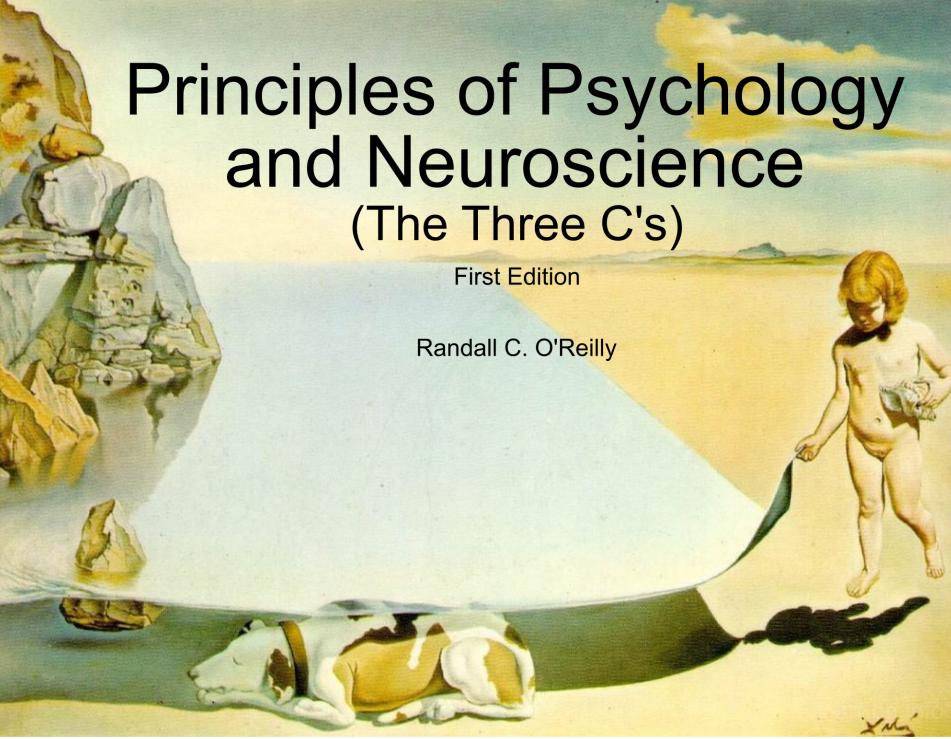


Principles of Psychology and Neuroscience, First Edition

Randall C. O'Reilly



Principles of Psychology and Neuroscience

(The Three C's)

First Edition

Randall C. O'Reilly

Xan

Contents

Preface	6
Chapter 0: Introduction	7
The Three C's	7
Compression	7
Contrast	8
Control	10
The Breakdown of Control	11
Other Principles and Perspectives	12
Where do we go from here?	12
Chapter 1: Science and Subjectivity – The Fundamental Challenge of Psychology	14
Subjectivity in Psychology: A Brief History	15
Fundamentals of Cognitive Neuroscience	16
Subjectivity and Science: Working with the Method	16
Research Methods in Psychology and Neuroscience	18
Neuroscience methods	20
Statistics	20
Conclusions	22
Summary of Key Terms	22
Chapter 2: Neuroscience	24
Simple Neurons Make Complex Work	25
The Tug-of-War in Your Brain	28
Large-Scale Brain Organization ("Gross" Anatomy)	30
The Big Brain Chunks	31
Functional Organization of the Neocortex	36
Hierarchical Organization	39
Neuroscience Methods	39
Functional Neuroimaging: fMRI, PET, EEG, MEG	40
Conclusions	41
Summary of Key Terms	41
Chapter 3: Consciousness, Drugs, Sleep, and Dreams	43
Neural correlates of consciousness	43
Features of Consciousness	43
Recurrent Connectivity as a Major NCC	44
Animal and AI Consciousness	46
Altered States	47
Neuromodulators and Drugs	47
Sleep and Dreams	49
Functions of Sleep	50
Sleep Stages	50
The Function of Dreams	52
Summary of Key Terms	53
Chapter 4: Sensation, Perception, and Attention	54
Perception is (Hierarchical) Compression	54
We See the "Real" World, not Raw Sensation	55
Sensory Systems	57
Vision	58
Compression and Contrast in Vision	60
Color Contrasts	63

Depth	63
Compression in Object Recognition	65
Time Contrast: The Novelty Filter	69
Audition	71
Attention	72
The Posner Spatial Cueing Task	72
Psychophysics	73
Summary	73
Summary of Key Terms	74
Chapter 5: Learning, Motivation, and Emotion	75
Synaptic Plasticity	75
Neocortical Learning	77
Dopamine-modulated Learning	78
Classical (Pavlovian) Conditioning	79
Operant / Instrumental Conditioning	83
Motivation	84
Goal-driven Behavior	86
Emotion and Arousal	87
Emotional / Motivational Encoding in vmPFC	90
Biological Grounding of Emotion and Arousal	90
Summary of Key Terms	92
Chapter 6: Memory	94
From Synapses to Memory	94
The Modal Model of Memory	96
The Hippocampus	98
Taxonomy of Long-Term Memory	100
Amnesia	102
Memory Capacity and the Importance of <i>Chunks</i>	103
Encoding and Retrieval Strategies (i.e., How to Study!)	104
Memory Retention and Interference	105
The Fallibility of Memory	107
Working Memory and the Prefrontal Cortex	108
Summary of Key Terms	108
Chapter 7: Thinking, Control and Intelligence	110
The Neural CPU in the Prefrontal Cortex and Basal Ganglia	111
What it takes to be a Computer	114
Individual Differences in Prefrontal Cortex / Basal Ganglia?	115
Strengths, Weaknesses, and Biases of our Neural Computer	117
Task Transfer and Education	118
Programs in the Mind: Problem Solving and Reasoning	119
Measuring Intelligence and its Implications	121
Multiple intelligences	122
Control	124
Summary of Key Terms	124
Chapter 8: Language	126
Chapter 9: Origins: Evolution, Genetics, and Development	127
Evolution	127
Genetics	129
Sexual Reproduction	131
Heritability and Individual Differences	132

Shared Environment and Parental Influences	136
Development	137
Piaget and the Development of the Neural CPU	137
The Development of the Prefrontal Cortex	139
Crystallized vs. Fluid Intelligence Development	140
Social, Personality and Moral Development	141
Lifespan development	143
Summary of Key Terms	144
Chapter 10: Personality	145
Acknowledgments	146
Glossary	147
About the Authors	148
References	149
...	

Published by Open Textbook, freely available

Copyright © 2018 Randall C. O'Reilly

All rights reserved. No part of this publication may be reproduced, distributed, or transmitted in any form or by any means, including photocopying, recording, or other electronic or mechanical methods, without the prior written permission of the publisher, except in the case of brief quotations embodied in critical reviews and certain other noncommercial uses permitted by copyright law. For permission requests, write to the author, addressed “Attention: Book Permissions,” at the address available from below.

<https://github.com/PsychNeuro/ed1>

To my family.

Preface

This is an in-progress experiment – feedback is more than welcome.

A number of the current figures have been shamelessly downloaded from the web, and sincere apologies are extended to those who would otherwise not wanted to have shared them in this way. Obviously if / when this book is officially published, this situation will be fixed.

Chapter 0: Introduction

Introductory Psychology textbooks typically provide a rather fragmented, fact-laden view of the field, relying on colorful graphics, exciting news stories, and personal anecdotes to generate interest in the material. This book represents a radical departure from that approach. Instead, the goal here is to provide a simple, succinct, principled account of the human mental world and how it emerges from our brains, that is coherent across the scope of phenomena covered in typical Intro Psych texts. There is a linear narrative story, intended to be read like a regular book, instead of jumping around looking at the figures and key words.

The overall portrait painted of you looks something like this (caution: it is not overly flattering, and you might even think this song is not about you, but go ahead and be vain – it is): You are obsessed with controlling your environment to satisfy a range of core desires and to mitigate strong fears. You are unlikely to be swayed by other people’s advice, but have no problem dishing it out. A challenge to your social standing or any other form of disrespect (the *diss*) is one of the worst offenses. You are willing to spin all manner of stories to maintain your sense of order in the world, *especially* when that sense is strongly challenged, often to the point of absurdity in the eyes of others.

You crave simple ways of understanding the world, to the point of massively oversimplifying the true complexities and ambiguities, preferring to think in terms of concrete anecdotes instead of broad abstractions, logical arguments, or, especially, statistics. You think you know how most stuff you use everyday works (bikes, cars, toilets..), but studies show that you are actually remarkably clueless (Keil 1981; Sloman and Fernbach 2018)– how exactly does that chain on a bike work? Perhaps most glaringly, you can’t help but think in terms of stereotypes, and inevitably focus on information that is consistent with your existing views, while ignoring all those nagging hints that all may not be as simple as you might like.

You only care about things that are new and unexpected, and are constantly comparing and evaluating yourself and others with a keen eye for who is doing better or worse along any number of important dimensions (wealth, beauty, smarts, athletic ability, popularity – you name it!) You are hypersensitive to who might be cheating or gaming the system, but are perhaps not so aware of unfair advantages you might have. More generally, you tend to think of yourself as being “your own person” and strongly underestimate how strong of an influence other people actually have over you. If you’re honest with yourself, you’ll admit that you spend way too much time thinking about what other people think of you – without recognizing that everyone else is doing the same thing, so that in fact the answer is a somewhat disappointing: “not much” (unless of course you do something embarrassing or strange or stupid, but even then, your memory of those events will typically far outlast those of others).

In other words, you are a *survivor*. You are a tough cookie. Your ancestors survived unbelievable hardships to get you here, to your relatively plush college-educated world. You are amazingly efficient. All those crazy details you don’t know about the world are largely irrelevant anyway. Seriously, does it really matter that you don’t know how the engine or transmission in your car works? You can drive, and get to where you need to go – and that is what really matters. Your brain is exquisitely tuned into what really matters, and despite over 60 years of attempts to recreate the magic of your brain in a computer, nothing has come even close (despite all the recent media hype to the contrary).

And yet, despite all your toughness and amazing abilities, you are very likely to have at least some level of significant mental dysfunction. You are more likely than not to suffer from depression, anxiety disorders (and often both of those together), drug dependence, etc. Unfortunately, the promise of a magic pill to cure these afflictions has turned out to be yet another disappointment. In fact, regular old “talking to another human being about your problems” (i.e., therapy, which is actually somewhat more involved and structured than that) is likely to be more effective than medication for most people.

The Three C’s

Surprisingly, we can make sense of all the above (and more!) using only three core principles:

Compression

Each neuron in the most important part of your brain (the *neocortex*) is wired for simplification, and the collective effect of the massive waves of electrical activity surging through your brain every millisecond is

to compress, reduce, and simplify information. Each neuron receives input signals from roughly 10,000 or so other neurons, but guess how much it can then say about that flood of information coming in? Almost nothing. First of all, it only has *one* output signal, the *spike*, which is an all-or-nothing affair. Furthermore, a typical neocortical pyramidal neuron will fire at most around 100 spikes in a second. And a second is a relatively long time in the inner loops of the brain – there is evidence that 1/10th of a second represents a kind of fundamental time-frame for information processing, so those 100 spikes reduce down to just 10 spikes within that critical window. And most neurons are firing far less rapidly than that. It's like when you tell your friend all your deepest thoughts, and they just say "huh". Neurons are the strong, silent type most of the time. But still waters run deep: when neurons *do* get excited about something, it is likely to be *important*, and most of what they are doing is *shielding you from constant TMI* (too-much-information – but you knew that already, so, kind of a meta thing we got going there...)

The raw scene coming into your eyeballs is truly gory: all jumbles of light, motion and color. When you were a tiny baby, you were overwhelmed by this “blooming buzzing confusion”, but now your neural networks have learned and developed to the point where you don't (can't!) even see that raw sensation anymore (unless you partake of various hallucinogenic substances, but even then, the level of disorder experienced is trifling compared to the pure chaos of the raw, unfiltered tidal wave of sensation coming in). We get small, fascinating hints of the magic power of our perceptual systems through illusions, and the occasional “viral gold / blue / brown dress” controversy, where people see or hear strikingly different things from the same stimulus. But overall, we really have absolutely no idea how much undercover cleanup work is going on inside our brains. If anyone was truly aware of the level of conspiracy operating in there, it would be scandalous. But, somehow, amazingly, we largely all end up converging on the same stable, boring illusions of simplicity. A table. A chair. Some french fries. People walking down the street. Cars driving by. Nothing strange going on here.

We would be utterly nonfunctional without this compression. For the same reason that those hallucinogenic drugs render people nonfunctional. If you want to do something useful with your time, you need to be able to make everything else in the world boring and irrelevant, so you can focus on *what matters*. If you're reading a book, or your tiny screen, it simply wouldn't work if every time you moved your eyes, the whole world was seen afresh, requiring you to reorient and rediscover what you were just reading and what you need to read next. Interestingly, this capacity for perceiving a stable, boring world seems to depend critically on a very active underlying process of *prediction* – your brain is stitching everything together in a seamless whole by filling in the gaps with what you *expect* or *predict* to see. You can easily see this, and relive some of your earliest experiences, by simply closing one eye, and then gently pushing on the bottom of the eyelid of your other, open eye. Suddenly, the world is a moving jumbly mess again! (Seriously, try it!)

Your brain's penchant for simplification (compression) does not stop with perception. Your highest levels of thought are similarly dominated by the same quest to render everything simple and predictable. Instead of recognizing the incredible high-dimensional diversity of our fellow beings, we inevitably reduce everyone to stereotypes. Even members of negatively stereotyped groups are caught in the evil maw of this process, exhibiting similar levels of stereotype-driven biases as everyone else. The ultimate expression of this compression process is the *anosognosia of everyday life* (aka the *Dunning-Kruger effect*; [New York Times Article](#)) – the lack of knowledge about our utter lack of knowledge. People can be remarkably unaware about what they don't know, and sometimes, this leads to funny situations. But, amazingly, most of the time, *it causes no obvious problems whatsoever*. We just keep getting on with our lives. And, as with perception, if we didn't, we'd never get anything done, because there is such a huge amount of stuff we routinely, safely ignore, that it would take many many lifetimes to process and understand it all.

Contrast

The next principle explains why we seem so fixated on comparing ourselves with others. Not just any others, but those certain people *who really get to you*. In that inexplicable, frustrating way. Why do I always have to be so jealous of those people? Can't I convince myself that the “grass is always greener?” Nope. As with compression, your brain is wired at the lowest level for magnifying contrasts, in this case via a special class of neurons called *inhibitory interneurons*, coupled with other important properties of all neurons that we'll cover in Chapter 2. The net effect is that your brain only sees things *relatively* (yep, we can have our own, special, relativity law in Psychology too – actually it is pretty general). A classic example of this is when you

come in from the bright sunny outdoors into a dimly-lit room. The difference in raw light energy coming into your eyeballs in these two situations is enormous, but, after a brief period of adaptation, you’re seeing things in the dim room that differ by a few photons here or there, whereas outside those few photons would be a minuscule drop in the bucket. In other words, our neurons *normalize* away the raw strength of whatever signal is coming into them, and remain sensitive to the *relative differences* compared to that overall signal. Those inhibitory neurons play the critical role of mathematically *dividing away* the raw signal strength, leaving the principal pyramidal neurons “in the zone” for responding to relative differences.

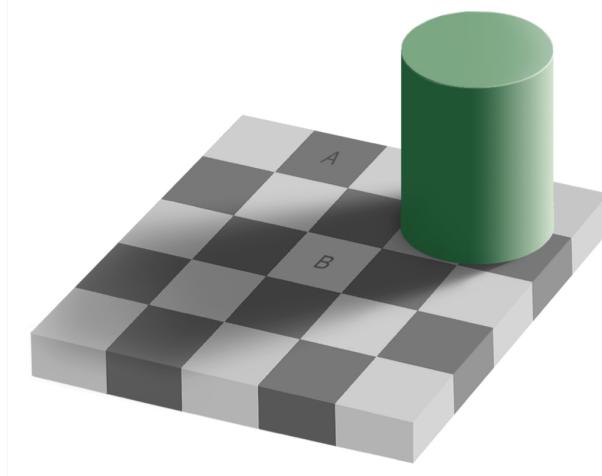


Figure 0.1: Illustration of the power of contrast in perception. Do you think the physical image-level color of square A is the same or different from that of B? Unbelievably, they are identical!

As with compression, perception provides some of the clearest windows into this phenomenon, for example Figure 0.1, showing the remarkable effects of contrast (and global scene understanding) on perception of color and brightness. Another remarkable example is the case of *perfect pitch* – why is it so unusual a skill for people to simply be able to recognize the absolute frequency of a sound? Mechanically, and mathematically, extracting such a frequency from a sound signal is trivial, and simple (a “Snark”-like guitar tuning device can be had for a few bucks). That this feat is so incredibly rare and difficult in humans just points to the pervasive power of the contrast relativity (most people can easily tell the relative pitch).

But contrast, like compression, is not restricted to the perceptual domain. It affects every level of thought, contributing to that insidious obsession with your relative standing among your peers. For example, studies routinely show that the absolute amount of money that people make is largely unrelated to various measures of their happiness – instead, what matters is their perceived level of income *relative to their peers*.

Contrast operates over time as well, in several important ways. First, at the perceptual level, we are highly sensitive to the rate of change over time of stimuli. The classic example here is the slow approach to boil being unnoticed by a hapless frog until it is too late (this is not exactly true – you can’t get all the way to boiling, but it is very likely true to at least some extent). Similarly, a cottage industry of amazing demonstrations of our inability to detect slow changes in visual scenes sprung up a few years ago: [YouTube video by Dan Simons](#). Once you become aware (upon repeated viewing or instruction) of the nature of these changes, it is truly astounding to realize how much you overlooked them the first time(s). If you rapidly flip between the start and end frames of these slow-moving videos, the changes pop immediately into view. Again, we see the *delta* (relative change), not the absolute value of things.

Nowhere is this more poignant, and pressing, than deep inside the *dopamine* system in the middle of your brain. As you already know (and would be annoyed to have me repeat, but I’m going to do it anyway, to prove the pertinent point), dopamine is widely believed to be the “pleasure drug” in your brain. It is associated with drugs of addiction, and actually most other major mental disorders in one way or another. However, this popular description of dopamine leaves out one of the most important points: dopamine is *not* about *raw* pleasure, but rather, about the *difference* between what you experience and what you *expected* to experience. Specifically, if you get exactly what you expected, your dopamine goes “meh”. This soul-crushing

response to your greatest accomplishments is exactly what critics do to performers, and indeed the dopamine system is best understood as being the *central critic* of your brain. Far from a center of epicurian delights, it is a hard-nosed bully that is never satisfied. And that dissatisfaction is what has driven us ever upward in all manner of exploits – many of them good, but many of them not so good.

Greed is really a byproduct of your dopamine system. Seriously, why in the world can't someone who already has *millions of dollars* just be happy with that, and give the rest away or do something else useful with it (and their all-too-brief lives). Because dopamine adapts quickly to that million-dollar feeling, and it keeps giving you back that critical "meh" response. You need more than that. You *deserve* more than that... It really is tough living with such an asshole critic in your head all the time. But then again, we really do owe every step of "progress", within our individual lives and as a society and a species, to that nasty little critic driving us ever upward and onward.

Finally, another obvious manifestation of the contrast principle is our collective obsession with the *news*. Especially with the advent of the 24hr news cycle and the constant updating of news information via electronic, online media, we are now living in a quickly-moving bubble of news that sweeps things up in its path and spits them out quicker than... yesterday's news. Or yesterhour's news. If you don't check in quickly enough (for fear-of-missing-out; FOMO), you'll likely miss huge swaths of news. In the "good old days", people used to read *weekly* news magazines (e.g., *Time*, *Newsweek*) – can you *imagine* reading news that might be an entire *week* old!? Everyone worries about this kind of thing, but really it is just what our brains are wired to do. Every conscious moment of our lives is driven by a thirst for knowing what has changed, what is different – anything that remains constant will quickly drift out of your mind, like that delicious aroma of dinner that I can no longer access (pro tip: go outside for just a minute or two – you'll be amazed as what you've been ignoring!), or, thankfully, that feeling of my butt sitting on this chair that I was thoroughly *so over with* until I just wrote that sentence..

Control

Last but certainly not least, is our obsession with *control*. Some of you may be thinking that you're not a control freak like those *other people*, but actually, every one of us is a crazy control freak at some level – it just differs in terms of what matters to us. Anyone want to have some stranger come pick you up and take you around to work with them all day? Or just invade your personal space? How would you feel if someone just started selling all your stuff on craigslist? Or how about those people who go door-to-door (or stop you on the street) and try to convince you to believe in some particular brand of religion? Or just your roommate who keeps nagging you about the dishes, or being too loud, etc. Yeah, there's definitely *something* for *everyone*, where it matters. Don't even start about the parents, whose only mission in life seems to be exerting unwanted control over yours. And usually, if you have two or more people living together, you quickly become aware of all that stuff that you didn't realize really matters to you. A lot.

Starting again in the brain, virtually every neuron in the brain is serving the master of control at one level or another. At the most basic level, it is about *motor control*, and a great example of the dedication of the brain to this particular function comes from a lowly sea squirt that starts off life as a mobile tadpole, and flits around in the ocean for a bit, looking for a good place to settle down. As soon as it finds its special place on the reef, it promptly eats its own brain! Because, the whole point of the brain in the first place, evolutionarily speaking, is to process sensory inputs *in the service of producing useful motor outputs to improve survival and the overall quality of life*. There's a reason nobody thinks highly of layabouts and 30-year-old's living in their parent's basement: progress requires action, and our brains are wired for action. In the brains of most species, there are big chunks devoted to the compression and contrast processing of sensory inputs, and the rest is devoted to using that information to figure out what kinds of opportunities and threats are out there in the world, and how to best optimize chances of survival within the repertoire of available motor actions. Not much space left over for cultivating arcane knowledge about civil war battles, or fantasy role-playing games, or whatever other weird, seemingly non-functional things people spend their time doing.

The human brain takes this obsession with motor control to the next level, by building an internal fortress / castle of the *self*. We're not quite sure to what extent any other beast even has a similar kind of thing inside their own mental worlds. The self is a model, a construct, built up over years, that helps us

predict how we are going to behave, and what we seem to really want (and not want). By having such a thing inside our own brains, we can use it to more accurately anticipate what kinds of motor actions are really going to get us what we want. This is especially important when dealing with other people, who are, compared to your average rock or tree, very complicated and unpredictable. I'm not saying you're a manipulative little jerk. I'm saying *everyone* is a manipulative little jerk, deep down. It is, again, just a logical extension of what brains are supposed to be doing. If they aren't good for maximizing pleasure and minimizing harm, then we might as well all just eat them for dinner!

This *self model* lying at the heart of our control system is like our secret nuclear power reactor inside our brains. It is the “nerve center” of our being. It does *not* take existential threats kindly. Anything that appears to threaten our internal sense of identity and control gets raised to the red alert level. This is why you can't just “mansplain” something to someone else, and expect them to instantly see the error of their ways, and instantly become a new, better self. We have a lot of investment in that *old* self, and it does not look kindly on being deposed from its despotic rule over its own internal kingdom.

Developmentally, the self emerges around age two, heralded by the onset of *tantrums*. Tantrums are the inevitable consequence of an emerging desire for control, coupled with an almost complete lack of *actual* control. This is really the defining battle of life, and it never really ends: the best you can hope for is some kind of truce as expressed in the Serenity Prayer of Reinhold Niebuhr: “God, grant me the serenity to accept the things I cannot change, Courage to change the things I can, And wisdom to know the difference.” The enduring power of this saying is another testament to the central importance of control in our overall life happiness.

Although the self is a despot at heart, it is also remarkably sensitive to external, social forces, creating one of the most fundamental and puzzling paradoxes of the human condition: We care deeply about what other people think of us, and are actually remarkably malleable in adapting our behavior under the influence of others. There are many demonstrations of the power of the social force, from the evil of Nazi Germany and controversial attempts to recreate those forces in the lab, to the seemingly more benign and amusing phenomenon of hypnosis. Biologically and ecologically, our very survival is utterly dependent on our ability to work together socially, and social motivations are undoubtedly wired directly into the depths of our brains, providing pathways that can be “hijacked” to get past the watchful eye of the self-model.

And therein lies the likely explanation for this paradox: these social forces can only act when delivered in ways that the self either does not recognize as threatening, or even endorses – like the immune system, we are highly sensitive to foreign invaders. The minute you are aware someone is trying to convince you of something, is the minute that it fails. But when a social virus is neatly packaged in a nice sugar coating, often in terms of reinforcing a sense of belonging with an identified *in-group*, then it can easily slip past the guards. These kinds of in-group / out-group (tribalism) dynamics are the strongest of social forces and underlie all our greatest evils (genocide, war, hate in all its forms). But they are also the basis of *love* and social cohesion of groups at many levels (family, religious, civic, sports fans, etc).

The prevalence of cults and the seemingly obvious insanity and self-harm of those ensnared in their traps, is testament to the powerful forces at work here. Unfortunately, some people have learned to tap into these powerful forces and abuse them to satisfy their own needs, of control over self and others. This phenomenon is also evident in nationalist or populist political movements, which mine this same deep need for social belonging, to support the leader's power in ways that clearly go beyond any kind of rational socioeconomic-level considerations. Understanding the nature of these powerful forces is thus something that transcends scientific disciplines and may be critical for the continued survival of our species – psychology and neuroscience have never been more relevant!

In short, the social extension of control is *power*, and many social-level dynamics, as well as personality variables, can be understood in terms of power and control (Hopwood et al. 2013).

The Breakdown of Control

Unfortunately, achieving *serenity now!* is very difficult. And all those challenges to the self can end up leading to a bout of depression, often coupled with anxiety or other unpleasant mental states. Although widely characterized in terms of *anhedonia* or the inability to experience pleasure, current research supports the idea that the core disorder of depression is really about *control*, or the perceived lack thereof. When

your self model is sufficiently challenged, it basically gives up on a lot of goals, and unfortunately, achieving those goals is a primary source of pleasure and satisfaction in life. So, yes, anhedonia is a consequence of depression, but the core of it is more about the inability to motivate yourself to get out of bed and do all those now-meaningless things that you used to find meaningful.

Consistent with this central role for control, one of the most promising components of modern therapy for treating depression is *behavioral activation*, which is essentially an attempt to reboot your core self-motivation control system. Indeed a major study found behavioral activation to be the most important element among a group of therapies, and as effective as medication (Dimidjian et al. 2014). And when you recognize the central role for control in depression, it is then less surprising that medications are relatively ineffective: for the vast majority of people, the problem is *not* about some kind of low-level imbalance in their brain chemistry: it is about their core mental power plant of control running out of steam. And it takes hard mental work, aided by effective therapeutic treatments, to reboot your own sense of mental self-control and efficacy.

For the smaller proportion of people who clearly do have a biologically-based mental disorder, it is still the case that the brain areas most centrally involved in self-control are the ones that are most likely to be affected. Schizophrenia and OCD for example involve the frontal cortex, basal ganglia, and dopamine systems of the brain, which are the main players in developing and sustaining our internal self control system. Thus, understanding how different parts of the brain function to support this critical self-model system is a major goal of current research in Psychology and Neuroscience, and this book is designed to get you started on a journey toward understanding this cutting-edge work.

Other Principles and Perspectives

There are many candidates for “the fourth C”, and different names could have been chosen to refer to the above “three C’s” (e.g., reduction, relativity, and... respect?), but being a slave to the simplifying force of *Compression*, it is useful to try to see as much as possible through the lens of these three principles. Furthermore, as briefly introduced above, these principles can be tied directly into the most fundamental properties of the nervous system, and thus provide a critical *bridging function* between Neuroscience and Psychology. Nevertheless, it is important to always remain aware of all the compressing taking place, and to acknowledge that this radical attempt at synthesis may strike many practicing scientists as overly simplistic or downright wrong-headed. However, my hope is that the benefits outweigh the costs overall, without attempting to overly minimize those costs.

Where do we go from here?

This question can be asked at two levels: the short-term question of where this book is headed, and the longer-term question of where our species is headed!? Although it may seem like our current cultural and political environment reflects an extreme magnification of many of the negative aspects of human mental function as described above, another perspective is that these truly are perennial battles and challenges that we have struggled with since the dawn of human history, and that they are borne of fundamental properties of the human brain that also have many positive aspects. Like everything it seems, double-edged swords abound. And the core premise, and promise, of science is that by understanding something deeply, we are better positioned to make the best of it. This contrasts with the idea that by somehow reifying “bad” features of the human brain, we are therefore justifying the bad ends they produce. That is not the aim here.

With those big picture questions out of the way, we can turn to the plot for the rest of this adventure story through the human brain. Unlike a good mystery story, we’re going to ruin the whole thing right up front, in the hopes of achieving a better understanding and mental roadmap in the bargain.

Chapter 1 will provide a big-picture overview of the challenges and promise of achieving a scientific understanding of the brain and the mind (particularly the mind). The main challenge here is the *subjective* nature of the subject matter, which Psychology has wrestled with since its inception. Indeed, subjectivity is, according to Rene Descartes (and many others), the only thing we know for certain, and it stands as a fundamental barrier for any attempt to transcend its bounds and achieve an *objective* understanding of the world around us, and, especially, the world inside our subjective minds.

Chapter 2 will cover the nuts and bolts of the brain, but always connected directly to the bigger picture

via the three-C's principles and their applications. We'll see in detail how each neuron functions as such an amazing "information compactor", compressing those 1000's of signals into its single spiky output. We'll then take an amazing "connected" voyage through the pathways of the neocortex, seeing how the great chain of neurons locked in their long-lasting embraces create channels where information flows in different ways. We'll wrestle with the central question of whether brain areas are truly "specialized" for different functions or not, and whether there is any "there" there, as in, "where *is* that memory anyway?"

Chapter 3 will dip into the fascinating shores of consciousness, sleep, and arousal, into the land of dreams and altered states of awareness, providing additional grist for the subjective nature of experience, and surprising insights into the way our brains work. We have powerful *neuromodulatory* systems that are altered by psychoactive drugs, helping to reveal the ways in which these systems function when not-so-altered as well.

Chapter 4 is all about the first two C's: Compression and Contrast, as manifest in sensation and perception, where we try to understand illusions like that shown earlier in Figure 0.1, and how it is that our brains render this most challenging of computational problems almost entirely *transparent* – we are almost completely unaware of the incredible feats of inference and insight that our brains are performing, effortlessly and efficiently.

Chapter 5 tackles one of the most important unsolved questions in the field: how does our brain *learn* everything that we as adults so blithely take for granted? Yeah, all that school was kind of hard maybe, but how exactly did you really learn everything? Especially in those critical times before school-age learning, when you mastered perception, language, and got a good start on motor control. Most intro textbooks focus almost exclusively on decades-old conditioning level learning – we'll cover that, but also give you a much more recent, cutting-edge sense of how new developments in AI and neuroscience may provide some insights into these fundamental puzzles.

Chapter 6 picks up where learning leaves off: the fascinating, misty land of memory. Where are all those childhood memories actually stored in your brain? What is the difference between "motor memory" for how to ride a bike, versus knowing the capital of India, and other types of memories? How can you improve your memory to get better grades in school, and also achieve a deeper understanding of things?

Chapter 7 examines the upper echelons of your brain, where "thinking", reasoning, planning, etc take place. What really is the difference between someone who is "smart" vs. not so much – are so-called "intelligence" tests for real, or just a biased cultural artifact?

Chapter 8 just scratches the surface of the deep topic of language, and its critical role in shaping human thought. Language is fundamentally a social activity, and we can learn a lot about ourselves through this lens.

Chapter 9 goes back to the beginning(s) and considers the evolution of the human brain and human behavior, how genetics shapes us relative to the role of experience (nature vs. nurture), and what we know about how we develop across the lifespan.

Chapter 10 expands our horizons into the domains of social interactions and personality, focusing on these critical questions of how social forces shape our thoughts and behavior, and how people differ in the ways that they interact with others.

Chapter 11 wraps up with the all-too-prevalent breakdowns in our amazing nervous systems in the form of the major mental disorders and how they are treated through therapy and pharmacological interventions. The most successful forms of modern therapy are based on basic, practical principles, and include a considerable focus on understanding and acceptance of how your brain functions. In this way, this textbook might be considered a form of therapy: hopefully you will come away from this with a much deeper and more satisfying understanding of what makes you tick, and what ticks you off, so that you might go on to lead a happier and fulfilling life!

Typography tip (as a reward for reading all the way through!): **bold** is for keywords that are likely to show up on tests, and *italics* is for important but less-likely-to-be-tested terms.

Chapter 1: Science and Subjectivity – The Fundamental Challenge of Psychology

Psychology is the science that attempts to understand the human mind. The human mind is the most fascinating and amazing “thing” in the known universe, and the idea that you can actually attempt to study it using the basic reductionistic approach of science may seem a bit of a stretch. And indeed it has been – but at this point in the development of the field, many practicing scientists are likely to feel at least somewhat confident that significant progress has been made, without fundamental, obvious limitations to how far we can go.

Despite all this progress and optimism, we will see in this chapter that there are fundamental boundaries to what science can penetrate, and these boundaries have shaped the field from its inception. Thus, understanding these limitations helps put the field of psychology and neuroscience into perspective in multiple ways, and in fact many of the limitations we discuss apply to science, and all human knowledge, more broadly.

The central issue we must confront head-on is the inescapable problem of *subjectivity*. By subjectivity we mean not just the fact that different people have different opinions or perspectives on things, though that is a big part of it. Instead, we need to step back a bit to look at the *really big picture* (i.e., Philosophy), starting with the fundamental problem of subjectivity as expressed by **Rene Descartes** (way back in 1637), in his famous statement: *Cogito Ergo Sum – I think therefore I am*.

There are two essential implications of this statement – we’ll explore the first one in depth before turning to the second. The first implication is that *subjective experience is primary*. If you put yourself into the mindset of a very skeptical, doubting philosopher, you might find yourself questioning just about everything, *except* this one, primary fact: you are *here* (wherever you are), *thinking*. If you really push it, you might appreciate that you can’t really be sure that the world itself exists outside of your mind! This very challenging train of thought is well-captured in several modern movies, perhaps most notably in the *Matrix* series, where, in fact (in the movies at least), there turns out to be every reason to have such doubts. In philosophical circles, this line of thinking is known as *solipsism*, and lest you think that this is just an irrelevant and obscure way of thinking, one of the great innovators of our time, Elon Musk, is apparently convinced that we’re all living in a giant simulation. He also apparently smokes entirely too much dope, but be that as it may.

This is the kind of all-encompassing subjectivity that we want to more fully understand and appreciate. What does this line of thinking mean for the study of psychology, or science more generally?

This is where we can usefully bring in Descartes’ second major implication from *Cogito Ergo Sum: dualism*. Dualism is the idea that there are two fundamentally different “substances” in the universe: the regular physical stuff of the everyday world, and this entirely separate, magical transcendent thing called *mind*, which lives apart from that other, regular stuff. The opposing view is called *materialism*, where the mind is seen as just a product of the material world like everything else, and in particular a product of the physical processes taking place within the *brain*, as widely embraced in modern neuroscientific approaches to psychology.

You might be somewhat surprised to hear that many modern-day philosophers still embrace dualism, and one of the most outspoken advocates is David Chalmers, who argues that understanding the nature of subjective experience, or *qualia*, is the *hard problem* of consciousness and simply cannot be explained in objective, materialistic, scientific terms (Chalmers 1995).

You might also be surprised to hear that, despite being one of those modern materialistic neuroscientists, I actually agree with Chalmers, and Descartes (in spirit at least, so to speak)! I think that there are two fundamentally different “somethings” in the universe, but, unlike Descartes and Chalmers, I don’t think the dividing line is between *mind* and *matter*, but rather, between *subjective* and *objective* perspectives (Nagel 1974).

Following Descartes (again), we can take subjective experience as primary – it is the only thing I am fully certain of. But it is also primary in another, essential way: it is uniquely, completely, definitionally, *mine*. It is literally impossible for *you* to experience *my* subjective experience, because, by definition, *my* subjective experience is exactly the sum-total of what it “feels like” to be me. If we somehow were to add *you* into my brain, my subjective experience would be irreparably altered. If you are somehow sharing in my subjective experience as it is happening, you would have to have direct access to every level of my brain, and

not just “objective” access as you might get from a super-hi-tech future brain scanner, but *direct, internal, subjective* access, “from the inside out”.

In other words, you would have to literally be inside my brain. And you can’t be inside my brain because I’m already here. From the materialist perspective, we can identify my subjective experience as emerging directly from my brain – it is what it feels like to be my brain. If you truly appreciate this equivalence, then it should be readily apparent that there can be only one “mind” for every brain (we’ll look into the fascinating phenomenon of multiple personality disorder later, but it doesn’t change this fundamental conclusion – all those personalities are just as irrevocably trapped inside the one brain as you and I are, and in fact we all have something like multiple personalities too).

Another way of thinking about this is in terms of identical twins. Let’s imagine we have the most identical of identical twins ever to exist. Their brains are *completely identical* in every way possible. Would those twins have the same subjective experience? No. They might have a great deal in common, but, fundamentally, they would not, and could not, directly experience exactly what the other is experiencing. Why not?

It all boils down to *perspective*. Each physical thing in the universe has its own unique perspective, if we take this term to mean a particular spatial location, and a particular trajectory through space and time in the past (and going onward into the future), that is fundamentally *unique* to that thing. This is why the twins cannot share their subjective experiences: they are two separate, distinct things, and, inevitably, they “see the world” from two different vantage points. The only way they could share experiences is if they could somehow superimpose themselves into exactly the same point in space, and do so over a sufficiently long time period to synchronize their history of experience, which plays such a critical role in our subjective life, in addition to the immediate sensations coming in from the outside world.

Anyway, the key point of all this is that *if* you allow that subjective experience can never be shared among different brains, *then* it follows that there is a fundamental divide between this inner subjective world, and the “regular” outside *objective* world. I believe this divide captures the essence of what Chalmers is talking about in terms of the irreducible nature of the qualia of consciousness – the impossibility of trying to explain in objective terms “what it feels like” to experience things in our subjective, inner world. Furthermore, it does so without introducing anything particularly magical or fundamentally at odds with materialism: subjective experience is not separate from the physical world in terms of some kind of magical “substance” that it is constituted from – it is just separate in terms of this notion of *perspective* – the unique point of view (literally, where they are standing / sitting / looking) that each subjective being has all to themselves.

Subjectivity in Psychology: A Brief History

Stepping back from this big philosophical abyss, what does it all mean for the attempt to study psychology as a science? The primary, obvious problem is that psychology is the study of *what it is like to be a human being*, and if this is fundamentally a subjective thing that can never be directly shared with any other human being, how can we possibly hope to arrive at some kind of objective, scientific understanding? Well, the first step is to follow Chalmers and attempt to *partition the problems* – we can carefully attempt to set aside the *hard problems* associated with the nature of subjective experience, and focus instead on the so-called *easy problems* that are left over. *If* there is enough interesting stuff left over in this space of easy problems, then it probably makes pragmatic sense to just see how far we can get in trying to understand that stuff, and then, once we seem to have exhausted that space, perhaps we could circle back and start reconsidering some of those hard problems.

This overall approach provides a reasonable narrative for the history of psychology as a scientific discipline. The person most widely credited with founding the science of psychology, **Wilhelm Wundt**, had the innovative idea in the late 1800’s that, after millenia of armchair speculation, you could actually apply the techniques of empirical science to understanding the human mind / brain. Wundt made many groundbreaking contributions, but his legacy, at least at the level of introductory psychology texts, is as a founder of the *introspectionist* school of psychology, which also includes **William James**, who also made major lasting contributions to the field. When the next major paradigm shift took place in the early 1900’s, it emerged as a strong reaction and rejection of this introspectionist approach, which was characterized as being overly concerned with all those hard problems of subjective experience. Introspectionists would try to systematize and characterize the contents of subjective experience, and the hard-nosed *behaviorists*

who came next regarded these investigations as insufficiently objective, rigorous, and replicable. Instead, they emphasized purely objective, externally-observable *behavior* as the only valid data in psychology (hence the term behaviorism). The main figures in this era (e.g., **John B. Watson, B. F. Skinner, and Ivan Pavlov**) focused on how external, objective factors such as reward and punishment affected subsequent behavior through *conditioning*.

Thus, these first two epochs of scientific psychology embody exactly this tension between the subjective and objective worlds. The next paradigm shift took place in the 1950's and 60's with the *Cognitive Revolution*, riding the wave of digital computers, which made it fashionable to start talking about internal mental operations in terms of the *information processing model* of the mind – i.e., the mind as a computational device. Scientists leading this new field, such as **Herbert Simon** and **Alan Newell**, started thinking about how the mind could perform complex mental operations such as scientific proofs, chess, and other challenging tasks (Newell and Simon 1972). People created running computer models of how these internal thought processes might work, which provided a compelling way to render that formerly “loosey-goosey” internal world in a much more rigorous, objectively-characterizable way.

