



PHYSICS/NEURO 141. THE PHYSICS OF SENSORY SYSTEMS IN BIOLOGY

Prof. Aravi Samuel, Department of Physics.

Dr. Alina Vrabioiu, Department of Physics.

Ariana-Dalia Vlad, Department of Physics.

Living organisms use sensory systems to inform themselves of the sights, sounds, and smells of their surrounding environments. Sensory systems are physical measuring devices, and are therefore subject to the laws and limits of physics. Here, we will consider the physics of sensory measurement and perception, and study ways that biological systems have solved their underlying physical problems. We will discuss specific cases in vision, olfaction, and hearing from a physicist's point of view. **N.B.:** This is one course that is taught as both Neuro 141 and Physics 141. Whether a student enrolls in Neuro 141 or Physics 141 is flexible, to best suit each student's academic planning and concentration requirements.

ARAVI SAMUEL grew up in upstate New York, and graduated with a BA in physics and PhD in biophysics from Harvard. For his PhD, he studied the biophysics of bacterial chemotaxis with Howard Berg, a pioneering biophysicist at Harvard. Berg established our mechanistic understanding of sensory perception in bacterial chemotaxis, one of the best understood systems in biology ¹ (Figs. 114, 3). Aravi thinks about the biophysics of circuits and behavior in bacteria, worms, and flies in much the same mechanistic way as his mentor. To find out more about Aravi's work, see [his website](#).

Email: samuel@physics.harvard.edu

ALINA VRABIOIU was a postdoc with Howard Berg, and is continuing her postdoctoral work with Aravi on the bacterial flagellar motor.

Email: alina_vrabioiu@fas.harvard.edu

ARIANA-DALIA VLAD was a student in this course in 2022, and is currently a senior studying quantitative biology and physics.

Email: avlad@college.harvard.edu

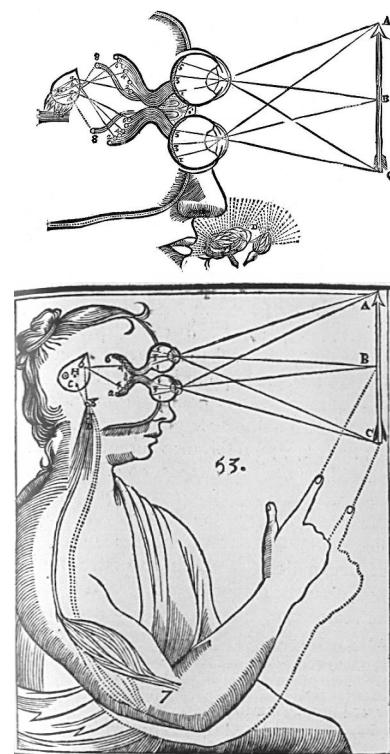


Figure 1: René Descartes (1596-1650) depicts sensory integration between two of Aristotle's five senses, vision and smell (above) and a sensorimotor pathway that mediates a behavioral response (below). Descartes was an early believer that sensory perception and animal behavior could be reduced to physical processes and the mechanical flow of information through an organism, the conceptual framework for modern neuroscience and this course. Descartes did make some mistakes. In these drawings, sensory perception and behavior are mediated by the pineal gland, not the brain.

LECTURES (Tu, Th 9-10:15 EST in NW255) will teach fundamental physics and math that illuminates sensory neuroscience. We will also discuss classic papers about vision, hearing, and chemical sensing that can be unlocked with physical reasoning.

SECTIONS will occur at the same time and same place as lectures, occasionally taking the place of an ordinary lecture but where the teaching fellows will go over problem-solving that is relevant to problem sets.

OFFICE HOURS will be held by Aravi, Alina, and Ariana by appointment.

COURSE MATERIALS will be distributed as this main PDF course packet, containing hyperlinks to directly download all other additional required and recommended reading material (chapters of other books and primary papers). There is no required course textbook. **Primary papers** that we cover each week will draw on both classic and current studies in vision, hearing, and chemical sensation.

GRADES will be based on Problem sets and coding assignments, 60%; Final student presentation (20%); Class participation (20%).

- **Problem Sets** We will have occasional problem sets, roughly every third week, that will include questions about the physics and biology of sensory systems. The problem sets will be longer than the typical weekly problem set given in most classes, so it is advisable to start early. In section, you will go over solving similar problems with the help of the TFs. We will occasionally include small coding assignments involving simulations or data analysis.
- **Final presentation** At the end of the course, students will make a final presentation. Students will
 1. Write their own review-style paper
 2. Present a short talk that describes the importance of their chosen paper and communicates any essential physics or math needed as background
 3. Construct their own coding project (e.g., data analysis or simulation) that illuminates their paper

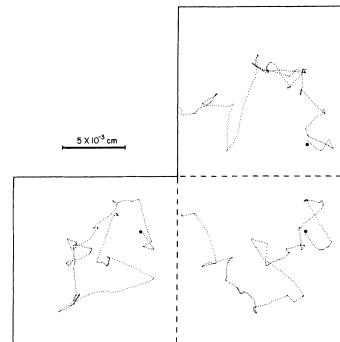


Figure 2: Bacteria perform chemotaxis by a random walk. The three-dimensional track of a single swimming bacteria viewed in xy, yz, and xz projections. The movement can be characterized as an alternating sequence of runs (periods of forward movement) and tumbles (periods of erratic rotational movement). When the bacteria is pointed in a direction it wants to go, runs get longer. The random walk becomes biased towards preferred environments.

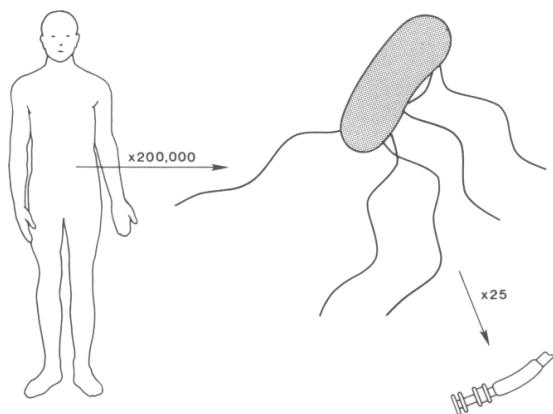


Figure 3: Man, *E. coli*, and its flagellar motor, a study in scale. Physics is different at different scales. No sensory or behavioral system is better understood than that of bacterial chemotaxis, where behavior can be reduced to the rotation of individual bacterial flagella, and every event from perception of sensory inputs (chemoreceptors binding to receptors) to motor control has been determined.

Contents

CALENDAR	5
WHY STUDY SENSORY SYSTEMS?	6
HOW BIOLOGISTS STUDY SENSORY SYSTEMS	16
SOME STATISTICAL MECHANICS	33
COLOR VISION AND THE DIMENSIONALITY OF PERCEPTION	39
THE STATISTICS OF PHOTON ABSORPTION BY PHOTORECEPTORS	48
HOW MANY PHOTONS CREATE VISION	54
SINGLE PHOTONS AND SINGLE ROD CELLS	62
RELIABLE SIGNAL TRANSDUCTION IN THE ROD CELL	75
OLFACTION AND DIFFUSION	80
THE DIFFUSION COEFFICIENT	103
COUNTING MOLECULES	110
COUNTING MOLECULES WITH RECEPTORS	117
HEARING WITH HAIR CELLS	121
NON-LINEARITIES IN HEARING	132
THE SENSITIVITY OF THE HAIR CELL	137

CALENDAR

Class Meeting	Lecture Topic	Slides
Sep 5	Why Study Sensory Systems?	Slides
Sep 7	How biologists study sensory systems	Slides
Sep 12	Some Statistical Mechanics	Slides
Sep 14	Vision	Slides
Sep 19	Vision	Zoom
Sep 21	Review Section	Slides
Sep 26	Vision	Problem Set One Due
Sep 28	Vision	Slides
Oct 3	Vision	Slides
Oct 5	Vision	Slides
Oct 10	Vision	Slides
Oct 12	Review Section	
Oct 17	Smell	Slides Problem Set Two Due
Oct 19	Smell	Slides
Oct 24	Smell	Slides
Oct 26	Smell	Slides
Oct 31	Smell	Slides
Nov 2	Review Section	Slides ,
Nov 7	Hearing	Slides , Problem Set Three Due
Nov 9	Hearing	Slides
Nov 14	Hearing	Slides
Nov 16	Hearing	Slides
Nov 21	Review Section	
Nov 28	Student Presentations	Problem Set Four Due
Nov 30	Student Presentations	
Dec 5	Prof. Samuel's Birthday	

WHY STUDY SENSORY SYSTEMS?

OUR BEHAVIORS begin with sensory perception. We gather information about the sights, sounds, smells, tastes, and textures of our worlds using our sensory systems. We respond to this sensory information using our brains and motor systems. We have wondered about our capacity for sensory perception as long as we have wondered about ourselves. In his treatise *de Anima*, Aristotle focused on sensory perception and motility, arguing that these were the most essential biological processes in making us alive. He thought sensory perception and motility would be windows to the soul.

We might argue with Aristotle about the “soul”, but a modern neuroscientist would accept that our sensory and motor systems are credible windows into the brain. Understanding how the human brain creates cognition and behavior is our major challenge in 21st-century science. To make progress, we can make a strong case to start with sensory systems. Every neuron is an “information processing unit” with the task of mapping incoming signals to outgoing signals. Incoming signals might come from other neurons (e.g., by synaptic communication) or from the rest of the body (e.g., by non-synaptic chemical communication). Outgoing signals are sent to other neurons, muscle cells, or the rest of the body through synapses and chemical messengers. The experimental advantage with sensory neurons is that these cells provide the scientist with a direct handle on their most salient incoming signals, environmental stimuli that can be controlled in the laboratory, like photons for photoreceptors or molecules for chemoreceptors. If we can understand the principles by which one sensory neuron maps an incoming signals to outgoing signals, we can illuminate the principles by which any neuron maps inputs into outputs.

Evolutionary conservation and homology allows us to leverage sensory neurons to understand common features of any neuron. “*Nothing in biology makes sense except in the light of evolution.*” – Dobzhansky. We will learn that sensory neurons share general computational principles in the mapping of incoming signals to outgoing signals that can be recognized throughout biology, from single-celled microorganisms to the photon-counting rod cells in the human retina. that abstract operating principles with other neurons, More concretely, any given sensory neuron in an animal can share homologous molecules and mechanisms, with other sensory neurons across Aristotle’s five sensory modalities and beyond. The molecular basis of olfaction, where olfactory neurons smell molecules from the environment (Fig. 5) is closely linked to the molecular basis of synaptic



Figure 4: Aristotle’s Five Senses. Vision, Hearing, Taste, Smell, Touch. Hearing and touch are both mediated by mechanosensory neurons. Smell and taste are both mediated by chemosensory neurons.

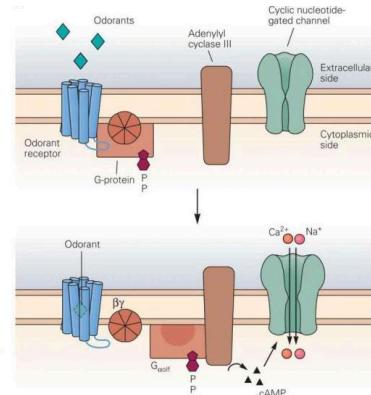


Figure 5: Odorant receptors. Binding an odorant causes an odorant receptor to interact with a G-protein signal transduction cascade that increases cAMP concentration. Elevated [cAMP] opens cyclic nucleotide-gated cation channels, causing a change in membrane potential, the neural signal for olfactory detection. Here, we illustrate a “metabotropic” olfactory receptor where stimulus energy is indirectly converted to neural activity through biochemistry that activates separate ion channels. Some olfactory receptors are metabotropic, but some, like insect olfactory receptors, are ionotropic (like the neurotransmitter receptor shown in Fig. 6 where the receptor contains its own ion channel).

communication, where postsynaptic neurons essentially “smell” neurotransmitters released by presynaptic neurons at synaptic clefts (Fig. 6).

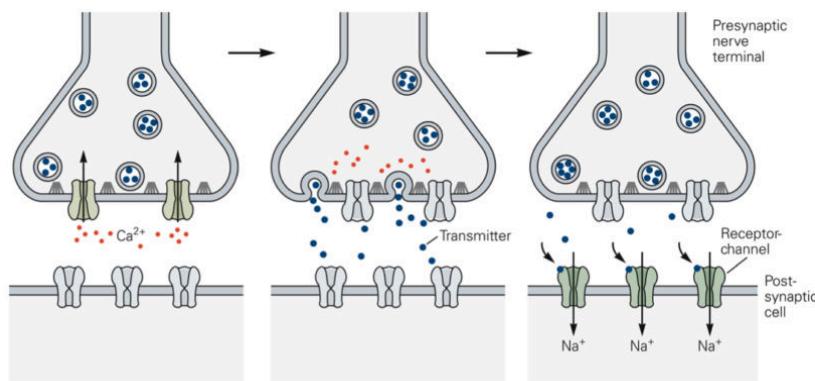


Figure 6: Synaptic transmission as olfaction. An action potential arriving at the terminal of a presynaptic neuron causes voltage-gated Ca^{2+} channels at the active zone to open. Channel opening and the rise in intracellular calcium levels causes vesicles containing neurotransmitter to fuse with the cell membrane and release their contents. Released neurotransmitter molecules diffuse across the synaptic cleft and are “smelled” by specific receptors on the postsynaptic membrane. Here, we illustrate an “ionotropic” receptor where neurotransmitter binding is directly converted to neural activity by opening an ion channel in the receptor itself. Some neurotransmitter receptors are ionotropic, but some are metabotropic (in the same way as the vertebrate olfactory receptor shown in Fig. 5).

Evolution creates bigger brains with added functions, but always by reusing and adapting a relatively small set of especially pliant information-processing mechanisms. The first recognizable nervous systems probably belonged to animals where only one or two cells carried out every information-processing step from perception to action. Larger nervous systems separate the many information-processing steps that were carried out by single primordial neurons, dividing functions among different neurons across multicellular circuits and brains, always by repeating, reusing, and adapting homologous building blocks. If we can understand how one sensory neuron responds to the outside world, we are well equipped to understand how any neuron responds to the inside world of the brain.

Another reason that physicists might want to study sensory systems is that sensory phenomena are particularly amenable to identifying and characterizing the primary physical triggers that turn environmental stimuli (photons, chemicals, mechanical force) into neural activity. A pattern of stimulus energy is translated into a pattern of neural activity that is then used by the brain to assess its surrounding world. Sensory systems are highly evolved, often working at the physical limits of operation. We see individual photons. Many organisms smell individual molecules. Understanding the physical challenges in counting photons and molecules requires understanding the physics of light, thermal noise, molecular diffusion, fluid mechanics, and statistical mechanics. Understanding sensory systems is a motivation for learning many areas of fundamental physics in the context of concrete biological application. We will learn this physics from the “ground up” in a way that is directly applied to biological

questions. Not only might we learn physics more deeply by *using* it, we will gain practice in using physics to build models of the natural world, an important skill beyond neuroscience.

Beyond Aristotle

THE FIVE SENSES that Aristotle cataloged in *de Anima* remain the five senses of conventional wisdom. But if the essence of sensory perception is gathering *information*, there are countless legitimate sensory systems throughout biology from microorganisms to man. Most animals sense pain, called nociception. Many animals sense the position and posture of their own bodies, called proprioception. Many animals sense internal states like visceral information about sickness or hunger that is communicated to the brain, called interoception. Many animals have sensory modalities that we do not. Weakly electric fish sense the perturbations of self-generated electric fields as a type of radar (Fig. 8). Bats use sound and hearing as another type of radar.

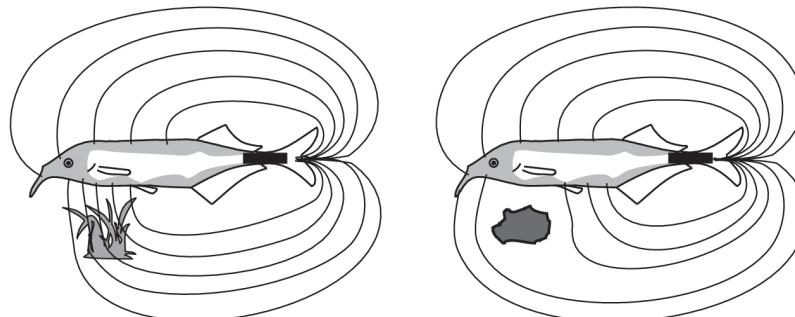


Figure 7: **Thermal imaging cameras.**
The pit organ of vipers is an exquisitely sensitive thermal imaging camera, albeit with lower resolution than this one formed by a man-made camera.

Sensory Receptors

We know how many environmental stimuli activate their corresponding receptor molecules. One large, highly conserved family of G-protein coupled receptors are common to all animals. GPCRs have proved especially pliant receptors across evolutionary history, having been adapted to detect photons (rhodopsin molecules in our retina), chemicals (olfactory receptors in our noses), and countless neurotransmitter and hormonal receptors that effectively “smell” molecules in our brains and bodies. GPCRs detect photons by virtue of a pigment (retinal) that changes shape upon photon absorption, triggering the conformational change in the surrounding protein that modulates its activity (Fig 38). GPCRs detect molecules when ligand binding

Figure 8: **Imaging with electric fields.**
We create visual maps of our surroundings our eyes, but other animals do the same with other sense. Schematic two-dimensional drawings of the electric fields of a *Gnathonemus petersii* distorted by a water plant (good conductor, left) and a stone (isolator, right). The fish is viewed from the side. Electrical field lines are drawn as thin lines. The electrorreceptive body surface of the fish is shown in grey.

directly triggers a conformational change that modulates protein activity.

The conformational change that accompanies the activation of an olfactory receptor has now been visualized. Unlike mammalian olfactory receptors which are GPCRs and *metabotropic* (meaning that they activate an intracellular signaling pathway that eventually changes the electrical activity of the cell) insect olfactory receptors are *ionotropic* (meaning that odorant binding directly opens an ion channel in the receptor itself). **Cryo-electron microscopy** allowed direct visualization of this structural change – opening and closing the ion channel within a receptor– by comparing bound and unbound states (Fig. 10). The evolutionary divergence between insect ionotropic olfactory receptors and vertebrate metabotropic receptors highlights the ancientness of a critical sensory modality that has evolved separately in different phyla for hundreds of millions of years. Even older are the olfactory receptors and signaling pathways that allow bacteria to ‘smell’ their environments, using “metabotropic” mechanisms without shared evolutionary history to animals.

Sometimes we know what receptors sense a particular stimulus, but not how they do it. Most organisms sense temperature. Temperature sensing is critical for small animals that do not regulate their own body temperatures – called poikilotherms or ectotherms or cold-blooded (if they have blood). Ectotherms rely on environmental temperature to regulate their own body temperatures by moving from place to place, think of reptiles sunning themselves to warm up. Thermosensation can also be used to gather information, like the pit organs of certain snakes that act like image-forming *pinhole cameras*, seeing warm objects by focusing infrared radiation onto highly sensitive thermoreceptors (Fig. 7). Thus, the most sensitive known thermoreceptors in biology can sense temperature changes as small as 0.001 °C/sec! Many thermoreceptors have been identified – like certain insect gustatory receptors that evolved to sense warming and cooling. However, the precise biophysical mechanism for temperature sensing remains poorly understood. While we know what receptors are needed, we do not know how thermal stimuli as small as millidegrees drive the specific changes in receptors that toggle their activities.

Receptor diversity within sensory modalities

Evolutionary divergence of receptors within a sensory modality can increase the sophistication and dimensionality of sensory coding. We have *trichromatic* vision because our retinas have three types of *cone cells* that are sensitive to long, medium, or short wavelengths of

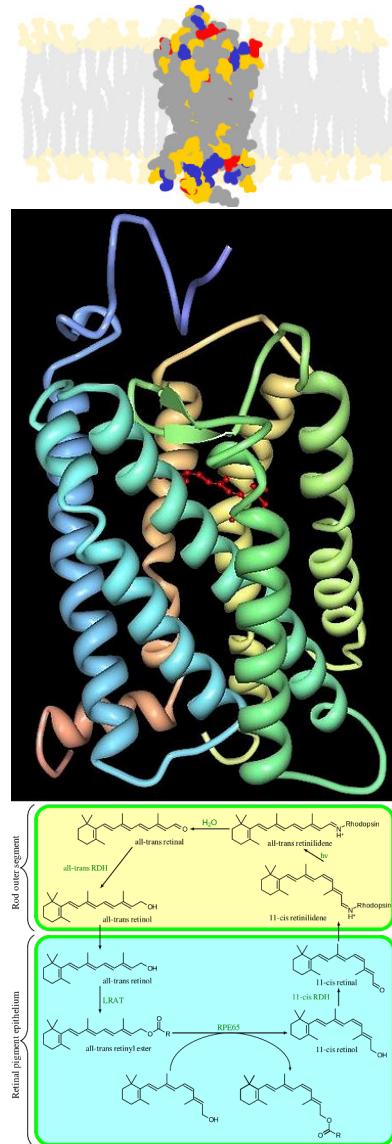


Figure 9: **Rhodopsin.** Three dimensional structure of bovine rhodopsin, a membrane bound protein. The chromophore, retinal, is embedded within the protein. Absorption of light energy, $h\nu$, causes a conformational change – 11-cis-retinal becomes all-trans-retinal – which alters protein structure to activate G-protein-coupled signal transduction. Biochemistry in the rod cell restores the cis configuration.

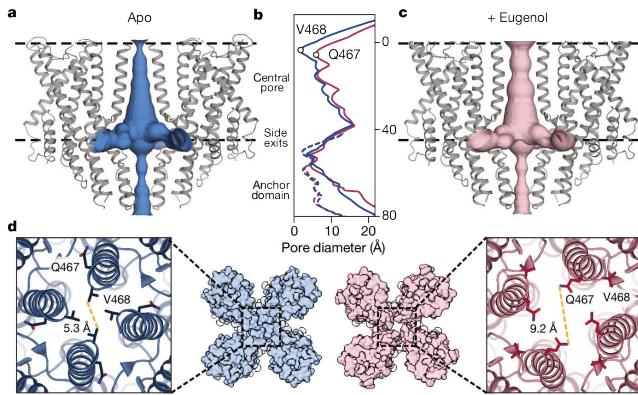


Figure 10: Odorant-evoked opening of an ionotropic olfactory receptor. a, c, The channel pores of unbound (a, blue) and eugenol-bound (c, pink). Black dashed lines, membrane boundaries. b, The diameter of the ion conduction pathway (solid lines) and along the anchor domain (dashed lines). d, Close-up view of the pore from the extracellular side. Marmol et al. (2021).

bright light by virtue of three slightly different rhodopsin molecules (Fig. 14). We have one type of rod photoreceptor that is sensitive to dim light. Because we only have one type of rod photoreceptor with one type of wavelength sensitivity, vision in dim light is *monochromatic*. But unlike the cone cells, our rod cells can detect single photons, a level of sensitivity that is needed to see the dimmest stars in the night sky.

Sensory discrimination is more complicated when different stimulus types must be separated and identified by different receptors. In our noses, diverse *olfactory receptors* in the nasal epithelium are tuned to different volatile chemicals. Mammalian noses contain many different types of olfactory receptor neurons, each with its own molecular olfactory receptor (Fig. 11). When we whiff a scent, blends of different types and concentrations of molecules enter the nasal cavity. These molecules bind to different olfactory receptors with different chemical specificity. Odorant molecules can exhibit relatively small differences in their chemical structure and properties (Fig. 12). Thus, one receptor can bind (and be activated) by many different kinds of molecule. One molecule can bind (and activate) many different kinds of receptors. Because typical smelly objects will emit many different odorant molecules, patterns of olfactory receptor activity can be expected to be complicated. A *combinatorial code* might be needed to discriminate odor molecules and identify smells (Fig. 13).

When many receptors in many cells are used to collect information along the different dimensions of a sensory modality (different colors in vision, different chemicals in smell), the information must be collected and analyzed by downstream circuits. Information travels from sensory receptors to our brains along *afferent* neuronal processes. Our modern circuit-level understanding of sensory processing

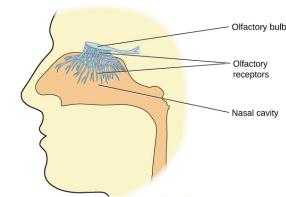


Figure 11: Olfactory receptor neurons innervate the olfactory epithelium in the human nasal cavity.

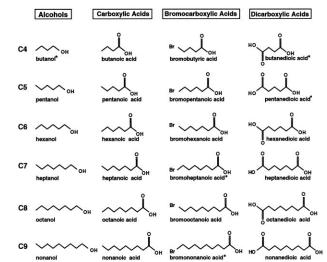


Figure 12: Volatile odorant molecules among the many thousands of chemicals that we can smell.

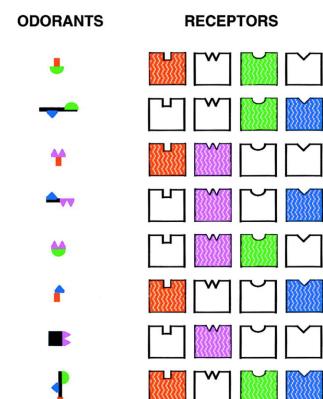
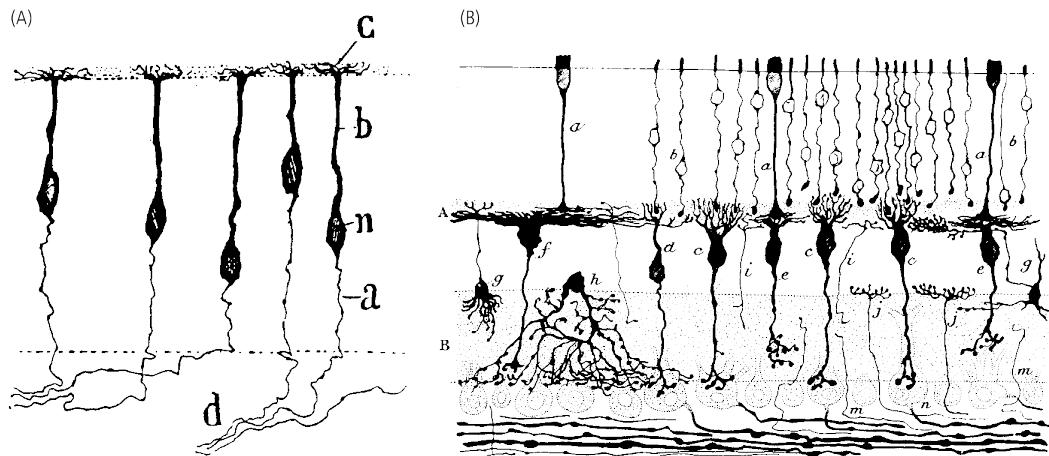


Figure 13: Cartoon of the combinatorial code where differently shaped molecules can bind to olfactory receptors with differently shaped binding pockets. Mammals can have 100s of different olfactory receptors that presumably detect many 1000s of different odorant molecules.

began with the anatomical studies of **Santiago Ramón y Cajal** (1852–1934) who visualized the detailed structure of many types of sensory neurons, as well as the wires and intercellular contacts (synapses) that carried information to the brain (Fig. 14). Cajal did this by perfecting Golgi's method for sparsely staining neural tissues with silver salts. The Golgi staining method made individual neurons visible at random. Each neuron could be categorized and reconstructed using light microscopy. Collecting neuronal structures across samples, Cajal worked out the basic structure and connectivity of many neural circuits, revealing many principles of information processing that remain valid today.



The physics in sensory perception

STIMULUS ENERGY is the initiating event of all sensory perception. Photoreceptors transduce the energy of photons. Chemoreceptors transduce chemical binding energy. Mechanoreceptors transduce mechanical stretch or movement. Information is generated by some ‘energy absorbing’ primary event in sensory detection – a photoreceptor molecule absorbs a photon, a chemoreceptor binds a molecule, a mechanoreceptor is tugged open.

Sensory information eventually reaches circuits that shape organism behavior. For this to happen, the primary sensory signal – the energy from absorbing single photons or binding single molecules – must be filtered, amplified, and transmitted through intervening pathways. These different stages of sensory perception can be inter-

Figure 14: Santiago Ramón y Cajal has claim to be the father of modern neuroscience through his detailed anatomical studies of neural tissue. (A) Bipolar sensory neurons from mammalian olfactory mucosa. a, Axon; b, peripheral process; c, sensory dendrite; d, axon; n, nucleus. (B) Section of retina of an adult dog. A, Outer plexiform (synaptic) layer; B, inner plexiform (synaptic) layer; a, cone fiber; b, rod cell body and fiber; c, rod bipolar cell with vertical dendrites; d, cone bipolar cell with vertical dendrites; e, cone bipolar cell with flattened dendrites; f, giant bipolar cell with flattened dendrites; g, special cells stained very rarely (perhaps inter-plexiform cells); h, diffuse amacrine cell; i, ascendant nerve fibers (probably processes of cell not well stained); j, centrifugal fibers coming from central nervous system; m, nerve fiber (probably again of poorly stained cell); n, ganglion cell. (from Fain, Chapter 1)

preted in engineering terms (Fig.). Every sensory system has some sort of ‘antenna’ that serve to detect primary signals. At the cellular level, the antenna for vision, for example, might be interpreted to be the photoreceptor cell. At the molecular level, the antenna might be molecular rhodopsin which directly absorbs each photon.

The primary sensory signal brings an amount of energy to the antenna that must be amplified. The energy of a single photon is calculated using Einstein’s relation $E = hc/\lambda$. The energy of one blue-green photon ($\lambda=500 \text{ nm}$) is $4 \times 10^{-19} \text{ J}$.

After detection, primary signals must be filtered or analyzed to create useful information. Filtering is needed to separate true signals from noise. Thermal energy constitutes a constant and unavoidable source of random noise that contaminates the sensory antenna. One of the most useful results from **Statistical Mechanics** is the *Equipartition Theorem* that states that all bodies in thermal equilibrium have $k_B T/2$ of average energy in every quadratic degree of freedom. For example, a gas molecule has an average kinetic energy of $k_B T/2$ along each axis of motion. Kinetic energy is quadratic in velocity components, three degrees of freedom in three-dimensional space:

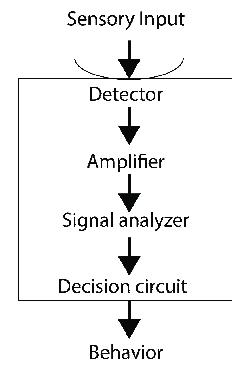
$$\left\langle \frac{mv_x^2}{2} \right\rangle = \left\langle \frac{mv_y^2}{2} \right\rangle = \left\langle \frac{mv_z^2}{2} \right\rangle = \frac{k_B T}{2} \quad (1)$$

Thermal energy at room temperature ($\sim 25^\circ \text{C}$), roughly the temperature of most biology, is $4 \times 10^{-21} \text{ J}$. By comparison, the energy of one blue photon is $\sim 100 \times$ thermal energy. Although thermal energy fluctuates about $k_B T/2$, because a single photon is so much more energetic, its corresponding signal is much larger than this fluctuating noise. If the energy threshold for single photon detection is set to be much larger than $k_B T/2$ and below $\sim 100k_B T/2$, the energy of ‘true signals’ caused by single photons can be reliably filtered from ‘false signals’ caused by thermal energy fluctuations.

THE BOLTZMANN DISTRIBUTION predicts that any detection threshold can be broached, albeit perhaps rarely, by thermal noise. For a body at thermal equilibrium, the probability of a thermal fluctuation with energy E_c is exponentially distributed:

$$P(E_c) \propto e^{-E_c/k_B T} \quad (2)$$

This means that the probability that a rhodopsin molecule reaches energies comparable to visible photons by thermal fluctuations, although vanishingly small, is also non-zero. Our eyes have hundreds of millions of photoreceptor cells (Fig. 21). Each photoreceptor cell is filled with billions of rhodopsin molecules. Although the likelihood



$$h = 6.6 \times 10^{-34} \text{ J s}$$

$$k_B = 1.4 \times 10^{-23} \text{ J K}^{-1}$$

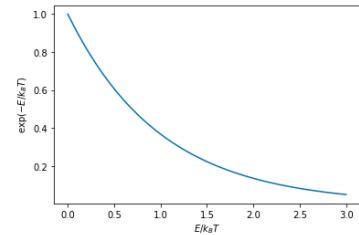


Figure 15: The Boltzmann Distribution

that a single rhodopsin molecule reaches detection threshold because of temperature might be small, so many rhodopsin molecules mean that ‘dark noise’ events – the spontaneous activation of a rhodopsin molecule by temperature, not by photon absorption – can occur with significant frequency. The activation of rhodopsins by true photons have to be discriminated from the static noise of spontaneous rhodopsin activation by temperature. It is impossible to say whether a single given rhodopsin activation was triggered by a photon and not by thermal fluctuation. Multiple simultaneous rhodopsin activations are required to discriminate a true flash of light from a dark noise event.

Seeing single photons

WHAT IS THE SMALLEST NUMBER OF PHOTONS that a human observer can reliably detect? In the early history of the photon, Lorentz realized that a ‘just detectable’ flash of light delivered ~ 100 photons to the cornea. Most photons that reach the cornea do not reach the retina. The deeper question is: What is the threshold number of photons that is absorbed by photoreceptor cells to produce ‘seeing’? Hecht, Shlaer, and Pirenne (1942) did the classic experiment that established that single photons absorbed by single rod cells could be integrated into a perception of a flash of light. The threshold number of photons individually absorbed by a group of photoreceptor cells at the ‘threshold of seeing’ was 5–7 photons (Fig. 16).

Smelling single molecules

BACTERIAL CHEMOTAXIS is driven by smell. Bacteria are covered with chemoreceptors for molecules like amino acids that signify food. They use these chemoreceptors to assess surroundings and swim to favorable places. *E. coli* swims by rotating helical flagellar filaments attached to its $\sim 1 \mu\text{m}$ size body (Fig.). When all flagella rotate counterclockwise, as viewed from outside the cell, a bundle forms that pushes the cell forward at $\sim 25 \mu\text{m}/\text{s}$. These ‘runs’ typically last $\sim 1 \text{ s}$. Occasionally, one or more flagellar motors switch from CCW to CW rotation. This ends the run by disrupting the flagellar bundle. The cell stays in place and ‘tumbles’ until all flagella return to CCW rotation, and the cell starts a new run in a new direction.

Bacterial chemotaxis works by counting molecules. If the bacteria counts more attractant molecules over time during a run, it postpones the next tumble. Runs in favorable directions are thus longer than runs in unfavorable directions. Although each tumble randomly

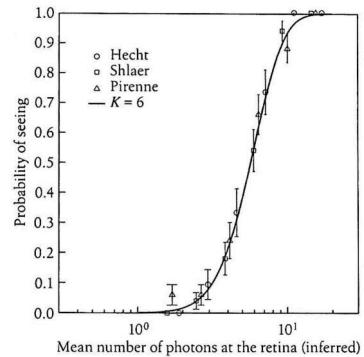


Figure 16: **The threshold of seeing** as the probability of seeing a flash plotted against the logarithm of the number of photons estimated to be absorbed by the retina at different flash strengths. Open symbols represent different measurements from three observers, the authors of the experiment.

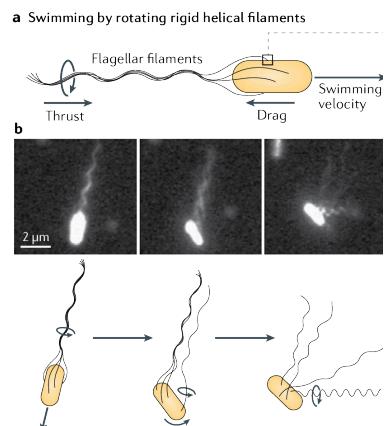


Figure 17: **Bacteria swim with a random walk** alternating periods of forward movement (counterclockwise rotation of flagellar filaments) with tumbles (clockwise rotation of one or more flagellar filaments).

reorients each run, a *biased random walk* ensues that inexorably drives the bacteria where it wants to go.

Noise in counting molecules during chemotaxis has a thermal origin, but in another sense than thermal activation of unbound receptors. Consider a bacterium-sized volume, $L=1 \mu\text{m}$ on each side. If the cell instantaneously counted all molecules inside this volume, it would count $\sim 600,000$ molecules if the mean concentration was 1 mM , 600 molecules if the concentration was $1 \mu\text{M}$, and 60 molecules if the concentration was 10^{-7} M . But molecules constantly move in and out of the measurement volume, and so the number in an instantaneous count will fluctuate in a way governed by **Poisson statistics** (Fig. 18). The standard deviation in the number of counted molecules will be the square root of the mean number of molecules. And so the relative error in estimating molecular concentration based on a single count within the $1 \mu\text{m}^3$ -sized measurement volume will be $600,000 \pm 800$ at 1 mM ($\sim 1.3\%$ error), 600 ± 24 at $1 \mu\text{M}$ ($\sim 4\%$ error), and 60 ± 8 molecules at 10^{-7} M ($\sim 13\%$ error). Berg and Purcell (1977) showed that bacteria do better than this by integrating measurements over time. During bacterial chemotaxis, measurable changes in bacterial behavior are caused by changing the occupancy of single receptors.

All biology is sensing molecules

Understanding the biophysics of chemoreception is not idiosyncratic to bacterial chemotaxis. Virtually every process in intracellular and intercellular biological signaling involves the diffusive movement of molecules. Intracellular biochemical pathways involve the binding of ligands and enzymes (Fig. 19). Synaptic communication between nerve cells involves the diffusion of neurotransmitters from a presynaptic cell to the receptors on a postsynaptic cell. Understanding the biophysics of chemoreception is fundamental to understanding life.

RECOMMENDED READING

- S M Block. Biophysical principles of sensory transduction. *Society of General Physiologists*, 47:1–17, 1992. ISSN 0094-7733 [Download paper](#)
- H C Berg. A physicist looks at bacterial chemotaxis. *Cold Spring Harbor Symposia on Quantitative Biology*, 53 Pt 1:1–9, 1988. ISSN 0091-7451 [Download paper](#)
- Chapter One. Gordon L Fain. *Sensory Transduction*. Sinauer Associates, Sunderland, Mass., 2003. ISBN 0878931716 [Download paper](#)

ADDITIONAL READING

- Josefina del Marmol, Mackenzie A. Yedlin, and Vanessa Ruta. The structural basis of odorant recognition in insect olfactory receptors. *Nature*, 597(7874):126–131, 2021. ISSN 0028-0836 [Download paper](#)

Updated: November 14, 2023

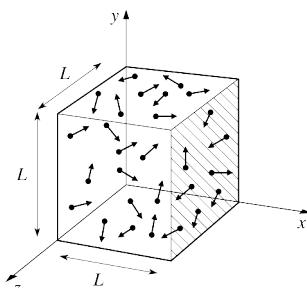


Figure 18: **Diffusion** is driven by the random and incessant motion of particles in and out of sampling volumes. At any point in time, the number of particles in a sampling volume fluctuates about a mean concentration owing to these ‘Brownian movements’.

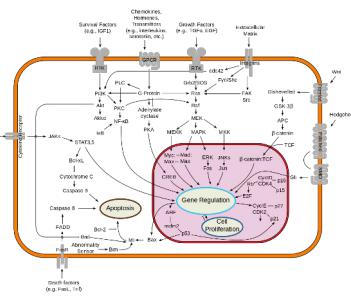


Figure 19: **Eukaryotic signal processing** is dominated by diffusion and the random binding and unbinding of molecules to receptors.

- S Hecht, S Shlaer, and M H Pirenne. Energy, quanta, and vision. *The Journal of General Physiology*, 25(6): 819–840, 1942. ISSN 0022-1295 [Download paper](#)

HOW BIOLOGISTS STUDY SENSORY SYSTEMS

Psychophysics

THE QUANTITATIVE STUDY of sensory perception, *psychophysics*, began with 19th century experimental psychology. Conscious human responses to well-controlled stimuli began to be measured and analyzed. With a human subject, a scientist can obtain meaningful perceptual measurements simply by asking questions. Did the subject detect a stimulus – did she or didn't she see the flash of dim light? Can the subject tell the difference (discriminate) between two stimuli – which is heavier, block A or block B? For the answers to such questions to be meaningful, one should use well-controlled stimuli and ask simple questions with clear answers. But even with rigorously designed experiments and simple questions, interpreting the results of psychophysical experiments can be subtle. *Psychometric curves*, mathematical functions that interrelate sensory response and stimulus intensity, must be carefully interpreted (Fig. 20).

Psychometric curves are shaped by the activities of underlying molecular and cellular sensors. An animal perceives a quantitative change in a given stimulus because some sensory molecules or cells undergo a quantitative change in activity. A stimulus might be perceived as stronger with an increase in stimulus intensity – brighter light, louder sound, higher odorant concentration – for different reasons. One possibility is that more sensors of a given type become activated – our low-light vision is mediated by one type of photoreceptor, when more rods within the retinal field of a dim object are activated, object appears brighter. Another possibility is that individual sensors become more strongly activated – the mechanosensory hair cells of the inner ear detect cochlear vibrations, the larger the vibration amplitude, the louder the sound. A third possibility is that different sensors with different intensity thresholds for the same stimulus become activated – our noses contain many different olfactory receptors with different activation thresholds for a given odor molecule, so the set of olfactory receptors that are activated by a given smell can contain information about the identity and intensity of odor molecules. Our goal as neuroscientists is to explain animal perception in terms of behavior with underlying mechanisms in terms of molecules and cells. To do this, sensory neuroscientists can study the relationships between psychometric curves that describe behavior and “dose-response” curves that describe molecular and cellular activity.

When studying sensory systems, it is useful to design experiments

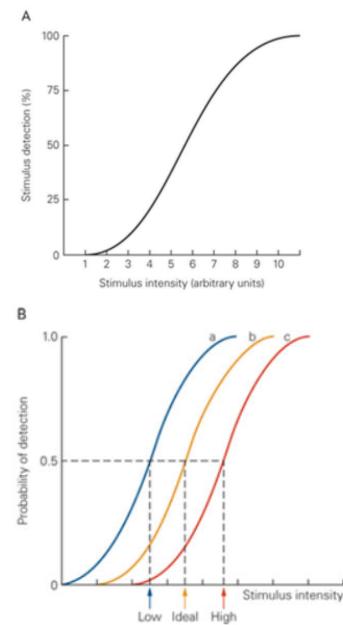


Figure 20: **The psychometric curve.** A. The psychometric function plots the percentage of stimuli detected by a human observer as a function of stimulus magnitude. Threshold is defined as the stimulus intensity detected on 50% of the trials. B. Detection and discrimination thresholds depend on the criteria used by individual subjects. Where an ‘ideal’ observer correctly detects the presence and absence of stimuli at the response threshold with equal probability (curve b), an observer who is told to respond to the slightest indication of a stimulus may report many false positives when no stimuli occur and has a lower response threshold (curve a). An observer who is told to respond only when very certain that a stimulus has occurred reports more hits than false positives and has a higher response thresholds (curve c). From *Principles of Neural Science*.

above *response thresholds* but below saturation, using stimulus intensities where the sensory system switches between ‘inactive’ and ‘active’ states with changes in stimulus intensity. One good reason is to study sensory systems where they are most relevant to natural behavior.

