then to allocate all receptors among the channels with maximal partial information gain $\partial I_0 / \partial r_i$, in any manner. This result is quite intuitive. If the channel information gains are fixed, we allocate all receptors to the channel with the highest gain. Since Equation 2.34 is valid for all non-negative random variable $X$, this result holds for arbitrary environmental statistics. Furthermore, since $\partial I_0 / \partial r_i = a_i \geq 0$ due to Jensen's inequality, the optimal solution remains unchanged if we replace the equality constraint by an inequality $\mathbf{1}^T \boldsymbol{r} \leq N$. Again, since we allow $P(\boldsymbol{X})$ to be arbitrary, this optimal solution holds for any environmental statistic and model of receptor activation. Figure 2.9 illustrates this result by optimizing across two Poisson channels for different number of receptors ($N$).

Figure 2.9: **Optimal receptor distribution across two Poisson channels for various values of $N$, Related to Figure 2.2 and Figure 2.7.** Given two Poisson channel with independent inputs $X_1$ and $X_2$, where $E[X_2] > E[X_1]$. Each plot corresponds to a particular choice of $N$ to be divided between the two channels. Solid curves represent the partial derivative the mutual information $I(X;Y)$ with respect to the two channel receptor number $r_1$ and $r_2$, evaluated for different receptor allocations. Dotted horizontal lines indicate the receptor distribution between the two channels that maximizes $I(X;Y)$.

In line with the KKT condition of Equation 2.27, Figure 2.9 shows the optimal receptor distribution across the two channels (dotted lines) occurs precisely where their partial derivatives (solid lines) are equal. Even though channel 2 (orange) experience an average ligand concentration that is only 5% higher than channel 1, the optimal solution allocates nearly all receptors (98%) to channel 2 when $N$ is small. As $N$ increases, this asymmetry of the optimal solution reduces significantly, resulting in a 5% difference in receptor number between the two channels when $N = 200$. In agreement with results we derived, strong receptor localization occurs when $N$ is small due to the fact that $\frac{d}{dr_i} I(X;Y)$ becomes nearly independent of receptor number (note the difference in y-range across the three plots in Figure 2.9).

**Absolute ligand concentration and dynamic range** In addition to receptor number, environmental factors can strongly influence receptor placement. Intuitively, one would expect the larger the difference between two channels' input ligand concentration, the larger the asymmetry should be in their receptor allocation. Figure 2.10A confirms this intuition, showing that as the relative difference in average ligand concentration sensed between two channels increase, receptor distribution between the two channels becomes more asymmetric. Figure 2.10A also suggests two additional features of the optimal strategy that are less intuitive,

1. As ligand concentration increases, optimal strategy switches from allocating more to allocating fewer receptors to region of higher ligand concentration

2. When ligand concentration are either high or low, optimal receptor placement can become highly localized, concentrating most receptors to a few channels



Figure 2.10: **Optimal receptor distribution across two Poisson channels for different absolute ligand concentration and relative difference in concentration, Related to Figure 2.2.** (A) Each curve represents the optimal receptor proportion $r_1^*/N$ for channel 1 for increasing levels of ligand input for channel 2 while keeping $E[C_1]$ fixed. (B) While keeping the relative difference between $E[C_1]$ and $E[C_2]$ fixed at 5%, black curve shows the optimal receptor proportion in channel 1 for different level of $E[C_1]$. Blue curve shows the approximation of the partial derivative shown in Equation 2.34. Black dotted line indicates peak of blue curve, red dotted line indicated $r_1^*/N = 0.5$. $K_d = 40nM$. (C) binary entropy function ($H(p) = -p \log p - (1-p) \log (1-p)$) in nats

The first feature can be seen by observing the fact that as $E[C_1]$ increases in Figure 2.10A, the optimal receptor distribution $r_1^*/N$ (gray to black) changes from being below 0.5 to above 0.5, even though $E[C_2] > E[C_1]$ for all cases plotted. The second case can be seen by observing the slope of the graphs. As $E[C_1]$ becomes either high (black) or low (gray), the optimal receptor distribution becomes more sensitive to $(E[C_2] - E[C_1])/E[C_1]$, the relative difference in average input level. For a minor difference in concentration of $< 5\%$, nearly all receptors become allocated to one of the two channels

Both of these observations are indeed general features of the optimal strategy and we can explain both using $\partial I_0 / \partial r_i$ defined in Equation 2.34. We can gain further intuition of both features from the shape of the binary entropy function (Figure 2.10C).

1. Figure 2.10B shows that when $E[X_1]$ is low, $\partial I_0/\partial r_1$ is an increasing function in $E[X_1]$, suggesting that more receptors should be allocated to channels with higher probability of receptor activation (i.e., ligand concentration). However, this monotonicity switches as $E[X_1]$ increases, with $\partial I_0/\partial r_1$ becoming a decreasing function in $E[X_1]$. The point at which monotonicity of $\partial I_0/\partial r_1$ switches (dashed black line) matches precisely with when the optimal strategy switches from allocating more receptors to region of higher probability of receptor activation to region of lower probability, as shown by the solid black curve passing the dashed red line. Thus feature #1 can be fully explained by the gradient of mutual information. This switch in strategy is intuitive when we consider the binary entropy function. Recall that $X_i$ maps to the probability of receptor activation in the $i$-th channel, so it shares a similar interpretation as the success probability $p$ of the binary entropy function. The entropy function $H(p)$ is maximized when the success probability (analogously receptor activation) is neither high nor low (Figure 2.10C). Thus it can be less useful to place more receptors at regions of higher ligand concentration, since those receptors will simply stay activated, being uninformative of the input.

2. Figure 2.10B shows that the slope of $\partial I_0/\partial r_1$ is maximized when $E[X_1]$ is either low or high. The larger the difference in the partial derivatives between two channels, the more receptors will need to be allocated before their partial derivative agree, a necessary condition for achieving optimality according to Equation 2.27. Therefore, a small relative difference in input concentration between channels can lead to large difference in information gain per receptor, when absolute ligand concentration is either high or low (relative to $K_d$), leading to strong localization of receptors. This behavior can also be explained using the binary entropy function, specifically the fact that the rate of change in entropy is maximized at low and high success probability (Figure 2.10C). Analogously, placing receptors in regions where likelihood of activation is 0 or 1 is useless from an information perspective (zero entropy/uncertainty in output), so all receptors should be allocated to a region with non-zero entropy, no matter how small the difference in likelihood of activation (thus ligand concentration) is.

**Modeling chemical microenvironment**

**Soil chemical microenvironment**    In soil, free-living unicellular eukaryotes can sense and respond to signaling ligands secreted by soil bacteria. We follow mathematical models described in (Raynaud and Nunan, 2014) and (Melke et al., 2010), modeling the spatial distribution of ligands in two steps: 1) model the spatial distribution of bacteria in soil, and 2) model each bacteria as an independent point sources of ligands.

**Modeling bacteria distribution in soil**    We follow the procedure outlined Raynaud and Nunan (2014), which allows us to generate realistic bacterial distributions found in soil. This procedure involves sampling from a spatial statistical model, based on Log Gaussian Cox Processes (LGCP) fitted to image data of observed bacterial distribution in soil. In a LGCP, the observed number of bacteria per unit area is modeled as a Poisson process in which the rate parameter is treated as being the exponential of a Gaussian process. Specifically, we consider Gaussian processes with an exponential covariance function,

$$C(r) = \sigma^2_{\text{bacteria}} e^{-r/\beta},$$ (2.41)

so the Gaussian process (and the LGCP) is fully determined by three parameters, its mean ($\mu$), variance ($\sigma^2_{\text{bacteria}}$), and scale ($\beta$). In the limit as $\sigma^2_{\text{bacteria}} \to 0$, we obtain a homogeneous Poisson process. The average intensity of a LGCP (number of bacteria per unit area) is given by,

$$\lambda = e^{\mu + \sigma^2_{\text{bacteria}}/2}$$ (2.42)

We used parameters reported in (Raynaud and Nunan, 2014), with $\mu = -7.52$, $\sigma^2_{\text{bacteria}} = 1.9$, and $\beta = 25$, to simulate a bacterial density of approximately $10^9$ cells/g on a $1000 \times 3000 \mu\text{m}^2$ rectangular domain (containing approx. 4000 cells). These are the default parameters unless otherwise stated in the main text. We used R 3.6.1 with packages *spatstat* (Baddeley, Rubak, and Turner, 2015) and *RandomFields* (Schlather et al., 2020) to generate all bacteria distributions.

**Modeling chemical distribution given bacteria distribution.**    Given a spatial distribution of bacteria, we model the distribution of secreted molecules using standard reaction-diffusion models (Melke et al., 2010). Specifically, such models treat each bacteria as a static, independent sources, producing ligands with rate $\alpha$

that diffuse ($D$) and degrade ($\gamma$). The resulting ligand concentration field is then the solution of the following partial-differential equation (PDE) ,

$$\frac{\partial c(x,t)}{\partial t} = \alpha|_{\text{bacteria}} + D\Delta c - \gamma c, \tag{2.43}$$

Rather than approximating each parameter of Equation 2.43, we model the ligand distribution produced by a bacteria using a 2D Gaussian density profile, and directly fit the Gaussian profile to empirical measurements. The concentration $c$ at a given position $x$ in the domain is then the sum over all such Gaussian profiles evaluated at $x$, which can be expressed mathematically as

$$c(x) = \sum_{q \in \mathcal{U}} \frac{C}{\sqrt{2\pi s^2}} \exp\left\{-\frac{||x - q||^2}{2s^2}\right\} \tag{2.44}$$

where $\mathcal{U}$ represents the set of bacterial positions generated using the LGCP model. $C$ represents the total concentration of each Gaussian profile, and $s$ determines the width of the profile, both of which are assumed to be uniform across all bacteria. We extract both parameters based on a geostatistical block kriging analysis of the spatial distribution of AHL in soil ($C$ chosen such that mean concentration (across the entire spatial domain) is approximately 0.6 nM, $s = 9\,\mu$m) (Burton et al., 2005; Gantner et al., 2006; Sheng et al., 2017; Y.-J. Wang and Leadbetter, 2005).

**Tissue chemical microenvironment** We follow mathematical models of ligand distribution in tissue outlined in (Rejniak et al., 2013; Milde, Bergdorf, and Koumoutsakos, 2008), simulating a tissue environment using a PDE model that incorporates four transport mechanisms: (1) free diffusion, (2) ECM binding, (3) fluid advection, (4) cellular uptake. The spatial domain is a rectangle of size $300\,\mu$m $\times$ $900\,\mu$m. We model ligands being supplied through fluid flows from the left boundary of the domain, and penetrate the interstitial space between immobilized cells. Soluble ligands are then transported by diffusion and fluid flow, and become immobilized upon binding to an extracellular matrix (ECM) made up of networks of interconnected fibers containing ligand binding sites. We explicitly represent both ECM-bound ($c_b$) and soluble forms of the ligand ($c_s$), so that the the total ligand concentration $c(x,t)$ at position $x$ and time $t$ is equal to,

$$c(x,t) = c_s(x,t) + c_b(x,t). \tag{2.45}$$

Mathematically, we can describe the dynamics of the soluble fraction $c_s(x, t)$ as follows,

$$\frac{\partial c_s}{\partial t} = \kappa|_{\text{boundary}} - u(x,t) \cdot \nabla c_s + D\Delta c_s - \beta_c c_s|_{\text{cells}} - k_{\text{ECM}}(e(x) - c_b)c_s - \gamma_s c_s. \quad (2.46)$$

1. The first term, $\kappa$, represents production/release of molecule at the left boundary.

2. The second term represents fluid transport, where $u(x, t)$ is the velocity field of the interstitial fluid with input flow speed $u^{\text{in}}$ at the left boundary. We impose zero-velocity condition on the top and bottom boundary.

3. The third term represents diffusion with $D$ as the ligand diffusion coefficient.

4. The fourth term represents cellular uptake with rate $\beta_c$, a process that only occurs near immobilized cells distributed across the domain.

5. The 5th term represents ECM binding. The concentration of ECM binding site $e(x)$ at position $x$ is generated using a minimal model of ECM protein distribution (see paragraph on "Generating ECM fiber network"). Binding occur with rate proportional to $e(x) - c_b(x, t)$, the level of available ECM binding site. Since the on-rate of ECM binding is much larger than the off-rate, we assume the off-rate to be zero.

6. The last term represents enzymatic degradation of ligand.

The dynamics of ECM-bound fraction $c_b(x, t)$ is much simpler, involving a term corresponding to ECM binding, a degradation term due to enzymatic decay .

$$\frac{\partial c_b}{\partial t} = k_{\text{ECM}}(e(x) - c_b)c_s - \gamma_b c_b. \quad (2.47)$$

To generate a ligand concentration field $c$, we take $\kappa$ to be non-zero for a brief period of time, representing a bolus of ligand release. Then, we simulate the combined dynamics of bound and soluble fractions for sufficiently long until the ligand distribution $c(x, t)$ is relatively stable. In practice, we observe that $c \approx c_b$ after a sufficiently long period of time, since the soluble fraction quickly become insignificant due to fluid flow. The resulting concentration field represents an interstitial gradient. The average concentration is set by setting the release rate $\kappa$ such that the concentration of the soluble fraction $c_s$ matches measured chemokine concentration found in interstitial fluids (1 pM to 10 pM) (Xiangdan Wang et al., 2008; Clark et al., 2015).

**Generating ECM fiber network**   To generate a distribution of ECM binding sites $e(x)$, we use a minimal computation model of fiber network (Harjanto and Zaman, 2013; Schlüter, Ramis-Conde, and Chaplain, 2012; Byoungkoo Lee et al., 2014). The model generates ECM fibers represented by line segments, which could represent fibronectin, collagen, laminin, or other fibrous matrix components. To position each fiber, one end of each segment is randomly positioned following a uniform distribution within the domain. The other end's position is determined by picking an angle, uniformly from $[0, 2\pi)$, and length sampled from a normal distribution with mean 75 µm and standard deviation of 5 µm (as measured for collagen by Friedl et al. (1997)). In total, 4050 fibers were placed in the domain. For the PDE simulation, the generated network is discretized by counting the number of fibrous proteins around each node in the simulation lattice. The density of fiber within each node is then converted to a concentration value representing the level of ECM binding sites, resulting in an average concentration of ECM binding site of 520 nM.

**Simple chemical gradient**

One of the simplest model of chemical gradient can be described by the following PDE,

$$\frac{\partial c(x,t)}{\partial t} = \alpha\delta(x_0) + D\Delta c - \gamma c, \tag{2.48}$$

where ligands are produced at rate $\alpha$ from a localized source at $x_0$, diffuses with diffusivity $D$ and undergoes first order degradation with rate $\gamma$. The steady-state solution of Equation 2.48 is a single exponential gradient,

$$c(x) = C_0 \exp{(-x/\lambda)}, \tag{2.49}$$

where the ligand concentration $c(x)$ only depends on distance $x$ from the source, the concentration at the source boundary $C_0 = \alpha/(2\sqrt{D/\gamma})$ and the decay length $\lambda = \sqrt{D/\gamma}$. By taking the source location $x_0$ to be the entire left boundary of the spatial domain, the stimulated interstitial gradient is well-described by the exponential model. Specifically, by first averaging the interstitial gradient (along the axis parallel to the ligand source) and fitting the resulting 1-D profile to Equation 2.49 using Matlab's fit function, we obtain an excellent fit with correlation coefficient $R^2 = 0.98$. This fitted exponential profile is the simple, monotonic gradient used in the paper.

| Environment class | Parameter | Symbol | Value | Ref |
|---|---|---|---|---|
| Soil | Domain size | — | 1000 µm × 3000 µm | — |
| | mean | $\mu$ | −7.52 | (Raynaud and Nunan, 2014) |
| | variance | $\sigma^2_{\text{bacteria}}$ | 1.9 | (Raynaud and Nunan, 2014) |
| | scale | $\beta$ | 25 | (Raynaud and Nunan, 2014) |
| | Concentration (per bacteria) | $C$ | 114 nM | (Burton et al., 2005; Sheng et al., 2017; Y.-J. Wang and Leadbetter, 2005) |
| | Spread of ligand (per bacteria) | $s$ | 9 µm | (Gantner et al., 2006) |
| Tissue | Domain size | — | 300 µm × 900 µm | — |
| | Diffusion coefficient | $D$ | $45\ \mu m^2\ s^{-1}$ | (Miller et al., 2018) |
| | Cellular uptake rate | $\beta_c$ | $10^{-2}\ s^{-1}$ | (Rejniak et al., 2013) |
| | Cellular uptake distance | — | 1.15× radius | (Rejniak et al., 2013) |
| | Fluid viscosity | $\mu$ | $2.5\ \mu g\ \mu m^{-1}\ s^{-1}$ | (Rejniak et al., 2013) |
| | Spatial discretization | $\Delta x$ | 2 µm | — |
| | Time step | $\Delta t$ | 0.0178 s | — |
| | Interstitial fluid input flow | $u^{\text{in}}$ | $0.1\ \mu m\ s^{-1}$ to $2\ \mu m\ s^{-1}$ | (Swartz and Fleury, 2007) |
| | Average [ECM binding site] | — | 520 nM | (Y. Wang and Irvine, 2013) |
| | ECM binding rate | $k_{\text{ECM}}$ | $9.3 \times 10^{-5}\ nM^{-1}\ s^{-1}$ | (Shields et al., 2007) |
| | Production/release rate | $\kappa$ | $7\ nM\ s^{-1}$ | (Xiangdan Wang et al., 2008) |
| | Number of cells | — | 3 × 8 cells | — |
| | Soluble ligand degradation | $\gamma_s$ | $1 \times 10^{-3}\ s^{-1}$ | (Milde, Bergdorf, and Koumoutsakos, 2008; Y. Wang and Irvine, 2013) |
| | Bound ligand degradation | $\gamma_b$ | $1 \times 10^{-5}\ s^{-1}$ | (Milde, Bergdorf, and Koumoutsakos, 2008) |
| | Mean ECM fiber length | — | 75 µm | (Schlüter, Ramis-Conde, and Chaplain, 2012) |
| | Variance in fiber length | — | 5 µm | (Schlüter, Ramis-Conde, and Chaplain, 2012) |
| | Total number of ECM fibers | — | 4050 | (Schlüter, Ramis-Conde, and Chaplain, 2012) |
| Simple gradient | Max concentration | $C_0$ | 0.7 nM | – |
| | Decay length | $\lambda$ | 60 µm | – |