However, as parallel work in the field of Neuroscience continued to advance, it gradually became clear that the brain really doesn't work anything like a standard digital computer. Instead, it is really a *massively parallel* computer with billions of computing elements (neurons) that combine the functions of computation and memory, which are otherwise separated in a standard digital computer. Psychologists **David Rumelhart** and **James McClelland** published a ground-breaking pair of books in the mid 1980's that popularized this new understanding of how information processing might work in the brain (Rumelhart and McClelland 1986; McClelland and Rumelhart 1986), and subsequent advances in the ability to take high-resolution pictures of the activity inside the human brain (*neuroimaging*) have led to the currently-dominant paradigm that integrates neuroscience and cognitive psychology (i.e., *cognitive neuroscience*) to come up with coherent understanding of how exactly the brain gives rise to the phenomena of the mind.

Fundamentals of Cognitive Neuroscience

This book is grounded squarely in this new paradigm of cognitive neuroscience, and attempts to provide a coherent set of core principles that connect directly from the basic processing carried out by individual neurons, all the way up to the highest levels of mental life. We are still largely avoiding significant consideration of the vast inner world of subjective life, but there is a robust field studying the *neural correlates of consciousness* (NCC) that we will discuss in depth in Chapter 3. Slowly but surely, we are building bridges between the objectively-identifiable properties of the human brain, and the subjective experiences that tend to co-occur with particular such brain states. Thus, we are developing a richer objective understanding about the kinds of neural mechanisms that give rise to our subjective mental life. But even with all of these advances, I don't think we could ever explain to a non-human-brain lifeform what it feels like subjectively to be a human brain. Thus, the subjective world remains our own private dominion, and literature, art, and movies provide the richest vehicles for sharing those experiences across the inevitable subjective gap between us all.

Subjectivity and Science: Working with the Method

The challenges imposed by the primacy of subjectivity have far-reaching implications beyond the field of psychology. First, given that some people can't even agree that there *is* an objective, external world outside the mind, how can we possibly even begin to start talking about *objective knowledge* and *facts*? This appreciation for the primary nature of subjective experience forces us to recognize that objective knowledge itself is entirely dependent on the subjective motivation of individuals to entertain a strong enough belief in this notion of objective reality, to put up with all the effort it takes to make any progress in understanding and advancing objective knowledge.

Those individuals are called “scientists”, and they follow a particular method, the **scientific method**, which has the following basic steps:

1. Come up with a general question or problem, e.g., based on an informal **observation** about something of interest (e.g., Newton observes the apple falling on his head, which gets him thinking..)
2. Form a specific **hypothesis** about how that something might work, which makes testable **predictions**

(e.g., there is an invisible force called *gravity* that causes all objects to experience the same acceleration, making the testable prediction that a feather and a hammer should fall at the same rate *in a perfect vacuum* so as to eliminate the “confound” of friction).

3. **Collect data** that could actually test the predictions of the hypothesis, in comparison to other possible hypotheses (e.g., measure how fast things fall, ideally in a vacuum if you happen to have one of those lying around). It is essential that the data be collected using a well-specified procedure that could be **replicated** by other scientists.
4. **Analyze the data** to determine whether any effects observed are strong enough to be clearly distinguishable from random chance and noise.
5. **Draw conclusions** – how compelling are the data, what holes are there in the data that would allow other hypotheses to explain the observed effects, etc?
6. Iterate! Plug the holes, think of other alternative explanations, test those, etc.

These steps can incrementally pull us out of our individual subjective fortresses through the critical lever of **consistency**. If you articulate a clear sequence of steps to perform an experiment, and tell me exactly what you observe as results, and I do the same thing to the best of my ability, and get *consistent* results, then it seems like there might be something *real* and *objective* going on, or at least the world isn’t completely random. As more and more people do the same thing, and continue to get consistent results, the odds that each one of us is just being individually tricked by some kind of subjective illusion would seem to go down.

As this scientific process continues, ever broader networks of interconnected hypotheses and associated empirical data accumulate, and if all of these remain somehow consistent with each other, it really starts to seem like there might be some kind of *laws* governing the behavior of the outside world. Furthermore, all this scientific knowledge makes its way into technology, which depends on those same laws, further bolstering the network of consistency. Fast forward to the modern world, and we now have the *standard model* of physics that provides a single consistent framework for understanding virtually all physical phenomena that have been subject to experiment, and drives incredible technology that would have been considered pure magic in times past.

Despite all this amazing progress made through the iterative application of the scientific method, you still have people like Elon Musk, one of the great *users* of physical laws, nevertheless concluding that it is all a giant simulation. And still plenty of people who believe that the Earth is flat, etc. And there is *nothing* you can do to convince these people otherwise. Such is the ultimate primacy of our subjective perspective on the world: the *only* porthole we have onto that supposed objective reality out there is through our very own, individual, subjective lenses. Because our subjective worlds are fundamentally uniquely our own, this also means that nobody can force anyone to believe anything that they aren’t otherwise prepared to believe. Objective reality really is a second-class citizen, and is entirely dependent on the patronage of the ruling, sovereign subjectivity, just as scientists are still to this day dependent on the hard work and wealth of others to have the luxury of time and resources to create this huge network of consistent hypotheses and data.

Even within the scope of the scientific method, subjectivity abounds. Where, exactly are these hypotheses, or conclusions, supposed to come from? How many scientists looking at the exact same empirical data draw the same conclusions? You’d be surprised how subjective and inconsistent cutting-edge science really is. History is full of examples where a visionary pioneer was ridiculed by their colleagues, until enough evidence accumulated, and enough old people in power died, to allow the new ideas to flourish (“science advances one funeral at a time”, according to Max Planck). The widely-accepted description of how science actually works, developed by Thomas Kuhn (Kuhn 1962), emphasizes this sociological, psychological reality of science, with one major consequence being the strong suppression of ideas that are inconsistent with the current paradigm.

We can understand this phenomenon in terms of the three C’s principles. Compression says that people crave simplicity, and the current paradigm embodies that: it is something that a large number of people know and agree about. Having that overturned requires confronting a high level of uncertainty and complexity. Control is paramount here: that challenge to a widely-believed paradigm is experienced as a direct, personal challenge to your entire mental fortress – psychologically, it is really the same as challenging someone’s belief in a particular religion. Furthermore, the uncertainty directly undermines the feeling of control as well. And control interacts with contrast – the “paradigm believers” constitute a social in-group, and anyone challenging the paradigm is immediately a strongly-contrasting out-group member, and all the deep tribal motivations

are aroused in this case, causing the challenger to be treated like a real outcast and pariah.

In other words, science is just people being people. However, despite all our limitations and inevitable subjectivity, there is some indication that following some approximation of the scientific method really does seem to work, at least over the longer arc of history.

Before we get more into the nuts and bolts of actual experiments and statistical analysis techniques in psychology and neuroscience, there is one further perspective on the problem of subjectivity in science that bears mentioning. This comes from Robert Pirsig, who wrote the famous book, *Zen and the Art of Motorcycle Maintenance*, which is actually more about philosophy of science and personal autobiography, rather than Zen per se. Pirsig literally went insane (as in, institutionalized, electroconvulsive shock therapy, etc) in the course of struggling with the question of where hypotheses come from – he realized that there was no rational explanation for how to come up with a good hypothesis, and it seems like there could easily be an infinite number of plausible hypotheses, so this throws a massive monkey wrench into the entire rational foundation of science.

Thus, subjectivity, creativity, and individual genius truly lie at the heart of science – most scientists are reasonably capable of evaluating hypotheses in terms of their consistency with data and with the larger network of other validated hypotheses, but relatively few scientists are responsible for coming up with the major hypotheses in the first place. Oh, and by the way, Pirsig suffered from Schizophrenia so that probably had more to do with his mental breakdown than the problem with hypotheses, but anyway it makes for a good story.

Research Methods in Psychology and Neuroscience

The one thing that everyone in science can agree upon is that *data* is essential! You can come up with cool-sounding new hypotheses all day long, but ultimately the data has the final say (despite the grumbling of many a theoretician). Nobel prizes strongly favor those who generate data, and even Einstein never got one for his general relativity theory! In this section, we'll discuss the specific types of data that psychologists and neuroscientists tend to collect, and what kinds of analyses are typically done with that data. We will cover this succinctly because it all *sounds* perfectly logical, but actually applying it requires a good deal of practice and experience, which is beyond the scope of this book, and likely the course you're currently taking.

In psychology, there are three major ways in which data is collected, each with complementary trade-offs:

- **Descriptive Methods** – these tend to be the least *invasive* techniques, involving various ways of capturing what is actually happening in human behavior, such as observation, case studies, and surveys. A modern version employs cell phones with apps that ping people at random times during the day and ask them what they're doing, or thinking about, etc. The disadvantage of these techniques is in their relative inability to inform you about *why* people might be behaving the way they are – the other two techniques improve on that aspect of things, but, particularly with the experimental method, tend to require more artificial, less naturalistic kinds of experiments.
- **Correlational Studies** involve measuring multiple different **variables** (something that can be measured which varies across people, such as weight, IQ, vocabulary, diet, etc) and determining the extent to which these variables **correlate** or vary systematically in relationship to each other. For example, people's weight and height tend to be positively correlated, because as one goes up, the other does too. Critically, as with most real-world data, this is not a **perfect** correlation – there are many exceptions in either direction – but overall, on average, there is a relationship. The single most important limitation of correlational studies, is that the presence of a **correlation does not imply causation**. Often, however, it does work the other way around: causation *does* typically produce correlations of some sort. Interestingly, psychologists have shown that we're particularly bad at overgeneralizing these kinds of relationships (Wason 1968), and thus we tend to falsely conclude that correlation does imply causation. In short, the human brain relies on correlation as a kind of “quick and dirty” shortcut for finding causal relationships in the world, and we find it remarkably difficult to recognize that the two are not equivalent. For example, most studies on the effects of diet on health are correlational, and yet the media and even scientific papers regularly interpret these as showing a causal link. “Drink more coffee because you'll live longer!” Well, what if in fact the observed correlation between coffee and longevity is due to the fact that more wealthy people drink more coffee, and it is really the wealth and all its

associated benefits that is driving the longevity. Coffee is just “along for the ride”. This is the **third variable problem** (in this case, the third variable is wealth), and it is the bane of correlational studies, because *there is always a third variable* (and a fourth, and a fifth, etc). And it is typically very difficult to rule out the possibility that everything is being caused by one of these unmeasured “third variables”.

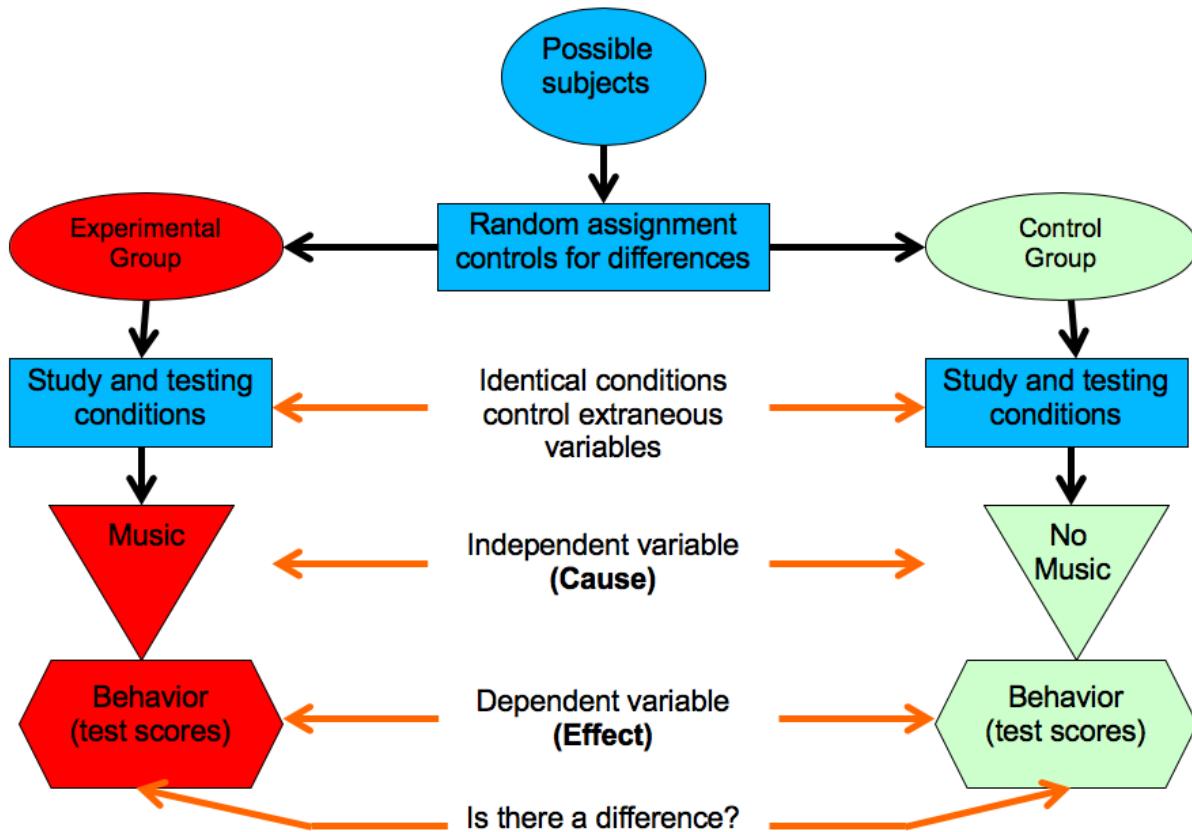


Figure 1.1: Logic of an experimental study, using random assignment to eliminate third variables from the study participants. It is also essential to minimize all other differences between the experimental and control conditions (i.e., *confounds*, or additional “third variables”), to more precisely identify the single *independent variable* (i.e., the *causal* variable) as truly being responsible for the differences measured in the *dependent variable*

- **Experimental Studies** are the only way to truly establish a causal relationship, and even then it is still a major challenge to really accomplish this feat. The key trick is to use *randomness* and careful designs to attempt to systematically eliminate all possible “third variables”. A huge source of third variables is each individual person participating in the study. Like all the bacteria on your skin, you are crawling with third variables. Your genes, your upbringing, your neighborhood, your schools, your friends, your... everything, is a teaming cesspool of third variables! The key trick in an experimental study is to use the cleansing power of randomness to wash away all those third variables, by **randomly assigning people to different conditions**. No third variable can withstand the incredible power of such random assignment – if we find a systematic difference between two completely random samples of the population, it cannot be due to their pre-existing conditions!

However, random assignment is also the achilles heel of experimental studies, because it is often impossible to use random assignment for many questions of interest. Can you really look at the effects of parenting style on subsequent emotional development, by randomly assigning kids to parents!? Same goes with any long-term study on things like diet and lifestyle – you can sometimes sorta force people to eat some particular diet over a period of a few months or so, but that just isn’t going to work for the decades it likely takes for most diet effects to really impact overall health outcomes. There are also

other important ways of eliminating further possible third variables (typically called **confounds** in this context) from experiments, but random assignment is the most important (see Figure 1.1 for a diagram of the overall logic).

In summary, each of these different techniques is most appropriate for different kinds of questions, given the different tradeoffs. The key thing as a student and a citizen is to understand the limitations of any given study, so you can make an informed decision about what it really means. And don't expect the media to do this for you. Seriously, look at *any* correlational study on health / diet / etc and see how clearly the story, or the original article, discusses the limitations on any kind of causal implications from the study.

Neuroscience methods

Methods in neuroscience (and cognitive neuroscience) tend to be either correlational or experimental. The vast majority of **neuroimaging** studies are purely correlational, measuring the neural correlates of various different tasks or other manipulations performed while participants are in the brain scanner (we'll learn more about these scanners in the next chapter). By now, the neural correlates of just about every possible human activity (yes, including sex) have been measured in a scanner. But because of the correlational nature of these results, it is difficult to know whether the recorded brain activity is just *epiphenomenal* (i.e., just along for the ride), or whether it is really causal and somehow *responsible* for the behavior in question.

To attempt to address this causality question, scientists have used various forms of electrical and magnetic stimulation, which can disrupt or enhance neural firing in a relatively localized region of the brain. For example, **transcranial magnetic stimulation (TMS)** applied over the primary motor cortex can cause your muscles to flinch. However, just as with other experimental studies, the resulting brain states after TMS are not very "naturalistic", and it becomes difficult to interpret whether any changes in observed behavior are due to the disruption of the "normal" functioning of that brain area, or whether they just reflect the weird stuff that happens when you tweak that brain area in a completely unnatural way.

In animal neuroscience, much more precise causal inferences can be made by employing "invasive" techniques, such as directly cutting out different parts of the brain, or using modern **optogenetic** techniques to instantly and reversibly activate or deactivate a given population of neurons. These optogenetic techniques allow very specific populations of neurons to be targeted, and have produced a powerful new wave of causal empirical data, showing that very precise manipulations to very specific neural populations can sometimes have impressive overall effects. However, often even these results are over-interpreted and one must look very carefully for confounds in the resulting activity of other neural populations. Virtually every neuron in the brain is within a few synapses of every other neuron (i.e., the "6 degrees of separation" (from Kevin Bacon) phenomenon), so it remains very difficult to isolate what each specific subset of neurons is uniquely contributing. Indeed, as we'll see in the next chapter, the very premise of isolating specific functions may be entirely misguided.

Finally, animal neuroscience also affords much higher-resolution neuroimaging techniques which can resolve the activity of individual neurons, while also recording many such neurons at the same time. Such techniques provide the most powerful descriptive methods for characterizing what neurons actually do, and historically have been some of the most important data for fueling our theorizing and understanding of how the brain works.

Thus, truly each different type of technique plays a critical role in the overall arsenal of science.

Statistics

Finally, it is useful to be aware of the most widely used statistical techniques in psychology and neuroscience. Here is a brief overview:

- **Descriptive Statistics** – like descriptive methods, descriptive statistics are used to describe data, and differ from **inferential** statistics which are used to *infer* causality or correlation, as described below. The primary descriptive statistics are probably familiar to you: *mean*, *median*, *mode*, *range* and *standard deviation*. For a *normal* (bell-shaped, *gaussian*) distribution, the mean, median, and mode are all the same, and they tell you where the *middle* of the distribution is (i.e., the "average" person, etc). It is only when the distribution is *skewed* that they differ, with the mode and median being less

“pulled” by the long-tailed side of the distribution. You may have heard of income being reported in terms of medians – this is because income is a skewed distribution, with progressively fewer people making a *lot* more money than the mass of the “middle class” and below. The median and the mode more accurately capture this “middle class” salary because they don’t get pulled upwards as much by all the rich people (Figure 1.2).

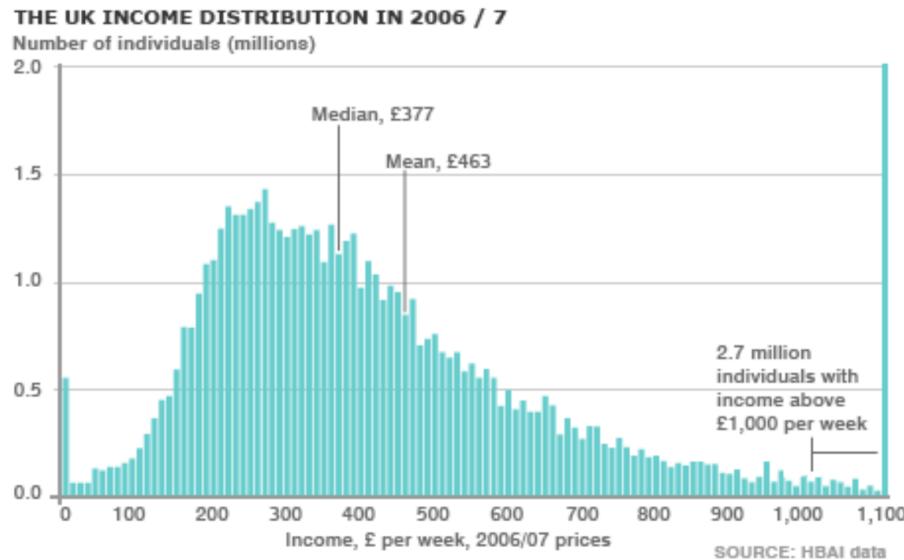


Figure 1.2: Mean, Median, and Mode tell different stories when the distribution is skewed (in this case, it is *right-skewed* – the skewer is the long tail to the right). The mean is pulled up by the tail much more than the median or mode, which do a better job of capturing the “middle class” income.

- **Correlation Coefficient and Scatterplots** – these are the primary tools for correlational studies. The correlation coefficient is a number, typically labeled r , which goes between -1 and 1, where -1 represents a perfect negative correlation, 0 is the complete absence of a correlation, and 1 is a perfect positive correlation. Importantly, both a strong negative and a strong positive correlation are equally important statistically, and indeed you can almost always just flip one of your variables around and turn one into the other (e.g., height vs. weight is positive, but “shortness” vs. weight is negative). A scatterplot simply plots the value along each variable (one on the X or horizontal axis, and the other on the Y or vertical axis), with each dot representing a different person (or whatever else is being measured). Thus, you can usually directly see the strength of the correlation in the shape of the “cloud” of such points (Figure 1.3).

One critical “pro tip” for looking at such scatter plots is finding “outlier” points that might be carrying a huge amount of weight. Just as a person sitting further out on a see-saw has more impact than one sitting further in, data points that are far away from the center of the cloud carry much stronger weight, and if they happen to lie along one of the positive or negative diagonals, they can produce a strong apparent correlation, even when all the rest of the points in the middle are clearly just milling about and going nowhere in relation to each other.

- **t-test, F-test (ANOVA) and the GLM** The “Student’s” t-test is the most basic of the *inferential* statistics used in experimental studies. It is *not* so-named because it is only for use by students, but rather it was the pen-name of the guy who invented it (William Gosset), to improve the quality of beer brewed by Guinness brewery in Ireland, no less! Too bad it isn’t called the “Stout” t-test. Anyway, it basically tells you if the difference between your experimental group and your control group is big enough to *not* be due to random chance. Thus, in applying this test, we “reject the null hypothesis” that our data is just random noise, but, critically, we’re not actually *proving* that our favored hypothesis is correct. We’re just saying it is relatively unlikely to be pure noise.

There are more “advanced” versions of this test, specifically the F-test used in the ANOVA (analysis

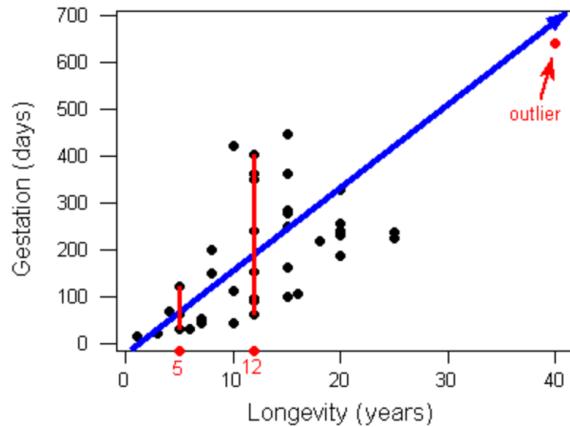


Figure 1.3: Scatterplot showing the positive correlation between length of gestation in the womb and overall lifespan, for different species of animals. The Elephant in the figure is the outlier, carrying undue amount of weight on the overall correlation coefficient. In this case, it is actually consistent with the rest of the data, but sometimes it is not, and yet the correlation still looks positive according to the r value. Thus, it is *essential* to always plot your raw data and ensure that the summary statistics are reflective of real aggregate effects!

of variance) procedure, and the full *generalized linear model (GLM)*, which can tell you about the importance of multiple different factors and their potential interactions.

You may have heard about the *replicability crisis* in various fields of science, including psychology, where many results that were thought to be “true” have “failed to replicate” – meaning that the original paper(s) reported a *significant* t-test result, and the subsequent ones did not (they instead found results consistent with pure noise). This is actually to be expected about 5% of the time, given the standard for publication is set at this 5% level. However, when you take into account how science is *actually* done, there are major systematic biases that enter into the process, which are not taken into account by these statistical tests, such that the actual effective probability of publishing garbage is closer to 50%!

There are now important changes afoot to combat the worst of these biases, and help ensure that this garbage probability goes back down to closer to 5%. But 5% itself is still a rather large number – in physics the standard is one in 3.5 million! And, amazingly, results that end up going into the “garbage” pile appear significant at levels below this standard, so randomness can sometimes be a challenging foe.

Conclusions

All of science, but especially psychology, is challenged by the primacy of human subjectivity. We are each sovereign nations unto ourselves, and science is more like policy making at the UN: it relies on slowly building up consensus among fundamentally capricious actors. The scientific method, however, provides a recipe for attempting to find the key *consistency* across people and across time that builds the foundation for the growing objective understanding of our world. In practice, if you don’t worry too much about all these philosophical issues, progress is being made. In the next chapter, we dive into the wealth of relatively new knowledge we have gained about how the brain works, building up a significant objective portrait to pair with the one painted by our own subjective worlds.

Summary of Key Terms

This is a checklist of key terms / concepts that you should know about from this chapter. As we’ll learn in the memory chapter, it is a great idea to test yourself on what was said about each of these terms, and then go back and double-check – that provides both beneficial repetition and also the *testing effect*.

- Subjective vs. Objective perspectives and the history of psychology
 - Rene Descartes, Cogito ergo sum, dualism

- Materialism, objective reality
 - Wundt, James and *introspectionism*
 - Watson, Skinner, Pavlov and *behaviorism*
 - Newell & Simon and *cognitive, information processing* approach
 - Rumelhart & McClelland and the *cognitive neuroscience* approach
 - Neural correlates of consciousness
- Scientific Method
 - Observation, hypothesis, data, replication, analysis, conclusions
 - Consistency over time, across people
 - Paradigms, social, psychological forces in science
- Methods:
 - Descriptive, Correlational, Experimental: pros and cons
 - Third variable problem
 - Power of random assignment
- Statistics:
 - Descriptive: mean, median, mode, range, standard deviation.
 - Correlational: r , scatterplot
 - Inferential: t-test, F-test (ANOVA)

Chapter 2: Neuroscience

From a materialist, neuroscientific perspective, *everything* that happens in your mind is due to underlying physical processes taking place in your brain. As we discussed in the last chapter, this does *not* mean that we can *reduce* your mind to the brain, but it does mean that there is a really huge mystery here (arguably the greatest mystery of all): how is it even remotely possible for a physical system to produce the amazing subjective delights (and terrors, and everything in between) that we all experience?

We start with a time-honored scientific approach: reduce the problem to the simplest possible system that exhibits the relevant behavior, and see if that makes it easier to understand. Consider the two gears as shown in Figure 2.1. As elaborated in the figure caption, there is something kind of “magical” that emerges out of the interaction between the two gears, which cannot be reduced directly to either gear separately. These **emergent** properties depend critically on the relationship and interaction between the two different parts – their relative sizes, rotational speeds, etc. If the larger gear interacted with a different, even larger gear, the overall system of interacting gears would exhibit very different emergent properties. Thus, you really can’t isolate these emergent properties to either gear in isolation. Furthermore, the actual material that the gears are made of is largely irrelevant, as long as it is reasonably solid. Thus, there truly is some kind of seemingly mysterious new “substance” being created out of this interaction, which can “transcend” its material basis. And yet, it nevertheless depends entirely and directly on having an actual material basis – e.g., the *picture* of those gears doesn’t work at all like two actual gears!

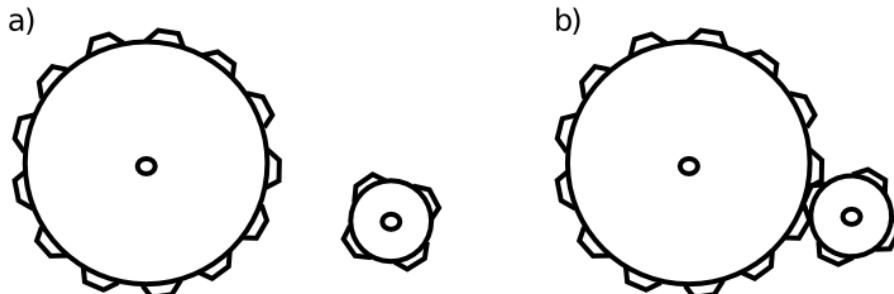


Figure 2.1: The principle of *emergence*, simply illustrated. The gears on the left do not interact, and nothing interesting happens. However, on the right, the interaction between the gears produces interesting, useful phenomena that *cannot* be reduced to the individual gears *separately*. For example, the little gear will spin faster, but the larger one will have higher torque at its axle – these properties would be entirely different if either gear interacted with a different sized gear. Furthermore, the material that the gear is made from really doesn’t matter very much – the same basic behavior would be produced by plastic, metal, wood, etc. Thus, even in this extremely simple case, there is something just slightly magical and irreducible going on – when two gears get together, something emerges that is more than the sum of the parts, and exists in a way independent of the parts, even while being entirely dependent on actually *having* those parts to make it happen. This seems like a good analogy for the relationship between the mind and the brain (but where the complexity and number of interactions is magnified billions of times over).

This simple two-interacting-gears scenario captures the strange relationship between mind and brain, where the mind depends entirely on the brain, and yet it fully transcends it. As we’ll see in a moment, the brain has billions of tiny, interacting parts (*neurons*), which, like the gears, interact in ways that produce emergent properties transcending their material substance. Moreover, there are so many neurons in the brain, and each one interacts with so many *other* neurons (receiving roughly 10,000 inputs and sending a similar number of outputs), that there is a vastly greater degree of emergent interactions taking place in the brain compared to our simple gear example. Thus, although it is essentially impossible for us to wrap our own minds around it, it should be possible to at least imagine in a vague way how something as fantastic and complex as the mind could indeed emerge out of all those billions and billions of interactions taking place every nanosecond, right inside your very own brain.

To try out another metaphor, you can also think about the brain as a massive LEGO set, with parts that *learn* to interconnect with each other in myriad ways. As you might have experienced in your youth, the number of different ways even a small pile of LEGOs can be combined to make different things quickly

exhausts the imagination. This *combinatorial explosion* of possibilities is an essential feature of the brain – our neurons can be interconnected in so unimaginably many different ways, that the possibilities are effectively infinite. Due to the explosive nature of combinations, even small numbers of elements can be combined in more different ways than there are atoms in the universe. To see this for yourself, type in $69!$ (factorial) on your calculator (or just google it), and you'll get a number that is $1.7\dots$ with 98 zeros! This factorial function lies at the heart of combinatorial explosion, and gives a rough sense of the number of different combinations of 69 parts. You can't even begin to conceive of (or even calculate) the value of 100,000,000,000! (i.e., the factorial of the 100 billion neurons in your brain).

Thus, your brain is one of the most complex systems in the universe – it doesn't have the most atoms, molecules or cells of course, but the key difference between your brain and an equivalent lump of gray jello is that the exact configurations of all those molecules and cells in your brain *mean something special* in terms of your unique memories, knowledge, etc. So all those different combinations and connections of neurons are *meaningfully distinct*, whereas you can rearrange the goo in the jello brain any which way and nobody would notice the difference.

Simple Neurons Make Complex Work

The magic of LEGO is that all the different parts interconnect using a single, simple principle (the stud fits tightly in a corresponding hole), so that you really can make all those different combinations work. The same is true of the brain: each neuron operates according to surprisingly simple, easily-understood principles, and the power emerges through all the interactions and combinations of these simple parts.

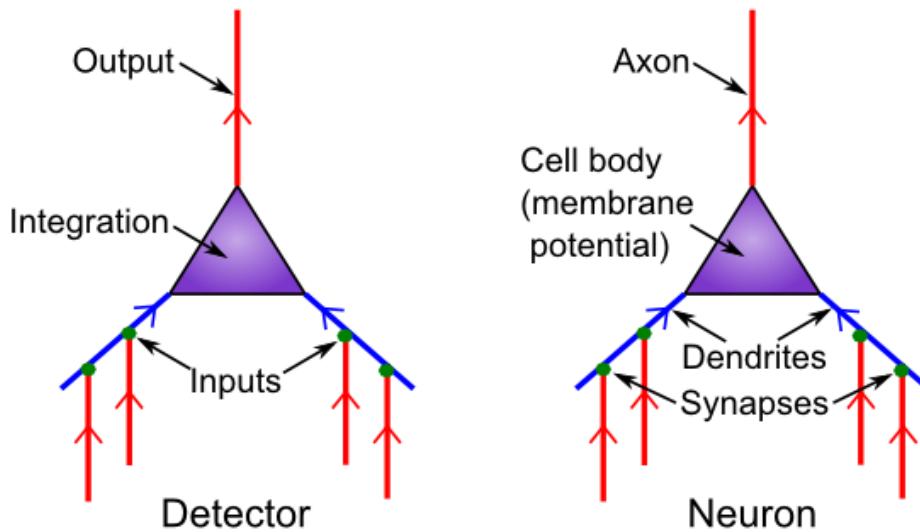


Figure 2.2: The neuron as a detector. Inputs come in via *synapses* connecting the *axons* of other neurons to the *dendrites* of a given neuron. This neuron *integrates* these inputs, resulting in an overall *electrical potential* (called the *membrane potential*, because it is the electrical difference between the inside and outside of the neuron's cell membrane), in the cell body. At the start of the axon (the *axon hillock*), a critical *go / nogo* “decision” is made – if the membrane potential is sufficiently elevated, then the neuron triggers an *action potential* (aka a “spike”), which races up the axon and delivers its signal to the many thousands of other neurons that are “listening” to this signal, via their own synaptic connections. Thus, the essence of neural function is *communication* – neurons are highly social little things, and our brain is really a huge social network of chattering naybobs.

The easiest way to understand what neurons are doing is in terms of **detection** (Figure 2.2). A neuron acts much like a smoke detector, constantly sampling its local environment (i.e., its inputs from other neurons), and looking for some set of incoming signals that indicate that something *important* might be going on (e.g., a fire in the case of the smoke detector). When it detects whatever it is looking for, it sends a signal out to other neurons, alerting them to the news, so they can incorporate this as one of many other pieces of information that they are themselves sampling in their own detection process. And so on, and so on...

Examples of the kinds of things different neurons have been shown to detect include: faces, specific people's faces (e.g., a famous case of a neuron tuned to Halle Barry, and another for Bill Clinton; (Quiroga et al. 2005)), eyes, letters, numbers, houses, different levels of visual depth, specific sounds, etc. Basically, anything that you can be aware of when looking out at the world is the result of neurons detecting those things from among all the possible configurations of visual features, including the words you're reading now, or your laptop, or your phone, or that pizza slice... everything!

Each of the major biological parts of the neuron take on a clear functional role within this overall detector model:

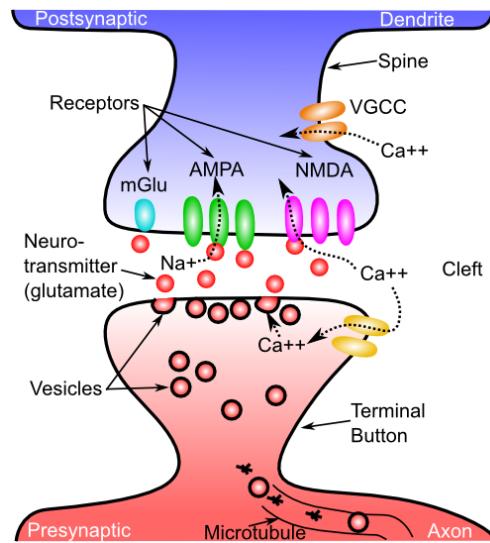


Figure 2.3: Details of the communication process across the synapse between a sending axon and a receiving dendrite. Neurotransmitter is released from the *terminal button* (pronounced French-style by those in the know), and binds to corresponding *receptors* on the dendrites, causing them to un-twist and thus open up small *channels* that allow electrically-charged *ions* to flow into the receiving neuron. Once neurotransmitter is released, it is taken back into the axon (*reuptake*) and is also broken down by enzymes, so that it tracks the rate of spiking by the sending neuron, and doesn't just hang out indefinitely. In addition to the primary excitatory *AMPA* channels that bind *glutamate* and allow Na^+ to enter, glutamate also binds to *NMDA* and *mGlu* receptors that are involved in learning, and other synapses use other neurotransmitters such as *GABA* which are inhibitory and allow Cl^- ions to enter.

- **Synapses:** are the tiny gaps between neurons, where the output signals from one neuron cross over and become the input signals to the next (Figure 2.3). Most synapses are *chemical*, involving the release of a *neurotransmitter* from the *presynaptic* axons, which then bind to *receptors* on the *postsynaptic* dendrites. These receptors twist open as a result of neurotransmitter binding, and allow *ions* (i.e., electrical charge) to flow into the dendrites, through the resulting open channels. The most common neurotransmitter in the neocortex is *glutamate*, and it opens up *AMPA* receptors, that allow sodium (Na^+) ions to flow into the dendrites, thus creating a *positive* electrical potential in the receiving neuron. After release, neurotransmitters are taken back into the terminal button (*reuptake*) and broken down for re-use by enzymes – this ensures that the amount of neurotransmitter binding accurately reflects the spiking rate of the sending neuron. All this complex-sounding machinery and terminology is actually very simple: Neurons like to excite other neurons by sending them exciting signals! The basic machinery is chemical and electrical, but the bottom line is just: how strongly do the input signals to a given neuron excite it? This is determined by the detailed properties of each of the roughly 10,000 synapses coming into a given neuron. This point bears emphasis, as we will return to it repeatedly: **the pattern of its synaptic connections determines what a given neuron detects, and thus, ultimately, what the brain knows.**
- **Dendrites:** provide a broad tree-like (dendrite literally means tree-like) "arbor" for all the synaptic inputs into a neuron, and they funnel the resulting electrical charges up into the cell body. This funnel-like property, illustrated in Figure 2.4, is the origin of the **compression** principle of brain

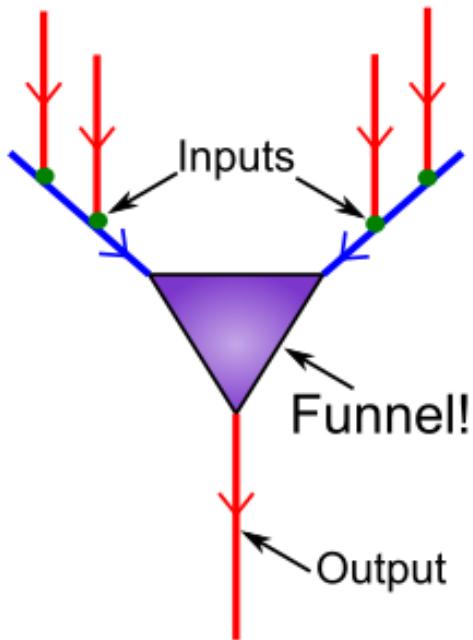


Figure 2.4: The neuron as a *funnel*, *compressing* its 10,000-odd inputs down into a *single* output signal, conveyed through its axonal output. This is the genesis of the *Compression* principle of brain function.

function, one of our three C's. As we noted in the introduction, each neuron is compressing its 10,000 different inputs into a single output signal, producing a roughly 10,000-to-1 compression factor. As an aside, one of the big debates in neuroscience is the extent to which these dendrites perform various kinds of more complex “processing” of their synaptic inputs, or simply convey the overall signal. There is evidence on both sides, and, as usual, the truth is likely somewhere in between.

- **Cell Body:** The neuron is a cell, and, despite its long tendrils, it has a cell body like other more compact kinds of cells, where the nucleus and other cellular machinery hangs out. It is here that all the dendritic signals converge, to produce the final compressed electrical potential that somehow summarizes everything coming into the cell at that moment. This electrical potential is called the *membrane potential* because the electrical signal is measured as a difference in electrical potential across the cell’s membrane (that fatty lipid bilayer that you might recall reading about in high school science). If this membrane potential is sufficiently excited, then special channels at the start of the axon (the *axon hillock*) will get *extra* excited and essentially flip a switch, causing the initiation of the *action potential* or *spike*. The details of this process were worked out by Hodgkin and Huxley in the 1950’s (Hodgkin and Huxley 1952), and have stood the test of time, forming the basis of modern detailed mathematical models of neuron firing.
- **Axon:** The spike propagates down the axon, effectively broadcasting this one signal out to the roughly 10,000 other neurons that it sends input to, continuing the great chain of communication among neurons. Note that the fact that the output is sent to roughly 10,000 other neurons doesn’t alter the *compression* property of neurons: it is just broadcasting the same signal, so there isn’t any new information there. Axons can have varying amounts of *myelin*, provided by helpful *glia* cells called *oligodendrocytes* (no you won’t be tested on those!), which serve to insulate the electrical “wire” that is the axon. Myelinated axons convey information more quickly, and *multiple sclerosis* is one of various disorders that involves the degeneration of this myelin, resulting in slowed signal conduction. There are many other forms of glia cells, but all of them are generally thought to play various supporting roles in the overall function of the brain, whereas the neurons are the “stars” of the show. These supporting roles are essential for keeping the brain functioning, and may affect various processes such as learning, but we’ll nevertheless

generally ignore them in this introductory treatment.