All biological systems *evolved* by natural selection. Sensory systems evolved to provide information about naturally occurring stimuli with certain characteristics in the real world. One learns less from sensory systems with stimuli far below threshold (i.e., too small to be detected) or far above saturation (i.e., too large to be discriminated).

Above response thresholds and below saturation, sensory systems have the most discriminatory power and provide the most information to the behaving organism. But in a regime where responses to changes in stimulus intensity are graded, sensory perception is not ‘all’ or ‘none’, binary performance that has the virtue of being easy for the experimenter to characterize on the basis of ‘yes’ and ‘no’ questions. Instead, sensory and behavioral responses will change in a gradual and typically probabilistic manner with changes in stimulus intensity, requiring more subtle (but also more informative!) analyses.

Let’s consider a concrete example. The rod cells of the human eye are specialized low-light (also called scotopic) vision. A good example of a naturally occurring stimulus that requires rod vision is the night sky. The dimmest visible stars will shower hundreds to thousands of photons per second on the dark-adapted pupil. Because most photons that reach the cornea will be scattered or absorbed before reaching the retina, rod cells must be able to detect small numbers of photons. The extraordinary sensitivity of the rod cell as photon counters was established at the turn of the last century. The astronomer Samuel Langley invented his *bolometer* in 1878, a device that allowed him to measure radiant energy with extraordinary precision – he could detect the infrared radiation from a cow at a quarter mile, quite a feat for a 19th century physicist. The bolometer allowed precise calibration of the radiant energy in visual stimuli that were just barely visible. With the realization that electromagnetic radiation is quantized as photons, Lorentz was able to estimate the “threshold for seeing” in terms of numbers of photons: just ~ 100 photons reaching the cornea could be reliably “seen”. Given photon losses due to reflection from the cornea or absorption within the eye, Lorentz’s estimate meant that single photons could trigger the activity of individual rod cells and induce a psychometric effect.

Indeed, the anatomy of the rod cell seems expressly engineered to catch photons (Fig. 21). A photon traveling along the long axis of a rod cell has to run a gamut of parallel disks loaded with visual pigment (rhodopsin molecules), but still only has a $\sim 2/3$ chance of being absorbed. Understanding the results of any photon-counting

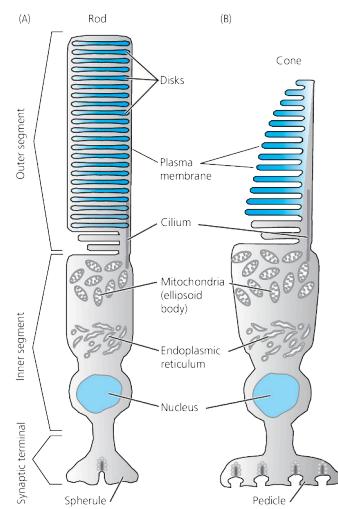


Figure 21: Vertebrate rods and cones. Principal structural features of vertebrate photoreceptors. (A) Rod. The outer segment is composed of disks detached from external plasma membrane. (B) Cone. The outer segment has membrane infoldings or lamellae instead of disks. The total area of sensory membrane is increased by these disk and lamellar structures, which are loaded with visual pigment. In a mammalian rod cell that is 0.025 cm long, roughly 2/3 of light is absorbed. (from Fain, Chapter 9).

experiment, whether at the level of animal perception or rod cell activity, requires thinking about probability and statistics. Not only is there intrinsic randomness in sensory response to a given stimulus, there is often intrinsic randomness in the delivery of any stimulus. Quantum mechanics does not allow a device that reliably deliver single photons – one at a time in a stream of identical pellets – to the eye of an observer or to a photoreceptor cell. The best we can do is understand the many sources of randomness that shape stimulus-response relationships, and develop probabilistic models that are consistent with behavioral and physiological measurements.

To make it easier to quantitatively analyze experimental results, studies are asked to make *comparative judgments* of a stimulus property, such as stimulus amplitude or frequency. To do this, we often use *two-alternative forced-choice protocol* with two observation intervals and a pair of stimuli. Subjects might be asked whether the second stimulus is stronger or weaker, higher or lower, faster or slower, same or different than the first stimulus. Or the subject might simply be asked whether the stimulus occurred in each interval. With binary questions and answers, there are only four outcomes – true positive, false positive, true negative, and false negative – making it easy to quantify and tabulate data (Fig. 22).

BAYES' THEOREM is naturally useful for interpreting psychophysical experiments in terms of *conditional probabilities*. Whether a stimulus is delivered in each trial has a well-defined probability determined by the experimenter. Whether a stimulus is detected has a probability that can be measured in the course of each experiment. Bayes' theorem is stated mathematically as the following:

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

where A and B are events and $P(B) \neq 0$. Here, the events are stimuli and responses.

- $P(A | B)$ is a *conditional probability*: the probability of event A given event B .
- $P(B | A)$ is also a conditional probability: the probability of event B given event A .
- $P(A)$ and $P(B)$ are the probabilities of events A and B respectively without conditions on other variables; they are also known as the marginal probability or prior probability.

In the experiment shown in Fig. 22, a stimulus event is true (A) if a red flash is delivered and false (\bar{A}) is a blue flash is delivered. In

		Response		Total stimuli
		Yes Red	No Blue	
Stimulus	Red	Hits (65)	Misses (35)	100
	Blue	False positives (20)	Correct rejections (80)	100
		85	115	200

Figure 22: **Two-alternative forced choice tasks.** The stimulus-response matrix for a stimulus detection task (yes-no) or a categorical judgment task (red-blue). Although there are two possible stimuli and two possible responses, the data represent conditional probabilities in which the experimenter controls the stimuli and measures the subject's responses. The numbers provide examples of behavioral data obtained from a strict observer who responds "yes" less often than the actual frequency of occurrence of the stimulus.

the same experiment, a response event is true (B) if a red flash is seen and false (\bar{B}) if a blue flash is seen. Red and blue flashes are delivered with equal probabilities so that $P(A) = P(\bar{A}) = 0.5$. However, a red flash is not always seen as red – of 100 red flashes, 65 are properly seen as red (hits or true positives) but 35 are falsely seen as blue (misses or false negatives). In terms of conditional probabilities, we write $P(B | A) = 0.65$ and $P(\bar{B} | A) = 0.35$. At the same time, a blue flash is not always seen as blue. Of 100 blue flashes, 20 are falsely seen as red (false positives) and 80 are properly seen as blue (true negatives). In terms of conditional probabilities, we write $P(B | \bar{A}) = 0.20$ and $P(\bar{B} | \bar{A}) = 0.80$. The probability of seeing a red flash is the sum of true positives and false positives: $P(B) = P(B | A)P(A) + P(B | \bar{A})P(\bar{A}) = 0.425$. The probability of seeing a blue flash is the sum of true negatives and false negatives: $P(\bar{B}) = P(\bar{B} | A) + P(\bar{B} | \bar{A})P(\bar{A}) = 0.575$.

The probability that a human subject sees a red or blue flash is an assessment of perceptual accuracy in the context of an experiment, but has less to do with perceptual accuracy in the real world. When we see something, we want to know the likelihood that the ‘something’ occurred. Using the data shown in Fig. 22, we can calculate the answer to this inverse question with the help of Bayes: what is the probability that a photon that is seen as red is actually red: $P(A | B)$.

Bayes' Theorem

BAYES' THEOREM MAY BE DERIVED from the definition of conditional probability:

$$P(A | B) = \frac{P(A \cap B)}{P(B)}, \text{ if } P(B) \neq 0$$

where $P(A \cap B)$ is the probability of both A and B being true. Similarly,

$$P(B | A) = \frac{P(B \cap A)}{P(A)}, \text{ if } P(A) \neq 0$$

Solving for $P(A \cap B)$ and substituting into the above expression for $P(A | B)$ yields Bayes' theorem:

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)} \text{ if } P(B) \neq 0$$

Interesting Bayesian results

Surprising results can emanate from Bayes' theorem. A classic example is testing for disease in a population, a pertinent modern example. Suppose, a particular test for whether someone has COVID-19 is 90% sensitive, meaning that the test is correctly positive result for 90% of infected persons (and incorrectly negative for the other 10%). The test is also 80% specific, meaning that it is correctly negative for 80% of uninfected persons (and incorrectly positive for the other 20%). Assuming 5% of the population is infected, what is the probability that a random person who tests positive is really infected? This conditional probability is: $P(\text{Infected} | \text{Positive})$.

$$\begin{aligned} P(\text{Infected} | \text{Positive}) &= \frac{P(\text{Positive} | \text{Infected})P(\text{Infected})}{P(\text{Positive})} \\ &= \frac{P(\text{Positive} | \text{Infected})P(\text{Infected})}{P(\text{Positive} | \text{Infected})P(\text{Infected}) + P(\text{Positive} | \text{Uninfected})P(\text{Uninfected})} \\ &= \frac{0.90 \times 0.05}{0.90 \times 0.05 + 0.20 \times 0.95} = \frac{0.045}{0.045 + 0.19} \approx 19\% \end{aligned}$$

In other words, even if someone tests positive, the probability that they are infected is only 19% - this is because in this group, only 5% of people are infected, and most positives are false positives coming from the remaining 95%.

THE MONTY HALL PROBLEM can also be solved by Bayes' Theorem.

Suppose you're on a game show, and you're given the choice of three doors: Behind one door is a car; behind the others, goats. You pick a door, say No. 1, and the host, who knows what's behind the doors, opens another door, say No. 3, which has a goat. He then says to you, "Do you want to pick door No. 2?" Is it to your advantage to switch your choice?

Stimulus variability and thresholds

How might the same stimulus give rise to a true positive or false negative? One reason is that any stimulus is rarely well-characterized as a binary variable, but is usually drawn from a probability distribution over a continuum of possible values. The observer applies a threshold to each stimulus to perform binary classification: ‘seen’ or ‘unseen’; ‘true’ or ‘false’; or ‘red’ or ‘blue’ as in Fig. 22. An ‘ideal’ observer would make the fewest mistakes, such as the sum of false positives and false negatives. But most observers might have different thresholds. An optimist (or lax observer) might be more likely to classify a stimulus as ‘true’ than the ideal observer, perhaps making more false positives but fewer false negatives. A pessimist (or strict observer) might be more likely to classify a stimulus as ‘false’, perhaps making fewer false positives but also more false negatives. The psychometric curves for different observers with different thresholds would accordingly shift along the stimulus axis (Fig. 20).

To model the effects of sensory threshold and stimulus variability on psychometric curves and error rates, we need to characterize the variability of the stimulus itself. Two stimuli might exist on a continuous range of measurable value, on the basis of which they must be discriminated. If the probability distributions of the measurable values for different stimuli exhibit overlap, it is not possible to make error-free judgment. Wherever a threshold can be set to delineate ‘true’ events from ‘false’ events, there is always a finite possibility for a ‘false’ event to exceed the threshold (leading to a false positive) or for a ‘true’ event to fall below threshold (leading to a false negative).

Stimulus variability around a mean measured value is most often characterized using Gaussian distributions. For example, whether a subject sees a ‘red’ flash (A , or true event) or blue flash (\bar{A} , or false event) is made on the basis of a continuously measured variable y . The mean value of y will be different for true and false events, $\langle y_A \rangle$ and $\langle y_{\bar{A}} \rangle$, respectively. We can then write the conditional probabilities of different values of y using Gaussian distributions with different means but (for simplicity) the same variance, as illustrated in (Fig. 23):

$$P(y | A) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(y - \langle y_A \rangle)^2}{2\sigma^2} \right]$$

$$P(y | \bar{A}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(y - \langle y_{\bar{A}} \rangle)^2}{2\sigma^2} \right]$$

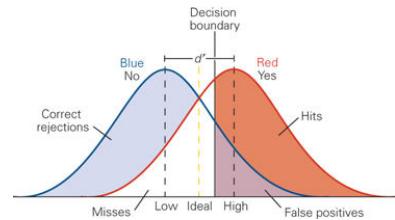


Figure 23: **Gaussian stimulus magnitudes.** Stimulus magnitudes can be represented by Gaussian curves with standard deviations that measure the fluctuation in sensations from trial to trial. The discriminability of a pair of stimuli is correlated with the distance between the two curves. When two stimuli are similar in magnitude, the two Gaussian curves overlap and no single criterion allows error-free responses. The frequency of true and false positives (and true and false negatives) is determined by the criteria used in the decision task. An ideal observer maximizes the number of correct responses and minimizes the total errors, setting the decision boundary at the intersection of the two curves. A strict observer minimizes the number of false positives but also reduces the total hits, setting the decision boundary to the right (solid line). A lax observer maximizes the number of hits but also increases the total false positives, setting the decision boundary to the left of the ideal subject.

Why do we use Gaussian distributions to characterize stimulus variability? A principled reason is the **Central Limit Theorem**. Suppose that a large sample of observations is obtained, each observation being randomly produced in a way that does not depend on the values of the other observations, and that the average (arithmetic mean) of the observed values is computed. If this procedure is performed many times, resulting in a collection of observed averages, the central limit theorem says that the probability distribution of these averages will converge to a normal distribution. An unprincipled reason is that the Gaussian distribution is simple-to-characterize with two just parameters (mean and variance) and an elegant mathematically function that facilitates calculation. Even when variability is non-Gaussian, the Gaussian is often used anyway to get useful approximate results without too much work or without too much danger.

Sensory Anatomy

FORM AND FUNCTION are complementary in biology, just like any in area of man-made engineering – designing electrical circuits, machinery, automobiles, or buildings – where structure dictates performance. The difference is that ‘engineering’ in biology was done by evolution and natural selection. The study of ‘form’ in biology occurs at many levels from the anatomy of body parts to cells to molecules.

The functional study of sensory systems began with gross anatomy. Long before magnifying glasses and microscopes were invented, analysis was limited to manual dissection and unaided human observation. For example, Galen’s early studies of the eye marveled at its unique and specialized structures. But without a proper understanding of physics of optics, the ‘model’ that Galen built was inherently flawed. Fascinated by the lens, an object like none other in the animal body, Galen made it the central structure in the vision mechanism (Fig. 24). When physical optics was discovered in the Renaissance, the properties of refraction became understood. On the practical side, optics led to useful inventions like the magnifying glass, telescope, and microscope. On the conceptual side, physical optics also clarified the role of the lens in eye. The lens is only a device for focuses images onto the retina at the back of the eye. The retina, a thin sheet of brain tissue, is the central mechanism for detecting images and relaying information to the brain.

Galen’s model is a significant distortion of eye anatomy, and it is hard to imagine how it held sway for a millennium. Scientists needed to understand the physics before they could properly see the parts of

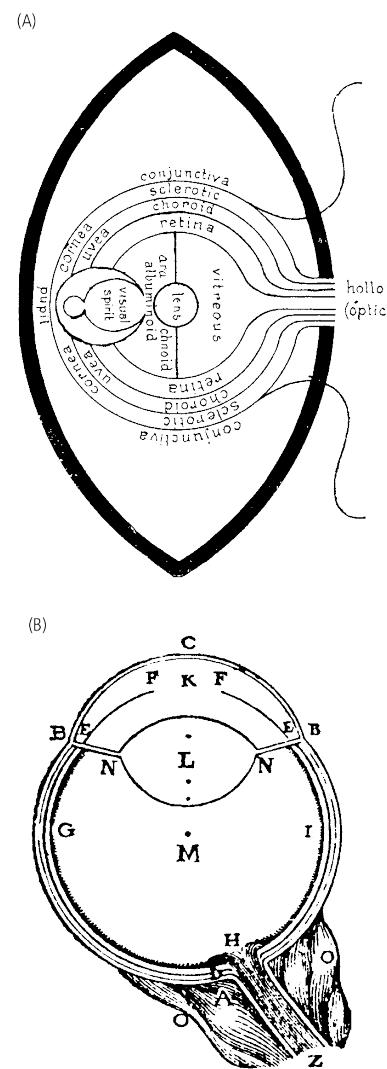


Figure 24: **Structure of the eye.** (A) Diagram of the eye from a ninth-century ad translation of Galen. (B) More anatomically correct diagram of cross-section of the eye made by Descartes. ABCB, Cornea and sclera; EF, iris (in actual fact closer to the lens than shown in Descartes’ diagram); K, aqueous humor; L, lens; EN, zonule fibers; M, vitreous humor; GHI, retina; H, optic nerve head; O, ocular muscles; and Z, optic nerve. From Fain, Chapter 1.

the eye in their proper places. Seeing is believing. Believing is also seeing. Descartes' model put the lens closer to its actual position and identified the muscles that change the shape of the lens during focal accommodation.

Microscopy

THE INVENTION OF THE OPTICAL MICROSCOPE led to anatomical investigations at the cellular level. Most of the sensory systems that we will study are in animals, vertebrate hearing, smell, and vision. The primary sensory receptors are specialized neurons for detecting sound, scent, and photons. The basic conceptual framework of neural circuit organization was largely established by Santiago Ramón y Cajal (1852-1934), the Spanish neuroscientist who was primarily a cellular-level anatomist (Fig. 25). Based on careful systematic analyses of sparsely labeled neurons in many brain tissues across animals and across developmental stages, Cajal identified basic principles about the organization of neural circuits in general and sensory circuits in particular. First, individual neurons interact with other neurons via contact or contiguity, synaptic contacts in today's language. Second, information flows with directionality through circuits, with dendrites and cell bodies on the input side and axons on the output side. These basic principles have been confirmed through anatomical investigations with higher spatial resolution using **electron microscopy**.

THE GOLGI STAINING METHOD was critical to Ramón y Cajal's success in resolving individual neurons. Neurons can be densely packed in brain tissue. Cell bodies are on the scale of micrometers. Synapses and nerve fiber thicknesses are on the scale of tens of nanometers. A theoretical limit to the resolution of any imaging system is the **Abbe diffraction limit**. Most optical systems, from light microscopes to electron microscopes to the human eye, suffer a spatial limit to resolution comparable to the wavelength due to diffraction. Even with the best compound microscopes, Cajal would not have been able to resolve neighboring structures closer than $\sim 300 - 500$ nm (corresponding to the wavelengths of visible light). But by individually labeling neurons that were separated by much greater distances, he could still perform exquisite anatomical reconstructions of each neuron, one cell at a time.

The Fly Eye

The wave nature of light is an important factor in determining the spatial limit of optical resolution, not just of microscopes but in vision. The "pixel sizes" of our retinas corresponds to the diameters of rod and cone photoreceptor cells, roughly 500 nm in the fovea where

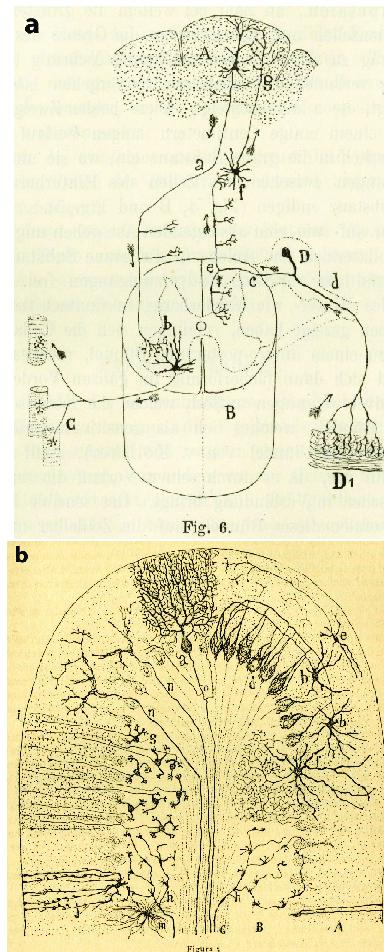


Figure 25: Cellular level anatomical analyses. (a) A nervous system-wide diagram of reflex and voluntary control of behavior by Cajal. Sensory information from the skin (D) is transmitted by dorsal root ganglion cells (d) to spinal cord (B) gray matter and to pyramidal neurons in the cerebral cortex (A), which in turn transmit impulses to motor neurons (b) in the spinal cord. For clarity, an interneuron between the spinal ending of (c) and an ipsilateral motoneuron are not shown. In this diagram, function (arrows) is predicted from structure. Neuron types in a gray matter region, the cerebellum. From Swanson and Lichtman (2016).

we obtain our highest spatial acuity. The visual angle of our highest spatial acuity, corresponding to the smallest letters on the eye chart, is about 1 minute of arc, spanning about 1-2 cone diameters in the fovea.

Insects have compound eyes, where each “pixel” of the field of view is provided by an ommatidium that points in a different direction, each gathering visual information from one angular field of view (Fig. 26). Feynman (the physicist) and Barlow (the neuroscientist) thought deeply about the anatomy of the compound eye. Take the insect eye to be a sphere of radius r divided into ommatidia (Fig. 27).

The larger the diameter of the ommatidium, the less angular resolution. This geometrical estimate of angular resolution is simply:

$$\Delta\theta_g = \frac{\delta}{r}$$

Shrinking δ would increase visual acuity, but diffraction puts a limit to how small δ can be. Light traveling through a thin slit will diffract. The thinner the slit, the wider the central peak of the diffracted wavefront (Fig. 28). Thus, if δ is too small, the ommatidium will “see” light at angles far from its axis:

$$\Delta\theta_d = \frac{\lambda}{\delta}$$

If δ is too large, enough ommatidia will see light because of poor angular resolution. So we adjust the δ to minimize poor angular resolution at excessively small or large values. We add the two effects and find the δ where the sum is minimized:

$$\frac{d(\Delta\theta_d + \Delta\theta_g)}{d\delta} = -\frac{\lambda}{\delta^2} + \frac{1}{r} = 0$$

This gives us an estimate for the optimum size of the ommatidium (Fig. 29):

$$\delta = \sqrt{\lambda r}$$

The diameter of the ommatidium should increase with the square root of the size of the eye. Barlow tested this idea by measuring insect eyes from 27 species of Hymenoptera (sawflies, wasps, bees, and ants), and thus experimentally verified his prediction.

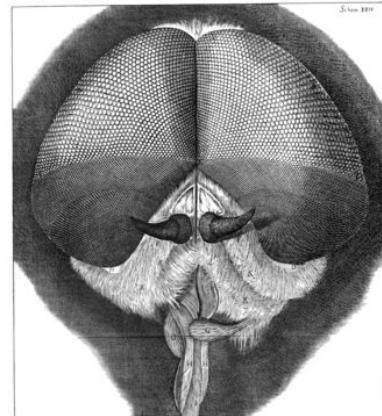


Figure 26: The compound eye.



Figure 27: Schematic view of packing of ommatidia in insect eye.

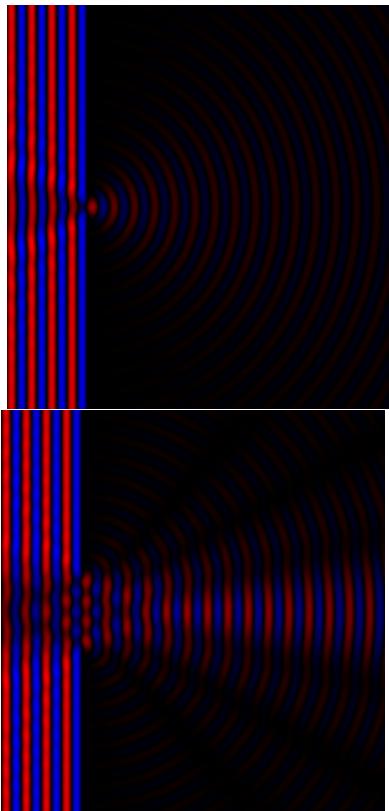


Figure 28: Diffraction.

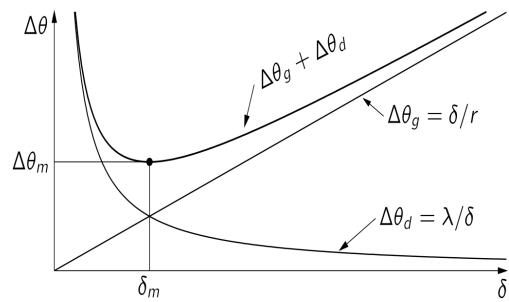


Figure 29: Optimum ommatidia diameters

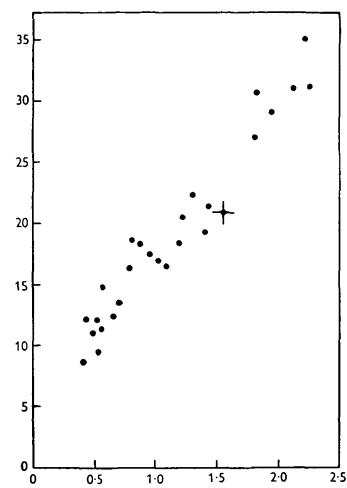


Figure 30: Actual ommatidia diameters
The ordinate axis is ommatidia diameter in micrometers. The abscissa is the square root of eye in millimeters.

ELECTRON MICROSCOPES have greater spatial resolution because the electron wavelength, defined by quantum mechanics by the de Broglie equation:

$$\lambda = \frac{h}{p}$$

where h is Planck's constant and $p = mv$ is momentum, is so much smaller in scanning or transmission electron microscopes. In typical microscopes, electron velocities reach 20%-70% the speed of light. The electron wavelength reaches ~ 12 picometers in a 10 kV SEM and ~ 2 picometers in a 200 kV TEM.

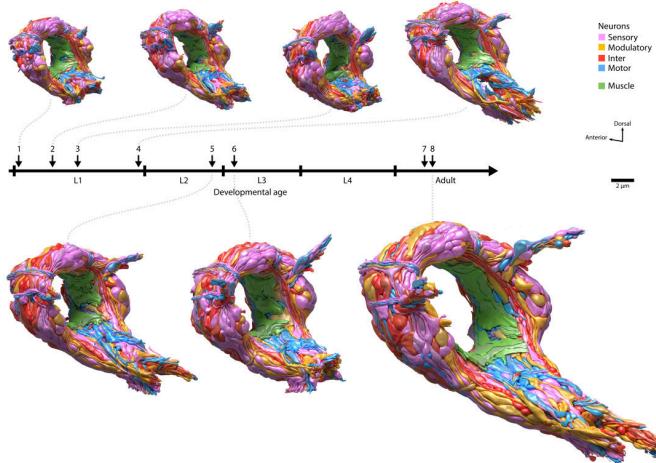


Figure 31: *C. elegans* *C. elegans* was the first animal to have its entire nervous system mapped by electron microscopy by John White and co-workers in the 1980s. More recently, eight *C. elegans* brains were reconstructed from birth to adulthood, mapping and comparing every chemical synapse and neuronal shape across whole-brain connectomes. From Witvliet et al. 2021

Electron microscopy imaging of brain sections has led to the characterization of synapses and neurons in diverse tissues. The main disadvantage of electron microscopy is that it requires ultrathin ($1\text{ }\mu\text{m}$) histological sections, which means that a single section is never adequate to generate full structural analysis of any cell. To overcome the problem of thin sections, scientists have used serial sectioning, in which each brain section is part of a sequence that transects a volume into hundreds or even tens of thousands of sections. Tracing objects from one section to the next reconstructs the geometry of the neurons and can also be used to identify the sites and cellular participants of each synapse. Connectomics has been used to reconstruct the entire nervous systems of small animals like nematodes and fruit flies (Fig. 31), and small parts of the nervous system of vertebrates (Fig. 32).

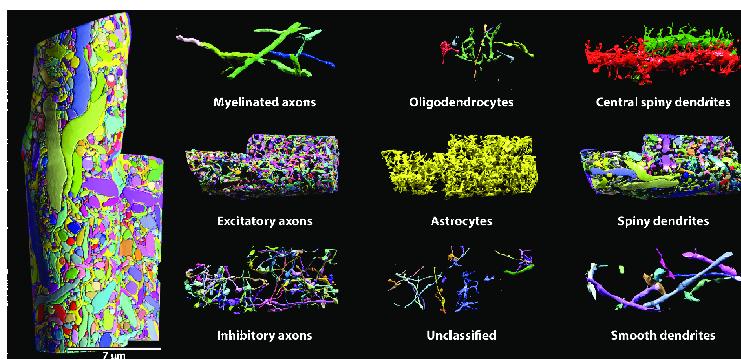


Figure 32: Reconstructions using modern, serial electron-microscopy connectomics approaches. On the far left is a fully reconstructed volume surrounding two apical dendrite segments of layer 5 pyramidal neurons in somatosensory cortex. Each colored object is a separate neuronal or glial cell process or extension. These processes are categorized in the images to the right. From Swanson and Lichtman, 2016.

ANATOMICAL ANALYSIS AT THE LEVEL OF INDIVIDUAL SENSORY RECEPTORS requires atomic resolution that greatly exceeds what is used for connectomics. X-ray crystallography has long been used to determine atomic-level structures of biomolecules. Crystallography only works for molecules that can be crystallized. Most sensory receptors are embedded in cell membranes, making them hard to crystallize in their native environments. Electron microscopy has atomic resolution, in principle, but sample degradation and poor contrast has long made it difficult to apply to molecules. Technical breakthroughs have led to cryo-electron microscopy (cryo-eM), which has allowed scientists to image many biomolecules for the first time, including membrane-bound, multimolecular sensory receptors. Cryo-EM, combined with the computer-assisted approach of obtaining numerous images of the molecule of interest in different orientations, followed by algorithmic reconstruction, enables 3D structures of biomolecules at Angstrom-level resolution, such as that of Piezo, a major mechanosensory channel in animals (Fig. 33).

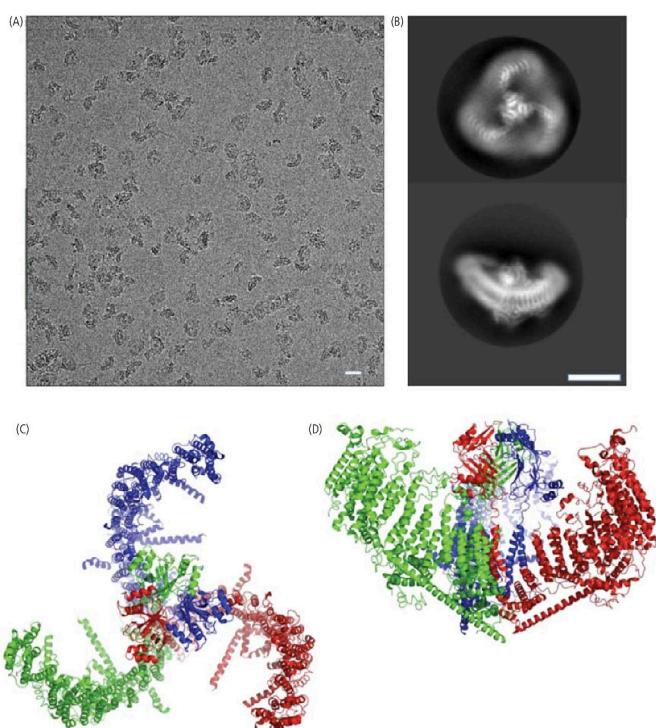


Figure 33: CryoEM of Piezo1 The gene for the Piezo1 protein was expressed in a cell line and purified. The protein was suspended on an electron-microscope grid and rapidly frozen by plunging the grid first into liquid ethane and then into liquid nitrogen. (A) Representative raw micrograph of the protein on the grid; scale bar is 20 nM. (B) Protein images like those in (A) were separated into groups according to their orientation, first manually to produce templates and then automatically under computer control. Images in each class were averaged, and representative averaged classes are shown viewed from the top (upper image) and side (lower image); scale bar is 100 Å. (C, D) Atomic model of the trimeric channel at an overall resolution of 3.7 Å shown as a ribbon diagram, viewed from the top (C) and side (D). The three subunits have been given different colors. From Fain, Chapter 1.

Sensory Physiology

ANIMAL SUBJECTS can be asked whether they perceived a stimulus, whether directly (if a human) or by experimental design (quantitative analysis of behavioral responses in monkey, mouse, or smaller organisms). Knowing whether internal mechanisms – molecules, neurons, or brain circuits – perceived a stimulus requires *physiological analysis*. We will mostly discuss perception in animals, where sensory mechanisms are located in neurons. Neuronal activity can be measured in terms of electrical signals that travel along their ‘wires’ in a way that is mediated and amplified by ion channels in their cell membranes. Sensory cells have specialized ion channels that are directly (in the case of ionotropic receptors) or indirectly (in the case of metabotropic receptors) activated by the environmental stimulus itself.

E. D. Adrian recorded some of the first *extracellular recordings* from sensory neurons by placing the axons of touch receptor cells near wire electrodes. Skin pressure changed the frequency of *action potentials* in different ways depending on the stimulus. Thus, action potential firing patterns are a mechanism for encoding and communicating sensory information to the brain. The first photosensory responses were made by Hartline in the horseshoe crab, where action potential patterns also depended on the intensity and duration of light stimulation. Patterns of neuronal activity in sensory neurons encode the incoming stimulus.

THE BRAIN SOLVES AN INVERSE PROBLEM. It must reconstruct a past incoming stimulus based on the present pattern of neuronal activity. As before, when we considered the statistics of stimulus and response at the level of organism behavior, Bayesian analysis is needed to think about the statistics of stimulus and response at the level of neuronal activity. Because neuron activity patterns can be intrinsically randomness, a physiologist must measure the *probability* of a specific neuronal response that might be conditioned on a specific stimulus: $P(\text{response} \mid \text{stimulus})$. The neuronal activity pattern is what the brain “knows” about the external environment. The animal must then perform a Bayesian inference, estimating the probability that a specific stimulus occurred in the external world conditioned on a specific neuronal response in its brain: $P(\text{stimulus} \mid \text{response})$.

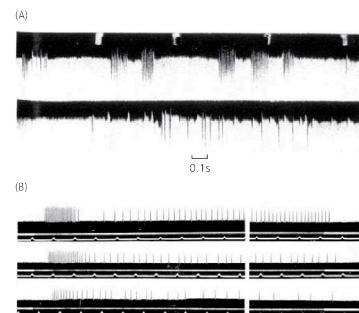
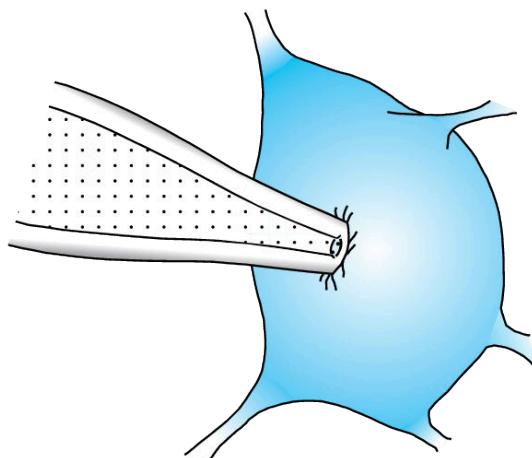


Figure 34: **Early electrical recordings of sensory responses.** (A) Action potentials recorded from single axons dissected from the cutaneous nerve of a frog. (B) Action potentials from the lateral eye of the horseshoe crab *Limulus*. Each trace gives the response to a different light intensity, which was systematically increased by an additional factor of ten from dimmest (bottom) to brightest (top). From Fain, Chapter One.

INTRACELLULAR RECORDINGS made it possible to record the detailed electrical responses of single neurons (Fig. 81). Electrodes are inserted into glass tubes that are melted and pulled to a fine point and filled with salt solution. Sharp electrode recordings can be made by penetrating individual cells. *Patch electrode recordings* can be made by polishing the tip of the glass electrode to be very smooth, such that when it is pressed against a cell membrane and a slight suction is applied, a very tight seal is formed, sometimes called a gigaseal with typical electrical resistances of $10\text{-}100\text{ G}\Omega$. High seal resistance ensures that most of the electrical current in each recording travels through the cell membrane that is being tested, not leaking through the seal.

(A)



(B)

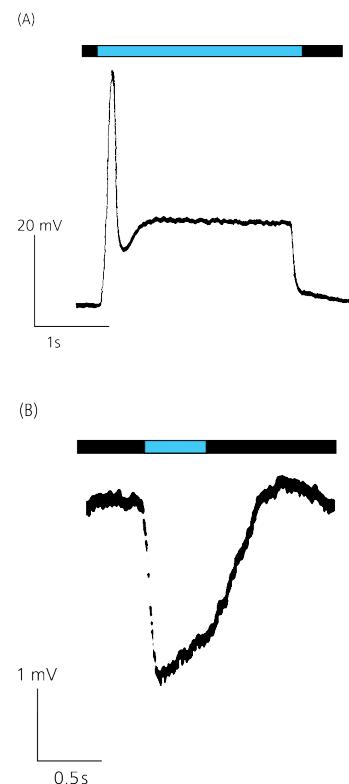
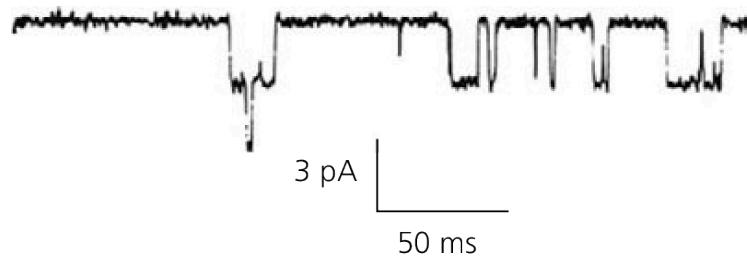


Figure 35: Intracellular recordings from sensory receptors. Bars above recordings show timing and duration of light flashes. (A) Depolarizing voltage response from photoreceptor of *Limulus* ventral eye. (B) Hyperpolarizing voltage response from photoreceptor (cone) of a fish. This is the first published recording of the response of a vertebrate photoreceptor. From Fain, Chapter One.

Figure 36: Patch-clamp recording from single channels. (A) The tip of a patch pipette is pushed against the cell body of a cell and slight suction is applied to form a seal. (B) Single-channel currents recorded from muscle acetylcholine receptors. The pipette contained $0.3\text{ }\mu\text{M}$ acetylcholine. Downward deflections indicate channel opening. At least two channels were present in this membrane patch. From Fain, Chapter One.

Molecular Biology

THE REVOLUTION IN MOLECULAR BIOLOGY deepened the analysis of sensory mechanisms (Fig. 37). Most sensory receptors are integral membrane proteins. Most known receptors were identified either by protein purification and sequencing, or by genetic analysis (finding a mutant that lacks a sensory modality, and working out the missing gene that was responsible for phenotype). Modern techniques to find sensory receptors include single-cell transcriptional and whole-genome sequencing and analysis. The genetic sequence of a putative receptor has characteristics that can betray its identity. Integral membrane proteins must have extensive sequences within the hydrophobic interior of the lipid bilayer. From amino acid sequences, which amino acids lie within the membrane and which face the cytoplasmic or extracellular solution can be inferred. Some amino acids (such as valine and isoleucine) are hydrophobic. Some amino acids (such as aspartate and lysine) are hydrophilic. Hydropathy analysis can be used to estimate the rough structure of putative receptors.

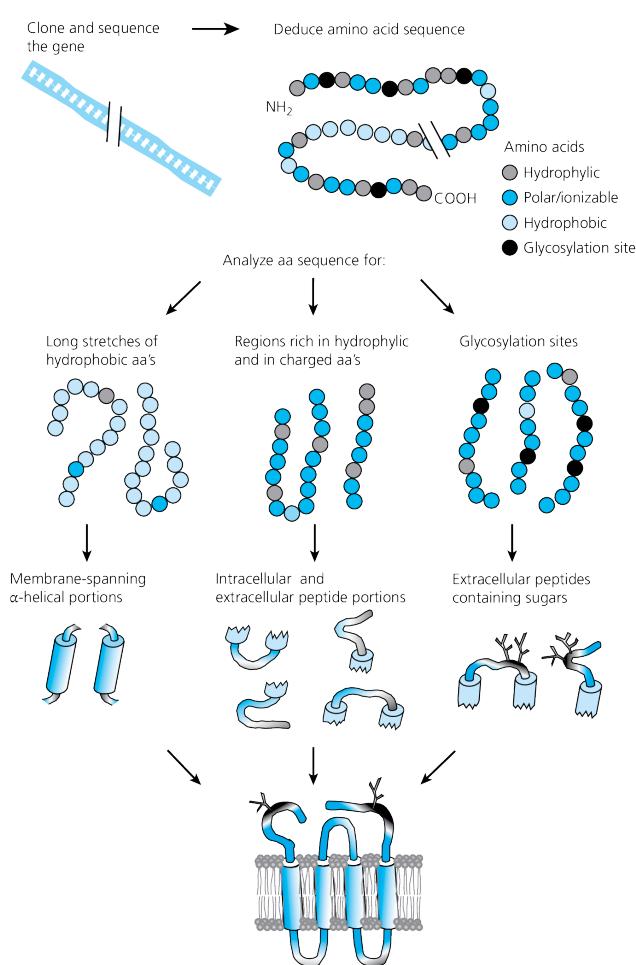


Figure 37: Analysis of hydropathy and the folding of membrane proteins. The amino acid sequence of a membrane protein can be used to make inferences about protein structure. From Fain Chapter One.