Table 2.1: **Parameters used for modeling all three environment classes: soil, tissue, and simple gradient, Related to Figure 2.2A.** Tissue simulation code was adopted from (Rejniak et al., 2013). Parameters for simple gradient obtained from fitting an exponential function to the spatially averaged profile of the tissue gradient

**Incorporate cost for receptor redistribution using the Wasserstein distance**

In a dynamically changing environment, receptor should redistribute in an efficient manner in order to maximize information acquisition. We extended our optimization problem of Equation 2.15 to incorporate a "cost" for changing receptor location. For a cell sensing a sequence of ligand profiles $\{c_t\}_{t=1}^{T}$ over time, the optimal receptor placement $r_t^*$ for $c_t$ now depends additionally on $r_{t-1}^*$, the receptor placement for the previous ligand profile,

$$r_t^* = \underset{\substack{r \geq 0 \\ \sum_i r_i = N}}{\mathrm{argmax}} \; I\left(\hat{c}_t; \hat{a} \mid r\right) - \gamma W_1\left(r_{t-1}^*, r\right). \tag{2.50}$$

Here, we model the cost for redistributing receptors using the Wasserstein-1 ($W_1$) distance. For completeness, we first introduce the formal definition of the $W_1$ distance before returning to a much simpler form that applies to our problem. Let $X \sim P$ and $Y \sim Q$ represent two random variables defined over $M \subset \mathbb{R}^d$. Further, let $\mathcal{J}(P, Q)$ denote all joint distributions $J$ for $(X, Y)$ that have marginal $P$ and $Q$. The $W_1$ distance between $P$ and $Q$ is,

$$W_1(P, Q) = \inf_{J \in \mathcal{J}(P,Q)} \int_{M \times M} \|x - y\|_1 dJ(x, y). \tag{2.51}$$

One way to understand the above definition is to consider different ways of transporting a distribution of mass $P(x)$ to a different distribution $Q(x)$. Given some cost function associated with each unit of mass transported, the $W_1$ distance is the minimum transport cost achievable. In this way, the $W_1$ distance assumes that the transformation from $P$ to $Q$ occurs in an optimal manner. Note that this distance function is non-negative and symmetric, and does not require $P$ and $Q$ to be probability distributions, it applies whenever the total mass is preserved between $P$ and $Q$.

Although Equation 2.51 is difficult to compute in general, it has a closed form for the special case of $d = 1$ which is the case we are considering. Instead of using the canonical form of the $W_1$ distance in 1-D, we need to use a generalized form that applies to distributions on a circle (Rabin, Delon, and Gousseau, 2011). For two receptor distributions on the 1-D surface of a 2-D cell, represented as non-negative vectors $a$ and $b$ of length $m$, the $W_1$ distance takes on the form,

$$W_1(a, b) = \sum_{i=1}^{m} |\phi_i - \mu|, \tag{2.52}$$

where $\phi_i = \sum_{j=1}^{i} \left( \frac{a_j}{\|a\|_1} - \frac{b_j}{\|b\|_1} \right)$ and $\mu$ is the median of the set of values $\{\phi_i, 1 \leq i \leq m\}$. We derive the gradient of Equation 2.52 as,

$$\frac{\partial}{\partial a_k} W_1(a, b) = \sum_{i=1}^{m} \mathrm{sgn}(\phi_i - \mu) \sum_{j=1}^{i} \left( \delta_{jk} - \frac{a_j}{\|a\|_1} \right). \tag{2.53}$$

We perform optimization with this gradient using the fmincon function (with sqp algorithm) in MATLAB.

## Numerical simulation of receptor feedback scheme

In our feedback scheme, receptor $r(x, t)$ is modeled by considering three redistribution mechanisms: (1) lateral diffusion of $r$ along the plasma membrane $(D\nabla^2_{\mathrm{memb}} r)$, (2) endocytosis of $r$ along the plasma membrane $(k_{\mathrm{off}} r)$, (3) incorporation of cytoplasmic pool of receptors, $R_{\mathrm{cyto}}$, to the membrane at rate proportional to local receptor activity $(hAR_{\mathrm{cyto}})$. $A(x, t)$ is a random variable that denotes receptor activity along the cell membrane, and is a function of local receptor number. Then, the equation describing the distribution of $r$ across the cell membrane can be expressed mathematically as,

$$\frac{\partial r(x, t)}{\partial t} = D\nabla^2_{\mathrm{memb}} r - k_{\mathrm{off}} r + hAR_{\mathrm{cyto}}, \tag{2.54}$$

where the total number of receptors $r_{\mathrm{tot}} = \int_{\mathrm{memb}} r + R_{\mathrm{cyto}}$ is fixed. We simulate receptor distribution by treating the cell membrane as a 1D space and the cytosol as a single, homogeneous compartment. This simplification allows us to simulate our PDE using the Crank-Nicolson method in one spatial dimension. Given space and time units $\Delta x$ and $\Delta t$, respectively, the Crank-Nicolson method with $R_i^j := r(i\Delta x, j\Delta t)$ and $A_i^j := A(i\Delta x, j\Delta t)$ is given by the difference scheme

$$\frac{R_i^{j+1} - R_i^j}{\Delta t} = \frac{D}{2\Delta x^2} \left( R_{i+1}^j - 2R_i^j + R_{i-1}^j + R_{i+1}^{j+1} - 2R_i^{j+1} + R_{i-1}^{j+1} \right)$$
$$- \frac{k_{\mathrm{off}}}{2} \left( R_i^j + R_i^{j+1} \right) + \frac{hA_i^j}{2} \left( R_{\mathrm{cyto}}^j + R_{\mathrm{cyto}}^{j+1} \right) \tag{2.55}$$

where, $i = 1, 2, 3, \ldots m$, representing $m$ discrete membrane compartments and $R_{\mathrm{cyto}}^j$ represents the additional cytosol compartment. Since the membrane is represented by a circle, we have the following pair of conditions,

$$R_0^j = R_m^j, \quad R_{m+1}^j = R_1^j. \tag{2.56}$$

Lastly, total receptor number across all compartments is conserved,

$$\sum_{i=1}^{m} R_i^j + R_{\text{cyto}}^j = \sum_{i=1}^{m} R_i^{j+1} + R_{\text{cyto}}^{j+1}. \tag{2.57}$$

Now, we can combined Equation 2.55–Equation 2.57 and rewrite everything in vector form. First, let

$$\alpha := \frac{D}{2\Delta x^2}, \quad \beta := \frac{k_{\text{off}}}{2}, \quad \kappa_i^j := \frac{hA_i^j}{2},$$

and rewrite equation (2.55) as,

$$\frac{R_i^{j+1}}{\Delta t} - \alpha\left(R_{i+1}^{j+1} - 2R_i^{j+1} + R_{i-1}^{j+1}\right) + \beta R_i^{j+1} - \kappa_i^{j+1} R_{\text{cyto}}^{j+1} = \frac{R_i^j}{\Delta t} + \alpha\left(R_{i+1}^j - 2R_i^j + R_{i-1}^j\right)$$
$$- \beta R_i^j + \kappa_i^j R_{\text{cyto}}^j \tag{2.58}$$

and define $U^j$ to be the $(m+1)$-dimensional vector with components $R_i^j$ for $i = 1, 2, 3, \dots m$ and $U_{m+1}^j = R_{\text{cyto}}^j$. The difference scheme is given in the vector form

$$PU^{j+1} = QU^j. \tag{2.59}$$

where,

$$P = \begin{bmatrix} \frac{1}{\Delta t} + 2\alpha + \beta & -\alpha & 0 & \cdots & 0 & -\alpha & -\kappa_1^{j+1} \\ -\alpha & \frac{1}{\Delta t} + 2\alpha + \beta & -\alpha & 0 & & \cdots & 0 & -\kappa_2^{j+1} \\ 0 & \ddots & \ddots & \ddots & & & \vdots \\ \vdots & & & \ddots & \ddots & \ddots & \\ 0 & \cdots & 0 & -\alpha & \frac{1}{\Delta t} + 2\alpha + \beta & -\alpha & -\kappa_{m-1}^{j+1} \\ -\alpha & 0 & \cdots & 0 & -\alpha & \frac{1}{\Delta t} + 2\alpha + \beta & -\kappa_m^{j+1} \\ 1 & 1 & \cdots & & & 1 & 1 \end{bmatrix} \tag{2.60}$$

$$Q = \begin{bmatrix} \frac{1}{\Delta t} - 2\alpha - \beta & \alpha & 0 & \cdots & 0 & \alpha & \kappa_1^j \\ -\alpha & \frac{1}{\Delta t} - 2\alpha - \beta & \alpha & 0 & & \cdots & 0 & \kappa_2^j \\ 0 & \ddots & \ddots & \ddots & & & \vdots \\ \vdots & & & \ddots & \ddots & \ddots & \\ 0 & \cdots & 0 & \alpha & \frac{1}{\Delta t} - 2\alpha - \beta & \alpha & \kappa_{m-1}^j \\ \alpha & 0 & \cdots & 0 & \alpha & \frac{1}{\Delta t} - 2\alpha - \beta & \kappa_m^j \\ 1 & 1 & \cdots & & & 1 & 1 \end{bmatrix} \tag{2.61}$$

Because $A$ is invertible, the Crank-Nicolson scheme reduces to the iterative process

$$U^{j+1} = P^{-1}QU^j. \tag{2.62}$$

The entire evolution of $r$ can be solved where at each time step, we update receptor activity $A_i^j$ across all membrane position $i$ according to the random process described by equation (2.13), (2.14), followed by solving equation (2.62) for $U^{j+1}$.

| Parameter | Symbol | Value | Ref. |
|---|---|---|---|
| Receptor endocytosis rate | $k_{\text{off}}$ | $0.06\,\text{s}^{-1}$ to $0.18\,\text{s}^{-1}$ | (Marco et al., 2007) |
| Feedback constant | $h$ | $2 \times 10^{-3}\,\text{s}^{-1}$ to $4 \times 10^{-3}\,\text{s}^{-1}$ | (Marco et al., 2007) |
| Average receptor feedback rate | $\langle hA_i \rangle_i$ | $10^{-3}\,\text{s}^{-1}$ to $2 \times 10^{-3}\,\text{s}^{-1}$ | (Marco et al., 2007) |
| Receptor dissociation constant | $K_d$ | $40\,\mu\text{M}$ | — |
| Cell radius | | $10\,\mu\text{m}$ | — |
| Basal receptor activity | $\alpha$ | $0.1$ | — |
| Spatial discretization | $\Delta x$ | $0.6283\,\mu\text{m}$ | — |
| Time discretization | $\Delta t$ | $1\,\text{s}$ | — |
| Total receptor | $r_{\text{tot}}$ | $1000$ | — |
| Number of membrane bins | $m$ | $100$ | — |

Table 2.2: **Parameter values for feedback scheme simulation, Related to Figure 2.5.** The average rate of receptor incorporation, $\langle hA_i \rangle_i$, depends on receptor activity which changes as the cell moves in a heterogeneous environment. Thus, the range shown represents the time-averaged value for a cell moving through simulated tissue environments. The value of $h$ was chosen to achieve a physiologically relevant range for $\langle hA_i \rangle_i$.

We set the value of the feedback constant $h$ using empirical measurements from Marco et al. (2007) (Marco et al., 2007). In Figure 3M of Marco et al., the authors report a quartile box plot showing estimated values for a parameter they call h (which we will refer to as $\bar{h}$), with a mean estimate of around $1.6 \times 10^{-3}\,\text{s}^{-1}$. Note $\bar{h}$ is equivalent in meaning as our $hA_i$. However, since $hA_i$ will be different across different membrane bins and across time, we simulate the feedback scheme for a cell in a given environment and set the value of h such that the mean rate $\langle hA_i \rangle$ (averaged across membrane and time) is approximately equal to the mean estimate of $1.6 \times 10^{-3}\,\text{s}^{-1}$ reported by Marco et al.. The value $\bar{h}$ reported by Marco et al. corresponds specifically to the transport rate of the Cdc42 to the membrane. The parameter value was obtained by analyzing fluorescence recovery of GFP-Cdc42 in membrane regions bleached with a laser pulse. Although the measured value corresponds to Cdc42, it has been used to model the effective exocytosis rate for receptors shown to undergo activity-dependent localization, showing good agreement with empirical data (Hegemann et al., 2015). Similar values around

$10^{-3}$ s$^{-1}$ to $2 \times 10^{-3}$ s$^{-1}$ have been measured for the recycling rate of a wide range of GPCRs (JJ, 1993; Pippig, Andexinger, and Lohse, 1995; Koenig and Edwardson, 1996; Koenig and Edwardson, 1994).

**Numerical simulation of cell navigating in interstitial chemical gradients**

**Chemotaxis algorithm**    At $t = 0$, initialize a cell at position $p_0 \in \Omega \subset \mathbb{R}^2$.

At each subsequent time step $t = t + \Delta t$ with the cell at position $p_t \in \Omega$:

1. Compute mean ligand profile $c \in \mathbb{R}^m$ at the cell's current position.

2. Independently sample $n$ ligand profiles $\{C^{(i)}\}_{i=1}^n$ where each element $C_j$ is distributed as a Poisson random variable with mean equal to $c_j$ ($n = 30$ used in main text, refer to Figure 2.16 for other values of $n$).

3. For each ligand profile $C^{(i)}$ sampled, sample a corresponding receptor activity profiles $A^{(i)}$,

$$A^{(i)}|C^{(i)} \sim \prod_{j=1}^m \text{Pois}(\lambda_j), \text{ where } \lambda_j = r_j \left( \frac{C_j^{(i)}}{C_j^{(i)} + K_d} + \alpha \frac{K_d}{C_j^{(i)} + K_d} \right). \quad (2.63)$$

4. Compute average receptor activity $\bar{A} = \frac{1}{n} \sum_{i=1}^n A^{(i)}$

5. Compute the an estimator of gradient direction $\hat{\theta}$ using one of three approaches

   - Optimal decoder + noise (Hu et al., 2010): $\hat{\theta} = \arctan \left( \frac{\sin(\phi)^T \bar{A}}{\cos(\phi)^T \bar{A}} \right) + \mathcal{N}(0, 0.1)$, where $\phi_i = 2\pi i/m$, $i = 1, ..., m$ corresponds to the angle where $A_i$ is measured on the cell surface.

   - Random: $\hat{\theta}$ is sampled uniformly from the set of $m$ angles/directions $\{\phi_i\}_{i=1}^m$

   - Maximal increase: $\hat{\theta} = \phi_{i^*}$ where

   $$i^* = \text{argmax}_{1 \leq i \leq m} f(i) \text{ and } f(i) = \begin{cases} \bar{A}_i - \bar{A}_{i+m/2} & i \leq m/2 \\ \bar{A}_i - \bar{A}_{i-m/2} & i > m/2 \end{cases}$$

   This decoder selects the direction of maximum change in receptor activity across the cell surface.

   (*) In addition to the three decoders above, we consider the possibility of temporal averaging, where $\bar{A}$ is a running mean over the past 5 minutes

of receptor activity profiles (a total of $300 \times 30 = 9000$ sample profiles). This running average is decoded with the optimal decoder + noise as described above.

6. Set new cell position $p_{t+\Delta t} = p_t + s\Delta t[\cos(\hat{\theta}), \sin(\hat{\theta})]$, with speed $s = 2\,\mu\text{m}\,\text{min}^{-1}$, $\Delta t = 1$ s.

7. Repeat from step 1.

**Tissue gradient simulation for cell navigation**

**Localization task** In addition to using the same tissue environment as the rest of the paper, we simulated additional tissue gradients using the same set of parameters but different (randomly generated) ECM fiber networks. This results in tissue gradients that have the same macroscopic features but different patterns of microscopic fluctuations.