To summarize, each neuron is receiving a huge amount of input through its roughly 10,000 synapses, and it then compresses this all down into a single discrete spiking signal that it then broadcasts back out to the roughly 10,000 other neurons listening to its little story. Only when a neuron detects something “interesting” does it get excited enough to send this spiky signal out, and this *thresholding* is really the defining characteristic of a detector, making it respond *selectively*. Thresholding is just as important in neurons as it is in people: it can quickly get tiresome listening to someone with a *low threshold* who is always blabbering on about the most uninteresting things. The advent of Facebook and other forms of social media has greatly magnified this problem.

We’ll explore more about the kinds of interesting conversations neurons might be having in a bit, but first we’ll examine how this electrical magic operates within the neuron in more detail. This level of depth goes beyond most introductory texts, but understanding how this works helps us understand the effects of many different drugs, including alcohol and valium (as we discuss in the next chapter). Furthermore, a really simple analogy helps make it accessible, and this machinery ends up producing many of the **contrast** effects that are so central to our overall framework, so we’re motivated to take this brief detour.

The Tug-of-War in Your Brain

There are two major classes of synaptic inputs converging on each neuron: the excitatory ones described above (via the neurotransmitter glutamate opening AMPA receptors), and separate *inhibitory* synaptic inputs that are driven by a neurotransmitter called *GABA* which activates... *GABA* receptors, which allow negatively-charged Cl^- (chloride) ions to enter the cell. Thus, your brain runs primarily on table salt: $NaCl$, dissolved in water – we carry the ancient ocean around in our heads!

The inhibitory inputs come from an entirely separate set of specialized neurons known as *inhibitory interneurons*, which are somewhere between the principal, excitatory neurons and the glia in overall status within the pecking-order of the brain. These interneurons only act relatively locally, like glia, and they also play a largely *regulatory* role, regulating the overall level of electrical excitation surging through the brain. In contrast, the main excitatory *pyramidal* neurons (which constitute roughly 85% of neurons in the neocortex) can broadcast their exciting messages over long distances to far-flung regions of the brain, and are regarded as the primary *information processing* neurons (i.e., they are primarily responsible for all the chatter and compression going on).

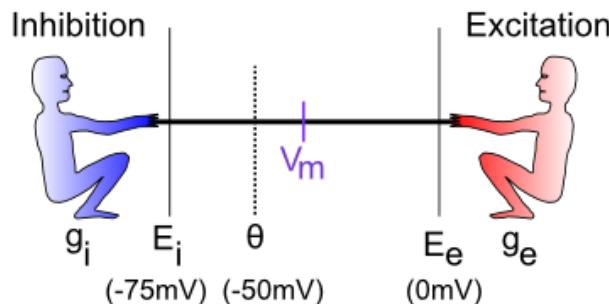


Figure 2.5: The Tug-of-War between excitation and inhibition, producing a beneficial balancing act, and a major source of **contrast** coding in the brain. Inhibition pulls the membrane potential (written as V_m , where V =voltage and m =membrane), down toward the *resting potential* of roughly -75 mV via the influx of negative Cl^- ions. Excitation pulls up toward roughly 0 to +55 mV (depending on type of neuron) via the influx of positive Na^+ ions. Thus, the components of ordinary table salt ($NaCl$) are driving this perpetual battle inside every one of your neurons. Theta Θ represents the *threshold* electrical potential, above which the neuron will fire a spike. The ability to do so depends on the *relative* balance between excitation and inhibition, not the absolute levels, and that is what makes us so sensitive to relative contrasts.

Inside each neuron, excitation and inhibition are forever locked in a pitched battle, which can be pictured as a tug-of-war, with each side pulling with varying strength, but each side always pulling in the same

direction (Figure 2.5). The “pitch” on which this battle is taking place is the amount of electrical charge in the cell, i.e., the membrane potential. The excitatory end is always pulling this potential upwards, while the inhibitory side is pulling it back down, and the actual potential represents the balance between these two forces.

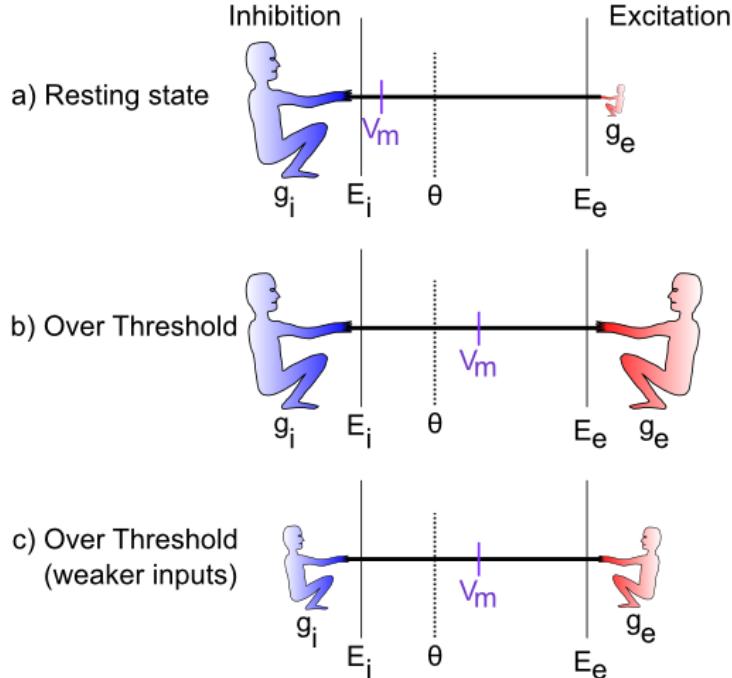


Figure 2.6: Illustration of the contrast or relative nature of the tug-of-war: only the relative strengths of excitation and inhibition matter, not their absolute values. Thus, neurons respond to the contrast between their inputs and overall levels of inhibition, which typically represent the the rough average of activity coming into a given brain area.

The essential point here is that **contrast**, or **relativity** emerges as the result of this tug-of-war battle. Specifically, it doesn’t matter how strong the two different sides on the tug-of-war are in absolute terms – all that matters is the *relative* strength of the two sides (Figure 2.6). Excitation could be relatively weak, but if inhibition is also weak, then the net balance between the two will be the same as if each was proportionally stronger.

Typically, the amount of inhibition is roughly proportional to the “average” amount of activity in the brain in any given area, so in effect, each neuron is effectively comparing how excited it is against this overall “average” level. Only those neurons that are getting *above average* level of excitation will actually get excited enough to fire spikes.

In real-world terms, this “average” inhibition is very much like the amount of money that your peers are making (or the amount of fun they appear to be having on their various social media accounts) – it forms the baseline or standard against which you measure yourself. Likewise, neurons are constantly comparing themselves against *their* peers, and all of the spiking going on in your brain is therefore always and inexorably *relative* to these peer-standards. For example, when you step outside into the bright sunlight, all the visual neurons suddenly get a huge wave of excitation relative to the dim indoor light from before. But you avoid suffering an epileptic seizure from all that excitement because those inhibitory interneurons are also getting this wave of excitation, causing them to send a proportional amount of damping inhibition on the party, keeping everyone in balance and on a more level keel. Yes, inhibition is the wet rag of the brain, but without it, you really would be suffering from seizures all the time.

In fact, this balancing act between excitation and inhibition is so important for overall brain function, that our brains are perched on a kind of “knife edge”, and the relatively high incidence of epilepsy in the population is likely a result of the fact that it is really hard to get this balance exactly right. And too much

inhibition has very bad consequences as well (indeed, it literally “depresses” your brain and makes it difficult for you to do anything).

The main treatments for epilepsy involve activating the GABA inhibitory system more strongly, thus altering this fundamental tug-of-war balance, and drugs that activate GABA are also used as anesthetics. The reason alcohol makes you sleepy and pass out if you have too much, is that it also activates GABA inhibition too – paradoxically that increased inhibition in your brain results in *behavioral disinhibition* – your anxieties get inhibited along with everything else. And people never take the drugs that go the other way, and reduce GABA inhibition – because they will definitely cause epilepsy and fry your brain!

Another key property of neurons is that they exhibit **adaptation / accommodation / fatigue** (all different terms for the same basic phenomenon), where their response to the same level of excitatory input decreases over time. This is not just some metabolic side-effect of neurons getting “tired” – it is a core computational property of how neurons process information. Specifically, it reflects a strong bias for emphasizing what is *new*, and discounting anything that has already been active and processed. It is an important form of *contrast**, contrasting the new against the old.

In summary, two out of the three of the core principles of this textbook, **compression** and **contrast**, emerge directly out of the basic function of neurons. As we discussed in the Introduction, we can trace the implications of these core neural properties all the way through the full scope of psychology and behavior. The Perception chapter will provide particularly compelling demonstrations of how compression and contrast play out in our perceptual lives – the story of perception is fundamentally the story of compression and contrast.

By thinking in these terms, we have managed to dramatically simplify our understanding of the brain, creating an almost transparent, level-spanning way of going from single neurons on up. However, none of this contradicts the emergence and complexity discussed at the outset. Instead, these principles just capture the overall general tendencies and propensities of the brain, but within that broader scope, there is a wild, complex, bubbling jungle of intertwined conversations and chatter constantly unfolding within your brain, thinking all manner of complex and ineffable thoughts. Next, we’ll move up from the level of individual neurons and start to think about how all these principles might play out in terms of how different brain areas are organized to facilitate effective overall behavior and cognition.

Large-Scale Brain Organization (“Gross” Anatomy)

We will attempt to answer two closely interrelated questions about the large-scale organization of the brain, one of which is relatively easy, and the other which remains rather more murky. The easy question is: “what are the obviously separate parts of the brain, which have a distinct evolutionary and structural basis?” The hard question is: “to what extent do any of these brain parts, or regions within these brain parts, support a distinct kind of overall function?” This latter question of *functional specialization* is challenging because neurons are so massively interconnected and interdependent on each other, that it is hard to clearly isolate any specific function. It is like any kind of team sport: we are tempted to focus on a few specific star players, but, really, the team depends on every player and the quality of their interactions (just like the gears in Figure 2.1, and any emergent system). For example, the quarterback on a football team can either look really good or bad as a function of how good the offensive line is, but nobody gives that line sufficient credit, focusing instead on the singular, more glamorous quarterback.

Historically, Karl Lashley in the 1920’s concluded from his extensive studies lesioning different parts of rat brains, that the brain is an *equipotential* system, operating according to some kind of *mass action* principle. That is, all areas contribute roughly equally to the overall function, and all that matters is how much overall neural tissue is intact – the more the better. This idea strongly conflicts with everything we know about mechanical systems, where each part has a specific, well-defined function. It also conflicts with the principle of *compression* as students of the brain: we want a simple, easy-to-understand picture about how things work. This is reflected in the ubiquitous drawings of the brain carved up into discrete functions (e.g., Figure 2.7), and in the discredited approach of *phrenology* which attempted to associate functions with different bumps on the skull. Critics have derided many recent neuroimaging studies as *neo-phrenology* because there is still this strong tendency to ascribe discrete functions to individual blobs of brain that “light up” when people are doing different tasks in the brain scanner.

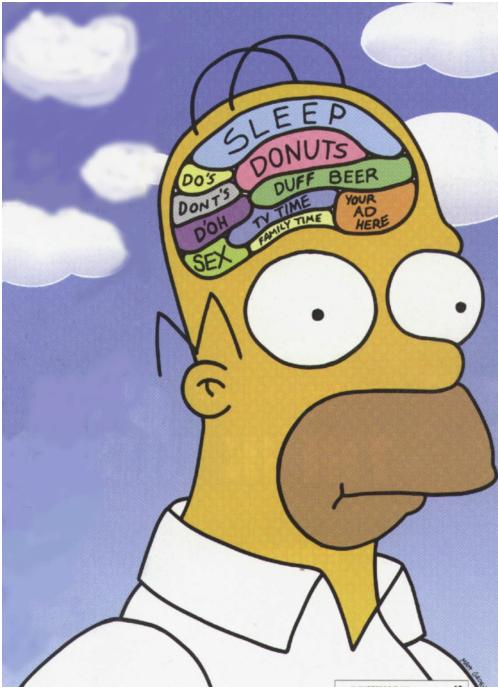


Figure 2.7: The brain would be a lot easier to understand if each part had an easily-labeled, distinct function.

As usual, the truth is somewhere in between these extremes. Even though neurons are massively interconnected and interdependent, and overall function emerges through these interactions, there is evidence that different brain areas are differentially important for different functions. But the level of functional specialization is much more *partial* and *overlapping* than completely distinct. One way of thinking about this is in terms of the saying that “All politics is local”, which is as true for the brain as it is for people, perhaps more so: neurons can’t pack-up and move to a different part of the brain. Instead, they are like the 85% of Pittsburghers who live their entire lives in the same little neighborhoods. This means that different neighborhoods can develop their own special “personalities” and focus on detecting particular kinds of signals.

On the other hand, the excitatory neurons in the brain also send out long-range connections. Network theorists characterize these as “small world” patterns of connectivity, such that, in the end, every neuron is only a few synapses away from every other neuron. This then limits how much “neighborhood funkiness” can develop. So, again, the brain is all shades of grey, not black-and-white: different areas are somewhat specialized, but also very interdependent.

The Big Brain Chunks

Figure 2.8 shows the different brain regions that can be easily distinguished based on evolutionary history and obvious structural differences. Keeping in mind the above qualifications, the overall functions of these areas are as follows:

- **Cerebral cortex (Neocortex):** This is the most important part of the human brain, supporting all of our special human abilities to think, read, talk, reason, plan, etc. We are exclusively conscious of activity in this part of the brain – this is where “you” live! Most of what we discuss throughout the textbook is focused on this part of the brain, and all of what we said above about neurons is focused specifically on this part of the brain (other parts may differ in various details, but the general properties are common across the brain). Often, we’ll just refer to this as “the cortex”, although technically “cortex” means “sheet-like” and other brain areas also have a “cortical” kind of organization. Anatomically, this sheet-like nature of the neocortex is more evident in smaller-brained animals – in humans, the neocortex is so greatly expanded that it is all folded over on itself. It is the wrinkled sheet upon which all of our hopes and dreams rest. We’ll go into more detail about the different lobes and

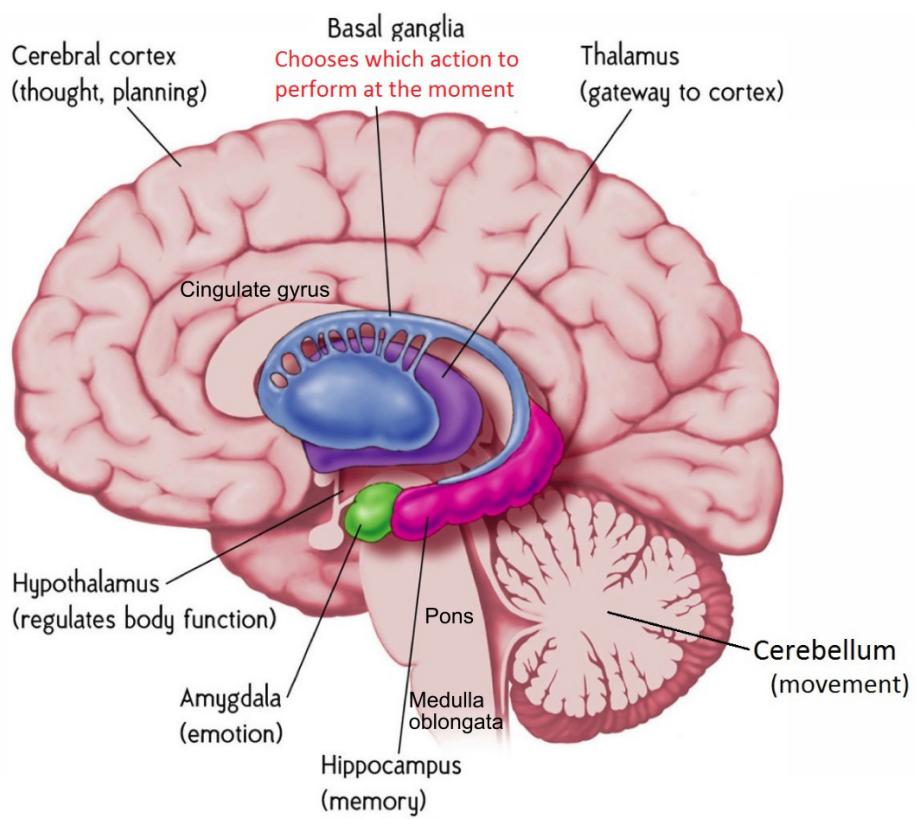


Figure 2.8: Large-scale (“gross”) brain structures and their overall specialized functions.

their relative functions in the next section.

- **Thalamus:** Functionally, the thalamus is completely intertwined with the neocortex, and neither can function without the other. There are massive, bidirectional interconnections between every part of the neocortex and the thalamus. Some parts of the thalamus “relay” information from the senses up into the primary sensory areas of the neocortex. For example, the *lateral geniculate nucleus (LGN)* receives most of the output from your eyes, and then sends that up into your *primary visual area (V1)* in the neocortex. However, it is not *just* a relay: V1 also sends massive “top-down” connections back into the LGN, and these serve to focus attention and organize all the low-level visual signals into a more coherent overall “picture”. As visual information processing proceeds up to higher levels in the cortex, the thalamus continues to play a critical role through a structure known as the *pulvinar*, which again has massive bidirectional interconnections with corresponding cortical areas.

The pulvinar has been implicated in attention, and also serves to coordinate different cortical areas by synchronizing brain activity in the *alpha* frequency (10 Hz or 10 cycles per second) – we’ll discuss these frequencies more later when we talk about sleep stages (the thalamus also plays an important role in sleep). Other areas of the thalamus are directly interconnected with both the frontal lobes and the basal ganglia, and are really inseparable from the overall function of those brain areas. Thus, overall, despite its anatomical separation, functionally it does not really make sense to think of the thalamus as a separate brain area from the neocortex – instead we should think of the *thalamocortical system* as a functional unit.

- **Basal Ganglia (Striatum):** This is a collection of different brain *nuclei* (chunks of neurons) that form a complete sequential pathway or loop from the neocortex and back up into different parts of the frontal lobes. Thus, like the thalamus (which is also a key part of this loop), it is hard to really separate the function of the basal ganglia from that of the frontal lobes of the neocortex, and damage to either of these areas produces very similar overall problems. Indeed, this *fronto-striatal* system is implicated in most of the major mental disorders that afflict us, including depression, anxiety disorders, ADHD, and OCD, as we’ll cover in detail in the chapter on mental disorders. As noted in the Introduction, this system is the most important player in the **control** component of our three C’s, and these disorders all involve a major element of control.

Anatomically, the input portion of the basal ganglia circuit is composed of the *Caudate Nucleus*, *Putamen*, and the *Nucleus Accumbens*, which collectively comprise the *Striatum*, which means “striped”. The striatum receives massive input from all over the neocortex, “digests” (compresses) it down into a basic “go” vs. “nogo” decision about whether to do something or not, and then sends that decision back up into the frontal lobes, by way of the thalamus. Thus, whereas the basal ganglia was previously thought to be more of a “habit learning” part of the brain, it is actually the real *decision maker* in your brain. Indeed, research shows that the basal ganglia makes a decision about what you’re going to do next about 1/3 of a second before you are consciously aware of it (Libet et al. 1983)! This reflects the fact that we’re only conscious of what is going on in the neocortex, and some people find it kind of unsettling that this “other” part of your brain is “making decisions on your behalf”. But really, you are *all* of your brain, not just the parts you’re subjectively aware of, and again, all of these areas are massively interconnected and interdependent, so don’t get too freaked out by this! Embrace your inner decision maker, which is responsible for those “gut feelings” that all so often end up being correct, even as they are often overridden by your over-analyzing conscious cortex.

The basal ganglia are unique in the brain by virtue of having by far the most *dopamine* receptors, and dopamine plays a critical role here by shaping the decision-making process according to what has worked, and not worked, in the past. Furthermore, the basal ganglia, particularly the nucleus accumbens in the *ventral* (bottom) part of the striatum, plays an essential role in controlling the firing of dopamine neurons, so that the overall dopamine signal reflects the *contrast* from your expectations, rather than raw reward or punishment itself. In the Learning chapter, we’ll see in more detail how this process works, in the context of *Classical Conditioning* and *Operant Conditioning*, which largely reflect the dopamine-driven learning processes taking place in the basal ganglia.

- **Amygdala:** The amygdala is a relatively small nucleus, which is named after the Greek word for almond (most anatomical labels describe either the shape, color, or texture of the brain structure), that plays an essential role in driving our emotional life. It is extensively interconnected with both the

basal ganglia and the dopamine system, and drives these systems to respond appropriately for positive and negative emotional events. For example, when a previously-neutral stimulus is associated with either a rewarding or punishing outcome in classical conditioning, the amygdala learns the association between the stimulus and this outcome, and drives dopamine firing and other behaviors to anticipate and prepare for the outcome (e.g., approaching yummy food and running away from fear-inducing scary stuff).

The Amygdala is also extensively bidirectionally interconnected with the neocortex, receiving sensory inputs and sending its emotional signals up to the medial and ventral regions of the frontal lobes, which are the emotion centers of your conscious world in the cortex. Thus, overall, the amygdala is a *hub* for emotional signals, interconnecting between lower-level brain stem systems such as the hypothalamus, and driving your high-level conscious emotional experiences. People with damage to this area don't necessarily have a complete absence of emotion, but they can't connect all the pieces together in an effective way, and often behave carelessly because they fail to anticipate the potential risks of their actions.

- **Hippocampus:** The hippocampus lives next door to the amygdala, and is essential for rapidly forming new memories of the daily events of your life (i.e., *episodic* memories). When you think of memory, mostly you're thinking of what the hippocampus does. Of all the brain areas we've considered so far, the hippocampus is the most strikingly specialized: highly selective damage to this brain structure can result in profound *amnesia* – particularly the inability to learn new episodic memories, but also the loss of at least a certain window of more recently-acquired memories. The famous patient H.M. (Henry Molaison) had his hippocampus lesioned surgically to alleviate epileptic seizures, and was unable to acquire new memories for the rest of his long life.

It is scary to contemplate H.M.'s case: so much of our sense of life's meaning is associated with making lifelong memories. Along with our emotions, our memories are the most cherished aspect of our subjective world, and it is truly horrifying to imagine losing this ability. Therefore, you should treat your hippocampus well: it is a bit of a "canary in the coal mine", and is often the first thing to go when you lose oxygen to the brain. Likewise, heavy drinking causes this area to lose function before others, resulting in memory blackouts.

As we'll explore in greater depth in the Memory chapter, the hippocampus has several biological specializations that enable its super-memorizer abilities, but these also result in it being a bit "hyper-sensitive", both biologically and functionally. Although the hippocampus is highly specialized, it nevertheless depends entirely on extensive input from the surrounding areas of the neocortex, which convey a massively *compressed* summary of everything going on in the rest of your brain. This high level of compression makes our memories relatively inaccurate and subject to many biases, but also extremely efficient.

In sum, the hippocampus essentially takes a "snapshot" of the current state of the brain, and later, when you want to recall some prior event, it can retrieve that snapshot and cause the rest of your brain to relive that moment. During recall, the hippocampus drives those same surrounding neocortical areas in the reverse direction from when the memory was initially encoded, again demonstrating the essential interdependence of all these different brain areas. Also, the hippocampus has somewhat separable "cognitive" and "emotional" components, with the emotional one extensively interconnected with the amygdala and those frontal emotional areas that strongly interconnect with the amygdala, and it plays a critical role in making your emotional responses appropriately responsive to different situations and contexts.

- **Cerebellum:** The cerebellum plays a critical role in learning to perform motor (muscle) movements in a smooth, efficient, and coordinated way. Anatomically, it is a kind of "mini brain" (that is what it's name means) tucked under the back of your brain, and it is also a "cortical" structure with a very distinctive sheet-like organization. In some ways, you can think of it as a kind of "hippocampus for motor learning", as suggested by the pioneering scientist David Marr in a pair of prescient papers (Marr 1969, 1971), which attempted to discern the functions of both the cerebellum and hippocampus based on their unique anatomical properties. Amazingly, his ideas have largely stood the test of time, and form the core of our modern conception of these areas.

Area	Learning Signal		
	Reward	Error	Self Org
<i>Primitive</i>			
Basal Ganglia	+++	- - -	- - -
Cerebellum	- - -	+++	- - -
<i>Advanced</i>			
Hippocampus	+	+	+++
Neocortex	++	+++	++
+ = has to some extent ... +++ = defining characteristic – definitely has - = not likely to have ... - - - = definitely does not have			

Figure 2.9: Learning Rules across the brain: some of the clearest differences between brain areas are in terms of the signals that drive learning in a given area. In particular, the basal ganglia and cerebellum each specialize on two of the most important types of learning signals: reward (and punishment) vs. error signals (which are *not* the same as punishment – these are instead detailed signals with specific information about exactly what didn't go according to plan in a motor action. The cerebellum is unique in having no dopamine innervation or receptors. More evolutionarily-modern areas incorporate multiple signals, and include self-organizing learning, which means learning that happens automatically all the time, as in the hippocampus automatically taking snapshots of cortical activity.

These brain areas are among the most functionally specialized, and both rely on a kind of “brute force” memorization strategy to achieve their special learning abilities. This brute force strategy requires a lot of neurons: half of the total neurons in your brain live in the cerebellum! An important consequence of this strategy is that it takes lots and lots of practice to really perfect any given motor skill (e.g., gymnastics, skiing, etc), because the cerebellum has to memorize each of the many different ways to perform a motor action.

The cerebellum learns to anticipate errors, awkwardness, and inefficiency in a given motor action plan, and sends well-timed corrective signals to prevent those from actually occurring. It receives error signals from a nucleus called the *inferior olfactory nucleus* (you can guess what it looks like), which drive a powerful error-correcting learning signal in the *purkinje* neurons that are one of the central actors in the cerebellar circuit. These purkinje neurons are truly amazing things, receiving over 100,000 different synaptic inputs (10 times as many as the typical neocortical neuron) – so many synapses are needed to be able to have distinct memories for each of those different motor action sequences. This form of error-driven learning is quite different from the “snapshot” memorization operating in the hippocampus, so these brain structures are also functionally distinct from each other, even though they both share the same overall brute-force memorization strategy.

Interestingly, the cerebellum and basal ganglia, which are both considered motor control systems, have almost no direct interconnections (a rarity in the brain, as we've seen) – but this actually makes good sense, because they each perform very different functions, at different time scales (Figure 2.9). The cerebellum deals with very fast “online” motor control at the scale of 10's of milliseconds, whereas the basal ganglia is more involved in the “outer loop” of deciding which of various possible motor plans to actually execute. Thus, the basal ganglia typically acts first to select the motor plan, and then the cerebellum takes over and ensures that the selected plan is executed to the best of your ability.

By analogy with the different roles in making a movie, the basal ganglia (together with the frontal cortex) is the *producer*, deciding what movie to make; the cerebellum is the *director*, who is there day-in-day-out on the set, dishing out detailed instructions to the actors to make it all look good; and motor circuits in the *pons* and other brainstem areas, on down into the spinal chord and the muscles, are the *actors*, actually carrying out the actions.

- **Hypothalamus:** This tiny structure plays a huge role in controlling your basic bodily functions, including eating, drinking, sleeping, arousal, sex, stress, immune response, etc. It is the kingpin in the *HPA axis* (hypothalamic-pituitary-adrenal), which is a system of interconnected structures that release hormones including *corticosteroids* in response to *stress*. The hypothalamus has many different nuclei,

each specialized for different domains, and some of these project up to the amygdala to drive emotional responses. For example, the positive reward feelings associated with eating and drinking come from the lateral hypothalamus, and these signals go into the amygdala and directly into the dopamine system, driving bursts of dopamine for (unexpected) positive events, like when a co-worker brings in leftover birthday cake to the office. The hypothalamus also receives top-down control signals from areas of ventral and medial frontal cortex, which can regulate the response to potentially stressful events, for example.

One fascinating line of research shows that rats who can *control* their exposure to mild electric shocks (by moving to another part of their cage), have significantly lower stress responses than the “yoked” rats that receive the exact same electric shocks, but have no control over them. Thus, the perception of control, which has been localized to those frontal cortical areas (consistent with the overall role of these areas in control more generally), is an essential factor in how the body responds to stressful situations (Maier and Watkins 2010). A clear real-world example of this is the difference between driving a car and riding along as a passenger – the driver typically experiences things as “under control” whereas the passenger is more likely to feel stress because the driver is going too fast otherwise being unsafe.

More generally, chronic exposure to negative, stressful situations over which a person has little perceived control can produce significant long-term mental health problems, leading to a kind of *learned helplessness* that is associated with depression (Maier and Seligman 1976). This kind of chronic stress exposure without perceived control is much more prevalent in people of lower socioeconomic status, and can produce a very unfortunate feedback loop of learned helplessness and significant health complications from stress. Understanding how these brain mechanisms work is thus of vital importance for addressing many pressing societal issues, again underlining the broader importance of psychology and neuroscience.

- **Brainstem Nuclei and Medulla Oblongata:** Finally, there are a number of different clumps of neurons in the brainstem that play critical roles in overall brain and body function. These are evolutionarily more ancient brain areas, like the hypothalamus, which have highly specialized functions. In computer terms, these are the core BIOS brain areas – the low-level hardware control areas. One group of such nuclei are collectively referred to as the *reticular activating system*, and include the sources of the major *neuromodulators* that “modulate” (alter) the functioning of neurons throughout the brain in various (often similar) ways:

- *ventral tegmental area* and *substantia nigra pars compacta*: dopamine – modulates learning in basal ganglia, other areas.
- *raphe nucleus (dorsal, median)*: serotonin – modulates arousal, sleep, mood.
- *locus coeruleus*: norepinephrine (noradrenaline) – modulates effort, engagement.
- *basal forebrain cholinergic nuclei*: acetylcholine (ACh) – modulates attention, arousal, learning (nicotine affects this system).

These core areas serve as master control knobs for the overall state of the brain, and are thus incredibly important and powerful. All of them receive extensive top-down projections from the frontal lobe, which thereby asserts its overall master control of these knobs, while also being subject to their effects. The mutual interdependence of all of these brain systems is evident even here at the lowest levels, and has many important implications for sleep, arousal, and other overall brain states.

Last but not least, the medulla oblongata wins the prize for the funniest name in the brain, but it is no laughing matter, providing essential low-level body control signals. Damage to this area often results in death. Enough said.

Functional Organization of the Neocortex

The neocortex is divided anatomically into four separate lobes (Figure 2-10), which can be given broad overall functional specializations that stem principally from the unique sensory inputs / motor output coming into / out of each lobe (each lobe gets one of the three major sensory input modalities, or drives motor output):

- **Occipital:** Receives the primary visual input from the LGN of the thalamus (in area V1 at the very back of the brain), and begins the processing of these inputs.
- **Temporal:** Extracts object identity information (e.g., face, pizza, laptop, etc) from visual signals

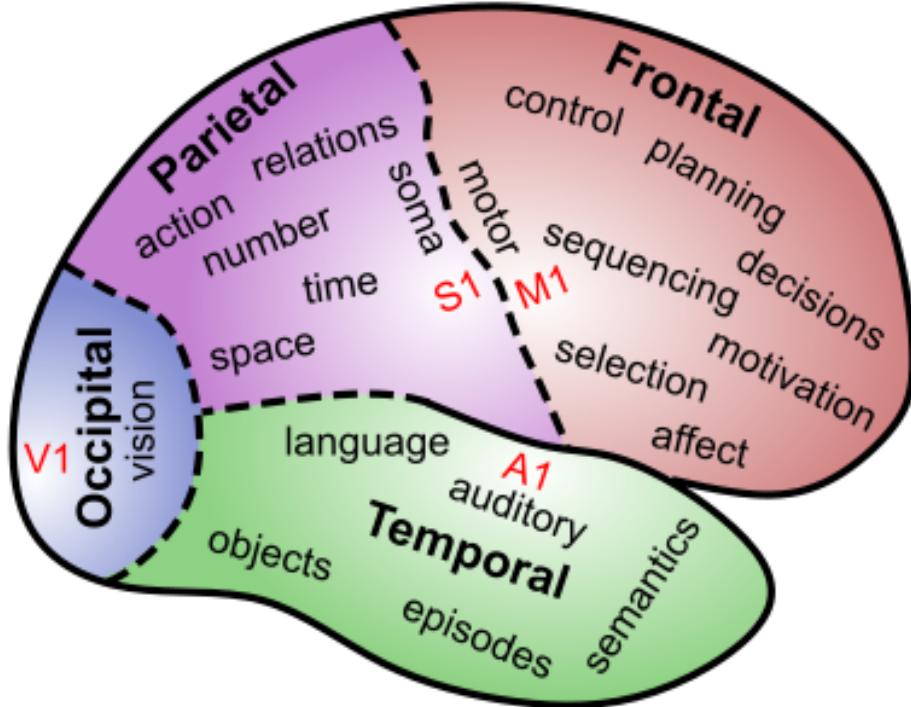


Figure 2.10: Functional specializations of the different lobes of the neocortex, stemming largely from their unique primary sensory / motor areas (V1 = primary visual cortex, A1 = primary auditory cortex, S1 = primary somatosensory cortex, M1 = primary motor cortex).

coming in from the occipital lobe, and connects those with auditory signals arising from primary auditory cortex (A1), which is in the upper (superior) portion of the temporal lobe. These connections form the initial basis for *language*, in terms of the ability to name objects recognized visually, and semantically understand the meaning of spoken words. The inner (medial) part of the temporal lobe connects up with the hippocampus, and is critical for assembling the *who-what-where* elements that define the episodes (events) of our lives, that the hippocampus takes snapshots of. The very tip of the temporal lobe (towards the front) is important for encoding our most abstract, high-level semantic knowledge (truth, justice, etc) (Lambon-Ralph et al. 2017).

- **Parietal:** Also feeds off of the occipital visual information, but in the service of guiding motor actions, by virtue of its position betwixt the occipital and frontal lobes, and the primary somatosensory inputs in area *S1*. *S1* is located just across the *central sulcus* (sulcus = groove) from the primary motor cortex, *M1* in the frontal lobe, and each has a matched *homunculus* (“little man”) representation of your entire body (Figure 2.11), which is distorted in its focus on the most important areas at the expense of others (e.g., your back doesn’t get a lot of neural space, whereas your fingers and mouth are very prominently represented). Because motor actions require proper positioning of your hands and body in space, the parietal lobe is where your understanding of spatial locations and relationships arises. Interestingly, these spatial representations are re-used for thinking about more abstract continuous quantities like time and number (Dehaene et al. 2004).
- **Frontal:** Is grounded by its *M1* primary motor outputs, which make this lobe focused on motor control across all levels of space and time. Progressively more frontal (i.e., *prefrontal*) areas encode progressively higher-level, extended action plans, to coordinate and organize the basic motor actions encoded back in *M1*. These higher levels of control require things like sequencing, planning, and decision-making, and as noted above, all of these functions depend critically on interactions between frontal cortex and the basal ganglia. You can think of the frontal cortex’s part of this interaction in terms of *generating possible action plans*, which the basal ganglia then evaluates according to its dopamine-driven learning

Sensory/Motor Homunculus

Copyright © The McGraw-Hill Companies, Inc. Permission required for reproduction or display.

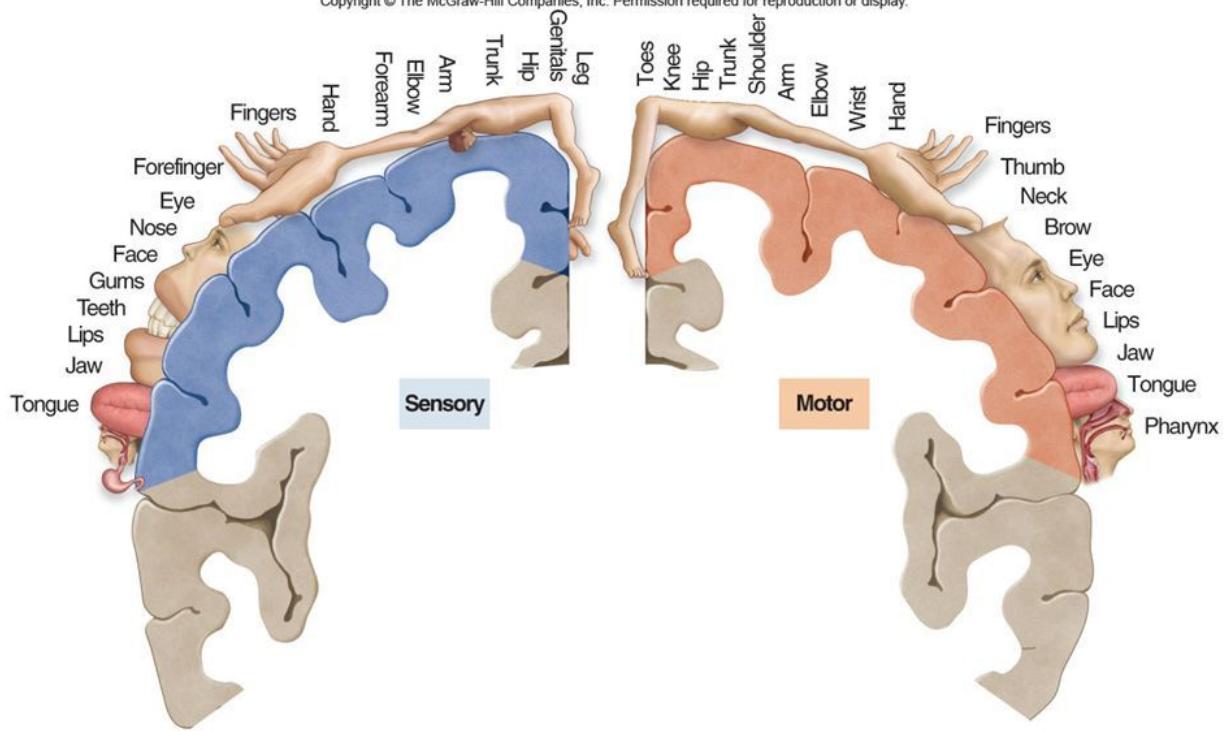


Figure 2.11: The coordinated *homunculus* ("little man") as represented across the primary somatosensory (S1) and motor (M1) cortex.

history, and it sends back up a strong signal selecting the plan most likely to maximize future reward and minimize punishment / cost.

If you have seen “Mad Men”, the frontal cortex is more of the “creative type” coming up with a new Ad pitch, and the basal ganglia is the tough customer critically evaluating the bottom line benefits and costs.

The frontal cortex has other “departments” that are also involved in representing these benefits and costs, in the inner (medial) and lower (ventral) parts, which are anchored by the inputs from the amygdala and other core visceral areas, including primary taste areas in the *insula*. These emotional (*affective*) and body-state inputs are essential for guiding the overall motor control and planning processes, to focus on the things that actually matter, thus giving the frontal lobes a primary role in goal-driven and motivated behavior.

Continuing the Ad agency analogy, the parietal lobe also plays a critical role, much like the art department, providing the sensory guidance needed to inform the action planning process (e.g., the initial sketches for the pitch). For example, the parietal lobe can provide a spatial map of a sequence of actions to be taken over time, to help in figuring out the best ordering of the individual steps in the sequence. Furthermore, it can represent the likely sensory outcomes of different possible action plans, in terms of both somatosensory and visual modalities (i.e., how my arm would feel if I moved it in a particular way, and where it would end up in space), which can then feed back to refine the overall motor plan.

Finally, the frontal connectivity with the temporal lobe (in the lower (ventral) parts of the outer (lateral) frontal lobe) also plays a critical role in representing all the important players and details of an overall plan: *who* and *what* will be affected, etc. The language functions (writing copy, communicating with others) and memory functions (recalling relevant past memories) of the temporal lobe are also engaged by this frontal planning system. In sum, the frontal lobes are characterized as the *central executive* of the brain (Baddeley and Hitch 1974), orchestrating all the other brain areas around the core actions and motivations that really matter.

Returning to the big question of functional specialization, we can clearly see the balance between different neighborhoods of neurons specializing on different kinds of information, but also depending critically on the work of other areas to get their own jobs done. In effect, the brain is just like any complex human organization (e.g., in a company, a university, the military, etc) – everyone depends to varying extents on the work that others are doing, but each person also performs some specific, specialized roles. Because neurons stay put in the brain, and neighboring neurons tend to be more strongly interconnected with each other, we can trace these networks of interaction and interdependency to help understand what each part is doing. Next, we’ll briefly examine the importance of another aspect of complex organizations: a hierarchical structure.

Hierarchical Organization

Figure 2.12 shows how this combination of interdependency and specialization plays out in the case of the visual pathway going from V1 up through the object recognition neurons in the inferior (bottom) part of the temporal lobe (*IT = inferotemporal cortex*, conveniently where *it* is recognized). There is an overall *hierarchical* organization to this pathway, such that the early stages detect simpler features (e.g., oriented edges in V1) while higher levels build on this to detect parts of objects in terms of collections of these features, and still-higher levels can then detect entire objects in terms of collections of features. Thus, by building up a cascade or hierarchy of detectors in this way, each performing their own part in a larger chain of *compression*, a very challenging overall problem can be broken down into simpler steps. Hierarchical organizations of this sort are ubiquitous and necessary for organizing, coordinating, and integrating the work of many individuals, whether it is people in the military or corporations, or neurons in the brain.