RECOMMENDED READING

- Eric R. Kandel, James H. Schwartz, and Thomas M. Jessell, editors. *Principles of Neural Science*. Elsevier, New York, third edition, 1991 [Download](#)
- Chapter One. Gordon L Fain. *Sensory Transduction*. Sinauer Associates, Sunderland, Mass., 2003. ISBN 0878931716 [Download paper](#)
- Larry W. Swanson and Jeff W. Lichtman. From Cajal to Connectome and Beyond. *Annual Review of Neuroscience*, 39(1):197–216, 2016 [Download](#)

ADDITIONAL READING

- Barlow's explanation of ommatidial size across insects.
 - H. B. Barlow. The size of ommatidia in apposition eyes. *Journal of Experimental Biology*, 29(4):667–674, 1952. ISSN 0022-0949 [Download](#)
- This paper describes the modern use of high-throughput connectomics to measure the developmental dynamics of an entire animal brain from birth to adulthood.
 - Daniel Witvliet, Ben Mulcahy, James K. Mitchell, Yaron Meirovitch, Daniel R. Berger, Yuelong Wu, Yufang Liu, Wan Xian Koh, Rajeev Parvathala, Douglas Holmyard, Richard L. Schalek, Nir Shavit, Andrew D. Chisholm, Jeff W. Lichtman, Aravinthan D. T. Samuel, and Mei Zhen. Connectomes across development reveal principles of brain maturation. *Nature*, 596(7871):257–261, 2021. ISSN 0028-0836 [Download](#)
- The discovery of the Piezo mechanosensory channels involved a remarkable integration of modern techniques in molecular biology and electrophysiology to yield one of the most elusive sets of proteins in sensory perception.
 - Bertrand Coste, Jayanti Mathur, Manuela Schmidt, Taryn J Earley, Sanjeev Ranade, Matt J Petrus, Adrienne E Dubin, and Ardem Patapoutian. Piezo1 and piezo2 are essential components of distinct mechanically activated cation channels. *Science*, 330(6000):55–60, 2010 [Download](#)

SOME STATISTICAL MECHANICS

Rhodopsin

RHODOPSIN, also called visual purple, is the light-sensitive receptor protein in the retina (Fig. 39). Rhodopsin is also a membrane-bound photoswitchable G-protein-coupled receptor (GPCR) with both a protein component and covalently-bound cofactor, a photoreactive chromophore called *retinal*. Isomerization of 11-cis-retinal into all-trans-retinal by light triggers a conformational change that causes rhodopsin to activate another G protein called *transducin*. Transducin activation begins a signal transduction cascade that eventually modulates cyclic guanosine monophosphate (cGMP) levels. Changing cGMP levels change the electrical activity of the rod cell through cGMP-gated ion channels in the cell membrane of the photoreceptor cell. Changes in the electrical activity of the photoreceptor cell membrane lead to changes in their chemical synaptic outputs that are communicated to the rest of the retina as visual events.

THE ENERGY of one blue-green photon ($\lambda=500 \text{ nm}$) is $4 \times 10^{-19} \text{ J}$. Because rhodopsin is in thermal equilibrium with its environment ($\sim 25^\circ\text{C}$), the chromophore will have, on average, an amount of thermal energy around $4 \times 10^{-21} \text{ J}$. The energy of a blue photon is much higher than thermal energy. On average, the thermal energy fluctuations of the chromophore can be expected to be much lower than a threshold needed for its activation. But how frequently will the chromophore exceed the threshold of activation by thermal fluctuations alone? Answering this question requires knowing about probability distribution of energy levels of an object in thermal equilibrium, not just its mean thermal energy.

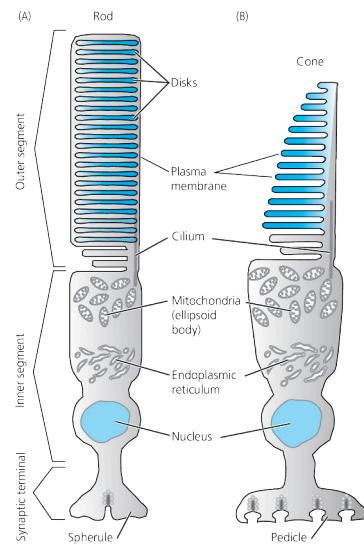


Figure 39: Rods and Cones.

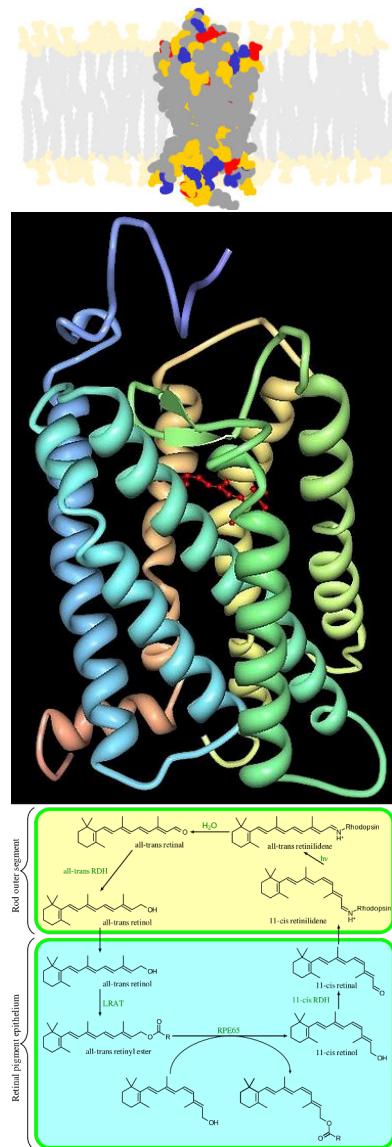


Figure 38: Rhodopsin. Three dimensional structure of bovine rhodopsin, a membrane bound protein. The chromophore, retinal, is embedded within the protein. Rhodopsin is embedded in detached discs in the outer segment of the rod photoreceptor cell and embedded in infolded lamellae in cones. Absorption of light energy, $h\nu$, causes a conformational change – 11-cis-retinal becomes all-trans-retinal – which thereby alters protein structure to activate G-protein-coupled signal transduction. Biochemistry in the rod cell restores the cis configuration.

The Boltzmann Distribution

IN STATISTICAL MECHANICS, a Boltzmann distribution is a probability distribution that gives the probability that a system will be in a certain state as a function of state energy and temperature:

$$p_i \propto e^{-\varepsilon_i/(kT)}$$

where p_i is the probability of the system being in state i , ε_i is the energy of that state, and a constant kT of the Boltzmann constant and temperature T . The symbol \propto denotes proportionality.

System can have broad meaning; an macroscopic ensemble of components (a molecule) or a single atom. What matters is that the system has a *degree of freedom* that allows it to enter different measurable states, and that energy is freely exchanged into and out of the system. An exponential distribution means that states with lower energy have a higher probability of being occupied (Fig. 40).

THE RATIO OF PROBABILITIES of two states is the Boltzmann factor and depends on their energy difference:

$$\frac{p_i}{p_j} = e^{(\varepsilon_j - \varepsilon_i)/(kT)}$$

If a system has M possible states, and the sum of the probabilities of being in each state is normalized, $\sum_{j=1}^M p_i = 1$, one can compute the probability of being in each state:

$$p_i = \frac{1}{Q} e^{-\varepsilon_i/(kT)} = \frac{e^{-\varepsilon_i/(kT)}}{\sum_{j=1}^M e^{-\varepsilon_j/(kT)}}$$

THE EXPONENTIAL ATMOSPHERE is a classic example of a continuous Boltzmann distribution (Fig. 41). The altitude of gas molecules in Earth's atmosphere is one spatial degree of freedom. The altitude of each gas molecule is associated with a gravitational potential energy: $E = mgh$. Assuming a uniform atmospheric temperature T , the probability that a molecule is at height h is:

$$p(h) \propto e^{-mgh/kT}$$

For this probability density to be properly normalized:

$$p(h) = \frac{e^{-mgh/kT}}{\int_{h=0}^{\infty} e^{-mgh/kT} dh}$$

The expectation value for the altitude of a gas molecule depends on its mass: $\langle h \rangle = \frac{kT}{mg}$.

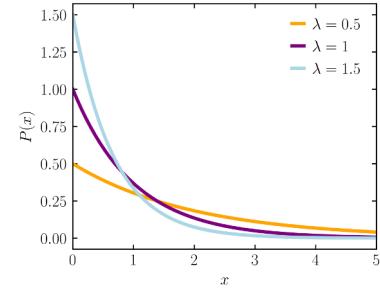


Figure 40: **Exponential Distributions.** Plot of the probability density function of the exponential distribution ($P(x) = \frac{e^{-\lambda x}}{\lambda}$) for rates $\lambda = 0.5, 1$ or 1.5 .

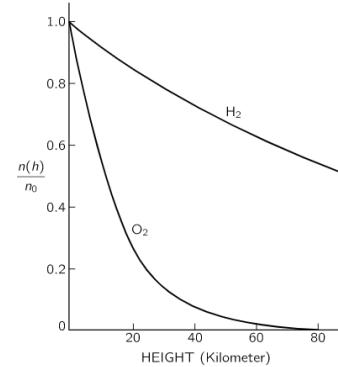


Figure 41: **Exponential atmosphere.** From Feynman's *Lectures in Physics*. The mass of one molecule of O_2 is 5×10^{-23} g. The mass of one molecule of H_2 is 3×10^{-24} g. Does Feynman's drawing make sense?

THE BOLTZMANN DISTRIBUTION governs the likelihood that a sensory receptor like rhodopsin has different states with different thermal energies. We describe the different conformational states of a molecule using a “reaction coordinate”, and assign an energy to each point along the reaction coordinate (Fig. 74).

Using Boltzmann factors, we can infer the relative likelihood of being in different states. Going from one stable state A to another state B might require transit through an activation state. The energy of this activation state creates a barrier that slows the reaction. At equilibrium, the relative likelihood of any two states is:

$$\frac{p_A}{p_B} = e^{-\frac{\Delta E}{kT}}$$

To estimate the reaction rate from A to B , we need the relative likelihood of being at the top of the activation barrier compared to state A . This estimates the fraction of particles in state A that are able to move over the barrier thanks to fluctuations in thermal energy:

$$k_{A \rightarrow B} \propto e^{-\frac{E_{act}}{k_B T}}$$

The Boltzmann distribution gives us insight into the fluctuating energies of sensory receptors before the arrival of stimulus energy, as well as the speed of the reactions that characterize sensory transduction.

Deriving the Boltzmann distribution by counting

Say that you have N particles that are given a total amount of energy E . These particles are able to freely exchange energy among them, but are constrained to quantal energy levels. How many particles can you expect to find at each energy level? In other words, what is the probability of a particle having a certain amount of energy? To find the distribution of the particles over their possible energy states, we enumerate the states ($s = 1, 2, \dots$), associate each state with a discrete energy ($\varepsilon_1, \varepsilon_2, \dots$), and assign a number of particles to each state (n_1, n_2, \dots). The Boltzmann distribution should tell us how many particles, n_s , of the N total particles that we can expect to find in the s state with energy ε_s .

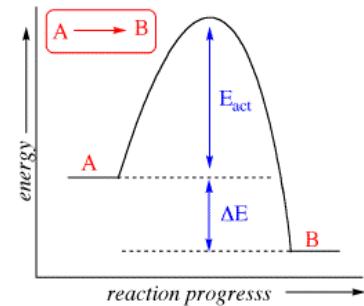


Figure 42: Reaction coordinate. From <http://butane.chem.uiuc.edu/pshapley/genchem2>

State	Energy	Number
1	ε_1	n_1
2	ε_2	n_2
3	ε_3	n_3
.	.	.
.	.	.
s	ε_s	n_s
.	.	.
.	.	.

Mass Conservation:
 $\sum_s n_s = N$.

Energy Conservation:
 $\sum_s n_s \varepsilon_s = E$.

The probability of any particle being in a given state s is n_s/N . Thus, the average energy of each particle is:

$$\langle \epsilon \rangle = \frac{\sum_s n_s \epsilon_s}{\sum_s n_s}$$

What is the probability of observing a specific distribution with say n_1 particles in state 1, n_2 in state 2, and so on? Our central assumption is that every possible distinct arrangement that satisfies conservation laws of mass and energy is equally possible. The probability of observing a particular distribution of n_s is thus proportional to the number of distinct arrangements that can be achieved with the N particles. The number of distinct particle arrangements that corresponds to the same distribution of particles among energy levels is a measure of the probability of that distribution. The number of ways a given distribution can be formed is a combinatorial problem:

$$W = \frac{N!}{n_1! n_2! n_3! \dots}$$

To find the distribution where W is largest, we need to maximize W with respect to n_s and with respect to the conservation laws. To do this, we need a convenient analytical expression for the factorial problem and we need to use Lagrange Multipliers.

We choose to maximize $\log W$ subject to the constraints of conservation of energy and mass. Thus,

$$d \left(\log W - \alpha \sum_s n_s - \beta \sum_s \epsilon_s n_s \right) = 0$$

The α and β are the Lagrange multipliers. Varying with respect to n_s and incorporating Stirling's Formula gives:

$$-\sum_s d n_s (\log n_s + \alpha + \beta \epsilon_s) = 0$$

Because this must hold true for every δn_s , every term in the sum must vanish. The values of n_s which do this are

$$\log n_s + \alpha + \beta \epsilon_s = 0$$

We can thus conclude that the occupancy of state s depends exponentially on its energy:

$$n_s \propto e^{-\beta \epsilon_s}$$

Actually, we have shown that the exponential dependence of state occupancy on energy is true only at the most probable W . We have not shown that every other arrangement can be ignored. The probability

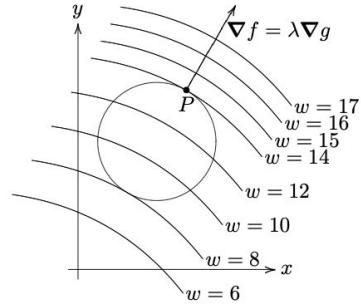


Figure 43: **Boltzmann's Tomb.** When you take Statistical Mechanics you will learn that $\log W = S/k_B$, the equation for entropy etched on his tomb.

Lagrange multipliers

Here, we sketch a geometric proof that builds intuition without worrying about rigor. For the function $w = f(x, y, z)$ constrained by $g(x, y, z) = c$, the maxima and minima are those points where ∇f is parallel to ∇g :

$$\nabla f - \lambda \nabla g = 0$$



For concreteness, we've drawn the constraint curve, $g(x, y) = c$, as a circle and some level curves for $w = f(x, y) = c$ with explicit (made up) values. Geometrically, we are looking for the point on the circle where w takes its maximum or minimum values.

Start at the level curve with $w = 17$, which has no points on the circle. Clearly, the maximum value of w on the constraint circle is less than 17. Move down the level curves until they first touch the circle when $w = 14$. Call the point where they first touch P. It is clear that P gives a local maximum for w on $g = c$, because if you move away from P in either direction on the circle you'll be on a level curve with a smaller value.

Since the circle is a level curve for g , we know ∇g is perpendicular to it. We also know ∇f is perpendicular to the level curve $w = 14$, since the curves themselves are tangent, these two gradients must be parallel. Q.E.D.

of alternative arrangements becomes negligible when N is very large, but this is hard to show and we will skip it.

We also have not shown that β , introduced here as a Lagrange multiplier, is related to temperature. For our purposes, we define $\beta = 1/k_B T$. When you take thermodynamics, you will learn that β behaves like $1/k_B T$ for various thermodynamic relationships and can have no other meaning. We are going to declare victory with the result that energy levels are exponentially distributed at thermal equilibrium.

Deriving the Boltzmann distribution Using Information Theory

CLAUDE SHANNON invented *information theory* starting with his own notion of the *entropy* of a probability distribution. Let X be a discrete random variable that can have different possible values, x . The probability density function, $p(x)$, is the likelihood of having different values.

SHANNON DEFINED HIS ENTROPY $H(x)$ of a discrete random variable X as:

$$H(X) = - \sum_x p(x) \log_2 p(x)$$

Say the random variable is the outcome of tossing a fair coin. X can either be heads ($p = 0.5$) or tails ($p = 0.5$). In this case, $H = 1$. Shannon's entropy is the amount of information in *bits* that you need to characterize the outcome of the toss.

What if the coin was biased and always came up heads? In this case, $H = 0$. You don't need any information to characterize the outcome of the coin toss (i.e., 0 bits), because the outcome is guaranteed. What if we vary p ? The entropy is minimum at $p = 0$ or $p = 1$ and maximum at $p = 0.5$ (Figure 44).

Roll a fair die. In this case, every outcome has $p = 1/6$ and $H = \log_2 6 \approx 2.5$. In Dungeons and Dragons, we have 8-sided die, and the number of bits needed to characterize one roll is 3.

The maximum entropy distribution is the one that minimizes the amount of prior information that is built into the distribution. If you know nothing about a probability distribution except the range of outcomes or the mean of the distribution, these constraints can be used to calculate the distribution that maximizes entropy using methods like Lagrange multipliers. This does not mean that the maximum entropy distribution is the right distribution, but it is a safe place to start. If the maximum entropy distribution isn't the distribution that correctly describes the system, then you lack prior

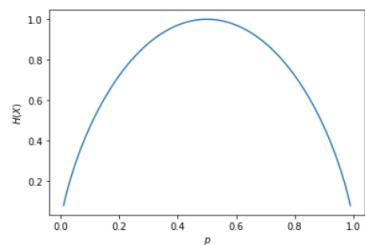


Figure 44: Entropy of a coin toss

and salient constraints about the system that would lead to different maximum entropy distribution.

Reconsider the problem of N particles exchanging a total amount of energy E among them. Each particle is constrained to discrete energy levels $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_s$

What is the probability distribution that governs the energies of each particle? This problem is fully described. We are not missing any salient facts needed to calculate a distribution. The correct probability distribution must be the one that maximizes entropy. We have phrased another Lagrange multiplier problem:

$$\delta \left(H - \alpha \sum_s p_s - \beta \sum_s \varepsilon_s p_s \right) = 0$$

We conclude that $p_s \propto e^{-\varepsilon_s \beta}$. If the distribution were anything else, we would have required additional salient constraints on the physical problem. But because we set up the problem without any other constraints, only the exponential distribution of energies, the Boltzmann distribution, is possible. *Q.E.D.*

REFERENCES

- This is Shannon's original paper where he invented Information Theory.
 - C. E. Shannon. A mathematical theory of communication. *Bell System Technical Journal*, 27(4):623–656, 1948. ISSN 0005-8580
[Download](#)
- Mehran Kardar. *Statistical physics of particles*. Cambridge University Press, Cambridge ; New York, 2007. ISBN 9780521873420

COLOR VISION AND THE DIMENSIONALITY OF PERCEPTION

VERTEBRATES HAVE TWO TYPES OF PHOTORECEPTORS, rods and cones. Our human eyes contain ~ 130 million rod cells for scotopic (low light) vision and ~ 7 million cone cells for color vision. All rod cells have one peak wavelength sensitivity, whereas our three cone sub-types are tuned to long, medium, and short wavelengths (Fig. 50). Thus, rod vision is monochromatic and cone vision is trichromatic. Trichromacy means that any color that we can see can be built from three different wavelengths of visible light such as red/green/blue or cyan/magenta/yellow. Thus, the stimulus space of human color vision has three dimensions. Let's make the idea of stimulus dimensionality more precise.

Human color vision provided early clues about the physical nature of light. Newton used a prism to separate white light into a spectrum of different colors (Fig. 46). He used a second prism to recombine spectrally-separated light to recreate white light. Somehow, light of one color (white) can be formed by summing light from different colors.

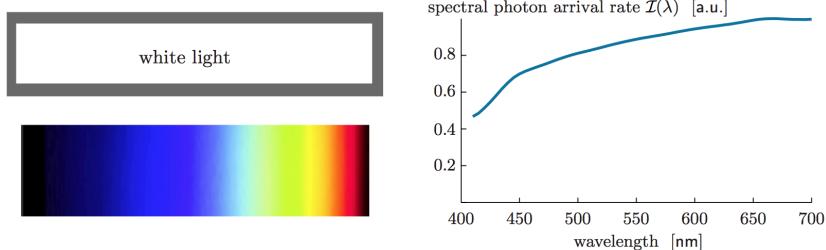


Figure 45: Rods and Cones.

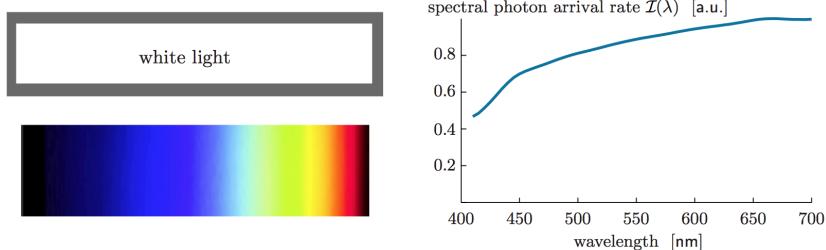


Figure 46: Sunlight consists of a broad spectrum of photon energies. When passed through a prism, it separates into a continuous distribution of colors. We can represent the spectrum by a photon arrival rate function that is nearly constant over the range of visible light.

In principle, electromagnetic radiation of one wavelength can carry information that can be completely separated from electromagnetic radiation of another wavelength. This is how frequency-modulated (FM) radio transmission works: each radio station sends information in a different electromagnetic ‘color’. Stimulus information that is carried on two different orthogonal axes (electromagnetic signals with different wavelengths) correspond to a 2-dimensional stimulus space. One can imagine a visual system that more fully exploits the human visual spectrum – from 400 nm (blue) to 700 nm (red) – by disentangling the visual information that is separately carried on a very large number of different wavelengths. But to read N -dimensional color information that is carried on N different wavelengths, the receiver needs at least N receptors that are separately tuned to different wavelengths. With only three cone sub-types, our

trichromatic retinas project all color information that might be carried by many different wavelengths onto a three-dimensional space. We discard any color information that might be contained in a high-dimensional stimulus space onto our much smaller three-dimensional internal representations of color. Most animals have fewer dimensions of color perception than we do (horses have two cone types, corresponding to dichromatic vision) but some more (mantis shrimps have 16 cone types, and might have the greatest power of color discrimination on the planet).

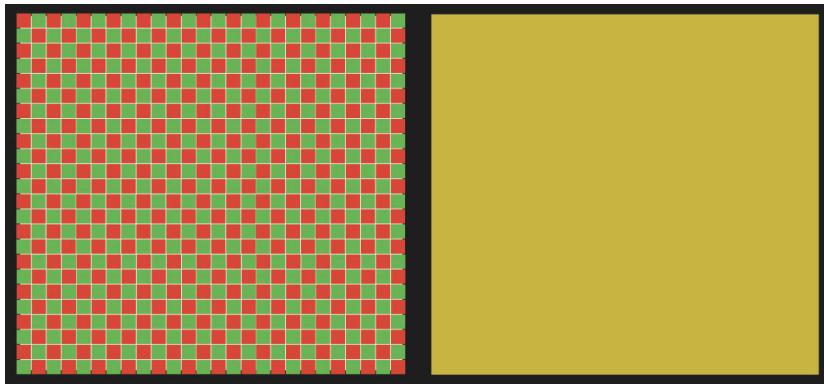


Figure 47: **Color illusion.** When viewed up close, the left box is seen to consist of small red and green squares. When viewed from afar, each photoreceptor receives light from both red and green squares, whose spectra merge. The resulting percept is closer to the color in the right box than red or green. Red + Green \sim Yellow.

We intuitively know that trichromatic vision discards spectral information. A mixture of spectrally pure red and green light looks yellow (Fig. 47). *Metamers* refer to physically different light spectra that lead to the same color perception. Long before biologists knew the molecular and cellular basis of color perception – the existence of three cone cells that express three different rhodopsin molecules – optical scientists used perceptual ‘color matching experiments’ to postulate trichromacy.

In color matching experiments, different sets of *basis lights* to illuminate a viewing screen (Fig. 48). Each basis light corresponds to a distinct spectrum of light wavelengths. Scientists discovered that only three basis lights, when varied in relative intensity, were sufficient to create a perceptual match to the full range of human color perception from red to violet. Moreover, three basis lights with different spectra can be chosen, as long as three spectra are different from one another and fall within the spectrum of visible light. The practical consequence of trichromatic vision when building computer or television screens is that the full gamut of color vision can be achieved with only three types of pixels (red, green, and blue). By mixing a relatively small palette of primary colors, a painter can endlessly create different hues. Chemical analysis of all of Vermeer’s paintings has revealed only twenty distinct pigments.

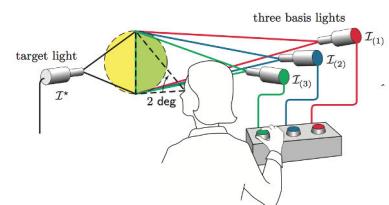


Figure 48: **Color matching.** A target light is projected onto the left half of a screen. A subject attempts to obtain a perceptual match by adjusting the intensities of three otherwise fixed basis lights that converge on the right half of the screen.

A quantitative model of color matching.

Suppose that we choose three basis lights with different spectra, meaning that they have different photon arrival rates as a function of wavelength: $\mathcal{I}_1(\lambda)$, $\mathcal{I}_2(\lambda)$, and $\mathcal{I}_3(\lambda)$. The units of these spectral light intensities are photons $s^{-1} nm^{-1}$, because light energy is distributed over time and over different wavelengths. To calculate the total photon arrival rate for each basis light (Φ_i), one must integrate over all wavelengths:

$$\Phi_i = \int \mathcal{I}_i(\lambda) d\lambda$$

In a color matching experiment, we must deliver these basis lights with adjustable intensities. To do this, we might fix the mean photon delivery rate of all basis lights to one baseline value, $\Phi = \Phi_1 = \Phi_2 = \Phi_3$, and use an adjustable scaling factor to individually change the intensities of each light: ζ_1 , ζ_2 , and ζ_3 . We give a human subject knobs to adjust the three ζ scaling factors to create a perceptual match to a monochromatic target (Fig. 48).

For simplicity, we can choose this monochromatic target to have a spectrum that is sharply peaked at one wavelength λ^* and fix the photon arrival rate at the same baseline as all basis lights, Φ . In the experiment corresponding to Fig. 49, the selected basis lights were also sharply peaked at single wavelengths – $\lambda_1=645$ nm, $\lambda_2=526$ nm, and $\lambda_3=444$ nm. In this experiment, a different combination of ζ_1 , ζ_2 , and ζ_3 were needed to create an experimentally observed perceptual match to the target light. Different target lights were chosen throughout the visible spectrum, each demanding a different set of ζ_1 , ζ_2 , and ζ_3 to create perceptual matches.

An interesting and subtle point for the experiment shown in Fig. 49 is the discovery of a range of target wavelengths that *cannot* be matched with three positive ζ values. There is a range of target wavelengths that require *negative* values for ζ_1 . Since negative light intensities are physically meaningless, how were experimental measurements made in this seemingly impossible range? When red basis light with negative intensity was needed to match the target wavelength, the experimenters simply added red basis light to the target light. Then, instead of matching the target light to the sum of red, green, and blue basis lights, the sum of the target and red basis light was matched to the sum of the green and blue basis lights. This works in a linear model of color perception, where perception follows from the simply addition of light at all wavelengths. If the stimulus is not too strong (such that cones are saturated), a linear model can be sufficient and is usually a good place to start.

In any case, every color in the spectrum of human perception can

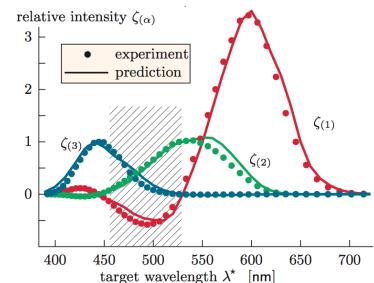


Figure 49: **Quantitative model of color matching.** Dots: experimental measurement of the relative intensities of basis lights needed to replicate target monochromatic wavelengths λ^* and results predicted from the measured spectral sensitivities of human cone cells using three basis lights with fixed wavelengths $\lambda_1=645$ nm, $\lambda_2=526$ nm, and $\lambda_3=444$ nm. It is impossible to create perceptual matches for monochromatic target lights in the hatched region that require negative values of long wavelength illumination (\mathcal{I}_1) in a linear model of color matching.

be represented with a stimulus space where every color is uniquely specified with only three free variables. A three-dimensional stimulus space saturates the full range of human color perception. .

The three-dimensional space of color representation.

The discovery that using three stimulus dimensions is sufficient to match any color led to the prediction that our internal representation of color also uses three dimensions. A three-dimensional internal representation of color would be formed at the beginning of ‘seeing’ if we used three different classes of photoreceptor for color – each with its own spectral sensitivity – to capture visual images. Now that we know that humans have three types of cone photoreceptor, we know that this is true, but let us build a quantitative model. First, we assign a spectral sensitivity function to all neurons of each photoreceptor class. This spectral sensitivity function is proportional to the probability that a photon of a given wavelength triggers the activity of the corresponding photoreceptor. The activity of each photoreceptor effectively counts all incoming and absorbed photons. But once a photon is counted, its color information (i.e., wavelength) is discarded. The rhodopsin molecule that absorbed the photon undergoes a binary change, from inactive to active, but no longer encodes the specific wavelength of the photon that elicited this change.

The three photoreceptor classes can be ordered by peak wavelength of their corresponding sensitivity functions (S, short, blue; M, medium, green; L, long, red): $\mathcal{S}_L(\lambda)$, $\mathcal{S}_M(\lambda)$, and $\mathcal{S}_S(\lambda)$. Although the color of incoming photons is lost in the activity of an individual photoreceptor, color information is preserved and encoded in the *relative* activities of the three photoreceptor types.

In each small wavelength range, $\Delta\lambda$, photons arrive at the retina with a mean rate $\mathcal{I}\Delta\lambda$. Photons stimulate the activity of each photoreceptor in proportion to the rate of photon arrival and its sensitivity function: $\mathcal{S}_i\mathcal{I}\Delta\lambda$. Collecting all wavelengths of light, the total number of photon absorptions by each photoreceptor is:

$$\begin{aligned}\beta_S &= \int d\lambda \mathcal{S}_S(\lambda) \mathcal{I} \\ \beta_M &= \int d\lambda \mathcal{S}_M(\lambda) \mathcal{I} \\ \beta_L &= \int d\lambda \mathcal{S}_L(\lambda) \mathcal{I}\end{aligned}$$

The brain concludes that two colors match when their spectra, or spectral arrival rates, cannot be distinguished – $\mathcal{I} \sim \mathcal{I}'$). The brain only has access to the activity of its three photoreceptor classes. Thus, the brain would conclude that two colors match when both colors

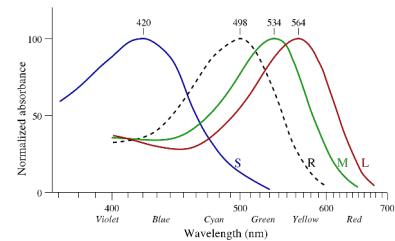


Figure 50: **Color vision. Normalized human photoreceptor absorbances for different wavelengths of light.**

evoke the same photon absorption rates in the three photoreceptors – $\beta_L, \beta_M, \beta_S$.

Given a target light spectrum (with intensity \mathcal{I}^*), how do we adjust the relative amounts of three basis lights with spectral intensities \mathcal{I}_i by setting values of ζ_i . What is the amount of each basis light in a sum that matches a target light when measured by photoreceptor activation? We need to solve the following equation for ζ_i :

$$\mathcal{I}^* \sim \zeta_1 \mathcal{I}_1 + \zeta_2 \mathcal{I}_2 + \zeta_3 \mathcal{I}_3$$

For the target light, the number of photon absorptions by the i th photoreceptor is by:

$$\begin{aligned}\beta_S^* &= \int d\lambda S_S(\lambda) \mathcal{I}^*(\lambda) \\ \beta_M^* &= \int d\lambda S_M(\lambda) \mathcal{I}^*(\lambda) \\ \beta_L^* &= \int d\lambda S_L(\lambda) \mathcal{I}^*(\lambda)\end{aligned}$$

Define $3 \times 3 = 9$ numbers, that describe the photon absorptions rates in each photoreceptor (S,M,L) evoked by each basis light ($i = 1, 2, 3$):

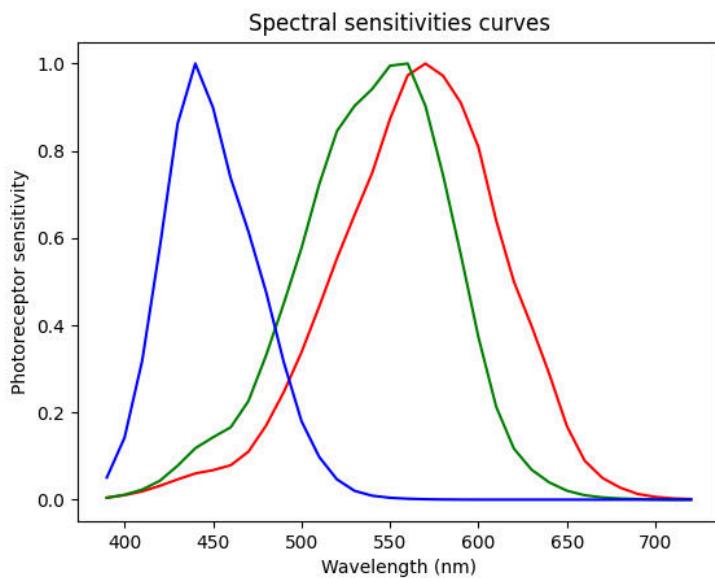
$$\begin{aligned}B_{i,S} &= \int d\lambda S_S(\lambda) \mathcal{I}_i(\lambda) \\ B_{i,M} &= \int d\lambda S_M(\lambda) \mathcal{I}_i(\lambda) \\ B_{i,L} &= \int d\lambda S_L(\lambda) \mathcal{I}_i(\lambda)\end{aligned}$$

To calculate the ζ_i that achieves color matching, we need to solve three linear equations that describe the activities of each photoreceptor in response to the target light and the to the summed basis lights. Rewriting the requirement for color matching in vector-matrix form :

$$\begin{bmatrix} \beta_S^* \\ \beta_M^* \\ \beta_L^* \end{bmatrix} = \begin{bmatrix} \zeta_1 B_{1,S} + \zeta_2 B_{2,S} + \zeta_3 B_{3,S} \\ \zeta_1 B_{1,M} + \zeta_2 B_{2,M} + \zeta_3 B_{3,M} \\ \zeta_1 B_{1,L} + \zeta_2 B_{2,L} + \zeta_3 B_{3,L} \end{bmatrix}$$

As homework, we will use the measured spectral sensitivities of human photoreceptors to calculate the basis lights needed to achieve color matching with different target lights.

Tabulated spectral sensitivities of photoreceptors.



[Download CSV file](#) corresponding to the measured spectral sensitivities of human photoreceptors.

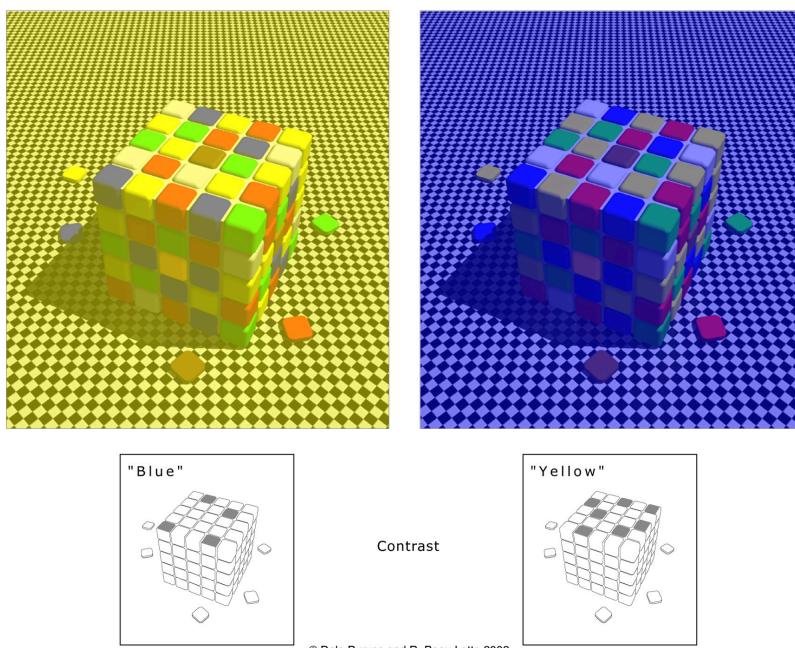
Wavelength	$\mathcal{S}_L(\lambda)$	$\mathcal{S}_M(\lambda)$	$\mathcal{S}_S(\lambda)$
390	0.00442	0.00440	0.0507
400	0.0105	0.0111	0.142
410	0.019	0.0231	0.319
420	0.0317	0.0434	0.580
430	0.0465	0.0778	0.862
440	0.0601	0.118	1
450	0.0677	0.143	0.899
460	0.0789	0.166	0.737
470	0.110	0.226	0.615
480	0.17	0.331	0.476
490	0.247	0.450	0.315
500	0.338	0.577	0.18
510	0.443	0.723	0.0979
520	0.553	0.845	0.0466
530	0.653	0.903	0.0202
540	0.75	0.942	0.00899
550	0.873	0.995	0.00438
560	0.973	1	0.00215
570	1	0.903	0.000969
580	0.972	0.744	0.000402
590	0.91	0.562	0.000157
600	0.808	0.372	0
610	0.64	0.213	0
620	0.500	0.117	0
630	0.398	0.0687	0
640	0.288	0.0395	0
650	0.168	0.0205	0
660	0.0893	0.0101	0
670	0.0499	0.00514	0
680	0.0276	0.00264	0
690	0.013	0.00133	0
700	0.00627	0.000626	0
710	0.00283	0	0
720	0.00134	0	0

Mysteries about color vision

In reality, color perception is much more than pixel-wise analysis of cone photoreceptor activity patterns. We perceive color by integrating information throughout the visual field. This is clearly demonstrated by various visual illusions about color.

One visual illusion is demonstrated when two targets with the same spectral intensities are surrounded by backgrounds with different spectral sensitivities. When this happens, the color of the target can appear different, even though the spectral composition of the light that is being absorbed by cone photoreceptors is exactly the same.

A more dramatic example of the context-dependence of color perception was devised when yellowish or bluish illumination was provided to a Rubik's cube tiled with different colors. Tiles that are spectrally gray when viewed in isolation can become blue when viewed with yellowish illumination or can become yellow when viewed with bluish illumination.



These examples suggest that measurements of human trichromatic color vision is strongly contingent on the experimental setup, involving the matching of large, homogeneous, context-free blocks of color. Color perception in the real world is a more subtle problem.

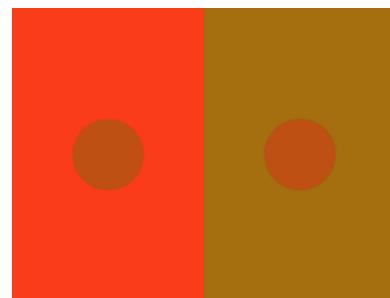


Figure 51: Spectrally identical patches can look differently colored when placed in spectrally different surrounds. The two central targets here are identical, as can be seen by masking out the surround.

Figure 52: Upper images show the cubes as if in yellowish (top left) or bluish (top right) illumination. The lower images show specific tiles of interest in the absence of these contexts. The yellow-looking tiles depicted as if under blue light and blue-looking tiles depicted as if under yellow light are actually a gray on their own.

REFERENCES

- Philip Charles Nelson. *From photon to neuron : light, imaging, vision*. Princeton University Press, Princeton, New Jersey, 2017. ISBN 9780691175188
- Dale Purves, R Beau Lotto, and Surajit Nundy. Why we see what we do. *American scientist*, 90(3):236–243, 2002. ISSN 0003-0996 [Download paper](#)

THE STATISTICS OF PHOTON ABSORPTION BY PHOTORECEPTORS

The anatomy of rod cells is specialized for low-light detection. Each rod cell is ~ 2 microns in diameter, but ~ 100 microns long. Incoming photons travel along the long axis, the better the chance of photon capture by visual pigment. Some animals have a shiny tissue layer behind the retina that acts as a retroreflector (Fig. 53). This *tapetum lucidum* reflects photons not captured by rods in a first pass, adding a chance of absorption in a second pass. Thus, the *tapetum lucidum* effectively doubles rod length without doubling visual pigment.

Starlight has photon fluxes of $<10^{-2}$ photons $\mu\text{m}^{-2} \text{ sec}^{-1}$. Bright sunlight has photon fluxes of $>10^8$ photons $\mu\text{m}^{-2} \text{ sec}^{-1}$. Cones mediate vision over the upper 7-8 log units of this range. Rods mediate vision at rates of photons per second. Rods are engineered to detect single photons. For rods to also *count* photons, they must reliably amplify the signal associated with each photon to be able to distinguish one photon from zero or two.

In photoreceptor cells, the visual pigments and signal transduction molecules that convert photon detection into electrical activity are in the outer segments. In rods, the visual pigments are rhodopsin molecules that are loaded into the membranes of packed disks (Fig. 54). Typical human rods have ~ 1000 disks. Each rod has about 10^5 rhodopsin molecules in each disk.

When a rod cell absorbs a photon, it creates and communicates a signal by changing synaptic output to downstream horizontal cells and bipolar cells in deeper retinal layers. Rods and cones lack the voltage-gated Na^+ needed to create fast all-or-none action potentials. Instead, photoreceptor cells have relatively slow, graded changes in membrane potential that modulate the rate of synaptic release. A biochemical signal transduction cascade connects changes in photon absorption to changes in membrane potential.

ROD PHOTORECEPTOR CELLS are designed to capture photons, but not *every* photon. The longer the rod, the more likely that an entering photon (at $x = 0$) is absorbed before exiting (at $x = L$). What is the probability, P , that a rod of length L captures a photon? This total probability is the sum over the probabilities that the photon is absorbed in each segment (between x and $x + dx$) along rod length.

Consider a thin segment of rod with thickness Δx . A photon that enters this segment has a small finite probability of being absorbed. This probability is proportional to both rhodopsin concentration, C , and an absorption cross-section, σ :

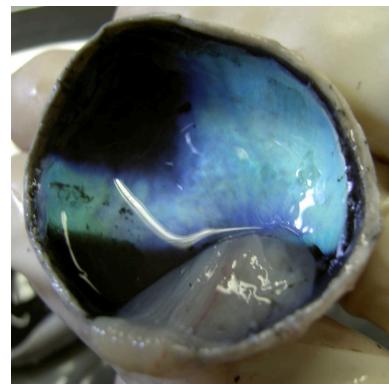


Figure 53: **Tapetum Lucidum**. Choroid dissected from a calf's eye appearing iridescent blue.