**Retention task** This task was motivated by the precision with which growth cones can retain themselves within specific regions of gradients of axon guidance cues. For this task, we used an ellipse-shaped cell with semi-major axis $= 5\,\mu\text{m}$, semi-minor axis $= 2\,\mu\text{m}$ to mimic the shape of a navigating growth cone. In-vivo observations of axon guidance cue gradients show very short decay length (Xiao and Baier, 2007; Xiao, Staub, et al., 2011), so we adjust several parameters to generate interstitial gradients with matching decay length. Below are the new parameter values that differ from Table 2.1,

| Parameter | Symbol | Value |
|---|---|---|
| Diffusion coefficient | $D$ | $5\,\mu\text{m}^2\,\text{s}^{-1}$ |
| Interstitial fluid flow speed | $u^{\text{in}}$ | $0.3\,\mu\text{m}\,\text{s}^{-1}$ |
| Production/release rate | $\kappa$ | $20\,\text{nM}\,\text{s}^{-1}$ |
| Soluble ligand degradation | $\gamma_s$ | $3 \times 10^{-1}\,\text{s}^{-1}$ |
| Bound ligand degradation | $\gamma_b$ | $3 \times 10^{-3}\,\text{s}^{-1}$ |

Table 2.3: **Parameter values used for tissue gradient generated for retention task that differs from the values in Table 2.1, Related to Figure 2.6.**

**Data on receptor placement, surface expression level, and binding affinity**

Empirical measurements of receptor cell surface expression level and binding affinity can be highly variable, depending on how the affinity was measured and the particular

cell type used. The data shown below simply represent a subset of values reported in literature.

| Cell type | Receptor/Ligand | $K_d$ (nM) | $N$ (cell surface) | Dynamic localization |
|---|---|---|---|---|
| T-cell | CCR2/CCL2 | 1.53 (Yoshimura, 2018) | 1100 (Yoshimura, 2018) | Yes (Nieto et al., 1997) |
| | CXCR4/CXCL12 | 5 (Qing et al., 2020) | 1572 (834-1961) (Benhur Lee et al., 1999) | Yes (Buul et al., 2003; Pelletier et al., 2000) |
| | CCR5/CCL5 | 4 (Qing et al., 2020) | 593 (167-1006) (Benhur Lee et al., 1999) | Yes(Nieto et al., 1997) |
| | IL-2R/IL-2 | $1.3 \times 10^{-2}$ (Fitzgerald et al., 2001) | 1000 (Höfer, Krichevsky, and Altan-Bonnet, 2012) | No (Nieto et al., 1997) |
| | TNFR-1/TNF | $1.9 \times 10^{-2}$ (Grell et al., 1998) | 1000 (Gehr et al., 1992; Thoma et al., 1990) | No(Nieto et al., 1997) |
| | TGFβR-2 | $5 \times 10^{-2}$ (Tripathi et al., 1993) | 10000 (Tucker et al., 1984) | No (Nieto et al., 1997) |
| Neutrophile | CR3/C3bi | 12.5 (Cai and Wright, 1995) | 40000 (Berger et al., 1984) | No (Pytowski, Maxfield, and Michl, 1990) |
| | C5aR/C5a | 2 (Huey and Hugli, 1985) | 50000 (Huey and Hugli, 1985) | No (Servant et al., 1999) |
| Neuron | GABA$_A$R/GABA | 12 (Jones et al., 1998) | 200 (Caruncho et al., 1995) | Yes (Bouzigues et al., 2007) |
| | Robo1/Slit | 235 ± 165 (Evans and Bashaw, 2010) | 22300* (Komatsu et al., 2017) | Yes (Pignata et al., 2019) |
| | PlxnA1 | NA | NA | Yes (Pignata et al., 2019) |
| Fibroblast | LPA-2/LPA | NA | NA | Yes (Ren et al., 2014) |

Table 2.4: **Receptor data used for Figure 2.7 of main text, Related to Figure 2.7.**. (*) receptor expression level data for Robo1 was taken from a non-neuronal cancer cell line.

**Extension of information theoretic framework to study other spatial sensing strategies such as modulation of cell shape**

We briefly illustrate an extension of our framework to study how cell shape can be tuned to improve cell sensing and navigation. Recall the generalized model of receptor activation from the Discussion,

$$\mathbb{E}[A_i \mid c_i] = f(\boldsymbol{\theta}_i)\Big(\frac{c_i}{c_i + K_d} + \alpha\frac{K_d}{c_i + K_d}\Big), \tag{2.64}$$

where $f$ is an unspecified function of an arbitrary set of variables $\boldsymbol{\theta}$, representing the "effective" number of receptors at position $i$.

Recent work has shown that given uniform membrane receptors sensing a uniform ligand field, membrane regions of higher curvature can exhibit higher receptor activity, due to higher local volume-to-surface ratio (Rangamani et al., 2013). Suppose we are interested in tuning membrane shape/curvature as a way to maximize information acquisition by cells. Assuming a constant, linear relationship between curvature at the i-th membrane position $\beta_i$ and "effective" receptor number $f$, and that receptors are uniformly distributed, then we have

$$\mathbb{E}[A_i \mid c_i] = \alpha\frac{r_{\text{tot}}}{N}\beta_i\Big(\frac{c_i}{c_i + K_d} + \alpha\frac{K_d}{c_i + K_d}\Big), \tag{2.65}$$

where $\alpha$ is a proportionality constant and $N$ is the number of membrane bins. This model is identical to our receptor model (2.19) up to a constant factor. Furthermore, if we assume a fixed total membrane area, the resulting optimization problem is nearly identical with (2.15), where total membrane area now play a similar role as total receptor number, and $\beta_i$ takes the place of $r_i$. Therefore, we expect general features of the optimal cell shape to match that of the optimal receptor placement. Namely, cells can maximize information acquisition by increasing membrane curvature at regions of high ligand concentration, by making narrow protrusions. One can derive a more accurate solution by considering a detailed model of the relationship between curvature and receptor activity outlined in (Rangamani et al., 2013).

By extension, a strategy to dynamically form narrow membrane protrusions at regions of high ligand concentration, without explicitly tuning receptor positions, should in principle boost navigation efficiency in a manner similar to the receptor feedback scheme we proposed, as the two strategies have qualitatively similar effects on the spatial distribution of receptor activity. Recent works show that indeed a feedback circuit that produces dynamic, narrow membrane protrusions is crucial for

neutrophil navigation. Cells that cannot form narrow protrusions can still move, but exhibit profoundly defective chemotaxis (Diz-Muñoz et al., 2016).

## 2.11 Supplemental figures



Figure 2.11: **Different metric assessing information gain offered by the optimal placement strategy over the uniform strategy, Related to Figure 2.2.** (A) versions of Figure 2.2C for different information metrics, where $I_{opt}$ and $I_{unif}$ is defined in Equation 2.5 in the main text, 1st row is absolute information gain between optimal and uniform receptors, 2nd row is absolute increase in the number of different classes to which the ligand profile can be subdivided after observing the receptor activity (with inset showing intermediate values of $\alpha$), 3rd row is relative information gain, 4th row is the average information obtained with uniform receptors; (B) versions of Figure 2.2E for different information metrics, $I_{opt,c} = I(\hat{c}; \hat{a} \mid \phi^*)$ and $I_{unif,c} = I(\hat{c}; \hat{a} \mid \phi^u)$; (C) versions of Figure 2.2D for different information metrics.

Figure 2.12: **Robustness of optimal efficacy to perturbation in receptor placement for other cell radius and efficacy metric, Related to Figure 2.3.** Colors of heat map represent ratio of perturbed efficacy $\eta(\phi^p)$ to optimal efficacy $\eta(\phi^*)$ for different combinations of shifting and flattening, computed for ligand profiles $\{c\}$ sampled from either (A) tissue or (B) soil; call-out boxes corresponds to different sets of perturbations, showing the average of the optimal $\{\phi^*(c)\}$ (gray) and perturbed $\{\phi^p(c)\}$ (red) receptor placements, after all profile peaks were centered; (C) same as (A) but for a cell of $10\,\mu\mathrm{m}$ radius and for efficacy metric, $2^{I_{\phi^*}} - 2^{I_{\phi^u}}$, corresponding to the increase in number of distinguishable input states between optimal and uniform placements, similarly for soil (D)

Figure 2.13: **Effect of redistribution cost $\gamma$ and constitutive receptor activity $\alpha$ on receptor redistribution according to the dynamic protocol, Related to Figure 2.4.** (A) Schematic showing a cell circling a ligand source which generates a stationary gradient (red). (B) Ligand profile experienced by the moving cell in Panel A at two different time points. (C) Optimal receptor profile computed using Equation 2.9 for the two time points shown in Panel B , for different values of redistribution cost $\gamma$. (D) Speed of the moving receptor cap (as shown in Panel C) for a wide range of $\gamma$, where speed is computed using the distance moved by the center-of-mass of the receptor distribution. (E) Different receptor redistribution dynamics for different degrees of constitutive receptor activity $\alpha$ (shown for two different pairs of ligand profiles), dotted lines represent ligand profile (red) and receptor profile (yellow patch) at one time step, while solid lines represent the ligand and receptor profile at the next time step. (F) Different receptor redistribution dynamics for different receptor redistribution cost $\gamma$, shown for the same pairs of ligand profiles as Panel E.

Figure 2.14: **Success rate of chemotactic cell navigating in simulated interstitial gradient for various values of $\gamma$ in dynamic protocol, Related to Figure 2.4.** Dynamic protocol of Equation 2.50 is solved step-wise for a cell simulated to navigate through an interstitial gradient, success rate is the proportion of simulated cells reaching gradient peak within 1 hour using the optimal + noise decoding method. $\gamma = 0$ corresponds to the case where the cost term is absent and receptors move simply to maximize mutual information. Parameter $\gamma$ determines the balance between maximizing information and minimizing receptor transport cost. Note the dynamic protocol should not be confused with the feedback scheme.

Figure 2.15: **Effect of rate parameter values on the effectiveness of receptor feedback scheme for information acquisition, Related to Figure 2.5.** (A) heat maps show the extent of receptor localization at the ligand peak for different choices of scheme parameters $k_{\text{off}}$ and $h$, for a cell in tissue (top) and soil (bottom), as measured by the fold change in receptor number near the ligand peak compared to a uniform distribution of receptors; call-out boxes show receptor morphology for different parameter values, star indicates default parameter values used in Figure 2.5. (B) ratio of scheme efficacy $\eta(\phi^s)$ to optimal efficacy $\eta(\phi^*)$ for static signals $\{c\}$ sensed by a 5 µm cell sampled from soil and tissue (C) ratio of scheme efficacy $\eta(\phi^s)$ to optimal efficacy $\eta(\phi^*)$ for a sequence of signals $\{c_t\}$ sampled by translating a 5 µm cell through soil and tissue environment at a speed of 2 µm min$^{-1}$.

Figure 2.16: **Cell navigation and retention performance in simulated interstitial gradient for various feedback scheme parameters, Related to Figure 2.6.** (A) heatmap showing success rate (proportion of simulated cells reaching gradient peak within 1 hour) for cells using receptor feedback scheme with different values of the endocytosis rate ($k_{off}$) and average incorporation rate ($\langle ha_i \rangle_i$), white star denotes parameter values used in Figure 2.6 of main text. (B) success rate for cells navigating with different sampling rate, which is the number of ligand profiles sampled per second. (C) histogram and corresponding success rate quantification for cells decoding gradient using three different decoding methods, see supplemental information for detail on each method (temporal averaging is done over a 5 minute window). (D) heatmap showing error rate (proportion of simulated time steps where cell was more than $5\mu m$ away from the gradient peak) for cells using receptor feedback scheme with different values of the endocytosis rate ($k_{off}$) and average incorporation rate ($\langle ha_i \rangle_i$), white star denotes parameter values used in Figure 2.6 of main text. (E) error rate for cells navigating with different sampling rate, which is the number of ligand profiles sampled per second. (F) histogram and corresponding error rate quantification for cells decoding gradient using three different decoding methods, see Section 2.10 for details

Figure 2.17: **Relative and absolute information gain for natural receptors of different surface expression level, constitutive activity, and binding affinity, Related to Figure 2.7.** (A) relative information gain for different values of $K_d$ and $N$; values computed using the tissue environment; $\alpha = 0.1$; red dots correspond to receptors that polarize in heterogeneous environments, white dots represent receptors that are constantly uniform, see Table 2.4 for receptor data. (B) relative (left) and absolute information gain (right) computed for the ten natural receptors shown in Panel A, across different values of $\alpha$, blue bars correspond to receptors that are constantly uniform, red bars correspond to receptors that polarize in ligand gradient.

## References

Baddeley, Adrian, Ege Rubak, and Rolf Turner (2015). *Spatial Point Patterns: Methodology and Applications with R*. London: Chapman and Hall/CRC Press. URL: `http://www.crcpress.com/Spatial-Point-Patterns-Methodology-and-Applications-with-R/Baddeley-Rubak-Turner/9781482210200/`.

Berg, Howard C, and Edward M Purcell (1977). "Physics of chemoreception". In: *Biophysical Journal* 20.2, pp. 193–219.

Berger, Melvin et al. (1984). "Human neutrophils increase expression of C3bi as well as C3b receptors upon activation". In: *The Journal of Clinical Investigation* 74.5, pp. 1566–1571.

Bouzigues, Cédric et al. (2007). "Asymmetric redistribution of GABA receptors during GABA gradient sensing by nerve growth cones analyzed by single quantum dot imaging". In: *Proceedings of the National Academy of Sciences* 104.27, pp. 11251–11256.

Buchwald, Peter (2019). "A receptor model with binding affinity, activation efficacy, and signal amplification parameters for complex fractional response versus occupancy data". In: *Frontiers in Pharmacology* 10, p. 605.

Burton, EO et al. (2005). "Identification of acyl-homoserine lactone signal molecules produced by Nitrosomonas europaea strain Schmidt". In: *Applied and Environmental Microbiology* 71.8, pp. 4906–4909.

Buul, Jaap D van et al. (2003). "Leukocyte-endothelium interaction promotes SDF-1-dependent polarization of CXCR4". In: *Journal of Biological Chemistry* 278.32, pp. 30302–30310.

Cai, Tian-Quan and Samuel D Wright (1995). "Energetics of leukocyte integrin activation". In: *Journal of Biological Chemistry* 270.24, pp. 14358–14365.

Candès, Emmanuel J and Michael B Wakin (2008). "An introduction to compressive sampling". In: *IEEE Signal Processing Magazine* 25.2, pp. 21–30.

Caruncho, HJ et al. (1995). "The density and distribution of six GABAA receptor subunits in primary cultures of rat cerebellar granule cells". In: *Neuroscience* 67.3, pp. 583–593.

Caselton, William F and James V Zidek (1984). "Optimal monitoring network designs". In: *Statistics & Probability Letters* 2.4, pp. 223–227.

Chau, Angela H et al. (2012). "Designing synthetic regulatory networks capable of self-organizing cell polarization". In: *Cell* 151.2, pp. 320–332.

Cheong, Raymond et al. (2011). "Information transduction capacity of noisy biochemical signaling networks". In: *Science* 334.6054, pp. 354–358.

Chou, Ching-Shan et al. (2011). "Noise filtering tradeoffs in spatial gradient sensing and cell polarization response". In: *BMC Systems Biology* 5.1, pp. 1–16.

Clark, Kristina EN et al. (2015). "Multiplex cytokine analysis of dermal interstitial blister fluid defines local disease mechanisms in systemic sclerosis". In: *Arthritis research & therapy* 17.1, pp. 1–11.

De Anna, Pietro et al. (2021). "Chemotaxis under flow disorder shapes microbial dispersion in porous media". In: *Nature Physics* 17.1, pp. 68–73.

Diz-Muñoz, Alba et al. (2016). "Membrane tension acts through PLD2 and mTORC2 to limit actin network assembly during neutrophil migration". In: *PLoS Biology* 14.6, e1002474.

Dlamini, Mcolisi, Timothy E Kennedy, and David Juncker (2020). "Combinatorial nanodot stripe assay to systematically study cell haptotaxis". In: *Microsystems & Nanoengineering* 6.1, pp. 1–12.

Dubuis, Julien O et al. (2013). "Positional information, in bits". In: *Proceedings of the National Academy of Sciences* 110.41, pp. 16301–16308.

Dytso, Alex, Michael Fauß, and H Vincent Poor (2020). "The Vector Poisson Channel: On the Linearity of the Conditional Mean Estimator". In: *IEEE Transactions on Signal Processing* 68, pp. 5894–5903.

Endres, Robert G and Ned S Wingreen (2008). "Accuracy of direct gradient sensing by single cells". In: *Proceedings of the National Academy of Sciences* 105.41, pp. 15749–15754.

Evans, Timothy A and Greg J Bashaw (2010). "Functional diversity of Robo receptor immunoglobulin domains promotes distinct axon guidance decisions". In: *Current Biology* 20.6, pp. 567–572.

Fitzgerald, Katherine A et al. (2001). *The cytokine factsbook and webfacts*. Elsevier.

Fowell, Deborah J and Minsoo Kim (2021). "The spatio-temporal control of effector T cell migration". In: *Nature Reviews Immunology*, pp. 1–15.

Gantner, Stephan et al. (2006). "In situ quantitation of the spatial scale of calling distances and population density-independent N-acylhomoserine lactone-mediated communication by rhizobacteria colonized on plant roots". In: *FEMS Microbiology Ecology* 56.2, pp. 188–194.

Gehr, Gisela et al. (1992). "Both tumor necrosis factor receptor types mediate proliferative signals in human mononuclear cell activation". In: *The Journal of Immunology* 149.3, pp. 911–917.

Grell, Matthias et al. (1998). "The type 1 receptor (CD120a) is the high-affinity receptor for soluble tumor necrosis factor". In: *Proceedings of the National Academy of Sciences* 95.2, pp. 570–575.

Guo, Dongning, Shlomo Shamai, and Sergio Verdú (2008). "Mutual information and conditional mean estimation in Poisson channels". In: *IEEE Transactions on Information Theory* 54.5, pp. 1837–1849.

Harjanto, Dewi and Muhammad H Zaman (2013). "Modeling extracellular matrix reorganization in 3D environments". In: *PLoS One* 8.1, e52509.

Hegemann, Björn et al. (2015). "A cellular system for spatial signal decoding in chemical gradients". In: *Developmental Cell* 35.4, pp. 458–470.

Hodge, A (2006). "Plastic plants and patchy soils". In: *Journal of Experimental Botany* 57.2, pp. 401–411.

Höfer, Thomas, Oleg Krichevsky, and Grégoire Altan-Bonnet (2012). "Competition for IL-2 between regulatory and effector T cells to chisel immune responses". In: *Frontiers in Immunology* 3, p. 268.