Neuroscience Methods

Finally, we conclude this chapter with a brief overview of some of the major techniques and methods used to understand how the brain works. We covered the issues of correlation vs. causation in neuroimaging and other such techniques in Chapter 1, so here we focus more on how these techniques actually work, and what

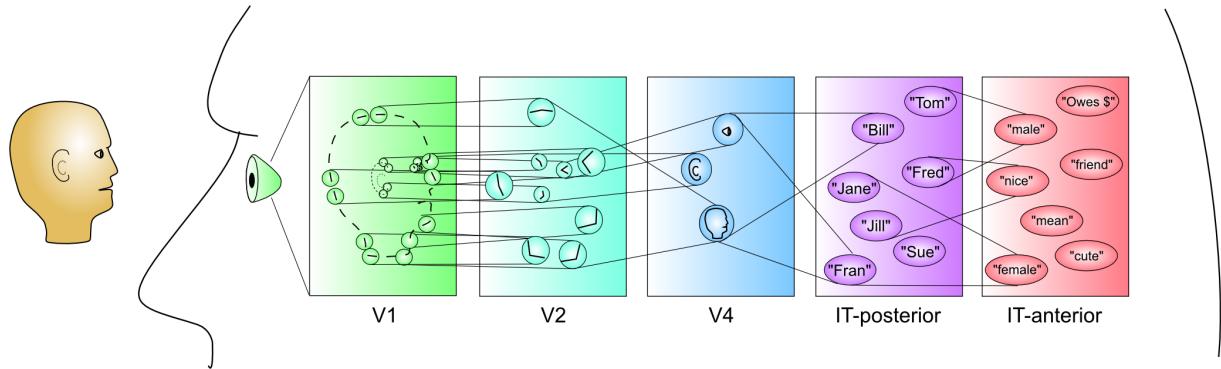


Figure 2.12: Hierarchical organization of detectors in the visual pathway going into the temporal lobe, supporting the ability to recognize (detect) entire objects, based on earlier levels detecting parts and features of parts. This shows the large-scale, cumulative effects of *compression* from very high-dimensional raw sensory inputs, to high-level, succinct interpretations of the world. Although a highly simplified cartoon, this roughly captures the nature of the process actually taking place in the brain.

their relative strengths and limitations are from a more practical perspective.

Functional Neuroimaging: fMRI, PET, EEG, MEG

The advent of practical techniques for imaging the activity of the living, breathing human brain has truly revolutionized the field of Psychology and Neuroscience. Initial pioneering work was done in the 1980's using the positron emission tomography (PET) scanner, which requires radioactive agents to be infused into the bloodstream. The PET scanner measures the decay of these radioactive labels, which can be formulated to bind to various different substances of interest in the brain, including different neurotransmitters such as dopamine, or glucose (sugar) to measure overall metabolic activity.

In 1992, several groups developed the ability to use magnetic resonance imaging (MRI) to measure the level of oxygen in the blood, known as the BOLD (blood-oxygen level dependent) signal, which varies as a function of overall neural activity within a given brain area. Interestingly, the brain over-reacts to neural activity, resulting in an over-supply of oxygen to the most active areas, rather than a depletion. This functional MRI (fMRI) technique has major advantages over PET, in not requiring an IV injection of radioactive tracers, and it has a much faster *temporal resolution* (i.e, the ability to resolve changes in activity over time). Furthermore, MRI machines are used in most clinical facilities of any reasonable size, so this technique made it possible for many scientists around the world to study how the brain responds to all manner of things inside the scanner.

In the ensuing years, fMRI techniques have improved to the point that remarkably small chunks of brain (called *voxels*, which are the volume analog of *pixels* in an image) about 1 mm on a side can be resolved, and in surprisingly many cases, these small voxels carry useful signals about what is going on in a given task. Current approaches typically focus on using the entire pattern of brain activity to understand how the brain works, which is consistent with our overall understanding about the way that many different neurons and brain areas work together to get the job done. Earlier, many scientists focused instead on identifying smallish blobs of activity that were particularly strongly activated by particular tasks (the *neo-phrenology* referred to earlier), but it has become evident that this only gives a small "porthole" view onto the full scope of brain function.

While fMRI can resolve relatively tiny voxels (i.e., it has good *spatial resolution*), its temporal resolution is still very limited (even though it is better than PET), because it is essentially measuring changes in blood flow, which take a while to react to changes in neural activity (about 6 seconds or so on average). A large number of different neural activity states can come and go within that 6 seconds, and all of these end up just getting blurred together in the overall fMRI signal.

To gain more insight into the detailed timecourse of cognition, there has been a continued and increasing

use of electroencephalography (EEG), which records real-time electrical signals using electrodes placed on the scalp, and has been around since the early 1900's. These signals immediately reflect changes in neural activity, providing excellent temporal resolution, but, alas, the remote recording of these signals from the scalp makes it very difficult to figure out exactly where the electrical signals are coming from within the brain. Thus, EEG has poor spatial resolution. Unfortunately, we do not yet have the perfect neuroimaging technique, which would have high resolution in both space and time. Nevertheless, advanced techniques in recording (using 100's of electrodes) and analysis have enabled EEG to achieve much better spatial resolution than before, and EEG can be combined with simultaneous fMRI recording to attempt to get the best of both worlds (though this remains challenging). You also may have heard about something called an *ERP* – this is just a way of averaging EEG signals together in a time-locked fashion, to create an *event related potential*, which has characteristic peaks and dips at different points in time, resulting from the waves of brain activation in response to a stimulus, or in preparation of a motor response.

Finally, there is a technique known as magnetoencephalography (MEG), which is the magnetic version of EEG. Recording these magnetic signals, which are much weaker than the electrical signals, requires advanced superconducting magnetometers, which in turn require complex cooling systems to get down to the superconducting realm. Thus, unlike EEG which is relatively inexpensive and portable, MEG is only available in a few labs around the world. However, it does have an advantage in spatial resolution over EEG, due to the way that the scalp distorts the electrical signal, but not the magnetic one.

Conclusions

Many scientists like to emphasize the popular sentiment that “the brain is a complete mystery” and we have barely scratched the surface in our understanding of it. However, you might get somewhat of a different impression from this chapter. In fact, we have a pretty good understanding of the large scale functional organization of the brain, which is consistent with all manner of data from neuroimaging and effects of brain damage, etc. As we'll see in the learning chapter, we have a remarkably good understanding of the details about how neurons learn at the synaptic level, and certainly we know a great deal about all the basic mechanisms underlying spiking. Detailed computer models incorporating all this data have been able to reproduce, at least at a coarse, approximate level, much of the actual human behavior observed in well-controlled laboratory studies, in domains such as perception, learning, memory, language, and cognitive control (O'Reilly et al. 2012).

Thus, while there certainly are many deep mysteries and major discoveries yet to be had, one could reasonably argue that we are at that stage in solving a jigsaw puzzle where a lot of the edges and key regions have already been filled in. Future editions of this textbook may not differ as much as you might think, as we start to fill in the rest of the picture. Only time will tell for sure, but there is at least room for optimism that we really are on the precipice of having a solid *science* of the brain and mind, and the goal of this textbook is to provide a coherent, comprehensive, and possibly a bit premature account of it.

Summary of Key Terms

This is a checklist of key terms / concepts that you should know about from this chapter. As we'll learn in the memory chapter, it is a great idea to test yourself on what was said about each of these terms, and then go back and double-check – that provides both beneficial repetition and also the *testing effect*.

- Neuron: synapse, dendrite, cell body, axon
 - neurotransmitter, reuptake, receptor, channel, ion, membrane potential, threshold, spike / action potential
 - 10,000 inputs to 1 output = *compression*
 - excitatory: glutamate, AMPA, Na⁺; inhibitory: GABA, Cl⁻
 - Tug-of-war creates *contrast* – neurons respond *relative to average*
- Brain:
 - Cerebral cortex / Neocortex:
 - * Occipital lobe: vision, V1
 - * Temporal lobe: objects, auditory, A1, language, episodes, semantics
 - * Parietal lobe: action, somatosensory, S1, homunculus, number, space, time, relations

- * Frontal lobe: motor, M1, control, planning, sequencing, decisions, motivation, affect
 - * Hierarchy of compression
- Thalamus: relay, thalamocortical system
- Basal ganglia: decision making, control, dopamine
- Amygdala: emotion
- Hippocampus: episodic memory
- Cerebellum: error-driven motor learning
- Hypothalamus: core body functions, HPA axis, stress
- Brainstem: reticular activating system, dopamine, serotonin, norepinephrine, acetylcholine
- Methods:
 - PET: radioactivity, slow (bad temporal resolution)
 - fMRI: blood oxygen (BOLD), faster than PET but still slow, good spatial resolution
 - EEG: electric signals, fast (real time, good temporal resolution), but poor spatial resolution
 - MEG: magnetic signals, fast, better spatial resolution, expensive

Chapter 3: Consciousness, Drugs, Sleep, and Dreams

Our entire subjective mental life is essentially synonymous with *consciousness* – hence its obvious fascination. As discussed in Chapter 1, psychology has wrestled with the challenge of subjective experience since its inception, and the stigma carried over from the behaviorist dogma against subjective experience still casts a shadow even to this day. When you start talking about dreams and the legacy of Sigmund Freud in talking about the mysterious workings of the subconscious, modern-day psychological scientists get even more uncomfortable!

As concluded in Chapter 1, we can view subjective and objective as two different perspectives, not two different magical substances, and while each individual has exclusive access to their own subjective perspective, we can all work together to build a consistent objective understanding through the scientific method. In this chapter, we build on the neuroscience foundation from Chapter 2, to understand more about the **neural correlates of consciousness** (*NCC*) and how we might somehow reconcile the subjective features of conscious experience with some objective facts about the brain.

Then, we delve into the intriguing world of altered states of consciousness induced by drugs and other factors, and then transition into the most universal state of altered consciousness that everyone can relate to: dreams! We'll see how dreams fit into the rest of what we know about the function of sleep and the different stages of sleep, and how the brain regulates activity between sleep and arousal.

Neural correlates of consciousness

As the chill cast over the scientific study of consciousness has thawed in the past few decades, a number of prominent theories of consciousness have emerged, most of which share a common core set of premises, even if the proponents may tend to emphasize their differences (this is how science works – you tend to gain a lot more attention by standing out than fitting in – *contrast* at work again). Before we get to these core ideas, we need to be clear about some terminology and ground rules.

First, as discussed in Chapter 1 (which you absolutely need to read before proceeding!), we can productively separate the **hard problem of consciousness** from the **easy problems** (Chalmers 1995). We do this by recognizing that the hard problem, associated with **qualia**, or the fundamental question of “what does it feel like?”, is inaccessible to objective science, because it is fundamentally *subjective* – objective science requires replicable data across many observers, but every subjective experience is an N of 1 (Nagel 1974). Nothing prevents us from speculating about the connections between subjective experience and objective understanding about brains, but it would essentially be impossible to *prove* anything at an objective, scientific level.

The best vehicles for conveying subjective experience are language and art, and as these are exclusively the province of humans, establishing which other types of brains might be conscious, and what their subjective experience might be like, is thus likely impossible. An important corollary of this is that our understanding of the term *consciousness* is inevitably shaped by our own subjective experience, and it must be acknowledged that there is no such thing as *generic* consciousness – we can only know about *human* consciousness, and, fundamentally, about our own singular subjective consciousness. Thus, the popular question, “which other animals, if any, are conscious?”, is really asking: “to what extent is the brain of another animal like that of the human, in the ways that might matter for shaping subjective experience?” While we can speculate about this, we will never really know for sure.

Even if we make an AI that seems “conscious” in every objective way, we can never inhabit its subjective world, and thus can never know for sure what its subjective conscious experience is like. However, if it starts writing heartfelt poetry, in a way that isn't just regurgitating statistical regularities in human writing (as is the case with the current generation of AI models, e.g., the otherwise impressive GPT-3 model (Brown et al. 2020)), we might be able to capture some glimpse of the subjective world of another type of being.

Features of Consciousness

So what is it *like* to be a human brain? Even though we all presumably have a strong feeling that we should be able to answer that question, actually doing so in a systematic, satisfying way has proven remarkably challenging. Early *introspectionists* like Wundt and his student Edward Titchener sought to enumerate and

categorize the subjective states of consciousness, but it was precisely the difficulty in doing so that motivated the behaviorists to completely reject the enterprise.

Modern treatments (e.g., Chalmer's list of the so-called "easy" problems) end up just listing what are essentially the main phenomenology studied by psychologists (and what we'll cover in the rest of this textbook):

- **Perception:** the ability to discriminate and categorize stimuli.
- **Attention:** the ability to focus processing on some subset, to the detriment of others.
- **Motor Behavior:** the ability to generate coordinated actions.
- **Reasoning:** the ability to integrate and organize knowledge.
- **Motivation & Control:** the goal-driven or *intentional* organization of behavior toward certain objectives.
- **Language:** the ability to verbally communicate or report on internal states, ideas, thoughts.
- **Self-awareness and metacognition:** the ability to somehow *access* internal states and be aware of being aware.
- **Wakefulness vs. sleep:** what essentially discriminates these states? Especially in the case of dreams, this is a widely recognized challenge!

While such lists do a good job of describing what human brains do, and are useful in the medical context, they don't quite seem to capture the essential role of consciousness *per se*. It seems like you could create a robotic system that could do all these things in a "zombie-like" manner that is essentially akin to sleep-walking: behaving and functioning in complex ways without anything that we would want to really identify as *consciousness* going on. Indeed, many such robots now exist, and nobody is ready to ascribe true consciousness to them.

One critical interim conclusion here is that *consciousness is not a single, monolithic construct* – it admits to many degrees or levels, and has many aspects or different facets. Likewise, existing evidence strongly suggests that no single brain region provides the unique source of consciousness – it cannot be simply reduced to the effects of one special magical entity (e.g., the pineal gland, as famously hypothesized by Descartes). Instead, consciousness is almost certainly an emergent phenomenon that depends critically upon neural interactions of some sort, but, just like the simple gears example from Chapter 2, it transcends this neural substrate and cannot be simply reduced to it.

In short, maybe a brain (or AI) that is just capable of basic perception, attention, and motor behavior has some minimal level of subjective conscious-like experience, but likely it is very different in overall character from our full-blown, super awesome human-level consciousness.

Recurrent Connectivity as a Major NCC

We can zero in on this key distinction between different *levels* of consciousness within your own subjective experience, to gain some insight into the critical properties of the more full-blown conscious states. Consider experiences when you were more "zombie-like" versus those where you were "very much aware" and "present" – what was the critical difference? For example, when driving down the freeway, thinking of other things in an "absent minded" manner, you may suddenly realize that you haven't been paying any attention to driving, and yet may have done fairly complex maneuvers such as changing lanes or even switching freeways etc.

This difference of being "absent" versus fully "present" seems to depend on the integrated combination of attention and awareness, with a *singular focus* of all the different parts of your brain working together. Somehow, each part of your brain, and whatever emergent entity we ascribe as the "self", is aware of what the other parts are doing, and all are organized around a common objective. This is often described as the **unitary** nature of consciousness.

It is this aspect of consciousness that many neuroscience-based approaches have focused upon, and where there seems to be an emerging consensus, with terms like *global workspace* (Baars 1988; Dehaene and Naccache 2001) and *integrated information* (Tononi 2004) used to describe this ability to integrate the disparate parts of the brain together. Perhaps the most neurally explicit theory along these lines is that **recurrent connectivity** in the brain is necessary for this integration aspect of consciousness (Lamme

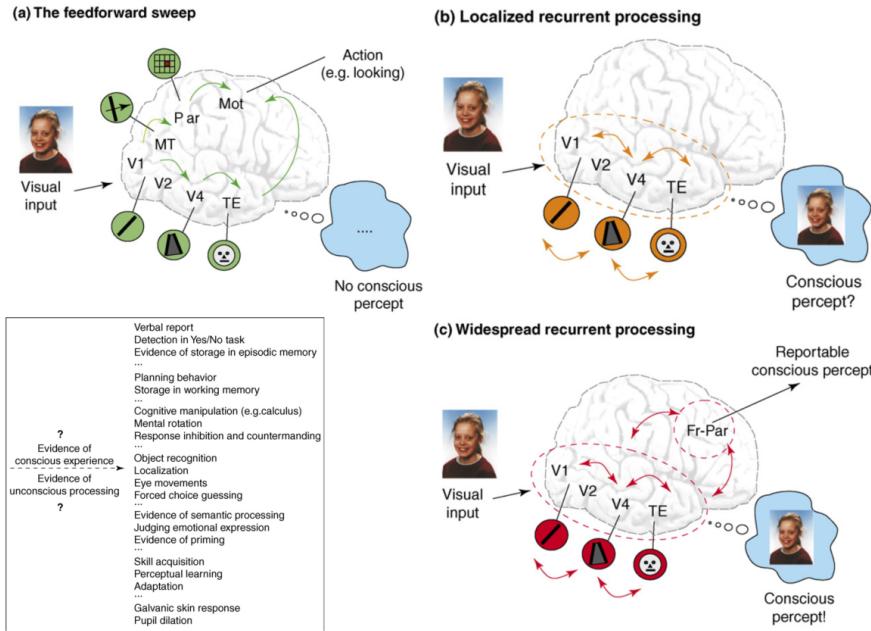


Figure 3.1: Consciousness is associated with robust recurrent processing across wide areas of the neocortex (Lamme, 2006). (a) With just the initial “feedforward sweep” of neural activity going from primary visual cortex (V1) on up, people generally do not report subjective, conscious awareness. (b) When more time has passed and bidirectional, recurrent connections can drive activity back down through these same visual areas, people are more likely to report conscious awareness. (c) With even more recurrent processing across more widespread areas of cortex, conscious awareness is strongly felt. Thus, the degree of recurrent processing provides a good measure of the graded strength of consciousness.

2006) (Figure 3.1). Intriguingly, this form of connectivity, where neurons in different areas have *bidirectional* connections (A excites B, and B also sends excitation back to A), is almost exclusively found in the neocortex. Perhaps not coincidentally, there is considerable evidence that we are also exclusively conscious of neural activity in the cortex – neural activity in subcortical areas such as the cerebellum or basal ganglia remains subconscious (as discussed in Chapter 2) (Koch et al. 2016).

The recurrent activity hypothesis also nicely accounts for **subliminal** stimuli, popularized by the [urban legend](#) of *subliminal advertising*, where briefly flashed messages in movie theaters are purported to have surprisingly large effects on purchases of popcorn and Coke. While the advertising part of this is not real (even blaringly *liminal* ads wish they could be so effective), many scientific studies have systematically varied the duration of stimulus presentation, and found that your ability to be consciously aware of the stimulus requires enough time to activate these recurrent circuits (Lamme 2006). And yes, there *are* some detectable effects of briefly flashed subliminal stimuli in these studies, e.g., being slightly faster at processing a second presentation of the same stimulus, known as *repetition priming*, but again nothing that would shape your overt behavior significantly.

In summary, it seems that our most fully conscious states depend on coordinating the firing of neurons across disparate areas of the neocortex around a common focus of attention, over a sufficiently long time period for the neurons to share their individual stories with each other. This account “resonates” with the idea that our subjective conscious feelings emerge directly out of neural firing in the brain, but also that they are not directly reducible to this neural firing – consciousness is perhaps the most emergent of all phenomena, requiring coordinated interactions among billions of neurons to generate enough raw emergent magic to give rise to this most amazing phenomenon in the known universe!

Animal and AI Consciousness

According to the recurrent activity account, we can look at the brains of other animals and determine the extent to which they have the recurrent connectivity necessary to support this form of consciousness, as a way of guessing whether they may share some of this same kind of subjective experience. Perhaps not surprisingly, our fellow primates are the closest on this score, and it would be difficult to argue that their subjective experience does not also include a similar form of subjective conscious experience. However, the fact that they lack language and all the amazing things that does for us and our cognitive abilities, will almost certainly make their subjective experience quite different than ours. Other mammals such as rodents have relatively little neocortex, and almost none of the critical prefrontal cortex that really broadly interconnects and coordinates cortical processing, so maybe they have a correspondingly diminished conscious awareness?

Interestingly, almost all modern neural-network-based AI avoids any significant amount of recurrent connectivity, and the very restricted form of “recurrent networks” that are used are not capable of the kind of broad coordination associated with human consciousness. This is consistent with the general impression that these models really lack any kind of “self awareness” or introspective capability – they are just really massive networks of pattern transformers that can imitate the behavior of our brains to some extent, but without that critical recurrent “light” that really makes us tick.



Figure 3.2: Illustration of the concept of *metacognition* and *self awareness* that is central to our subjective sense of consciousness, and can plausibly emerge from massive recurrent connectivity.

Figure 3.2 illustrates this key notion of “metacognition” or self awareness that plausibly depends on massive recurrent activation, so that our emergent conscious state can directly access neural activity across the brain. In later chapters we will revisit this notion of metacognition and its various important roles in human cognition (e.g., in knowing how confident you are about a memory you’ve retrieved), and we will also discuss the important functional benefits of being able to integrate information across many different brain areas, e.g., for solving challenging novel problems. Thus, consciousness is not just an interesting sideshow in the theatre of the mind – it is very likely the main stage upon which our most advanced cognitive abilities depend, and future research on its properties and functions may be critical for making the kinds of advances needed to finally capture human-level intelligence in artificial systems.

Altered States

Like our consideration of the different levels of consciousness experienced in normal waking, various forms of altered states of consciousness can provide important insights into the overall nature of consciousness. First, well outside the realm of normal experience, a number of studies of patients with various forms of brain damage have led to the conclusion that no particular part of the brain is essential for consciousness, except that at least some part of the neocortex must be intact. Of course, certain brainstem areas are critical for life and for the neocortex to even function properly in the first place, but clearly the neocortex is essential. However, no specific part of the cortex itself is essential, although higher-level association areas are likely differentially important (Koch et al. 2016).

A more widely accessible source of altered states can be had by consuming various psychoactive substances, and the challenges and potential for sharing subjective conscious states across the interpersonal void can be experienced by reading literature written by those under the influence of such substances, such as “The Doors of Perception” by Aldous Huxley (namesake of the band *The Doors*), or watching psychedelic movies (e.g., the classic *Altered States*, which gives a good impression of psychologists in the 1970’s :) In general, these altered states can involve changes in the whole gamut of psychological functions, from perceptual hallucinations to heightened emotional states, and disordered cognition. These effects contrast with the less qualitatively dramatic consequences of drugs like benzodiazepines, which have more “linear” effects on arousal and sedation.

Some research has attempted to connect changes in neurochemistry and receptor function to the nature of the subjective psychoactive experiences, including changes in effective inhibitory and excitatory connectivity that can give rise to characteristic features of visual hallucinations (Bressloff et al. 2002). Thus, again, our subjective conscious states can be understood in terms of direct effects on the underlying neural substrates of the brain, but understanding these effects from the objective outside cannot directly tell us what it feels like subjectively.

Neuromodulators and Drugs

Given their sociological, psychological and medical importance, we devote some time here to understanding the way that various well-known drugs affect the brain. As discussed in Chapter 2, there are specific nuclei in the brain stem *reticular activating system* that release *neuromodulators* that have broad overall effects on wide areas of the brain. Perhaps not surprisingly, these are the major targets of psychoactive drugs, because they have such significant modulatory effects on the brain. By contrast, *glutamate*, which is the main *neurotransmitter* in the strict sense of transmitting detailed signals from one neuron to the next (in the neocortex and many other areas), is not directly affected by most drugs, and its effects are much more local and content-specific. GABA, the primary inhibitory neurotransmitter, can be considered more of a neuromodulator in that it has broader regulatory effects and is directly affected by psychoactive drugs, as we mentioned earlier. To be clear, this difference between neurotransmitter and neuromodulator is strictly a functional distinction – they are all just chemicals released by the axons of neurons, but it is useful to distinguish the transmission vs. modulation roles.

Drugs can affect the brain in two opposing ways:

- **Agonists** are drugs that mimic or amplify the effect of a given neuromodulator. This term is a bit “agonizing” because it doesn’t exactly sound like what it means, but you can perhaps remember it better in relation to its opposite (antagonist). Scientists typically reserve the term *agonist* to more precisely refer to chemicals that specifically bind to the same receptors as the endogenous neuromodulator, but we’ll adopt a looser definition that includes anything that has a net “positive” effect on the effect of the neuromodulator. For example, *Valium* and other *benzodiazepines* are direct GABA agonists by binding to the GABA receptor and enhancing the amount of Cl^- that enters the cell, whereas *Ritalin* (*Methylphenidate*) enhances dopamine effects by inhibiting the reuptake of dopamine after it is released, so it is a kind of agonist but acts more indirectly. There are many different biochemical mechanisms that can lead to a net agonist effect.
- **Antagonists** are drugs that suppress, inhibit, or otherwise work against a given neuromodulator. They “antagonize” that poor neuromodulator. *Curare* poison is a classic competitive antagonist for acetylcholine (ACh) at the synapses of nerve fibers onto muscles, thus acting to paralyze muscles. It acts

by binding directly to the same receptors that ACh normally binds to (and it does so more effectively, i.e., with greater *affinity*), but it does *not* actually open those receptor channels. *Botulinum toxin (botox)* is also an overall ACh antagonist, but it works by preventing the release of ACh.

Interestingly, the neuromodulators are biologically ancient chemicals that have very different effects throughout the body, which explains why drugs often have many side-effects. For example, ACh drives the most basic function of muscle contraction throughout the body, but in the brain it is one of those high-level control knobs affecting attention, arousal, and learning. Dopamine receptors are also involved in lactation. Evolution is very pragmatic in re-purposing existing technology. Furthermore, the major players of serotonin, dopamine, and norepinephrine are all chemically very similar *monoamines*, so many drugs affect all of them to varying extents. Thus, overall, understanding the full effects of any given drug can be very complicated.

- **Caffeine** is a direct antagonist for *adenosine* receptors, which in turn are antagonistic against dopamine, and overall lead to sedation (drowsiness). Thus, consistent with its widely known and appreciated subjective effects, it directly inhibits drowsiness, and also leads to a net increase in dopamine, producing pleasurable effects and leading to its addictive properties.
- **Nicotine** is an agonist for a type of ACh receptor (the *nicotinic* ACh receptor) that drives the attention and arousal effects of ACh in the cortex. This is consistent with the stereotypical chain-smoking author or detective using nicotine to enhance their ability to focus and concentrate.
- **Alcohol** (ethanol) has complex effects on neurons, that vary with dose and over time. It acts as a GABA agonist, increasing levels of inhibition, which accounts for its psychological effects in reducing anxiety, causing sedation, and reducing “behavioral inhibition” as discussed earlier. It also antagonizes the binding of glutamate to the NMDA receptor, which is involved in learning as well see in the Learning chapter. Both the GABA and NMDA effects combine to impair learning in the hippocampus, leading to memory blackouts.
- **Benzodiazepines** (Valium, Xanax, Midazolam, etc) are widely-used GABA agonists, which, like alcohol, reduce anxiety, cause sedation, and generally turn off the brain to varying extents. If you’ve ever had surgery, you’ve likely had Midazolam, which knocks you out and prevents you from remembering anything. In low doses, Midazolam has been used in scientific studies to produce a reversible hippocampal amnesia-like condition, due to the heightened sensitivity of the hippocampus to the effects of GABA.
- **Amphetamine** (speed, Adderall) is an agonist for both *norepinephrine (NE)* and dopamine, increasing release and actually reversing the reuptake process so that there is more of these neuromodulators in the synapse. Both of these neuromodulators affect attention and learning, consistent with the observed behavioral and cognitive effects. Adderall is used for treating people with ADHD, which is somewhat paradoxical given the “hyperactive” component of this disorder. However, it is likely that NE acts to keep people actively engaged for a longer time, “locking in” a given set of frontal control signals and preventing the characteristic distractability of ADHD (Aston-Jones and Cohen 2005).
- **Cocaine** is similar overall to amphetamine in both biochemical and psychological effects. It has a specific inhibitory effect on *dopamine transporter (DAT)* that is responsible for reuptake of dopamine, thus producing an overall agonist effect on dopamine (leaving more in the synapse). These direct effects on dopamine likely play a critical role in its addictive properties, as it simulates the effects of rewarding outcomes, in a way that circumvents the natural *contrast* mechanisms that discount rewards in proportion to expectations (Redish 2004).
- **SSRIs** (Prozac et al) affect *serotonin* function by inhibiting the process of reuptake of the neurotransmitter after it has been released (i.e., *serotonin-specific reuptake inhibitors*). This allows serotonin to linger longer, and potentially have a larger overall effect. However serotonin is so incredibly complex at multiple levels, that nobody really understands exactly what is going on, and we really can’t be sure if it is an agonist or an antagonist. For example, serotonin (and all the other neuromodulators) have negative feedback mechanisms that strongly regulate the amount released, and it is possible that blocking reuptake causes these feedback mechanisms to over-react, thus leading to a net reduction in serotonin release over time. Furthermore, different serotonin sub-nuclei within the raphe have contradictory effects, with some promoting positive emotional states and others having the opposite effects.
- **Psychedelics** (LSD, psilocybin, peyote, etc) all have primary effects on the serotonin system, which,

among its many talents, is important for regulating sleep. The simplest explanation for the effects of these substances is that they effectively produce a waking dream state – dreams have similar characteristics to psychedelic experiences in terms of heightened emotional states and disorganized cognition driven by loose associative connections instead of the more controlled, willful, goal-directed states that characterize waking cognition. The reasons for these effects can be directly connected to the effects of serotonin on different brain areas, as we discuss in the next section.

- **Cannabis** (Marijuana) is a unique case where the drug activates receptors and associated endogenous neurotransmitter systems that were previously unknown, and have only relatively recently been discovered as a direct result of studying the effects of the drug. Thus, the receptors and endogenous neurotransmitters are known as *cannabinoid* receptors and *endocannabinoids*, and now that we have the tools to identify these things, they turn out to be found all over the body, like all the other more well-known neuromodulatory systems. However, unlike these other systems, the reason we never knew of these cannabinoid systems before is that they don't have a central nucleus that releases them – instead they are produced locally in cell membranes, and have a very localized signaling role, by sending messages *backward* across the synapse (i.e., from dendrite back to axon, instead of the usual other way around). The detailed function of these systems is still relatively unknown, and represents an exciting frontier in current research, well-timed with the recent legalization of this interesting substance in a number of different US states and other countries.
- **Narcotics** (heroin, morphine, fentanyl, opiates) are agonists for the endogenous opioid system, involved in regulating the neural response to pain stimuli. Opioid receptors are found in the amygdala, basal ganglia, hypothalamus, and thalamus, and this explains the strong emotional, euphoric effects of these drugs. These substances are widely believed to be the most addictive of all drugs.

In summary, these drugs are jacking right into those global control knobs in the brain, and provide a critical window into understanding how our brains function normally: your endogenous states of arousal, excitement, sedation, etc are all controlling these very same knobs. Some of these psychoactive drugs, such as caffeine, alcohol, and nicotine are very widely used (and abused), and there have been increasingly urgent discussions about the ethics of performance-enhancing drug use in schools, which is on the rise. To what extent is this like doping in sports, or should we instead consider it more like using a calculator or a computer: something that augments our native biological abilities to the general betterment of society, etc? What is your opinion?

Sleep and Dreams

Even beyond the mystery of dreams, sleep itself is a fascinating topic that raises so many questions, and also presents critical everyday challenges for many people. A large number (35-40%) of people in the USA do not get enough sleep (Liu et al. 2016), and yet research increasingly shows how important sleep is to health and longevity. But I worry that this very research, and its popular dissemination (e.g., Matthew Walker's popular [TED talk](#)), may actually be having a deleterious effect, by making people even more anxious about sleep, and thus leading to more insomnia! Sleep is one of those few ancient, biological things that our conscious mind cannot directly control, and thus it presents a frustrating paradox: the more you *try* to fall asleep, the less capable you are of actually doing so. Nevertheless, the positive side of all the current sleep media is the sharing of many *indirect* ways of enhancing sleep, such as exercise, diet, and meditation / mindfulness strategies that really do work.

One of the most effective sleep strategies in my family was due to an article I read many years ago, suggesting that **neuropeptide S** (*NPS*, a more localised, specialized neuromodulator in the brainstem) may be responsible for the circadian fluctuations in how brave vs. afraid different animals are. We noticed how our cats were such “fraidy cats” during the day, but always wanting to go out and explore at night, and speculated that this must be due to neuropeptide S, with cats being nocturnal. And then we told our young boys about this, and how when they were trying to fall asleep but experiencing all these anxious thoughts and feelings, they were just like the cats during the day – when they wake up in the morning, all those nighttime fears will melt away in a flood of NPS! Just by knowing the neural basis, and circadian timing, of these fearful thoughts, it provides a significant measure of reassurance and mindfulness-like ability to put those thoughts in their proper perspective. It is just your brain trying to make sure you don't go out at night

and do something foolish like getting eaten by a saber-tooth tiger! And it looks like this NPS story is holding up, with it playing important roles in both sleep regulation and anxiety (in the amygdala, seat of both fear and positive emotions) (Chauveau et al. 2020; Tillmann et al. 2019; Jüngling et al. 2008).

In addition to **insomnia**, which most people experience at least a few times, and can be debilitatingly severe, there are a number of other more rare sleep disorders, such as **narcolepsy** (sudden onset of sleep during the day – famously prevalent in dogs), persistent **sleepwalking** or **night terrors**, and **REM sleep behavior disorder** which involves acting out violent dreams physically while remaining asleep. These sleep disorders are recognized in a special category of the **DSM-5** (*Diagnostic and Statistical Manual of Mental Disorders*, 5th edition, of the American Psychiatric Association), when persistent and impairing quality of life (see Chapter 12 for more discussion).

Functions of Sleep

Perhaps the most persistent and puzzling question about sleep is why organisms do it in the first place? It seems like such a colossal waste of time! As usual, everyone wants to find *the* single reason for sleep, but there are clearly many different essential functions of sleep, including:

- Recharging metabolic batteries and building proteins and other molecules that form the body's infrastructure. These are the *anabolic* pathways, where molecules are built up, as compared to the subsequent *catabolic* breakdown of these molecules, e.g., when bursts of energy are required.
- Repair and learning in the immune system – reductions in sleep are associated with impaired immune function.
- Protein synthesis and other anabolic processes in neurons, which can “solidify” or *consolidate* information learned during the day. Interestingly, we'll see in Chapter 5 that learning in the synapses connecting neurons (which is where knowledge is encoded in the brain), requires the buildup of *actin* proteins, which are also a critical ingredient in muscles. Thus, your brain is very much a muscle (this is not just a popular metaphor), and sleep helps stabilize its critical synaptic connections. There is a large and growing body of literature showing that memory can significantly improve after a good night's sleep (Klinzing, Niethard, and Born 2019), as we'll see in Chapter 6 on memory!

The importance of sleep has been starkly demonstrated in experiments where some rats kept continuously awake eventually died (after 5, 13, and 33 days), while all other sleep deprived rats suffered obvious effects of lack of the above processes, looking very sickly and shriveled (Rechtschaffen et al. 1983). However, contrary to popular mythology, people do not die from sleep deprivation, and the world record is somewhere around 11 or 18 days depending on the standard of authenticity. Furthermore, sleep is found in most animals, indicating its importance.

Sleep Stages

In addition to there being many benefits from sleep, sleep itself is a multifaceted process consisting of 4 different stages, each with deeper, slower brain waves (Figure 3.3). These stages unfold over a roughly 90 minute cycle progressively through the night. It is easiest to wake during the higher stages, and much more difficult and disorienting to wake during the deepest, slow-wave stages. It is during these deep slow-wave stages when most of the above repair and rebuilding steps take place, and so these are “prioritized” at the start of sleep, with progressively longer REM phases through the later sleep cycles.

Doing the math, a 5-cycle night is about 7.5 hours, and it may take roughly 30 min to actually fall asleep, so that accounts for the typical 8 hours of recommended sleep for adults. You can get a device that synchronizes your alarm to your brain waves or other sleep signals such as motion sensors, to avoid waking you in deep sleep, producing a better wakeup feeling. You can also just set your alarm in 90 minute multiples relative to when you expect to actually fall asleep (if you happen to know that!). Entraining your body to regular sleep cycles by going to bed at roughly the same time every night is also a good way to get solid sleep, and if you plan that to work with the 90 minute cycles, you should feel consistently well rested. Isn't that better than doing that one more thing before bed!?

The progressive slowing of brain waves across stages of sleep was nicely leveraged in the must-see Christopher Nolan movie *Inception*, with the idea that the slower brain waves translate into a compression of

Cycles of Sleep

DREAMITALL

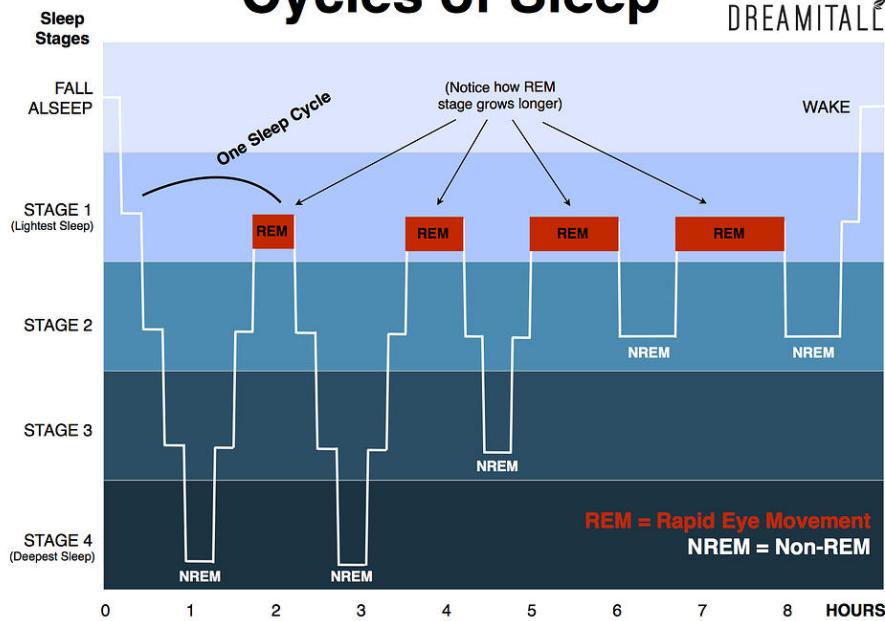


Figure 3.3: The four stages of sleep and their timecourse over a typical night, unfolding over roughly 90 minute cycles. Deeper stages are characterized by slower brain waves (slow wave sleep) and greater difficulty waking. Bain waves during REM (rapid eye movement) are similar to waking. Progressively less deep sleep and progressively more REM sleep occur over the course of a night.

time. This is an oft-reported feature of dreams, that many events seem to be compressed into a relatively short period of “waking time”. Technically speaking, you would actually need *faster* internal brain waves to explain this phenomenon, but don’t let that detract from your enjoyment of the movie! A more plausible explanation for these time dilation effects, suggested by [Nolan himself](#) is that dreams tend to jump around quickly, so it may seem that much has happened, without having all the intervening time actually passing (just like in movies!).

The brain mechanisms driving sleep and arousal involve key roles for the **hypothalamus** down low in the brainstem (which is also responsible for most other basic bodily functions) and the **thalamus** as it interconnects with the cortex. In effect, the thalamus shuts off the flow of sensory inputs, and motor outputs to / from the cortex, and it is also important in driving the slower rhythms associated with sleep stages. This cutting off of input / output connections is important for preventing sensory inputs from overly activating, and thus awaking, the cortex, and for preventing all those crazy dream states from actually being acted out (unless you suffer from REM sleep behavior disorder or sleepwalking, which are associated with incomplete blockage by the thalamus).

While there are 4 distinct sleep stages, most cognitive neuroscience research focuses on the differences between **slow-wave-sleep (SWS)** (stages 3 & 4) versus **rapid-eye-movement (REM)** sleep, which are where most of the time is spent. People woken from SWS report dreaming about 50% of the time, while in REM that figure is about 80%, so dreaming happens in both of these stages, contrary to popular belief that REM is the exclusive province of dreaming. However, dreams in REM sleep tend to be more vivid, bizarre and “dreamlike” (Hobson and Pace-Schott 2002), so REM is still the primary dreamtime.

In terms of the memory and learning effects of sleep, various sources of data suggest that, during SWS, your **hippocampus** is effectively teaching the neocortex about things it has learned, while REM sleep may drive more localized, cortically-specific synaptic changes (while your subjective self is off somewhere free-associating) (Hobson and Pace-Schott 2002; Diekelmann and Born 2010; Klinzing, Niethard, and Born 2019). This idea that the hippocampus trains the cortex makes sense according to computational models of learning (McClelland, McNaughton, and O'Reilly 1995), and it may depend in part on sharp-wave ripple

events taking place during SWS that originate in the hippocampus and propagate outward from there (Buzsáki 1989).