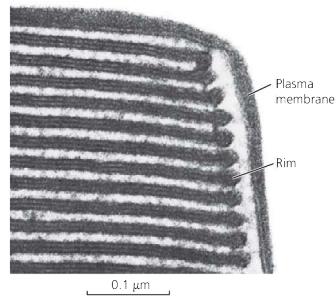
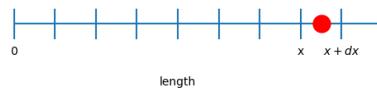


Figure 54: **Rod disks**. Low power electron micrograph of rod outer segment from Fain, Chapter 9.



$$p = \sigma C \Delta x$$

For a photon to be absorbed in the segment between x and $x + \Delta x$, it must have *not* been absorbed in any prior segment between 0 and x . What is the probability that the photon is not absorbed in any prior segment? Divide the distance from 0 to x into N segments. The probability that a photon entering any prior segment is absorbed by that segment is $\sigma Cx/N$. So the probability that a photon is not absorbed by each prior segment is $q = 1 - \frac{\sigma Cx}{N}$. These probabilities must be multiplied to achieve the total probability of not being absorbed by all N prior segments before reaching x :

$$\left(1 - \frac{\sigma Cx}{N}\right)^N$$

Hence, the probability that a photon that enters the rod at $x = 0$ is *first* absorbed between x and $x + dx$ is:

$$dP = \left(1 - \frac{\sigma Cx}{N}\right)^N \sigma C dx$$

In the limit of large N , this differential probability is $dP = e^{-\sigma Cx} \sigma C dx$. The total probability that an entering photon is absorbed by a rod within its length L is the sum of the probabilities that the photon is first absorbed between 0 and L :

$$P = \int_0^L dP = 1 - e^{-\sigma CL}$$

A definition of the exponential function in terms of a limit:

$$\lim_{N \rightarrow \infty} \left(1 - \frac{x}{N}\right)^N = e^{-x}$$

AS A ROD CELL INCREASES IN LENGTH the probability that each photon is absorbed will exponentially converge to 1. Longer rod cells require more rhodopsin molecules. Each rhodopsin molecule has a finite probability of spontaneous thermal isomerization. These spontaneous isomerizations create “dark noise” events. There must be a trade-off between the amount of signal (true absorbed photons) and the amount of noise (spontaneous isomerizations) when increasing rod length.

If the rate of spontaneous isomerization is r_{dark} , the number of dark-noise events in a given interval of time, τ , is proportional to the total number of rhodopsin molecules in the cell, N_{rh} : $\langle n_{dark} \rangle = r_{dark}\tau N_{rh}$. If the number of spontaneous isomerizations were exactly the average number of spontaneous isomerizations, the cell could subtract this number from the total number of isomerizations to calculate the fraction due to real photons. But spontaneous isomerizations occur randomly. A true photon is seen when its signal is larger than the fluctuating noise due to spontaneous isomerizations.

Binomial Statistics

Say that each rhodopsin molecule has a probability, p , of spontaneous isomerization in each trial. Say that we have N rhodopsin molecules. The mean number of ‘successful’ isomerizations in each trial is pN . What are the fluctuations around this mean? The answer requires binomial statistics in the limit of small p and large N . Before flipping rhodopsin molecules, let’s flip coins.

Start with a biased coin, which gives heads with probability p and tails with probability q . What is the probability of k heads from N flips? The probability that a single flip is heads is p . Since each flip is independent, I multiply together the probabilities of results from each set of flips: p^k for the heads, and q^{n-k} for the tails. I have to account for all permutations —HTTHT and HHTTT are distinct permutations with two heads, where each permutation has equal probability . The probability of k heads from N flips is:

$$P(k; N, p) = \binom{N}{k} p^k q^{N-k}, \quad (3)$$

$$\binom{N}{k} = \frac{N!}{k!(N-k)!}. \quad (4)$$

This is the **binomial distribution**. Is the binomial distribution normalized? In algebra, we use the binomial distribution to expand $(a + b)^N$. So:

$$\sum_{k=0}^N P(k; N, p) = \sum_{k=0}^N \binom{N}{k} p^k q^{N-k} \quad (5)$$

$$= (p + q)^N = 1^N = 1. \quad (6)$$

What is the mean of the binomial distribution. In other words, what is the expected value of k , the number of successes?

$$\langle k \rangle = \sum_{k=0}^N k P(k; N, p) \quad (7)$$

$$= \sum_{k=0}^N k \frac{N!}{k!(N-k)!} p^k q^{N-k} \quad (8)$$

$$= \sum_{k=1}^N \frac{N!}{(k-1)!(N-k)!} p^k q^{N-k} \quad (9)$$

$$= Np \sum_{k=1}^N \frac{(N-1)!}{(k-1)!(N-k)!} p^{k-1} q^{N-k}. \quad (10)$$

A change of variables to $m = N - 1$ and $s = k - 1$ makes it clear that the sum is the binomial distribution, which we know is normalized. Therefore,

$$\langle k \rangle = Np \sum_{s=0}^m \frac{(m)!}{(s)!(m-s)!} p^m q^{m-s} \quad (11)$$

$$= Np. \quad (12)$$

The expected number of successes is the probability of success multiplied by the number of trials. However, the mean is only one parameter of a probability distribution. Another parameter is **variance**, a measure of the spread of a distribution about the mean. Variance is defined as:

$$\text{Var}(k) = \sigma_k^2 = \langle k^2 \rangle - \langle k \rangle^2. \quad (13)$$

The square root of variance, σ_k , is the **standard deviation**. To solve for the variance of the binomial distribution, we must compute another expected value, $\langle k^2 \rangle$.

$$\langle k^2 \rangle = \sum_{k=0}^N k^2 P(k; N, p) \quad (14)$$

$$= Np \sum_{k=1}^N k \frac{(N-1)!}{(k-1)!(N-k)!} p^{k-1} q^{N-k} \quad (15)$$

$$= Np \sum_{s=0}^m (s+1) \frac{(m)!}{(s)!(m-s)!} p^m q^{m-s} \quad (16)$$

$$= Np (\langle s \rangle + 1) \quad (17)$$

$$= Np ((N-1)p + 1) \quad (18)$$

$$= N^2 p^2 + Np(1-p) \quad (19)$$

$$= N^2 p^2 + Npq. \quad (20)$$

$$(21)$$

So the variance of the binomial distribution is:

$$\sigma_k^2 = \langle k^2 \rangle - \langle k \rangle^2 \quad (22)$$

$$= N^2 p^2 + Npq - (Np)^2 \quad (23)$$

$$= Npq. \quad (24)$$

Poisson Statistics

Let us now try to take the limit of the binomial distribution as $N \rightarrow \infty$, but $\mu = Np$ is constant. Then:

$$P(k; N, p) = \frac{N!}{k!(N-k)!} p^k (1-p)^{N-k} \quad (25)$$

$$= \frac{N(N-1)(N-2)\dots(N-k+1)}{k!} \left(\frac{\mu}{N}\right)^k \left(1 - \frac{\mu}{N}\right)^{N-k} \quad (26)$$

$$= \frac{\mu^k}{k!} \frac{N(N-1)(N-2)\dots(N-k+1)}{N^k} \left(1 - \frac{\mu}{N}\right)^N \left(1 - \frac{\mu}{N}\right)^{-k}. \quad (27)$$

Taking the limit of each of the three N -dependent terms as $N \rightarrow \infty$, we get 1, $e^{-\mu}$, and 1, respectively. So:

$$P(k, \mu) = \frac{\mu^k}{k!} e^{-\mu}. \quad (28)$$

Since in the limit of small $p, q \rightarrow 1$, the variance $\sigma^2 = Np = \mu$. Poisson processes often show up in biology, where many phenomena are variations on the counting problem (binding processes, photon detection, random excitation, molecule synthesis etc.). Like the binomial distribution, the Poisson distribution is a discrete distribution. Note that the probability of having no events is:

$$P(0, \mu) = e^{-\mu}. \quad (29)$$

It follows that the probability of having at least one success is:

$$P(k \geq 1, \mu) = 1 - e^{-\mu}. \quad (30)$$

The Poisson distribution often shows up describing what is known as a Poisson process. A Poisson process is one where an event occurs with some constant probability per unit time λ , such that over some time t the mean number of events $\mu = \lambda t$. Thus, the probability that k events occur in time t is:

$$P(k, \lambda) = \frac{(\lambda t)^k}{k!} e^{-\lambda t}. \quad (31)$$

Optimal rod length

We are now ready to calculate the optimal length of a rod cell in terms of signal to noise. If the mean number of spontaneous isomerizations in a cell is $\langle n_{dark} \rangle = r_{dark} \tau N_{rh}$, the standard deviations of fluctuations in this number will be:

$$\delta = \sqrt{r_{dark}\tau N_{rh}}$$

A fraction of photons in each flash of light that enters a rod cell are absorbed by that rod cell. This constitutes the true signal:

$$N_{flash} (1 - e^{-\sigma CL})$$

The total number of rhodopsin molecules in the rod cell is a function of rhodopsin concentration, C , and rod volume, AL . The ratio of signal to noise is thus:

$$SNR = \frac{N_{flash} (1 - e^{-\sigma CL})}{\sqrt{CALr_{dark}\tau}}$$

This function has a maximum at an intermediate value of L between 0 and ∞ . Its maximum is reached when $CL \sim 1.26/\sigma$. This means that the probability of an incident photon not being absorbed when signal to noise is maximum is:

$$1 - P = e^{-CL\sigma} \sim e^{-1.26} \sim 0.28$$

Thus, to maximize signal-to-noise ratio, nearly 30% of photons should pass through the rod without being absorbed.

REFERENCES

- Chapter Nine. Gordon L Fain. *Sensory Transduction*. Sinauer Associates, Sunderland, Mass., 2003. ISBN 0878931716 [Download paper](#)
- F Rieke and DA Baylor. Single-photon detection by rod cells of the retina. *Reviews of Modern Physics*, 70 (3):1027–1036, 1998. ISSN 0034-6861 [Download paper](#)

HOW MANY PHOTONS CREATE VISION

RODS DETECT DIM LIGHTS, but how many photons are needed to create light that can be seen. The human eye has so many rod cells (~ 100 million) and each rod cell has so many rhodopsin molecules (~ 100 million) that even if spontaneous isomerizations were exceedingly rare per rhodopsin, they would occur with appreciable rate in the retina as a whole. Unless you knew you were in an absolutely dark room, the brain has no way of telling the difference between a rod cell that is activated by thermal isomerization from a rod cell that is activated by single photon absorption. To see a flash of photons, its signal must be larger than baseline activity caused by thermal isomerizations.

Scientists were interested in the smallest see-able light before the photon was discovered. In 1881, Langley reported a *bolometer*, a device capable of measuring temperature changes as small as 0.00001°C . The bolometer detected the temperature-induced change in electrical resistance of a metal conductor. Increasing temperature increases metal resistance. Langley's bolometer compared the tiny differences in resistance of an illuminated and un-illuminated reference metal. His bolometer could detect thermal radiation from a cow from a quarter mile. Langley estimated the minimum energy of a see-able flash: 3×10^{-16} Joules.

By 1905 when Einstein explained the photoelectric effect with photons, the lower bound on see-able energy for vision was 10-fold smaller. Lorentz used the new equation for photon energy, $E = h\nu$, to estimate the lower bound on see-able *photons*. He estimated that at least ~ 100 photons must be delivered to the cornea to see a flash. Different experiments in different conditions made similar estimates. Hecht, Shlaer, and Pirenne (1934) noted the most reliable:

Hecht et al. knew that rhodopsin (then called visual purple) was the molecular photosensor. They asked a new question: how many rhodopsins must be activated to be seen? This lower bound might be estimated using the measurable corneal reflectance, scattering of the vitreous humor, and rhodopsin absorption. They guessed that the number of photons that might be absorbed by rhodopsin is about one-tenth the number that reach the cornea. They sought a direct measure.

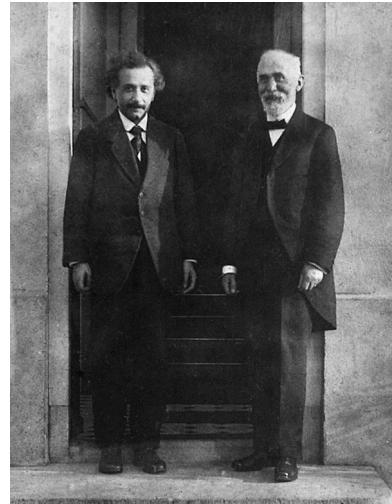


Figure 55: Einstein and Lorentz. In 1905, Einstein (left) published his paper on the photoelectric effect, the emission of electrons when electromagnetic radiation, such as light, hits a material. In classical electricity and magnetism, continuous light waves transfer energy to electrons, which would then be emitted when they accumulate enough energy. Einstein explained that the kinetic energy of emitted electrons was the difference in energy of single photons and a threshold voltage. There was no dependence on the number of photons, only the energy of single photons. The photoelectric effect could be fully explained to be a function of the frequency (energy) of single photons, quanta of light. His friend Lorentz (right) used the new equation for the energy of single photons to make the first estimate of the smallest see-able number of photons.

Wavelength	No. of quanta	Source
505	17-30	Chariton and Lea, 1929
507	34-68	von Kries and Eyster, 1907
530	40-90	Barnes and Czerny, 1932

TO ESTIMATE THE SMALLEST SEE-ABLE LIGHT, Hecht, Shlaer, and Pirenne maximized the odds of seeing a flash by optimizing conditions. They measured scotopic, rod-dominated vision after complete dark adaptation. In the dark, pupils dilate to allow more incoming photons to reach the retina. This happens in seconds and increases sensitivity by ten-fold. In bright light, all rhodopsin molecules in a rod cell are “bleached” and . All retinal, the visual pigment, is non-absorptive in its all-*trans* configuration. In bright light, all rod cells are hyperpolarized with low rates of synaptic release. The full light sensitivity of scotopic vision is reached after biochemistry restores all visual pigment to *cis* configuration. This takes about 30 minutes.

Rods and cones are not evenly spatially distributed in the retina (Fig. 56). Foveal vision – where we have highest-spatial resolution at the center-point of our visual field – is cone-dominated. Rods mostly contribute to peripheral vision. To see a dim light, don’t look straight at it.

Photons delivered to the retina are spread over space and time. Empirically, a larger test area that is illuminated weakly is seen as well as a smaller test area that is illuminated strongly. But this reciprocal relationship between intensity and area is not perfect. There is a maximum sensitivity corresponding to a circular retinal area that spans about 500 rod cells.

Small numbers of photons are most easily seen when arriving in one visual ‘moment’ when all photons can be counted at once. The visual ‘moment’ in both rod and cone photoreceptor cells is surprisingly long. Metabotropic receptors (like rhodopsin) require biochemical signal transduction to open and close ion channels. So metabotropic receptors are often slower than ionotropic receptors (like mechanosensory channels) where stimulus energy directly opens and closes the ion channel. This is why televisions operate at 30 or 60 frames per second. Any slower would create motion jitter between frames. Any faster would not improve the smoothness of movement in video. The reciprocal relationship between intensity and exposure time in scotopic vision holds for stimuli <0.01 sec. So Hecht, Shlaer, and Pirenne used 0.001 s flashes.

Finally, the color of the visual stimulus should be optimized for scotopic vision. The wavelength sensitivity of purified rhodopsin and scotopic vision peak near 510 nm (Fig. 57).

THE HSP EXPERIMENT involved a human observer who triggered a flash of light that would fall on their retina in an area spanning ~ 500 rods at about 20° from the center of vision. Another person

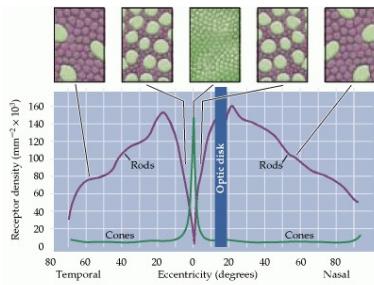


Figure 56: **Rod distribution.** Distribution of cones and rods in a typical human retina

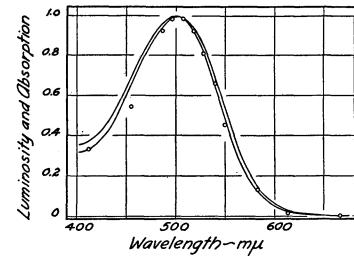


Figure 57: **Rod absorption.** Comparison of scotopic luminosity at the retina with visual purple absorption. The curves are the percentage absorption spectra of visual purple; the upper curve represents 20 per cent maximal absorption, and the lower one 5 per cent maximal absorption. All curves have been made equal to 1 at the maximum, 500 nm, for ease in comparison.

manipulated the filters and wedges that controlled flash intensity. Because the observer triggered each flash, the observer knew when to pay attention. Maximizing the readiness and receptivity of the observer also optimizes sensitivity to minimal flashes. But honesty mattered. The observer had to be truthful when they did not see a flash of light that they knew they had self-administered.

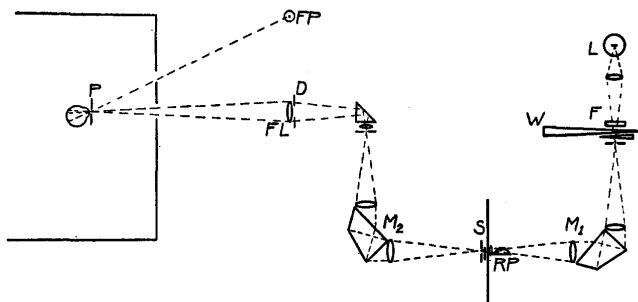


Figure 58: **Apparatus** for measuring minimum energies necessary for vision. The eye at the pupil P fixates the red point FP and observes the test field formed by the lens FL and the diaphragm D . The light for this field comes from the lamp L through the neutral filter F and wedge W , through the double monochromator M_1M_2 and is controlled by the shutter S .

A meaningful threshold for seeing might be when an observer successfully perceives a majority of stimuli. The HSP experiment did not allow *false positives* – the observer cannot see a flash that was not delivered. But every failure to see a flash is a *false negative* – the observer must admit to not seeing a self-administered light. The observer must not pretend to see flashes that they knew had occurred. HSP defined the ‘threshold for seeing’ as one probability, performance at a rate of 60% true positive/40% false negative.

HECHT, SHLAER, AND PIRENNE obtained remarkably consistent results for the threshold for seeing, the number of photons that are seen 60% of the time after reaching the cornea. Seven observers saw flashes that delivered between 54 and 148 blue-green photons to the cornea. These numbers are comparable to the earlier 'most reliable' measurements.

Observer	Quanta
S.H.	126
	135
	107
	87
	79
	123
S.S.	148
	79
	54
	56
	62
	96
C.D.H.	99
	104
	65
	76
M.S.	58
	58
S.R.F.	81
A.F.B.	112
M.H.P.	120
	83
	79
	83
	138

A DEEPER QUESTION is how many photons need to be absorbed by rod cells for the retina to 'see' a flash? One approach is to estimate the fraction of photons that arrive at the cornea that reach the retina and are absorbed by rhodopsin. The cornea reflects $\sim 4\%$ of incident photons. The vitreous humor absorbs $\sim 50\%$ of transmitted photons. Even rods of optimal length will only absorb $\sim 70\%$ of photons that reach the retina. The retina absorbs a small fraction of the photons that arrive at the cornea. If the number of photons that arrive at the cornea is N , the number of absorbed photons is $a = \alpha N$, where α is unknown and takes into account all losses by reflection and absorption and transmission.

Say that flashes result, on average, in one absorbed photon at the retina: $a = 1$. Not every flash will result in one photon absorption. Some will result in zero absorptions. Some will result in two or more absorptions. Poisson statistics describes the complete probability distribution of the number of absorbed photons in each trial, $P(k)$

$$P(k) = \frac{a^k}{k!} e^{-a}$$

, where a is the mean arrival number of photons across trials. Why Poisson statistics? No light source reliably delivers a fixed number of photons to the retina in every trial. Each photon that is released by a light source has a small probability of making its way to the retina and being absorbed. Large numbers of photons and small probabilities that a given photon is absorbed results in Poisson statistics. The probability of seeing curve can be plotted as a function of the mean stimulus size.

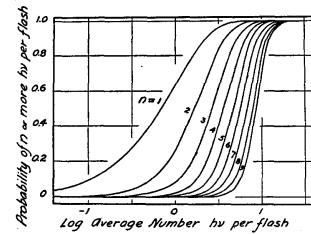


Figure 59: Probability of seeing, the theory. For any average number of quanta per flash, the y-axis gives the probabilities that the flash will deliver n or more quanta to the retina, with different values assumed for n .

If the threshold of seeing is one absorbed photon at the retina, then the ‘probability of seeing’ curve will still be a continuous sigmoidal function because of Poisson statistics. This sigmoidal curve is calculated by summing the probabilities of absorption of one, two, three photons and so on (Fig. 59). If the threshold of seeing is two absorbed photons, then Poisson distributions must be summed from the absorption of two photons and up. For a threshold of n photons, the probability of seeing curve is:

$$P_{see} = \sum_{k=n}^{\infty} \frac{a^k}{k!} e^{-a}$$

As the threshold increases, the sigmoidal curve shifts towards more photons being absorbed. As the curve shifts towards larger stimuli, it becomes steeper when plotted as probability versus the *logarithm* of stimulus size (Fig. 59).

There is merit in plotting P_{see} against the logarithm of stimulus size. Plot P_{see} as a function of the logarithm of the mean number of photons that arrive at the retina, $\log a$. Plot P_{see} as a function of the logarithm of the mean number of photons that arrive at the cornea, $\log N$. Because $a = \alpha N$, $\log a = \log \alpha + \log N$. The two P_{see} curves will have identical shapes except for a horizontal shift of $\log a$.

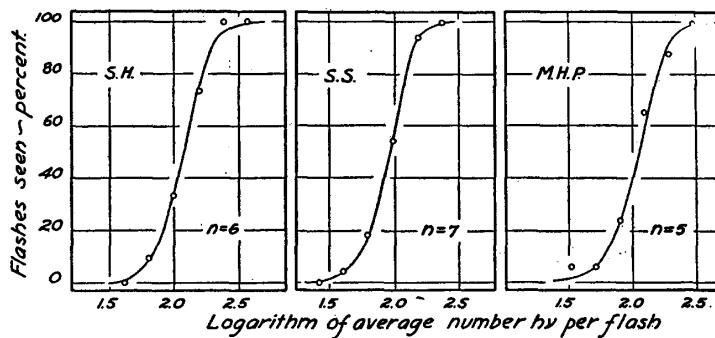


Figure 60: **Probability of seeing, the experiment.** Relation between the average energy content of a flash of light (in number of photons) and the frequency with which it is seen by three observers. Each point represents 50 flashes, except for S.H. where the number is 35. The curves are cumulative Poisson distributions with different thresholds.

Another merit in plotting against the logarithm of stimulus size (whether stimulus size is measured at the cornea or retina) is that the slope of the curve is a direct measure of the threshold:

$$\frac{dP_{see}}{d\log N} \approx \sqrt{n}$$

To derive this square root dependence, one needs to differentiate the probability of seeing, use Stirling’s formula, use the chain rule, and evaluate the slope near the inflection where $a \approx n$. As n increases, the steepness of the probability of seeing curve will increase. Whether the probability of seeing curve is plotted as the logarithm of photons at the cornea or the logarithm of photons at the retina does not matter to this steepness. Steepness is always a function of n , the threshold number of photons needed to see.

The steepness of the probability of seeing curve resembles the **signal-to-noise ratio** of sensory perception. The size of the signal – the numerator of the signal-to-noise ratio – is the number of absorbed photons. Because this signal obeys Poisson statistics, the standard deviation in stimulus size – the noise in the denominator of the signal-to-noise ratio – is its square root: $\sigma = \sqrt{n}$. For the three subjects for whom they had the most data – with initials S.S., S.H., and M.H.P. – HSP derived thresholds of 5,6, and 7 photons from the steepness of their respective probability of seeing curves.

WHY ARE MULTIPLE PHOTONS NEEDED TO SEE? HSP argued that probability of seeing curves depended on the statistical variability of any flash stimulus. But the probability of seeing curve can also depend on the statistical variability of internal noise. A weak signal corresponding to 5-8 photons must be distinguished from a background of spurious signals that will occur in the retina in total darkness. The retina cannot know whether a given rhodopsin activation was due to thermal activation or photon activation. The retina can only count rhodopsin activations. The threshold decision – whether a subject sees a flash of light in a given trial – must be based on the sum of weak and spurious signals. If the weak signal is ‘seen’, it is because this summed signal is larger than typical spurious signals in darkness.

The threshold for a perceptual decision is linked to the reliability of the response. If the threshold is lowered, then response reliability will be lowered. This is because more spurious signals can be interpreted as ‘seeing’, increasing the rate of false positives. If the threshold is raised, some weak signals will not be large enough to be distinguished from spurious signals alone, increasing the rate of false negatives. The probability of seeing and response reliability will be functions of signal size, noise, and the internal detection threshold used by the observer.

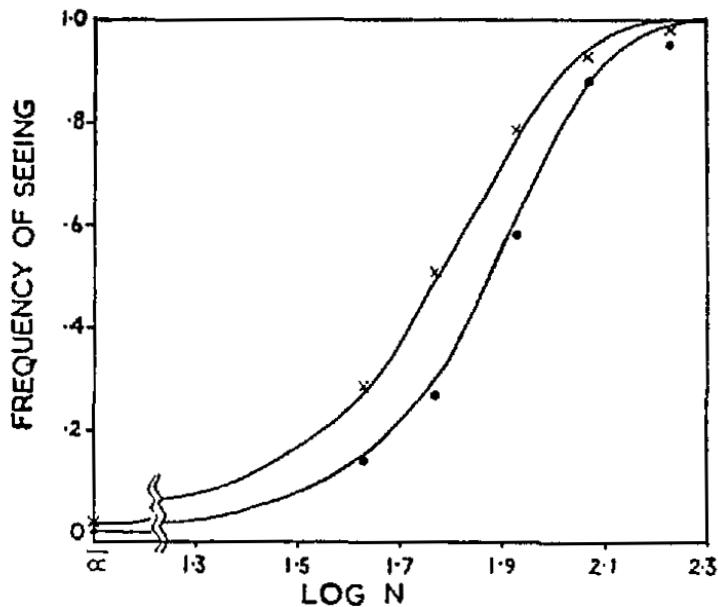


Figure 61: Poisson probability distributions. For any average number of quanta per flash, the y-axis gives the probabilities that the flash will deliver n or more quanta to the retina, with different values assumed for n .

HORACE BARLOW wanted to test the effect of response reliability on detecting weak signals. Human observers are conscious of the certainty of their own perceptions and the reliability of their own decisions. Barlow took advantage of this perceptual awareness to test the responses of one human observer, a fellow named Roy Rumble, at two different detection thresholds. Barlow delivered flashes of different intensities, including blanks where no flash was given, to Rumble. Rumble was encouraged to signal when a flash was “seen” and also when a flash was “possible”. Thus, the threshold number of absorbed photons “possible” would be lower than the threshold for “seen”. Because Barlow knew when blanks were delivered, he knew

when Rumble reported “false positives” and “false negatives”.

Indeed, the probability of seeing curve for “seen” events had a higher threshold (shifted to higher flash intensities) and higher steepness (required more photons) than probability of seeing for “possible” plus “seen” events. Response reliability differed for the two curves. The subject never said a blank was “seen”, but reported 3 of the 300 blanks as “possible”. Lowering the response threshold increased the rate of false positives to 1%. From these curves, Barlow extracted somewhat different parameters from his experiment than those of HSP.

- N , average number of photons at the cornea
- n , average number of photon absorptions, i.e., *bona fide* rod excitations
- x , average number of confusable events, i.e., spurious rod excitations
- $a = x + n$, total average of events (real photons plus noise)
- c , threshold number of events to “see”

Just as for HSP, the probability of seeing is a cumulative probability distribution involving these parameters.

$$P_{\text{see}} = \sum_{k=c}^{\infty} \frac{a^k}{k!} e^{-a}$$

The HSP experiment used the frequency of seeing curve as a function of the logarithm of flash intensity to estimate two values, 1) the fraction of incident photons at the cornea that activated rod cells, and 2) the threshold number of activated rod cells needed to see. These two values can be inferred from one frequency of seeing curve, because this sigmoidal curve has two quantifiable parameters, its horizontal shift and steepness. Barlow’s experiment had more parameters. Like the HSP experiment, a fraction of incident photons would activate rods (one unknown value). But two frequency of seeing curves in Barlow’s experiment meant two separate thresholds (two more unknown value). Lastly, every flash stimulus would be accompanied by a certain number of spurious rod activations (a fourth unknown, noise that has the same average amplitude in every trial). Barlow’s experimental measurement of two frequency of seeing curves contained enough information to calculate the four unknowns.

The steepness of the probability of seeing curve in Barlow’s formulation is somewhat different than the steepness in HSP’s formulation:

$$\frac{dP_{\text{see}}}{d \log N} \approx \frac{c - x}{\sqrt{c}}$$

As before, the steepness evokes the **signal-to-noise ratio**. In this case, the size of the signal in the numerator is the distance of the threshold c from confusable events x , and the size of the noise in the denominator is the expected fluctuations in that signal given by Poisson statistics, \sqrt{c} .

Parameters	Best fits
n/N	0.14
x	8.9
c for Possible or Seen	17
False positive rate for Possible or Seen	0.01
c for Seen	19
False positive rate for Seen	0.002

For “seen” events, meeting the threshold requires 10 events above the number of confusable events (spontaneous rhodopsin activation). In terms of signal-to-noise, events corresponding to “seen” *bona fide* flashes are ~ 2.3 standard deviations from the confusable events. Events corresponding to “possible” plus “seen” flashes are ~ 2 standard deviations from confusable events.

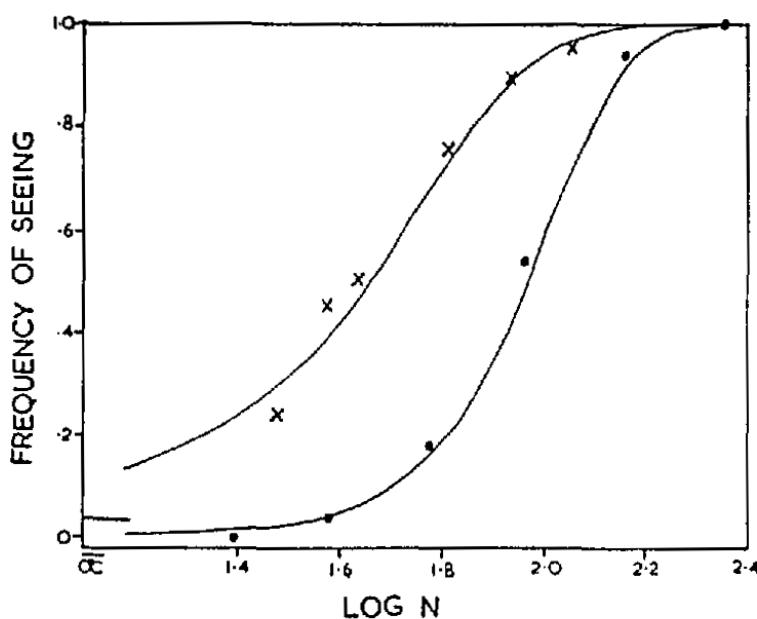


Figure 62: **Other experiments.** Data from Hecht et al. (dots) and van der Velden (crosses) fitted using Barlow's theoretical curves and parameters.

Is Barlow's formulation consistent with earlier measurements? From one probability of seeing curve in HSP (or one probability of seeing curve in another measurement by van der Velden), it is impossible to extract all the parameters that Barlow used. Earlier measurements did estimate the number of photons at the cornea, N , from which they inferred the fraction of photons that were lost and the threshold number of rhodopsin activations needed to see. If the number of spontaneous rhodopsin activations in earlier measurements was similar to that measured by Barlow, $x = 8.9$, these numbers could be inferred (Fig. 62). Fitting Barlow's parameters to HSP's experiments revealed roughly consistent numbers for the fraction of absorbed photons ($n/N = 0.13$) and the threshold of rhodopsin activations ($c = 21$). Fitting Barlow's parameters to another experiment by van der Velden revealed an unrealistically large fraction of absorbed photons ($n/N = 0.9$) and lower threshold of rhodopsin activations ($c = 15$).

REFERENCES

- S Hecht, S Shlaer, and M H Pirenne. Energy, quanta, and vision. *The Journal of General Physiology*, 25(6): 819–840, 1942. ISSN 0022-1295 [Download paper](#)
- HB Barlow. Retinal noise and absolute threshold. *Journal of the Optical Society of America*, 46(8):634–639, 1956. ISSN 0030-3941 [Download paper](#)

SINGLE PHOTONS AND SINGLE ROD CELLS

ARE SINGLE UNITS OF STIMULUS TRANSFORMED INTO SINGLE UNITS OF ROD CELL RESPONSE? The single stimulus unit in rod vision is one photon. This leads to a cellular response that can be measured as a change in cation current through the rod outer membrane. The predominant cation in the extracellular solution is sodium. The predominant cation inside the intracellular cytoplasm is potassium. In darkness, Na^+ channels in the outer segment are open, producing inward ionic current. At the same time, K^+ channels in the inner segment are open, producing outward ionic current. When a rhodopsin molecule is activated by absorbing a photon in the outer segment, it triggers a biochemical signal transduction cascade. This cascade closes Na^+ channels in the outer segment. By Ohms Law ($V = IR$), a drop in inward cation current means hyper-polarization of the membrane potential. Membrane potential is typically reported as internal voltage minus external voltage. Hyper-polarization means that the rod cell membrane become even more negative than baseline. Hyper-polarization lowers chemical synaptic transmission from the rod cell to downstream cells. Illumination quietens rod synaptic communication.

The anatomy of the rod cell is ideal for recording the change in light-evoked currents using a suction electrode. The outer segment is pulled into a tightly-fitting glass pipette. Any change in the current loop between an electrode inside the pipette and an electrode in the bath can be recorded. When a rod cell is sucked into the pipette, the change in this current becomes a direct measure of light-evoked changes in the opening and closing of cation channels in the outer segment.

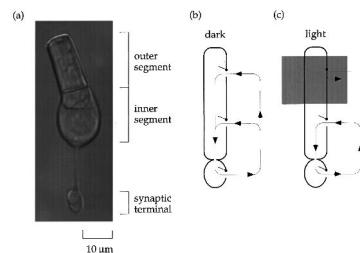


Figure 63: **Rod cells.** (a) Isolated rod photoreceptor from salamander. The cell membrane separates intracellular and extracellular salt solutions. The outer segment contains rhodopsin and transduces photons into rod activation. The inner segment keeps the cell alive. The synaptic terminal communicates signals to bipolar and horizontal cells. (b) In darkness Na^+ ions enter the outer segment through cGMP-gated channels. A current loop is completed by outward movement of K^+ ions through channels in the inner segment. Separate Na/K exchange pumps maintain the relatively lower intracellular concentration of Na^+ (which might be 15 mM inside a typical cell and 140 mM outside) and relatively higher intracellular concentration of K^+ (which might be 120 mM inside and 3 mM outside). (c) When a rod cell is exposed to light, some channels in the outer segment close. The cell hyperpolarizes, which reduces neurotransmitter release at the synaptic terminal.

CURRENTS THROUGH THE OUTER SEGMENT MEMBRANE WITH SUSTAINED ILLUMINATION. The bars beneath the traces in Fig. 64 indicate the duration of sustained illuminations and the numbers indicate intensity in photons $\mu\text{m}^{-2} \text{ sec}^{-1}$. The dimmest sustained lights evokes spontaneous transient 1 pA bumps in outward membrane current. Each spontaneous bump resembles flash-evoked responses with dim lights, also 1 pA bumps in outward current. Current fluctuations increase at higher light intensities until saturation at the brightest lights. Do the single spontaneous bumps evoked by dim sustained illumination represent single unit responses?

Recording rod currents with flash-evoked responses is like asking whether the rod saw a light. Instead of answering 'yes' or 'no', the rod responds with a measurable current. When a rod is subjected to consecutive dim flashes, some flashes evoke no apparent response. Some flashes evoke 1 pA bumps that resemble the spontaneous 1 pA bumps that occur during sustained dim illumination. Some flashes evoke larger currents. Flash responses appear quantized, consistent with representing quantal numbers of unit responses. The amplitude histogram for current responses with fixed dim flash intensity has sharp peaks at 0 pA and near 1.2 pA. The peaks are well-separated, but there is still significant variability in response amplitude around each peak, fit to normal distributions with small standard deviations around each peak.

If the 0 pA peak of the histogram is taken to correspond to zero responses and the 1 pA peak is taken to correspond to one single unit response, a threshold current of 0.5 pA is a good criterion for success (one or more unit responses) and failure (no unit response). A threshold near 0.5 pA is a criterion with few false positives and false negatives. With a threshold criterion of 0.5 pA, fifty-eight of the ninety-nine trials in Fig. 66A were failures. When the same cell was stimulated with lower light intensity, shown in Fig. 66B, forty-four of fifty-two trials were failures.

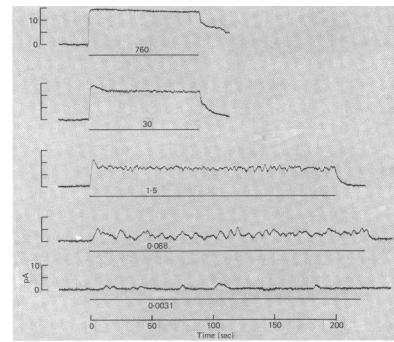


Figure 64: **Rod cells.** Response of a rod outer segment to steady lights. Ordinate is outward change in membrane current from level in darkness. Bars traces indicate duration of light stimuli; numbers give intensities in photons per μm^2 per second. From Baylor et al. 1979.

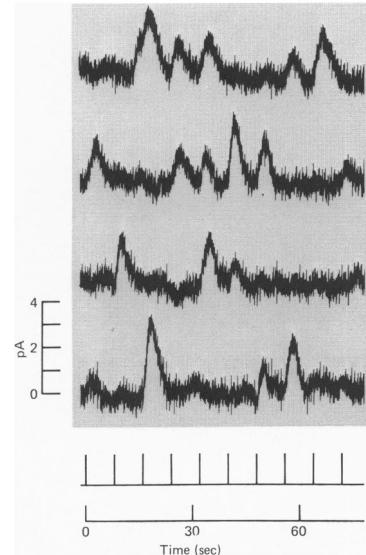


Figure 65: **Rod cells responses to flashes.** Response of outer segment to a series of forty consecutive dim flashes. Local illumination; 20 msec flash delivering 0.029 photons per μm^2 per second. From Baylor et al. 1979.

THE HISTOGRAMS OF STIMULUS-EVOKED CURRENTS (FIG. 66) CONTAIN MORE INFORMATION than the probabilities of success and failure. The stronger stimulus evokes more 1 pA responses than 0 pA responses, but also more 2 pA responses. Both 1 pA and 2 pA responses represent success, but do 2 pA responses really two unit responses? Beyond the absolute height of the first peak in the histogram (the probability of failure), can we predict the relative height of all other peaks in the histogram (the probability of one, two, or more unit responses per trial)?

Poisson statistics should apply to these physiological recordings. On average, many photons are delivered to each rod cell per flash – N is large. Each photon has a small probability of evoking a unit response – p is small. The mean number of unit responses per flash is small but non-zero – $\langle k \rangle = pN$. In this regime of Poisson statistics, knowing only the mean number of unit responses per trial, we know the probability of any number of unit responses. Hence, the probability $P(k)$ of k responses is:

$$P(k) = \frac{e^{-\langle k \rangle} \langle k \rangle^k}{k!}$$

The probability of failure is the probability that $k = 0$: For the data in Fig. 66A, $P(k = 0) = e^{-\langle k \rangle} = 58/99$, so $\langle k \rangle = 0.53$. The probability of success per flash ($P(k > 0)$) is less than one-half. The mean number of unit responses is more than one-half because success includes flashes that evoked two or more unit responses.

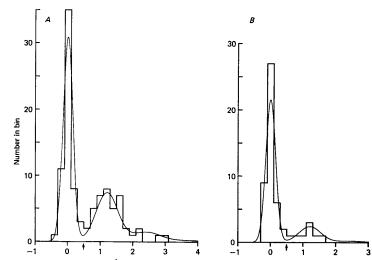


Figure 66: Rod cells. Amplitude histograms for flash responses. Ordinate is number of occurrences of amplitude shown on abscissa, bin width 0.2 pA. A, flash intensity (0.029 photons per μm^2 per sec. Response amplitude had mean $\mu = 0.56$ pA and variance $\sigma^2 = 0.58$ pA 2 . B, Same cell subjected to lower intensity flashes (0.014 photons per μm^2 per sec. Mean response: $\mu = 0.17$ pA, variance $\sigma^2 = 0.21$ pA 2). Arrows mark the criterion level of 0.5 pA used to separate responses from failures. From Baylor et al. 1979.

WE WANT TO DISTINGUISH DIFFERENT NUMBERS OF UNIT RESPONSES evoked by each flash from a continuous distribution of recorded currents. To do this, the unit response caused by photon absorption must be highly stereotyped. Many molecular events occur between the stimulus (photon absorption by rhodopsin) to the measured response (current changes in the outer segment membrane). After absorbing a photon, rhodopsin activates a signal transduction cascade to regulate the opening and closing of ion channels. This signal transduction cascade is a signal-amplifying network of enzymes and diffusible molecules (Fig. 67). In its activated state, rhodopsin activates copies of the protein transducin. Both rhodopsin and transducin are membrane-bound proteins, diffusing within the discs of the rod cell, interacting by random collision. One active rhodopsin activates ~ 1000 transducin in one second. One transducin molecule activates one membrane-bound enzymatic protein, cGMP phosphodiesterase (PDE). Each activated PDE molecule destroys ~ 50 cGMP molecules before PDE is inactivated. cGMP is the cytoplasmic messenger from the rod discs to the rod outer segment membrane. cGMP binds to ion channels in the outer segment membrane and stabilizes their open state. When cGMP concentration is lowered by the activation of one rhodopsin molecule, hundreds of ion channels close. In one second, closing one ion channel blocks $\sim 10,000$ cations from entering the rod.

The total “gain” of the system is the multiplication factor from single photon absorption to the change in large numbers of ions flowing through the outer segment. Trial-to-trial variability in gain creates a problem in counting the number of unit responses within each flash-evoked change in current. If the amplitude of a unit response is typically 1 pA, a 2 pA response might represent the combination of two unit responses with typical gain or might represent one unit response with atypically twice gain. For rods to reliably count photons, not just detect photons, the gain cannot be so variable. In fact, the gain is highly stereotyped. Later, we will discuss potential molecular mechanisms of gain invariance. For now, we assume that gain is fixed. What is the stereotyped change in membrane current corresponding to the single unit response?