Hu, Bo et al. (2010). "Physical limits on cellular sensing of spatial gradients". In: *Physical Review Letters* 105.4, p. 048104.

Huey, R and TE Hugli (1985). "Characterization of a C5a receptor on human polymorphonuclear leukocytes (PMN)". In: *The Journal of Immunology* 135.3, pp. 2063–2068.

Huston, Stephen J et al. (2015). "Neural encoding of odors during active sampling and in turbulent plumes". In: *Neuron* 88.2, pp. 403–418.

Iida, Fumiya and Surya G Nurzaman (2016). "Adaptation of sensor morphology: an integrative view of perception from biologically inspired robotics perspective". In: *Interface Focus* 6.4, p. 20160016.

JJ, Lauffenburger DA Linderman (1993). *Receptors: Models for binding, trafficking, and signaling*.

Jones, Mathew V et al. (1998). "Defining affinity with the GABAA receptor". In: *Journal of Neuroscience* 18.21, pp. 8590–8604.

Kennedy, Timothy E et al. (2006). "Axon guidance by diffusible chemoattractants: a gradient of netrin protein in the developing spinal cord". In: *Journal of Neuroscience* 26.34, pp. 8866–8874.

Kicheva, Anna et al. (2007). "Kinetics of morphogen gradient formation". In: *Science* 315.5811, pp. 521–525.

Kinoshita-Kawada, Mariko et al. (2019). "A crucial role for Arf6 in the response of commissural axons to Slit". In: *Development* 146.3, dev172106.

Koenig, Jennifer A and J Michael Edwardson (1994). "Kinetic analysis of the trafficking of muscarinic acetylcholine receptors between the plasma membrane and intracellular compartments". In: *Journal of Biological Chemistry* 269.25, pp. 17174–17182.

– (1996). "Intracellular trafficking of the muscarinic acetylcholine receptor: importance of subtype and cell type". In: *Molecular Pharmacology* 49.2, pp. 351–359.

Komatsu, N et al. (2017). "Enhancement of anti-robo1 immunotoxin cytotoxicity to head and neck squamous cell carcinoma via photochemical internalization". In: *Archives in Cancer Research* 5.4.

Krause, Andreas, Ajit Singh, and Carlos Guestrin (2008). "Near-optimal sensor placements in Gaussian processes: Theory, efficient algorithms and empirical studies". In: *Journal of Machine Learning Research* 9.2.

Lee, Benhur et al. (1999). "Quantification of CD4, CCR5, and CXCR4 levels on lymphocyte subsets, dendritic cells, and differentially conditioned monocyte-derived macrophages". In: *Proceedings of the National Academy of Sciences* 96.9, pp. 5215–5220.

Lee, Byoungkoo et al. (2014). "A three-dimensional computational model of collagen network mechanics". In: *PLoS One* 9.11, e111896.

Lim, Kihong et al. (2015). "Neutrophil trails guide influenza-specific CD8+ T cells in the airways". In: *Science* 349.6252.

Macara, Ian G and Stavroula Mili (2008). "Polarity and differential inheritance—universal attributes of life?" In: *Cell* 135.5, pp. 801–812.

Majzner, Robbie G et al. (2020). "Tuning the antigen density requirement for CAR T-cell activity". In: *Cancer Discovery* 10.5, pp. 702–723.

Marco, Eugenio et al. (2007). "Endocytosis optimizes the dynamic localization of membrane proteins that regulate cortical polarity". In: *Cell* 129.2, pp. 411–422.

Martinez, Marina and Edmund Kyung Moon (2019). "CAR T cells for solid tumors: new strategies for finding, infiltrating, and surviving in the tumor microenvironment". In: *Frontiers in Immunology* 10, p. 128.

"MATLAB Optimization Toolbox" (2021a). In: *The MathWorks, Natick, MA, USA*.

McClure, Allison W et al. (2015). "Role of polarized G protein signaling in tracking pheromone gradients". In: *Developmental Cell* 35.4, pp. 471–482.

Mehta, Pankaj et al. (2009). "Information processing and signal integration in bacterial quorum sensing". In: *Molecular Systems Biology* 5.1, p. 325.

Melke, Pontus et al. (2010). "A cell-based model for quorum sensing in heterogeneous bacterial colonies". In: *PLoS Computational Biology* 6.6, e1000819.

Milde, Florian, Michael Bergdorf, and Petros Koumoutsakos (2008). "A hybrid model for three-dimensional simulations of sprouting angiogenesis". In: *Biophysical Journal* 95.7, pp. 3146–3160.

Miller, Helen et al. (2018). "High-speed single-molecule tracking of CXCL13 in the B-follicle". In: *Frontiers in Immunology* 9, p. 1073.

Monti, Michele, David K Lubensky, and Pieter Rein Ten Wolde (2018). "Optimal entrainment of circadian clocks in the presence of noise". In: *Physical Review E* 97.3, p. 032405.

Mossman, Kaspar D et al. (2005). "Altered TCR signaling from geometrically repatterned immunological synapses". In: *Science* 310.5751, pp. 1191–1193.

Mugler, Andrew, Andre Levchenko, and Ilya Nemenman (2016). "Limits to the precision of gradient sensing with spatial communication and temporal integration". In: *Proceedings of the National Academy of Sciences* 113.6, E689–E695.

Nieto, Marta et al. (1997). "Polarization of chemokine receptors to the leading edge during lymphocyte chemotaxis". In: *The Journal of Experimental Medicine* 186.1, pp. 153–158.

Nunan, Naoise et al. (2001). "Quantification of the in situ distribution of soil bacteria by large-scale imaging of thin sections of undisturbed soil". In: *FEMS Microbiology Ecology* 37.1, pp. 67–77.

Pelletier, Anthony J et al. (2000). "Presentation of chemokine SDF-1$\alpha$ by fibronectin mediates directed migration of T cells". In: *Blood, The Journal of the American Society of Hematology* 96.8, pp. 2682–2690.

Petkova, Mariela D et al. (2019). "Optimal decoding of cellular identities in a genetic network". In: *Cell* 176.4, pp. 844–855.

Philipsborn, Anne C von et al. (2006). "Growth cone navigation in substrate-bound ephrin gradients". In: *Development* 133.13, pp. 2487–2495.

Pignata, Aurora et al. (2019). "A spatiotemporal sequence of sensitization to slits and semaphorins orchestrates commissural axon navigation". In: *Cell Reports* 29.2, pp. 347–362.

Pippig, Susanne, Sabine Andexinger, and Martin J Lohse (1995). "Sequestration and recycling of beta 2-adrenergic receptors permit receptor resensitization". In: *Molecular Pharmacology* 47.4, pp. 666–676.

Pytowski, Bronislaw, Frederick R Maxfield, and Josef Michl (1990). "Fc and C3bi receptors and the differentiation antigen BH2-Ag are randomly distributed in the plasma membrane of locomoting neutrophils". In: *The Journal of Cell Biology* 110.3, pp. 661–668.

Qing, Rui et al. (2020). "Non-full-length Water-Soluble CXCR4QTY and CCR5QTY Chemokine Receptors: Implication for Overlooked Truncated but Functional Membrane Receptors". In: *Iscience* 23.12, p. 101670.

Rabin, Julien, Julie Delon, and Yann Gousseau (2011). "Transportation distances on the circle". In: *Journal of Mathematical Imaging and Vision* 41.1, pp. 147–167.

Rangamani, Padmini et al. (2013). "Decoding information in cell shape". In: *Cell* 154.6, pp. 1356–1369.

Raynaud, Xavier and Naoise Nunan (2014). "Spatial ecology of bacteria at the microscale in soil". In: *PloS One* 9.1, e87217.

Rejniak, Katarzyna Anna et al. (2013). "The role of tumor tissue architecture in treatment penetration and efficacy: an integrative study". In: *Frontiers in Oncology* 3, p. 111.

Ren, Aixia et al. (2014). "Asymmetrical macromolecular complex formation of lysophosphatidic acid receptor 2 (LPA2) mediates gradient sensing in fibroblasts". In: *Journal of Biological Chemistry* 289.52, pp. 35757–35769.

Russo, Erica et al. (2016). "Intralymphatic CCL21 promotes tissue egress of dendritic cells through afferent lymphatic vessels". In: *Cell Reports* 14.7, pp. 1723–1734.

Sarris, Milka et al. (2012). "Inflammatory chemokines direct and restrict leukocyte migration within live tissues as glycan-bound gradients". In: *Current Biology* 22.24, pp. 2375–2382.

Schlather, Martin et al. (2020). *RandomFields: Simulation and Analysis of Random Fields*. R package version 3.3.8. URL: `https://cran.r-project.org/package=RandomFields`.

Schlüter, Daniela K, Ignacio Ramis-Conde, and Mark AJ Chaplain (2012). "Computational modeling of single-cell migration: the leading role of extracellular matrix fibers". In: *Biophysical Journal* 103.6, pp. 1141–1151.

Seifert, Roland and Katharina Wenzel-Seifert (2002). "Constitutive activity of G-protein-coupled receptors: cause of disease and common property of wild-type receptors". In: *Naunyn-Schmiedeberg's Archives of Pharmacology* 366.5, pp. 381–416.

Servant, Guy et al. (1999). "Dynamics of a chemoattractant receptor in living neutrophils during chemotaxis". In: *Molecular Biology of the Cell* 10.4, pp. 1163–1178.

Sheng, Hongjie et al. (2017). "Determination of N-acyl homoserine lactones in soil using accelerated solvent extraction combined with solid-phase extraction and gas chromatography-mass spectrometry". In: *Analytical Methods* 9.4, pp. 688–696.

Shields, Jacqueline D et al. (2007). "Autologous chemotaxis as a mechanism of tumor cell homing to lymphatics via interstitial flow and autocrine CCR7 signaling". In: *Cancer Cell* 11.6, pp. 526–538.

Shimonaka, Mika et al. (2003). "Rap1 translates chemokine signals to integrin activation, cell polarization, and motility across vascular endothelium under flow". In: *The Journal of Cell Biology* 161.2, pp. 417–427.

Slack, RJ and DA Hall (2012). "Development of operational models of receptor activation including constitutive receptor activity and their use to determine the efficacy of the chemokine CCL17 at the CC chemokine receptor CCR4". In: *British Journal of Pharmacology* 166.6, pp. 1774–1792.

Sokolowski, Thomas R and Gašper Tkačik (2015). "Optimizing information flow in small genetic networks. IV. Spatial coupling". In: *Physical Review E* 91.6, p. 062710.

Suchkov, Dmitry V et al. (2010). "Polarization of the yeast pheromone receptor requires its internalization but not actin-dependent secretion". In: *Molecular Biology of the Cell* 21.10, pp. 1737–1752.

Swartz, Melody A and Mark E Fleury (2007). "Interstitial flow and its effects in soft tissues". In: *Annu. Rev. Biomed. Eng.* 9, pp. 229–256.

Thoma, Bettina et al. (1990). "Identification of a 60-kD tumor necrosis factor (TNF) receptor as the major signal transducing component in TNF responses". In: *The Journal of Experimental Medicine* 172.4, pp. 1019–1023.

Tkačik, Gašper, Curtis G Callan, and William Bialek (2008). "Information flow and optimization in transcriptional regulation". In: *Proceedings of the National Academy of Sciences* 105.34, pp. 12265–12270.

Tkačik, Gašper and Thomas Gregor (2021). "The many bits of positional information". In: *Development* 148.2, dev176065.

Tkačik, Gašper, Aleksandra M Walczak, and William Bialek (2009). "Optimizing information flow in small genetic networks". In: *Physical Review E* 80.3, p. 031920.

Tripathi, Rarnesh C et al. (1993). "Trabecular cells express receptors that bind TGF-beta 1 and TGF-beta 2: a qualitative and quantitative characterization". In: *Investigative Ophthalmology & Visual Science* 34.1, pp. 260–263.

Tucker, Ronald F et al. (1984). "Specific binding to cultured cells of 125I-labeled type beta transforming growth factor from human platelets". In: *Proceedings of the National Academy of Sciences* 81.21, pp. 6757–6761.

Ueda, Masahiro et al. (2001). "Single-molecule analysis of chemotactic signaling in Dictyostelium cells". In: *Science* 294.5543, pp. 864–867.

Vergassola, Massimo, Emmanuel Villermaux, and Boris I Shraiman (2007). "'Infotaxis' as a strategy for searching without gradients". In: *Nature* 445.7126, pp. 406–409.

Vicente-Manzanares, Miguel and Francisco Sánchez-Madrid (2004). "Role of the cytoskeleton during leukocyte responses". In: *Nature Reviews Immunology* 4.2, pp. 110–122.

Wang, Ya-Juan and Jared Renton Leadbetter (2005). "Rapid acyl-homoserine lactone quorum signal biodegradation in diverse soils". In: *Applied and Environmental Microbiology* 71.3, pp. 1291–1299.

Wang, Liming et al. (2014). "A Bregman matrix and the gradient of mutual information for vector Poisson and Gaussian channels". In: *IEEE Transactions on Information Theory* 60.5, pp. 2611–2629.

Wang, Xiangdan et al. (2008). "Multiplexed cytokine detection of interstitial fluid collected from polymeric hollow tube implants—a feasibility study". In: *Cytokine* 43.1, pp. 15–19.

Wang, Xin et al. (2019). "Mating yeast cells use an intrinsic polarity site to assemble a pheromone-gradient tracking machine". In: *Journal of Cell Biology* 218.11, pp. 3730–3752.

Wang, Yana and Darrell J Irvine (2013). "Convolution of chemoattractant secretion rate, source density, and receptor desensitization direct diverse migration patterns in leukocytes". In: *Integrative Biology* 5.3, pp. 481–494.

Weber, Michele et al. (2013). "Interstitial dendritic cell guidance by haptotactic chemokine gradients". In: *Science* 339.6117, pp. 328–332.

Xiao, Tong and Herwig Baier (2007). "Lamina-specific axonal projections in the zebrafish tectum require the type IV collagen Dragnet". In: *Nature Neuroscience* 10.12, pp. 1529–1537.

Xiao, Tong, Wendy Staub, et al. (2011). "Assembly of lamina-specific neuronal connections by slit bound to type IV collagen". In: *Cell* 146.1, pp. 164–176.

Yang, Bo-Gie et al. (2007). "Binding of lymphoid chemokines to collagen IV that accumulates in the basal lamina of high endothelial venules: its implications in lymphocyte trafficking". In: *The Journal of Immunology* 179.7, pp. 4376–4382.

Yokosuka, Tadashi et al. (2005). "Newly generated T-cell receptor microclusters initiate and sustain T-cell activation by recruitment of Zap70 and SLP-76". In: *Nature Immunology* 6.12, pp. 1253–1262.

Yoshimura, Teizo (2018). "The chemokine MCP-1 (CCL2) in the host interaction with cancer: a foe or ally?" In: *Cellular & Molecular Immunology* 15.4, pp. 335–345.

Zhu, Meng et al. (2020). "Developmental clock and mechanism of de novo polarization of the mouse embryo". In: *Science* 370.6522.

*C h a p t e r   3*

# MORPHEUS: GENERATING COUNTERFACTUAL EXPLANATIONS OF TUMOR SPATIAL PROTEOMES TO DISCOVER THERAPEUTIC STRATEGIES FOR ENHANCING IMMUNE INFILTRATION

## 3.1 Abstract

Immunotherapies can halt or slow down cancer progression by activating either endogenous or engineered T-cells to detect and kill cancer cells. For immunotherapies to be effective, T-cells must be able to infiltrate the tumor microenvironment. However, many solid tumors resist T-cell infiltration, challenging the efficacy of current therapies. Here, we introduce Morpheus, an integrated deep learning framework that takes large-scale spatial omics profiles of patient tumors, and combines a formulation of T-cell infiltration prediction as a self-supervised machine learning problem with a counterfactual optimization strategy to generate minimal tumor perturbations predicted to boost T-cell infiltration. We applied our framework to 368 metastatic melanoma and colorectal cancer (with liver metastases) samples assayed using 40-plex imaging mass cytometry (IMC), discovering cohort-dependent, combinatorial perturbations, involving CXCL9, CXCL10, CCL22 and CCL18 for melanoma and CXCR4, PD-1, PD-L1 and CYR61 for colorectal cancer, predicted to support T-cell infiltration across large patient cohorts. Our work presents a paradigm for counterfactual-based prediction and design of cancer therapeutics using spatial omics data.

## 3.2 Introduction

The immune composition of the tumor microenvironment (TME) plays a crucial role in determining patient prognosis and response to cancer immunotherapies (Fridman et al., 2017; Binnewies, Mikhail, et al., 2018; Bruni, Angell, and Galon, 2020). Immunotherapies that alter the immune composition using transplanted or engineered immune cells (chimeric antigen receptor T-cell therapy) or remove immunosup-

pressive signaling (checkpoint inhibitors) have shown exciting results in relapsed and refractory tumors in hematological cancers and some solid tumors. However, effective therapeutic strategies for most solid tumors remain limited (Hegde and D. S. Chen, 2020; Choe, Williams, and Lim, 2020; Pitt et al., 2016). The TME is a complex mixture of immune cells, including T-cells, B cells, natural killer cells, and macrophages, as well as stromal cells and tumor cells (Fridman et al., 2017). The interactions between these cells can either promote or suppress tumor growth and progression, and ultimately impact patient outcomes. For example, high levels of TILs in the TME are associated with improved prognosis and response to immunotherapy across multiple cancer types (Haslam and Prasad, 2019; Lee and Ruppin, 2019). Conversely, an immunosuppressive TME characterized by low levels of TILs is associated with poor prognosis and reduced response to immunotherapy (Pittet, Michielin, and Migliorini, 2022). Durable, long-term clinical response of T-cell-based immunotherapies are often constrained by a lack of T-cell infiltration into the tumor, as seen in classically "cold" tumors such as triple-negative breast cancer or pancreatic cancer, which have seen little benefit from immunotherapy (Bonaventura, Paola, et al., 2019; Savas et al., 2016; Tsaur et al., 2021). The precise cellular and molecular factors that limit T-cell infiltration into tumors is an open question.