The Function of Dreams

Another ubiquitous question about sleep and dreams is whether dreams have any specific function, outside of the above-listed functions of sleep more generally? You have probably heard about Sigmund Freud's famous obsession with dreams as a way of working through unresolved subconscious tensions (typically of a sexual nature), but there is not much direct modern scientific evidence about such hypotheses. One attractive modern idea is that dreams are just the reflection of the memory replay effects described above (Wamsley 2014). Interestingly, the jumbled, random, bizarre nature of dreams may actually be beneficial in this memory consolidation process, to transform specific episodic memories into more flexible semantic knowledge where all the different parts of an episode are broken down and made available as separable mix-and-match elements (McClelland, McNaughton, and O'Reilly 1995).

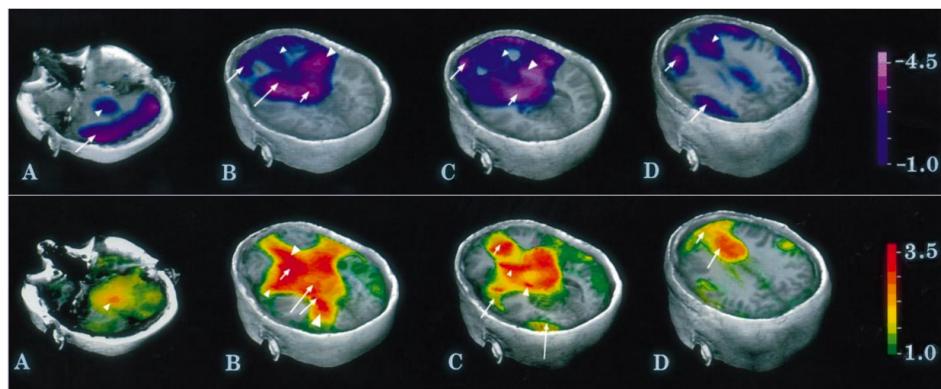


Figure 3.4: The frontal cortex goes deeply to sleep when you're dreaming, while the amygdala and hippocampus are highly active. This explains why your dreams are so incoherent and disorganized, and it is impossible to ever catch an airplane or show up for a test on time, etc. Also, dreams are emotionally charged and incorporate recent memories, thanks to the amygdala and hippocampus.

One of the most important neuroscience findings about brain activity during dream states is shown in Figure 3.4, demonstrating that your *prefrontal cortex* is the cortical brain area most deactivated during REM sleep (Hobson and Pace-Schott 2002). Thus, every night, we all get the benefit of experiencing a full frontal lobotomy! This then makes a lot of sense in terms of the jumbled bizarre nature of dreams: without the frontal executive task-master keeping the rest of the brain focused on a single current objective, the posterior cortex is left to its own devices, freely associating through its dense webs of interconnected associative networks of knowledge.

As is often thought to be the case with “artistic types” who seem to be more like this dream state even when awake this freedom from executive control is probably important for enabling creativity to flow. Indeed one of the most salient functions ascribed to dreams is exactly this ability to come up with novel solutions to challenging problems, as in the famous case of the chemist Kekule who discovered the ring shape of benzene from a daydream of a snake biting its own tail.

Another feature of the dreaming brain state is that the amygdala is typically more active (e.g., due to the deactivation of neuropeptide S), which would make sense of the increased level and intensity of emotional states typically experienced during dreaming. Furthermore, as we discussed earlier, psychedelic drugs act by driving some of these same brain changes seen during dreaming, and they also feature increased amygdala activation, and frontal deactivation, producing similar increases in emotional significance from everyday things and events, and an overall disorganization and inability to maintain coherent trains of thought.

Finally, we complete our own circle of ideas here and return to the question of consciousness in the context of the dream state. What can the state of consciousness during dreaming tell us about our normal waking conscious state, and about the nature of consciousness more generally? First, we can see that the

subjective properties typical of dream states are well correlated with known changes in brain activity during the different sleep stages, further establishing this direct connection between subjective, conscious states and the underlying neural substrate. However, despite some promising initial reports, it is not clear if more detailed electrical signatures of neural activity recorded in the scalp EEG can predict when a person has been dreaming or not during SWS (Wong et al. 2020). Thus, the relationship between brain states and dreaming may be more emergent and complex than can be effectively measured using these EEG signals.

So how do we know whether “life is but a dream” anyway? Well, if you try to use the scientific method while dreaming, and conduct replicable scientific experiments, you’ll likely fail (spectacularly). If you’re like me, you’ll just completely forget about the experiments after a bit, and wonder why those people from high school showed up so unexpectedly. I’m always failing to even catch a flight, which I lapse into and out of remembering that I’m late for. Thus, the fact that, during the waking state, we *can* successfully do science, seems like at least a pretty sensible indicator that these two states differ in important ways.

Whether the waking state is all just within some crazy big simulation remains a further, probably unanswerable question. But having created many virtual worlds for my computer models of the brain to inhabit, creating a simulation that can withstand the extreme distances and scales penetrated by modern astronomy and physics experiments would be challenging to the point of extreme implausibility, in my opinion. Maybe you and Elon Musk have a different opinion, but anyway, that’s why subjectivity is primary.

Summary of Key Terms

This is a checklist of key terms / concepts that you should know about from this chapter. As we’ll learn in the memory chapter, it is a great idea to test yourself on what was said about each of these terms, and then go back and double-check – that provides both beneficial repetition and also the *testing effect*.

- Neural correlates of consciousness
 - hard problem vs. easy problems
 - unitary nature of consciousness
 - recurrent connectivity / processing
 - subliminal processing
- Drugs:
 - Agonist: activates, enhances neurotransmitter / receptor function
 - Antagonist: inhibits, suppresses neurotransmitter / receptor function
 - Caffeine: adenosine
 - Nicotine: ACh
 - Alcohol: GABA
 - Benzodiazepines: GABA
 - Amphetamine: norepinephrine
 - Cocaine: dopamine
 - SSRI: serotonin
 - Psychedelics: serotonin
 - Cannabis: cannabinoid
 - Narcotics: endogenous opioid
- Sleep & Dreams
 - neuropeptide S
 - functions: recharging, rebuilding, immune system, learning consolidation
 - disorders: insomnia, narcolepsy, sleepwalking, night terrors, REM sleep disorder
 - stages: slow-wave-sleep vs. rapid-eye-movement (REM)
 - hypothalamus, thalamus role

Chapter 4: Sensation, Perception, and Attention

Channeling my hard-boiled teenage son, he would say: "Dad, why are you writing a whole !@#\$%ing chapter about *seeing* – you just look and you *see sht* – what's the big deal!" This is, actually, quite accurate. I mean, it is what he would say (this has been confirmed), and it is also what makes this chapter difficult to write: we all just take perception for granted, because from our subjective perspective, it does just happen, preconsciously, and we are only aware of the results. From a compression* standpoint, that's about all you need to know, right!? It just works, so get on with it!

Nevertheless, despite the potential futility of the exercise, I will persist in trying to convince you that perception is amazing and fascinating, and give you some sense of how it works, and why roughly 50% of your massive neocortex is required to solve this “trivial” problem. This tendency to underestimate the complexity of perception has been around for a long time: there is a famous story about how, at the dawn of AI research in the 1950’s, a random graduate student was tasked with solving the vision problem over the summer, so they could plug it into the rest of the system next year. Needless to say, it didn’t happen, and in fact it is only in the last 5 years or so that AI systems finally have semi-functional perceptual front-ends. You have likely experienced this in the speech recognition domain, when talking to Siri or other similar digital assistants: they often work but still make some basic mistakes, and they *definitely* don’t seem to really understand what you are saying at any deeper semantic level, but that’s another issue.

Perception is (Hierarchical) Compression

The trick to getting these AI systems to finally work was to adopt the strategy that the brain uses, by employing large networks of simulated neuron-like processing elements, organized over many hierarchical layers (i.e., “deep” neural networks), and trained by a learning mechanism known as *error backpropagation*, which was developed by psychologists in the 1980’s to better understand many properties of human cognition and learning (Rumelhart, Hinton, and Williams 1986; McClelland and Rumelhart 1986). A few important computational tricks made these networks work better and faster (LeCun et al. 1990), and the advent of fast computer chips developed for video gamers enabled these networks to be dramatically scaled up in size, resulting in significant performance improvements (Krizhevsky, Sutskever, and Hinton 2012; LeCun, Bengio, and Hinton 2015).

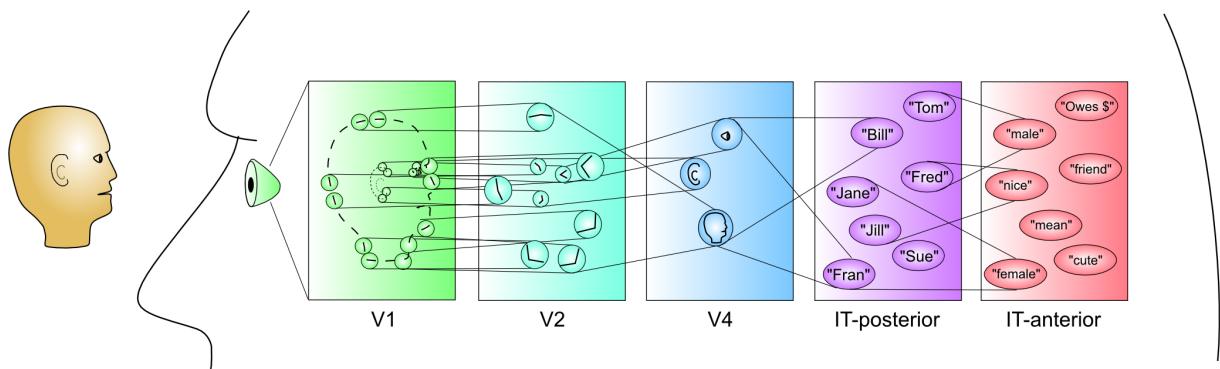


Figure 4.1: Hierarchical organization of detectors in the visual pathway going into the temporal lobe, supporting the ability to recognize (detect) entire objects, based on earlier levels detecting parts and features of parts. This shows the large-scale, cumulative effects of *compression* from very high-dimensional raw sensory inputs, to high-level, succinct interpretations of the world. Although a highly simplified cartoon, this roughly captures the nature of the process actually taking place in the brain.

The essential strategy learned by these deep neural networks, and the brain, is shown in Figure 4.1 (we already saw this figure in Chapter 2), where each layer **compresses** the complexity of the patterns on the layer before: *getting rid of irrelevant differences, while extracting the important ones that the system actually cares about*. This is the essential function of perceptual systems, served by the 10,000-to-1 compression

property of individual neurons that are detecting relevant patterns and ignoring irrelevant ones.

To use another metaphor, you can think of perception as a *filter*, filtering out irrelevant “junk” from the perceptual input signal, and purifying the relevant, important stuff. It takes multiple layers of such filters because one step of filtration can only do so much purification, and each such layer builds on the partially-purified outputs of the layer before. This is why it takes so many neurons in the brain, and in the AI models, to do a good job at perception – each individual neuron can only do a small part of the overall job.

We See the “Real” World, not Raw Sensation

Our subjective, conscious experience is dominated by the higher levels of this hierarchical perceptual filtering system, because that is what most strongly and directly interconnects with all the other brain areas, and this process of bidirectional communication and influence across brain areas is what drives consciousness as we discussed in the previous chapter. These higher layers are called **association cortex** because they “associate” with all these other brain areas.



Figure 4.2: The doors of perception, as represented by the author in a computer painting from college, depicting the hidden elemental nature of visual input and other potentially mysterious elements.

This is the reason for my son’s hard-boiled attitude about perception: by the time we’re aware of it, perception has already done all the hard filtering and compression work for us, and we just experience this nice simple impression of what is out there in the world! Interestingly, we nevertheless seem to retain a bit of a sneaking feeling that our perceptual systems might be hiding something interesting from us, e.g., in the popular notion of the “doors of perception” being blown wide open by psychedelic drugs and other such experiences (e.g., Figure 4.2). Somehow, we feel like we want to be able to see through all those filters, and see the world as it “truly is”.

However, the truth (represented by the mirror on the white table in Figure 4.2) is that our perceptual systems do a really amazing job of delivering a highly accurate representation of the world – there is no greater truth than what your eyes deliver to you, tirelessly, every moment. So look in the mirror and behold the truth! And, if you really want to get a different perspective, try a magnifying glass or, better yet, a microscope. That will reveal the separate Red, Green, Blue (RGB) dots in your cell phone screen, and give you a hint of all the cells in your skin, etc. That is a perspective that might blow some doors open, if you really think about it!

In perceptual science, the fact that our perception is of the world, not the raw sensory signals that our perceptual systems receive from our sense organs (more on those later), is referred to as **perceptual constancy** in general, with more specific versions such as **color constancy**, **size constancy**, etc. Figure 4.3 (which we saw already in Chapter 0) provides a powerful demonstration of color constancy, where the exact same raw pixel RGB values are perceived as strikingly different shades, based on our ability to integrate the various elements of the scene into a coherent overall interpretation, including the effects of lighting and shadows, and the regularity of the checkerboard pattern, etc.

Perhaps the most striking demonstration of the divergence between raw sensory input and subjective perception is from “The Dress”, which was a viral internet sensation in 2015 (Figure 4.4), because people

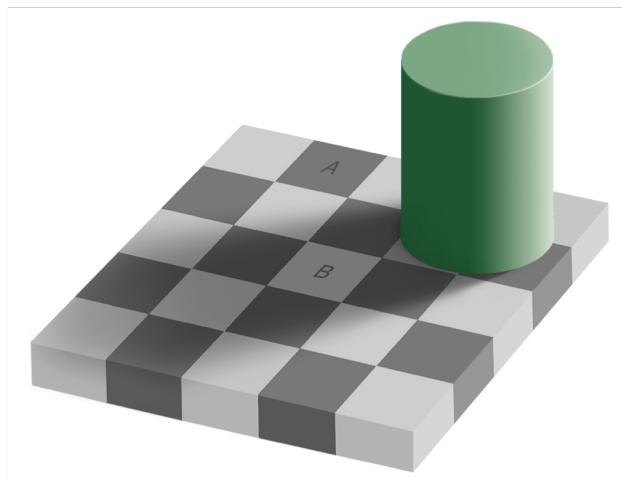


Figure 4.3: Illustration of color constancy, and more generally how we see the world, not the raw sensory signals. The raw pixels in A and B are identical, and yet we see them as strikingly different shades, based on all the “contextual” information about shadow and the regular checkerboard pattern, etc.

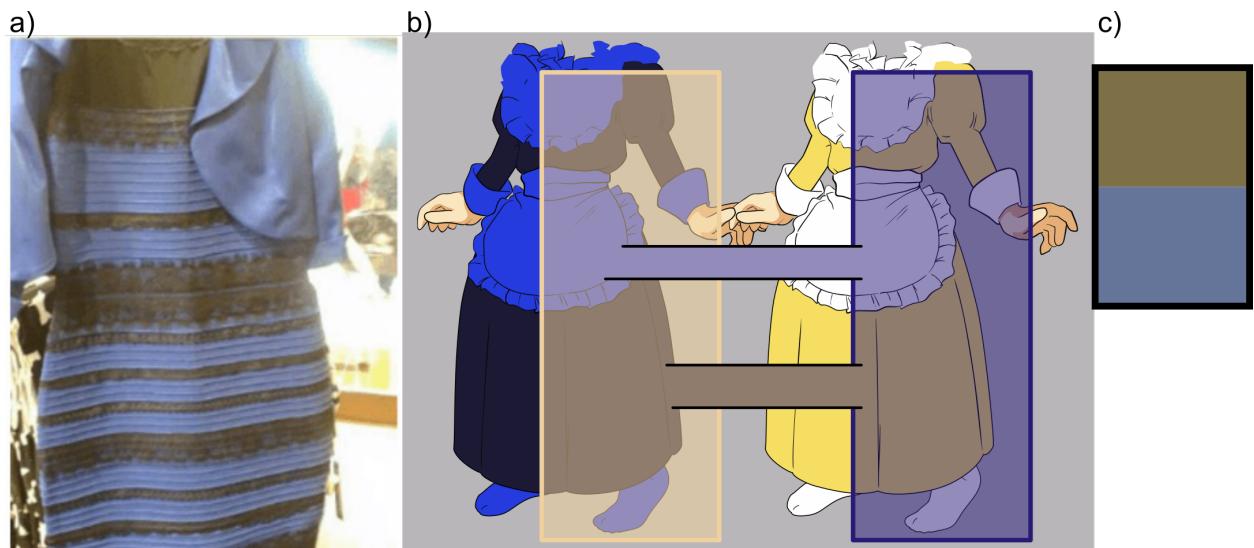


Figure 4.4: a) “The Dress”, which sparked a viral internet sensation in February, 2015, because different people have strongly divergent perceptions of the dress colors. Some see it as black and royal blue, and others as white and gold. b) Shows an explanation for the two different interpretations, with a gold vs. blue filter over the two different colors of underlying dress. Similar to Figure 4.3, it is hard to convince yourself that the colors within these filtered boxes are actually identical across the left and right boxes, but they are. (Figure design by Kasuga-jawiki; vectorization by Editor at Large; “The dress” modification by Jahobr, CC BY-SA 3.0, [wikimedia.org](https://commons.wikimedia.org)) c) The “objective” RGB colors sampled from the original image, which in isolation are clearly gold and blueish.

experienced very divergent yet strongly felt percepts of the dress colors. From a raw RGB perspective (Figure 4.4c), the dress colors are gold and blue, but studies show that 57% of people see it as black and blue. Figure 4.4b shows how two different color filters (gold on the left, blue on the right) produce identical stimulus-level color values for the two different dress colors that people report seeing.

Thus, the visual system is just making different **assumptions** about the lighting conditions across individuals, and, critically, your subjective percept is irrevocably “colored” by those assumptions, in an attempt to tell you what the *real* underlying color of the dress is, independent of the specific lighting conditions. This is the way in which your perception is more real than your raw sensation – most of the time it is amazingly accurate in telling you what the *real* materials (and shapes, sizes, etc) of things are in the world.

The fact that our perceptual systems need to make these assumptions sets them up for the inevitable *ass-u-me* situation (i.e., making an ass out of u and me): **illusions** reveal the nature of these assumptions, and are thus a fun and informative way to understand how the perceptual system works, under the hood. Some lucky researchers basically spend their entire day just coming up with new illusions (although this is increasingly less viable of a career path these days, as the chances of finding truly new such illusions is dwindling). Anyway, this is certainly the approach we’re going to take for the rest of the chapter, so buckle up and get ready to see some crazy stuff!

Before we go there, however, it is important to point out that the current generation of AI models discussed above almost certainly do *not* “see” the world in the way we do. There is nothing in the way that these models are trained that would cause them to make the same kinds of assumptions necessary for the color constancy demonstrated in the above figures. Furthermore, they are typically strictly *feedforward* in their processing of the sensory input – raw image pixels go straight up the hierarchy, without any higher-level interpretations of the scene coming back down to influence the way these lower levels process things. By contrast, the mammalian visual system is massively bidirectionally connected, and the *top-down* connections from higher levels to lower levels is critical for enabling the overall scene-level information to so strongly affect our basic perception of the elements of the scene, as so well demonstrated in Figure 4.3.

Extensive neural evidence shows how these top-down processes shape the firing of neurons throughout the visual system, all the way down in primary visual cortex (V1) (Lee et al. 2002; Angelucci and Bressloff 2006). And there are some neural network models with bidirectional connectivity, which demonstrate simple versions of these kinds of top-down phenomena (O’Reilly et al. 2013), but considerable more work needs to be done to capture something like Figure 4.3. In terms of the critical question of what kind of learning signal might cause the model to encode the stable features of the world instead of the raw sensory features, one idea is that the brain may learn by trying to predict what will happen next, and the stable, predictable things in the world are therefore what is learned (more on this in the learning chapter). By contrast, existing AI models are typically trained to label the category of objects in a scene, using massive human-labeled datasets – in effect, they are just learning to imitate one small part of the human perceptual filtering process.

Sensory Systems

Now that we have a sense of some of the big-picture issues and challenges in perception, we will dig into some of the more specific details about different sensory systems, so we can understand how it all works. As you undoubtedly already know, we have 5 major sensory modalities, which we have reliable and salient conscious awareness of. However, there are two other important sensory pathways that are critical for motor control, that we are not as aware of, likely because they are typically activated by our own motor actions, so we don’t experience them separately from them. A number of other sensory modalities exist in other animals, including echolocation in bats, whiskers in a range of animals, and magnetic sensing in a range of animals, especially birds.

Figure 4.5 shows the basic facts about each of these sensory modalities, including what physical **stimulus** activates them, the names of the receptor(s) that **transduce** this stimulus into a neural signal, and how that neural signal makes its way into the cortex by way of the thalamus. Interestingly, our sense of smell, which is evolutionarily the most ancient sense, present in even the most primitive beasts in the ocean, bypasses the thalamus and jacks straight into the cortex, in a brain area that is close to the hippocampus. This may explain why odors can be such powerful memory triggers, as famously captured in Marcel Proust’s *Remembrance of Things Past*.

Modality	Stimulus	Receptors	Thal → Cortex	Absolute threshold
Vision (sight)	Light (photons)	Rods & Cones in Retina	LGN → V1	Candle flame from 30 miles (on clear night)
Audition (hearing)	Sound (variation in air pressure)	Hair cells in Cochlea	MGN → V1	Tick of watch at 20 feet (in a quiet room)
Olfaction (smell)	Airborne molecules	Hair cells in Olfactory epithelium	(none) → Olfactory cortex	1 drop of perfume diluted in air of 6 room house
Gustation (taste)	Food molecules	Taste Buds in Papillae	VPN → Insula	1 teaspoon of sugar in 2 gallons of water (try it!)
Somesthesia (touch)	Touch, pressure, temperature, pain	Free nerve endings in Skin	PMN, VPN → S1	Wing of a fly falling on cheek from 1cm
Proprioception (self movement)	Muscle stretch	Muscle Spindle fibers	VPS → S1	
Vestibular (balance)	Head rotation, acceleration	Semicircular Canals & Otoliths	VPN → S1	

Figure 4.5: The 5+2 main sensory modalities for people (first 5 are consciously salient, last 2 less so), and the stimulus that activates the corresponding receptors, and the pathway through the thalamus and into cortex. LGN = lateral geniculate nucleus; MGN = medial...; VPN = ventral posterior nucleus; PMN = posterior medial nucleus; VPS = ventral posterior, superior nucleus. The absolute threshold suggests how sensitive our receptors are (most animals are much more impressive in the olfaction department)

In computational, *signal processing* terms, the **subcortical** part of the sensory pathway, prior to the arrival of the neural signal in the thalamus and cortex, performs significant *preprocessing* of the signal using often complex, sophisticated, and evolutionarily pre-wired circuits. For example, the subcortical auditory pathways have very fast spiking neurons that can process the tiny time differences between when a sound arrives in one ear versus the other, to extract the angle of the sound source relative to the head. The cochlea also performs the rough equivalent of a *fourier transform* on the auditory signal, transforming the time-varying sound wave into a *spectrogram* organized according to the different *frequencies* of sound present. This format is then much easier for the cortex to process.

These subcortical processing steps also perform the first pass on **compression** and **contrast** processing of the overall sensory signal – this is very well understood in the visual processing pathway, as we'll see in a moment.

The role of the thalamus in the sensory pipeline is somewhat less clear. A unique feature of the thalamus is that the neurons there have essentially no direct connections amongst themselves, and therefore they would seem to be incapable of significantly transforming the sensory signal (which depends on these connections given the way that neurons process information, as covered in Chapter 2). However, the thalamus receives massive numbers of top-down connections from its corresponding sensory cortex (e.g., more than 90% of connections in the LGN come from V1), suggesting that a major function of the thalamus is in supporting *top-down modulation* of the incoming sensory signal. For example, this can support **attention**, where these top-down cortical signals “shine a spotlight” on a subset of incoming sensory signals – we return to this topic later in the section on Attention.

Vision

Figure 4.6 shows the overall pathway of visual information, from light rays reflecting off of objects in the world, that are then focused by the lenses of your eyes onto the **retina** at the back of the eyeball, where the light is transduced (converted) from photons into electrons by **photoreceptor** cells (rods and cones), in very much the same way that the camera in your cell phone does it. Speaking of which, your eye has roughly 120 megapixels (i.e., 120 million photoreceptors), but they are not distributed uniformly as they are in your camera (Figure 4.7). Instead, there are many more **cones** (color sensitive photoreceptors) in the

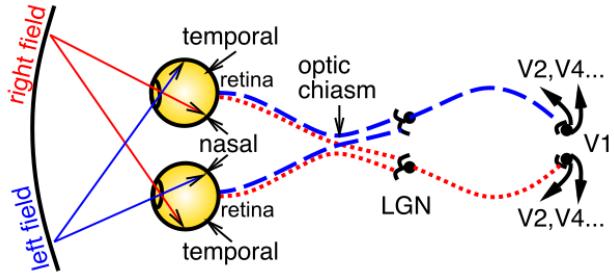


Figure 4.6: Visual pathway from light rays reflecting off of objects, through the lens of the eyes, and onto the retina at the back of the eyeball. The photoreceptors there transduce the light into electrical signals, which are then transformed into the firing of the retinal ganglion cells, that send their axons to the LGN of the thalamus, crossing over at the optic chiasm so the full left visual field ends up going to the right hemisphere, and vice-versa. The LGN neurons then communicate the visual signal up to area V1 in the very back of the brain, in the occipital lobe, and from there it makes its way back forward up the hierarchy of compression and filtering as shown in Figure 4.1.

center of your retina, called the **fovea**, with the density falling off rapidly as you go out from this center into the **periphery**, where the monochromatic, motion-sensitive **rods** predominate. If you do some math about *visual angles* as shown in Figure 4.7, the fovea can resolve about 300 dots per inch (dpi), which is the resolution of current high-res displays like those found on most higher-end cell phones. Thus the “retina” marketing from Apple is actually accurate.

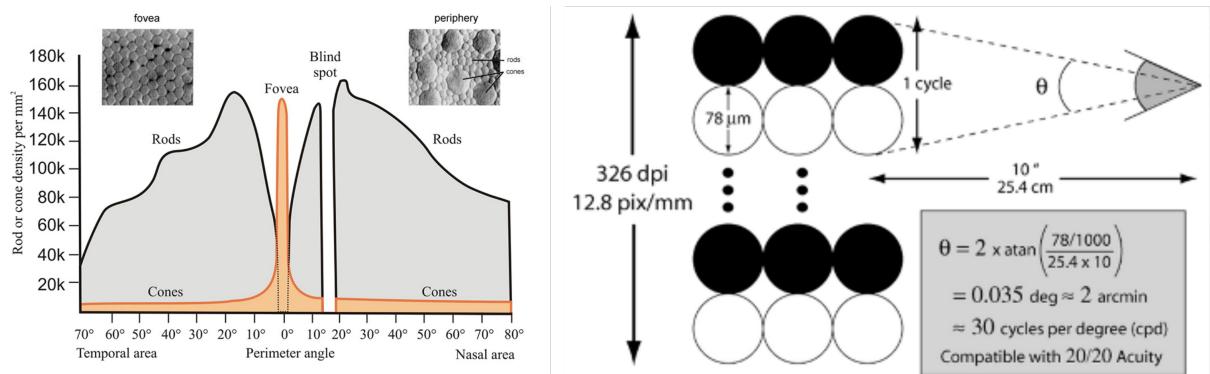


Figure 4.7: How many megapixels is your eye? First, unlike a camera, the photoreceptors are concentrated in the very center of your retina, called the **fovea**, which is also where most of your color-sensitive **cones** are (left panel). Out in the **periphery**, monochromatic **rods** predominate – they are also better at detecting motion, so you are actually better at seeing things move when you’re not looking right at them. The right panel shows that we can resolve about 300 dots per inch (dpi) at a distance of 10 inches, which is not coincidentally the “retina” screen resolution on modern high-res displays, and even on old-fashioned laser printers.

While the rods do not resolve multiple different colors, and have lower resolution, they are much better at detecting **motion**, and thus it can be useful to actually look away from something to detect motion better. Cats seem to know this trick, and will look away from the mouse playing dead (or the cat toy playing dead, more likely), to better detect when it starts to move. Or maybe they are just communicating disdain. Hard to know with cats.

There are several fairly crazy-but-true things about the visual system:

- The rod and cone photoreceptors are at the very *back* of the retina, behind all the extra pre-processing circuitry that does the compression and contrast enhancement. There are also blood vessels all over the place. Apparently they don’t end up blocking the light too much, or distorting its path, but still. Also, as shown in Figure 4.7, there is a huge chunk of visual space, about 1 degree of visual space in diameter (same as the fovea), which has no photoreceptors at all! This is the so-called **blind spot** where the axons for the optic ganglion cells (i.e., the *optic nerve*) are all gathered and head back to the thalamus.

You can find this spot (if you haven't already), by moving your outstretched arm with thumb pointed up, out to an angle of about 12:30 or 1 o'clock (where 12 is straight ahead), and noticing where, now that you're paying attention, you see it disappear! You don't notice it normally because, as we've said already, your perception is about the world, not your raw sensory signals, so your brain just papers over that little hole there for you, using all the signals surrounding it.

- The photoreceptors are constantly active (*depolarized*), and light actually *inhibits* them, instead of turning them on, as you would otherwise expect. Why this is the case seems to still be a mystery – it is not true across the animal kingdom, so it could go the other way. Interestingly, there are *many* instances like this across the brain, where neurons are *tonically* (continuously) active and then get inhibited by relevant signals. Although you would think this would cost a lot of energy, and the brain does consume about 1/3 of the human body's energy budget, the actual amount of energy required to sustain neural firing is likely a small fraction of this total cost (all the maintenance and upkeep and building of synapses, etc is likely much more expensive).
- Everything in the retina is upside-down and backward relative to the world (Figure 4.8). This is a result of basic optics, as evident in Figure 4.6 tracing the light rays from the world, through the lens, and into the retina. Do you need any more convincing that we see the world and not our raw sensory stimulus? Through experience in the world, we quickly learn the very systematic relationship between patterns of neural firing in the retina, and the direction of gravity, etc, and this is then what we perceive. This lesson is really important more generally: all neurons communicate with spikes that are essentially identical to those from any other neuron – they do not come with extra “annotations” indicating things like “up” or “down” – so each neuron has to learn *de novo* the meaning of each of its inputs in relation to all the others. This is also why it is actually quite possible that you are a brain in a vat (as in the *Matrix* movie) – everything we know is encoded in these patterns of neural spiking, and we only make plausible inferences about the most “reasonable” interpretation of all these patterns.

Also, if you want to know why mirrors only flip left and right, but not top and bottom, check out this video by [Physics Girl](#) who provides a comprehensive explanation. Basically, mirrors only reflect, they don't flip anything – we just get confused when looking at a reflection, because it shows us as we would look if we were facing ourselves from a point on the other side of the mirror. Careful not to fall into the looking glass!

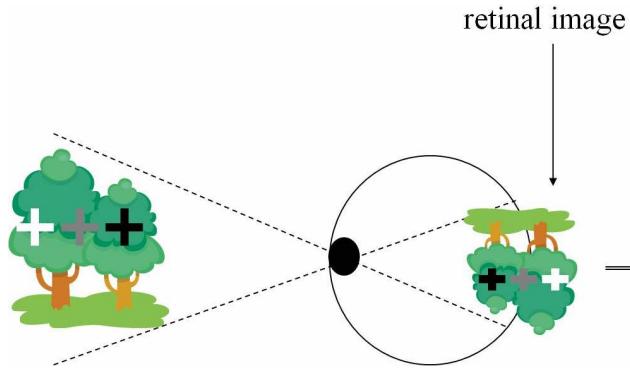


Figure 4.8: Our raw visual inputs are upside-down and backward relative to the real world, due to the basic optics of lenses – this is also true of a camera. Because we see the world, not the raw sensory signal, we learn the relationships between these visual signals and other reliable features of the world, like gravity. Indeed, all neural signals are completely arbitrary in the first place – there is nothing “up” or “down” about a pattern of neural firing – it is just spikes!

Compression and Contrast in Vision

Despite all the weirdness, the retina works really well, and one of the most important things it does is to **compress** and enhance **contrast** in the visual signal, before it is converted into neural spikes and sent through the **retinal ganglion** cells up to the LGN of the thalamus. Figure 4.9 shows some of the

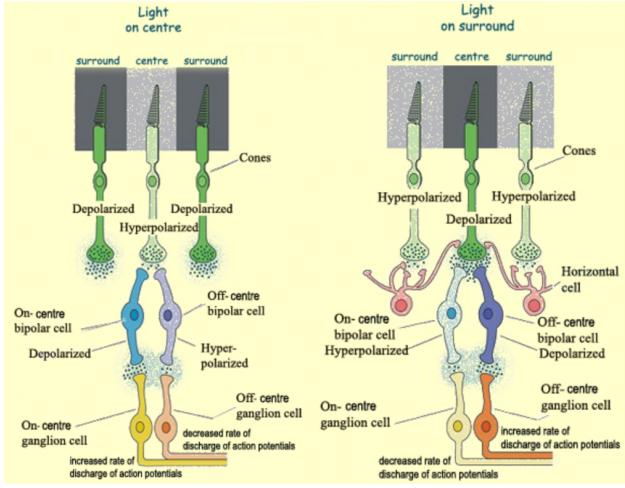


Figure 4.9: Circuits of dynamic pre-processing within the retina itself, which transform the raw transduced light signals by center-surround contrast coding.

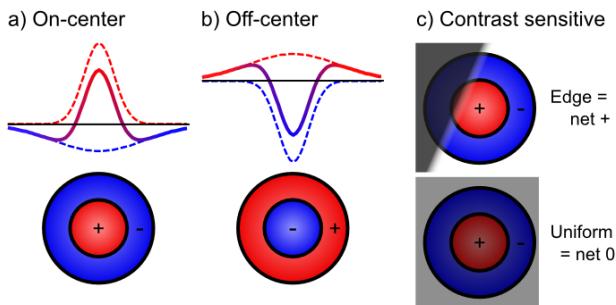


Figure 4.10: Diagrams of the receptive fields for on-center vs. off-center versions of center-surround cells. The receptive field (RF) refers to the pattern of illumination across the retina (or out in the world, for a given fixed position of the eyes) that directly influences the firing of a given neuron. For the on-center neuron, it is excited by light in a smaller central region, and inhibited by light in the wider surround region, while the off-center is the opposite. This center-surround organization is critical for *compression* by reducing or eliminating firing where the light is constant across the entire RF (the excitation and inhibition cancel out). Such neurons only respond to regions of *contrast* in the image, where there is a transition (edge) between a darker and lighter level of illumination. Thus, our visual systems mainly encode the regions where light and dark are changing – i.e., regions of *contrast* in the image – not the relatively uninteresting uniform regions in between.

preprocessing circuits that achieve this compression and contrast enhancement, through a feature known as **center-surround** contrast, which is illustrated in more abstract form in Figure 4.10.

Specifically, there are two main types of these center-surround signals that emerge out of the retina: **on-center** and **off-center**. The on-center type of retinal ganglion cell gets excited by light hitting the central region of the target-like **receptive field** shown on the left of Figure 4.10, and they are inhibited when light hits the wider outer-ring (the *surround*). The off-center is the opposite: it is excited by light in the surround, inhibited by light in the center.

Figure 4.10 illustrates the main consequence of these opposing patterns of center-surround excitation and inhibition: if there is a uniform, consistent amount of light across both the center and surround, the excitation and inhibition cancel out, and the cell will not fire! Only when there is some kind of transition or *edge* where more light falls on the center relative to the surround, or vice-versa, will these cells fire. Thus, the retina is already doing a huge amount of *compression* before sending signals up to the cortex, by filtering out regions of uniform illumination. These are the blank walls and blue sky in our visual inputs, and they are not where the interesting stuff is, so it is a good idea to filter this stuff out, and focus processing specifically on the *edges*.

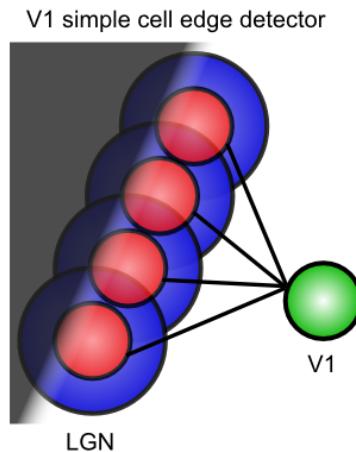


Figure 4.11: Simple cells in primary visual cortex (V1) combine multiple center-surround inputs from the LGN to form *edge detectors* that encode a consistent, elongated edge or transition in illumination across an image. This is referred to as the *classical receptive field* for V1 neurons, and its discovery by Hubel and Weisel in the 1950's and 60's won them the Nobel prize! The JPEG compression technique for pictures works by extracting these same kinds of oriented edges of contrast in images, and it greatly reduces the number of bits of information needed to encode large images.

The primary visual cortex (V1) builds on the center-surround signals coming from the LGN (which passes them along from the retina more-or-less intact), to detect more elongated, oriented edges of light / dark *contrast* in the image. Hubel and Wiesel discovered this property of V1 neurons by recording from V1 neurons in anesthetized cats shown simple images of oriented bars of light, and won the Nobel prize in 1981 for this and other discoveries about the visual system. Here's a [YouTube Video](#) about their work, and a more modern approach using [reverse correlation](#). These oriented edge detectors are shown in Figure 4.1 as the starting point for even more complex patterns detected at higher layers in the visual system.

Thus, in addition to seeing upside-down and backward, we mainly see the outlines or edges of things, and essentially assume the continuation of surfaces in between these edges. This may explain why we so readily process line drawings, which provide a good “illustration” of what our higher visual areas are largely processing. The efficiency of encoding the visual world in this way is demonstrated also by the widely-used JPEG compression standard, which greatly reduces the size of image files.

In summary, vision provides clear, well-understood examples for how the brain compresses incoming signals to extract features that will be of greatest value for subsequent stages of processing. In this and many other cases, this compression occurs by focusing on the key points of *contrast* – where things are changing. This principle of contrast also applies to the time domain as well – all stages of visual processing are also

particularly sensitive to changes over time, such as the onsets and offsets of illumination. This is particularly true of the rod-driven motion processing pathways, which dominate in the cat visual system, and explain why cats are so captivated by moving laser pointer dots.

Color Contrasts

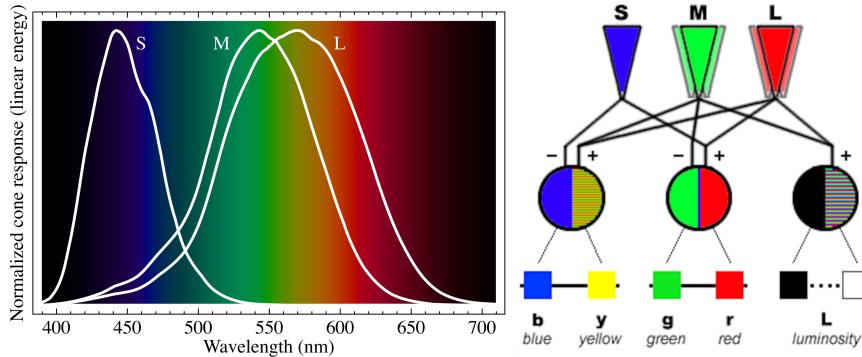


Figure 4.12: Color is encoded by receptors sensitive to different wavelengths of light (just barely different in the case of M and L) which are paired into two sets of opponents: red (L = long wavelength) vs. green (M = medium wavelength) and blue (S = short wavelength) vs. yellow, which is created by the combination of red and green, and is not provided by its own separate photoreceptor. Luminosity (black vs. white) is coded by integrating across all receptors, including rods which have a more blue-shifted tuning.

In addition to coding illumination contrasts at edges, the visual system is also tuned to contrasts between different wavelengths of light, which forms the basis for color vision (Figure 4.12). Color is efficiently encoded in the brain by using only three different types of cone photoreceptors, which can span the entire spectrum of visible light by mixing these three elements in different relative strengths, just as a painter can mix a wide range of colors from a small set of primary colors. Although we commonly think of these **primary colors** as Red, Green and Blue (RGB), Figure 4.12 shows that the Red and Green detectors are surprisingly overlapping in their response to different frequencies. For this reason, and because of their wide tuning, scientists refer to these photoreceptors as L = long wavelength (Red), M = medium wavelength (Green), and S = short wavelength (Blue), but we'll stick to the more familiar primary color names.

For the same reason that center-surround coding produces effective compression and contrast-enhancement, the brain also encodes color using color contrasts, using two **opponent** pairs of colors: Red vs. Green, and Blue vs. Yellow. Indeed, these color contrasts are often superimposed with the center-surround contrasts, such that e.g., the center responds to Red and the surround Green.

A number of visual illusions reveal the underlying color opponency at work in our visual systems, such as those in Figures 4.13 and 4.14, which interact with the constant motion of our eyes to create apparent motion where there is none (see [Akiyoshi Kitaoka's Web Page](#) for many other such demonstrations). Other illusions involve staring at one set of colors and then looking at a white page, wherein the fatigue caused in the one pair of the opponent allows the other to get a bit more active, driving its perception. [This illusion](#) provides a compelling set of illusory percepts based on red vs. green opponency.