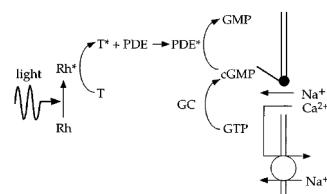


Figure 67: Schematic of signal transduction cascade. The light-sensitive current is carried by Na^+ and Ca^{2+} ions, which enter the outer segment through channels in the surface membrane. These channels are held open by the binding of cyclic guanosine monophosphate (cGMP). A decrease in cGMP in response to photon absorption permits channels to close, reducing the current. Rhodopsin (Rh) is activated by photon absorption. Active rhodopsin catalyzes the activation of transducin (T), which in turn activates a cGMP phosphodiesterase (PDE). Activated PDE hydrolyzes cGMP, causing its concentration to fall. Recovery of the light response requires the restoration of the cGMP concentration. This is accomplished by guanylate cyclase (GC), which synthesizes cGMP from guanosine triphosphate (GTP). Ca^{2+} is removed from the outer segment by an exchange protein.

THE SIMPLEST VISUAL AND INTUITIVE ESTIMATE of the single unit response is locating the first peak with non-zero current in the histogram of flash-evoked current changes. In Fig. 66, this peak occurs near ~ 1.2 pA.

THE MEAN RESPONSE AMPLITUDE of all membrane currents after all flashes allows another estimate. Calculate μ , the mean current response after all flashes, collecting both ‘successes’ with non-zero current and ‘failures’ with near-zero current. From Poisson statistics, estimate $\langle k \rangle$, the mean number of unit responses per flash. The constant of proportionality between μ and $\langle k \rangle$ will be the current amplitude of the unit response, a .

$$\mu = \langle k \rangle a$$

For Fig. 66, $\mu = 0.56$ pA and $\langle k \rangle = 0.53$, so a is ~ 1.06 pA.

THE VARIANCE OF MEASURED RESPONSES allows a third estimate. Consider the time-varying measurements of membrane current before and after each flash. From current measurements at each time point, calculate the mean current change from baseline, $\langle i(t) \rangle$. One can also calculate the variance in current measurements at each time point, $\langle (i(t) - \langle i(t) \rangle)^2 \rangle$.

Before each flash, the mean current change from baseline will fluctuate around zero. The variance in current measured at each time point will also fluctuate, but around a stationary non-zero value, a ‘noise floor’. This noise floor is not caused by photon absorptions, but by stationary fluctuations intrinsic to the experimental setup in total darkness. These stationary fluctuations have zero mean but non-zero variance. These stationary fluctuations are also uncorrelated with any fluctuations in stimulus-evoked current caused by photon absorption. The total variance in current measurements after a flash will then be the sum of the variance caused by stationary fluctuations and the variance caused by the statistics of stimulus-evoked responses.

Consider current measurements after the flash at $t = 0$. We need to describe the unit response as a time-varying but stereotyped waveform, $a(t)$. We describe the noise floor at each time point with a_n , a random number with zero mean $\langle a_n \rangle = 0$ and non-zero variance a_n^2 . The measured current at each time point after each flash is caused by the sum of k unit responses and the noise floor:

$$i(t) = ka(t) + a_n$$

The mean current at each time point after each flash is:

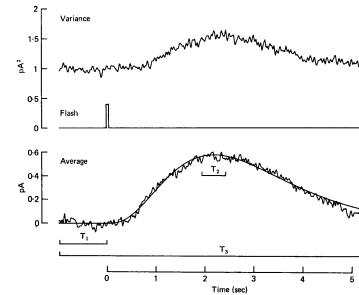


Figure 68: **Mean and variance of the flash response** Mean (below) and variance (above) were computed at 12 msec intervals before and after each flash. Flash timing is shown in middle.

$$\langle i(t) \rangle = \langle k \rangle a(t)$$

The variance in current measurements after each flash is:

$$\langle (i(t) - \langle i(t) \rangle)^2 \rangle = \langle (k - \langle k \rangle)^2 \rangle a^2(t) + a_n^2$$

Poisson statistics dictates that the mean number of unit responses will equal the variance in the number of unit responses:

$$\langle k \rangle = (\langle (k - \langle k \rangle)^2 \rangle)$$

Putting these equations together, we find that $a(t)$ is the proportionality between the variance (after subtracting the noise floor) and the mean of measured currents after each flash.

$$\langle (i(t) - \langle i(t) \rangle)^2 \rangle - a_n^2 = \langle i(t) \rangle a(t)$$

At the peak rise in variance and peak rise in mean current, $a(t)$ is ~ 1 pA, consistent with other estimates.

Does the unit response correspond to the activation of a rod cell by a single photon? Or might each unit response require multiple photon absorptions? Hecht, Shlaer, and Pirenne used ‘probability of seeing’ curves as a function of light intensity to establish 5-8 absorbed photons as the visual threshold for human perception. We can do the same thing for a rod cell. The mean number of photon isomerizations will be proportional to light intensity:

$$m = AI$$

where I is flash intensity in photons per μm^2 and the constant A is the effective collecting area of the outer segment (μm^2). If one photon is needed to trigger a unit response, then the probability of success, P_S , is one minus the probability of failure where failure is the probability of absorbing zero photons, $P(0) = e^{-AI}$:

$$P_S = 1 - e^{-AI}$$

This curve has only one free parameter, A , to fit the probability of seeing curve of a single rod cell as a function of flash intensity. For the data shown in Fig. 69, A is $12.7 \mu\text{m}^2$. This cross-section is in reasonable agreement with cell dimensions, pigment density, and quantum efficiency of isomerization.

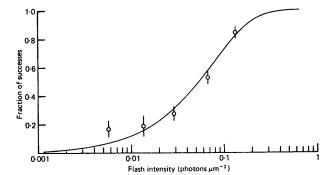


Figure 69: **Frequency of seeing experiment with single rods** Fraction of responses exceeding the criterion for single unit responses at five intensities plotted against flash intensity on a logarithmic scale. The curve is $P_S = 1 - e^{-AI}$ where $A = 12.7 \mu\text{m}^2$.

SINGLE UNIT RESPONSES, recorded in the outer segments of individual rod cells in response to dim illumination, are consistent with single photon absorptions. Similar electrical events also occur in rod cells in complete darkness. Absent photon energy, thermal energy should also be able to spontaneously activate rhodopsin, albeit much more rarely since thermal energy fluctuations are much smaller than the energy of a visible photon. Are the single unit responses in darkness consistent with thermal rhodopsin isomerization?

Outer membrane currents of single rod cells from the toad were recorded using electrophysiology in complete darkness and in light (Fig. 70). In both darkness and light, membrane current exhibits continuous, small amplitude fluctuations that represent baseline noise from the experimental setup. But in darkness, an additional fluctuation consists of 1 pA discrete events. These discrete electrical events are rare (~ 0.03 per sec) and have the same amplitude and time-course as photon-triggered events in dim light.

If the rare 1 pA events in darkness represent thermal fluctuations, they should exhibit Poisson statistics. The rod outer segment contains 2×10^9 rhodopsin molecules. The probability per unit time of spontaneous thermal isomerization must then be $\sim 10^{-11}$ per second. On average, the waiting time for spontaneous thermal isomerization for each rhodopsin molecule will be 1,000 years. One way to test Poisson statistics is to count successes in an observation interval. With discrete electrical events occurring at 0.03 per sec, the average number of events in one minute will be $\langle k \rangle = 1.8$. The probability of k events will then be $P(k) = \frac{e^{-\langle k \rangle} \langle k \rangle^k}{k!}$. Another way to test Poisson statistics of events that occur over time is to measure intervals between successes.

Consider a process in which a specific event occurs with a constant probability per unit time λ . Our example is the probability per unit time that a single rhodopsin undergoes isomerization inside a rod cell, observed as an electrical event with a recording pipette. A canonical example is the probability that a single atom undergoes radioactive decay from a lump of radioactive material, observed as a click on a Geiger counter. In either case, start observing the system at $t = 0$ and ask: what is the probability that the next event happens between t and $t + dt$? We can compute just the probability of event occurrence within the $(t, t + dt)$ interval by multiplying the probability per unit time (λ) and interval duration (dt). The probability that this is the *first* event since the beginning of observation at $t = 0$ is λdt times the probability that *no* events occurred in the $(0, t)$ interval. To calculate this probability of non-occurrence divide the timeline from 0 to t into N intervals, $\Delta t = t/N$. The probability of non-occurrence

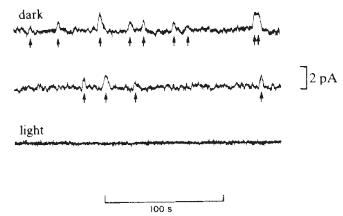


Figure 70: Spontaneous electrical events in rods in darkness Sample records of frog outer segment current in darkness (two upper traces) and bright light (bottom trace). Arrows below traces indicate times at which the electrical event exceeded a criterion level of 0.5 pA above baseline. 82 events were counted over 2631 seconds.

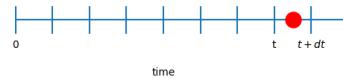


Figure 71: Line extending from 0 time $t + dt$ can be divided into a large number of increments of length t/N to derive the Poisson interval distribution.

in each of the N intervals is $1 - \frac{\lambda t}{N}$. The probability of non-occurrence in all of the N intervals is $\left(1 - \frac{\lambda t}{N}\right)^N$. In the limit of large N , the binomial expansion converges to an exponential function:

$$\lim_{N \rightarrow \infty} \left(1 - \frac{\lambda t}{N}\right)^N = e^{-\lambda t}$$

The Poisson interval distribution is the probability of non-occurrence in the $(0, t)$ interval and occurrence in the $(t, t + dt)$ interval:

$$P(t)dt = \lambda e^{-\lambda t} dt$$

The Poisson interval distribution is properly normalized:

$$\int_{t=0}^{\infty} P(t)dt = \int_{t=0}^{\infty} \lambda e^{-\lambda t} dt = 1$$

The mean of the Poisson interval distribution also gives the intuitively correct answer for the mean waiting time for the next event ($\langle t \rangle = \frac{1}{\lambda}$):

$$\langle t \rangle = \int_{t=0}^{\infty} t \lambda e^{-\lambda t} dt = \frac{1}{\lambda}$$

The expectation value of t^n is

$$\langle t^n \rangle = \int_{t=0}^{\infty} t^n \lambda e^{-\lambda t} dt = \frac{n!}{\lambda^n}$$

The mean-square interval is $\langle t^2 \rangle = \frac{2}{\lambda^2}$. Therefore, the standard deviation is:

$$\sigma_t = \sqrt{\langle (t - \langle t \rangle)^2 \rangle} = \langle t \rangle$$

The mean and the standard deviation of the Poisson interval distribution are the same. A common normalized measure of the width of a probability distribution is the **Coefficient of Variation**, defined as the ratio of the standard deviation σ to the mean μ : $CV = \frac{\sigma}{\mu}$. The Coefficient of Variation of the Poisson interval distribution is 1. Expected deviations from the mean waiting time are as long as the mean waiting time.

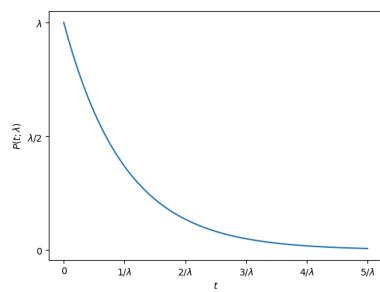


Figure 72: The Poisson interval distribution plotted in units of time $\tau = 1/\lambda$.

THE POISSON DISTRIBUTION AND THE POISSON INTERVAL DISTRIBUTION ARE INTERCONNECTED. We can use the Poisson interval distribution to predict the number of events that occur in a fixed time t' . The probability that an event does *not* occur in the interval t' ($P(k = 0)$) is the probability that the waiting time for the next event is longer than t' :

$$P(k = 0) = \int_{t'}^{\infty} \lambda e^{-\lambda t} dt = e^{-\lambda t'}$$

The probability that one event occurs in the interval t' ($P(k = 1)$) is the probability that the next event occurs between t_0 and $t_0 + dt_0$ times the probability that zero events occur in the rest of the interval, $t' - t_0$, integrated over all values of t_0 from 0 to t' :

$$\int_0^{t'} \lambda e^{-\lambda t_0} e^{-\lambda(t' - t_0)} dt_0 = \lambda t' e^{-\lambda t'}$$

The probability that two events occur in the interval t' ($P(k = 2)$) is the probability that the next event occurs between t_0 and $t_0 + dt_0$ times the probability that one event occurs in the rest of the interval, $t' - t_0$, integrated over all values of t_0 :

$$\int_0^{t'} \lambda e^{-\lambda t_0} \lambda (t' - t_0) e^{-\lambda(t' - t_0)} dt_0 = \frac{(\lambda t')^2}{2} e^{-\lambda t'}$$

By this process of iteration, we can calculate the probability that k events occur in the interval t' :

$$P(k; \mu) = \frac{\mu^k}{k!} e^{-\mu}$$

where $\mu = \lambda t'$. This is the Poisson distribution, where the mean number of events is $\langle k \rangle = \mu$.

IF THE SPONTANEOUS ELECTRICAL EVENTS IN DARKNESS CORRESPOND TO THERMAL ISOMERIZATION, time intervals between events should follow the Poisson interval distribution. The rod cell shown in Fig. 70 exhibited 82 events. The cumulative number of events less than an interval T should then be the integral of the Poisson interval distribution from 0 to T times the total number of events:

$$\text{No. intervals} \leq T = N \int_0^T \lambda e^{-\lambda t} dt = N (1 - e^{-\lambda T})$$

where $\frac{1}{\lambda} = 32$ sec is the average interval between successive events for the data shown in Fig. 73.

THERMAL ENERGY IS HIGHER AT HIGHER TEMPERATURES. The rate of spontaneous thermal isomerizations should increase systematically with higher temperature. One advantage of using rod cells from amphibians is that they are robust to temperature variations. Unlike mammals, toads and their sensory cells survive across a broad range of temperatures, from 0 to 30°C, a large enough range to measure changes in statistical mechanics with changes in $k_B T$.

Spontaneous rhodopsin isomerization is a first-order chemical reaction characterized by two low-energy stable states (the inactive *cis* state and the active *trans* state) separated by a high-energy transition state. Before photon absorption, almost all rhodopsin are in the *cis* state. The rate of the forward reaction will be proportional to the fraction of these rhodopsin molecules that spontaneously reach the transition state from the *cis* state. These molecules that reach the transition state can progress to the *trans* state or return to the *cis* state. Knowing the activation energy, the difference in energy E_{act} between the transition state and the *cis* state, the Boltzmann distribution tells us the fraction of molecules at the transition state, P_{trans} , where all molecules start from the *cis* state:

$$P_{trans} = e^{-\frac{E_{act}}{k_B T}}$$

The rate of the spontaneous isomerization is this probability times a frequency factor v . The frequency factor, with units of sec^{-1} , is a catch-all heuristic that captures the number of attempts per second of the forward reaction. The logarithm of the rate of the *cis* to *trans* state is thus

$$\log k_{cis \rightarrow trans} = \log v - \frac{E_{act}}{k_B T}$$

The Arrhenius plot, the logarithm of reaction rate versus inverse temperature $1/T$, tells us the activation energy E_{act} . From the rate

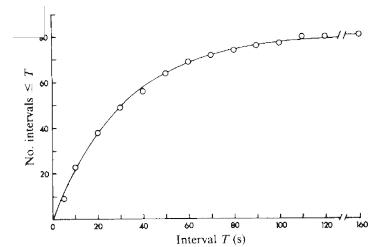


Figure 73: Cumulative distribution of intervals between successive events in the same cell shown in Fig. 70.

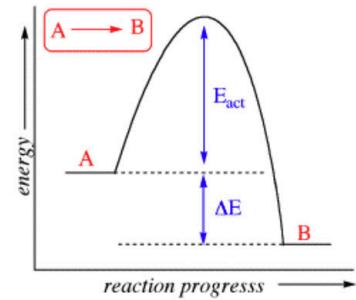


Figure 74: Cartoon of the reaction coordinate for rhodopsin from the A (*cis*) to B (*trans*) state via an intermediate high-energy transition state.

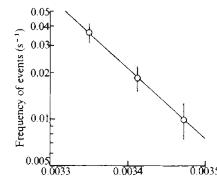


Figure 75: Arrhenius plot of frequency of occurrence of dark events in a rod (log scale) against inverse absolute temperature.

of spontaneous electrical events at different temperatures, we can estimate the rhodopsin activation energy, $E_{act} \approx 30k_B T$. This activation energy is much larger than $1k_B T$, so that spontaneous thermal isomerizations are large, but below the $\sim 100k_B T$ energies of visible photons that usually trigger rhodopsin isomerization.

THE WEAKEST PULSES OF LIGHT THAT A HUMAN CAN DETECT corresponds to 10-20 total rhodopsin isomerizations. The threshold for seeing pulses of light has to be larger than the number of spontaneous thermal isomerizations for the signal-to-noise ratio to be larger than one. If the rate of thermal isomerizations in the rod cells of the toad can be reduced by lowering temperature, the threshold for seeing pulses of light might also be reducible.

We know a dark-adapted starved toad in a light-tight box “sees” a flash of light when it snaps at a dimly-illuminated “worm-dummy”. The number of rod cells that are illuminated by the worm-dummy is proportional to its size (~ 4500 rods in the experiment shown in Fig. 76). Toads increase snapping frequency when the rhodopsin isomerization rate due to illumination (measured in isomerizations per rhodopsin per sec) exceed a threshold comparable to intrinsic thermal isomerization rate ($4.9 \times 10^{-12} \text{ Rh}^{-1} \text{ sec}^{-1}$). All toads responded promptly in the intensity range from 3.0×10^{-11} isomerizations $\text{Rh}^{-1} \text{ sec}^{-1}$. A significant increase in snapping occurs at 3.0×10^{-12} isomerizations $\text{Rh}^{-1} \text{ sec}^{-1}$.

Toads start detecting objects that add only 3.0×10^{-12} isomerizations $\text{Rh}^{-1} \text{ sec}^{-1}$ to the ongoing rate of $4.9 \times 10^{-12} \text{ Rh}^{-1} \text{ sec}^{-1}$. To understand why this is possible, we need to consider signal-to-noise ratios. Say that a signal is collected from 440 rods, each containing 3.25×10^9 molecules of rhodopsin, over a 1.9 sec interval (parameters corresponding to the most sensitive ganglion cell recorded by the same researchers using electrophysiology). At the behavioral threshold, this signal would correspond to 8.4 photon-triggers isomerizations and 13.3 thermal isomerizations. The mean signal amplitude is 8.4 photoisomerizations, but the magnitude of Poisson fluctuations from the mean signal amplitude will be $\sqrt{8.4 + 13.3}$. The signal-to-noise ratio is ~ 1.8 , which suggests low reliability. But the brain of the toad probably improves on its signal-to-noise ratio by summing the signal from many ganglion cells activated by the worm-dummy spanning ~ 4500 rods.

The threshold for seeing for amphibians at different temperatures and human (trying to see the same worm-dummies in the same setup at different light intensities) reveals a monotonic dependence with absolute temperature. Thermal isomerization of rhodopsin sets the ultimate limit on threshold intensity for visual detection.

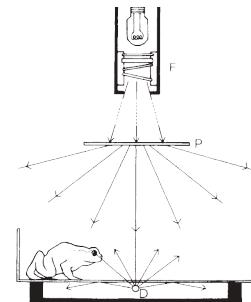


Figure 76: Schematic of experimental setup of snapping experiment. The distance of the toad from the worm-dummy will illuminate a region of about 4,500 rods, triggering a certain number of isomerizations per rhodopsin per second. For reference, the thermal isomerization rate is 4.9×10^{-12} per rhodopsin per second.

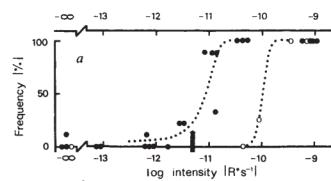


Figure 77: Schematic of experimental setup of snapping experiment. The distance of the toad from the worm-dummy will illuminate a region of about 4,500 rods, triggering a certain number of isomerizations per rhodopsin per second. For reference, the thermal isomerization rate is 4.9×10^{-12} per rhodopsin per second.

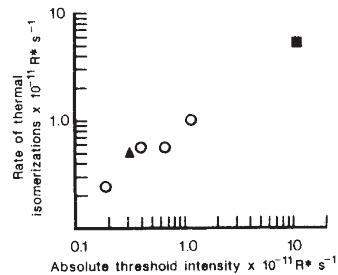


Figure 78: Correlation between rates of thermal rhodopsin isomerizations and absolute threshold intensities, expressed as rates of isomerization per rhodopsin per second in the retina of the toad (▲), frog (○), and human (■).

REFERENCES

- Greg D. Field and Fred Rieke. Mechanisms regulating variability of the single photon responses of mammalian rod photoreceptors. *Neuron*, 35(4):733–747, 2002. ISSN 0896-6273 [Download paper](#)
- D A Baylor, T D Lamb, and K W Yau. Responses of retinal rods to single photons. *The Journal of Physiology*, 288(1):613–634, 1979. ISSN 0022-3751 [Download paper](#)
- K.-W YAU, G MATTHEWS, and D. A BAYLOR. Thermal activation of the visual transduction mechanism in retinal rods. *Nature*, 279(5716):806–807, 1979. ISSN 0028-0836 [Download paper](#)
- A.-C Aho, K Donner, C Hydén, L. O Larsen, and T Reuter. Low retinal noise in animals with low body temperature allows high visual sensitivity. *Nature*, 334(6180):348–350, 1988. ISSN 0028-0836 [Download paper](#)
- J CHEN, CL MAKINO, NS PEACHEY, DA BAYLOR, and MI SIMON. Mechanisms of rhodopsin inactivation in vivo as revealed by a cooh-terminal truncation mutant. *Science*, 267(5196):374–377, 1995. ISSN 0036-8075 [Download paper](#)
- Thuy Doan, Ana Mendez, Peter B Detwiler, Jeannie Chen, and Fred Rieke. Multiple phosphorylation sites confer reproducibility of the rod's single-photon responses. *Science*, 313(5786):530–533, 2006. ISSN 0036-8075 [Download paper](#)

RELIABLE SIGNAL TRANSDUCTION IN THE ROD CELL

Photon absorption leads to a stereotyped change in the membrane current of a rod outer segment (Fig. 79). This current change is the output of a signal-amplifying cascade that blocks the entry of millions of cations with the absorption of single photons (Fig. 67). Amplification is essential to detect single photons, effectively a conversion of single-photon energies into perceptual signals registered by the brain. This amplification is strikingly reliable as the mean amplitude of the elementary electrical response is much larger than its standard deviation. The coefficient of variation (the unit-less, normalized ratio between mean and standard deviation) is low. The electrical activity of the rod cell does not only reliably indicate whether photons have been absorbed but also reliably indicates how many photons have been absorbed.

Every signal-amplifying molecular event between the activation of single rhodopsin molecules and the closing of ion channels is inherently stochastic. Molecules internally vibrate with thermal energy. Molecules diffuse laterally along membranes and three-dimensionally through the rod cytoplasm to interact with other molecules and propagate signals. Every stochastic molecular event is a source of variability in the amplification of single-photon detection to rod activity. With greater variability in signal-amplification comes less reliability in discriminating different stimulus amplitudes. Consider the variability of the first step in signal-amplification. A single photon-activated rhodopsin molecule can activate a large number of transducin molecules before rhodopsin itself becomes inactivated. Activated rhodopsin and transducin interact by lateral diffusion and random collision in the disks of the rod cell. The first step in amplification is the multiplicative number of transducin molecules that are activated by one activated rhodopsin. The longer the lifetime of activated rhodopsin, the more transducin can be activated.

A simple mechanism that might reduce the variability in the number of activated transducin would be saturation. If an activated rhodopsin activated *every* transducin molecule in its neighborhood and the number of available transducin molecules were constant, every photon would activate the same number of transducin molecules. Experiments suggest that the signal transduction machinery of the rod cell is not locally saturated by single photons (Fig. 80). If saturation does not explain low trial-to-trial variability, we need to consider other mechanisms.

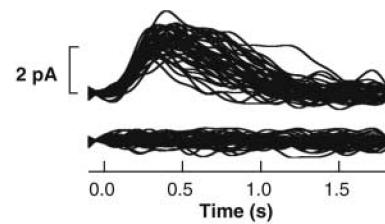


Figure 79: **Rod reliability.** Graph of 50 consecutive single-photon responses and responses to zero absorbed photons isolated from the same mammalian rod cell (from Doan, 2006).

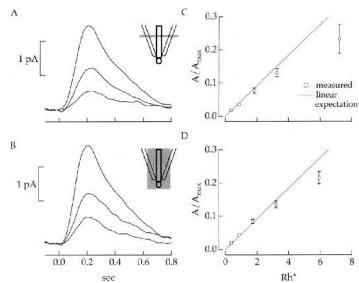


Figure 80: **Single photons do not locally saturate signal transduction inside rod cells.** The average responses of a guinea pig rod to dim flashes of 3 different strengths delivered either locally (A) or uniformly (B). The increase in response amplitude with increasing flash strength was similar in each case. Figures C and D show average results from ten guinea pig rods. Responses to both uniform and local illumination summed nearly linearly up to three activated rhodopsin molecules. The response per activated rhodopsin was essentially identical for the two stimuli. From Field and Rieke (2002).

SINGLE MOLECULES OFTEN SHOW LARGE TRIAL-TO-TRIAL FLUCTUATIONS caused by the variable duration of a molecule's active lifetime. Single ion channels can stationary and continuous variability in opening and closing, which can be recorded with a patch pipette (Fig. 81). After opening, a typical ion channel has a constant probability per unit time of closing. The stochastic lifetimes of such a channel is exponentially distributed. The mean lifetime of the "open" state would be equal to the standard deviation of the open state, resulting in a coefficient of variation equal to one.

The lifetime of activated rhodopsin should be proportional to the magnitude of signal amplification. If the pool of available transducin is not saturated, the number of activated transducin molecules should increase with rhodopsin lifetime. Variability in rhodopsin lifetime would directly cause variability in signal amplification. If the lifetime of activated rhodopsin were exponentially distributed like a typical ion channel, the coefficient of variation in the number of activated transducin in this lifetime would be one. If activated rhodopsin lifetime were more narrowly distributed, so might the coefficient of variation of this step in signal amplification.

Rhodopsin is a G protein-coupled receptor with molecular structure, signaling, and shutoff mechanisms common to many other GPCRs. These signal transduction and shut-off mechanisms may have evolved to produce signal amplification that reliably indicates the number of active receptors. Rhodopsin is inactivated by phosphorylation near its carboxyl (COOH)-terminus and subsequent binding of the molecule arrestin. Genetically deleting phosphorylation sites at the COOH-terminus of rhodopsin yields a molecule that is capable of activation but resistant to shutoff.

When normal rods that express normal rhodopsin are subjected to dim flashes – photoisomerizing one rhodopsin per trial – the electrical response in these control experiments are highly stereotyped, 1 pA in amplitude and 0.3 sec in duration. When rod cells that express both normal and truncated rhodopsin were subjected to dim flashes, many responses were the same as the control experiments, but some individual electrical responses were much higher in amplitude and much longer in duration (Fig. 83). Unlike the stereotyped normal responses, prolonged responses were variable in duration. The probability distribution of prolonged response durations resembles an exponential, suggesting that shut-off is a stochastic transition with constant probability per unit time (Fig. 84). The COOH-terminus must play an important role in regulating the amplitude and duration of the flash response. Without the COOH-terminus, rhodopsin

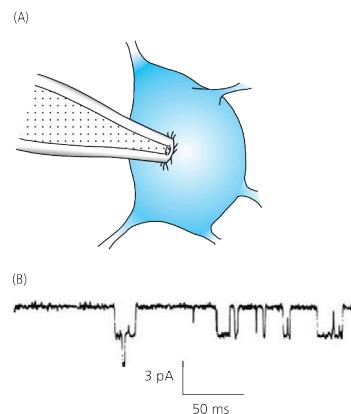


Figure 81: Patch recording single ion channels

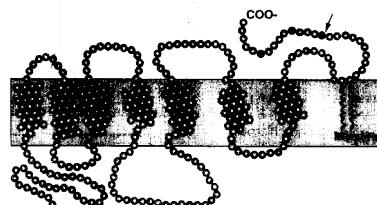


Figure 82: Rhodopsin. A cross-sectional model for rhodopsin in the disc membrane where each circle represents an amino acid. Sites of rhodopsin kinase phosphorylation include Ser³³⁴, Ser³³⁸, and Ser³⁴³ shown by filled circles.

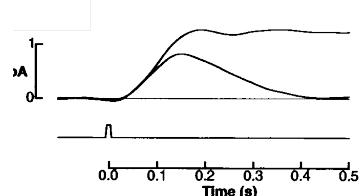


Figure 83: Rhodopsin lifetimes in truncation mutants. Electrical recording of single-photon responses in rod cells that co-express truncated and normal rhodopsin either follow their stereotyped time course or an aberrant prolonged time course.

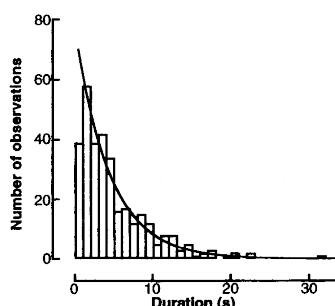


Figure 84: Rhodopsin lifetimes in truncation mutants. The histogram of prolonged rod responses in the truncation mutant reveals an exponential distribution resembling the Poisson interval distribution.

shut-off is dominated by one first-order rate-limiting step.

The COOH-terminus has multiple potential phosphorylation sites. In addition to three serine sites (Ser³³⁴, Ser³³⁸, and Ser³⁴³), rhodopsin also has three threonine sites (Thr³³⁶, Thr³⁴⁰, and Thr³⁴²). If shut-off involves sequential, step-wise phosphorylation of these sites, total rhodopsin lifetime might be the sum of the intervals between phosphorylations.

$$T = t_1 + t_2 + t_3 + \dots$$

For simplicity, assume that the mean interval between successive phosphorylations is the same, $\langle \tau \rangle$. If shut-off requires N successive phosphorylations, the mean time for shut-off will be the sum of individual mean intervals, $\langle T \rangle = N \langle \tau \rangle$. Because each phosphorylation step is independent, the variance in shut-off times, $\langle T^2 \rangle - \langle T \rangle^2$, will be the sum of the variances in interval times. If each phosphorylation step occurs with the same kinetics and, apart from having to occur in sequence, has a waiting time that is independent of other steps, the summed variances of interval times is N times the variance of one interval time:

$$\langle T^2 \rangle - \langle T \rangle^2 = N (\langle \tau^2 \rangle - \langle \tau \rangle^2)$$

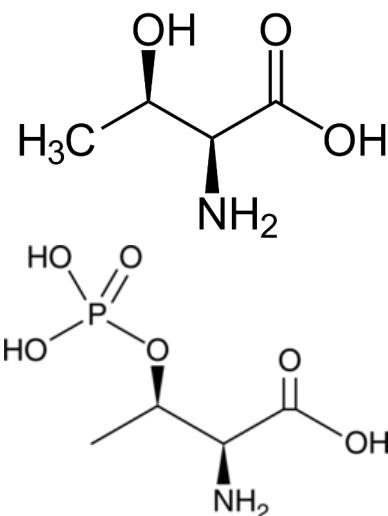


Figure 85: Threonine can be phosphorylated.

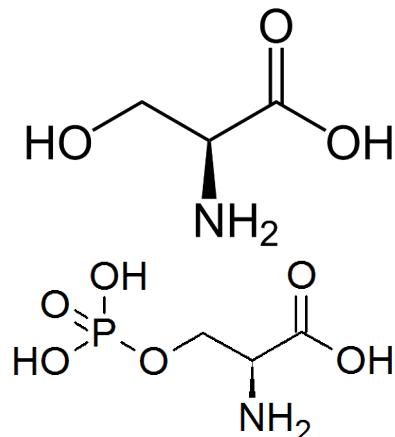


Figure 86: Serine can be phosphorylated.

IF RHODOPSIN SHUT-OFF WERE A FIRST-ORDER STOCHASTIC PROCESS, the coefficient-of-variation of rhodopsin lifetime and the integrated rhodopsin activity that is transduced into electrical events are one.

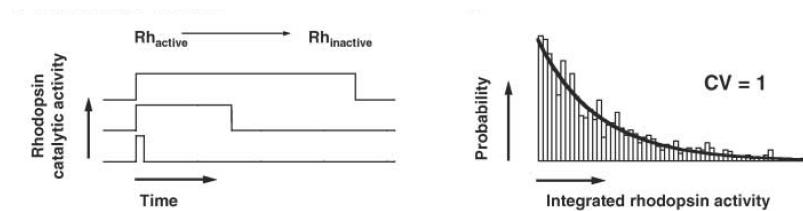


Figure 87: Simulated activity of rhodopsin after a single stochastic shutoff step

IF RHODOPSIN SHUT-OFF WERE A MULTI-STEP SEQUENCE WHERE ONLY THE FINAL STEP AFFECTS RHODOPSIN ACTIVITY, the distribution of rhodopsin lifetimes as well as integrated rhodopsin activity would exhibit smaller coefficients of variation.

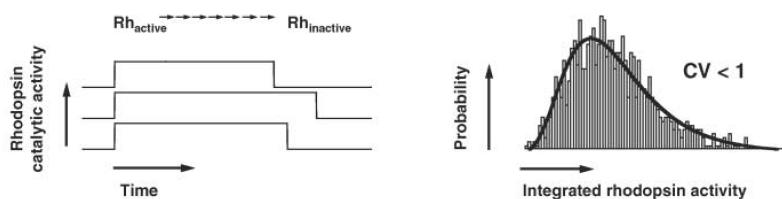


Figure 88: Simulated activity of rhodopsin with seven independent steps with equal rate constants.

Single-photon responses in mutant rod cells where the COOH-terminus of rhodopsin has been mutated to have different numbers of phosphorylation sites are both more variable and more prolonged.

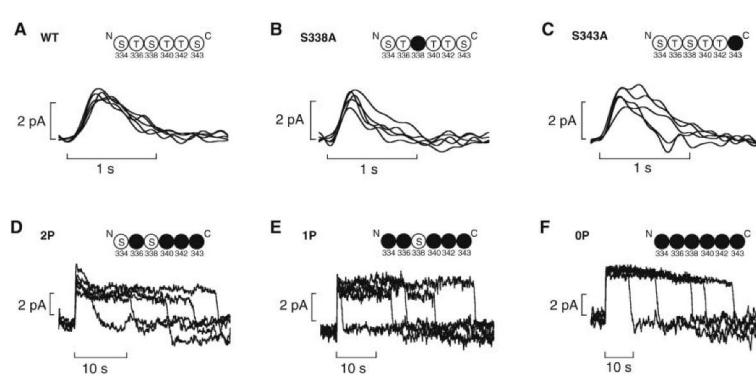


Figure 89: Examples of single-photon responses produced by wild-type and mutated rhodopsin. Five identified single-photon responses from a mouse rod expressing (A) wild-type (WT), (B and C) rhodopsin with five phosphorylation sites (S338A and S343A), (D) rhodopsin with two sites (2P), (E) rhodopsin with one site (1P), and (F) rhodopsin with zero sites (0P). Response variability increases as the number of remaining phosphorylation sites decreases.

FOR A MODEL INVOLVING MULTI-STEP SHUTOFF where each step has identical first-order kinetics, the coefficient-of-variation in shut-off times should decrease as \sqrt{N} where N is the number of steps. If the number of steps involves $N_p=6$ phosphorylations followed by arrestin binding, $N = 7$. Indeed, the coefficient-of-variation of single-photon response variability systematically depends on the number of phosphorylation sites with C.V. predicted by $1/\sqrt{N_p + 1}$.

REFERENCES

- Greg D. Field and Fred Rieke. Mechanisms regulating variability of the single photon responses of mammalian rod photoreceptors. *Neuron*, 35(4):733–747, 2002. ISSN 0896-6273 [Download paper](#)
- D A Baylor, T D Lamb, and K W Yau. Responses of retinal rods to single photons. *The Journal of Physiology*, 288(1):613–634, 1979. ISSN 0022-3751 [Download paper](#)
- K.-W YAU, G MATTHEWS, and D. A BAYLOR. Thermal activation of the visual transduction mechanism in retinal rods. *Nature*, 279(5716):806–807, 1979. ISSN 0028-0836 [Download paper](#)
- A.-C Aho, K Donner, C Hydén, L. O Larsen, and T Reuter. Low retinal noise in animals with low body temperature allows high visual sensitivity. *Nature*, 334(6180):348–350, 1988. ISSN 0028-0836 [Download paper](#)
- J CHEN, CL MAKINO, NS PEACHEY, DA BAYLOR, and MI SIMON. Mechanisms of rhodopsin inactivation in vivo as revealed by a cooh-terminal truncation mutant. *Science*, 267(5196):374–377, 1995. ISSN 0036-8075 [Download paper](#)
- Thuy Doan, Ana Mendez, Peter B Detwiler, Jeannie Chen, and Fred Rieke. Multiple phosphorylation sites confer reproducibility of the rod's single-photon responses. *Science*, 313(5786):530–533, 2006. ISSN 0036-8075 [Download paper](#)

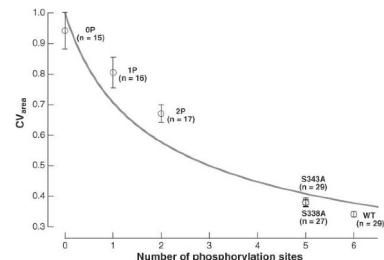


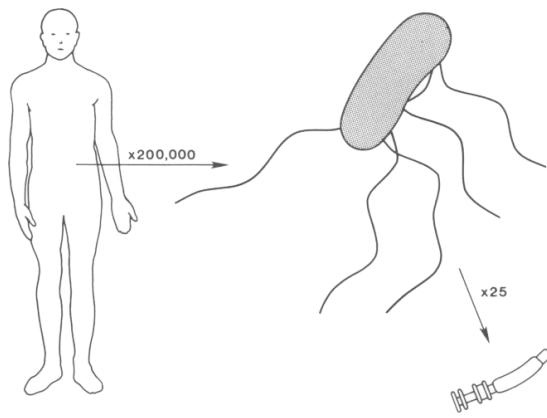
Figure 90: Correlation of single-photon response variability with number of rhodopsin phosphorylation sites.

OLFACTION AND DIFFUSION

MOTILE BACTERIA ACTIVELY SEEK FOOD BY PERFORMING CHEMOTAXIS. To perform chemotaxis, *E. coli* assesses the local concentration of food molecules by smell. Chemoreceptors on the cell surface transiently bind small molecules like amino acids. The activity of these metabotropic receptors communicates olfactory information into the cell through a signal processing cascade. This signal processing cascade regulates swimming behavior.

E. coli swims by rotating 5-6 flagellar filaments with molecular motors located at their base. The rod-shaped bacterial cell is $\sim 2 \mu\text{m}$ long and $\sim 1 \mu\text{m}$ wide. A flagellar filament is a left-handed helix with $\sim 2.3 \mu\text{m}$ pitch and $\sim 2.3 \mu\text{m}$ radius. When all flagellar motors turn counterclockwise (as viewed from outside the cell), the filaments form a helical bundle that propels the cell in a roughly straight 'run'. When one or more motors turn clockwise, the bundle falls apart and the cell tumbles in place until starting a new run in a new direction.

How does the bacterium use olfactory information to modulate CW and CCW flagellar rotation to reliably ascend spatial gradients of chemical attractants? This chemotactic feat requires solving significant physical problems. The bacterium generates net chemotactic movement up gradients by counting odor molecules with riotous and random thermal movements, while the bacterium swims with its own random walk with alternating runs and tumbles.



In an isotropic environment without chemical gradients, *E. coli* wanders in all directions. Run durations are randomly chosen from an exponential distribution (Fig. 93). This is expected with a Poisson interval distribution where the cell has a constant probability per unit time of ending each run with a tumble. During a tumble, one

Updated: November 14, 2023

Figure 91: **A sense of scale.** Comparisons of man, *E. coli*, and the flagellar rotary motor that drives swimming

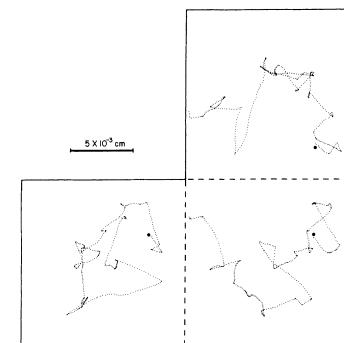


Figure 92: **Biased random walks.** A digital plot (12.6 points/sec) of the displacement of a wild-type cell (*E. coli* strain AW405) executing a random walk in a homogeneous, isotropic medium. These are planar projections of a three-dimensional track: If the left and upper panels are folded out of the page along the dotted lines, the projections appear in proper orientation on three adjacent faces of the cube. Tracking began at the large dot and continued for about 30 sec. The cell swam at the speed of $2 \times 10^{-3} \text{ cm/sec}$. There were 26 runs and tumbles.

or more flagellar rotate in the CW direction, the flagellar bundle flies apart, and the cell moves erratically until all flagella resume CCW rotation and form a new bundle. Runs typically last ~ 1 sec and tumbles last ~ 0.1 sec.

TO PERFORM CHEMOTAXIS IN A SPATIAL CHEMICAL GRADIENT, bacteria count surrounding molecules during each run and modulate the probability per unit time of starting a tumble. When *E. coli* swims in a spatial gradient of attractive molecules, runs that carry it up the gradient are extended, runs that carry it down the gradient are not. Although each tumble randomly determines the direction of the next run, the ‘biased random walk’ strategy produces net movement in the favorable direction over time. Runs are pointed in all directions. Runs that are pointed in a favorable direction are selectively lengthened.

For the biased random walk strategy to be effective, the bacterium has to assess ambient odor concentrations within ~ 1 sec, the duration of a typical run. Within ~ 1 sec, the bacterium has to ‘decide’ whether it wants to prolong its current run. One problem is that the molecules that the bacterium is trying to count are also moving randomly. Even if the cell were a ‘perfect monitor’, capable of an instantaneous and error-free census of all molecules in its vicinity, the number of counted molecules will vary from moment to moment.