Spatial omics technologies capture the spatial organization of cells and molecular signals in intact human tumors with unprecedented molecular detail, revealing the relationship between localization of different cell types and tens to thousands of molecular signals (Moffitt, Lundberg, and Heyn, 2022). T-cell infiltration is modulated by a rich array of signals within the tumor microenvironment (TME) such as chemokines, adhesion molecules, tumor antigens, immune checkpoints, and their cognate receptors (Lanitis et al., 2017). Recent advances in *in situ* molecular profiling techniques, including spatial transcriptomic (Rodriques et al., 2019; Eng et al., 2019) and proteomic (Giesen et al., 2014; Goltsev et al., 2018) methods, simultaneously capture the spatial relationship of tens to thousands of molecular signals and T-cell localization in intact human tumors with micron-scale resolution. IMC is one such technology that uses metal-labeled antibodies to enable simultaneous detection of up to 40 antigens and transcripts in intact tissue (Giesen et al., 2014).

Recent work on computational methods as applied to multiplexed tumor images have primarily focused on predicting patient-level phenotypes such as survival, by identifying spatial motifs from tumor microenvironments (Bhate, Salil S, et al.,

2022; Z. Wu et al., 2022; Schürch et al., 2020; Aoki, Tomohiro, et al., 2023). These methods have generated valuable insights into how the complex composition of TMEs influences patient prognosis and treatment response, but they fall short of generating concrete, testable hypotheses for therapeutic interventions that may improve patient outcomes. Given the prognostic significance of T-cell infiltration into tumors, we need computational tools that can predict immune cell localization from environmental signals and systematically generate specific, feasible tumor perturbations that are predicted to alter the TME to improve patient outcomes.

Counterfactual explanations (CFEs) can provide important insight in image analysis applications (Chang et al., 2019), but have not been applied to multiplexed imaging data. Traditionally, CFEs help clarify machine learning model decisions by exploring hypothetical scenarios, showing how the model's interpretation would change if a feature in an image were altered slightly (Wachter, Mittelstadt, and Russell, 2017). For instance, slight pixel intensity variations or minor edge alterations in a tumor's appearance on an X-ray might lead a diagnostic model to classify the scan differently. Numerous CFE algorithms exist to elucidate a model's decision boundaries and shed light on its sensitivity to specific image features (Verma et al., 2020). In multiplexed tissue images where each pixel captures detailed molecular information, variations in pixel intensity directly correspond to specific molecular interventions. Thus, spatial omics data enables the extension of CFEs from understanding to predicting actionable interventions.

In this work, we introduce Morpheus, an integrated deep learning framework that first leverages large scale spatial omics profiles of patient tumors to formulate T-cell infiltration prediction as a self-supervised machine learning (ML) problem, and combines this prediction task with counterfactual optimization to propose tumor perturbations that are predicted to boost T-cell infiltration. Specifically, we train a convolutional neural network to predict T-cell infiltration using spatial maps of the TME provided by IMC. We then apply a gradient-based counterfactual generation strategy to the infiltration neural network to compute changes to the signaling molecule levels that increase predicted T-cell abundance. We apply Morpheus to melanoma (Hoch et al., 2022) and colorectal cancer (CRC) with liver metastases (Zhijun Wang et al., 2023) to discover tumor perturbations that are predicted to support T-cell infiltration in tens to hundreds of patients. We provide further validation of ML-based T-cell infiltration prediction using an additional breast cancer data set (Danenberg et al., 2022). For patients with melanoma, Morpheus predicts

combinatorial perturbation to the CXCL9, CXCL10, CCL22 and CCL18 levels can convert immune-excluded tumors to immune-inflamed in a cohort of 69 patients. For CRC liver metastasis, Morpheus discovered two cohort-dependent therapeutic strategies consisting of blocking different subsets of CXCR4, PD-1, PD-L1 and CYR61 that are predicted to improve T-cell infiltration in a cohort of 30 patients. Our work provides a paradigm for counterfactual-based prediction and design of cancer therapeutics based on classification of immune system activity in spatial omics data.

## 3.3 Results

**Counterfactual optimization for therapeutic prediction**

The general logic of Morpheus (Figure 3.1A) is to first train, in a self-supervised manner, a classifier to predict the presence of CD8+ T-cells from multiplexed tissue images (Figure 3.1B). Then we compute counterfactual instances of the data by performing gradient descent on the input image, allowing us to discover perturbations to the tumor image that increases the classifier's predicted likelihood of CD8+ T-cells being present (Figure 3.1C). The altered image represents a perturbation of the TME predicted to improve T-cell infiltration. We mask CD8+ T-cells from all images to prevent the classifier from simply memorizing T-cell expression patterns, guiding it instead to learn environmental features indicative of T-cell presence.

We leverage IMC profiles of human tumors to train a classifier to predict the spatial distribution of CD8+ T-cell in a self-supervised manner. Consider a set of images $\{I^{(i)}\}$, obtained by dividing IMC profiles of tumor sections into local patches of tissue signaling environments, where $I^{(i)} \in \mathbb{R}^{l \times w \times c}$ is an array with $l$ and $w$ denoting the pixel length and width of the image and $c$ denoting the number of molecular channels in the images (Figure 3.1B). Each image shows the level of $c$ proteins across all cells within a small patch of tissue. From patch $I^{(i)}$, we obtain a binary label $s^{(i)}$ indicating the presence and absence of CD8+ T-cells in the patch and a masked copy $x^{(i)}$ with all signals originating from CD8+ T-cells removed (see Supplemental methods). The task for the model $f$ is to classify whether T-cells are present ($s^{(i)} = 1$) or absent ($s^{(i)} = 0$) in image $I^{(i)}$ using only its masked copy $x^{(i)}$. Specifically, $f(x^{(i)}) \in [0, 1]$ is the predicted probability of T-cells, and then we apply a classification threshold $p$ to convert this probability to a predicted label $\hat{s}^{(i)} \in \{0, 1\}$. Since we obtain the image label $s^{(i)}$ from the image $I^{(i)}$ itself by unsupervised clustering of individual cells, our overall task is inherently self-supervised.

Figure 3.1: **An integrated counterfactual optimization framework for discovering therapeutic strategies predicted to drive CD8+ T-cell infiltration in human tumors.** (A) Overview of the Morpheus framework, which consists of first (B) training a neural network classifier to predict the presence of CD8+ T-cells from multiplexed tissue images where CD8+ T-cells are masked. (C) The trained classifier is then used to compute an optimal perturbation vector $\delta^{(i)}$ per patch by jointly minimizing three loss terms ($L_{\text{pred}}$, $L_{\text{dist}}$, $L_{\text{proto}}$). The perturbation $\delta^{(i)}$ represents a strategy for altering the level of a small number of signaling molecules in patch $x_0^{(i)}$ in a way that increases the probability of T-cell presence as predicted by the classifier. The optimization also favors perturbations that shift the image patch to be more similar to its nearest T-cell patches in the training data, shown as proto. Each perturbation corresponds to adjusting the relative intensity of each imaging channel. Taking the median across all perturbations produces a whole-tumor perturbation strategy, which we assess by perturbing *in silico* tumor images from a test patient cohort and examining the predicted T-cell distribution after perturbation.

Given a set of image patches, we train a model $f$ to minimize the following T-cell prediction loss, also known as the binary cross entropy (BCE) loss,

$$L = -\frac{1}{N} \sum_{i=1}^{N} \left[ s^{(i)} \log\left(\hat{s}^{(i)}\right) + \left(1 - s^{(i)}\right) \log\left(1 - \hat{s}^{(i)}\right) \right], \tag{3.1}$$

where

$$\hat{s}^{(i)} = \begin{cases} 1 & \text{if } f(x^{(i)}) \geq p \\ 0 & \text{if } f(x^{(i)}) < p \end{cases} \tag{3.2}$$

and $p$ is the classification threshold. We select $p$ by minimizing the following root mean squared error (RMSE) on a separate set of tissue sections $\Omega$,

$$\text{RMSE}^2 = \frac{1}{|\Omega|} \sum_{j \in \Omega} \left| \frac{1}{N_j} \sum_{i=1}^{N_j} s^{(i)} - \frac{1}{N_j} \sum_{i=1}^{N_j} \hat{s}^{(i)} \right|^2. \tag{3.3}$$

The RMSE is a measure of the differences between the observed and predicted proportions of T-cell patches in a tissue section averaged across a set of tissues $\Omega$, which we take to be the validation set.

We evaluated the performance of various classifiers, including both traditional convolutional neural networks (CNNs) and vision transformers. In all cases, we observed similar performance (Table 3.6). We settled on a U-Net architecture because of ease of extension of the model to multichannel data sets. Our U-Net classifier consists of a standard U-Net architecture (Buda, Saha, and Mazurowski, 2019) and a fully connected layer with softmax activation (Supplemental methods). To increase the number of samples available for training, we take advantage of the spatial heterogeneity of TMEs and divide each tissue image into 48 µm × 48 µm patches upon which the classifier is trained to predict T-cell presence (Supplemental methods).

Using our trained classifier and IMC images of tumors, we employ a counterfactual optimization method to predict tumor perturbations that enhance CD8+ T-cell infiltration (Figure 3.1C). For each image patch $x_0^{(i)}$ that does not contain CD8+ T-cells, our optimization algorithm searches for a perturbation $\delta^{(i)}$ such that our classifier $f$ predicts the perturbed patch $x_p^{(i)} = x_0^{(i)} + \delta^{(i)}$ as having T-cells, hence $x_p^{(i)}$ is referred to as a counterfactual instance. Ideally, we want each perturbation to involve perturbing as few molecules as possible, and realistic in that the counterfactual instance is not far from image patches in our training data so we can be more confident of the model's prediction. We can obtain a perturbation $\delta^{(i)}$ with these desired properties

by solving the following optimization problem adopted from (Looveren and Klaise, 2021),

$$\delta^{(i)} = \min_{\delta} L_{\text{pred}}(x_0^{(i)}, \delta) + L_{\text{dist}}(\delta) + L_{\text{proto}}(x_0^{(i)}, \delta), \tag{3.4}$$

such that

$$\begin{aligned} L_{\text{pred}}(x_0^{(i)}, \delta) &= c \max(-f(x_0^{(i)} + \delta), -p), \\ L_{\text{dist}}(\delta) &= \beta\|\delta\|_1 + \|\delta\|_2^2, \\ L_{\text{proto}}(x_0^{(i)}, \delta) &= \theta\|x_0^{(i)} + \delta - \text{proto}^{(i)}\|_2^2 \end{aligned} \tag{3.5}$$

where $\delta^{(i)}$ is a 3D tensor that describes perturbation made to each pixel of the patch.

The three loss terms in Equation (3.4) each correspond to a desirable property of the perturbation we aim to discover. The term $L_{\text{pred}}$ encourages validity, in that the perturbation increases the classifier's predicted probability of T-cells to be larger than $p$, so the network will predict the perturbed tissue patch as having T-cells when it previously did not contain T-cells. Next, the term $L_{\text{dist}}$ encourages sparsity, in that the perturbation does not require making many changes to the TME, by minimizing the distance between the original patch $x_0^{(i)}$ and the perturbed patch $x_p^{(i)} = x_0^{(i)} + \delta$ using elastic net regularization. Lastly, the term $\text{proto}^{(i)}$ in the expression for $L_{\text{proto}}$ refers to the nearest neighbor of $x_0^{(i)}$ among all patches in the training set that are classified as having T-cells (see Supplemental methods). Thus the term $L_{\text{proto}}$ explicitly guides the perturbed image $x_p^{(i)}$ to lie close to the data manifold defined by our training set, making perturbed patches appear similar to what has been observed in TMEs infiltrated by T-cells.

Since drug treatments cannot act at the spatial resolution of individual micron-scale pixels, we constrain our search space to only perturbations that affect all cells in the image uniformly. Specifically, we only search for perturbations that change the level of any molecule by the same relative amount across all cells in an image. We incorporate this constraint by defining $\delta^{(i)}$ in the following way,

$$\delta^{(i)} = \gamma^{(i)} \odot_3 x_0^{(i)}, \tag{3.6}$$

where $\gamma^{(i)} \in \mathbb{R}^c$ defines a single factor for each channel in the image and the circled dot operator represent channel-wise multiplication, so that within each channel, the scaling factor is constant across the spatial dimensions of the image. In practice, we directly optimize for $\gamma^{(i)}$, where $\gamma_j^{(i)}$ can be interpreted as the relative change to the mean intensity of the $j$-th channel. However, given our classifier does have fine spatial resolution, we can search for targeted therapies such as perturbing only

a specific cell type or restricting the perturbation to specific tissue locations by changing Equation (3.6) to match these different types of perturbation.

Taken together, our algorithm obtains an altered image predicted to contain T-cells from an original image which lacks T-cells, by minimally perturbing the original image in the direction of the nearest training patch containing T-cells until the classifier predicts the perturbed image to contain T-cells. Since our strategy may find different perturbations for different tumor patches, we reduce the set of patch-wise perturbations $\{\delta^{(i)}\}_i$ to a whole-tumor perturbation by taking the median across the entire set.

**Convolutional neural networks predict T-cell distribution**



Figure 3.2: **U-Net classifiers accurately predict T-cell distribution in IMC images of melanoma, metastatic liver, and breast tumor.** (A) Histograms showing the distribution of tumor cores per patient and CD8+ T-cell fractions per core across all three data sets and data splits. (B) Predicted and actual T-cell distribution of tissue sections from test cohorts in melanoma, liver tumor, and breast tumor data set. (C) Predicted and true proportion of patches with T-cells within a tissue section, each dot corresponds to a tissue section, diagonal black line indicates perfect prediction. (D) The RMSE (Equation (3.3)) across all (test) tissue sections for three different classes of models.

We applied Morpheus to two publicly available IMC data sets of tumors from patients with metastatic melanoma (Hoch et al., 2022) and colorectal cancer (CRC) with liver metastases (Zhijun Wang et al., 2023) (Figure 3.2A). We validate the infiltration prediction on an additional breast cancer data set (Danenberg et al.,

2022). While this breast cancer data focuses on cell-type markers over functional modulators of T-cell infiltration, making it unsuitable for therapeutic prediction, it serves to further validate our ML-based prediction of T-cell infiltration.

The melanoma data set (Hoch et al., 2022) was obtained by IMC imaging of 159 tumor cores from 69 patients with stage III or IV metastatic melanoma. Each tissue was imaged across 39 molecular channels, consisting of markers for tumor, immune, and stromal cells, as well as 11 different chemokines (RNA) (Supplemental methods). The CRC data set (Zhijun Wang et al., 2023) consists of 209 tissue sections taken from 30 patients imaged across 42 channels, including 60 sections from primary CRC tumors, 89 sections CRC metastases to the liver and 60 "healthy" liver sections obtained away from the metastases (Supplemental methods). The breast cancer data set (Danenberg et al., 2022) was obtained by IMC imaging of 749 breast tumor cores from 693 patients. The tissues were imaged across 37 channels, consisting of markers for tumor, lymphoid, myeloid and stromal cells (Supplemental methods).

For each of the three tumor data sets, we trained a separate U-Net classifier that effectively predicts CD8+ T-cell infiltration level in unseen tumor sections (Supplemental methods). The two classifiers trained on melanoma and CRC data sets achieved the best performance with an AUROC of 0.77 and 0.8, respectively, whereas the classifier trained on breast tumors achieved a AUROC of 0.71 (Table 3.5). Figure 3.2B shows examples of actual and predicted T-cell distributions in tumor sections, demonstrating that our classifiers accurately predict the general distribution of T-cells. For each tissue section of a cancer type, the predictions were obtained by applying the corresponding U-Net classifier to each image patch independently. Comparing the true proportion of T-cell patches in a tissue section against our predicted proportion also shows strong agreement (Figure 3.2C). The true proportion of patches with T-cells is calculated by dividing the number of patches within a tissue section that contain CD8+ T-cells by the total number of patches within that section. We quantify the performance of our U-Nets on the entire test data set using the RMSE (Equation 3.3), which represents the mean difference between our predicted proportion and the true proportion per tumor section (Figure 3.2D). Our classifiers performs well on liver tumor and melanoma, achieving a RMSE of only 6% and 8%, respectively, and a relatively lower performance of 11% on breast tumor. Taken together, these results suggest that our classifier can accurately predict the T-cell infiltration status of multiple tumor types.

In order to gain insight into the relative importance of non-linearity and spatial information in the performance of the U-Net on the T-cell classification task, we compared the U-nets' performance to a logistic regression model (LR) and a multi-layer perceptron (MLP). Both the LR and MLP model are given only mean channel intensities as input, so neither have explicit spatial information. Furthermore, the LR model is a linear model with a threshold whereas the MLP is a non-linear model. Figure 3.2D shows that across all three cancer data sets, the MLP classifier consistently outperforms the LR model, reducing RMSE by $20-40\%$ to suggest that there are significant nonlinear interactions between different molecular features in terms of their effect on T-cell localization. The importance of spatial features on the T-cell prediction task, however, is less consistent across cancer types. Figure 3.2D shows that for predicting T-cells in breast tumor, the U-Net model offers negligible boost in performance relative to the MLP model ($< 2\%$ RMSE reduction), whereas for liver tumor, the U-Net model achieved a RMSE $50\%$ lower compared to the MLP model. This result suggests that the spatial organization of signals may have a stronger influence on CD8+ T-cell localization in liver tumor compared to breast tumor.