Despite the presence of these “oopsie” cases in carefully-crafted illusions, the color opponent system normally enables us to make highly accurate inferences about the “real” colors of different objects, compensating for the impact of different lighting conditions. This is known as *white balance* in photographic terms, and in the brain it depends on the ability of e.g., a preponderance of yellow activation from yellow lighting producing a compensatory accentuation of the blue opponent, which we could see at work in “The Dress” illusion discussed at the outset.

Depth

Our retinas only provide a 2D window onto the full 3D world, so we have to rely on a number of assumptions to reconstruct that missing 3rd dimension – these assumptions then provide a treasure trove of illusions that

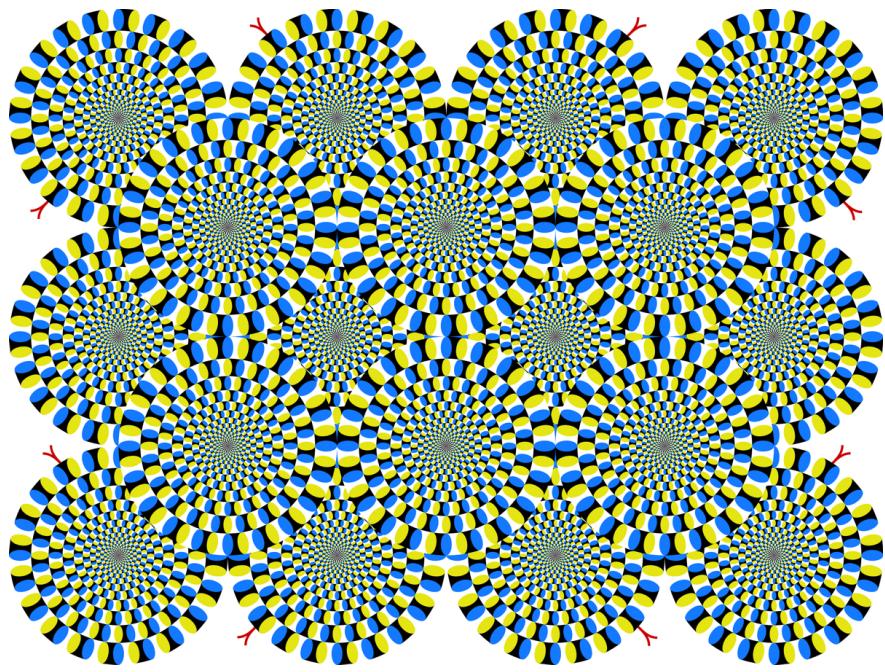


Figure 4.13: Rotating snakes illusion from Akiyoshi Kitaoka, which depends on eye movements interacting with blue-yellow opponent color coding.

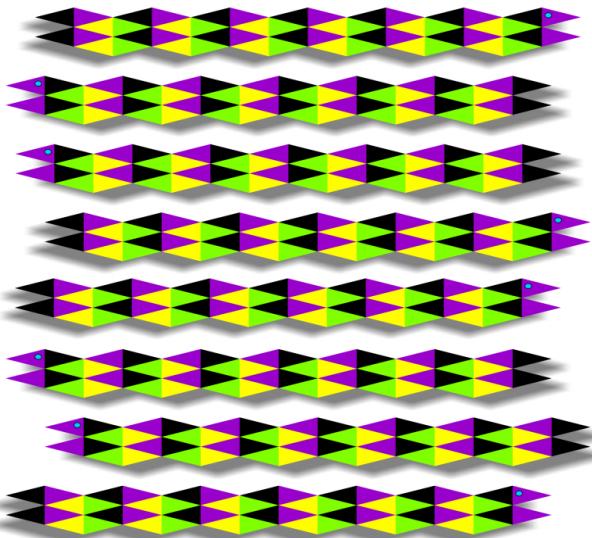


Figure 4.14: Moving snakes illusion from Akiyoshi Kitaoka.

can keep you entertained for hours! This missing 3rd dimension of depth is an important example of the *ill posed* nature of perception, which is a mathematical description of a situation where you have more unknown variables than data points available, so you have to rely on extra constraints or assumptions.

There are two major categories of **depth cues** used to reconstruct the missing 3rd dimension:

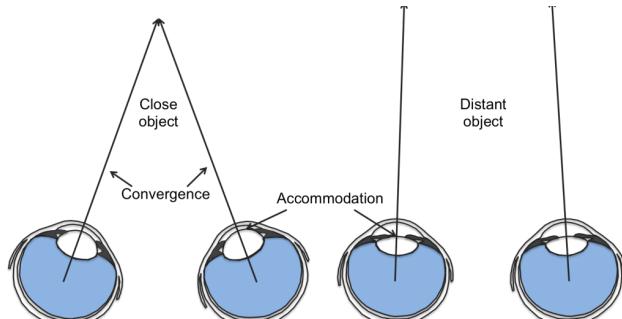


Figure 4.15: Convergence of eyes on a common point provides one of the two binocular depth cues. The other is different offsets of features across the two eyes, produced by parallax.

- **Binocular** cues that depend on the **parallax** effect, where the left and right eyes receive slightly different views of the world based on their different horizontal positions. The use of these two views to extract depth is known as **stereopsis** (same root as “stereo”, and same effect as stereo sound). There are two forms of these binocular cues: the **retinal disparity** (differences) across the left and right eye views of the same general region of visual space (enabled by the cross-over and combination of different-eye views into the same visual hemisphere, as shown in Figure 4.6), and the **convergence** of the two eyes to focus the center of vision at the same point in depth (Figure 4.15).

These binocular cues give the most vivid sense of 3D depth, and are what 3D movie and TV technologies provide, by being able to project different images to the two eyes, either by using different color filters (red vs. blue, as in the old 1950’s 3D glasses) or different polarization of light (the current generation of 3D glasses). 3D technology has generally remained a bit of a “fad”, not essential like color, because we have so many other cues to depth (and the glasses can be uncomfortable, and they cut the image brightness in half).

Figure 4.16 shows an *autostereogram* or single-image stereogram, which is a version of a *random dot stereogram*, popularized by the *Magic Eye* books several years ago, which requires you to adjust your eye convergence outward (“wall eyed”, focusing well beyond the image plane itself), causing each eye to see a different offset of the random dots in the image – these offsets across the two eyes are the *retinal disparity* signal, and the fact that you can see depth with *only* this retinal disparity signal was a big deal when first discovered by Bela Julesz in 1959. When you (eventually) get just the right eye position, a previously-hidden 3D world gradually materializes out of the sea of random dots! This takes a lot of practice and patience – see [wikipedia page](#) for more info and examples.

- **Monocular** cues, which operate strictly within a single 2D image, include **occlusion** (one object in front of another), **relative size** (larger = closer), **texture gradients** (also including local surface texture indicating convex and concave shapes), and **linear perspective**, as shown in Figure 4.17. There are many others, especially if you include motion, which can give strong depth signals. One can trace a progression in art over centuries in terms of the use of these cues to create a perception of depth, with occlusion, relative size, and texture gradients being among the earlier ones, while the use of linear perspective was revolutionary in the renaissance period. Modern-day sidewalk artists employ these techniques to great effect (Figure 4.18).

Compression in Object Recognition

As illustrated in the hierarchical filtering process from Figure 4.1, every stage of processing in the perceptual pathways (and everywhere in the cortex more generally) produces more and more compression and contrast effects like those we’ve seen already, starting right in the retina itself. Thus, you should not be surprised



Figure 4.16: Autostereogram – the 2D image gives a hint as to the magic 3D world that lies just under its surface – you just have to focus your eyes out in the distance, and let your brain “settle” to see it.



Figure 4.17: Monocular depth cues including linear perspective and texture enable us to see depth in otherwise flat, 2D images. The left panel also shows how our brains automatically use these depth cues to infer the size of different objects in the scene, causing you to automatically perceive the top bar as much larger than the bottom one, even though they are identical in size. This is another example of how we perceive the world, not the raw sensation, and is known as *size constancy* in this case.



Figure 4.18: Sidewalk art that takes advantage of monocular depth clues to provide a compelling illusion of depth (but only when viewed from the right point).



Figure 4.19: We are biased to see objects, especially faces, even where none exist. This reflects the higher-level compression of visual scenes into known object categories.

to learn that the brain has a strong bias toward organizing the features in an image into a much simpler, compressed encoding in terms of *objects*. Saying “it’s a dog” represents a massive degree of compression relative to all the visual information that goes into the image of a typical dog. This bias toward seeing known objects in images is behind the common “mild hallucination” of perceiving objects where none exist, for example in the shapes of clouds, wood grain, toast – whatever has enough raw material to organize into objects (Figure 4.19).



Figure 4.20: An image of a plate of spaghetti transformed by Google’s Deep Dream neural network algorithm that progressively enhances features of the image consistent with what it has learned across a large number of “normal” images. This produces hallucinogenic images similar to those seen on LSD and other psychedelic drugs, reflecting the imposition of our simplifying biases onto images – just a more extreme form of Figure 4.19.

Figure 4.20 (yep) shows the output of a neural network model trained on millions of photographs, when the input image (a plate of spaghetti in this case) was progressively altered in a way that better fit the internal biases of the model. These images resemble the kinds of hallucinations produced by psychedelic drugs such as LSD, suggesting that these drugs have the effect of enhancing the influence of internal representations over the raw input stimuli. Thus, we are all mildly hallucinating all the time, and the waking dream state produced by these drugs just accentuates these processes.



Figure 4.21: Is there something more than just dots in this image? The ability to extract the underlying shapes from this image represents the benefit of top-down activation in driving our perceptual system – even though it has a mild hallucinatory side-effect, it is often beneficial to be able to see partially-occluded or otherwise hard to see objects.

Figure 4.21 demonstrates the important benefits of this bias to “see things” in images – more often than not, there actually are things there, which might be obscured in various ways, and having the ability

of top-down expectations and biases to influence our perception helps pull these things out (O'Reilly et al. 2013).

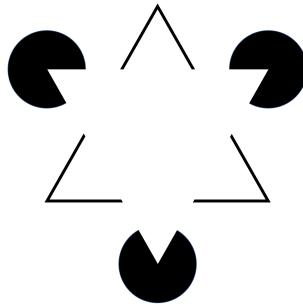


Figure 4.22: The classic Kanisza triangle, demonstrating how we interpret the “suspicious coincidences” of the wedges and apparent occlusion of the background triangle in terms of the simpler percept that there is a white triangle in front of the other items.

The classic **Kanisza triangle** (Figure 4.22) provides another demonstration of this bias toward seeing a simpler, more compressed encoding of the world. These kinds of effects have traditionally been explained in terms of **gestalt principles**, developed by a school of influential German school of psychologists in the early 1900's to explain how we tend to impose higher-level *gestalt* groupings onto images (Figure 4.23). From the modern perspective, the attempt to explicitly enumerate long lists of such “principles” seems like a mismatch relative to the way that top-down and bottom-up *constraints* or *biases* interact in a more graded, emergent way in actual perception. These biases and constraints are much “softer” and more “fluid” than things you might articulate as “principles”.

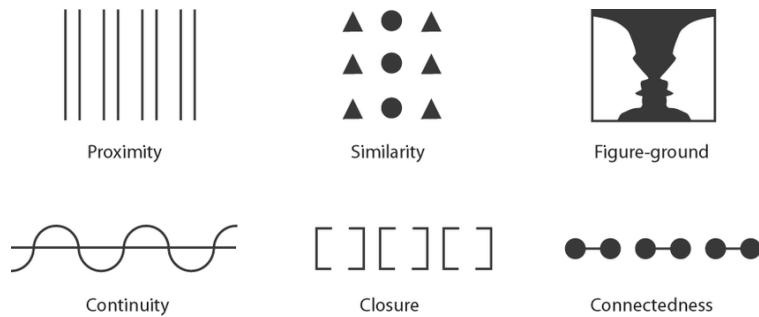


Figure 4.23: Some of the many gestalt principles for how items in an image are grouped and organized.

And speaking of soft, fluid images, the artwork of Salvador Dali provides an excellent illustration of the hallucinogenic top-down biases at work, trying to organize and simplify the world (Figure 4.24).

Time Contrast: The Novelty Filter

Another very important form of **contrast** that drives perception is *contrast over time* – i.e., “the news” – the visual system prefers new stimuli and new ways of seeing things. In other words, it functions as a **novelty filter** – filtering out the old and focusing on the new. Like compression, this happens at all levels in the system, from the retina on up. In the retina, if you prevent the eye from moving at all, and present a static image, the firing of retinal neurons will slowly start to fade away, and the world will go black (this experiment has actually been conducted by paralyzing the eye muscles! The various color opponency illusions mentioned earlier also demonstrate this novelty-filter property, where staring at one color causes the opponent color to be relatively more activated).

Figure 4.25 demonstrates a higher-level version of the neural adaptation that drives these novelty filter effects. If you stare at these images for long enough, you will find your brain spontaneously switching to a

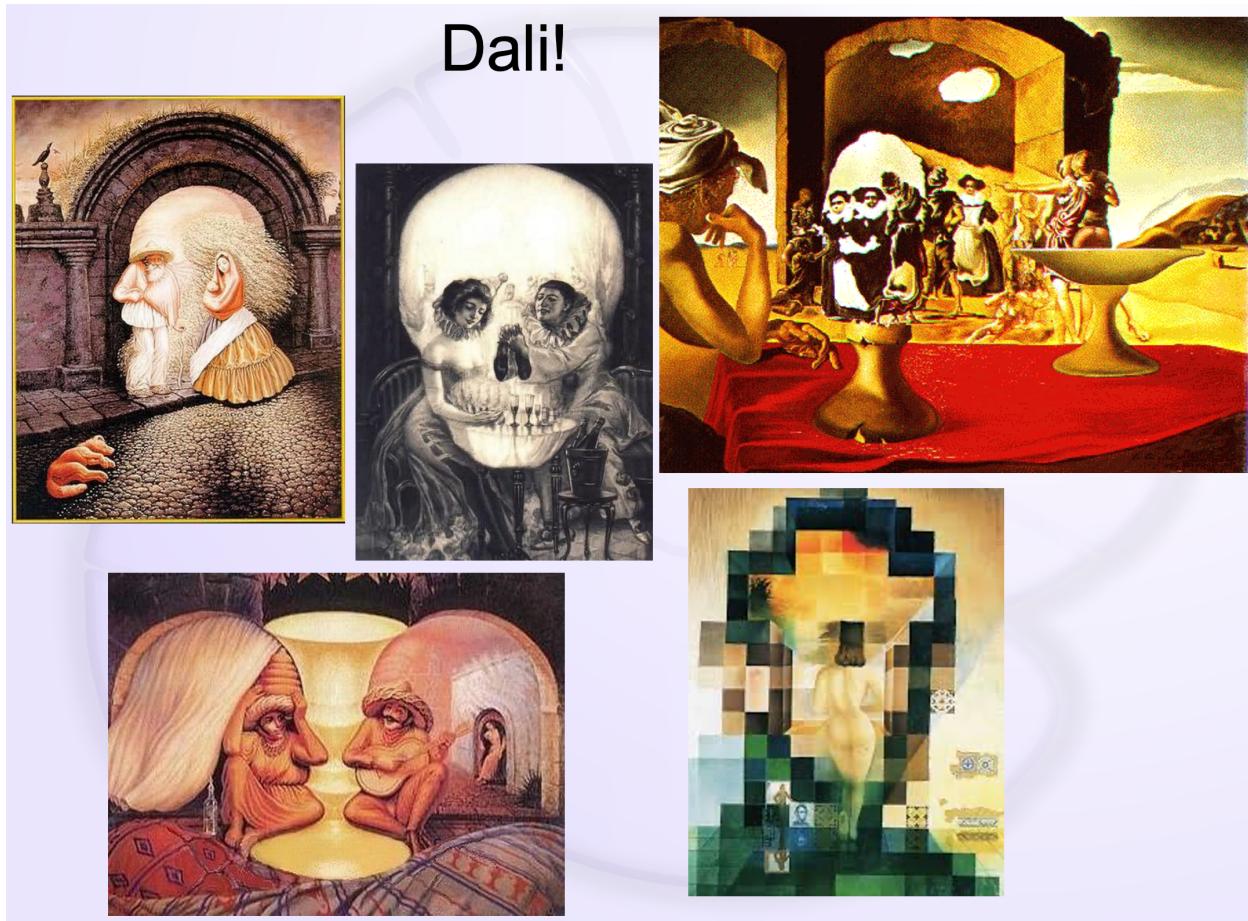


Figure 4.24: The art of Salvador Dali demonstrates the power of object (especially face) biases in perception.

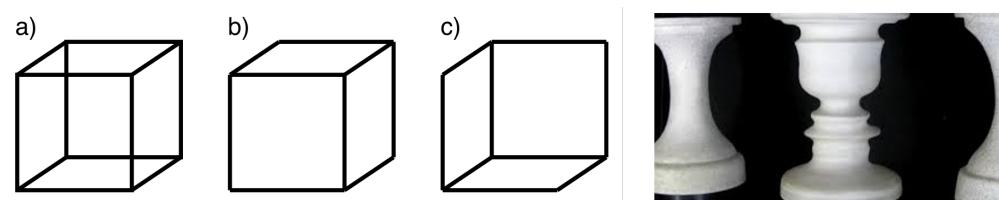


Figure 4.25: The visual system also enhances *contrast* over *time* – active neurons experience *accommodation* or *adaptation*, thus favoring novel perspectives instead of continuing to see the same thing over time. Stare at the *necker cube* on the left, or the ambiguous figure-ground image on the right, and you'll find your brain spontaneously switching between the two different ways of seeing them.

new way of seeing the image, without any explicit attempt on your part to do so. In fact, it is difficult to prevent your brain from switching in this way – because it is built right into the neural hardware.

Audition

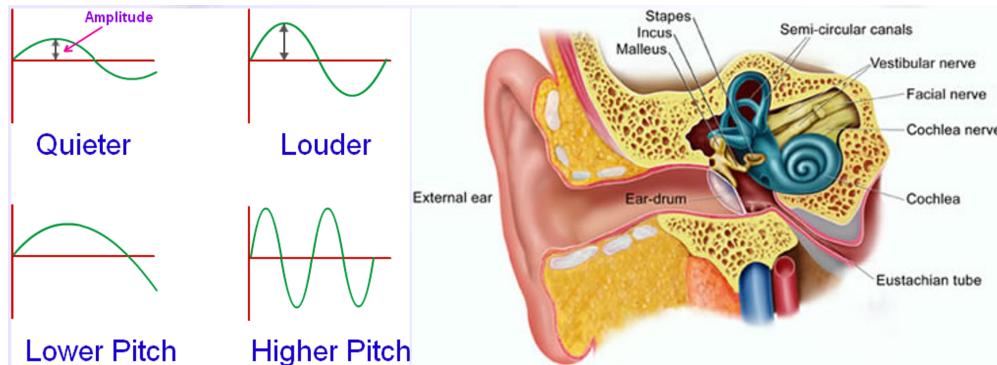


Figure 4.26: Overview of the auditory system, where rapid changes in air pressure (sound) are amplified by the tiny bones in your inner ear, causing the hair cells in the cochlea to bend – this mechanical force is then transduced into neural firing and processed by many stages of subcortical networks before making its way up to the primary auditory cortex (A1) by way of the MGN nucleus of the thalamus.

Because the same principles just explored for the visual system apply to all of the other sensory modalities, and it is harder to share the relevant perceptual experiences via a book, we will give these other modalities the short shrift that my son recommended at the outset. Figure 4.26 shows the essential features of the auditory transduction process, where sound waves are converted into neural firing, via hair cells in the cochlea.

First, sound travels as *waves*, and thus can be described in terms of its **amplitude** (intensity) and **frequency or wavelength** (pitch). These vibrating sound waves cause the **ear drum** to vibrate (like a drum!) and this vibration is then amplified by tiny bones (**ossicles**: Malleus, Incus, Stapes – no you don't need to memorize these!) that cause your inner squid (the liquid-filled **cochlea**) to vibrate in tune with the sound. Inside the cochlea are only about 3,500 inner hair cells that are responsible for transducing the liquid vibrations into neural firing signals. This is an incredibly tiny number of *anything* at the cellular level – your hearing is precious and you should do everything you can to preserve those priceless cells!

The cochlea transforms the time-varying sound waves into a much more useful kind of representation, remarkably similar to what is typically done in artificial signal processing approaches, by effectively performing a *fourier transform*. This occurs by the **place coding** of frequency, such that different subsets of hair cells are activated for different frequencies of sound (for low frequency sounds, however, the frequency remains as a time-varying neural firing signal). This is called a **spectrogram**, typically plotted with time on the X axis and frequency on the Y axis. By splitting out different frequencies across different neurons, it becomes easier to recognize patterns that span across these frequencies – e.g., the distinctive patterns of human speech, called *formants* are characteristic patterns of frequency changes over time.

Interestingly, in the spectrogram representation, these patterns look like contrast edges at different orientations (increasing, decreasing), which are exactly the same patterns processed at the lowest levels of the visual system. Indeed, the very same neural networks that do a good job of recognizing visual patterns can recognize these auditory patterns as well. This then makes sense of the amazing experiment where the visual signals in a ferret were re-routed to its auditory cortex, and the cells there developed very similar firing patterns as are typically found in visual cortex (Angelucci et al. 1997).

In short, the auditory cortex does the same kind of compression and contrast processing of the auditory signals, extracting simpler ways of summarizing all of the sound information through a series of hierarchical layers. At the upper levels, which we are consciously aware of, we have things like “Max wants a banana”, summarizing a long complex auditory stimulus with a relatively few bits of information.

Attention

In addition to all the compression described so far, there is another critical driver of compression effects in perception, known as **attention**. Subjectively, attention is often described as a **spotlight**, shining bright mental light on one, or at most a few, items in the current *attentional focus*, which also has the consequent effect of pushing everything else off into the shadows. Like compression, attention operates everywhere in the cortex, and can be understood in terms of the interactions between bidirectional excitation and inhibition among neurons, where inhibition is what pushes everything else out into the shadows, and bidirectional excitation reinforces the attentional focus (Cohen et al. 1994). Not coincidentally, this bidirectional excitation is the same central ingredient in consciousness that we discussed in Chapter 3 (aka recurrent processing) – the focus of attention typically corresponds to what we are consciously aware of.

There is a special part of the brain, in the *parietal lobe*, that seems to be particularly important for **spatial attention** (i.e., paying attention to different parts of space), which has been extensively studied. Although attention itself is ubiquitous, spatial attention is particularly important in perception because discrete objects tend to occupy different regions of space, and thus we tend to use this spatial attentional focus as a way of directing attention at different objects of interest. This special role for spatial attention is also directly tied to the fact that the very same neural circuits in the parietal lobe are used for deciding where to move our eyes. Thus, attention is typically synonymous with *looking*, and the motor action of looking (moving the eyes) requires working with the spatial coordinates of where to look.

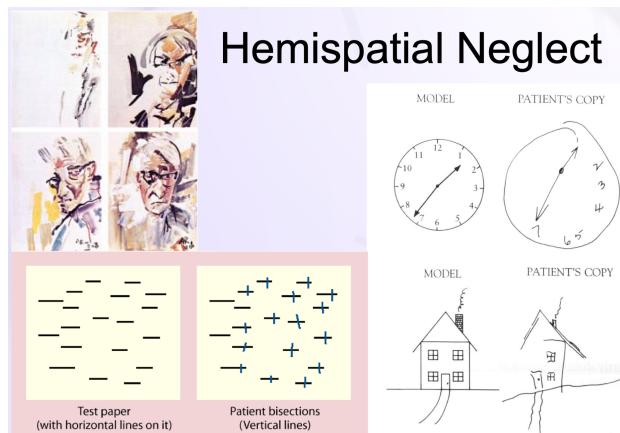


Figure 4.27: Demonstrations of hemispatial neglect. Upper left: Progression of self portraits by an artist with hemispatial neglect, showing gradual remediation of the neglect over time. Right: Drawings of given target objects by patients with hemispatial neglect, showing profound neglect of the left side of the drawings. Lower left: Results of a line bisection task for a person with hemispatial neglect. Notice that neglect appears to operate at two different spatial scales here: for the entire set of lines, and within each individual line

Some of the most striking evidence that the parietal cortex is important for spatial attention comes from patients with **hemispatial neglect**, who tend to ignore or neglect one side of space (Figure 4.27). This condition typically arises from a stroke or other form of brain injury affecting the right parietal cortex, which then gives rise to a neglect of the left half of space (due to the crossing over of visual information). Interestingly, the neglect applies to multiple different spatial *reference frames*, as shown in the line bisection task, where lines on the left side of the image tend to be neglected, and also each individual line is bisected more toward the right, indicating a neglect of the left portion of each line (Figure 4.27).

The Posner Spatial Cueing Task

One of the most widely used tasks to study the spotlight of spatial attention is the Posner spatial cueing task, developed by Michael Posner (Posner 1980) (Figure 4.28). One side of visual space is cued, and the effects of this cue on subsequent target detection are measured. If the cue and target show up in the same side of space (*valid* cue condition), then reaction times are faster compared to when they show up on different sides of space (*invalid* cue condition). This difference in reaction time (RT) suggests that spatial attention is

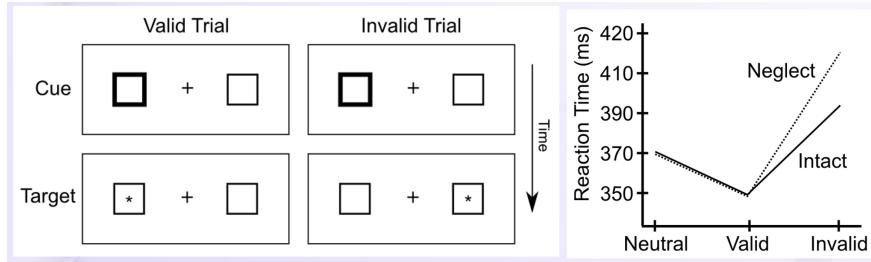


Figure 4.28: The Posner spatial cueing task, widely used to explore spatial attention effects. The participant is shown a display with two boxes and a central fixation cross – on some trials, one of the boxes is cued (e.g., the lines get transiently thicker), and then a target appears in one of the boxes (or not at all on catch trials). The participant just presses a key when they first detect the target. The typical data plotted on the right, showing that the target detection reaction time is quicker for valid cues vs. invalid ones. This suggests that spatial attention was drawn to that side of space. Patients with hemispatial neglect exhibit slowing for targets that appear in the neglected side of space, particularly when invalidly cued.

drawn to the cued side of space, and thus facilitates target detection. The invalid case is actually worse than a neutral condition with no cue at all, indicating that the process of reallocating spatial attention to the correct side of space takes some amount of time. Interestingly, this task is typically run with the time interval between cue and target sufficiently brief as to prevent eye movements to the cue – thus, these attentional effects are described as *covert attention*, while eye movements constitute *overt attention*.

Patients with hemispatial neglect show a disproportionate increase in reaction times for the invalid cue case (Figure 4.28), specifically when the cue is presented to the good visual field (typically the right), while the target appears in the left.

Psychophysics

Last, and frankly, likely least for most readers, we are required by an unwritten psychology rule to cover the field of *psychophysics*, which is an attempt to measure the most basic aspects of perception with a degree of precision that is intended to impress our colleagues in physics (who likely remain unimpressed). These basic aspects of perception center around finding the very weakest stimulus intensity that can be detected, across different modalities.

The **absolute threshold** is the name of this very weakest stimulus, and it is listed in Figure 4.5 for each of the main sensory modalities. A key factor here is defining what it means that a stimulus can be detected – what if you can't quite detect it all the time, but still most of the time – does that count? The convention is actually to push all the way down to a 50% probability of detection.

We can also determine the **discrimination threshold** or **just-noticeable difference (JND)**, which is how big of a *difference* between two different stimuli that can be reliably detected (again typically at the 50% probability level). One of the most exciting results in psychophysics is that this JND is a function of the intensity of the stimuli, and the famous **Weber's Law** says that it is a constant proportion of the overall intensity. For example, if people can detect differences in fairly dim lights of a few percent, then to detect differences in much brighter lights, the raw differences must also be much larger, exactly in proportion to the overall intensity of the lights. This proportion is known as the **Weber fraction**.

One of the reasons we have to cover this (btw, it is all over – not so bad after all!) is that psychophysics was one of the earliest examples of scientific psychology, starting with Ernst Weber's work in the 1830's and fully established by Gustav Fechner in 1860.

Summary

This chapter has a lot of detailed information, but the overarching theme of *compression* and *contrast* hopefully comes through. Your brain is wired to be a *simplicity filter* and a *novelty filter*, delivering the simplest interpretation of an complex pattern of sensory input, and focusing on what is new and different.

These same processes, which can be tied directly back to the properties of neurons as explained in Chapter 2, operate throughout your brain, at all levels, shaping how you perceive other people (in terms of simplifying stereotypes) and the world (always seeing out news, and quickly discounting the past). Thus, perception truly is a window onto the soul, and we return to these lessons throughout the remainder of the textbook.

Summary of Key Terms

This is a checklist of key terms / concepts that you should know about from this chapter. As we'll learn in the memory chapter, it is a great idea to test yourself on what was said about each of these terms, and then go back and double-check – that provides both beneficial repetition and also the *testing effect*.

- Hierarchical compression
 - Filtering
 - Raw sensation vs. subjective perception
- Perceptual constancy
 - Color constancy
 - Size constancy
 - Assumptions, illusions
- Sensory systems
 - physical stimulus, transduction, subcortical preprocessing
 - vision, rods, cones, retina, LGN, V1
 - audition, hair cells, cochlea, MGN, A1
 - olfaction, hair cells, olfactory epithelium, olfactory cortex
 - gustation, taste buds, papillae, VPN, Insula
 - somesthesia, free nerve endings
- Vision
 - retina, periphery, fovea, photoreceptors, blind spot
 - cones for color, rods for motion
 - center-surround contrast detectors
 - R,G,B primary color photoreceptors (L,M,S)
 - red-green, blue-yellow opponent color coding
 - binocular depth cues: retinal disparity, convergence
 - monocular depth cues: occlusion, relative size, texture, linear perspective
 - object, top-down biases, hallucinations, Kanisza triangle
 - gestalt principles
 - time contrast: novelty filter
- Audition
 - amplitude, frequency, wavelength
 - ear drum, cochlea, hair cells
 - place coding, spectrogram
- Attention
 - spotlight
 - spatial attention in parietal lobe
 - hemispatial neglect
 - Posner spatial cuing task
- Psychophysics
 - absolute threshold
 - discrimination threshold, just-noticeable-difference (JND)
 - Weber's law, Weber fraction

Chapter 5: Learning, Motivation, and Emotion

Learning is the single most important process taking place in the brain. Without learning, nothing else is possible. All of our focus on the three-C's of compression, contrast, and control presumes a brain with sensible patterns of synaptic connectivity, that produce *useful* forms of each of these phenomena. Without learning, neurons would randomly compress incoming sensory information, detecting irrelevant, bizarre features that don't have any behavioral relevance. Contrast would compare these random things against each other, producing equally meaningless relative comparisons. Control would drive us toward random goals, and our behavior would be just a jumble of strange impulses.

Learning is essential because there are *way* (way, way, way...) too many synapses for any kind of genetic process to shape in a detailed way. There are only about 20,000 different protein-coding genes in the human genome, which is only 2 times the number of synaptic inputs on a *single* neuron. It is inconceivable that genes could code for any sensible fraction of the 100 *billion* times that amount of information that would be required to configure the full human brain. This genetic argument accords with the obvious fact that we learn the vast majority of our abilities over an extremely protracted developmental window, in a way that depends critically on the experiences and education that we are exposed to.

Thus, the brain (specifically the neocortex) is fundamentally a *self-organizing* system, which somehow magically transforms raw sensory inputs into *knowledge* encoded in its billions of synaptic connections. The mystery of this process has long perplexed philosophers, who have explored the opposing ideas of **empiricism** vs. **rationalism** and positions in between. Empiricists embrace the idea that learning proceeds directly from sensory experience, while rationalists argue that there is no way that raw experience by itself is sufficient to create the sophisticated level of knowledge an adult human (philosopher) has. Modern scientific approaches to this question retain much of this ancient debate, with some favoring a generous amount of innate knowledge, and others arguing that almost everything is learned.

We'll return to these issues in the Development chapter, but the quick summary is that neither extreme view is likely to be correct, with genetic and experiential factors each playing critical roles. In particular, there is ample evidence that genes establish broad patterns of initial connectivity and orchestrate developmental transitions, such as synaptic pruning, which in turn strongly influence an experience-driven learning process operating at synapses throughout the neocortex.

Our objective in this chapter is to first understand the nature of these synaptic learning processes, which have been figured out in spectacular detail at this point, and explore some broader ideas about how they might result in this magical self-organization of knowledge over development. Then, we turn to the forms of learning that were the focus of behaviorism: *classical and operant conditioning*. These both depend on similar dopamine-driven learning mechanisms operating in the basal ganglia, amygdala, and related areas, which are now very well understood. These forms of learning shape our core decision-making process to select actions that are likely to be rewarding, and not punishing.

Finally, we broaden our perspective beyond the limited world-view of the behaviorists, and consider the possibility that *internal* factors such as *goals*, *drives*, and, ultimately, *emotions*, might play a central role in driving both our learning and decision-making behavior. This perspective, long embraced by social psychologists, is only recently beginning to be explored from the neuroscientific angle, which has been perhaps overly-enamored with the remarkable alignment between the classic externally-driven behaviorist conditioning processes and the function of dopamine in the basal ganglia.

Synaptic Plasticity

If learning is the most important thing in the brain, then the most important thing about learning is that it takes place in the synapses interconnecting neurons. This idea goes back at least to Santiago Ramon y Cajal in the late 1800's, the pioneering Spanish neuroscientist who advanced the idea that interconnected networks of neurons are doing most of the work in the brain. Logically, the strength of the connections between neurons should alter the patterns of information flow through these networks, and thus makes sense as the primary locus of learning, and knowledge.

Donald Hebb cemented this idea with a compelling, well-specified proposal that memories are formed when neurons that are active at the same time increase the strength of their synaptic connections, so that

they are then more likely to co-activate each other in the future (Hebb 1949). In effect, learning is “gluing together” the different elements of a memory. This idea has been captured with the pithy statement that “neurons that fire together, wire together”.

However, it was not until 1966 that this **Hebbian** form of learning was actually demonstrated in the brain, by Bliss and Lomo (Bliss and Lomo 1973). They described a form of **Long Term Potentiation (LTP)** of the synaptic strengths between well-defined groups of neurons, where potentiation means “getting stronger” and the “long-term” aspect of it was critical to distinguish from earlier discoveries of synaptic potentiation that only lasted for a few minutes. If synaptic changes are really the basis for learning and knowledge in the brain, they had better last for more than a few minutes, because clearly our memories and knowledge can last a very long time.

The field of LTP research expanded rapidly from that point onward, and progressively more detailed questions were addressed about the exact nature of what is changing in the synapses, and what specific factors in the activity of the sending and receiving neurons on either side of the synapse were critical for causing it to change. After many controversies and twists and turns in this amazing story of scientific discovery, we now have a very solid and detailed understanding of how this process works, at least in terms of the underlying biochemical mechanisms. It is a fabulous success story for the power of the scientific method, to drill down and figure out exactly how some complex system actually works. Perhaps most remarkably, Hebb’s original idea seems to have been nicely supported, by a remarkable interaction of different moving parts: changing the strength of the synapse requires *both* the sending and receiving neurons to be active.

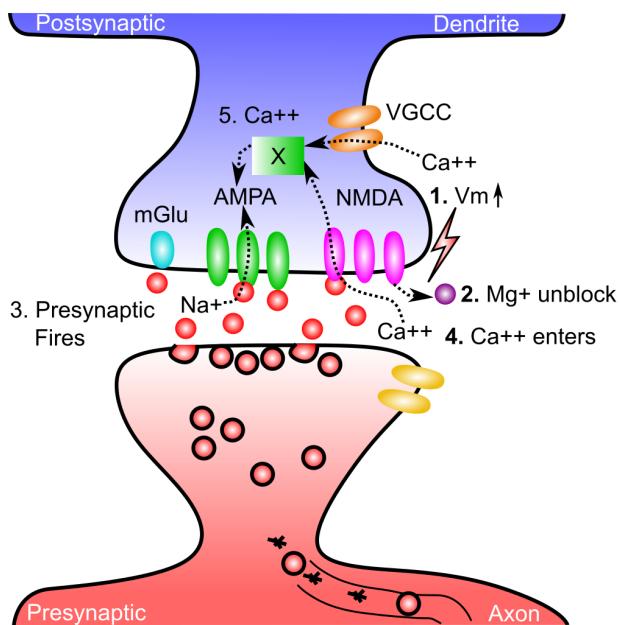


Figure 5.1: Mechanisms of synaptic plasticity, resulting in changes in the overall strength of the synaptic connection between the sending axon and the receiving dendrite. 1. The receiving neuron must be active, so that: 2. its elevated membrane potential (V_m) kicks out the positively-charged Mg^+ ions from the NMDA receptors. 3. The sending neuron must fire and release glutamate, which then binds to the NMDA receptors. 4. Causing them to open and Ca^{++} ions to enter. 5. Ca^{++} then triggers complex chemical pathways that ultimately result in changes in the numbers of AMPA receptors poking out across the membrane, which thus changes the overall amount of Na^+ that can enter for any given firing of the sending neuron.

Figure 5.1 shows the major steps in the process of synaptic change. The receiving neuron must be active enough so that its elevated membrane potential pushes out positively-charged magnesium ions (Mg^+), which are otherwise blocking the opening of the *NMDA* receptors. And the sending neuron must be actively releasing glutamate neurotransmitter, as a result of spiking, because glutamate binding to the NMDA receptors (in addition to the AMPA receptors) is necessary to cause them to open. Whereas AMPA receptors allow Na^+ ions to flow into the cell, NMDA allows *calcium* (Ca^{++}) ions to enter, and these Ca^{++} ions then trigger a

cascade of chemical reactions that ultimately leads to the change in synaptic plasticity. This critical role for Ca^{++} is consistent with many other similar such biochemical processes throughout the body – evolution often reuses existing mechanisms.

The main consequence of Ca^{++} entry is a change in the number of AMPA receptors in the synapse, which then changes the overall amount of Na^+ that can enter when the sending neuron spikes. Much more can be said about the details of these Ca^{++} driven chemical pathways (Rudy 2013), and the other associated changes that take place in the synapse, but the core logic remains the same as Hebb envisioned it: both neurons must be active for the synapse to change.

However, Hebb overlooked one *essential* aspect of learning, which was also neglected in the early days of research on LTP. This is the fact that you can't only *increase* the strength of synapses. Eventually, all the synapses would get ever-stronger, and the brain would blow up in a huge epileptic seizure. Instead, it is equally if not more important that synapses also *decrease* in synaptic strength, which has been named **Long Term Depression** or **LTD**. Decreases may be more important than increases, from the perspective of the *compression* function of neurons: each neuron has to essentially throw away a huge amount of information in order to compress its 10,000 inputs into a single output signal, and LTD makes synapses weaker and thus facilitates this information filtering process.

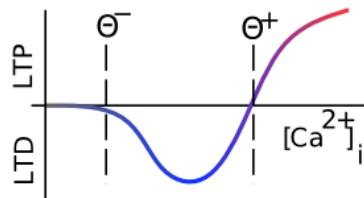


Figure 5.2: Direction of synaptic change as a function of the amount of calcium entering the dendritic spine. Lower amounts of Ca^{++} result in LTD = long-term depression or decrease in synaptic strength, whereas higher amounts result in LTP = potentiation or increase in synaptic strength.

In any case, Figure 5.2 shows that the balance between LTP and LTD is a function of the overall amount of calcium entering the dendrite – lower amounts result in LTD, while higher amounts result in LTP. This behavior emerges from a competition between two different chemical pathways, one which drives LTP and the other LTD, and their relative dependence on Ca^{++} levels. This is yet another tug-of-war taking place within neurons – this competitive dynamic is a very commonly-found mechanism at all levels of the brain.

One intriguing finding that makes sense in terms of this balance between LTP and LTD, is that weak activation of perceptual inputs seems to make those things harder to see, while strong activation makes them easier to see (Newman and Norman 2010). Thus, the weak activation leads to the lower levels of Ca^{++} , and causes LTD, whereas the stronger activation drives higher levels and LTP.

Neocortical Learning

Now that we know in detail how learning operates at the synaptic level, you might think that all of the mysteries of brain function should be solved, given what we said about the essential role of learning. Unfortunately, this is not the case. There are a number of challenges here, but chief among them is that there are so many synapses and neurons involved in learning any given bit of knowledge, that it is essentially impossible to go directly from behavior of the individual synapse up to this *emergent* behavior of learning in the larger neural network. The major tool that can be used to bridge this gap are computer simulations of neural networks, with equations capturing something like the function shown in Figure 5.2, operating within networks of neurons that behave something like real neurons in response to stimulus inputs.