Consider a perfect monitor with measurement volume V . The mean number of molecules inside the volume is $\langle N \rangle = c_0 V$, where c_0 is molecular concentration. Each census will fluctuate in number as molecules move in and out of V . These fluctuations can be described by Poisson statistics. The criteria of Poisson statistics are a large number of trials (here the number of molecules in the whole environment), low probability of success p per trial (here, the probability that each molecule is ‘successfully’ inside V , given by V divided by the volume of the whole environment), and a finite mean number of successes per trial (here, the number of molecules N inside V). With Poisson statistics, the variance in the number of counted molecules is equal to the mean:

$$\langle (N - \langle N \rangle)^2 \rangle = \langle N \rangle$$

Each census will typically differ from the mean number of molecules inside V by about the standard deviation in N . The fractional error in estimating the mean number of molecules inside V with one census is the coefficient of variation, the normalized ratio of variation from the mean δN and the mean:

$$\frac{\delta N}{N} = \sqrt{\frac{\sigma_N}{\langle N \rangle}} = \sqrt{\frac{1}{\langle N \rangle}} = \sqrt{\frac{1}{c_0 V}}$$

One census of molecules inside V provides a direct estimate of c_0 : $c_0 = N/V$. The fractional error in estimating c_0 based on one census inside V is:

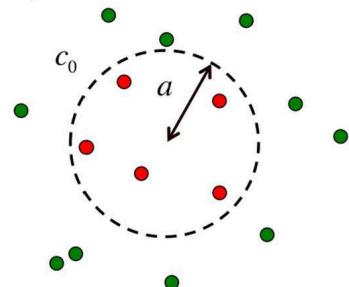


Figure 94: **The perfect monitor.** The perfect monitor is permeable to ligand molecules and estimates the concentration c_0 by counting the molecules in its volume during time T .

$$\frac{\delta c}{c_0} = \frac{1}{\sqrt{c_0 V}}$$

If a μm -sized bacterium in a 1 mM ambient concentration of molecules conducted an error-free census inside a μm -sized volume, it would count $\sim 600,000$ molecules. Census results would fluctuate around $\sim 600,000$ by $\pm \sqrt{600,000}$ or roughly ± 800 . The fractional error in measuring 1 mM molecular concentrations with one census would be $\frac{\delta c}{c_0} \sim 0.1\%$. One census by a μ -sized perfect monitor provides a reliable estimate of molecular concentrations near 1 mM. *E. coli* is effective at chemotaxis at lower molecular concentrations.

If the ambient concentration were 1 μM , the bacteria would count 600 ± 25 molecules with one census, a fractional error of 4%. If the concentration were 10 nM, the bacteria would count 6 ± 2 molecules with one census, a fractional error of 25%. Error rates at low concentrations are problematic.

How does the bacterium do better? If the bacterium made multiple independent measurements of ambient concentration across an interval of time, fractional error would decrease. This is possible because the molecules that the bacterium is trying to count are constantly moving, randomly walking in and out of the measurement volume V . After a sufficient turnover time after each census, molecules inside V will diffuse out, a new sample of molecules will diffuse in, and a new and independent census can be taken. This turnover time τ is set by diffusion. In a time interval T , the bacterium can take $M \approx T/\tau$ independent censuses. If the bacterium takes the average of M independent censuses, it would reduce uncertainty in estimating c_0 by $1/\sqrt{M}$:

$$\frac{\delta c}{c_0} = \frac{1}{\sqrt{M c_0 V}} = \frac{1}{\sqrt{T c_0 V / \tau}}$$

To calculate τ , we need the physics of diffusion at the scale of molecules and cells.

Diffusion: Microscopic Theory

Diffusion is the random migration of particles in suspension from motion due to thermal energy. A particle with absolute temperature T will have kinetic energy for each independent degree of freedom drawn from the Boltzmann distribution. The Boltzmann distribution specifies that the probability that a particle will be in state i associated with energy E_i is exponentially distributed:

$$p_i \propto \exp\left(-\frac{E_i}{kT}\right)$$

For three-dimensional movement of classical particles, we partition the state space of velocities into infinitesimal volume elements $dv_x dv_y dv_z$. The occupation probability of each element is specified by its associated energy. Considering only the x direction, the normalized probability distribution function for the x velocity of one particle depends on kinetic energy $E = \frac{mv_x^2}{2}$:

$$p(v_x) dv_x = \sqrt{\frac{m}{2\pi kT}} \exp\left(-\frac{mv_x^2}{2kT}\right) dv_x$$

Particle velocities along the x -axis are Gaussian distributed with standard deviation $\sqrt{\langle v_x^2 \rangle} = \frac{kT}{m}$. A mean kinetic energy of $kT/2$ is associated with each x , y , and z translational degree of freedom. The mean total kinetic energy will be

$$m \left\langle \frac{v^2}{2} \right\rangle = m \left\langle \frac{v_x^2}{2} \right\rangle + m \left\langle \frac{v_y^2}{2} \right\rangle + m \left\langle \frac{v_z^2}{2} \right\rangle = \frac{3kT}{2}$$

Knowing particle mass m and absolute temperature T , we know its distribution of v and the distribution of component velocities along each axis.

The probability distribution of v is derived from the probability distributions of v_x , v_y , and v_z . Movement in x , y , and z are statistically independent. The composite probability distribution for these three variables is:

$$p(v_x, v_y, v_z) dv_x dv_y dv_z = \left[\frac{m}{2\pi kT} \right]^{3/2} \exp\left(-\frac{m(v_x^2 + v_y^2 + v_z^2)}{2kT}\right) dv_x dv_y dv_z$$

The probability distribution of speed, $v = \sqrt{(v_x^2 + v_y^2 + v_z^2)}$, is obtained by converting from Cartesian to polar coordinates:

$$p(v) dv = 4\pi \left[\frac{m}{2\pi kT} \right]^{3/2} \exp\left(-\frac{mv^2}{2kT}\right) v^2 dv$$

This is called the Maxwell-Boltzmann Distribution (Fig. 95).

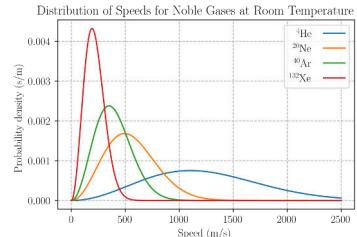


Figure 95: The Maxwell-Boltzmann distribution. The probability density functions of v of noble gases at a $T = 298.15$ K (25°C). The y-axis is in s/m so that the area under any section of the curve (which represents the probability of v being in that range) is dimensionless.

Instantaneous velocity is fast, diffusion is slow

We want to understand the movement of biological particles in and around cells. Lysozyme is an antimicrobial protein that is part of our innate immune system. The mass of one lysozyme molecule is $m = 2.3 \times 10^{-20} \text{ g}$. kT at $300 \text{ }^{\circ}\text{K}$ is $4 \times 10^{-14} \text{ g cm}^2/\text{sec}^2$. The root-mean-square velocity of lysozyme in the x direction is $\langle v_x^2 \rangle^{1/2} = 1.3 \times 10^3 \text{ cm/sec}$. The mean overall velocity v is $2.2 \times 10^3 \text{ cm/sec}$.

Unhindered at room temperature, lysozyme moves quickly across the room. In water, lysozyme quickly collides with other molecules that cause tiny and rapid changes in direction. Lysozyme does not ballistically move far in any direction. Instead, lysozyme undertakes a random walk with successive small steps in all directions. Diffusion describes the cumulative effect of these microscopic random walks. The ‘speed’ of diffusion by random walks is much slower than the instantaneous ballistic speed between collisions specified by the Boltzmann distribution.

Diffusion: microscopic theory

A MICROSCOPIC STUDY OF DIFFUSION begins by considering movement in only one dimension (Fig. 97). Particles start at time $t = 0$ at position $x = 0$ and execute random movements by a set of rules.

- Each particle steps to the right or to the left once every τ seconds, moving at velocity $\pm v_x$ a distance $\delta = \pm v_x \tau$. We treat τ and δ as constants, but they depend on particle size, the surrounding liquid, and absolute temperature T .
- At each time point, the probability of stepping to the right is $1/2$ and the probability of stepping to the left is $1/2$. The walk is not biased.
- With each step, particles forget what they did at the previous step. Each step is statistically independent.
- Each particle moves independently. Particles do not interact. This is reasonable when particles are reasonably dilute.

THESE RULES OF RANDOM MOVEMENT have two consequences. The first consequence is that particles go nowhere on average. The second consequence is that their root-mean-square displacement is not proportional to elapsed time, but to the square-root of time. Consider an

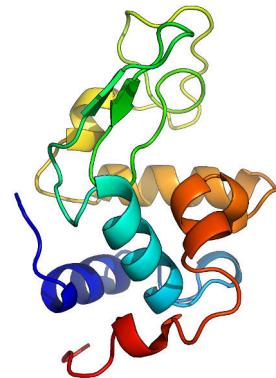


Figure 96: **Lysozyme**. An antimicrobial enzyme produced by animals that forms part of the innate immune system.



Figure 97: **1-D random walk**. Particles executing a one-dimensional random walk start at the origin, o , and move in steps of length δ , occupying positions $o, \pm\delta, \pm 2\delta, \pm 3\delta, \dots$

ensemble of N particles. Let $x_i(n)$ be the position of the i th particle after the n th step. According to the first rule, particle position after the n th step differs from position after the $(n - 1)$ th step by $\pm\delta$:

$$x_i(n) = x_i(n - 1) \pm \delta \quad (32)$$

The $+$ sign applies to roughly half of the particles. The $-$ sign applies to the rest. The mean displacement of the particles after the n th step is found by summing over particles (index $i = 1$ to N) and dividing by N :

$$\langle x(n) \rangle = \frac{1}{N} \sum_{i=1}^N x_i(n) \quad (33)$$

When we express $x_i(n)$ in terms of $x_i(n - 1)$, we find:

$$\langle x(n) \rangle = \frac{1}{N} \sum_{i=1}^N [x_i(n - 1) \pm \delta] \quad (34)$$

$$= \frac{1}{N} \sum_{i=1}^N x_i(n - 1) \quad (35)$$

$$= \langle x(n - 1) \rangle \quad (36)$$

The second term in the brackets ($\pm\delta$) averages to zero, because its sign is positive or negative with equal probability. Because $\langle x(n) \rangle = \langle x(n - 1) \rangle$, the mean position does not change from step to step. If all the particles start at the origin, their mean position stays at the origin. Particles spread symmetrically from the origin.

How quickly do particles spread? One measure of spread is root-mean-square displacement from the origin $\langle x^2(n) \rangle^{1/2}$, which is the same as the standard deviation because the mean displacement is zero. To find $\langle x^2(n) \rangle$, we write $x_i(n)$ in terms of $x_i(n - 1)$ and take the square:

$$x_i^2(n) = x_i^2(n - 1) \pm 2\delta x_i(n - 1) + \delta^2 \quad (37)$$

Then we compute the mean by averaging over all particles from $i = 1$ to N :

$$\langle x^2(n) \rangle = \frac{1}{N} \sum_{i=1}^N x_i^2(n) \quad (38)$$

After substitution,

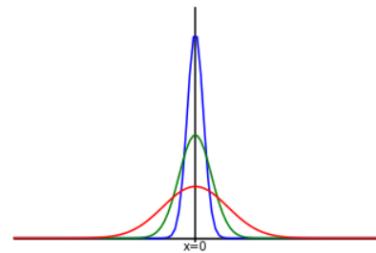


Figure 98: **Gaussian Distributions.**
The probability of finding particles at different points x at times $t = 1, 4$, and 16 sec. The standard deviations increase with the square root of time.

$$\langle x^2(n) \rangle = \frac{1}{N} \sum_{i=1}^N [x_i^2(n-1) \pm 2\delta x_i(n-1) + \delta^2] \quad (39)$$

$$= \langle x^2(n-1) \rangle + \delta^2 \quad (40)$$

As before, the second term in brackets averages to zero. Since $x_i(0) = 0$ for all particles, $\langle x^2(0) \rangle = 0$. We now use Eq. 40 iteratively to calculate the mean square position at all times. Thus, $\langle x^2(1) \rangle = \delta^2$, $\langle x^2(2) \rangle = 2\delta^2$, $\langle x^2(3) \rangle = 3\delta^2$, and so on. We conclude that the mean-square-displacement increases with the step number n .

$$\langle x^2(n) \rangle = n\delta^2$$

Root-mean-square displacement increases with \sqrt{n} . Because particles execute n steps in a time $n = t/\tau$, mean-square displacement is proportional to t and root-mean-square displacement is proportional to \sqrt{t} .

We write $x(t)$ rather than $x(n)$ to express x as a function of t .

$$\langle x^2(t) \rangle = (t/\tau)\delta^2 = (\delta^2/\tau)t \quad (41)$$

We define a new physical constant, $D = \delta^2/2\tau$, from the mathematical constants of the rules of the random walk. D has units of length²/time. The diffusion coefficient characterizes particle spread according to:

$$\langle x^2 \rangle = 2Dt \quad (42)$$

and

$$\langle x^2 \rangle^{1/2} = (2Dt)^{1/2} \quad (43)$$

There is no simple way to calculate D for a real molecule in a real liquid at a specific temperature from the mathematics of lattice random walks with step size δ and step time τ . Einstein discovered how to calculate D from the physics of statistical mechanics and fluid mechanics of real particles in real liquids at temperature T .

For a small molecule in water at room temperature, like the amino acids that bacteria try to smell, $D = 10^{-5}\text{cm}^2/\text{sec}$. Such a particle will diffuse a distance $x = 10^{-4}\text{cm}$ in a time $t \approx x^2/2D = 5 \times 10^{-4}\text{sec}$. This particle diffuses a distance $x = 1\text{ cm}$ in a time $t \approx 5 \times 10^4\text{sec}$ or 14 hours. But this particle diffuses a distance $x = 1\mu\text{m}$, the length of a bacterial cell, in $t \approx 1\text{ms}$. The ‘speed’ of diffusion depends on

the distance that needs to be traveled. The rapidity of diffusion at the size scale of the bacterial cell has pluses and minuses.

One plus from rapid diffusion is improved signal-to-noise in counting molecules with rapid molecular turnover. It takes $\tau \sim 1$ msec for molecules to diffuse in and out of a cell-sized volume surrounding the cell. In $T = 1$ sec, the duration of a typical run of a swimming *E. coli*, the bacterium as a ‘perfect monitor’ of a local cell-sized volume can take $M \approx T/\tau = 1000$ independent censuses. Taking the average of M independent censuses by a perfect monitor, uncertainty in estimating the ambient concentration c_0 is reduced by the factor $1/\sqrt{M}$.

With bacterium size a , perfect monitor volume $V = a^3$, turnover time $\tau \approx a^2/D$, integration time T , and ambient concentration c_0 of molecules with diffusion coefficient D , fractional error in estimating concentration by the perfect monitor becomes:

$$\frac{\delta c}{c_0} = \frac{1}{\sqrt{Tc_0aD}}$$

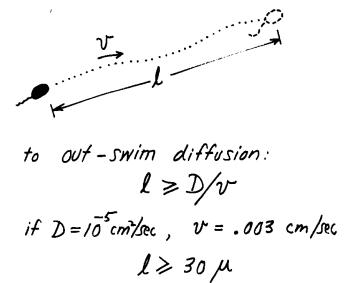
Error is reduced with longer integration times. Error is reduced when estimating higher ambient concentrations. Error is reduced with larger perfect monitors. And error is reduced with higher diffusion coefficients that allow more independent measurements.

Another plus is that signaling within the cell is fast and efficient. Intracellular signaling molecules must diffuse inside the cell from the chemoreceptors that count ambient molecules to the flagellar motors that drive locomotion. These signaling molecules will diffuse from end-to-end of the cell in a millisecond. Signal transduction inside micrometer-sized cells, not just *E. coli* but any living cell, can be fast because diffusion is rapid at short distances.

The speed of diffusion directly affects how bacteria assess ambient gradients. The speed of diffusion makes it difficult to make instantaneous spatial comparisons from end-to-end decide whether it is pointed up gradients. A bacterium swimming up a gradient will have a hard time deciding that there are instantaneously more molecules at its front end than its back end when the ambient molecules that it is trying to count move from back to front in just a millisecond. It would be difficult to make spatial comparisons on the size scale of micrometers because diffusion can erase spatial gradients on this scale within milliseconds.

Instead of spatial comparisons, bacteria make temporal comparisons. They monitor changes in the ambient concentration of molecules over time. The bacterium must decide whether the local concentration near the end of a run is higher than the local concentration was at the beginning of the run. For the bacterial cell to count ‘new’ molecules by the end of each run, it has to out-swim the

molecules it counted at the beginning of the run (Fig. 99). Because the displacement of the swimming bacteria grows linearly with time $\Delta x_{\text{swimming}} \sim v_{\text{speed}} t$, it eventually outruns the diffusive displacement of ambient molecules that grows as the square root of time $\Delta x_{\text{diffusion}} \sim \sqrt{Dt}$. For a bacterium that swims at $v_{\text{speed}} \sim 30 \mu\text{m/sec}$, this happens when $t > D/v_{\text{speed}}^2 \sim 1 \text{ sec}$, which seems to explain why runs are at least as long as 1 sec on average (Fig. 100).



"If you don't swim that far you haven't gone anywhere."

Figure 99: **Greener pastures.** For a bacterial cell to tell whether it has found a greener pasture, it has to move to that pasture. The cell has to 'out-swim' diffusion. This happens when it travels a distance $l \sim D/v$.

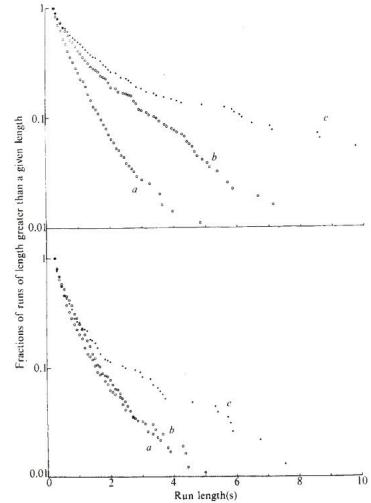


Figure 100: Runs are longer in gradients, but runs up gradients are even longer than runs down gradients. This difference is dramatic when run-length distributions are plotted for different attractants, serine and aspartate. A log-linear plot reveals that run length distributions are exponential in time, as one would expect from a Poisson interval distribution. Mean run lengths are $\sim 1 \text{ sec}$.

Two- and three-dimensional random walk

If motions in the x , y , and z directions are independent, $\langle x^2 \rangle = 2Dt$ and $\langle y^2 \rangle = 2Dt$ and $\langle z^2 \rangle = 2Dt$. In two dimensions, the square of the distance from the origin to the point (x, y) is $r^2 = x^2 + y^2$; therefore

$$\langle r^2 \rangle = 4Dt \quad (44)$$

In three dimensions, $r^2 = x^2 + y^2 + z^2$, and

$$\langle r^2 \rangle = 6Dt \quad (45)$$

The Recurrence Theorem

George Polýa discovered a striking difference between one-, two-, and three-dimensional random walks. In one and two-dimensions, a random walker on a lattice is guaranteed to visit and revisit every point in space in infinite time, no matter how far that point might be from the origin. In one and two-dimensions, a random walker that leaves the origin is guaranteed to eventually return.

In three-dimensions or more, a random walker *can* escape. A particle that leaves the origin has a finite probability of returning to origin. Thus, the particle also has a finite probability of escaping to infinity and never returning. Even with infinite time, every particle has a finite probability of reaching any particular point in space. The probability that a particle reaches a specific point in space can be calculated as its ‘capture’ probability. If the particle ever gets to that specific point, it gets captured. The probability of capture at a specific point is the probability that the particle ever wanders to that point.

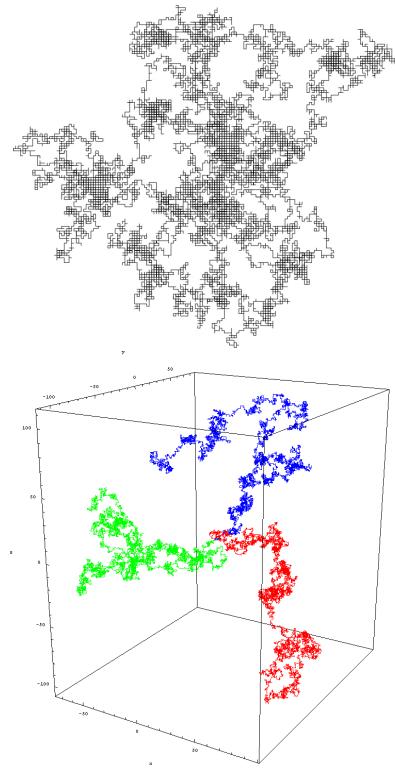


Figure 101: Random walk simulations.
Top. One random walker in two dimensions. Bottom. Three random walkers in three dimensions.

The binomial distribution

We have learned that particles undergoing free diffusion have a zero mean displacement and a root-mean-square displacement that is proportional to the square root of time. What can we say about the shape of the distribution of particles? To find out, we have to work out the probabilities that the particles step different distances to the right or to the left. It is convenient to generalize the one-dimensional random walk and suppose that a particle steps to the right with a probability p and to the left with a probability q . The probability that such a particle steps exactly k times to the right in n trials is given by the binomial distribution

$$P(k; n, p) = \frac{n!}{k!(n-k)!} p^k q^{n-k} \quad (46)$$

The displacement of the particles in n trials, $x(n)$, is equal to the number of steps to the right minus the number of steps to the left times the step length, δ :

$$x(n) = [k - (n - k)] \delta = (2k - n)\delta \quad (47)$$

Since we know the distribution of k , we know the distribution of x . The two distributions have the same shape.

The mean displacement of the particle is:

$$\langle x(n) \rangle = (2 \langle k \rangle - n)\delta \quad (48)$$

where

$$\langle k \rangle = np. \quad (49)$$

The mean-square displacement is

$$\langle x^2(n) \rangle = \langle [2k - n]\delta]^2 \quad (50)$$

$$= (4 \langle k^2 \rangle - 4 \langle k \rangle n + n^2)\delta^2 \quad (51)$$

where

$$\langle k^2 \rangle = (np)^2 + npq \quad (52)$$

For the case $p = q = 1/2$, $\langle x(n) \rangle = 0$ and $\langle x^2(n) \rangle = n\delta^2$ as expected.

The Gaussian Distribution

When n and np are both very large, the binomial distribution, $P(k; n, p)$ is equivalent to:

$$P(k)dk = \frac{1}{(2\pi\sigma^2)^{1/2}} e^{-(k-\mu)^2/2\sigma^2} dk \quad (53)$$

where $P(k)dk$ is the probability of finding a value of k between k and $k + dk$, $\mu = \langle k \rangle = np$, and $\sigma^2 = npq$. This is the Gaussian or normal distribution. By substituting $x = (2k - n)\delta$, $dx = 2\delta dk$, $p = q = 1/2$, $t = n/\tau$, and $D = \delta^2/2\tau$,

$$P(x)dx = \frac{1}{(4\pi Dt)^{1/2}} e^{-x^2/4Dt} dx \quad (54)$$

where $P(x)dx$ is the probability of finding a particle between x and $x + dx$. The variance of this distribution is $\sigma_x^2 = 2Dt$. Its standard deviation is $\sigma_x = (2Dt)^{1/2}$.

Dirac delta function

In the limit of small time ($t \rightarrow 0$), the probability distribution $P(x) = \frac{1}{(4\pi Dt)^{1/2}} e^{-x^2/4Dt}$ has the properties of the Dirac delta function. The Dirac delta function, $\delta(x)$ is a function on the real line which is zero everywhere except at the origin, where it is infinite. In our random walk derivation, all particles started at $x = 0$ at $t = 0$. Thus, the probability distribution in the limit of small time exhibited the two essential characteristics of the Dirac delta function:

$$\delta(x) \simeq \begin{cases} +\infty, & x = 0 \\ 0, & x \neq 0 \end{cases}$$

and normalization

$$\int_{-\infty}^{\infty} \delta(x) dx = 1$$

Thus, $\lim_{t \rightarrow 0} P(x, t) = \delta(x)$.

$$\int_{-\infty}^{\infty} \delta(x - a) f(x) dx = f(a)$$

There are many different heuristic definitions of the Dirac delta function for different contexts. It is useful to know that the Gaussian or normal distribution in the limit of small time works as a Dirac delta function.

Diffusion: macroscopic theory

By considering the microscopic random movements of individual particles on lattices, we derived the continuous and time-varying probability distribution of the spatial positions of particles spreading by diffusion. In one dimension, particles released from $x = 0$ at $t = 0$ initially exhibit a probability distribution of spatial positions that resembles a Dirac delta function. This probability distribution evolves over time as a Gaussian or normal distributions with standard deviation $\sigma_x = (2Dt)^{1/2}$. Instead of measuring the probability distribution of individual particles, one might observe the bulk concentration profile of many particles exhibiting diffusion, all released from the origin $x = 0$ and $t = 0$, not interacting with one another, each particle exhibiting a random walk within the surrounding fluid. Particle concentration corresponds to the number of particle per unit length (in one dimension) or number of particles per volume (in three dimensions), such that:

$$N = \int_{-\infty}^{\infty} C(x) dx$$

or

$$N = \int_V C(x, y, z) dx dy dz$$

The small number of particles, dN , in a small volume, dV , will be $dN = C(x, y, z)dV$.

NOW, WE DESCRIBE DIFFUSION BY STARTING FROM A MACROSCOPIC PERSPECTIVE, considering bulk particle flows and concentration profiles from the outset. We start with one dimension. Suppose we know the number of particles at each point along the x axis at time t , $N(x, t)$. How many particles move across unit area in unit time from the point x to the point $x + \delta$? In other words, what is net particle flux in the x direction, J_x ? After the next step, time will be $t + \tau$. Half the particles that were at x will have stepped across the dashed line from left to right. Half the particles that were at $x + \delta$ will have stepped across the dashed line from right to left.

The net number crossing to the right will be:

$$-\frac{1}{2} [N(x + \delta) - N(x)] \quad (55)$$

To obtain the net flux, we divide by the area normal to the x axis and by the time interval, τ :

$$J_x = -\frac{1}{2} [N(x + \delta) - N(x)] / A\tau \quad (56)$$

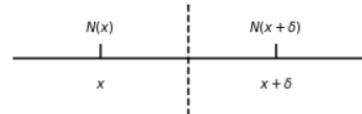


Figure 102: At time t there are $N(x)$ particles at position x , $N(x + \delta)$ particles at $x + \delta$. At time $t + \tau$, half of each set will have stepped to the right and half to the left.

Multiplying by δ^2/δ^2 and rearranging, we obtain:

$$J_x = -\frac{\delta^2}{2\tau} \frac{1}{\delta} \left[\frac{N(x + \delta)}{A\delta} - \frac{N(x)}{A\delta} \right] \quad (57)$$

The quantity $\delta^2/2\tau$ is the diffusion coefficient, D . $N(x + \delta)/A\delta$ is the number of particles per unit volume at the point $x + \delta$, i.e., the concentration $C(x + \delta)$. $N(x)/A\delta$ is concentration $C(x)$. With these substitutions,

$$J_x = -D \frac{1}{\delta} [C(x + \delta) - C(x)] \quad (58)$$

But δ is very small. In the limit $\delta \rightarrow 0$,

$$J_x = -D \frac{\partial C}{\partial x} \quad (59)$$

THIS IS FICK'S FIRST EQUATION. It states that the net flux (a vector field over x and t) is proportional to the slope of the concentration function (a scalar quantity over x and t). The constant of proportionality is D . If the particles are uniformly distributed, the slope is 0, i.e., $\partial C/\partial x = 0$ and $J_x = 0$. If the slope is constant, i.e., if $\partial C/\partial x$ is constant, J_x is constant. This occurs when C is a linear function of x .

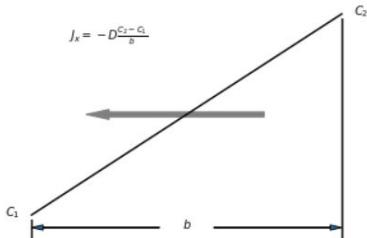


Figure 103: The flux due to a linear concentration gradient $(C_2 - C_1)/b$. There is net movement of particles from right to left solely because there are more particles at the right than at the left.

FICK'S SECOND EQUATION follows from the first, provided that the total number of particles is constant. Consider the box shown in Fig. 104. In a period of time τ , $J_x(x)A\tau$ particles will enter from the left and $J_x(x + \delta)A\tau$ particles will leave from the right. The volume of the box is $A\delta$. If particles are neither created nor destroyed, the number of particles per unit volume in the box must increase at the rate

$$\frac{1}{\tau} [C(t + \tau) - C(t)] = -\frac{1}{\tau} [J_x(x + \delta) - J_x(x)] \frac{A\tau}{A\delta} \quad (60)$$

$$= -\frac{1}{\delta} [J_x(x + \delta) - J_x(x)] \quad (61)$$

In the limit $\tau \rightarrow 0$ and $\delta \rightarrow 0$, this means that

$$\frac{\partial C}{\partial t} = -\frac{\partial J_x}{\partial x} \quad (62)$$

When we combine Fick's First and Second Laws:

$$\frac{\partial C}{\partial t} = D \frac{\partial^2 C}{\partial x^2} \quad (63)$$

Fick's second equation states that the time rate of change in concentration (at x and t) is proportional to the curvature of the concentration function (at x and t); the constant of proportionality is D . The diffusion equation tells us how a nonuniform distribution of particles will redistribute itself over time. If we specify the boundary conditions of our problem, we can calculate all later concentration distributions from an initial concentration distribution.

Here, it is useful to remind ourselves that that Gaussian distribution that we derived for individual particle probability density functions will simply be proportional to the bulk concentration of all particles over space and time. The Gaussian distribution satisfies the diffusion equation, Eq. 63, with an initial condition $t = 0$ where all particles are at the origin:

$$C(x, t) = \frac{N}{(4\pi Dt)^{1/2}} e^{-x^2/4Dt} \quad (64)$$

Diffusive flux in three dimensions

In three dimensions, we have $J_x = -D\partial C/\partial x$, $J_y = -D\partial C/\partial y$, and $J_z = -D\partial C/\partial z$. These are components of a flux vector:

$$\mathbf{J} = -D \nabla C \quad (65)$$

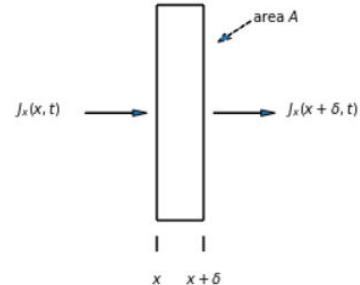


Figure 104: Fluxes through the faces of a thin box extending from position x to position $x + \delta$. The area of each face is A . The faces are normal to the x axis.

The concentration changes with time as

$$\frac{\partial C}{\partial t} = D \nabla^2 C \quad (66)$$

where ∇^2 is the three dimensional Laplacian.

If the problem is spherically symmetric, the flux is radial,

$$J_r = -D \partial C / \partial r \quad (67)$$

and

$$\frac{\partial C}{\partial t} = D \frac{1}{r^2} \frac{\partial}{\partial r} \left(r^2 \frac{\partial C}{\partial r} \right) \quad (68)$$

Diffusion with reflecting and absorbing barriers

So far, we have considered particles moving without restriction. How is particle motion affected by reflecting or absorbing walls? To build intuition about such boundary conditions affect particle motion, we return to lattice models of random walks. Suppose that particles are released at $x = 0$ on a lattice with unit steps $\pm\delta$, but with a reflecting barrier at x' (Fig. 105). For $x' > 0$, the reflecting wall is to the right of the start position. When the particle reaches $x = x'$, it has to step back to $x' - 1$ with unit probability. Without the reflecting wall, the probability distribution for particle position – in the limit of long time t , step time τ , and diffusion coefficient $D = \delta^2/2\tau$ – is:

$$P(x)dx = \frac{1}{(4\pi Dt)^{1/2}} e^{-x^2/4Dt} dx \quad (69)$$

With the wall, we need a new and different probability distribution, defined for $x < x'$ where probability of a trajectory ending at any point is augmented by trajectories that involved reflection from the wall. In Fig. 105, the time course of the position of a particle is represented in the (x, t) -plane. In this figure, the spatial displacement of the particle by a step $\pm\delta$ at every time interval τ means that its point moves upward one unit while moving laterally one unit. The probability of every path that reaches a point $x < x'$ after n reflections from the wall is larger than the probability of paths that do not reach the wall by a multiplying factor of 2^n – this is because the probability of a leftward step at each each reflection is $p = 1$ instead of $p = 1/2$. We calculate the probability of arriving at any point $x < x'$ in the presence of the wall by adding to Eq. 69 (the probability of being at $x = x'$ in the absence of the wall) the probability of reaching an “image” point at $2x' - x$ in the absence of the wall:

$$P(x, t; x') = P(x, t) + P(2x' - x, t) \quad (70)$$

This is true based on counting individual paths. Consider the path OED (Fig. 105), which involves one reflection at x' . Reflecting this path about the vertical line through x' yields a trajectory to the image point $2x' - x$. For every trajectory leading to the image point when the wall is absent – crossing the line through x' once – there is one trajectory which leads to x with one reflection. Instead of double-counting each trajectory that is reflected once at the wall to calculate the probability of arriving at x , we can instead add a unique trajectory that leads to $2x' - x$.

What if the trajectory involves two reflections? Consider the trajectory OABCD which leads to x after two reflections. The probability

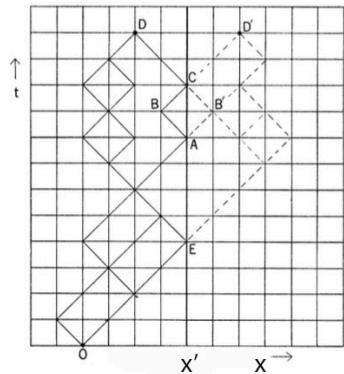


Figure 105: Trajectories in the presence of a reflecting or absorbing walls at $x = x'$.

of this trajectory should be $4\times$ larger because of the reflections. There are two trajectories (OAB'CD' and OABCD') leading to the image point and one trajectory (OAB'CD) that leads to x which we should exclude because of the barrier. These three additional trajectories together with OABCD give four trajectories leading to x or $2x' - x$ in the absence of the reflecting wall. Instead of quadruple-counting each trajectory that is reflected twice at the wall to calculate the probability of arriving at x in the presence of the reflecting wall, we can sum the trajectory that leads to x or $2x' - x$ in the absence of the wall. This argument generalizes to any number of reflections, validating Eq. 70.

In the limit of long t , the probability distribution of particle positions with a reflecting wall at $x = x'$ becomes:

$$P(x, t; x') = \frac{1}{(4\pi Dt)^{1/2}} e^{-x^2/4Dt} + \frac{1}{(4\pi Dt)^{1/2}} e^{-(2x' - x)^2/4Dt} \quad (71)$$

This probability distribution is only valid for $x < x'$. Evidently $P(x, t; x') = 0$ for $x > x'$.

It is useful to note that the according to Eq. 71,

$$\left(\frac{\partial P}{\partial x}\right)_{x=x'} = 0 \quad (72)$$

The gradient of the probability distribution (and thus the gradient of the distribution of particle concentration) is zero at a reflecting wall. The flux J_x through a reflecting wall is zero.

NOW CONSIDER AN ABSORBING WALL at $x = x'$. Whenever the particle arrives at x' , it stops moving. The probability of arriving at $x < x'$ after time t in the presence of an absorbing wall must be smaller than the probability in the absence of the wall by the removal of trajectories that crossed the wall.

Consider the probability of arriving at x after time t in the presence of the wall at $x = x'$. When counting the number of distinct trajectories that lead to x , we need to exclude all trajectories that involved one or more arrivals at $x = x'$. If we consider all sequences which lead to x in the absence of the absorbing wall, we should then exclude a number of "forbidden" trajectories. Every forbidden trajectory uniquely defines a trajectory leading to the image $2x' - x$ of x reflected across the line $x = x'$ in the (x, t) plane. Conversely, every trajectory leading to the image point $2x' - x$ in the absence of the absorbing wall yields a forbidden trajectory leading to x that crosses $x = x'$. Hence, the probability distribution of particle positions with an absorbing wall at $x = x'$ becomes:

$$P(x, t; x') = P(x, t) - P(2x' - x, t) \quad (73)$$

In the limit of long t ,

$$P(x, t; x') = \frac{1}{(4\pi Dt)^{1/2}} e^{-x^2/4Dt} - \frac{1}{(4\pi Dt)^{1/2}} e^{-(2x' - x)^2/4Dt} \quad (74)$$

According to these equations,

$$P(x = x', t; x') = 0 \quad (75)$$

The boundary condition for an absorbing wall is that particle concentration vanishes at the wall, $C = 0$.

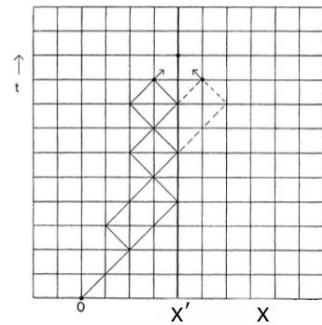


Figure 106: Particles absorbed by a wall.

Diffusion to a spherical adsorber

Consider a spherical adsorber of radius a in an infinite medium. Every particle reaching the surface of the sphere is gobbled up, so the concentration at $r = a$ is 0. The concentration at $r = \infty$ is C_0 . With these boundary conditions, the diffusion equation has the solution:

$$C(r) = C_0 \left(1 - \frac{a}{r}\right) \quad (76)$$

The flux is

$$J_r = -DC_0 \frac{a}{r^2} \quad (77)$$

The net migration of molecules is radially inward, as shown by the dashed arrows in Fig. 107. The particles are adsorbed by the sphere at a rate equal to the area, $4\pi a^2$ times the inward flux $-J_r(a)$:

$$I = 4\pi DaC_0 \quad (78)$$

We will refer to this adsorption rate, I , as a diffusion current. Note that this current is proportional not to the area of the sphere, but its radius. As the radius increases, the area increases as a^2 but the concentration gradient decreases as $1/a$.

Probability of capture

Suppose a particle is released near a spherical adsorber of radius a at a point $r = b > a$? What is the probability that the particle will be adsorbed at $r = a$ rather than wander away for good?

Consider a spherical shell source of radius b between a spherical adsorber of radius a and a spherical shell adsorber of radius c as shown in **Figure 7**. The concentration rises from 0 at $r = a$ to a maximum value C_m at $r = b$ and then falls again to 0 at $r = c$. With these boundary conditions, the diffusion equation has the solution:

$$C(r) = \begin{cases} \frac{C_m}{1-a/b} \left(1 - \frac{a}{r}\right) & \text{if } a \leq r \leq b, \\ \frac{C_m}{c/b-1} \left(\frac{c}{r} - 1\right) & \text{if } b \leq r \leq c \end{cases} \quad (79)$$

The radial flux is

$$J_r(r) = \begin{cases} \frac{-DC_m}{1-a/b} \frac{a}{r^2} & \text{if } a \leq r \leq b, \\ \frac{DC_m}{c/b-1} \frac{c}{r^2} & \text{if } b \leq r \leq c \end{cases} \quad (80)$$

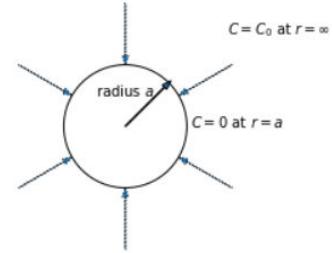


Figure 107: A spherical adsorber of radius a in an infinite medium containing particles at an initial concentration C_0 . The dashed arrows are lines of flux.

Thus, the diffusion current from the spherical shell source to the inner adsorber is

$$I_{in} = 4\pi DC_m \frac{a}{1 - a/b} \quad (81)$$

and the diffusion current from the spherical shell source to the outer adsorber is

$$I_{out} = 4\pi DC_m \frac{c}{c/b - 1} \quad (82)$$

The ratio

$$\frac{I_{in}}{I_{in} + I_{out}} = \frac{a(c - b)}{b(c - a)} \quad (83)$$

is the probability that a particle released at $r = b$ will be adsorbed at $r = a$. In the limit $c \rightarrow \infty$, this probability is just a/b . This is the probability of capture for the sphere of radius a immersed in an infinite medium. As b increases, this probability decreases as $1/b$.

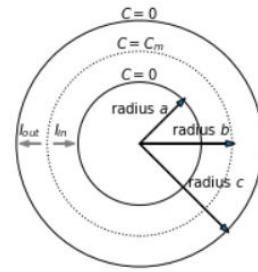


Figure 108: A spherical shell source, radius b , between a spherical adsorber of radius a and a spherical shell adsorber of radius c . Particles released at $r = b$ move inward and are adsorbed at $r = a$ at rate I_{in} or move outward and are adsorbed at $r = c$ at rate I_{out} . Their steady-state concentration rises from 0 at $r = a$ to C_m at $r = b$ and then falls again to 0 at $r = c$.

REFERENCES

- H C Berg. A physicist looks at bacterial chemotaxis. *Cold Spring Harbor Symposia on Quantitative Biology*, 53 Pt 1:1–9, 1988. ISSN 0091-7451 [Download paper](#)
- H.C. Berg and E.M. Purcell. Physics of chemoreception. *Biophysical journal*, 20(2):193–219, 1977. ISSN 0006-3495 [Download paper](#)
- Gerardo Aquino, Ned S Wingreen, and Robert G Endres. Know the single-receptor sensing limit? think again. *Journal of Statistical Physics*, 162(5):1353–1364, 2015. ISSN 0022-4715 [Download paper](#)
- Howard C. Berg. *Random walks in biology*. Princeton University Press, Princeton, N.J., rev. ed. edition, 1993. ISBN 0691000646
- S. Chandrasekhar. Stochastic problems in physics and astronomy. *Reviews of Modern Physics*, 15(1):1–89, 1943. ISSN 0034-6861

THE DIFFUSION COEFFICIENT

THE MICROSCOPIC THEORY OF DIFFUSION builds intuition about the dynamics of random walks. A key idea is that ‘root mean square’ displacement of a particle grows as the square root of time, as opposed to the ballistic movements of particles in free space where displacement is proportional to time and speed).

$$\langle x^2 \rangle^{1/2} = (2Dt)^{1/2} \quad (84)$$

For random walks on a lattice, the diffusion coefficient is $D = \delta^2/2\tau$. This diffusion coefficient, defined in terms of a unit step δ and unit interval τ , is difficult to relate to the physics of real particles in real liquids at finite temperatures. In 1905, Einstein discovered how to interrelate Brownian movements, fluid dynamics, and statistical mechanics with the mathematics of random walks. Einstein derived what is now called the Einstein-Smoluchowski relation that explains how to derive the diffusion coefficient with physical quantities:

$$D = \frac{k_B T}{f_r} \quad (85)$$

$k_B T$ is from statistical mechanics, the product of Boltzmann’s constant and temperature. The denominator is the *frictional drag coefficient*. An object under the influence of a steady force resisted by viscous drag will move with velocity proportional to applied force. The frictional drag coefficient is the constant of proportionality between this velocity and force.