# Applying Morpheus to metastatic melanoma samples



Figure 3.3: **Combinatorial chemokine therapy predicted to drive T-cell infiltration in patients with metastatic melanoma** (A) Whole-tumor perturbations optimized across IMC images of patients (row) from the training cohort, with bar graph showing the median relative change in intensity for each molecule.

Figure 3.3: **(continued)** (B) Distribution of cancer stages among patients within two clusters, gray indicates unknown stage, chance probability from hypergeometric distribution. (C) Volcano plot comparing chemokine level and cell-type abundance from patient cluster 1 and 2, computed using mean values and Wilcoxon rank sum test. Gray indicates non-statistical significance. (D) Patch-wise chemokine profile (left); 1-D heatmap (right): infiltration status (light/dark = from infiltrated/deserted tumor), tumor cell (light/dark = present/absent), CD8+ T-cells (light/dark = present/absent). (E) Patch-wise correlation between chemokine signals and the presence of CD8+ T-cells. (F) (Top) UMAP projection of tumor patches (chemokine channels) show a clear separation of masked patches with and without T-cells. (Bottom) colored arrows connect UMAP projection of patches without T-cells and their corresponding counterfactual (perturbed) patch, where the colors correspond to k-nearest neighbor clusters (i-iv) of the counterfactual patches, highlighting the minimal nature of the perturbations. Pie charts (i-iv) shows the distribution of patients whose original tumor patches are found in the corresponding cluster regions in the UMAP. (G) Cell maps computed from a patient's IMC image, showing the distribution of T-cells before and after perturbation. (H) Original vs. perturbed (predicted) mean infiltration level across all patients (test cohort) with 95% confidence interval (only shown for patients with more than two samples). Stage IV patients received perturbation strategy 1 (yellow), stage III patients received perturbation strategy 2 (green). (I) Mean infiltration level across all patients (test cohort) for optimized perturbation strategies of varying sparsity, error bar represents 95% CI.

Applying our counterfactual optimization procedure using the U-Net classifier trained on melanoma IMC images, we discovered a combinatorial therapy predicted to be highly effective in improving T-cell infiltration in patients with melanoma. We restricted the optimization algorithm to only perturb the level of chemokines, which are a family of secreted proteins that are known for their ability to stimulate cell migration (Hughes and Nibbs, 2018) and have already been harnessed to augment T-cell therapy (Foeng, Comerford, and McColl, 2022). By optimizing over multiple chemokines, Morpheus opens the door to combinatorial chemokine therapeutics that has the potential to more effectively enhance T-cell infiltration into tumors. Figure 3.3A shows that patients from the training cohort separate into two clusters based on hierarchical clustering of perturbations computed for each patient. Taking median across all patients in cluster 1, the optimized perturbation is to increase CXCL9 level by 370%, whereas in patient cluster 2, the optimized perturbation consists of increasing CXCL10 level by 280% while decreasing CCL18 and CCL22 levels by 100% and 70%, respectively (Figure 3.3A). Both CXCL9 and CXCL10 are well known for playing a role in the recruitment of CD8+ T-cells to tumors. On

the other hand, CCL22 is known to be a key chemokine for recruiting regulatory T-cells (Kohli, Pillarisetty, and T. S. Kim, 2022) and CCL18 is known to induce an M2-macrophage phenotype (Schraufstatter et al., 2012), so their expression likely promotes an immunosuppressive microenvironment inhibitory to T-cell infiltration and function.

Figure 3.3B shows that the choice of which of these two strategies was selected for a patient appears to be strongly associated with the patient's cancer stage, with strategy 1 being significantly enriched for patients with stage IV metastatic melanoma and strategy 2 being significantly enriched for patients with stage III cancer, with a probability of 0.053 of such difference being due to chance. Probing deeper into the difference between these two patient clusters, we find that all chemokines have lower mean expression in the tumors of patients in cluster 1 compared to cluster 2, while there are no significant differences between the two groups in terms of the cell-type compositions within tumors (Figure 3.3C). Since the levels of CCL22 and CCL18 is 37% and 31% higher in patients from cluster 2 and both chemokines have been implicated in having an inhibitory effect on T-cell infiltration, it is reasonable that the optimization algorithm suggests inhibiting CCL18 and CCL22 only for patients in cluster 2. However, the switch from boosting CXCL9 to CXCL10 is not as straightforward. A possible explanations is that boosting CXCL10 is important when blocking CCL18 and CCL22 in order for the perturbed patches to stay close to the data manifold, leading to more realistic tissue environments.

Morpheus selected perturbations that would make the chemokine composition of a TME more similar to T-cell rich regions of immune-infiltrated tumors. Figure 3.3D shows that melanoma tissue patches can be clustered into distinct groups based on their chemokine concentration profile. One cluster (highlighted in blue) contains exactly the patches from immune-infiltrated tumors that contain both tumor and T-cells, which likely represents a chemokine signature that is suitable for T-cell infiltration. Alternately, a second cluster (highlighted in red) which contains patches from immune-desert tumors that have tumor cells but no T-cells likely represents an unfavorable chemokine signature. In comparison to the cluster highlighted in red, Figure 3.3D shows the cluster highlighted in blue contains elevated levels of CXCL9, CXCL10 and reduced levels of CCL22 which partially agrees with the perturbation strategy (Figure 3.3A) discovered by Morpheus. Lastly, Figure 3.3E shows that our four selected chemokine targets cannot simply be predicted from correlation of chemokine levels with the presence of CD8+ T-cells, as both CCL18

and CCL22 are weakly correlated ($< 0.1$) with CD8+ T-cells even though the optimized perturbations requires inhibiting both chemokines, suggesting the presence of significant nonlinear effects not captured by correlations alone.

We can directly observe how Morpheus searches for efficient perturbations by viewing both the original patch and perturbed patches in a dimensionally reduced space. Figure 3.3F (top) shows a UMAP projection where each point represents the chemokine profile of an IMC patch. T-cell patches (with their CD8+ T-cells masked) are well separated from patches without CD8+ T-cells. The colored arrows in the bottom UMAP of Figure 3.3F illustrate the perturbation for each patch as computed by Morpheus, and demonstrate two key features of our algorithm. First, optimized perturbations push patches without T-cells towards the region in UMAP space occupied by T-cell-infiltrated patches. Second, the arrows in Figure 3.3C are colored to show that optimized perturbations seem efficient in that patches are perturbed just far enough to land in the desired region of space. Specifically, red points that start out on the right edge end up closer to the right after perturbation (region ii and iii), while points that start on the left/bottom edge end up closer to the left/bottom (region i and iv), respectively. We make this observation while noting that UMAP, though designed to preserve the topological structure of the data, is not a strictly distance-preserving transformation (McInnes, Healy, and Melville, 2018). Furthermore, the pie charts (i-iv) are colored by the patient of origin to show the region of space where points are being perturbed to are not occupied by tissue samples from a single patient with highly infiltrated tumor. Rather, these regions consist of tissue samples from multiple patients, suggesting that our optimization procedure can synthesize information from different patients when searching for therapeutic strategies.

After applying the second perturbation strategy from Figure 3.3A *in silico* to IMC images of a tumor, Figure 3.3G shows that T-cell infiltration level (defined as the proportion of tumor patches with T-cells) is predicted to increase by 20 fold. We applied our two perturbation strategies on patients in our test cohort *in silico* after stratifying by cancer stage, using strategy 1 on patients with stage IV melanoma and strategy 2 on patients with stage III melanoma.Figure 3.3H shows that this predicted improvement holds across nearly all 14 patients from the test group, boosting T-cell infiltration level from an average of 23% across samples to a predicted 63% post perturbation. For the three test patients with multiple tumor sections (patient 64, 57, 89), we see small to moderate variation in predicted improvement across samples.

The combinatorial nature of our optimized perturbation strategy is crucial to its predicted effectiveness. We systematically explored the importance of combinatorial perturbation by changing parameter $\beta$ of Equation (3.4) which adjusts the sparsity of the strategy, where a more sparse strategy means fewer molecules are perturbed. Figure 3.3I shows that perturbing multiple targets is predicted to be necessary for driving significant T-cell infiltration across multiple patients, with the best perturbation strategy involving two targets predicted to generate only 60% of the infiltration level achieved by the best perturbation strategy involving four targets. In conclusion, within the scope of the chemokine targets considered, combinatorial perturbation of the TME appears necessary for improving T-cell infiltration in metastatic melanoma.

**Applying Morpheus to CRC with liver metastases samples**



Figure 3.4: **Blocking subsets of PD-L1, CXCR4, PD-1, and CYR61 predicted to drive T-cell infiltration in CRC cohort.**

Figure 3.4: **(continued)** (A) Optimized tumor perturbations aggregated to the patient (row) level (train cohort). Bar graph shows the median relative change in intensity for each molecule across all patients within their cluster. (B) Patch-wise correlation between the levels of different molecules and the presence of CD8+ T-cells. (C) Pie charts show proportion of patients in each cluster that have fatty liver disease (FLD), chance probability from hypergeometric distribution. (D) Volcano plot comparing molecule levels and cell-type abundance between the two patient cluster using tumor tissues, computed using mean values and Wilcoxon rank sum test with Bonferroni correction. (E) Optimized perturbations aggregated to the level of tissue samples (row). (F) UMAP projection of IMC patches, left UMAP shows T-cell patches colored by the tissue samples they are taken from. right UMAP shows counterfactual (perturbed) instances optimized for tumor patches without T-cells (red). (G) Line plots shows T-cell infiltration level for each tissue section from the test cohort, before and after perturbation. Bar plots show predicted mean T-cell infiltration level for each test patient. (H) Mean infiltration level across all test patients using perturbation strategies of varying sparsity, obtained by varying $\beta$ in Equation (3.4), error bar represents 95% CI.

Applying Morpheus to IMC images from the CRC cohort, we discovered two patient-dependent therapies predicted to be highly effective in improving T-cell infiltration. Figure 3.4A shows the optimal perturbations computed for every patient from the training cohort, aggregated over all tumor samples for each patient. Our method consistently discovered two distinct patient-dependent strategies for improving T-cell infiltration, as revealed by hierarchical clustering of all patient-level perturbations (Figure 3.4A). Taking median over patients in the first cluster, the optimized strategy involves completely inhibiting PD-1, PD-L1, and CXCR4. While for the second group of patients, the optimized strategy involves completely inhibiting CYR61, PD-1, PD-L1, and CXCR4 (Figure 3.4A). Interestingly, all four of the perturbation targets correlated poorly with the presence of CD8+ T-cells compared to the other proteins that were not selected as perturbation targets (Figure 3.4B), suggesting the presence of significant spatial and nonlinear effects not captured by correlations alone.

All perturbation targets identified by our optimization procedure have been found to play crucial roles in suppressing T-cell function in the TME, and treating patients with inhibitors against subsets of the selected targets have already improved T-cell infiltration in human CRC liver metastases. Regulatory T-cells (Tregs) are recruited into tumor through CXCL12/CXCR4 interaction (Ghanem et al., 2014), and the PD-1/PD-L1 pathway inhibits CD8+ T-cell activity and infiltration in tumors.

In addition, CYR61 is a chemoattractant and was recently shown to drive M2 TAM infiltration in patients with CRC liver metastases (Zhijun Wang et al., 2023). Inhibition of both PD-1 and CXCR4, which were consistently selected by Morpheus as targets, have already been shown to increase CD8+ T-cell infiltration in both patients with CRC and mouse models (Biasci, Daniele, et al., 2020; Y. Chen et al., 2015; Steele et al., 2023). Finally, Figure 3.4A shows that the fifth most common proposed perturbation involves inhibiting IL-10. Indeed, blockade of IL-10 was recently shown to increase the frequency of non-exhausted CD8+ T-cell infiltration in slice cultures of human CRC liver metastases (Sullivan et al., 2023).

The emergence of the two distinct perturbation strategies may be explained by variation in liver fat build-up among patients. Patient cluster 1 is made up of significantly more patients with fatty liver disease (70% FLD) compared to patient cluster 2 (22%), where the probability of this due purely to chance is 0.047 (Figure 3.4C). Furthermore, Figure 3.4D shows that both YAP and CYR61 levels are significantly higher in tumors from patient cluster 1, by 50% and 15%, respectively. Indeed, CYR61 is known to be associated with non-alcoholic fatty liver disease (FLD) (Zhijun Wang et al., 2023) and YAP is a transcription coregulator that induces CYR61 expression (Zhang, Pasolli, and Fuchs, 2011). However despite patients in cluster 1 having higher levels of CYR61, it is only for patients in cluster 2 where the optimal strategy involves blocking CYR61. We postulate that this seemingly paradoxical finding may arise because removing CYR61 from patients in cluster 1 represents a more pronounced perturbation, given their inherently higher concentration. A perturbation of this magnitude would likely shift the tumor profile significantly away from the data manifold, where the classifier's prediction about the perturbation's effect becomes less reliable, hence such a perturbation would be heavily penalized during optimization due to the $L_{\text{proto}}$ term.

Using only raw image patches, Morpheus discovers tissue-dependent perturbation strategies (Figure 3.4E). As depicted in Figure 3.4E, by aggregating perturbations at the individual tissue level, we observe that the optimized perturbation for "healthy" liver sections is straightforward, necessitating only the inhibition of CXCR4. Recall "healthy" sections are samples obtained away from sites of metastasis. In contrast, promoting T-cell infiltration into primary colon tumors is anticipated to involve targeting a minimum of three signals. Our method finds that liver metastases appears to fall between these two tissue types. The optimized perturbation strategy for some liver metastases samples is to block CXCR4, while requiring the inhibition of the

same set of signals as primary tumors for others. Furthermore, direct comparison between perturbations optimized for metastatic tumor and primary tumor samples does not reveal a significant difference in strategy (Figure 3.6). We can partly understand the discrepancy between tissues by plotting a UMAP projection of all T-cell patches from the three tissue types (Figure 3.4F, left). The clear separation between T-cell patches from "healthy" tissue and those from primary tumors underscores that the signaling compositions driving T-cell infiltration likely differ substantially between the two tissue types. This distinction is likely what prompted our method to identify markedly different perturbation strategies. Furthermore, some patches from metastatic tumors co-localize with "healthy" tissue patches in UMAP space, while other patches co-localizes with primary tumor patches. This observation again aligns with our previous result, where optimized perturbation strategies for metastases samples share similarities with strategies for either "healthy" tissue or primary tumor (Figure 3.4E).

Despite the CRC data set comprising a complex blend of healthy, tumor, and hybrid metastatic samples, Morpheus targets the most pertinent tissue type when optimizing perturbations. During both the model training and counterfactual optimization phases, we did not make specific efforts to segregate the three tissue types. Furthermore, we did not provide tissue type labels or any metadata. Despite these nuances, Figure 3.4F shows that the counterfactual instances for tumor patches (dark blue) from primary and metastases samples are mostly perturbed to be near T-cell patches from primary (cyan) and metastatic tumor (gold), instead of being perturbed to be similar to T-cell patches from "healthy" tumors (purple). This result is partly a consequence of our prototypical constraint which encourages patches to be perturbed towards the closest T-cell patch. For a patch from a metastatic tumor without T-cells, the closest (most similar) T-cell patch is likely also from a metastatic tumor than from a "healthy" tissue. However, there are occasional exceptions where T-cell patches from "healthy" tissues can influence the optimization of tumor tissues, as outlined by the dashed ellipse in Figure 3.4F, especially if they share similar features as tumor regions.

The two therapeutic strategies we discovered generalize to patients in our test cohort (Figure 3.4G,H). Given that we have two therapeutic strategies, one enriched for patients with FLD and another for patients without FLD, we apply different perturbation strategies *in silico* across all test patients depending on their FLD status. Aggregated to the patient level, Figure 3.4G shows that CD8+ T-cell infiltration

level is predicted to increase for nearly all patients, with the exception of patient 28. Furthermore, aggregating to the entire test cohort, Figure 3.4H shows a statistically significant boost to mean infiltration level from 15% to a predicted 35% post perturbation. However, when comparing individual tissue samples, Figure 3.4G reveals significant variation in the predicted response to perturbation among samples from the same patient and tissue types. In patient 7, one primary tumor sample is predicted to see a nearly threefold increase in T-cell infiltration after perturbation, yet almost no change is expected for patient 7's other two primary and three metastatic samples. Similar patterns are observed in patients 14 and 17. This marked variability in response among a significant portion of test patients underscores the challenges posed by intra-tumor and inter-patient heterogeneity in devising therapies for CRC with liver metastases. This result further implies that, for studying CRC with liver metastases, collecting numerous tumor sections per patient could be as crucial as establishing a large patient cohort. Lastly, combinatorial perturbation is again predicted to be necessary to drive significant T-cell infiltration in large patient cohorts. By increasing $\beta$ in Equation (3.4), we generated strategies with between one and four total targets, where our four-target perturbation is the only strategy predicted to produce a statistically significant boost to T-cell infiltration (Figure 3.4H).

## 3.4 Discussion

Our integrated deep learning framework, Morpheus, combines deep learning with counterfactual optimization to directly predict therapeutic strategies from spatial omics data. One of the major strengths of Morpheus is that it scales efficiently to deal with large diverse sets of patients samples including metachronous tissue from the same patients but different sites, which will be crucial as more spatial transcriptomics and proteomics data sets are quickly becoming available (A. Chen et al., 2022). Larger data sets could allow us to train more complex models such as vision transformers, capturing long range interactions in tissues to improve prediction of T-cell localization. Furthermore, a large set of diverse patient samples will more accurately capture the extent of tumor heterogeneity, enabling Morpheus to discover therapeutic strategies for different sub-classes of patients.