Extensive work with such models has repeatedly shown that the known Hebbian-like learning mechanisms described above does *not* result in the kinds of larger-scale learning that people are clearly capable of. The reasons for this are well understood, but beyond the scope of this discussion. Furthermore, the kind of learning that *does* work reliably in these neural models, and is used in the recent powerful AI (artificial intelligence) models currently powering the speech recognition and other advanced capabilities in your cell phone and

other gadgets (as discussed in the previous chapter), is called **error backpropagation** (Rumelhart, Hinton, and Williams 1986), and it makes some additional demands on the biology that some influential people have argued are implausible (Crick 1989).

This problem has been my specific area of research for over 20 years, and my colleagues and I have developed progressively more biologically plausible models of how this error-driven learning process could work within the neocortex (O'Reilly 1996; O'Reilly et al. 2012). Our latest idea is that the brain is constantly making predictions about what will be seen next, at a rate of about 10 times per second (i.e., the *alpha* frequency), and very specific patterns of neural connectivity in the neocortex and thalamus provide a “ground truth” correct answer against which those predictions are compared. Thus, the difference between these predictions and what actually happens provides the error signals driving learning, and we have shown how these error signals, which exist as differences in the activity states of neurons over time, could drive learning in synapses throughout the neocortex. Furthermore, our computer models show that this form of learning can indeed acquire the kinds of sophisticated knowledge that people do, for example the ability to recognize different categories of objects (O'Reilly et al. 2020).

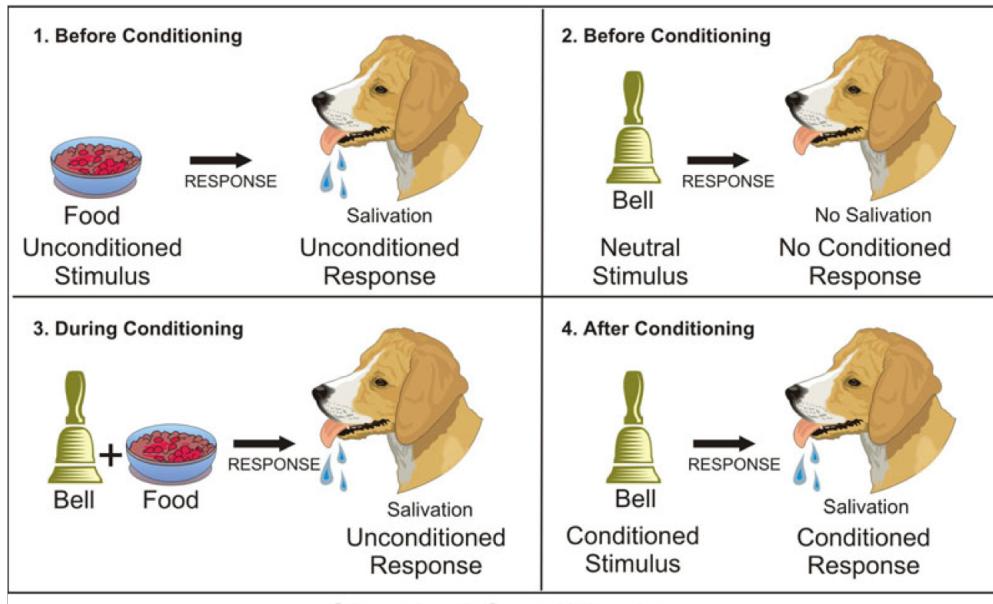
There is just one problem with all this: while our proposed synaptic plasticity mechanisms are consistent with the existing body of detailed knowledge, they also make a few extra demands that have not been tested empirically. So we do not yet know if this theory all goes through or not. Furthermore, there are various other different theories about how all of this could work, which make different, testable predictions. Thus, hopefully we'll get some answers in the not-too-distant future, and then we can potentially connect the dots all the way from the beautifully detailed biochemical level up to the high-level effects of these mechanisms in forming new knowledge representations within the neural networks of the neocortex. For now, we have to live with a bit of a hole in our overall understanding of this most important process of learning in the brain.

Before moving on to dopamine, another way of understanding learning taking place in the neocortex is in terms of **imitation learning**, also known as **observational learning**, where somehow we are able to observe other people's behavior, and then turn around and produce some approximation of that behavior ourselves. The popular phrase “monkey see, monkey do” suggests that this form of learning is widespread in primates, but the actual behavioral data across a range of species suggests that people are by far the most likely to engage in true imitation, while other species exhibit a range of socially-influenced learning that often falls short of direct imitation of actions (Carcea and Froemke 2019). This imitative capacity is closely related to *cultural transmission* of behavior across individuals, and the best non-human examples of this comes from chimpanzees who learn techniques for getting termites using sticks, or using moss as a sponge (Lamont et al. 2017).

Although imitation may sound relatively simple, upon closer examination, the process of turning the perception of behavior into your own motor program requires a highly sophisticated perceptual and motor control system. Thus, the fact that even young infants appear to be capable of this is quite remarkable (Meltzoff and Moore 1994; Ferrari et al. 2006). An important neural substrate for this form of learning has been found, in the form of **mirror neurons** that appear to achieve this feat of mapping observed behavior into the same patterns of neural firing that are active when you perform the same behavior (Iacoboni, Woods, and Rizzolatti 1999). However, it is not known how these neurons learn this mapping in the first place, so it remains a phenomenon in search of a deeper explanation. Nevertheless, there is an intriguing suggestion that these mirror neurons might be affected in autism spectrum disorders, which could potentially account for the difficulties in empathy in this population (Gallese, Keysers, and Rizzolatti 2004).

Dopamine-modulated Learning

Most introductory textbooks do not address any of the above topics in learning, and focus exclusively on the relatively well-understood domain of conditioning, which has been studied since the days of Pavlov and the behaviorist school in the early 1900's. This has become an area of renewed interest in neuroscience, since the discovery that dopamine activity almost perfectly accounts for the nature of these conditioning phenomena (Montague, Dayan, and Sejnowski 1996; Schultz, Dayan, and Montague 1997).



Classical Conditioning

Figure 5.3: The classical conditioning paradigm.

Classical (Pavlovian) Conditioning

The classical conditioning paradigm (Figure 5.3) centers around learning the connection between a previously *neutral* stimulus (the **conditioned stimulus** or **CS**) and a biologically-established, affectively significant *outcome*, known as the **unconditioned stimulus** or **US**. In the classic experiments by Pavlov, the ringing of a bell served as the CS, and food reward as the US, and the subjects were dogs, who learned over a few repetitions of the CS followed by the US to salivate after hearing the bell, in anticipation of receiving the food. The salivation is somewhat confusingly labeled the **un/conditioned response** (U/CR), where it is *un*-conditioned (UCR) prior to learning in response to the food US, and *conditioned* (CR) after learning in response to the CS. So, the same response has two different labels depending on what is driving it.

Ecologically, this simplified lab experiment is thought to capture the real-world learning about different stimuli that help us anticipate and prepare for important upcoming outcomes. For example, when you are hungry and driving down the highway on a road trip, the sight of a McDonald's sign alerts you to the availability of food there. Thus, the McDonald's sign is effectively a CS, and indeed this conditioning paradigm applies well to the goal of advertising, which is to establish a solid connection between a brand logo and desirable US outcomes. In many ads, the use of sexual imagery or famous faces directly activates the brain's reward pathways – they are literally replicating the classical conditioning paradigm to associate the CS (brand / logo) with the US driven by these rewarding stimuli.

Although Pavlov and the behaviorists were exclusively concerned with overt behavior such as salivation, we now know the internal biology that drives this form of learning. Figure 5.4 shows how dopamine neurons in the *ventral tegmental area (VTA)* of the brainstem reticular activating system respond in a classical conditioning experiment (Schultz 1986; Schultz, Dayan, and Montague 1997). When the US (labeled R = reward, a juice drop) is presented without any prior CS, dopamine responds with robust firing above its “tonic” steady base rate of firing. This is consistent with the naive idea that dopamine encodes raw reward signals.

However, when the very same reward is presented after a CS (which has been reliably paired with the reward in prior conditioning trials), *dopamine no longer fires to the reward!* Furthermore, when the CS is presented and the reward is *withheld*, dopamine neurons show a suppression or *dip* in firing below their tonic baseline. Psychologically, you would feel disappointed if you didn't get the reward you expected, and indeed that is exactly what the dopamine neurons are signaling.

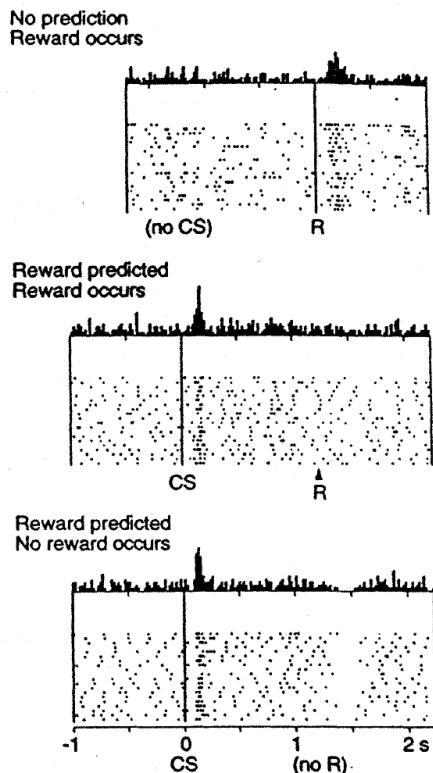


Figure 5.4: Dopamine neuron firing in a classical conditioning paradigm of CS followed by US (labeled R for reward – it was actually a drop of juice). Top: Unexpected rewards (at time point R) drive dopamine firing. Middle: Trained CS followed by R shows dopamine firing at the CS, but *not* for the reward. Bottom: Trained CS followed by *omission* of R shows reduction of firing at R. Each row of dots shows when a dopamine neuron fired a spike on a given recording trial, and the bars at the top show the accumulated histogram of all the spikes at that corresponding point in time across all such trials. Dopamine does *not* respond to raw reward input, because it fails to fire in when the reward is accurately predicted by the CS, in the middle panel. Furthermore, it directly signals “disappointment” by reducing dopamine firing when an expected reward is not received. These and many similar results show that dopamine responds to the *contrast* or difference between predicted and actual rewards. From Schultz et al, 1997.

These results, from the pioneering work of Wolfram Schultz and colleagues, have profound, far-reaching implications, and represent one of the most exciting and important findings in neuroscience. They are also one of the most important examples of the **contrast** principle, as we emphasized in the introductory chapter. Specifically, these results show that dopamine neurons respond to the contrast or difference between an expectation of reward, and what is actually received, *not* to the raw reward input itself. This contrast property of dopamine is what drives insatiable greed, dissatisfaction, and apathy, because once we learn to expect any given positive reward-like outcome, we no longer receive dopamine for it!

This property of dopamine is what causes kids to be so entitled and spoiled: they come to expect all that coddling from their overprotective parents, and all the excitement they get from playing video games, so that when they finally get out into the “real world”, it is all so difficult and filled with disappointing drudgery. It is also why so many famous people, especially rock stars it seems, turn to dopamine-activating drugs of abuse – once they adapt to their new amazing famous lifestyle, their dopamine system no longer gives them that amazing feeling of unexpected reward. Drugs like cocaine artificially bypass the expectation-driven contrast mechanisms of the dopamine system, so they continue to drive dopamine bursts. However, even here the system slowly adapts and more and more drug is required to achieve the same effects, so really there is no escape from the evil maw of the dopamine contrast effect!

From a hard-nosed learning theory perspective, there is a very good reason why the dopamine system must work in this contrast-based way: *learning is most efficient when it is focused on what is not yet learned.* Learning something you already know simply doesn’t make much sense. Thus, in the case of classical conditioning, continuing to learn about the fact that the CS predicts the reward after the system has already acquired this association isn’t very useful. And this logic suggests that dopamine is fundamentally a *learning* signal, not a reward signal. In particular, as we’ll see in the next section, dopamine directly affects learning in the basal ganglia and other brain areas, including the areas that are learning about the CS – US association in the first place. Thus, as dopamine stops firing at the time of the US (R, reward), it stops the further learning of this association (because it has already been learned).

This basic theory of how learning should function was systematized by Robert Rescorla and Allan Wagner in a seminal paper (Rescorla and Wagner 1972), where they proposed a simple mathematical “learning rule” that says that the amount of new learning should be proportional to the contrast or difference between the actual reward you receive and what you already expect the reward should be. This is also known as a **reward prediction error (RPE)**.

Roughly a decade later, Rich Sutton and Andy Barto published an important extension to this idea (Sutton and Barto 1981), known as the *temporal differences (TD)* learning rule, which can also account for the fact that dopamine learns to fire at the onset of the CS, even as it stops firing for the expected US. Furthermore, this work led to the development of many advanced mathematical techniques in a field collectively known as **reinforcement learning (RL)**, which is a branch of *machine learning* that deals specifically with learning from overall reward / punishment signals. These RL techniques have been used in many different AI technologies, and play a central role in the recent advances from the Google DeepMind group, in their models that learn to play Go and challenging video games (Silver et al. 2017; Mnih et al. 2015).

Although the TD learning rule provides an elegant and powerful mathematical description of classical conditioning, the brain networks actually involved in this form of learning are considerably more complex. Figure 5.5 (Mollick et al., n.d.; Hazy, Frank, and O'Reilly 2010) shows a summary diagram of these networks. Conditioning learning involves interactions between two “vertically” organized sub-systems, one involved in forming associations between CSs and USs and another that drives the contrast-with-expectations differences in the dopamine system. This first, associative system depends on nuclei within the *amygdala*, and the second depends on the ventral (bottom) areas of the striatum in the basal ganglia. This framework can account for a wide range of data about the biology and function of dopamine-driven learning in the brain, and, given the overall complexity of the system, the ability to simulate it all in a computer model is essential for understanding how it all works.

Extinction and Context in Conditioning

As was the case with LTD (long-term-depression), figuring out how the associations between CS and US are *unlearned* is just as important as figuring out how they are learned. This involves the phenomenon of

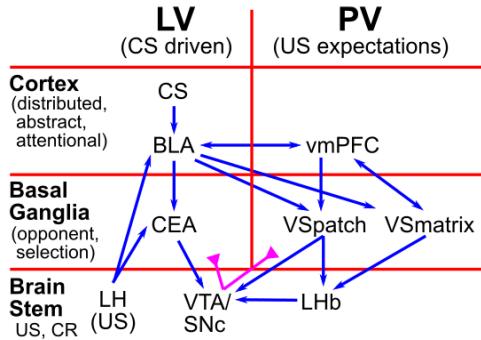


Figure 5.5: Biological systems involved in classical conditioning. There are two major learning systems, LV and PV (columns), each of which has cortex-like, basal-ganglia-like, and brainstem components (rows). LV (*learned value*) reflects the contributions of the amygdala to forming CS – US associations. PV is the ventral striatum, which drives reward-prediction-error (RPE) firing in the dopamine system (VTA / SNC). Together these constitute the PVLV system, named in honor of Pavlov. The BLA (basolateral amygdala) within the LV system learns to associate CS's with corresponding US's, while the CEA (central amygdala) reduces these higher-level associations down to specific “Go” vs “NoGo” signals in a basal-ganglia-like fashion, and directly drives dopamine firing and core behavioral responses (*conditioned responses*) appropriate for different US's. The PV system likewise has a cortical component in the ventral and medial areas of the prefrontal cortex (vmPFC), and a basal-ganglia component in the ventral striatum (VS). Dopamine firing in the VTA / SNC drives learning throughout all of these areas.

extinction, where the CS is repeatedly presented while withholding the US. From Figure 5.4, this should produce repeated *dips* in dopamine levels (much disappointment), which in turn should drive LTD in synaptic connections, causing the association between the CS and US to be unlearned.

While this all does occur, the situation turns out to be considerably more complicated, in ways that make sense ecologically. To make a long story short, the brain actually learns *new associations* during these extinction events, in addition to weakening (somewhat) the existing ones. These new associations effectively encode **context-specific exceptions** to the original association – e.g., “in this particular situation, you’re not going to get the food, but you might still get it in other situations”. Furthermore, the nature of this new learning is under top-down control from the ventral / medial frontal cortex, which can play a critical role in interpreting the nature of what is going on: has the world really changed, or is it just kind of random (Quirk and Mueller 2008; Gershman, Blei, and Niv 2010)?

The advantage of all this is that the initial CS – US association is relatively preserved, and especially if this was something learned through a painful, dangerous experience, it is probably a good idea to keep these memories around. Better safe than sorry. The disadvantage is evident in the phenomenon of PTSD (post-traumatic stress disorder), where traumatic memories cannot be extinguished, and keep intruding into normal life. There are significant individual differences in the extent of PTSD, and a major factor reflects the ability to exert top-down control and establish a strong new context to override the traumatic situation.

Another important consequence of this type of learning is that people will tend to hold onto these associations even in the face of repeated disconfirmation, explaining each new failure as another “special case” or circumstance – this is evident for example in doomsday cults, which respond to each absence of predicted doomsday by reinforcing their core beliefs, while attributing the failures to unforeseen contingencies (Boudry and Braeckman 2012).

In the lab, these extinction phenomena are observed in the phenomena of **spontaneous recovery**, **reinstatement**, and **renewal**, which are typically observed in aversive conditioning situations (i.e., when the US is a negative outcome, like getting shocked). Spontaneous recovery refers to the re-emergence of the CS – US association after extinction (typically after some time has passed), without any further training, clearly showing that extinction learning did not erase this original memory. Reinforcement occurs after a single US presentation without the prior CS, after which the CS – US association is reinstated. The US reactivates the associated memories and this is enough to overcome the extinction learning. Renewal is particularly revealing of the important role of context. In this case, the subject is conditioned in one environment (A) and

extinguished in a second, novel context (B). When put back into the original context (A), the original CS – US association is *immediately* effective without any further learning. In other words, the subject learned a context-specific exception (“when in context B, I won’t get shocked”) instead of unlearning the original association.

Operant / Instrumental Conditioning

Classical conditioning is the sensory front-end to the other major form of learning studied by the behaviorists: *operant* or *instrumental* conditioning. This form of learning occurs through the reinforcement or punishment of *actions*, instead of stimuli. The central idea is captured in **Thorndike’s law of effect**: *actions that lead to good outcomes are more likely to be taken, while those that lead to bad outcomes are less likely* (Thorndike 1911). This is so intuitive that it is difficult to imagine it being otherwise, but nevertheless, it captures a considerable amount of behavior in humans and animals.

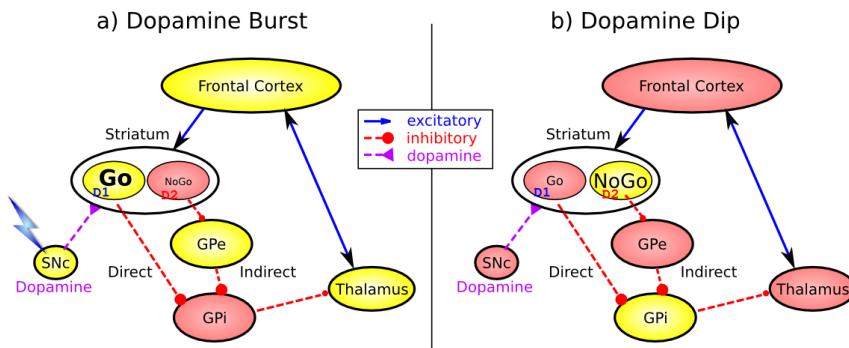


Figure 5.6: How increases in dopamine (bursts) and decreases in dopamine (dips) drive learning in opposing Go vs. NoGo pathways in the basal ganglia. Through the complicated basal ganglia circuitry, the firing of Go (aka direct pathway) neurons leads to a net excitation of motor plans in the frontal cortex. The NoGo pathway has the opposite effect, preventing the frontal activation that would otherwise occur from Go activation. When an action leads to an unexpected positive outcome, the resulting dopamine burst activates a special type of dopamine receptor (the D1 receptor), which drives LTP learning in the input synapses to the Go neurons. This makes those neurons more likely to fire again under similar circumstances, achieving Thorndike’s law of effect. The opposite happens when dopamine dips occur for unexpectedly bad outcomes, which interestingly has a net LTP effect on the NoGo neurons via D2 receptors, and an LTD effect on the Go neurons.

We now know how this type of learning works, in terms of dopamine’s effect on the basal ganglia (Figure 5.6) (Frank 2005; Gerfen and Surmeier 2011). Unexpected positive outcomes following a given action result in a burst of dopamine (as we saw in Figure 5.4), and this dopamine burst acts on D1 receptors located on the “Go” neurons of the basal ganglia to drive LTP of the synapses into the neurons that decided to trigger that action. Thus, these stronger synaptic inputs make it more likely that the same action will be triggered again in the future, when similar inputs are driving the basal ganglia (i.e., in similar situations), thereby achieving Thorndike’s law of effect.

The opposite pattern of changes occurs when unexpectedly bad outcomes arise, which drive dips in dopamine firing, and end up strengthening inputs to the “NoGo” neurons that compete against the Go pathway and prevent an action from being triggered. Thus, actions that lead to bad outcomes are less likely to be triggered, consistent with the other half of Thorndike’s law of effect.

The overall relationship between dopamine and the basal ganglia is summarized in Figure 5.7, where classical conditioning processes train the **critic** what kinds of rewards or punishments to expect, and the resulting differences between these expectations and actual rewards / punishments, reflected in the dopamine signal, then drives learning in the **actor** (basal ganglia). This image of dopamine as a critic fits with our overall conception of the *contrast* nature of this signal: it is never satisfied and quick to criticize, just like a critic. Of course, the poor actor has to do all the hard work of coming up with stuff for the critic to critique, but, as you may have experienced, it is often hard to be properly critical of your own behavior, whereas it is much easier to see what is wrong with other people. Thus, separating the critic and actor components in the

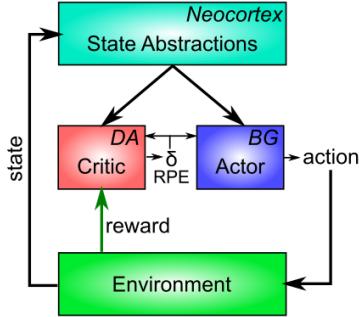


Figure 5.7: Actor-Critic schematic for the relationship between the dopamine (DA) signal driven by principles of classical conditioning (the critic), and action decisions triggered in the basal ganglia (BG) (the actor). The environment state is represented by various abstract, higher-level *compressed* representations in neocortex, which feeds into both the critic and actor. The actor decides on actions to take as a function of these neocortical inputs, and the critic generates predictions about the kinds of US outcomes that are likely to result. Learning in both the critic and actor is a function of the dopamine signal, which is symbolized as a delta, or *reward prediction error* (RPE). Thus, classical and operant / instrumental conditioning are connected through this actor – critic relationship.

brain makes sense, and is another example where two fully interdependent systems can nevertheless be seen as performing distinct functions.

Partial Reinforcement, Gambling, and Shaping

One of the topics that the behaviorists explored extensively was *reinforcement schedules* – different rates and patterns of delivering rewards. The most interesting and relevant finding from this work is that **partial reinforcement** can have surprisingly strong effects compared more reliable reinforcement schedules. In a partial reinforcement schedule, rewards are only delivered randomly on a fraction of successful action trials. In effect, it is just like gambling, where there is a relatively infrequent, random payout. The net effect of this is to confound the critic system, which can no longer accurately predict what kind of outcome to expect. Therefore, when a positive reward is received, it is not *discounted* like would have been if it was perfectly predictive. You will get that burst of dopamine for the reward! This is why gambling can be so addictive – it works just like addictive drugs in disabling the stingy, harsh dopamine critic.

Another important discovery in instrumental conditioning was that more complex behaviors can be built up from simpler elements through the process of **shaping**. This is the technique used to get circus animals to perform their complex tricks, for example, and is often used in scientific research with animals to study more difficult cognitive tasks.

Finally, it is important to recognize the difference between a **primary** vs. **secondary reinforcer**. A primary reinforcer directly satisfies a biological need (e.g., food or water), while a secondary reinforcer is indirect, and must be learned. Money, points, and gold stars are common examples of secondary reinforcers, which are effective for motivating people to do things. Interestingly, animals typically require primary reinforcers, but people readily learn to value secondary reinforcers. This ability to value initially arbitrary stimuli is essential for modern economic life – it would be rather inconvenient to have to directly exchange food, water, or other items of direct value.

Motivation

Despite the satisfying modern synthesis between dopamine and the behaviorist-era conditioning phenomena, this overall view of behavior focuses almost entirely on **external / extrinsic** factors (reward / punishment) to the exclusion of **internal / intrinsic** factors such as goals, drives, desires, etc. This is consistent with the behaviorist-era prohibition on considering internal factors more generally, but we should have no such constraints on our modern thinking about this topic. Nevertheless, the current neuroscience-based research still carries some of this extrinsic bias, with the central role of internal factors having been somewhat less emphasized. By contrast, researchers in the field of social psychology have a long tradition of thinking about

the central role of goals, desires, emotions and mood on behavior.

Before exploring some of these ideas, it is interesting to ponder the state of mind of a behaviorist from the 1920's: did they really think that their *own* personal behavior was fully determined by external rewards and punishments? Were they not aware of having internal goals that drove them to torture rats for long hours, day after day, in pursuit of such ineffable, remote rewards as scientific understanding and a chance of prestige and fame? The tangible rewards associated with scientific research are sufficiently distant and improbable, while the immediate working conditions involve relative poverty and extreme hard work, that it is really hard to understand why people would do such a thing in terms of purely external rewards. Instead, there must be some significant long-term internal forces driving such "crazy" pursuits, which are evident across many domains of human endeavor.



Figure 5.8: Drive reduction theory according to Hull, 1943. Basic needs create drives when those needs are not satisfied, and behavior is then recruited to satisfy those drives.

One of the few types of internal state that behaviorists did consider was the notion of a *drive* or state of internal discomfort (e.g., due to lack of food or water) that then motivates behavior toward reducing that discomforting state (Hull 1943) (Figure 5.8). But this **drive reduction** theory has trouble accounting for motivations such as our desire to learn and work, which don't really seem to be associated with discomfort-reduction processes.



Figure 5.9: Maslow's hierarchy of needs. Higher-level needs are only considered once lower-level ones are satisfied.

A more comprehensive theory of motivation was developed by *Abraham Maslow*, at around the same time as Hull (Maslow 1943). Maslow's **hierarchy of needs** (Figure 5.9) captures the intuitive idea that

higher-level needs are not relevant unless the more basic needs essential for survival are satisfied. The two lowest levels in the hierarchy are physiological needs (breathing, food, water, etc) and safety. Once those are satisfied, then higher-level needs such as love and belonging and esteem become relevant. Finally, at the highest level, Maslow put *self actualization*, which includes things like morality, creativity, and lack of prejudice. Interestingly, this highest level resembles the Buddhist notions of enlightenment, where one transcends lower-level attachments and needs, and can act in a more principled, rational, and yet spontaneous manner. These frameworks capture the subjective feeling that we are controlled by our basic needs, and yet we yearn to be free from these low-level demands.

One problem with Maslow's theory, shared with any theory that attempts to articulate universal features of human behavior, is that people are rarely so compliant, and regularly violate his strict hierarchy. For example, people have been known to literally work themselves to death, including recent cases of video gamers playing to death as a result of neglecting basic bodily needs. Furthermore, teenagers routinely risk their personal safety in order to show off and otherwise enhance their social belonging and perceived self-esteem. Nevertheless, as a general tendency, the hierarchy makes sense, and certainly the numerous cases of cannibalism in the face of extreme hunger suggest the power of these more basic physiological needs.

Goal-driven Behavior

A more general motivational framework is based on the notion of *goals* and the idea that people are specifically motivated to achieve their goals. These goals can be highly diverse in their specifics, but they share the common property of delivering positive reward signals upon goal completion (e.g., “the satisfaction of a job well done”), or even progress toward goal completion (“almost there, just around the corner.”), and corresponding negative states associated with failure (disappointment, embarrassment, lack of self-esteem). Many aspects of goal-driven behavior have been studied over the years (Tolman 1948; Miller, Galanter, and Pribram 1960; Powers 1973; Klinger 1975; Gollwitzer 1993; Carver and Scheier 1990).

In the animal behavioral tradition, goal-driven behavior has been studied in the context of paradigms such as **satiety** and **devaluation** (Balleine and Dickinson 1998). In these cases, an animal is instrumentally conditioned to press one lever for food while in a state of hunger, and is then given as much food as they want. They are then put back into the box with the lever – if behavior is driven by purely *habitual* stimulus – response associations, they should push the lever even if they are no longer hungry. However, if they are actually thinking about the outcome produced by pressing the lever, and recognizing that they don't want that outcome, then they should not press the lever. Interestingly, results show that damage to the ventral and medial areas of prefrontal cortex (**vmPFC**) cause rats to press the lever even when they are full. The same kinds of results have been shown in the devaluation studies, where the food is subsequently paired with a bitter taste outside of the lever-pressing context, so that it is no longer desirable. If the animal still presses the lever, then they aren't clearly representing the outcome of the lever press.

The importance of the vmPFC brain areas for goal-driven cognition is consistent with neural data showing that these areas (in particular the *orbital frontal cortex, OFC* and *anterior cingulate cortex, ACC*) have many neurons that anticipate the possible US outcomes associated with a given situation and actions taken within that context, and impair goal-directed behavior when damaged (Rudebeck et al. 2006; Wallis and Kennerley 2011). More generally, this is consistent with the overall role of the prefrontal cortex in driving goal-driven controlled behavior, which requires the tight coordination between plans and their potential outcomes in order to decide on the plans that will lead to the most desirable potential outcomes.

As discussed in the Neuroscience chapter, these vmPFC areas are directly interconnected with the basal ganglia, amygdala, and dopamine system, forming the overall *control* and *decision-making* system of the brain, and each of these areas plays a critical role in supporting the overall emergent ability to behave in a goal-driven, controlled manner. Furthermore, as we'll see in the clinical disorders chapter, these are the brain systems that are implicated in most of the major clinical disorders.

Finally, one of the most fascinating and important demonstrations of the importance of intrinsic motivation comes from studies showing that giving people extrinsic rewards can actually *undermine* intrinsic motivation (Deci, Koestner, and Ryan 2000)! For example, giving kids awards for drawing actually caused them to draw less than kids who did not receive these awards. These results are controversial, however, and systematic reviews of the literature have reached opposite conclusions (Cameron, Banko, and Pierce 2001).

One of the most important factors appears to be whether the task in question is actually reasonably strongly intrinsically motivating in the first place: there is stronger evidence of the undermining effect when the task has stronger intrinsic interest, compared to more “boring” tasks, for which external rewards might be useful.

Emotion and Arousal



Figure 5.11: Valence (positive vs. negative) vs. arousal (high vs. low activation) *circumplex model*.

The fact that the same brain areas involved in goal-driven motivated behavior are also the primary areas associated with emotion raises the important question as to the relationship between emotion and motivation. It is somewhat difficult to provide a crisp, principled definition of *emotion*, which thus makes it difficult to arrive at a clear understanding of its relationship with motivational states. Some widely-recognized properties of emotion are that it has some kind of distinctive, characteristic subjective feeling, is associated with physiological arousal at least to some extent, and that it drives associated behavioral responses. It is also generally agreed that emotion should be biologically grounded, at least in the more “primitive” or basic level of emotions.

All of these properties are consistent with the idea that emotional states and motivations have a strong connection (Cardinal et al. 2002). For example, one’s overall *happiness* is most strongly associated with feelings of personal self-efficacy and control (along with interpersonal connectedness and belonging). Likewise, feelings of *sadness* are strongly associated with disappointment, failure, and lack of control. Thus, a simple overall hypothesis is that emotions are the subjective states associated with our core motivational systems. Let’s see how far this idea can take us, in understanding the full spectrum of emotional states.

Figure 5.11 shows the simplest standard model of emotion, known as the **circumplex model**, which distinguishes between two separate dimensions of *valence* vs. *arousal*. Valence refers to the “sign” of the emotion, positive vs. negative, while arousal refers to the intensity of the emotion. Anger and exhilaration are opposite valences but the same high level of arousal. These two valences are also associated with opposing *approach* vs. *avoid* behavioral orientations, which have been identified as core opponent aspects of emotional / motivational states and corresponding personality dimensions (Carver and White 1994; Read et al. 2010).

While this simple framework captures the most essential dimensions of affective / emotional states, it is likely that the valence aspect of emotion is considerably more complex than the *bivalent* (two valences) nature of the circumplex model. For example, *Paul Ekman* found that there are 6 **basic emotions** that have clearly recognizable facial expressions, which are universal across cultures: anger, disgust, fear, happiness, sadness, and surprise (Ekman and Friesen 1976). Later work added other emotions based on vocal and facial



Figure 5.12: Six different basic emotions as represented by facial expressions: anger, disgust, fear, happiness, sadness, and surprise.

expressions, and Plutchik proposed a systematic wheel of emotions based on 8 emotion categories arranged in opponent pairs, with an arousal dimension as well (Figure 5.13).

In addition to the basic happy / sad elements of emotion, which may be more closely related to goal-driven motivational states, some of these other emotional states are more clearly social in nature, and the role of facial expressions and vocalization clearly implicates a strong social communication role for emotions. Thus, we can potentially organize emotional states in terms of a set of distinct functional domains, where these states can serve to motivate people toward appropriate patterns of behavior. We roughly organize these according to Maslow's hierarchy of needs, with slightly different groupings.

- **Physiological States:** **Hunger**, **thirst**, **pain**, **tiredness**, **lust**, and the **need to excrete** are all basic motivational states associated with core body functions necessary for survival, and correspond with the physiological level in Maslow's hierarchy. These may not be considered "emotional" states per se, but they share the same properties of being strongly biologically determined, varying in level of intensity or arousal, and capable of driving appropriate behaviors to mitigate negative states and approach positive ones. While hunger and thirst may not typically need to be communicated socially using basic facial expressions, tiredness and lust likely do, and have clear social cues in the form of yawning and flirting behavior.
- **Safety states:** **Fear** is the emotional correlate of Maslow's safety level, and has both a direct internal motivational role (driving you to avoid scary situations), and an important social communication role for alerting others of potentially dangerous situations, which is facilitated by the presence of the unique fear facial expression. **Disgust** is an interesting case which may be more strongly social in nature: it is important to communicate to others that food might be rotten and disgusting, and it seems likely that this original function has been extended to apply to labeling the behaviors of others in the group as dangerous or otherwise something to be avoided. **Hate** is an emotional state that is also clearly negative and social in nature, and associated with disgust: it is the emotional state and social communication associated with labeling others as belonging to the out-group.
- **Social states:** **Love** is the opposite of hate, and is the positive social affective state associated with

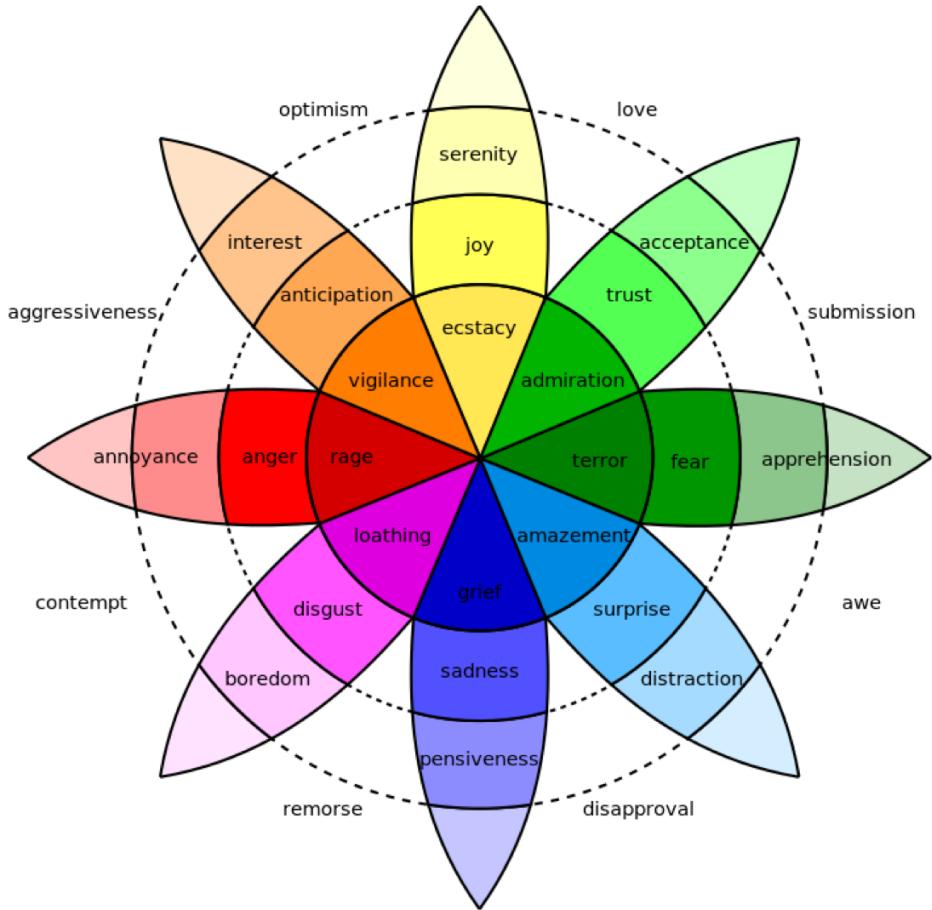


Figure 5.13: Plutchik's wheel of emotions, with arousal (intensity) represented as distance from the center along any of 8 different categories of opponent emotions.

members of the in-group. It obviously corresponds with Maslow's 3rd level of love and belonging. Other important social states include **dominance** and **submission** dynamics, along with trust and admiration, which have to do with establishing and perceiving relative status within the social order. These correspond with Maslow's esteem level, and are very strong and often-overlooked motivational states for social beings, from dogs to monkeys to humans. We do not necessarily have clear terms for these as emotional states (e.g., the feeling of being dominated by, or of dominating, a social other), but there is evidence that they are important factors in personality and interpersonal interactions (Hopwood et al. 2013), and certainly we have terms such as "diss" = disrespect and "pissing contest" that refer to such interactions.

- **Goal-associated states:** many of the remaining states are associated with goal-driven behavior, including **happiness** and **sadness** (and their varying levels of intensity or arousal) as noted above, but also **anger** and **frustration** which are associated with impediments to progress toward achieving one's goals, and **curiosity**, **interest**, and **surprise** which are associated with recognition of new interesting avenues to pursue. **Boredom**, **distraction**, **optimism**, and **anticipation** are also other states that clearly seem to be goal-related. **Grief** and loss are perhaps not so obviously goal-related, but in some ways they reflect a profound disruption of one's sense of overall control and order in the universe, in addition to the basic feelings of missing a loved one.

Thus, overall, it does seem that emotional states can generally be understood as corresponding to biologically-determined motivational states, which provides a clear functional story for why we have emotional states in the first place (Cardinal et al. 2002). As such, this raises important questions about the standard "Hollywood" story about the special status of emotion as a unique aspect of human beings. Under this motivational framework, many of our emotional states are common across all mammals at least, and represent a genetically-coded, low-level aspect of our brains, not something special and unique about humans. On the other hand, because our emotional states are so strongly felt, and provide dramatic color to our lives, we regard them as special.

Also, emotion is what keeps us from harming each other (except when it is what drives us to harm each other, in the case of hate and anger), and the lack of basic emotional connections in psychopaths enables them to do horrible things that "normal" people would never do. So, from a survival perspective, we really depend on everyone sharing these protective emotional responses, and anything that doesn't is immediately scary and foreign. Furthermore, we do have a large portion of our vmPFC devoted to emotional processing, and these emotional representations are likely novel combinations of more basic, lower-level emotional states, shaped over our personal histories, and thus likely provide a much richer and elaborated emotional tapestry than found in other animals.

Emotional / Motivational Encoding in vmPFC

Figure 5.14 shows a map of what the vmPFC emotional / motivational tapestry might look like, based on tracing the inputs and outputs of these areas relative to lower-level emotional and motivational areas in subcortical areas (Ongür and Price 2000). Consistent with the circumplex model (Figure 5.11), there are separable areas for positive vs. negative valence, and arousal. Interestingly, the negative valence area, known as area 25 or subgenual ACC, has been implicated in major depressive disorder through the work of *Helen Mayberg* and colleagues, and electrical stimulation in this area is a promising treatment (Riva-Posse et al. 2014).

Also, consistent with the anatomical principles from the Neuroscience chapter, these areas relate to nearby areas in terms of the ACC areas at the top relating to motor plans coded in surrounding PFC areas, and OFC areas toward the bottom being driven by visual, olfactory, taste, and visceral inputs coded in nearby areas. Thus, we can think of ACC as being more associated with action planning, including things like effort and difficulty costs, while OFC is more important for representing outcomes in terms of their relevant sensory features (taste, appearance, etc).