TO CALCULATE THE FRICTIONAL DRAG COEFFICIENT FOR A GIVEN OBJECT, we need to use fluid mechanics. Fluid mechanics involves a set of differential equations called the Navier-Stokes Equations. Even for an incompressible fluid – i.e., with fixed fluid density and fluid viscosity – the Navier-Stokes Equations are daunting:

$$\rho_0 \left(\frac{\partial \mathbf{u}}{\partial t} + (\mathbf{u} \cdot \nabla) \mathbf{u} \right) = \eta \nabla^2 \mathbf{u} - \nabla p \quad (86)$$

The Navier-Stokes equations are a sum of all forces on each fluid volume. Inertial terms, forces due to mass and acceleration, are on the left. The terms on the right are forces due to viscous shear and pressure gradients. Solving the Navier-Stokes equations means calculating the scalar pressure field $p(x, y, z, t)$ and the fluid velocity vector field $\mathbf{u}(x, y, z, t)$. These solutions can require substantial effort with

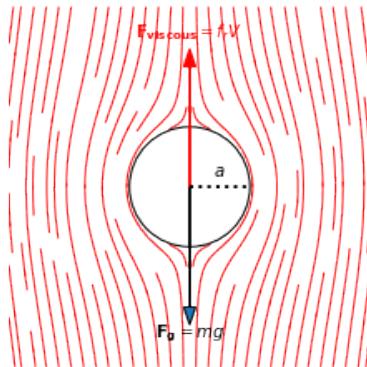


Figure 109: **Viscous drag.** An object of mass m falling in a fluid with low density and high viscosity experiences a gravitational force (mg) balanced by viscous resistance, and achieves a terminal velocity V . The proportionality between force and terminal velocity is the frictional drag coefficient, f_r that depends on the size and shape of the object and fluid viscosity. Calculating the frictional drag coefficient requires solving Navier-Stokes equations for fluid movements outside the cell. At low Reynolds numbers, when viscous forces are much larger than inertial forces, the equations of fluid dynamics are easier to solve. The frictional drag coefficient of a sphere of radius a at low Reynolds number is $f_r = 6\pi\eta a$, where η is fluid viscosity. For water, $\eta = 0.01 \text{ g cm}^{-1} \text{ s}^{-1}$

differential equations or numerical calculation. One typically solves a boundary value problem. Velocities on a bounding surfaces are specified, and the Navier-Stokes equations are used to calculate fluid movements between boundaries. Then, the forces on the boundaries, like on the surfaces of a bacterial cell and flagella, can be calculated from fluid movements.

The magnitude of forces due to viscous shear and inertia can be estimated within order of magnitude for a given problem. A given problem will be characterized by a size scale (an order of magnitude estimate a , like $\sim 1 \mu\text{m}$ for a swimming bacterium or $\sim 1 \text{ m}$ for a swimming human), velocity scale (an order of magnitude estimate v), fluid viscosity (η), and fluid density (ρ). The ratio between inertial forces and viscous forces in a given situation is a dimensionless number called the *Reynolds number*:

$$\text{Re} = \frac{\rho v a}{\eta} \quad (87)$$

For bacteria, Reynolds numbers are very small. Inertial forces are negligible in comparison to forces due to viscous shear. When inertial terms are excluded from the Navier-Stokes equations, they become linear in velocity and instantaneous (no time-derivatives):

$$\eta \nabla^2 \mathbf{u} - \nabla p = 0 \quad (88)$$

Solving these equations can be challenging, but their form allows intuitive inferences about the physics of fluid dynamics at low Reynolds number. First, these are ‘force balance’, and the equations imply that all forces are proportional to velocities \mathbf{u} and viscosities η . How much force is needed to move a sphere of radius a at velocity v ? Given that the units of viscosity are [$\text{g cm}^{-1} \text{s}^{-1}$] and the units of velocity are [cm s^{-1}], the only way to get the dimensions of force in such a way that force is proportional to velocity and viscosity is for the frictional drag coefficient of the sphere to be proportional to its radius:

$$f_r \sim \eta a \quad (89)$$

You need to properly solve the differential equations to calculate that the proportionality constant is 6π . By dimensional analysis alone, we got a pretty close estimate. The exact force we need to tow a sphere with radius a at velocity v is thus:

$$F = 6\pi\eta av \quad (90)$$

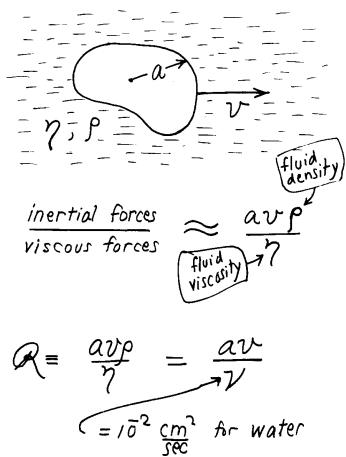


Figure 110: **The Reynolds number.**
An object moves through a fluid with velocity v . It has dimension a . In Stoke's law, the object is a sphere, but it can be anything. η and ρ are the viscosity and density of the fluid. The ratio of the inertial forces to the viscous forces, as Osborne Reynolds pointed out, is $\rho v a / \eta$. For water, $\eta \sim 0.01 \text{ g cm}^{-1} \text{s}^{-1}$ and $\rho \sim 1 \text{ g cm}^{-3}$.

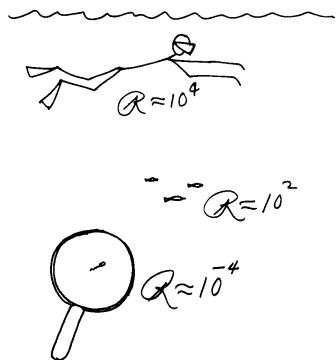


Figure 111: **Viscous drag.** We swim at high Reynolds numbers. We accelerate water behind us, which accelerates us forward. Bacteria cannot accelerate the fluid around them. They swim at low Reynolds number. They swim by taking advantage of viscous shear.

Earlier, we considered lysozyme, a small protein of typical size. We can estimate its diffusion coefficient by estimating its frictional drag coefficient. The size of a small protein like lysozyme is $\sim 1 \text{ nm}$ (Fig. 112). At room temperature, the diffusion coefficient of a small biological molecule can be estimated using the Einstein relation: $D \sim 10^{-5} \text{ cm}^2 \text{ s}^{-1}$.

THE DIFFUSION COEFFICIENT was first derived using statistical mechanics and fluid mechanics by Einstein in 1905. Here, we present a derivation by Langevin in 1908, who approached the problem by adding stochastic dynamics to Newton's equations of motion. Consider the 1-dimensional motion of a particle in fluid, with no external forces except inertia and viscosity. Absent other forces, inertial and viscous forces will sum to zero.

$$F = m \frac{d^2x}{dt^2} + f_r \frac{dx}{dt} = 0 \quad (91)$$

Substituting $v = dx/dt$ creates a simple first-order differential equation for velocity:

$$\frac{dv}{dt} + \frac{f_r}{m} v = 0 \quad (92)$$

If the initial velocity is v_0 , this differential equation is easily solved:

$$v(t) = v_0 \exp\left(\frac{-f_r t}{m}\right) \quad (93)$$

Velocity will exponentially decay with a time constant that depends on particle mass and the frictional drag coefficient. For molecules, the time constant of exponential decay is on the order of picoseconds, far shorter than any reasonable observation time. Viscous damping rapidly quenches all inertial movements.

And yet, small particles exhibit incessant Brownian movement.

The force-balance equation requires an additional force:

$$m \frac{d^2x}{dt^2} + f_r \frac{dx}{dt} = X(t) \quad (94)$$

where $X(t)$ fluctuates over time with the thermal movements of the fluid. Multiplying this force-balance equation by x and rearranging gives:

$$\frac{m}{2} \frac{d^2x^2}{dt^2} - mv^2 = -\frac{f_r}{2} \frac{dx^2}{dt} + xX(t) \quad (95)$$

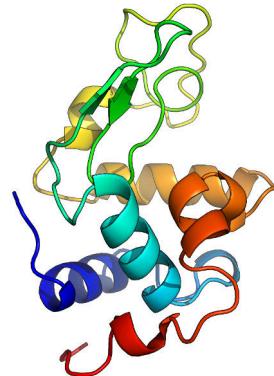


Figure 112: **Lysozyme**. An antimicrobial enzyme produced by animals that forms part of the innate immune system.

The equipartition theorem tells us the average kinetic energy of the particle.

$$\langle mv^2 \rangle = k_B T \quad (96)$$

Because the random thermal force is indifferently positive or negative, and because the diffusing particle is as likely to drift towards positive or negative x , $\langle xX(t) \rangle = 0$. Now take the ensemble average of the force balance equation. The ensemble average is an average over virtual copies of the system, each representing a possible state of the system over time. After an ensemble average, we get:

$$\frac{m}{2} \frac{d^2 \langle x^2 \rangle}{dt^2} - k_B T = -\frac{f_r}{2} \frac{d \langle x^2 \rangle}{dt} \quad (97)$$

Simplify the notation with the substitution $\frac{d \langle x^2 \rangle}{dt} = z$:

$$\frac{m}{2} \frac{dz}{dt} + \frac{f_r}{2} z = k_B T \quad (98)$$

The general solution of this “first-order inhomogeneous differential equation” is:

$$z(t) = \frac{2k_B T}{f_r} + C e^{-\frac{f_r}{m} t} \quad (99)$$

The second term exponentially decays with time. For small molecules with typical f_r and m , the time constant of this exponential decay, as we have already argued, is on the order of picoseconds. Neglecting this term gives and replacing z :

$$\frac{d \langle x^2 \rangle}{dt} = \frac{2k_B T}{f_r} \quad (100)$$

Hence, $\langle x^2 \rangle = 2 \frac{k_B T}{f_r} t$. Mean square displacement grows linearly with time, just as one expects with a random walk. The proportionality constant is the diffusion coefficient, $\langle x^2 \rangle = 2Dt$, so that:

$$D = \frac{k_B T}{f_r} \quad (101)$$

EINSTEIN DERIVED THE EINSTEIN RELATION by starting with Fick's Law. Any concentration gradient of diffusing particles will create fluxes that gradually erase those gradients. But consider the exponential atmosphere. In a gravitational field at non-zero temperatures, particles with mass m are distributed in an exponential gradient in the vertical dimension, z :

$$\rho(z) = \rho_0 e^{-\frac{mgz}{k_B T}} \quad (102)$$

Fick's First Law requires an upward diffusive flux because of Brownian movement:

$$J_z^{diff} = D \times \frac{mg}{k_B T} \times \rho_0 e^{-\frac{mgz}{k_B T}} \quad (103)$$

At the same time that gas particles are diffusing upward, they are falling downward with Earth's gravity. This sedimentation velocity is the ratio between gravitational force (mg) and frictional drag coefficient (f_r) on each particle. Downward flux (particles per unit area per unit time) due to sedimentation is the product of particle density and sedimentation velocity:

$$J_z^{sed} = -\frac{mg}{f_r} \times \rho_0 e^{-\frac{mgz}{k_B T}} \quad (104)$$

At steady-state, downward sedimentation flux must balance upward diffusive flux, $J_z^{diff} + J_z^{sed} = 0$:

$$D \times \frac{mg}{k_B T} \times \rho_0 e^{-\frac{mgz}{k_B T}} - \frac{mg}{f_r} \times \rho_0 e^{-\frac{mgz}{k_B T}} = 0 \quad (105)$$

After canceling terms:

$$\frac{D}{k_B T} - \frac{1}{f_r} = 0 \quad (106)$$

which yields the *original* derivation of the Einstein relation:

$$D = \frac{k_B T}{f_r} \quad (107)$$

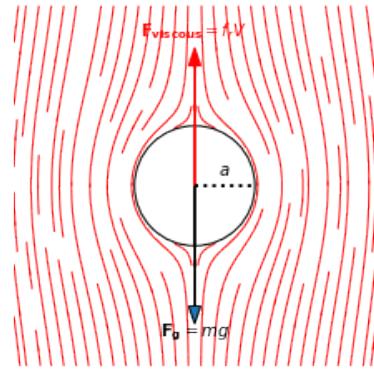


Figure 113: Sedimentation velocity.

Rotational diffusion

We now know why the runs of *E. coli*'s biased random walk are at least one second long (to outrun diffusion and to make more accurate measurements of ambient chemical concentrations). But why are runs not much longer than a second? Runs are not perfectly straight, but gently meandering. An average run drifts in orientation about 27° . This is because random thermal movements apply to rotation as well as translation. An object at thermal equilibrium will have a mean rotational kinetic energy along each axis of $kT/2$. The analysis of rotational diffusion and rotational random walks is analogous to the analysis of translation. In a rotational lattice, the random walker will step an angle $\pm\phi$ every τ seconds, yielding a mathematical diffusion coefficient $D_r = \frac{\phi^2}{2\tau}$. Instead of using Newton's laws with linear forces and translational movements in x , Langevin would have used the rotational equivalents with torques and angular movements in θ to derive the rotational diffusion coefficient $D_r = \frac{kT}{f_\theta}$. The rotational frictional drag coefficient f_θ is the constant of proportionality between applied torque N_θ and angular velocity $\Omega = d\theta/dt$.

$$\Omega = \frac{N_\theta}{f_\theta} \quad (108)$$

The mean-square angular deviation in time t is

$$\langle \theta^2 \rangle = 2D_r t \quad (109)$$

For a sphere of radius a ,

$$f_\theta = 8\pi\eta a^3 \quad (110)$$

For a two-dimensional angular random walk (Fig. 116), $\langle \theta^2 \rangle = 4D_r t$. Approximating an *E. coli* cell as a sphere of radius 10^{-4} cm, $D_r \sim 0.062$ radians² per sec. The cell wanders off course about 30° in 1 sec, as observed. This is why run duration should not be extended much more than 1 sec when swimming up gradients, because the cell will lose its way by rotational diffusion.

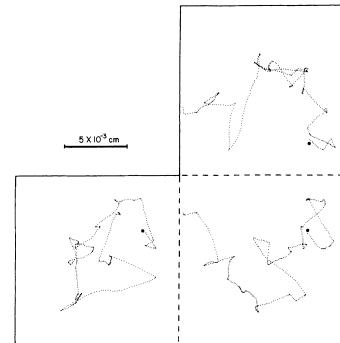


Figure 114: Bacteria perform chemotaxis by a random walk. The three-dimensional track of a single swimming bacteria viewed in xy, yz, and xz projections. The movement can be characterized as an alternating sequence of runs (periods of forward movement) and tumbles (periods of erratic rotational movement). When the bacteria is pointed in a direction it wants to go, runs get longer. The random walk becomes biased towards preferred environments.

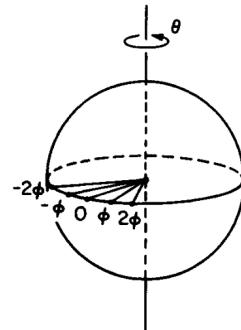


Figure 115: Rotational random walk of a sphere about one axis in steps $\pm\phi$.

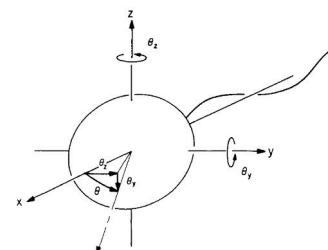


Figure 116: A microorganism propelled along the x axis by a flagellum can wander off course by rotating about either the y or z axis. The total angular deviation, θ , has components θ_y and θ_z . $\theta^2 = \theta_x^2 + \theta_y^2$.

REFERENCES

- Don S. Lemons and Anthony Gythiel. Paul Langevin's 1908 paper "On the Theory of Brownian Motion". *American Journal of Physics*, 65(11):1079–1081, 1997. ISSN 0002-9505 [Download paper](#)
- Albert Einstein. *The Motion of Particles Suspended in Liquids at Rest, According to the Kinetic Theory of Heat*. United States. Government Printing Office, 1905 [Download paper](#)
- E. M. Purcell. Life at low reynolds number. *American Journal of Physics*, 45(1):3–11, 1977 [Download paper](#)

COUNTING MOLECULES

THE PERFECT COUNTER AND THE PERFECT ADSORBER are simple models to understand how cells count molecules. The perfect counter is characterized as immediately and accurately counting all molecules in a small cell-sized volume. The perfect adsorber captures all the molecules that are brought to a cell by diffusive current. The rate at which the adsorber captures molecules is calculated with Fick's Laws.

$$I = 4\pi DaC_0 \quad (111)$$

where D is the diffusion coefficient, a is the sphere radius, and C_0 is the steady concentration at infinity. By measuring the rate of diffusive current, the cell can estimate the far-field molecule concentration.

Real cells are neither perfect counters nor adsorbers. *E. coli* counts molecules by smell. Odorant molecules transiently bind on and off receptors distributed across the cell surface. Receptor activity is coupled to both ambient molecular concentration and to a sensory transduction network within the cell.

To be a perfect adsorber, the cell needs to be completely covered by molecular receptors that capture odorant molecules without letting go. This is not what happens. Most cells – whether *E. coli* trying to do chemotaxis or any cell trying to count signaling molecules – uses a distribution of surface-bound receptors that bind and unbind molecules. The adsorber model is a simplification that allows analytical solutions to build intuition. How many receptors are needed to assess the environment? How should these receptors be distributed? To shed light on these questions, we will build an adsorber model of a patchy cell with a distribution of receptors.

Each molecule around the patchy cell exhibits a random walk. When the molecule is near the cell, it randomly bumps along the cell surface. Sometimes it will bump into a receptor. Not every molecule that is near the cell will reach the cell. In three dimensions, random walking molecules sometimes *never* bump into the cell, wandering to infinity without ever being counted. This result is from George Polya's **Recurrence Theorem** regarding one-, two-, and three-dimensional random walks. For molecules and adsorbers in one or two dimensions, every molecule eventually finds the adsorber. In three-dimensions, each molecule has a finite probability of finding the adsorber that depends on the size of the adsorber and starting distance.

The random trajectories of molecules near the cell provide in-

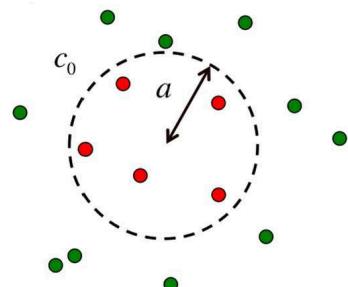


Figure 117: **The perfect monitor.** The perfect monitor is permeable to ligand molecules and estimates the concentration c_0 by counting the molecules in its volume during time T .

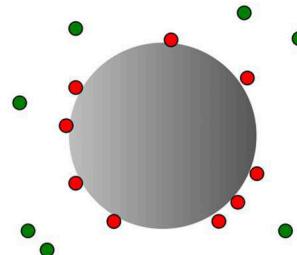


Figure 118: **The perfect adsorber.** The perfect adsorber estimates the ligand concentration from the number of molecules incident on its surface during time T .

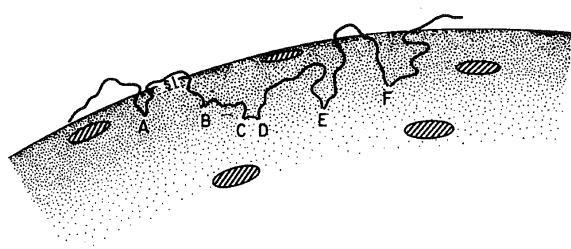


Figure 119: Random walk near a patchy sphere Path of a diffusing molecule that touches the surface of a cell at a sequence of points A, B, ... F. The cell has radius a . The receptor patches, shown shaded, are of radius s . A and B constitute independent tries at hitting a patch, but C and D do not. Note between A and B the excursion of distances perpendicular to the surface of the sphere.

tuition about why small numbers of receptors can be remarkably effective in detecting randomly walking molecules (Fig. 120). Say each receptor has a radius s . A molecule that starts its random walk a distance s from the cell has high probability of reaching the surface. This probability is the “probability of capture” that we calculated earlier using concentric spheres:

$$P_s = \frac{a}{a+s} \quad (112)$$

After a molecule that starts at radius $r = a + s$ reaches the cell at $r = a$, it wanders back to its starting radius. When it returns to radius $a + s$, it would have roughly explored an s -sized lateral region of the cell surface. If the molecule does not encounter a receptor in its first exploration, it would explore a new neighboring s -sized region of cell surface on its next exploration. Thus, each time that the molecule diffuses back and forth from the cell to radius $a + s$ will represent a new exploration of the cell surface.

Sometimes, molecules that start at $a + s$ never get to the surface. Instead of exploring for receptors, these molecules wander away to infinity, never to be seen again. The probability of failure to reach the surface in each excursion is $1 - P_s$. How many times will a molecule explore the cell surface before escaping to infinity? How many s -sized regions on the cell surface will be explored? The more regions on the cell surface that are explored, the more likely that the molecule will be recognized by cellular receptors. How does this likelihood change with the number of receptors?

The probability that a molecule at $r = a + s$ executes exactly n excursions to the surface separated by reappearances at $r = a + s$ before escaping to infinity is $P_s^n(1 - P_s)$. The average number of excursions is thus:

$$\langle n \rangle = \sum_{n=0}^{\infty} n P_s^n (1 - P_s) = \frac{P_s}{1 - P_s} = \frac{a}{s} \quad (113)$$

For a molecule exploring nanometer-sized patches on a micrometer sized cell starting a nanometer away from the surface, the average

number of excursions is ~ 1000 . The molecule has 1000 chances to find a receptor and be counted.

What is the probability that such molecule is eventually counted? To estimate this probability, we assume that each receptor is a perfect adsorber. After a molecule is adsorbed, it has been counted (whatever happens after the molecule leaves the receptor and eventually escapes to infinity; it might return to the cell and be counted again, but it only needs to be counted at least once to be counted at all). To calculate the probability that the molecule is counted at all, we start by calculating the probability that it escapes without being counted at all. In each excursion, the molecule must (a) travel to the surface with probability P_s and (b) not get adsorbed with probability $\beta = 1 - (Ns^2/4a^2)$. The probability of not getting adsorbed, β , is the fraction of the surface that lacks receptors. N receptors, each with area πs^2 , cover a total surface area of $N\pi s^2$ on a sphere with surface area $4\pi a^2$. A molecule can escape without being counted after 0, 1, 2, or more excursions. It does not matter how many excursions, n , are made, only that the molecule eventually escapes without being counted. This total escape probability is:

$$P_{esc} = \sum_{n=0}^{\infty} \beta^n P_s^n (1 - P_s) = \frac{1 - P_s}{1 - \beta P_s} \quad (114)$$

$$= \frac{4a}{4a + Ns} \quad (115)$$

We now estimate the fraction of molecules that arrive at the surface that *are* counted by at least one receptor before diffusing away to infinity, $1 - P_{esc}$. Let J_{max} be the maximum rate of counting molecules by a sphere that is completely covered by receptors. The rate at which molecules are counted by the patchy sphere is a fraction of this saturated rate:

$$\frac{J}{J_{max}} = \frac{Ns}{4a + Ns}$$

IS THERE ANOTHER WAY TO CALCULATE THE COUNTING EFFICIENCY OF A PATCHY SPHERE? Modeling the trajectory of a diffusing particle near the surface has intuitive appeal, but we did make some simplifying assumptions when making this “back of the envelope” calculation. Is there a more formal way to calculate diffusive fluxes to differently shaped adsorbers – spheres, disc-like adsorbers, patchy spheres?

Steady-state flux to adsorbing objects, where one fixes the far-field concentration of the diffusing particles to c_∞ , is calculated by solving the time-independent diffusion equation for $c(x, y, z)$ (i.e., $\partial c / \partial t = 0$):

$$\nabla^2 c = 0$$

The flux of particles is a vector field based on this concentration field:

$$\mathbf{J} = -D \nabla c$$

The total diffusive current entering a closed surface is calculated by integrating the flux around the surface:

$$I = \int_S \mathbf{J} \cdot d\mathbf{s}$$

So far, we have solved this problem in one special case of a spherical adsorber, exploiting the spherical symmetry to simplify the calculation. Developing de novo solutions for other geometries would be arduous. But others have done the work for us.

Consider the equations for electrical potential and electric field in charge free space, using cgs units where these equations are simpler (in SI units, Gauss's law is $\nabla \cdot \mathbf{E} = \frac{\rho}{\epsilon_0}$; in cgs units, Gauss's law is $\nabla \cdot \mathbf{E} = 4\pi\rho$). Laplace's Equation gives the electrical potential in charge-free space:

$$\nabla^2 \phi = 0$$

The gradient of the electrical potential gives a simple form for the electric field in charge-free space:

$$\mathbf{E} = -\nabla \phi$$

Gauss's Law tells us that the total electric charge on any closed surface is a surface integral of the electric field:

$$Q = \frac{1}{4\pi} \int_S \mathbf{E} \cdot d\mathbf{s}$$

Consider the electrical field and electrical potential surrounding a charged conductor, grounding the conductor $\phi = 0$ and setting the far-field potential to ϕ_∞ . The linearity of our equations means that the far-field electrical potential is proportional to the charge on the conductor:

$$Q = C\phi_\infty$$

where C is the capacitance of the object. In cgs units, the capacitance of a sphere is its radius, $C = a$.

Note side-by-side analogues between the equations for diffusive flux and electrostatics. I is analogous to Q . J is analogous to E . ϕ_∞ is analogous to c_∞ . The two sets of equations differ by factors of 4π and D , but are otherwise the same. Thus, solutions in electrostatics can be used as solutions for diffusive flux after correcting multiplicative factors.

The flux of diffusing particles to an adsorber is proportional to far-field concentration, $I \sim c_\infty$. With the capacitance of the same-shaped object in electrostatics, we use its capacitance to solve for the diffusive flux to the object. Correcting multiplicative factors, the flux to this object is:

$$I = 4\pi C D c_\infty$$

where C is its electrical capacitance in cgs units.

Thus, the flux to a complete adsorbing sphere (where capacitance is $C = a$) is:

$$I = 4\pi a D c_\infty$$

From electrostatics, the capacitance of a thin, conducting disk of radius b is $2\pi/b$. This gives the flux to both sides of the disk as $I = 8b D c_\infty$. If we model one adsorbing receptor as one one-sided disk on an insulating sphere, we can estimate the flux to a receptor:

$$I = 4s D c_\infty$$

If more disks are added randomly to the surface (far enough apart so that they do not affect lines of flux to one another), total flux will increase linearly with the number of disks, N . Flux can never exceed the flux to a totally adsorbing sphere. As the number of receptors increases, it must asymptotically reach the flux to a complete adsorbing sphere. What is this functional relationship?

Purcell and Berg carefully calculated the capacitance of an insulating sphere of radius a covered with a random distribution of conducting disks of radius s interconnected by vanishingly thin conducting wires (i.e., the capacitance of the patchy sphere). The capacitance of this object was:

$$C = \frac{Ns a}{Ns + \pi a}$$

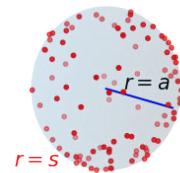


Figure 120: The patchy sphere.

The diffusive flux to the patchy sphere can be calculated using its capacitance and relative to maximum flux to a completely adsorbing sphere ($I_{max} = 4\pi Dc_\infty a$):

$$I = 4\pi Dc_\infty \frac{Ns a}{Ns + \pi a}$$

$$= I_{max} \frac{Ns}{Ns + \pi a}$$

How quickly does this function reach the asymptote? With patch radius 1 nm (an estimate of molecular size) and sphere radius 1 μm (an estimate of cellular size), flux reaches half-maximum when $N = a\pi/s$ or ~ 3100 . At half-maximum, receptors cover $\sim 0.1\%$ of cell surface. The average separation of randomly distributed patches is $\sim 60 \times$ patch radius.

A sparse distribution of receptors allows a cell to count molecules nearly as well a saturating distribution of receptors. Because diffusion gives each molecule many chances to visit the surface, sparse receptors are as likely to count each molecule closely-packed receptors. If the receptor patches were condensed into one patch that covers the same surface area on the cell as distributed receptors that reach half-maximal efficiency, counting efficiency would actually be reduced, from $I_{max}/2$ to $I_{max}/\sqrt{3100}$. Many cells count many different molecules using many different receptors. Because a cell only has to commit a small fraction of surface area to counting any one type of molecule with one type of receptor, it has plenty of room for other types of receptors. This is a deep insight into the efficiency of molecular-sensing by any micrometer-sized cell.

It turns out, however, that *E. coli* does not actually distribute its receptors randomly on its cell surface! When electron microscopy was finally used to characterize distribution of *E. coli* chemoreceptors, they were discovered to be tightly packed into single patches! The physics is not wrong. These cells with single receptor patches must have lower counting efficiency. It was later discovered that receptor patches provide other benefits that presumably offset the loss of counting efficiency.² Tightly-packed receptors can communicate with one another, amplifying chemoreceptor signals in ways that Berg and Purcell did not imagine.

There are more things in heaven and Earth, Horatio,
Than are dreamt of in your philosophy.
— Hamlet , Shakespeare

² J. R Maddock and L Shapiro. Polar location of the chemoreceptor complex in the escherichia coli cell. *Science*, 259 (5102):1717–1723, 1993. ISSN 0036-8075

REFERENCES

- H.C. Berg and E.M. Purcell. Physics of chemoreception. *Biophysical journal*, 20(2):193–219, 1977. ISSN 0006-3495 [Download paper](#)
- Gerardo Aquino, Ned S Wingreen, and Robert G Endres. Know the single-receptor sensing limit? think again. *Journal of Statistical Physics*, 162(5):1353–1364, 2015. ISSN 0022-4715 [Download paper](#)
- Steven M. Block, Jeffrey E. Segall, and Howard C. Berg. Impulse responses in bacterial chemotaxis. *Cell*, 31(1):215–226, 1982. ISSN 0092-8674 [Download paper](#)
- JE SEGALL, SM BLOCK, and HC BERG. Temporal comparisons in bacterial chemotaxis. *Proceedings of the National Academy of Sciences*, 83(23):8987–8991, 1986. ISSN 0027-8424 [Download paper](#)
- S Leibler and N Barkai. Robustness in simple biochemical networks. *Nature (London)*, 387(6636):913–917, 1997. ISSN 0028-0836 [Download paper](#)
- S Leibler, U Alon, M. G Surette, and N Barkai. Robustness in bacterial chemotaxis. *Nature (London)*, 397(6715):168–171, 1999. ISSN 0028-0836 [Download paper](#)

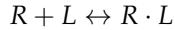
COUNTING MOLECULES WITH RECEPTORS

THE PERFECT COUNTER AND THE PERFECT ADSORBER are simplified models for how a bacterial cell counts molecules. An actual cell counts molecules with receptors. A more realistic model must consider the binding and unbinding of ligand molecules to receptors as the signal by which the cell analyzes its environment. A typical receptor has a binding site that accommodates one molecule. The occupancy of this binding site is a time-varying quantity $p(t)$, which might be '1' in the bound state and '0' in the unbound state (Fig. 121).

THE LAW OF MASS ACTION allows us to calculate the time-averaged occupancy of a receptor characterized by a single dissociation constant, K . In equilibrium at ligand concentration c , the expected average occupancy of a receptor is:

$$\langle p \rangle = \frac{c}{c + c_{1/2}}$$

Where $c_{1/2}$ is the concentration at which average occupancy is 50%. To derive this occupancy probability, we can use the Law of Mass Action to calculate the concentrations of bound and unbound receptors. Free receptors R and free ligand L bind to form bound receptors $R \cdot L$.



The rate of the forward binding reaction is proportional to both the concentration of receptor $[R]$ and ligand $[L]$. The reason for this manner of concentration dependence is probability. The probability that a given point in space contains a receptor or a ligand is proportional their concentrations. The probability that a given point in space contains both a receptor and ligand is proportional to the product of their concentration. Because receptor and ligand have to be co-localized to bind to one another, the rate of the forward binding reaction is proportional to the product of their concentrations. On the other hand, the rate of the reverse unbinding reaction is proportional to the concentration of bound receptor $[R \cdot L]$:

$$\frac{d[R \cdot L]}{dt} = k_{on}[R][L] - k_{off}[R \cdot L]$$

The mean probability \bar{p} that a given receptor is bound is the fraction of receptor in the bound state relative to total number of re-

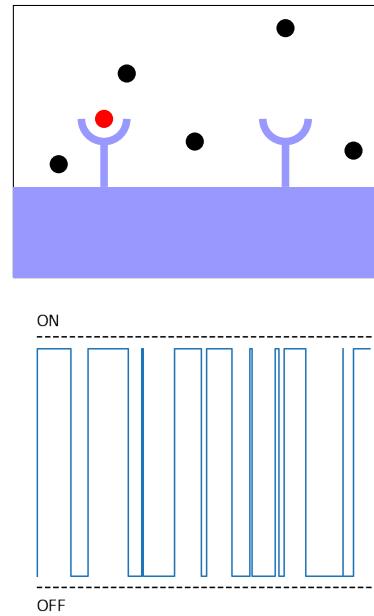


Figure 121: **Poisson receptor.** The receptor binds and unbinds ligands over time.

ceptors in the bound and unbound states. At steady-state, where $\frac{d[R \cdot L]}{dt} = 0$, we can rewrite \bar{p} in terms of ligand concentration $[L]$.

$$\begin{aligned}\bar{p} &= \frac{[R \cdot L]}{[R] + [R \cdot L]} \\ &= \frac{[L]}{[L] + k_{off}/k_{on}}\end{aligned}$$

The dissociation constant – $K_D = k_{off}/k_{on}$ – is also the concentration at which the receptor is bound with 50% probability.

AN ENTROPIC VIEW OF RECEPTOR OCCUPANCY affords another way to calculate the probability of receptor binding. We characterize an environment that contains L ligand molecules as a lattice with Ω sites. A receptor can bind one of these molecules, and when this happens only $L - 1$ ligand molecules remain on the lattice. Each of these situations can be characterized with an energy. The energy associated with no bound receptors is the total energy of L ligand molecules associating with solute, $L\epsilon_{sol}$. The energy associated with one bound receptor is the total energy of $L - 1$ ligand molecules associating with solute and one molecule bound to the receptor, $(L - 1)\epsilon_{sol} + \epsilon_b$. Each situation corresponds to a different number of equally possible arrangements of ligand molecules among lattice sites. The multiplicity of the unbound states M_u corresponds to the number of ways that L molecules can be arranged among Ω sites, which, from combinatorics, is:

$$M_u = \frac{\Omega!}{L!(\Omega - L)!}$$

For $\Omega \gg L$, this simplifies to $M_u \approx \frac{\Omega^L}{L!}$.

The multiplicity of the bound states M_b corresponds to the number of ways that $L - 1$ molecules can be arranged among Ω sites:

$$M_b = \frac{\Omega!}{(L - 1)!(\Omega - L + 1)!} \approx \frac{\Omega^{L-1}}{(L - 1)!}$$

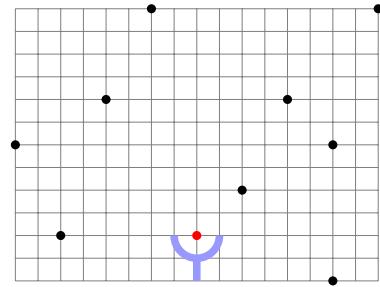
The statistical weight of the unbound state is its multiplicity M_u weighted by its Boltzmann factor: $M_u e^{-\beta L\epsilon_{sol}}$. The statistical weight of the bound state is *its* multiplicity M_b weighted by *its* Boltzmann factor: $M_b e^{-\beta[(L-1)\epsilon_{sol} + \epsilon_b]}$. The mean probability of the bound state is its statistical weight divided by the sum of the statistical weights of the bound and unbound state:

$$\bar{p} = \frac{\frac{\Omega^{L-1}}{(L-1)!} e^{-\beta[(L-1)\epsilon_{sol} + \epsilon_b]}}{\frac{\Omega^{L-1}}{(L-1)!} e^{-\beta[(L-1)\epsilon_{sol} + \epsilon_b]} + \frac{\Omega^L}{L!} e^{-\beta L\epsilon_{sol}}}$$

Multiply top and bottom by $(L!/\Omega^L)e^{\beta L\epsilon_{sol}}$:

$$\bar{p} = \frac{(L/\Omega)e^{-\beta\Delta\epsilon}}{1 + (L/\Omega)e^{-\beta\Delta\epsilon}}$$

where $\Delta\epsilon = \epsilon_b - \epsilon_{sol}$.



The overall volume of the environment is the volume of each lattice point V_{box} times the number of lattice sites Ω , so we can write ligand concentration as $c = L/\Omega V_{box}$. This results in:

$$\bar{p} = \frac{(c/c_0)e^{-\beta\Delta\varepsilon}}{1 + (c/c_0)e^{-\beta\Delta\varepsilon}}$$

where $c_0 = 1/V_{box}$.

Comparing this result with what we obtained with the Law of Mass Action provides an alternative statistical mechanical interpretation of the dissociation constant K_D .

HEARING WITH HAIR CELLS

HEARING IN VERTEBRATES is mediated by a sensory neuron for motion called the hair cell. Hair cells are used to detect movement in three contexts. Fish use hair cells along their lateral lines to sense water movement. The two other contexts are in the vertebrate ear. The vertebrate inner ear is a hollow system of ducts in skull bone. Hair cells located in ducts of the vestibular system detect head movement. Hair cells located in the cochlea, a spiral canal about 3 cm long making 2.75 turns, detect sound.

Hair cells are uniquely distinguished by the bundle of stereocilia that projects from their outward facing (or apical) surface. Deflection of these stereocilia is how movement is sensed and transformed into electrical and synaptic activity. Hair cells are depolarized (activated) by bundle deflection toward the tallest stereocilium and hyperpolarized (inactivated) by bundle deflection in the opposite direction.

FISH USE LATERAL LINES to detect water movement. Each lateral line organ is a protrusion from the skin. Each protrusion contains several hair cells with their stereocilia embedded in a gelatinous “cupula”. Half of these hair cells are oriented to be activated by forward water flow along the body. Half are oriented to be activated by backward flow.

THE ORGANS OF THE VESTIBULAR SYSTEM are in the inner ear. These organs are inside three semicircular canals oriented to three the orthogonal axes of space. A swelling in each semicircular canal contains hair cells with stereocilia embedded in a gelatinous “cupula”. Angular acceleration causes the fluid inside the semicircular canals to push or pull the cupula, activating or inactivating hair cells. Each semicircular canal preferentially detects angular acceleration around one spatial axis.

THE ORGAN OF CORTI, inside the cochlea of the inner ear, the snail-shaped canal next to the vestibular system, detects sound. The ear drum is vibrated by waves of sound pressure. These vibrations travel inward to the movable bones of the middle ear (the malleus, incus, and stapes), which vibrate the oval window at the entrance of the cochlea, setting fluid inside the cochlea in motion. Pressure waves in the cochlear fluid travel its entire length. These oscillating movements are transduced into neuronal activity by the organ of Corti, a spiral strip of epithelial cells with hair cells along its length (Fig. 126). The organ of Corti rests on a basilar membrane that divides the

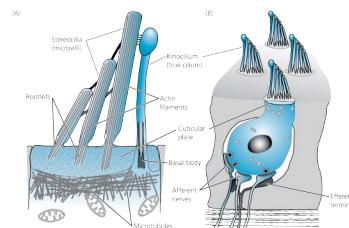


Figure 122: **Hair Cells.** (A) Top of hair cell, showing cuticular plate, stereocilia, and kinocilium. (B) Hair cell with synapses onto afferent nerves and from efferent nerves.

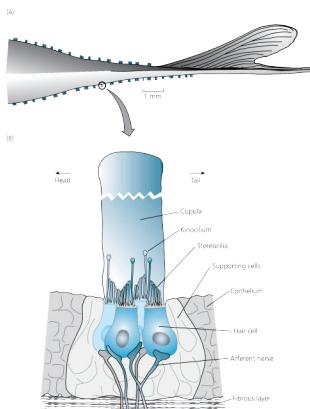


Figure 123: **Lateral line organs.** Distribution of lateral line organs alongside a minnow. Anatomy of a single organ. Two populations of hair cells have bundles oriented in opposing directions.

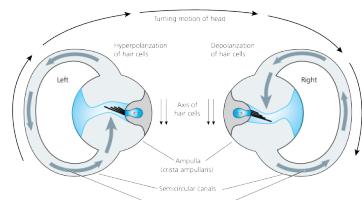


Figure 124: **Vestibular system** Organ in the inner ear that detects head position and movement

cochlea. Pressure waves push against the basilar membrane, causing it to oscillate with the frequencies of sound. The basilar membrane vibrates only if it is mechanically compliant at the frequency of sound. The stiffer the basilar membrane, the higher the frequency of its vibrations. The mechanical compliance of the basilar membrane changes from base to apex. The basilar membrane is narrow and stiff near the oval window (for high frequency sound) but wider and more compliant near the apex (for low frequency sound).

HOW DO HAIR CELLS TRANSDUCE DEFLECTION INTO ELECTRICAL SIGNALS? The 20–300 stereocilia that form hair bundles are arranged differently in different hair cells. In all hair cells, the stereocilia are arranged in order of increasing height forming a beveled shape. Stereocilia have an actin cytoskeleton, the same type of filament found in the contractile apparatus of muscle cells. At the base of each stereocilium, the actin cytoskeleton narrows into a rootlet that inserts into a plate of actin at the apical surface. Thus, when the stereocilia are deflected, they move as nearly rigid rods pivoting about their insertion point

Many hair cells also have a single largest cilium called a kinocilium, which expands into a ball at its top. The kinocilium is a ‘true’ cilium because its cytoskeleton is made from microtubules, not actin like the stereocilia. All hair cells initially have a kinocilium, which plays a role in morphogenesis and establishing hair bundle polarity during development. In the hair cells of the mammalian cochlea, the kinocilium degenerates by adulthood.

The biophysical understanding of the hair cell was pioneered by A.J. Hudspeth. Hudspeth isolated hair cells from the bullfrog and recorded their membrane potentials with an intracellular microelectrode while deflecting the stereocilia. Deflection of the stereocilia bundle toward the kinocilium caused membrane depolarization. Deflection of the stereocilia bundle away from the kinocilium caused membrane hyperpolarization. Thus, movement was directly transduced into neural activity.