## 3.5 Code availability

Code for model training, perturbation optimization and analysis are publicly available at `https://github.com/neonine2/morpheus-spatial`. Our optimization code was implemented in Python and was built upon the open source Python

library Alibi (Klaise et al., 2021).

## 3.6 Data availability

All data sets used in this study are published and publicly available.

## 3.7 Acknowledgements

We would like to thank Inna Strazhnik for her support with figure illustrations. We would like to thank Akil Merchant, Alma Andersson, Aviv Regev, Long Cai, Barbara Wold, Michal Polonsky, Jonathan Fox, Yujing Yang, Abdullah Farooq and all members of the Thomson lab for insightful discussion that significantly improved this work. We gratefully acknowledge the support of the National Institutes of Health's Information Technology for Cancer Research (ITCR) program and the Merkin Institute for Translational Research.

## 3.8 Supplemental methods

### IMC data sets

All data sets used in this paper are publicly available. Metastatic melanoma data set from Hoch et al. (Hoch et al., 2022) contains 159 images or cores taken from 69 patients, collected from sites including skin and lymph-node. CRC liver metastases data set from Wang et al. (Zhijun Wang et al., 2023) contains 209 images or cores taken from 30 patients. Breast tumor data set from Danenberg et al. (Danenberg et al., 2022) contains 693 images or cores taken from 693 patients. The RNA and protein panels used for each of the three data sets are listed in Table 3.1.

### Data split

For all three IMC data sets, we followed the same data splitting scheme to divide patients into three different groups (training, validation, testing) while ensuring similar class balance across the groups, which in our case means that the proportion of image patches with and without T-cells are roughly equal across the three groups for each data set. Specifically, each image within a data set was divided into $48\,\mu m \times 48\,\mu m$ patches and the number of patches with and without CD8+ T-cells was computed for each image. Furthermore, each patch was downsampled from $48 \times 48$ pixels to $16 \times 16$ pixel dimension where each pixel now represents a $3\,\mu m \times 3\,\mu m$ region. We applied spectral analysis to study the effect of using different patch size to predict T-cell infiltration and found that our selected patch size remains highly informative of T-cell presence (Figure 3.5). Next, the patients

| Metastatic melanoma | | CRC with liver metastases | | Breast tumor | |
|---|---|---|---|---|---|
| Vimentin | **DapB** | CD45 | Glnsynthetase | Histone H3 | SMA |
| CD163 | **CCL4** | CD163 | NKG2D | CK5 | CD38 |
| B2M | **CCL18** | CCR4 | PD-L1 | HLA-DR | CK8-18 |
| CD134 | **CXCL8** | FAP | CD11c | CD15 | FSP1 |
| CD68 | **CXCL10** | LAG3 | HepPar1 | CD163 | ICOS |
| GLUT1 | **CXCL12** | FOXP3 | $\alpha$SMA | OX40 | CD68 |
| CD3 | **CXCL13** | CD4 | CD105 | HER2 (3B5) | CD3 |
| LAG3 | **CCL2** | CD68 | VISTA | Podoplanin | CD11c |
| PD-1 | **CCL22** | CD20 | CD8$\alpha$ | PD-1 | GITR |
| HistoneH3 | **CXCL9** | TIM3 | CXCR4 | CD16 | c-Caspase3 |
| CCR2 | **CCL19** | PD-1 | iNOS | CD45RA | B2M |
| PD-L1 | **CCL8** | CD31 | CYR61 | CD45RO | FOXP3 |
| CD8 | SMA | CDX2 | CAIX | CD20 | ER |
| SOX10 | CD31 | CD3 | CD44 | CD8 | CD57 |
| Mart1 | pRB | CD15 | CD11b | Ki-67 | PDGFR$\beta$ |
| cleavedPARP | MPO | HLA-DR | IL10 | Caveolin-1 | CD4 |
| CD15 | CK5 | CXCL12 | HLA-ABC | CD31-vWF | CXCL12 |
| CD38 | HLA-DR | GranzymeB | Ki67 | HLA-ABC | panCK |
| S100 | Cadherin11 | HistoneH3 | CXCR3 | HER2 (D8F12) | |
| FAP | | Galectin9 | YAP | | |
| | | CD14 | CK19 | | |

Table 3.1: Protein and RNA panels imaged for each of the IMC data sets, with RNA targets bolded

are shuffled between the three groups until three criteria are met: 1) the number of patients across the three groups follow a 65/15/20 ratio, 2) the difference in class proportion between any two of the three groups is less than 2%, and 3) the training set contains at least 65% of total patches. The actual data splits used in the paper are described in Table 3.2.

| Data set | Group | Patient count | Patch count | Percent patches with CD8+ T-cells |
|---|---|---|---|---|
| Metastatic melanoma | Training | 102 | 23741 | 29.6% |
| | Validation | 28 | 6045 | 30.3% |
| | Testing | 29 | 5950 | 30.4% |
| CRC with liver metastases | Training | 19 | 44449 | 15.9% |
| | Validation | 4 | 6957 | 14.4% |
| | Testing | 7 | 14907 | 15.9% |
| Breast cancer | Training | 485 | 41104 | 23.7% |
| | Validation | 113 | 9015 | 23.4% |
| | Testing | 151 | 12987 | 23.8% |

Table 3.2: Data split for Melanoma, CRC cohort, and breast tumor IMC data set

**Single-cell phenotyping**

For each data set, we used the cell-type classification (tumor and CD8+ T-cells) from the original paper. For the melanoma data set, cell phenotyping was performed using the Shiny application of the R package cytomapper (Eling et al., 2020), which allows labeling of cell populations using multiple gates. CD8+ T-cells were defined using CD3 and CD8, tumor cells are positive for any or multiple of SOX9, SOX10, MITF, Mart1, S100A1, and p75. For the CRC and breast cancer data set, cell-type labeling was performed using PhenoGraph (Levine et al., 2015).

**Classifier training**

In this work, we trained three classes of models to perform our T-cell prediction task. All models presented in this paper were trained with early stopping based on the validation Matthews Correlation Coefficient (MCC) for 10–20 epochs. All models were trained on an NVIDIA GeForce RTX 3090 Ti GPU using PyTorch version 1.13.1 (Paszke et al., 2019). More details about hyperparameters and implementations can be found in our GitHub repository.

**T-cell masking strategy**

The purpose of model training is for the model to learn molecular features of a tissue environment that supports the presence of CD8+ T-cell, so it is important for us to remove features of the image that are predictive of CD8+ T-cell presence but are not part of the cell's environment, for example, the expression profile of T-cells themselves. We devised a non-trivial cell masking strategy in order to remove T-cell expression patterns without introducing new features that are highly predictive of T-cell presence but are not biologically relevant. A simple masking strategy of zeroing out all pixels belonging to CD8+ T-cells will introduce contiguous regions of zeros to image patches with T-cells, which is an artificial feature that is nonetheless highly predictive of T-cell presence and hence will likely be the main feature learned by a model during training. To circumvent this issue, we first apply a cell "pixelation" step to the original IMC image where we reduce each cell to a single pixel positioned at the cell's centroid. The value of this pixel is the sum of all pixels originally associated with the cell, representing the total signal from each channel within the cell. We then mask this "pixelated" image by zeroing all pixels representing CD8+ T-cells. Since there are usually at most two T-cell pixels in an image patch, zeroing them in a $16 \times 16$ pixel image where most ($> 90\%$) of the pixels are already zeros

is not likely to introduce a significant signal that is predictive T-cell presence. We show that our strategy is effective at masking T-cells without introducing additional features through a series of image augmentation experiments ( Supplemental Note 1 Assessment of T-cell masking strategy).

## Logistic regression models

We trained a single-layer neural network on the average intensity values from each molecular channel to obtain a LR classifier, predicting the probability of CD8+ T-cell presence in the image patch. This model represents a linear model where only the average intensity of each molecule is used for prediction instead of their spatial distribution within a patch.

## MLP models

Similar to a LR model, the MLP also uses averaged intensity as input features for prediction but is capable of learning nonlinear interactions between features. The MLP model consists of two hidden layers (30 and 10 nodes) with Rectified Linear Unit (ReLU) activation.

## U-Net models

To train networks that can make full use of the spatial information, we used a fully convolutional neural network with the U-Net architecture. The U-Net architecture consists of a contracting path and an expansive path, which gives it a U-shaped structure (Buda, Saha, and Mazurowski, 2019). The contracting path consists of four repeated blocks, each containing a convolutional layer followed by a ReLU activation and a max pooling layer. The expansive path mirrors the contracting path, where each block contains a transposed convolutional layer. Skip connections concatenates the up-sampled features with the corresponding feature maps from the contracting path to include local information. The output of the expansive path is then fed to a fully connected layer with softmax activation to produce a predicted probability. The model was trained from scratch using image augmentation to prevent over-fitting, including random horizontal/vertical flips and rotations, in addition to standard channel-wise normalization. We train our U-Net classifiers using stochastic gradient descent with momentum and a learning rate of $10^{-2}$ on mini-batches of size 128.

**Counterfactual optimization**

Given an IMC patch $x^{(i)}$ without T-cells, and a classifier $f$, our goal is to find a perturbation $\delta^{(i)}$ for the patch such that $f$ classifies the perturbed patch as having T-cells. For CNN models, $\delta^{(i)} \in \mathbb{R}^{w \times l \times d}$ is a 3D tensor that describes changes made for every channel, at each pixel of the patch.

Given a CNN classifier $f$ and a IMC patch $x^{(i)}$ such that $f(x_0^{(i)}) = \mathbb{P}(\text{T-cells present}) < p$, where $p > 0$ is the classification threshold below which the classifier predicts no T-cell, we aim to obtain a perturbation $\delta^{(i)}$ such that $f(x_0^{(i)} + \delta^{(i)}) > p$, by solving the following optimization problem adopted from (Looveren and Klaise, 2021),

$$\delta^{(i)} = \min_\delta L_{\text{pred}}(x_0^{(i)}, \delta) + L_{\text{dist}}(\delta) + L_{\text{proto}}(x_0^{(i)}, \delta), \tag{3.7}$$

such that

$$L_{\text{pred}}(x_0^{(i)}, \delta) = c \max(-f(x_0^{(i)} + \delta), -p), \tag{3.8}$$

$$L_{\text{dist}}(\delta) = \beta \|\delta\|_1 + \|\delta\|_2^2, \tag{3.9}$$

$$L_{\text{proto}}(x_0^{(i)}, \delta) = \theta \|x_0^{(i)} + \delta - \text{proto}^{(i)}\|_2^2, \tag{3.10}$$

$$\delta^{(i)} = \gamma^{(i)} \odot_3 x_0^{(i)} \tag{3.11}$$

where proto$^{(i)}$ is an instance of the training set classified as having T-cells, defined by first building a k-d tree of training instances classified as having T-cells and setting the $k$-nearest item in the tree (in terms of Euclidean distance to $x_0^{(i)}$) as proto. We use $k = 1$ for all counterfactual optimization. For all other parameters, we list their values in Table 3.3. During optimization, the weight $c$ of the loss term $L_{\text{pred}}$ is updated for $n$ iterations, starting at $c_0$. If we identify a valid counterfactual for the present value of $c$, we will then decrease $c$ in the subsequent optimization cycle to increase the weight of the additional loss components, thereby enhancing the overall solution. If, however, we do not identify a counterfactual, $c$ is increased to put more emphasis on increasing the predicted probability of the counterfactual. The parameter $s_{\text{max}}$ sets the maximum number of optimization steps for each value of $c$.

## 3.9 Supplemental information

**Supplemental Note 1 Assessment of T-cell masking strategy**

By masking T-cells prior to model training, we may have inadvertently introduced a new signal that is predictive of T-cells, specifically that fewer cells in an image increases the probability of T-cells being present. To assess the possibility of such an effect, we study the impact of random cell masking on the predicted probability

| Parameters | Melanoma | CRC |
|:---:|:---:|:---:|
| $\beta$ | 2 | 80 |
| $\theta$ | 60 | 40 |
| $p$ | 0.5 | 0.43 |
| $c_0$ | 1000 | 1000 |
| $n$ | 5 | 5 |
| $s_{\max}$ | 1000 | 1000 |

Table 3.3: Parameter values used for counterfactual optimization

of T-cell presence. For each patch, we generate three randomized versions, where we randomly select one of the cells to mask for each version. We then compute the difference in predicted probabilities and predicted label between the randomized patches and the original patch (Table 3.4). Across all three data sets, we do not see a statistically significant difference in the predicted labels when comparing randomly masked patches to original patches. We do see a significant increase in the predicted probability value for CRC in the randomized images, although the mean change is very small at $1.38 \times 10^{-4}$. These results suggest our T-cell masking strategy did not introduce an artificial signal whereby simply removing cells at random will increase the chance that T-cells are predicted to be present.

Table 3.4: Difference in predicted probability and predicted positive labels between randomly masked patches and original patches, p-value obtained from a one-sample T-test.

| Cancer type | Predicted label | | Predicted probability | |
|:---|:---:|:---:|:---:|:---:|
| | Mean difference | p-value | Mean difference | p-value |
| Melanoma | $-2.22 \times 10^{-4}$ | 0.132 | $-1.45 \times 10^{-4}$ | $2.28 \times 10^{-5}$ |
| Breast tumor | $-1.93 \times 10^{-5}$ | 0.739 | $6.96 \times 10^{-6}$ | 0.588 |
| CRC | $8.95 \times 10^{-5}$ | 0.433 | $1.38 \times 10^{-4}$ | $2.48 \times 10^{-13}$ |

**Supplemental Note 2 Choice of IMC patch size**

In order to obtain enough TME samples to train our classifier models, we took advantage of the inherent heterogeneity in the TME and divided each tissue image into 48 µm × 48 µm patches and treated each patch as an independent sample during training. Here, we perform spectral analysis to study the relationship between spatial patterning of proteins at various length scales and CD8+ T-cell infiltration. Specifically, we compute the power spectral density (PSD) of each breast tumor image. The PSD shows the relative importance of patterning at various length scales ("wavelengths") in the expression map. We then compute the Pearson correlation between patterning of a protein at a given length scale and T-cell infiltration (Figure 3.5). Figure 3.5 shows there are significant information pertaining to T-cell infiltration at our selected patch size. For certain proteins such as HLA-DR and FSP1, we see significantly more information is present at longer length scales of around 200 µm. This result suggests that for sufficient amount of IMC data, the performance of our classifier model may be improved by increasing the size of image patches.

Figure 3.5: **Correlation between each frequency band of each protein channel and T-cell infiltration level (proportion of CD8+ T-cell patches) across all IMC images for the breast cancer data set**. Red dotted line indicates the patch size of 48 µm used in this work.

**Optimized perturbation for primary CRC and liver metastases**



Figure 3.6: Optimized perturbations for tissue sections from primary colorectal tumor and liver metastases aggregated to the patient (row) level (train cohort).

**Evaluation metric for all trained classifiers**

Table 3.5: Performance of different classifier models trained on melanoma, CRC, and breast tumor IMC images to predict the presence of T-cells ($p = 0.5$)

| Cancer type | Model | Accuracy | Precision | Recall | F1 | AUROC | MCC |
|---|---|---|---|---|---|---|---|
| Melanoma | Linear | 0.84 | 0.83 | 0.37 | 0.51 | 0.67 | 0.48 |
| | MLP | 0.85 | 0.79 | 0.48 | 0.59 | 0.72 | 0.53 |
| | U-Net | 0.86 | 0.72 | 0.59 | 0.64 | **0.76** | 0.56 |
| Breast | Linear | 0.82 | 0.57 | 0.16 | 0.24 | 0.57 | 0.23 |
| tumor | MLP | 0.86 | 0.71 | 0.42 | 0.52 | 0.69 | 0.47 |
| | U-Net | 0.87 | 0.70 | 0.50 | 0.58 | **0.73** | 0.52 |
| CRC | Linear | 0.86 | 0.71 | 0.19 | 0.29 | 0.59 | 0.31 |
| | MLP | 0.88 | 0.67 | 0.50 | 0.57 | 0.73 | 0.51 |
| | U-Net | 0.88 | 0.68 | 0.61 | 0.64 | **0.77** | 0.57 |

Table 3.6: Performance of different CNN and ViT models trained on IMC image patches of melanoma

| Model Name | U-Net | ResNet-18 | EfficientNet-B0 | MedViT |
|---|---|---|---|---|
| Number of parameters | 12.8M | 11.2M | 4.1M | 31.3M |
| Test accuracy | 0.86 | 0.86 | 0.851 | 0.853 |

*Chapter 4*

# FUTURE WORK

## 4.1 Adapting information-theoretic framework to optimize other cell properties with respect to environmental statistic

One can easily adapt our information-theoretic framework introduced in Chapter 2 to understand how variables other than receptor placement affects spatial sensing. Although this work is about optimizing receptor placement, the key quantity being tuned is the spatial distribution of receptor activity, hence our result is relevant to any variable that 1) affects receptor activity and 2) redistributes across space. To illustrate, consider a generalized model of receptor activation,

$$\mathbb{E}[A_i \mid c_i] = f(\boldsymbol{\theta}_i)\left(\frac{c_i}{c_i + K_d} + \alpha\frac{K_d}{c_i + K_d}\right), \tag{4.1}$$

where $f$ is an unspecified function of an arbitrary set of variables $\boldsymbol{\theta}_i$, and $f(\boldsymbol{\theta}_i)$ represents the "effective" number of receptors at position $i$. In this work, we considered the case where $\theta_i = r_i$ and $f(r_i) = r_i$, but other factors such as phosphorylation level and membrane curvature also affect local receptor activity $A_i$ (Rangamani et al., 2013). In this way, one can optimize spatial sensing by tuning variables other than receptor placement, by specifying alternative forms of $f$. For example, it is known that given uniformly distributed receptors, those found in membrane regions of higher curvature can exhibit higher activity (Rangamani et al., 2013). Suppose we want to know the optimal way to adjusT-cell shape to maximize information acquisition, by assuming a linear relationship between local curvature $\beta_i$ and "effective" receptor number, i.e., $f(r_i, \beta_i) = \beta_i r_i$. Given uniform receptors and a constraint on total contour length (or area) of the membrane, we quickly arrive at the optimal solution since this problem is now identical to our original formulation. The optimal strategy is to increase membrane curvature at regions of high ligand concentration, by making narrow protrusions (for details see Supplement, Section 2.10).