Changes in the membrane potential of neurons are caused by the opening and closing of ion channels. The rapid changes in membrane potential caused by sound waves (we hear vibrations from 20 Hz to 20,000 Hz) points to a direct gating mechanism involving force-sensitive ionotropic channels. How is the movement of the stereocilia bundle motion coupled to the force acting on ion channels? The location of these channels was discovered using microscopes to image ion flow into stereociliary tips. The mechanosensitive ion channels of hair cells are both permeable to depolarizing inward flow of Na^+ and permeable to Ca^{++} that can play other roles by binding

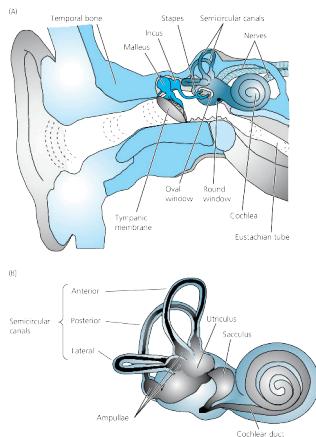


Figure 125: The cochlea. Structure of the human ear. Enlargement of components of the mammalian inner ear, showing semicircular canals, otolith organs, and cochlea.

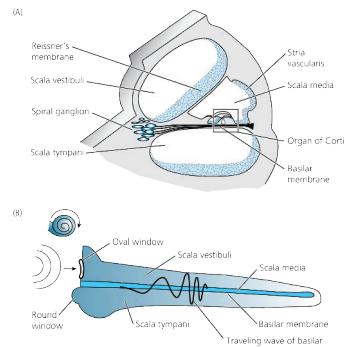


Figure 126: The cochlea (A) Cross-section of the cochlea showing the organ of Corti. (B) The cochlea imagined as uncoiled to illustrate the motion of the basilar membrane. Pressure on the oval window is communicated through the fluids of the scala vestibuli and scala tympani to the round window, producing a traveling wave of basilar membrane oscillation.

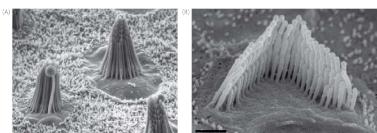


Figure 127: Scanning electron micrographs of hair bundles. (A) Bullfrog sacculus. Note kinocilium with ball at top. (B) Stereocilium bundle of an outer hair cell from the mouse cochlea. Scale bar, 1 μm .

to calcium-sensitive signaling proteins inside cells. The opening of mechanosensitive ion channels can be detected by monitoring the flow of Ca⁺⁺ into cells loaded with fluorescent Ca⁺⁺ indicating dye. These dyes change their fluorescence with changes in free-Ca⁺⁺ concentration. (Because many neurons exhibit increases in intracellular calcium concentration with increases in neuronal activity, calcium-indicating fluorescent dyes (and more recently genetically-engineered calcium indicating proteins) are used as proxies for neuronal activity involving microscopes instead of electrodes.) When stereocilia are stimulated by pushing on the hair bundle, Ca⁺⁺ influx begins at the stereociliary tips, revealing the location of the mechanosensitive channels.

How is bundle deflection turned into a force that opens the ion channels? The tops of stereocilia are interconnected by an extracellular network of cross-links that are parallel to the direction of bundle deflection that depolarizes the hair cell. These connecting fibrils are called “tip links”. The anatomy and arrangement of these tip links suggests that hair bundle deflection towards the largest stereocilia will stretch the tip links (creating a force that might pull channels open), whereas motion in the opposite direction will slacken the tip links (allowing channels to close). We now know that the location of these ion channels is near the lower insertion of each tip link, on the top of the shorter stereocilia, not the side of the taller stereocilia.

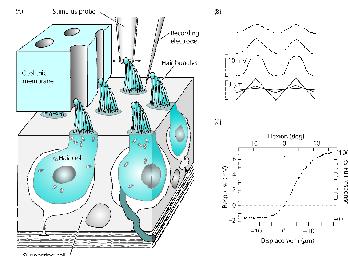


Figure 128: Scanning electron micrographs of hair bundles. (A) Bullfrog sacculus. Note kinocilium with ball at top. (B) Stereocilia bundle of an outer hair cell from the mouse cochlea. Scale bar, 1 μm .

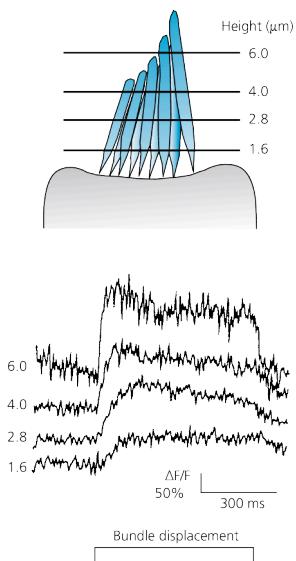


Figure 129: Location of stretch-sensitive channels in the hair cell. Change in indicator-dye fluorescence produced by the entry of Ca⁺⁺ through mechanoreceptive channels in different regions of a hair bundle, as measured with two-photon confocal microscopy. Cell had been filled with a fluorescent Ca⁺⁺ indicator dye. The plane of section of the measurement is indicated in the drawing above (note numbers next to traces below). The amplitude of fluorescence is given as a ratio of fluorescence change (ΔF) to resting fluorescence (F).

WHEN FORCE IS APPLIED TO THE HAIR BUNDLE, the bundle moves like a spring according to Hooke's law:

$$F = K_B X \quad (116)$$

where F is the force applied to the hair bundle, X is the distance the hair bundle moves, and K_B is the spring constant of the hair bundle. The stiffer the spring and the larger the spring constant K_B , the smaller the hair bundle movement to an applied force. The spring constant of the hair bundle K_B has multiple components. A passive component is provided by the stiffness of the stereocilium rotating about its base. If the combined spring constant of the rootlets, K_S , were all that mattered, the force-displacement curve of the hair bundle would be a Hooke's relationship with a stable point X_S with zero tension on the rootlets:

$$F = K_S(X - X_S)$$

But the tip links add tension that affects the force-displacement curve of the stereocilia bundle. Elastic elements connect the tops of stereocilia and pull on the mechanosensitive channels. These elastic elements are called gating springs, and representing the elasticity of the tip links and other proteins in the stereocilium. Finally, the channels themselves can contribute to the force-displacement curve with their own opening and closing movements. The contribution of channels to hair bundle stiffness is called gating compliance.

GATING COMPLIANCE was first measured by attaching a flexible glass fiber to the kinocilium of a bullfrog hair cell. When a thin glass rod is heated and rapidly pulled, it sharpens to a delicate needle point. The end of the needle point will move with external force, which can be carefully calibrated. The shape of the bent needle reveals how much force is being applied to it. When a needle is attached to a hair cell bundle and precisely moved, the bend in the needle conveys the force on the hair bundle and the position of the tip of the needle conveys displacement of the hair bundle. By systematically measuring changes in force ΔF with changes in hair bundle position ΔX , the spring constant of the bundle can be mapped out.

$$K_B(X) = \frac{\Delta F}{\Delta X} \quad (117)$$

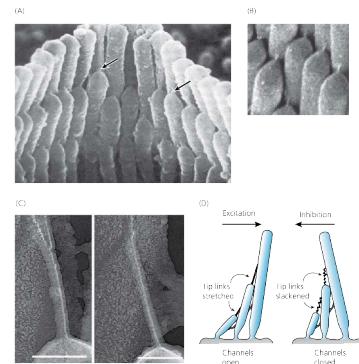


Figure 130: Tip links. (A) Scanning electron micrograph showing tip links (arrows) from outer hair cell of the guinea pig cochlea. (B) Scanning electron micrograph of a chicken auditory receptor. (C) Freeze-etch image of upper insertions of tip links of hair cells from guinea pig cochlea. Scale bars, 100 nm. (D) Proposed role of tip links in channel gating.

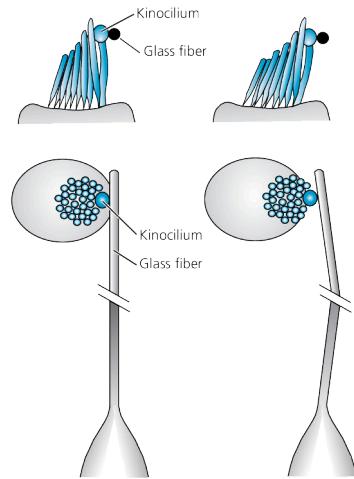


Figure 131: Measuring the gating spring. The hair bundle is moved with a fine glass fiber adhering to the kinocilium.

At the same time that the hair bundle is precisely deflected with calibrated glass needles, the membrane potential of the hair cell is measured with a separate electrode. When the hair bundle is moved toward the kinocilium, the membrane potential depolarizes, saturating with displacements of about 100 nm. Movements in the opposite direction caused the membrane potential to hyperpolarize.

Comparing the measurements of glass needle position and shape with electrical recordings reveals that changes in bundle stiffness coincide with the opening and closing of ion channels. Bundle stiffness as a function of bundle displacement, $K_B(X)$, shows that the bundle is most stiff when the channels are either entirely open (with positive deflections toward the tallest stereocilia) or entirely closed (with negative deflections toward the shortest stereocilia). Bundle stiffness is minimum when the probability of channel opening is about 0.5.

Because mechanical energy is needed to open channels, tension on the hair bundle determines the probability of channel opening – conversely, channel opening will affect hair bundle tension. When all channels are closed or open, the stiffness of the gating spring/channel apparatus will only be from the stiffness of the gating spring. If channels are actively opening and closing, the stiffness of the apparatus will be lowered by extra “give” caused by the changing shape of the channel. A quantitative model of total hair bundle stiffness K_B should incorporate the passive stiffness of the stereociliary pivots K_S , the stiffness of the gating springs K_G , and the active opening and closing of the mechanosensitive ion channels. Because the opening and closing of these ion channels is much more rapid than hair bundle movements or changes in membrane potential, their activity can be characterized as a mean probability p of the open state. The mean probability p will be a continuous function of hair bundle force and displacement. Below, we will derive this model that predicts hair bundle stiffness K_B as a function of channel open probability p . For now, we give the result of the modeling:

$$K_B = K_S + K_G - \frac{Nz^2 p(1-p)}{kT} \quad (118)$$

where N is the number of channels and z is a “gating force” for a single channel. The steady-state probability p of the channels’ being open is a sigmoidal function of bundle displacement (X) and the single-channel gating force:

$$p = \frac{1}{1 + \exp \left[-\frac{z(X-X_0)}{kT} \right]} \quad (119)$$

where X_0 is the displacement at which half the channels are open,

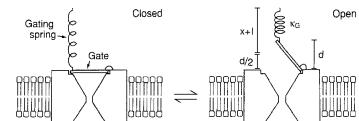


Figure 132: The conductance of a transduction channel is regulated by a molecular gate that assumes two positions, *open* and *closed*. Positive displacement of the hair bundle increases tension in the gating spring, of stiffness κ_G . When the channel is closed, the spring is extended by a distance $x + d/2$ beyond its natural length, l . Opening the channel shortens the spring by a distance d .

$p = 0.5$. At this position, the model predicts that the bundle stiffness should decrease by $Nz^2/4kT$ from the level of $\mathbf{K}_S + \mathbf{K}_G$, which is attained when the bundle is extensively deflected in either the positive or negative direction. N can be estimated from the number of stereocilia. \mathbf{K}_S can be estimated by measuring the force-displacement curve of a hair bundle without gating springs – this can be done by removing calcium from the extracellular buffer which destroys the tip links. z can be measured by fitting to the compliance and membrane potential curves.

Before we derive our model, we need to discuss some important aspects of experiment measurements and interpretations.

THE MECHANICAL AND ELECTRICAL RESPONSE of the hair bundle was simultaneously measured using a calibrated glass fiber attached to the kinocilium and an electrode in the hair cell. The base of the fiber was moved in steps of varying distance. A force of ~ 70 pN on the hair cell towards the tallest stereocilia will immediately evoke rapid positive displacement of the bundle and depolarization. These are followed by a more complex mechanical and electrical response. Immediately after the initial displacement and depolarization, the hair bundle exhibits a rapid and transient rebound in the negative direction. This rebound is accompanied by rapid termination of the electrical response. This process lasting only milliseconds is called rapid adaptation. Rapid adaptation is followed by slow relaxation of bundle position in the positive direction to a new steady-state position. Surprisingly, the new steady-state position is even *further* in the positive direction than the original rapid force-evoked displacement. This relaxation is called slow adaptation because it lasts tens of milliseconds. When the force pulse is removed, the bundle displays a rapid negative displacement, and then slowly relaxes back to the original resting position before the pulse.

The mechanisms of fast and slow adaptation are still debated. One theory is that rapid inflow of Ca^{++} causes rapid closing of the mechanosensitive channels, a mechanism for fast adaptation. Another theory is that Ca^{++} contributes to slow adaptation by mediating an adaptive change in tip link tension. An alternate mechanism for closing mechanosensitive channels is by reducing tension in the gating springs. Tension would be reduced by allowing the upper insertion of tip links to slide downward along the side of the taller stereocilium. (Similarly, the hair cell might adapt to hyperpolarizing deflection that closes mechanosensitive ion channels by increasing the tension in the gating spring by moving the tip link insertion upward).

The molecular events related to the opening and closing of ion channels, fast adaptation, and slow adaptation will be resolved when the relevant molecules are conclusively identified. There has been growing consensus that the mechanosensitive ion channels are Tmc1 and Tmc2, ion channels that have been implicated in genetic forms of human deafness. Myosin isoforms have been studied for roles in slow adaptation and active movement of the upper tip link insertion. Myosin isoforms are good candidates, being a class of motor protein known to interact with actin filaments, most notably as the contractile mechanism in muscle cells. But the identities and roles of myosin isoforms critical to hearing are still debated.

We seek a phenomenological model of the mechanoelectrical trans-

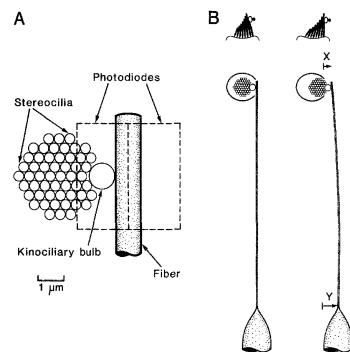


Figure 133: A hair cell was stimulated by displacing a glass fiber adhering to kinocilium. Displacement of the fiber base by Y displaces the hair bundle by X . The extent to which the calibrated fiber was bent ($Y - X$) provided a measure of stiffness

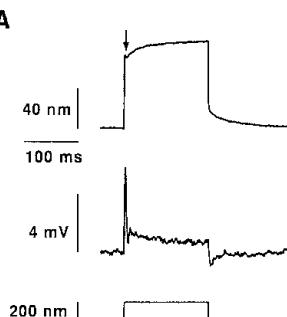


Figure 134: Apply a force of about 70 pN. Bundle moves rapidly about 74 nm. A transient reversal or rebound quickly occurs. Bundle relaxes to a new steady state with a time constant of 31 ms. After force pulse, bundle relaxes to old steady state with a time constant of 43 ms.

duction in the hair cell, how force and displacement cause direct gating of ionotropic channels in the context of the known anatomy of the hair cell. While the identity and inventory of molecules in hearing are important, our models are aimed at the physics and physiology of the hair cell, not their detailed biochemistry or molecular interactions. We seek a conceptual understanding of the mechano-electrical events that explain the measured force-displacement and displacement depolarization curves. First, we will focus on the immediate mechanical and electrical response after a force step, measured on millisecond time scale before rapid or slow adaptation. We will then discuss other interesting consequences of slow adaptation, an active process that is vital to frequency tuning and nonlinear amplification.

TO MAP THE RELATIONSHIP BETWEEN BUNDLE STIFFNESS AND DISPLACEMENT, mechanical responses were measured for positive and negative deflections of various amplitudes. Force, displacement, and electrical measurements were made between 0.75 - 1.25 ms after moving the fiber. This measurement interval represents the immediate response of the hair cell, before rapid or slow adaptation occurs. This measurement interval is also after the mechanical relaxation of the fiber tip to its new position, which requires at least 250 μ s based on the viscous drag of the fiber tip and the stiffness of glass. Because rate constants of channel closing and opening occur much faster, the open probability of mechanosensitive channels has time to reach statistical equilibrium within the observation interval and can be modeled using the Boltzmann distribution.

Bundle stiffness depends explicitly on position. As predicted by the gating spring model, an intermediate range of displacements over 125 nm is characterized by lower stiffness, reaching a minimum at a displacement 25 nm positive to rest position. For the cell illustrated in Fig. 136, the bundle minimal stiffness was 290 μ Nm $^{-1}$ smaller than the stiffness of 1150 μ Nm $^{-1}$ measured immediately after the bundle was pushed more than 50 nm in the negative direction or more than 75 nm in the positive direction.

The simultaneous recording of membrane potential in the hair cell indicates that transduction current is most sensitive to displacement over where the bundle is least stiff. Fitting the gating spring model to the transduction current as a function of displacement – using the Boltzmann relation of Eq. 119 – provides an estimate of the single channel gating force as 170 fN. This estimate is consistent to the value obtained for the single channel gating force obtained by fitting the theory to stiffness as a function of displacement using Eq. 118.

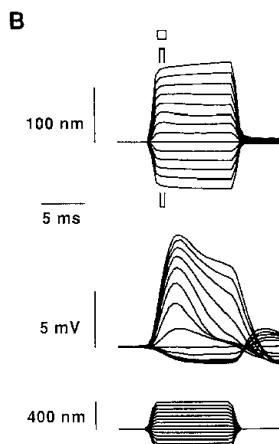


Figure 135: Apply a force of about 70 pN. Bundle moves rapidly about 74 nm. A transient reversal or rebound quickly occurs. Bundle relaxes to a new steady state with a time constant of 31 ms. After force pulse, bundle relaxes to old steady state with a time constant of 43 ms.

THE THEORETICAL BASIS FOR THE GATING COMPLIANCE DESCRIBED BY EQ. 118 IS DERIVED FROM STATISTICAL MECHANICS.

Suppose that the gating spring has a stiffness κ_G . At a given displacement of the hair bundle, let x be the extension of the gating spring midway between the open and closed states of the channel. Let the difference in the length of the gating spring between the open and closed state be d . When the channel is closed, the energy of the gating spring and channel (representing both the mechanical energy of the spring from Hooke's Law and the energy of the molecular configuration of the closed state μ_c^0) is:

$$g_c^0 = \frac{1}{2}\kappa_G(x + d/2)^2 + \mu_c^0$$

When the channel is open, the energy of the gating spring and channel is:

$$g_o^0 = \frac{1}{2}\kappa_G(x - d/2)^2 + \mu_o^0$$

The energy difference between the open and closed states is:

$$\Delta g^0 = g_o^0 - g_c^0 = -\kappa_G \times d \times x + \mu_o^0 - \mu_c^0$$

The ratio of the probabilities of the open and closed state is an exponential function of this energy difference according to the Boltzmann distribution. The probability (p) of finding the channel in the open state in equilibrium is thus:

$$p = \frac{1}{1 + \exp[-z(X - X_0)/kT]}$$

where the single channel gating force is $z = \kappa_G \times d \times \gamma$. And γ is a geometric factor (much smaller than one) that interrelates the change in the length of the nearly vertical gating spring x to a horizontal change in the displacement of the hair bundle X .

The steady-state transduction current will be proportional to p . The force exerted by a gating spring on its insertions is $f_c = \kappa_G(x - d/2)$ when the channel is open and $f_o = \kappa_G(x + d/2)$ when the channel is closed. The time average of the force exerted by one transduction element is thus

$$\begin{aligned} f &= pf_o + (1 - p)f_c \\ &= \kappa_G(x + d/2) - \kappa_Gd \times p \end{aligned}$$

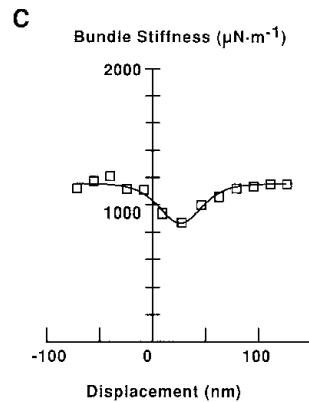


Figure 136: Stiffness smaller over a roughly 125 nm range of displacements and was minimal at 26 nm positive to resting position Minimal stiffness was 290 $\mu\text{N}/\text{m}$ smaller than the stiffness of 1150 $\mu\text{N}/\text{m}$ when pushed far negative or far positive.

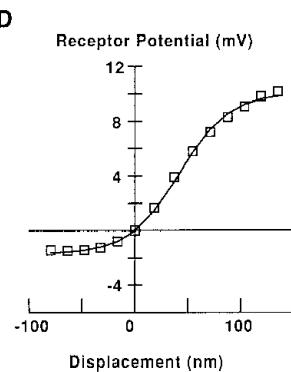


Figure 137: Sensory transduction is most sensitive for deflections where bundle is least stiff

Because γ represents the lever ratio between gating spring elongation and bundle displacement, the force exerted at the tip of the bundle (and in the direction of bundle displacement) is $N\gamma f$. N is the number of stereocilia in each bundle.

The steady-state force required to hold the hair bundle at position X

$$F = \mathbf{K}_S(X - X_S) + N\kappa_G\gamma(\gamma X + x_r + d/2) - Nz p$$

in which \mathbf{K}_S is the stiffness of the elastic components in parallel with the transduction elements such as the basal tapers and X_S is the steady-state position without gating springs.

Differentiation of F with respect to X gives the bundle stiffness:

$$\mathbf{K}_B = \mathbf{K}_S + N\kappa_G\gamma^2 - Nz^2 p(1 - p)/kT$$

Adaptation

THE HAIR CELLS ADAPTS TO CONSTANT STIMULI. After sustained deflections, the range of positions where the hair bundle is most sensitive will shift to new holding positions. The gating-spring model would predict that a shift in the position of maximum compliance should also shift the region of mechanoelectrical sensitivity. After measuring the mechanical and electrical response of a hair cell at rest position, the hair bundle was subjected to a sustained positive force. Force-displacement curves and electrical responses were measured around the new resting position. The bundle was then subjected to a negative force and the same physiological characterizations were repeated. As expected, the relation between stiffness and displacement changed with the relation between membrane potential and displacement. The position of greatest compliance closely followed the position at which membrane potential changes reached half-maximal value.

REFERENCES

- J. Howard and A.J. Hudspeth. Compliance of the hair bundle associated with gating of mechanoelectrical transduction channels in the bullfrog's saccular hair cell. *Neuron*, 1(3):189–199, 1988. ISSN 0896-6273
- Gordon L Fain. *Sensory Transduction*. Sinauer Associates, Sunderland, Mass., 2003. ISBN 0878931716

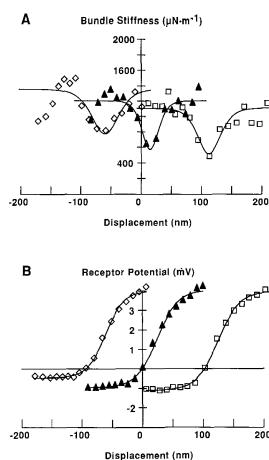


Figure 138: The position of increased compliance determines the region of mechanosensitivity. A hair cell was stimulated with force pulses. The experiment was repeated while the bundle was offset by positively or negatively directed stimuli that produced steady displacements of 103 nm and -93 nm. Gating compliance changes its position during adaptation.

NON-LINEARITIES IN HEARING

Sound stimulates a hair cell by deflecting mechanically-sensitive hair bundles. Gating springs at the tops of hair bundles direct part of this deflecting force toward opening mechanically-sensitive ion channels. Because mechanical energy is used to both deflect the hair bundle to a distance X and modulate the opening probability of ion channels p , the force-displacement curve of the hair bundle $F(X)$ depends partly on the open probability of the channels, p :

$$F(X) = \mathbf{K}_s(X - X_S) + N\kappa_G\gamma(\gamma X + x_r + d/2) - Nz p$$

A non-linear Boltzmann relationship interrelates the open probability of channels p and bundle displacement X , thereby introducing a non-linearity in the force-displacement curve.

$$p = \frac{1}{1 + \exp[-z(X - X_0)/kT]}$$

With simple springs with linear force-displacement curves ($f = kx$), sinusoidal forces at any frequency drive oscillations at the same frequency. For the sake of brevity in mathematical analysis, sinusoidal oscillations at a frequency f (cycles per second) can be described using complex numbers and Euler's formula:

$$e^{i2\pi ft} = \cos(2\pi ft) + i\sin(2\pi ft)$$

The real part of the complex number $e^{i2\pi ft}$ describes a sinusoidal oscillation in a real coordinate system.

A linear spring with linear force-displacement curve, $f = kx$, will exert sinusoidal forces in phase with any sinusoidal oscillation. Linearity means that when the spring is made to oscillate with a waveform that contains a set of sinusoidal frequencies, the spring will exert forces at the same frequencies. For example, if the spring is made to oscillate at frequencies f_1 and f_2 with amplitudes a_1 and a_2 – where the spatial coordinate is the real part of $X(t) = a_1 e^{i2\pi f_1 t} + a_2 e^{i2\pi f_2 t}$ – the spring will exert sinusoidal forces at the same frequencies – the real part of $F(t) = K a_1 e^{i2\pi f_1 t} + K a_2 e^{i2\pi f_2 t}$.

BUT WHAT IF THE SPRING WERE NON-LINEAR? Before addressing the non-linearity of the force-displacement curve of a real hair bundle, we consider a simpler non-linear spring. Near equilibrium, the force-displacement curve of a general non-linear spring is described with just the first terms of a Taylor Series:

$$F(X) = K_1 X + K_2 X^2 + K_3 X^3$$

If this spring is driven by a sinusoidal oscillation with a pure frequency, $X(t) = ae^{i2\pi ft}$, it will exert forces with sinusoidal oscillations at frequency f as well as harmonic frequencies $2f$ and $3f$ because of the additional quadratic and cubic terms. However, when this spring is simultaneously oscillated at two frequencies, $X(t) = a_1 e^{i2\pi f_1 t} + a_2 e^{i2\pi f_2 t}$, it will exert forces with additional oscillation frequencies. These additional 'combination' frequencies arise in cross-terms when expanding the quadratic and cubic terms. Combination frequencies include $f_1 + f_2$, $f_1 - f_2$, $2f_1 + f_2$, $2f_1 - f_2$, $f_1 + 2f_2$, and $f_1 - 2f_2$.

DISTORTION PRODUCTS ARE AN AUDITORY ILLUSION discovered by the composed Tartini in the eighteenth century. Tartini discovered that sounds delivered as a combination of two pure frequencies f_1 and f_2 would be heard as a combination of tones at many frequencies, not just f_1 , f_2 , and their harmonics but also various sum and difference tones. Combination tones are comparatively more faint. The most commonly heard combination frequency is the difference tone, $f_1 - f_2$, at a lower frequency than the delivered tones.

Direct measurement of the vibrations of the basilar membranes of anaesthetized chinchillas in response to pure tones and combinations of tones revealed a large range of distortion products. The magnitudes of these distortion products in direct measurement of basilar membrane motions are comparable to effects in human psychophysical studies. This suggests that the mechanism for two-tone distortion products is located in the cochlea. The essential non-linearity in the basilar membrane responses that give rise to two-tone distortion products are absent in dead or damaged cochlea, suggesting that the non-linear is an active process from live cells.



Figure 139: **Tartini tones.** Differences tones (lower red notes) between Yankee doodle (top melodic line) and a drone at high C.

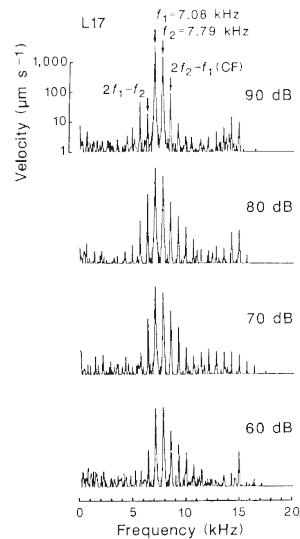


Figure 140: **Distortion products.** Frequency spectra of basilar membrane responses to two tone stimuli measured using laser velocimetry at the basal turn of the chinchilla cochlea. The spectra were obtained by Fourier analysis of responses to tone pairs at different amplitudes from 60-90 dB. The primary frequencies were chosen such that $2f_2 - f_1$ coincided with the characteristic frequency (the most sensitive frequency of the cochlear location under observation). The spectra have several peaks owing to two-tone distortion. The spectral peaks corresponding to the most prominent distortion products in psychophysical studies $2f_2 - f_1$ and $2f_1 - f_2$ are indicated by arrows.

Distortion products in the hair cell

The cochlea contains a mechanism for creating distortion products. Might the hair cell itself constitute a mechanism at the single-cell level? The approach that was used to originally map the force-displacement curve of a hair bundle – attaching a calibrated flexible glass fiber whose bend reveals the applied force and whose position reveals applied bundle displacement – was used to measure the force that a hair bundle exerts when driven with different sinusoidal oscillations and sums of different oscillations. When a hair bundle is moved back and forth at a single frequency, the hair bundle exerts force at that frequency and its harmonics (Fig. 141). When a second frequency is added to the applied displacement waveform, the hair bundle responds with a more complex pattern. In addition to forces at the primary frequencies, the hair bundle also produces forces at difference and sum frequencies ($f_2 - f_1$ and $f_2 + f_1$) and at combination frequencies ($2f_2 - f_1$ and $2f_1 - f_2$). When the hair bundle was damaged (gating springs were destroyed), the distortion products disappeared. These distortion products are a product of the active process of live hair cells.

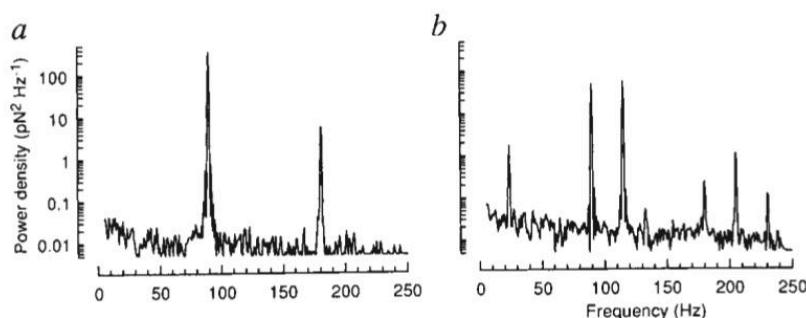


Figure 141: Distortion products. When moved back and forth, a hair bundle produced force at the stimulus frequency and the second harmonic. When a second frequency was added, the bundle produced forces at the difference and sum frequencies ($f_2 - f_1$ and $f_1 + f_2$) and at combination frequencies ($2f_1 - f_2$ and $2f_2 - f_1$).

If the non-linearity in gating compliance is responsible for distortion products, it should be possible to predict the pattern of distortion products based on an independent measurement of the non-linearity. After measuring the force-displacement curve of a hair cell, distortion products were elicited by subjecting the same hair cell to sinusoidal stimulation at two frequencies. Then, the fitted force-displacement equation was used to predict the force exerted by the hair bundle during oscillatory displacement. The theoretically-predicted power spectrum with peaks corresponding to distortion products closely matched the experimentally-measured power spectrum.

If the gating compliance is a source of distortion products, hair

cells are not only a target of traveling waves of different frequencies that travel along the basilar membrane but also a source of these waves. When a sound wave is analyzed by the cochlea, each frequency component evokes a traveling wave whose amplitude peaks at a frequency-dependent position along the basilar membrane. If hair cells account for distortion products, they must also be able to exert force on the basilar membrane, thereby activating hair cells tuned to different frequencies at different locations in the cochlea.

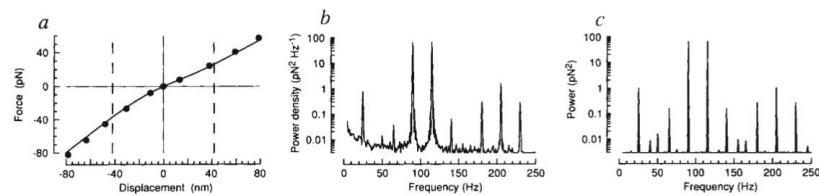


Figure 142: The relation between bundle displacement and force was determined by deflecting a bundle with stimulus pulses and measuring the flexion of the stimulus fiber. Stimulation of the hair bundle at 90 Hz and 115 Hz elicited distortion products including $f_2 - f_1 = 25$ Hz, $2f_2 - 2f_1 = 50$ Hz, $2f_1 - f_2 = 65$ Hz, $2f_2 - f_1 = 140$ Hz, $2f_1 = 180$ Hz, $f_1 + f_2 = 205$ Hz, and $2f_2 = 230$ Hz. Distortion products of similar frequency and magnitude arose in the power spectrum calculated from the bundle's measured stiffness and gating compliance. The primary frequencies and stimulus amplitudes in (c) were the same as those in (b).

REFERENCES

- F Jaramillo, V. S Markin, and A. J Hudspeth. Auditory illusions and the single hair cell. *Nature*, 364(6437):527–529, 1993. ISSN 0028-0836
- Luis Robles, Mario A Ruggiero, and Nola C Rich. Two-tone distortion in the basilar membrane of the cochlea. *Nature*, 349(6308): 413–414, 1991. ISSN 0028-0836

THE SENSITIVITY OF THE HAIR CELL

The mechanical properties of hair bundles were originally measured with glass fibers of calibrated stiffness. Brownian movement and active processes in the hair bundle can affect bundle stiffness. The viscous drag on the fiber can impair the cell's responsiveness to rapid movements. A non-invasive means of measuring the mechanical properties of hair bundles would be valuable. Winfried Denk and colleagues measured the spontaneous motion of unencumbered hair bundles with laser interferometry. An unencumbered hair bundle is subject to random forces, being surrounded by molecules exhibiting thermal movement. The hair bundle is not free to move. The simplest model of hair bundle position is a mass on a spring. If we can understand the random thermal movements of a mass on a spring, we can understand how to extract its physical properties from careful measurements on those movements (not unlike being able to measure the frictional drag coefficient of a small particle based on its Brownian movement). In the case of the hair bundle, we want to know the effective spring constant that constraints its motion as well as its frictional drag coefficient. Both quantities can be extracted from sensitive measurements of hair bundle movement over time at their smallest meaningful level caused by thermal motion.

The passive response of an object to thermal movement can be determined without stimulating the object by applying the equipartition theorem. The mechanical stiffness of a structure, such as the hair bundle, can be determined from the variance of its position fluctuations by the equipartition theorem which states that each degree of freedom with a quadratic energy term has, on average, $k_B T/2$ of energy. If x is the deflection from the average position, the mechanical stimulus of a hair bundle (κ) is

$$\kappa = \frac{k_B T}{\langle x^2 \rangle}$$

Measuring γ , the frictional drag coefficient of the hair bundle, requires going back to the Langevin equations of motion with a random forcing term. For an unrestrained particle performing Brownian movement, γ is derived from its movement fluctuations. We need to understand the movement fluctuations of a particle trapped in a harmonic potential energy well near its average position ($\kappa x^2/2$) provided by its attached spring. We need to make a detour into Fourier analysis.

Fourier Analysis

Say we measure the position of a particle constrained to one dimension over time $x(t)$. If the particle is in a harmonic potential well ($\kappa x^2/2$), it will not travel far from equilibrium at $x = 0$. From the movements of the particle over time, we can calculate the Fourier transform of those same movements. We define the Fourier transform in frequency space ($\hat{x}(f)$ where f is in cycles per second):

$$\hat{x}(f) \triangleq \int_{-\infty}^{\infty} e^{-i2\pi ft} x(t) dt$$

The Parseval–Plancherel states that the integral of a function's squared modulus is equal to the integral of the squared modulus of its frequency spectrum.

$$\int_{-\infty}^{\infty} |x(t)|^2 dt = \int_{-\infty}^{\infty} |\hat{x}(f)|^2 df$$

The squared modulus of the frequency spectrum is called the "power spectrum" and signifies the amount of energy attributed to the signal $x(t)$ at different frequencies. Another important theorem is that the Fourier transform of the time convolution of the signal is equal to this power spectrum:

$$|\hat{x}(f)|^2 = \mathcal{F}\{x^*(-t)*x(t)\} = \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} x^*(t-\tau)x(t)dt \right] e^{-i2\pi f\tau} d\tau$$

LANGEVIN INTRODUCED A RANDOM FORCING TERM, $\eta(t)$ to Newton's equations in one-dimension x (with velocity v) to derive Einstein's diffusion relationship, $D = kT/\gamma$.

$$m\dot{v} + \gamma v = \eta(t)$$

In the original paper, the effect of $\eta(t)$ was to maintain the thermally-driven movements of the particle to satisfy the equipartition theorem, but $\eta(t)$ was erased from analysis in its time-average. We know more about the mathematical structure of $\eta(t)$ than that its zero time-average. Its value at all time-points is uncorrelated ($\eta(t)$ has no memory) and so its time convolution should be zero at all times except zero. At zero, its time convolution should have positive amplitude (whose dependence on temperature and the properties of the particle we would like to determine):

$$\langle \eta(t')\eta(t) \rangle = \sigma^2 \delta(t - t')$$

If we were to solve the first-order homogeneous differential equation describing the velocity of a particle in one-dimension (i.e., where $\eta(t) = 0$), we would find:

$$v(t) \sim e^{-\gamma t/m}$$

A standard technique for solving the first-order inhomogeneous differential equation with non-zero $\eta(t) = 0$ is to guess that the solution is the product of the solution to the homogeneous different equation times an unknown function:

$$v(t) = e^{-\gamma t/m} \times f(t)$$

Substituting the guess into the original differential equation gives:

$$v(t) = \frac{e^{-\gamma t/m}}{m} \int_{-\infty}^t \eta(\tau) e^{\gamma \tau/m} d\tau$$

with rearrangement giving:

$$v(t) = \frac{1}{m} \int_{-\infty}^t \eta(\tau) e^{-\frac{\gamma}{m}(t-\tau)} d\tau$$

Changing variables $t' = t - \tau$ gives:

$$v(t) = \frac{1}{m} \int_0^\infty \eta(t-t') e^{-\frac{\gamma}{m}t'} dt'$$

The equipartition theorem dictates mean square velocity:

$$\langle v^2 \rangle = \frac{k_B T}{m}$$

The magnitude of the random forcing term (σ^2) has to be such that equipartition is satisfied. Multiplying $v(t)$ by itself gives:

$$v^2(t) = \frac{1}{m^2} \int_0^\infty \int_0^\infty \eta(t-t') \eta(t-t'') e^{-\frac{\gamma}{m}(t'+t'')} dt' dt''$$

and taking the expectation value of $v^2(t)$ gives:

$$\langle v^2(t) \rangle = \frac{1}{m^2} \int_0^\infty \int_0^\infty \langle \eta(t-t') \eta(t-t'') \rangle e^{-\frac{\gamma}{m}(t'+t'')} dt' dt''$$

But $\langle \eta(t-t') \eta(t-t'') \rangle = \sigma^2 \delta(t'-t'')$. Exploiting the properties of the δ function to calculate the integral, requiring the equipartition theorem to be satisfied, and rearranging gives:

$$\sigma^2 \int_0^\infty e^{-\frac{2\gamma t'}{m}} dt' = m k_B T$$

We conclude that $\sigma^2 = 2\gamma k_B T$ and that

$$\langle \eta(t) \eta(t') \rangle = 2\gamma k_B T \delta(t-t')$$

WE ARE NOW READY TO TACKLE THE OBJECT IN A HARMONIC POTENTIAL.

Consider the movements of an over-damped particle in a harmonic potential well, $E = \kappa x^2/2$ driven by a random thermal forces $\eta(t)$. Mass and acceleration do not matter, the force due to viscous drag is the frictional drag coefficient times velocity $-\gamma \dot{x}$, and the restoring force from the spring is $-\kappa x$:

$$\gamma \dot{x} + \kappa x = \eta(t)$$

The Fourier-transform of this equation of motion gives us a direct relationship between the Fourier transform of position and the Fourier transform of the random forcing term:

$$2\pi i \gamma f \tilde{x}(f) + \kappa \tilde{x}(f) = \tilde{\eta}(f)$$

Energy and power spectra

The *energy* of a continuous, real-valued signal $x(t)$ is defined as

$$\|x(t)\|^2 \triangleq \int_{-\infty}^{\infty} |x(t)|^2 dt = \int_{-\infty}^{\infty} |\hat{x}(f)|^2 df \quad (120)$$

where the second equality is from Plancherel's theorem. Eq. 120 will be infinite for signals that are periodic or non-localized in time or frequency. In such cases, we define the average *power* as

$$\overline{x(t)^2} \triangleq \lim_{T \rightarrow \infty} \frac{1}{T} \int_{-T/2}^{T/2} x(t)^2 dt = \lim_{T \rightarrow \infty} \frac{1}{T} \int_{-\infty}^{\infty} [x_T(t)]^2 dt \quad (121)$$

where

$$x_T(t) \triangleq \begin{cases} x(t), & |t| < T/2, \\ 0, & |t| \geq T/2, \end{cases}$$

and the notation $\overline{f(t)}$ denotes an average over all time. By analogy to (120), we may write

$$\overline{x(t)^2} = \int_{-\infty}^{\infty} S_x(f) df, \quad (122)$$

where $S_x(f) \triangleq \lim_{T \rightarrow \infty} \frac{1}{T} |\hat{x}_T(f)|^2$ is called the *power spectral density*. Letting $x_T^-(t) \triangleq x_T(-t)$, we observe

$$\begin{aligned} |\hat{x}_T(f)|^2 &= \hat{x}_T(f) \hat{x}_T^*(f) \\ &= \hat{x}_T(f) \hat{x}_T(-f) \text{ (conjugation)} \\ &= \mathcal{F}\{x_T(t)\} \mathcal{F}\{x_T(-t)\} \text{ (time reversal)} \\ &= \mathcal{F}\{x_T * x_T^-\} \text{ (convolution).} \end{aligned}$$

We therefore have

$$\begin{aligned}
 S_x(f) &= \lim_{T \rightarrow \infty} \frac{1}{T} \mathcal{F}\{x_T * x_T^-\} \\
 &= \mathcal{F}\left\{\lim_{T \rightarrow \infty} \frac{1}{T} (x_T * x_T^-)(\tau)\right\} \\
 &= \mathcal{F}\left\{\lim_{T \rightarrow \infty} \frac{1}{T} \int_{-T/2}^{T/2} x(t)x(t+\tau)dt\right\} \\
 &= \mathcal{F}\{\overline{x(t)x(t+\tau)}\} \tag{123}
 \end{aligned}$$

(124)

where $\overline{x(t)x(t+\tau)}$ is the autocorrelation function. In other words, the Fourier transform of the autocorrelation function equals the power spectral density. This result is called the Wiener-Khinchin Theorem.

REFERENCES

-