## 4.2 Optimizing spatial organization at different stages of information processing

Optimizing information transmission by organizing effectors in space can happen at all stages of signal processing within the cell, but is likely most effective at the receptor level. The most obvious reason is due to the data processing inequality,

which states that post-processing cannot increase information. Therefore, only optimization at the level of receptor activation can increase the total amount of information that is available to the cell. The second reason is due to the "hourglass" topology of cell signaling networks, which represent the fact that a large number of signaling inputs converge onto a small number of effectors internal to the cell (Csete and Doyle, 2004). For example, G-protein-coupled receptors, one of the largest group of cell surface receptors, drive downstream signaling through the same G-proteins. This feature makes optimizing spatial organization at later stages of information processing very difficult, since information can be easily lost by diffusion of effector molecules activated by different inputs, which ends up "mixing" different spatial signals.

## 4.3 Extensions of Morpheus

For future work, we would like to apply Morpheus to spatial transcriptomics data sets with hundreds to thousands of molecular channels. Although spatial transcriptomics can profile significantly more molecules compared to spatial proteomic techniques (Rodriques et al., 2019; Eng et al., 2019), the number of spatial transcriptomic profiles of human tumors is currently limited due to the cost, with most public data sets containing single tissue sections from 1–5 patients which is far too small to apply Morpheus. However, spatial transcriptomics is likely to be more standardized compared to proteomics, which use customized panels. As commercial platforms for spatial transcriptomics start to come online (Janesick et al., 2022), we will likely be seeing large scale spatial transcriptomics data sets in the near future, with $\sim 70$–$90\%$ of the same probes shared between experiments.

A technical extension of Morpheus involves incorporating prior knowledge of gene-gene interactions to model the causal relations between genes. Molecular features in tissue profiles can exhibit strong dependencies, therefore, changing the level of one molecule can affect the expression of others. For example, increased levels of interferon-gamma (IFN-$\gamma$) in the tumor microenvironment, can upregulate the expression of PD-L1 on tumor cells (Qian et al., 2018). In order to be more realistic and actionable, a counterfactual should maintain these known causal relations. We can apply a regularizer to penalize counterfactuals that are less feasible according to established gene interactions from knowledge graphs, such as Gene Ontology (Consortium, 2004).

Other extensions of Morpheus includes predicting cell-type specific perturbations,

which can be done by directly restricting the perturbation to only alter signals within specific cell types. Additionally, although we applied Morpheus to the specific problem of driving T-cells to infiltrate solid tumors, we can generalize our framework to predict candidate therapeutics that alter the localization of other cell types. For example, Morpheus can train a classifier model to predict localization of TAMs and compute perturbations predicted to reduce their abundance in the TME.

In this work, we focused on identifying generalized therapies by pooling predictions across multiple patient samples, but we can also apply Morpheus to find personalized therapy for treating individual patients. The variation in the optimized perturbations we observe among patients in both melanoma and liver data sets suggest personalize treatments could be significantly more effective compared to generalized therapies (Figure 3.3A, Figure 3.4A). Furthermore, Figure 3.4G shows that a therapeutic strategy could have highly variable effect even across different tissue samples from the same patient. This variability suggests that to generate therapy for an individual patient, it may be necessary to acquire significant quantities of biopsy data. We can then apply our optimization procedure to a random subset of the samples, and then test the resulting perturbation strategy on the remaining samples to see how well the strategy is predicted to perform across an entire tumor or other primary/secondary tumors.

Incorporating Morpheus in a closed loop with experimental data collection is another promising direction for future work. Data can be collected from patients or animal models with perturbed/engineered signaling context, and this data can be easily fed back into the classifier model to refine the model's prediction. The perturbation could be based on what the model predicts to be effective interventions, as is the case with Morpheus. We can also study tissue samples on which the model tends to make the most mistake and train the model specifically using samples from similar sources, such as similar patient strata or disease state.

# BIBLIOGRAPHY

Aoki, Tomohiro, et al. (2023). "The spatially resolved tumor microenvironment predicts treatment outcome in relapsed/refractory Hodgkin lymphoma". In: *bioRxiv*.

Berg, Howard C, and Edward M Purcell (1977). "Physics of chemoreception". In: *Biophysical Journal* 20.2, pp. 193–219.

Bhate, Salil S, et al. (2022). "Tissue schematics map the specialization of immune tissue motifs and their appropriation by tumors". In: *Cell Systems* 13.2, pp. 109–130.

Biasci, Daniele, et al. (2020). "CXCR4 inhibition in human pancreatic and colorectal cancers induces an integrated immune response". In: *Proceedings of the National Academy of Sciences* 117.46, pp. 28960–28970.

Binnewies, Mikhail, et al. (2018). "Understanding the tumor immune microenvironment (TIME) for effective therapy". In: *Nature Medicine* 24.5, pp. 541–550.

Bonaventura, Paola, et al. (2019). "Cold tumors: a therapeutic challenge for immunotherapy". In: *Frontiers in Immunology* 10, p. 168.

Bouzigues, Cédric et al. (2007). "Asymmetric redistribution of GABA receptors during GABA gradient sensing by nerve growth cones analyzed by single quantum dot imaging". In: *Proceedings of the National Academy of Sciences* 104.27, pp. 11251–11256.

Bruni, Daniela, Helen K Angell, and Jérôme Galon (2020). "The immune contexture and Immunoscore in cancer prognosis and therapeutic efficacy". In: *Nature Reviews Cancer* 20.11, pp. 662–680.

Buda, Mateusz, Ashirbani Saha, and Maciej A Mazurowski (2019). "Association of genomic subtypes of lower-grade gliomas with shape features automatically extracted by a deep learning algorithm". In: *Computers in Biology and Medicine* 109. DOI: 10.1016/j.compbiomed.2019.05.002.

Candès, Emmanuel J and Michael B Wakin (2008). "An introduction to compressive sampling". In: *IEEE Signal Processing Magazine* 25.2, pp. 21–30.

Chang, Chun-Hao et al. (2019). "Explaining Image Classifiers by Counterfactual Generation". In: *International Conference on Learning Representations (ICLR)*.

Chen, Ao et al. (2022). "Spatiotemporal transcriptomic atlas of mouse organogenesis using DNA nanoball-patterned arrays". In: *Cell* 185.10, pp. 1777–1792.

Chen, Yunching et al. (2015). "CXCR4 inhibition in tumor microenvironment facilitates anti-programmed death receptor-1 immunotherapy in sorafenib-treated hepatocellular carcinoma in mice". In: *Hepatology* 61.5, pp. 1591–1602.

Choe, Joseph H, Jasper Z Williams, and Wendell A Lim (2020). "Engineering T cells to treat cancer: the convergence of immuno-oncology and synthetic biology". In: *Annual Review of Cancer Biology* 4, pp. 121–139.

Chou, Ching-Shan et al. (2011). "Noise filtering tradeoffs in spatial gradient sensing and cell polarization response". In: *BMC Systems Biology* 5.1, pp. 1–16.

Consortium, Gene Ontology (2004). "The Gene Ontology (GO) database and informatics resource". In: *Nucleic Acids Research* 32, pp. D258–D261.

Csete, Marie and John Doyle (2004). "Bow ties, metabolism and disease". In: *TRENDS in Biotechnology* 22.9, pp. 446–450.

Danenberg, Esther et al. (2022). "Breast tumor microenvironment structures are associated with genomic features and clinical outcome". In: *Nature Genetics* 54.5, pp. 660–669.

Eling, Nils et al. (2020). "cytomapper: an R/Bioconductor package for visualization of highly multiplexed imaging data". In: *Bioinformatics* 36.24, pp. 5706–5708.

Endres, Robert G and Ned S Wingreen (2008). "Accuracy of direct gradient sensing by single cells". In: *Proceedings of the National Academy of Sciences* 105.41, pp. 15749–15754.

Eng, Chee-Huat Linus et al. (2019). "Transcriptome-scale super-resolved imaging in tissues by RNA seqFISH+". In: *Nature* 568.7751, pp. 235–239.

Foeng, Jade, Iain Comerford, and Shaun R McColl (2022). "Harnessing the chemokine system to home CAR-T cells into solid tumors". In: *Cell Reports Medicine*.

Fridman, Wolf H et al. (2017). "The immune contexture in cancer prognosis and treatment". In: *Nature Reviews Clinical Oncology* 14.12, pp. 717–734.

Ghanem, Ismael et al. (2014). "Insights on the CXCL12-CXCR4 axis in hepatocellular carcinoma carcinogenesis". In: *American Journal of Translational Research* 6.4, p. 340.

Giesen, Charlotte et al. (2014). "Highly multiplexed imaging of tumor tissues with subcellular resolution by mass cytometry". In: *Nature Methods* 11.4, pp. 417–422.

Goltsev, Yury et al. (2018). "Deep profiling of mouse splenic architecture with CODEX multiplexed imaging". In: *Cell* 174.4, pp. 968–981.

Haslam, Alyson and Vinay Prasad (2019). "Estimation of the percentage of US patients with cancer who are eligible for and respond to checkpoint inhibitor immunotherapy drugs". In: *JAMA Network Open* 2.5, e192535–e192535.

Hegde, Priti S and Daniel S Chen (2020). "Top 10 challenges in cancer immunotherapy". In: *Immunity* 52.1, pp. 17–35.

Hoch, Tobias et al. (2022). "Multiplexed imaging mass cytometry of the chemokine milieus in melanoma characterizes features of the response to immunotherapy". In: *Science Immunology* 7.70.

Hu, Bo et al. (2010). "Physical limits on cellular sensing of spatial gradients". In: *Physical Review Letters* 105.4, p. 048104.

Hughes, Catherine E and Robert JB Nibbs (2018). "A guide to chemokines and their receptors". In: *The FEBS Journal* 285.16, pp. 2944–2971.

Huston, Stephen J et al. (2015). "Neural encoding of odors during active sampling and in turbulent plumes". In: *Neuron* 88.2, pp. 403–418.

Iida, Fumiya and Surya G Nurzaman (2016). "Adaptation of sensor morphology: an integrative view of perception from biologically inspired robotics perspective". In: *Interface Focus* 6.4, p. 20160016.

Janesick, Amanda et al. (2022). "High resolution mapping of the breast cancer tumor microenvironment using integrated single cell, spatial and in situ analysis of FFPE tissue". In: *Biorxiv*, pp. 2022–10.

Klaise, Janis et al. (2021). "Alibi Explain: Algorithms for Explaining Machine Learning Models". In: *Journal of Machine Learning Research* 22.181, pp. 1–7.

Kohli, Karan, Venu G Pillarisetty, and Teresa S Kim (2022). "Key chemokines direct migration of immune cells in solid tumors". In: *Cancer gene therapy* 29.1, pp. 10–21.

Krause, Andreas, Ajit Singh, and Carlos Guestrin (2008). "Near-optimal sensor placements in Gaussian processes: Theory, efficient algorithms and empirical studies". In: *Journal of Machine Learning Research* 9.2.

Lanitis, E et al. (2017). "Mechanisms regulating T-cell infiltration and activity in solid tumors". In: *Annals of Oncology* 28, pp. xii18–xii32.

Lee, Joo Sang and Eytan Ruppin (2019). "Multiomics Prediction of Response Rates to Therapies to Inhibit Programmed Cell Death 1 and Programmed Cell Death 1 Ligand 1". In: *JAMA Oncology* 5.11, pp. 1614–1618.

Levine, Jacob H. et al. (July 2015). "Data-Driven Phenotypic Dissection of AML Reveals Progenitor-like Cells that Correlate with Prognosis". In: *Cell* 162.1, pp. 184–197.

Looveren, Arnaud Van and Janis Klaise (2021). "Interpretable counterfactual explanations guided by prototypes". In: *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 650–665.

Martinez, Marina and Edmund Kyung Moon (2019). "CAR T cells for solid tumors: new strategies for finding, infiltrating, and surviving in the tumor microenvironment". In: *Frontiers in Immunology* 10, p. 128.

McInnes, L., J. Healy, and J. Melville (Feb. 2018). "UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction". In: *ArXiv e-prints*. arXiv: 1802.03426.

Moffitt, Jeffrey R, Emma Lundberg, and Holger Heyn (2022). "The emerging landscape of spatial profiling technologies". In: *Nature Reviews Genetics* 23.12, pp. 741–759.

Mugler, Andrew, Andre Levchenko, and Ilya Nemenman (2016). "Limits to the precision of gradient sensing with spatial communication and temporal integration". In: *Proceedings of the National Academy of Sciences* 113.6, E689–E695.

Nieto, Marta et al. (1997). "Polarization of chemokine receptors to the leading edge during lymphocyte chemotaxis". In: *The Journal of Experimental Medicine* 186.1, pp. 153–158.

Paszke, Adam et al. (2019). "PyTorch: An imperative style, high-performance deep learning library". In: *Advances in Neural Information Processing Systems* 32.

Pignata, Aurora et al. (2019). "A spatiotemporal sequence of sensitization to slits and semaphorins orchestrates commissural axon navigation". In: *Cell Reports* 29.2, pp. 347–362.

Pitt, JM et al. (2016). "Targeting the tumor microenvironment: removing obstruction to anticancer immune responses and immunotherapy". In: *Annals of Oncology* 27.8, pp. 1482–1492.

Pittet, Mikael J, Olivier Michielin, and Denis Migliorini (2022). "Clinical relevance of tumour-associated macrophages". In: *Nature Reviews Clinical Oncology* 19.6, pp. 402–421.

Qian, Jiawen et al. (2018). "The IFN-$\gamma$/PD-L1 axis between T cells and tumor microenvironment: hints for glioma anti-PD-1/PD-L1 therapy". In: *Journal of neuroinflammation* 15.1, pp. 1–13.

Rangamani, Padmini et al. (2013). "Decoding information in cell shape". In: *Cell* 154.6, pp. 1356–1369.

Rejniak, Katarzyna Anna et al. (2013). "The role of tumor tissue architecture in treatment penetration and efficacy: an integrative study". In: *Frontiers in Oncology* 3, p. 111.

Rodriques, Samuel G et al. (2019). "Slide-seq: A scalable technology for measuring genome-wide expression at high spatial resolution". In: *Science* 363.6434, pp. 1463–1467.

Savas, Peter et al. (2016). "Clinical relevance of host immunity in breast cancer: from TILs to the clinic". In: *Nature Reviews Clinical Oncology* 13.4, pp. 228–241.

Schraufstatter, Ingrid U et al. (2012). "The chemokine CCL18 causes maturation of cultured monocytes to macrophages in the M2 spectrum". In: *Immunology* 135.4, pp. 287–298.

Schürch, Christian M et al. (2020). "Coordinated Cellular Neighborhoods Orchestrate Antitumoral Immunity at the Colorectal Cancer Invasive Front". In: *Cell* 182.5, pp. 1341–1359.

Servant, Guy et al. (1999). "Dynamics of a chemoattractant receptor in living neutrophils during chemotaxis". In: *Molecular Biology of the Cell* 10.4, pp. 1163–1178.

Steele, Maria M et al. (2023). "T-cell egress via lymphatic vessels is tuned by antigen encounter and limits tumor control". In: *Nature Immunology* 24.4, pp. 664–675.

Sullivan, Kevin M et al. (2023). "Blockade of interleukin 10 potentiates antitumour immune function in human colorectal cancer liver metastases". In: *Gut* 72.2, pp. 325–337.

Tsaur, Igor et al. (2021). "Immunotherapy in prostate cancer: new horizon of hurdles and hopes". In: *World journal of urology* 39, pp. 1387–1403.

Verma, Sahil et al. (2020). "Counterfactual explanations and algorithmic recourses for machine learning: A review". In: *arXiv preprint arXiv:2010.10596*.

Wachter, Sandra, Brent Mittelstadt, and Chris Russell (2017). "Counterfactual explanations without opening the black box: Automated decisions and the GDPR". In: *Harv. JL & Tech.* 31, p. 841.

Wang, Zhijun et al. (2023). "Extracellular vesicles in fatty liver promote a metastatic tumor microenvironment". In: *Cell Metabolism*.

Weber, Michele et al. (2013). "Interstitial dendritic cell guidance by haptotactic chemokine gradients". In: *Science* 339.6117, pp. 328–332.

Wu, Zhenqin et al. (2022). "Graph deep learning for the characterization of tumour microenvironments from spatial protein profiles in tissue specimens". In: *Nature Biomedical Engineering*, pp. 1–14.

Zhang, Haiying, H Amalia Pasolli, and Elaine Fuchs (2011). "Yes-associated protein (YAP) transcriptional coactivator functions in balancing growth and differentiation in skin". In: *Proceedings of the National Academy of Sciences* 108.6, pp. 2270–2275.