Biological Grounding of Emotion and Arousal

Finally, there is a somewhat strange history of thinking about emotion that is typically emphasized in introductory textbooks, and seems to reflect the desire to understand emotional states as special, biologically-

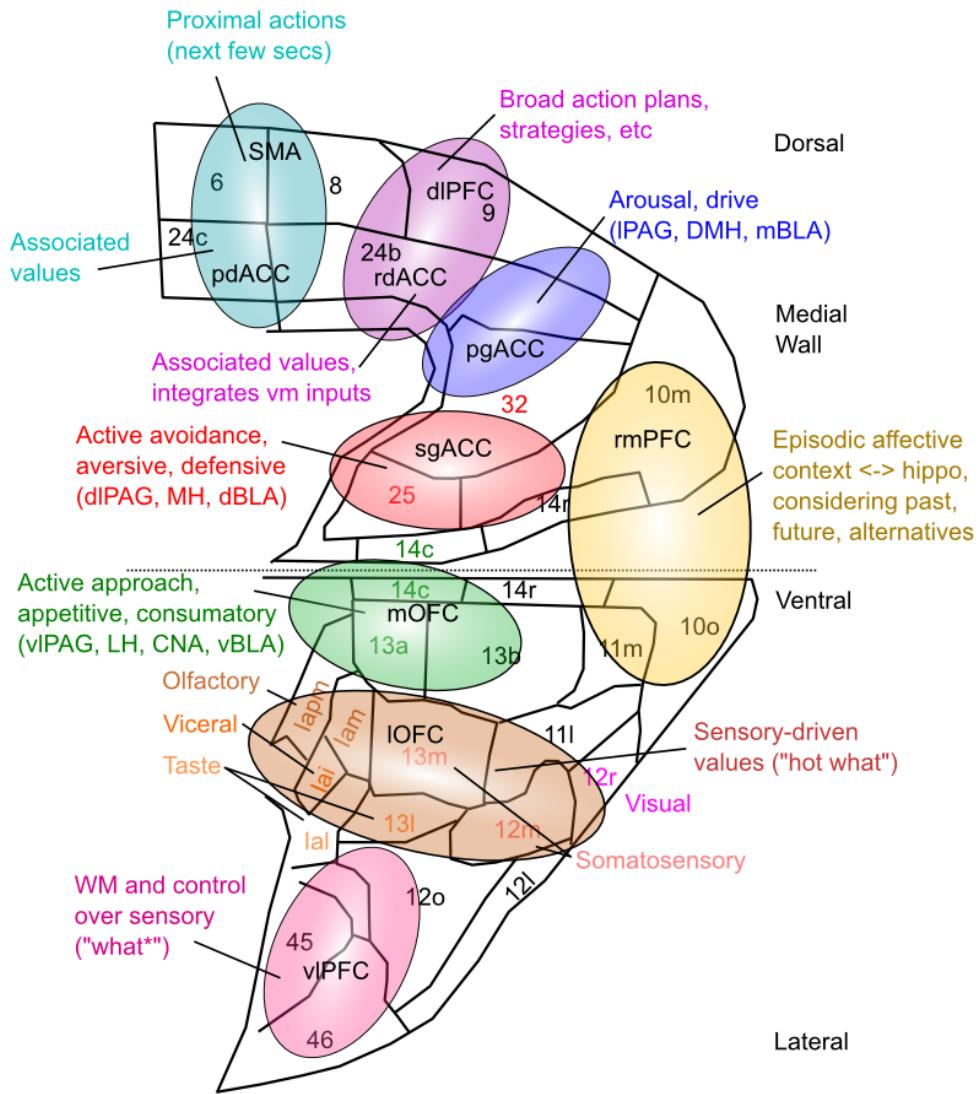


Figure 5.14: Map of ventral / medial frontal cortex (vmPFC) areas and their associated roles in emotional / motivational states, as a function of connectivity with subcortical areas that have established emotional / motivational valences. The broad organization is consistent with the circumplex model, with separate positive (appetitive) and negative (aversive) areas, and a separate arousal area, along with other forms of specialization. BLA = basolateral amygdala; CNA = central amygdala; PAG = periaqueductal grey; LH / MH / DMH = lateral / medial / dorsomedial hypothalamus.

grounded, important states. Specifically, William James and Carl Lange each independently proposed that emotion arises first in our bodily responses such as sweating, heart racing, etc, and is only later recognized as an emotional response as a direct result of these initial physiological responses. In contrast to this **James-Lange** theory, Walter Cannon and Phillip Bard proposed that higher-level processes in the brain play a critical role in driving our emotional experiences. Finally, Stanley Schacter and Jerome Singer argued that both physiological and higher-level interpretational processes were both essential, with their **two-factor theory**.

Ultimately, all of these theories still emphasize that emotional states have both physiological and higher-level interpretational aspects, and the unique, interesting aspect of emotion is that it can activate the body in ways that purely abstract mental states do not. From a modern perspective, it is clear that many different brain and body responses occur essentially in parallel, producing our rich, complex, and fascinating subjective experiences of emotion.

One final issue concerns the optimal level of arousal for driving motivated behavior. The **Yerkes-Dodson law** (Yerkes and Dodson 1908) established the principle that there is an optimal level of arousal somewhere in the middle between low and high levels, following an **inverse-U-shape** curve. This same curve has been found for levels of dopamine as well. You may have experienced this in experimenting with different levels of caffeine – too much is actually not productive, as you get too hyper and unable to focus.

Summary of Key Terms

This is a checklist of key terms / concepts that you should know about from this chapter. As we'll learn in the memory chapter, it is a great idea to test yourself on what was said about each of these terms, and then go back and double-check – that provides both beneficial repetition and also the *testing effect*.

- Synaptic Plasticity
 - Hebbian Learning
 - Long Term Potentiation (LTP)
 - Long Term Depression (LTD)
 - 5 Steps of NMDA / Ca++ synaptic plasticity
 - What determines LTP vs. LTD direction of synaptic plasticity?
 - Error backpropagation
- Classical Conditioning
 - Conditioned stimulus (CS)
 - Unconditioned stimulus (US)
 - Un/conditioned response (U/CR)
 - Dopamine responses to CS, R, no-R in conditioning expt
 - Rescorla-Wagner learning rule / reward prediction error model of dopamine
 - Extinction, and its context sensitivity: spontaneous recovery, reinstatement, renewal
- Operant / Instrumental Conditioning
 - Thorndike's law of effect
 - Implementation thereof in terms of dopamine effects on Go / NoGo
 - Actor / Critic model
 - Partial reinforcement and gambling
 - Shaping to build up complex behaviors
 - Primary and secondary reinforcers
- Motivation
 - External (extrinsic) vs. internal (intrinsic) motivation
 - Drive reduction
 - Maslow's hierarchy of needs (levels of hierarchy)
 - role of vmPFC in satiety / devaluation effects
- Emotion and Arousal
 - Circumplex model
 - Six basic emotions according to Ekman's original faces
 - Relationship between emotion and motivation

- Importance of both physiological and higher-level interpretations for emotion
- Yerkes-Dodson law

Chapter 6: Memory

Memory is the direct product of learning, so everything we learned in the previous chapter will help us understand how memory works. If you can remember it, of course. Some of the major questions that have been the focus of memory research include:

- What different kinds of memory are there?
 - Are there specialized brain areas for different kinds of memory?
- How long does memory last (for each different type)?
- What factors determine how well memories are encoded and recalled?

Thus, the study of memory has been focused on fairly practical and descriptive questions, befitting the essential role that memory plays in everyday life (and especially for students). Our memories are also a core aspect of our sense of self, and movies such as *Total Recall* (based on a Philip K. Dick short story, as so many good movies are) have explored this function of memories in provocative and interesting ways. Furthermore, by now most people have heard about the profound amnesia caused by damage to the *hippocampus*, e.g., from the famous case of Henry Molaison (H.M.) who had his hippocampus surgically removed and lost the ability to form new memories for the rest of his long, exceptionally well-studied life. The movie *Memento* artistically and accurately captures the subjective nature of this condition, and is required viewing for anyone interested in memory (I personally have *two* copies, each a gift to my wife – memory certainly can be fallible). What makes the hippocampus so important for memory? What kinds of memory do *not* depend on the hippocampus? These are some of the important questions we will address in this chapter.

From Synapses to Memory

If memory is the direct product of learning, and learning is the direct product of synaptic plasticity as we learned in the previous chapter, then in principle *memory should be found in every synapse in the brain*. In fact, this is *true*, but it is also true that some synapses are more important than others. A deep understanding of memory requires reconciling these two perspectives on memory, and integrating some additional properties of neurons beyond their synapses.

First, it is useful to compare and contrast the nature of memory in the brain with memory in a computer. More generally, it is tempting to try to use the computer as a model for the brain, as was especially popular in the early days of cognitive psychology, but it turns out that the brain doesn't work anything like a computer at the hardware level. Nevertheless, some features of computers do emerge out of the brain, despite the fundamental differences in their underlying hardware – this may make the computer analogy more confusing than not, but given the prevalence of these computer analogies, it is important to get this all straight!

In a computer, there are two major types of memory: RAM (random access memory), and a “hard disk”, which these days is typically just a different kind of solid-state chip, rather than the spinning platters used in actual hard drives, but it still plays the same functional role. RAM is where *active* memories reside – the stuff the computer is currently working on. Elements from RAM can be very quickly read into the central processing unit (CPU), processed, and then written back into RAM, often many times (e.g., when you are editing a document in your word processor, those words live in RAM, and are accessed many times to redraw the screen as you scroll and edit). When you are done working, you save the memories from RAM to the hard drive, where they can reside essentially permanently. If the power goes off before you save, the RAM is lost – it is active (fast) and *temporary*, whereas the hard disk is slower but permanent. Computers typically have much less RAM than hard disk storage – the amount of information needed for active processing is typically much less than the sum total of all information that has been processed and stored.

You have likely had the experience of suddenly forgetting what you were just talking about, or entering a room with a clear purpose, which then just vanishes into thin air. These seem like distinctly RAM-like properties: temporary, and associated with what you were currently thinking about / working on, and more capacity-limited – you can't juggle too many things at the same time. On the other hand, we clearly have access to a vast storehouse of long-term memories, which often require some specific cognitive effort to retrieve (and as you get older, this effort seems to grow), and clearly cannot all be active at the same time – this seems like a hard disk type of memory.

Thus, this distinction between fast, temporary, limited memory (RAM) vs. slower, permanent, high-

capacity memory (hard disk) somehow applies in the brain too, except the brain doesn't have a CPU, nor does it have discrete hardware modules like RAM and a hard disk. As we saw in the previous chapters, processing in the brain is *distributed* across all of the billions of neurons in the brain, with each neuron playing a small role within larger networks. Each neuron is detecting some particular patterns as a function of its synaptic connections, helping to compress and simplify the vast stream of information flowing through the brain. Learning operates by changing these synapses (i.e., long-term potentiation, LTP, and long-term depression, LTD), such that memory is also fully distributed across the brain, instead of being concentrated in a separate device like a hard drive.

The long-term nature of synaptic learning, and the vast numbers of such synapses, provides a nice fit with the hard-disk like properties, but what could be the neural equivalent of RAM? Without a CPU, there is no need for quickly reading and writing information from a RAM-like memory system. Instead, everything the neuron needs to carry out its detection and compression function is right there in its synapses, and learning directly modifies these synaptic connections.

The answer is **neural activity** – the ongoing spiking activity of the vast numbers of neurons in your brain that are currently above-threshold and sending their signals to other neurons. Indeed, this neural activity is an essential additional contributor to memory, because once excited, this activity tends to persist over time – and persistence is the essence of memory. However, unlike synaptic changes, and like RAM, neural activity is definitely transient – once some new pattern of activity sweeps over your neurons, whatever was there before is effectively lost (unless it somehow got recorded via synaptic changes). Furthermore, it is more limited in capacity – you have many fewer neurons than synapses, and because all these neurons are bidirectionally communicating with each other, the actual number of distinct, coherent memory states is drastically smaller than the number of neurons (only about 4, as we'll see below).

In summary, we'll start our investigation of memory with the following principles derived from neuroscience:

- Memory can be broadly defined as *any* form of persistence of information over time in the brain – any trace of some prior event can be considered a type of memory.
- Neurons have two primary sources of such persistent information:
 - **Activity** in the form of ongoing spiking, electrical potentials underlying that spiking, and the chemical states of other parts of the neuron, which are *transient* – once a neuron stops firing and its other electrical and chemical states dissipate, a memory trace is no longer actively present in that neuron.
 - **Synaptic changes** from learning, which are relatively *long-lasting*, and change what kinds of input signals will activate the neuron in the future (i.e., what it *detects*).
 - These two aspects of neural memory directly influence each other, because learning is driven by neural activity, and changes in synapses result in different patterns of neural activity. Despite this interdependence, these different types of memory have different functional properties and can be usefully distinguished.
- The specific **content** of the memory supported by any given neuron and its synapses is a direct function of its role within the larger neural networks of the brain – memory happens everywhere in the brain at all times, directly within content-specific processing areas (e.g., visual memories in visual cortex, etc). In addition, there are two brain areas that play an outsized role in memory, due to their specific neural properties and location within the brain's networks – both of these areas are situated at the *top* of the overall cortical hierarchy, so they have broad *access* and broad *influence* over everything going on in the brain:
 - The **hippocampus** is specialized for performing very *fast* encoding of synaptic changes in a way that avoids the massive *interference* effects that such synaptic changes would otherwise cause in other brain areas – this enables it to rapidly encode **episodic memory** of ongoing daily events, as H.M.'s case demonstrated.
 - The **frontal cortex / basal ganglia** system is uniquely capable of sustaining patterns of neural activity over longer durations and in the face of other distractions, supporting **working memory**, which is important for maintaining the current information that you happen to be “working” on. The relative deactivation of the frontal cortex during REM sleep demonstrates what cognition would

be like without this robust form of active memory – you would become much more distractable and unable to stay focused on a given task. This is characteristic of people with damage to the frontal cortex, and to some extent in people with ADHD (though their basic frontal function is typically indistinguishable from people without ADHD, as we'll explore in the Disorders chapter).



Figure 6.1: The *telephone* game, which is an apt metaphor for how information is transferred in the brain. Neurons, like people, take a given signal, interpret it in their own particular way (as a function of their synapses), and send out their own interpretation. The idea of direct symbolic information transfer as in a digital computer does not apply in the brain.

Finally, there is one more critical constraint from neuroscience, having to do with the widely-used concept of *transferring* information from one part of the brain to another. As noted above, this is how everything works in a computer (information is constantly being transferred among the different components of RAM, CPU, and hard disk), but information in the brain is not encoded *symbolically* as it is in a computer, and therefore cannot be so easily moved around. Instead, as we've emphasized repeatedly, each neuron has learned to detect patterns of activity in its inputs, and thus information can only be transferred by neurons in another brain area detecting their own version of the information encoded in a given brain area.

In other words, information transfer in the brain is much more like the game of *telephone*, where a given message is passed from one person to another, often resulting in amusing misunderstandings (Figure 6.1). The same thing happens in the brain: information transfer is *always* accompanied by fundamental transformations of the content, with each area adding its own *spin* or interpretation, with important consequences for understanding the relative veracity of memory.

The Modal Model of Memory

Figure 6.2 summarizes the **modal model** of memory, which is so-named because it summarizes the common elements of many different models of human memory that had been developed in the early part of the cognitive revolution (Atkinson and Shiffrin 1968). It does a good job of capturing many different phenomenological aspects of memory, and we can use it to see how the above neural principles play out in practice. It involves three separable components, *sensory memory*, *short-term memory (STM)*, and *long-term memory (LTM)*, with information flowing from one to the next dependent on relevant active processes including *attention* (sensory memory to STM) and *encoding* (STM to LTM).

First, sensory input activates **sensory memory**, which is characterized as a transient, high-capacity memory system that represents the sensory input at various levels of abstraction. Sensory memory corresponds largely to the *activity* of neurons that have been stimulated by the sensory input, at various levels along the

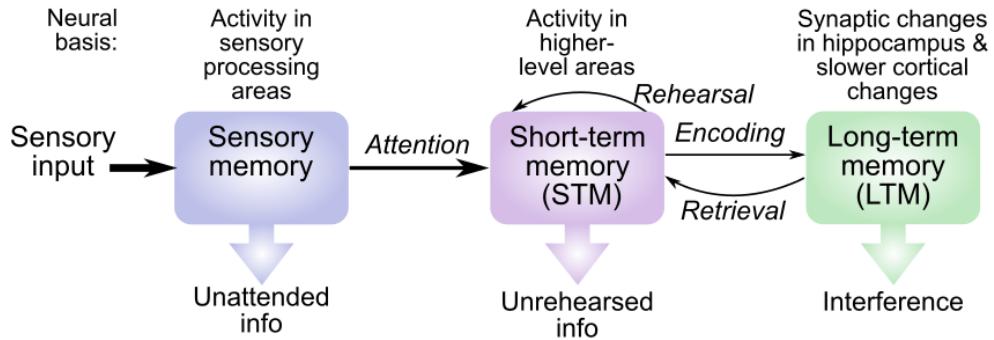


Figure 6.2: The modal model of memory and its neural basis, in terms of neural activity and synaptic changes. At each step, processing is required to transition to the next: only attended sensory items enter STM (and the rest is lost), and actively encoded STM information enters LTM. Active rehearsal sustains information in STM. Information in LTM can be retrieved back into STM, and is lost primarily via interference.

kinds of hierarchically-organized sensory processing pathways discussed in previous chapters.

There are different names for this activity within each modality, including **iconic** memory in the visual pathways, and **echoic** memory in the auditory pathway. Iconic memory generally persists for less than a second, whereas echoic memory lasts longer, up to about 4 seconds. These differences reflect the extent to which the neural activity in associated visual or auditory brain areas can persist. Because auditory information is inherently transient and evolving over time, the brain has extensive subcortical mechanisms that integrate and preserve these auditory signals over time, resulting in its longer persistence.

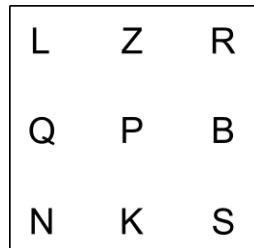


Figure 6.3: Sperling's sensory memory task. In the *full report* condition, participants attempted to retrieve all items, and typically only recalled about 4.5 on average. In the *partial report* condition, an auditory cue presented after a variable delay indicated which of the three rows to recall. For delays less than a second, they could accurately recall the letters within the cued row, indicating the presence of a high-capacity sensory memory trace (iconic memory) that decays within a second if not activated into STM via attention.

Classic experiments by George Sperling and others established these duration values, by flashing a display with 3 rows of 3 letters each (Figure 6.3), and probing people to report a particular row from the display after variable delays (Sperling 1960). In this *partial report* condition, people were generally able to report the information within about a second, but not longer. Critically, the relatively large amount of information in the full display was above people's capacity to encode in its entirety (as established through other *full report* conditions where they had to try to recall all of the letters), so the partial report cue allowed them to focus attention on one row, resulting in the activation of corresponding representations in STM. However, once the sensory memory trace fades, it is gone, and cannot be "transferred" to STM.

Experiments such as these also established the next step of the modal model, which is more strongly capacity-limited, but longer-lasting, and is referred to as **short-term memory (STM)**. Only information within the focus of **attention** makes the jump from iconic or echoic sensory memory into STM, and given the capacity constraints, attention can only grab about 3-4 "items" into STM from sensory memory (corresponding to a single row from the Sperling task). From a neural perspective, STM corresponds to neural activity in higher levels of the neocortex (in temporal and parietal lobes) that have more highly compressed encodings of

the sensory input. Thus, as noted above, the “transfer” of information from sensory memory to STM results in a significant compression and abstraction of the original signal. The ability to uniquely activate these compressed, abstract detector neurons in higher brain areas requires attention to filter the lower-level sensory input, thus explaining both the need for attention and the lower capacity of STM relative to sensory memory.

Furthermore, the smaller capacity of STM enables it to persist for longer periods of time, because more neurons across multiple of these higher-level areas can participate in representing this information, resulting in a more redundant and robust collection of such neurons. In the terminology from the chapter on consciousness, STM corresponds to the *fully recurrent* activated state, which is highly likely to be the subject of conscious awareness. Indeed, one of the defining characteristics of STM is that you are consciously aware of it.

Thus, the overall picture of STM is that the underlying neurons are mutually activating each other via bidirectional excitatory connections, causing a bit of an “echo chamber” as these spiking signals pass back and forth among these neurons, resulting in a longer-lasting activation trace compared to sensory memory. Rough estimates of the duration of STM extend up to about 30 seconds, but this is strongly dependent on the process of **maintenance rehearsal**, which involves the deliberate attempt to keep those neurons firing robustly by continuously focusing attention on them.

Interestingly, up to this point, the modal model only includes memory mechanisms based on *neural activity*. This reflects the fact that the synapses in the sensory pathways have been very well-trained by the time anyone is participating in Sperling-style experiments, so the synaptic changes there typically don’t make much of a noticeable difference. However, there is an extensive literature on *priming* and *perceptual learning* which can reveal the effects of these ongoing synaptic changes. Thus, as noted above, memory really is happening at every synapse in the brain, whenever activity is sufficient to drive synaptic changes. However, you sometimes have to try pretty hard to see the effects of these changes, and the modal model only covers the most obvious forms of memory.

Finally, the last component of the modal model introduces a form of memory that does depend on synaptic changes, in the form of **long-term memory (LTM)**. In the terms of the modal model, memories are “transferred” into LTM from STM through the process of **encoding**. They can also be recovered back from LTM into STM via **retrieval** processes. This model was developed during the 1960’s, when the computer metaphor was at its height, and this encoding process was typically envisioned as transferring “files” between the RAM-like STM and the hard-disk of LTM. But what does this correspond to in the brain, given that we don’t think the concepts of RAM, hard-disk, or transfer really apply in the brain?

The Hippocampus

This is where the *hippocampus* makes its grand entrance on the memory scene: in most cases, the initial encoding of information from the active state of the cortex (i.e., STM) into a form that can be later retrieved (i.e., LTM) depends on the hippocampus. Because the hippocampus sits at the top of the neocortical hierarchy of areas, it can quickly take a “snapshot” of the current pattern of activity across the upper layers of the cortex (Figure 6.4). Thus, the unique anatomical position of the hippocampus, plus some important special properties of the hippocampus itself, enable it to play such a critical role in the encoding and retrieval of memories.

In brief, you can think of the hippocampal neurons as *detecting* the elements of a memory (e.g., the *who, what, where* elements of an event or *episode*). Synaptic changes in these neurons then enable even a subset of those elements (e.g., the query “what did you have for dinner last night?”) to re-activate these same neurons in the hippocampus. When these neurons fire, they act in turn to re-activate the memory out in the neocortex (i.e., the *retrieval* arrow between LTM and STM in the modal model, Figure 6.2). Thus, whereas neurons in the visual pathways are detecting objects and object features, neurons in the hippocampus are detecting *memories*, and that is why they play such a central role in our mnemonic life.

Figure 6.5 shows one of the magic tricks used by the hippocampus to be able to rapidly encode new memories without overwriting other existing memories, known as **pattern separation**. This is the key idea developed by David Marr as mentioned in the Neuroscience chapter, which applies to both the hippocampus and the cerebellum (Marr 1969, 1971). The idea is that if you simply reduce the number of neurons firing in the hippocampus compared to the cortex (i.e., make them *sparse*), then the patterns of activity in the

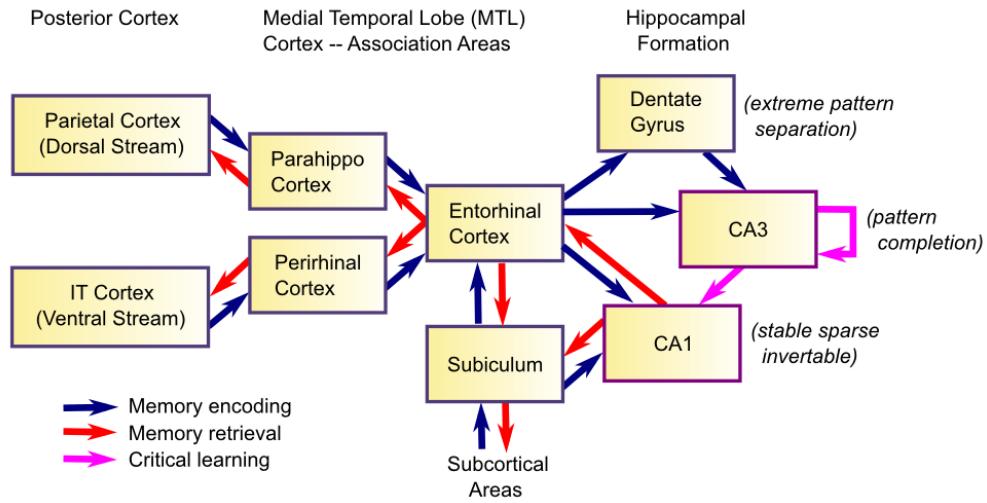


Figure 6.4: Connectivity and structure of the hippocampus. Sensory memory and STM are supported by activity in the posterior cortex areas, which then feed into two cortical areas in the *medial* (middle) region of the temporal lobe, the *parahippocampal* and *perirhinal* cortex. These then feed into the *entorhinal* cortex, which thus has a maximally compressed encoding of everything active in the rest of the brain. The areas of the hippocampal formation then effectively take a snapshot of this cortical activity.

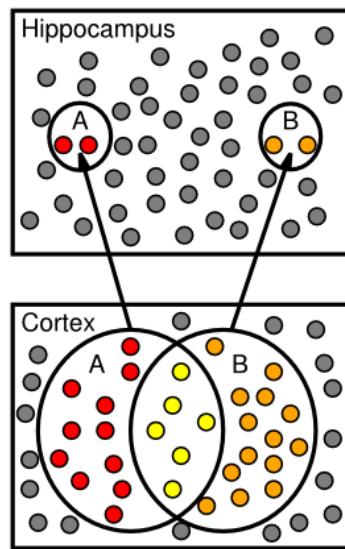


Figure 6.5: Pattern separation in the hippocampus: overlapping patterns of neural activity in the cortex result in separate, non-overlapping patterns in the hippocampus, because it has *sparse* activity (i.e., very few neurons active).

hippocampus will overlap much less than those in the cortex, and therefore, there will be less overlap or interference in the synaptic changes involved in memory encoding. Mathematically, this derives from the fact that squaring a small number, such as .01, results in a *much* smaller number (.0001) – the small number (.01) is the probability of a neuron getting active, and the square is the resulting probability that it would be active in *two* different memories. More realistic, detailed simulations of the hippocampal circuit confirm this basic principle (O'Reilly and McClelland 1994).

There are many important implications of this pattern separation property. First, as we noted in the Neuroscience chapter, this results in a kind of *brute force* memorization strategy in the hippocampus. It doesn't try to make any direct connections between related memories – instead it just effectively sticks each memory in its own separate “box” (i.e., a highly distinct neural pattern with no systematic relationship to other memories). This is great for quickly finding a place to stick a new memory, but it means that the hippocampal version of those memories is a completely disorganized, haphazard pile of these separate boxes.

Thus, there is a much slower process of *organizing* and *systematizing* all those memories, known as **systems consolidation** (which is distinct from synaptic-level consolidation processes involved in the LTP / LTD mechanisms, as we'll clarify below). Specifically, memories that are initially encoded in the hippocampus are gradually incorporated into synaptic changes among neurons in the neocortex, resulting in the formation of more systematic, well-organized **semantic knowledge** (McClelland, McNaughton, and O'Reilly 1995). Some of this consolidation may take place during slow-wave sleep, as discussed in Chapter 3 (Wilson and McNaughton 1994; Buzsáki 1989; Roumis and Frank 2015), and much of it certainly depends on the usual retelling and ruminative replaying of memories throughout the course of daily life.

Hippocampal pattern separation and memory consolidation have major implications for educational learning and expertise. Everything you learn in class is initially encoded through hippocampal brute-force memorization, and only over a relatively long period of repeated learning and practice does a systematic and *productive* form of semantic knowledge emerge. This is consistent with how much experience it takes to become an expert in a given domain: roughly 10,000 hours or 10 years (Ericsson and Lehmann 1996). Thus, if you really want to master something, be prepared to spend a long time slowly shaping your neocortical synapses to develop the necessary systematic knowledge base.

Another important implication of pattern separation is the canary-in-a-coal-mine nature of the hippocampus. Driving down the activity level of the hippocampus requires an extensive amount of GABA inhibition, and thus the hippocampus is extra sensitive to the effects of alcohol and benzodiazepines (e.g., valium, midazolam), which are GABA agonists as discussed in Chapter 3. Furthermore, the rapid rate of learning in the hippocampus requires high levels of NMDA receptors, which makes this system susceptible to epileptic seizures due to the development of over-strong excitatory synaptic connections (recall that H.M. had his hippocampus removed due to epilepsy, which often has a hippocampal source). Both of these factors may contribute to a heightened sensitivity to oxygen deprivation.

Pattern separation also has important implications for the retrieval of memories from the hippocampus. To the extent that it is always trying to keep different patterns separate, it is then hard to take a partial retrieval cue (e.g., the “what did you have for dinner?” question) and have that re-activate the original pattern of neural activity that was present when the memory was originally encoded. This retrieval process is called **pattern completion**, as it involves filling-in or completing the partial cue pattern. Instead of doing pattern completion, the hippocampus might just end up encoding a retrieval attempt as a brand new experience, and activate entirely new neurons as a result of pattern separation.

Thus, pattern separation and pattern completion are essentially opposing forces within the hippocampus. Pattern completion is supported by special connections within one of the main areas of the hippocampus (the CA3), which effectively “glue” together the different elements of a memory. Detailed analyses of the battle between pattern separation and pattern completion suggest that the specific anatomy of the hippocampus is particularly well-suited for balancing between these competing demands (O'Reilly and McClelland 1994).

Taxonomy of Long-Term Memory

Although the hippocampus plays a dominant role in the process of encoding new memories into LTM from STM (in the terms of the modal model), synaptic changes occur everywhere there is activity in the brain. Thus, a major focus of memory research has been attempting to document and organize all these different

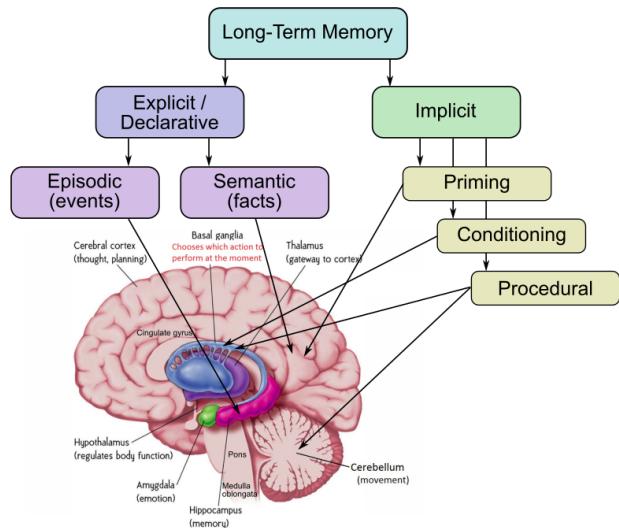


Figure 6.6: A standard Long-Term Memory taxonomy, and associated neural substrates. The broadest distinction is between consciously accessible memories supported by the cortex and hippocampus, versus non-conscious memories largely in subcortical areas. Priming is an interesting case of a cortical memory effect that is not directly accessible to consciousness.

“types” of long-term memory within overall memory **taxonomies** (akin to the taxonomies used to organize different species of animals, for example). Because memory and processing are both occurring within each and every neuron, these taxonomies are really just descriptions of the different kinds of processing taking place in the brain, which we reviewed in detail in the Neuroscience chapter. Nevertheless, we will briefly review the most popular taxonomy, and the research that went into its construction.

Figure 6.6 shows perhaps the most widely adopted LTM taxonomy, proposed initially by Endel Tulving (Tulving 1972) and refined by Larry Squire (Squire 1992). It features a top-level division between **explicit** or **declarative** memory, as contrasted with **implicit** or non-declarative memory. Explicit memories are those we can have direct conscious access to (and declarative means you can declare it verbally), while implicit memories are not consciously accessible. Given what we know about consciousness from Chapter 3, explicit memories are therefore those in the neocortex. Interestingly, we likely are not directly conscious of hippocampal memories, given the requirement of recurrent / bidirectional connectivity for consciousness, which is only partially present in the hippocampus. Instead, we become conscious of hippocampally-supported memories when they are recalled back into cortex.

The two major subtypes of explicit memory are hippocampal **episodic** memories (i.e., the memories of all the daily events and episodes in our lives, and those we read about or watch in movies or on TV), and **semantic** memory, which is a summary term for all of the facts and knowledge we have, which has been integrated into our cortical synapses over many years of memory consolidation. During this consolidation process, the episodic character of the knowledge gets winnowed away, leaving only the bare knowledge devoid of the **source** or **context** information about where we learned these facts. Newly learned facts (e.g., much of what you are learning in this course) still retain their episodic trace – you can probably recall when you heard about something interesting for the first time in lecture, or read about it in a book. Sometimes, people feel like they have a particularly clear sense of where on the page they read something, but in my experience this has proven illusory more often than not.

Within the much more diverse umbrella of implicit memories, there are *procedural*, *conditioning*, and *priming* memory traces. The separability of **procedural** memory from hippocampal episodic memory was vividly demonstrated by H.M., who was able to learn a challenging new procedural task such as learning to trace a picture when looking in a mirror (try it – it is hard!) at the same rate as neurologically intact control participants. This is because procedural learning depends on circuits through the frontal cortex, cerebellum and basal ganglia, not on the hippocampus. Likewise, as we reviewed in the Learning chapter, **conditioning**

depends on the amygdala, basal ganglia, and dopamine systems, and is thus separable from hippocampal and cortical memories (and was also intact in H.M.).

The value of this memory taxonomy is debatable. Really, it is just assigning new labels for the functions of brain areas, which can be much more richly and accurately described (e.g., as in the Neuroscience chapter) than in such a broad taxonomy. Furthermore, it is missing many important parts of the brain. Perhaps most importantly, the central division according to the criterion of conscious access is problematic at many levels. Consciousness is inherently subjective, and putting a subjective construct at the center of a major theoretical framework jeopardizes the entire enterprise. Furthermore, it immediately eliminates application to animal memory (Morris 2001), as the notion of consciousness in animals is certainly fraught with controversy. It also unnecessarily complicates any kind of straightforward understanding of memory in terms of underlying neural mechanisms.

For example, given our detailed understanding of how the hippocampus works, it is highly likely that even rats (which have a large hippocampus relative to the rest of their brain) encode something like episodic memories of all the different experiences in their lives. Rats likely don't sit around idly reminiscing as people do, but that doesn't mean they don't re-activate their episodic memories in response to relevant stimulus cues – indeed, this has been demonstrated in many experiments recording from hippocampal neurons in rats. Thus, it is more productive to find the many parallels in brain systems across species, so that we can integrate a much broader scope of data into our theories of memory.

The case of **priming** is particularly illustrative of the limitations of a consciousness-based framework. Priming is the measurable facilitation in processing information that was previously processed (i.e., "priming the pump"). It results from small synaptic changes throughout the neocortex, driven by neural activity. Thus, although we are not directly conscious of priming itself (e.g., we don't know that our responses are faster by about 10 msec), we *are* typically conscious of much of the activity that drives priming. And these are the very same synaptic changes that add up over time to produce new semantic memory learning. So does it really make sense to put this in the implicit memory category? Another example is the considerable contributions of the frontal and parietal cortex to procedural tasks: we can certainly be aware of activity in these brain areas, and yet they are put in the implicit category.

Another interesting case similar to priming is the difference between **recognition** and **recall**, which has been studied extensively in the memory literature (Jacoby, Toth, and Yonelinas 1993). Recognition memory is characterized as using the overall feeling of *familiarity* with a given stimulus to decide whether it was on a given memory list, whereas recall involves the explicit, conscious *recollection* of episodic details from the time of study. Recollection generally depends on the hippocampal pattern completion process to re-activate those episodic details, whereas familiarity can be supported by differences in neocortical activity patterns reflecting synaptic weight changes, which are not strong enough to drive full recollection (Norman and O'Reilly 2003). Thus, familiarity is similar to priming, but interestingly, H.M. and some other severe amnesics were impaired at familiarity-based recognition memory, but their priming was intact. This is because the familiarity signal is likely driven by the neocortical areas surrounding the hippocampus (e.g., perirhinal and entorhinal cortex) that were damaged along with the hippocampus proper, whereas most priming tests probe lower-level semantic or visual cortical areas.

Amnesia

Patients with hippocampal damage such as H.M. have also shown us that two different types of **amnesia** (loss of memory function) can be *dissociated* (i.e., separated, do not always co-occur): **retrograde** vs. **anterograde** amnesia. Retrograde refers to memories of the past (like "retro" styles etc), while anterograde refers to the ability to form new memories. H.M. was profoundly impaired in his ability to form new memories, and thus suffered from severe anterograde amnesia. However, many of his memories from his more distant past were largely intact, meaning that he had comparatively mild retrograde amnesia. Furthermore, his basic semantic knowledge of facts etc was largely intact.

We can understand this dissociation in terms of the basic explanation of hippocampal function given above. The hippocampus is critical for rapidly learning new episodic memories, because of its unique position at the top of the cortical hierarchy, and its special properties including pattern separation and a relatively fast learning rate. Thus, damage to the hippocampus almost always produces significant impairments in

encoding new episodic memories. However, because of the gradual incorporation of episodic memories into the neocortex through the consolidation process, older memories from the past can still be recalled even without the help of the hippocampus.

Interestingly, consolidation predicts that more recent memories leading up to the point of hippocampal damage should be most impaired, as they have had less time to be consolidated into the neocortex. This *gradient* of retrograde amnesia is often observed at least to some extent in human amnesics. Interestingly, extensive investigations of retrograde gradients and memory consolidation in rats have produced inconsistent results, likely reflecting the variability in the extent to which rats actually recall prior episodes, across different experimental paradigms (Sutherland, O'Brien, and Lehmann 2008; Anagnostaras, Maren, and Fanselow 1999).

Another fascinating form of amnesia is **childhood amnesia**, which is the widely-documented phenomenon that people cannot generally remember anything before about 3 years of age. Go ahead, give it a try – can you? Many studies have attempted to understand the reasons for this amnesia (Hayne 2004). Overall, it is likely a result of the fact that the neocortex is not sufficiently well organized before that age, to support the ability of earlier hippocampal snapshots to be translated back out into the cortex. In effect, the “language” that the earlier snapshots were recorded in is no longer something that the more mature brain can understand (and indeed language learning itself likely plays a significant role).

Although the neocortex continues to develop and learn in significant ways beyond the age of 3, there is presumably just enough stability for those earliest memories to persist. And those early memories that you can still recall have likely been recalled, reinforced, and elaborated many times in the ensuing years, so they are well-consolidated and may not actually be very accurate anymore. Nevertheless, in my own case, I feel like I do have vivid, first-person memories of my 3-year-old-self living for 6 months on the island of Grenada in the Caribbean, including a scary encounter with a large crab behind the house. Thus, as is the case with memory in general, emotional arousal and the relative novelty of experiences play a large role in one’s ability to later recall them.

Memory Capacity and the Importance of *Chunks*

One of the dimensions along which memory systems vary is in terms of their capacity, with sensory memory being high capacity, STM having a strongly limited capacity of around 3-4 items, and LTM being essentially unlimited in its overall storage capacity. But any consideration of capacity raises the central question of *what counts as an item for the purposes of measuring capacity?* In the Sperling experiments (Figure 6.3), items were individual letters, but what if we instead put words where the letters were in the 3x3 grid display? Memory capacity will be about the same, now measured in words instead of letters, but that represents a considerable increase in overall *letter* memory capacity!

The answer to this puzzle is to introduce the concept of a **chunk**, which is somewhat circularly defined as an element that acts like a single item with respect to memory capacity measurements. If the stimuli are *random* letters, then each letter is a chunk, but if the letters can be formed into words, then the word becomes the chunk. Likewise, if words can be combined into sensible sentences, then those sentences become the chunks. In short, a non-circular definition of a chunk is *anything that we have an existing stable neocortical semantic representation of*. This is still not very precise, but it will do for now.

Based on an influential and provocative article by *George Miller* (Miller 1956), many textbooks incorrectly cite the capacity of STM as *7 plus or minus 2*. However, Sperling’s original data, and data from many other tasks and domains, strongly suggests that it is actually **the magic number 4** (Cowan 2001; Luck and Vogel 1997). Although overly simplistic, one way of thinking about this is that each of your two cerebral hemispheres can hold 2 items when pushed to the limit, such that $2 \times 2 = 4$, with 1 per hemisphere being much more comfortable (Buschman et al. 2011). The higher capacity of 7 applies only to verbal memory that can be sustained by a rehearsal mechanism known as the *phonological loop*, where you repeatedly verbalise (in your mind, but also using your actual vocal muscles at a subthreshold level) the to-be-remembered material (Baddeley, Gathercole, and Papagno 1998). Our extensive experience with verbal material presumably produces this larger capacity beyond what is generally available with the “default” neural mechanisms.

Encoding and Retrieval Strategies (i.e., How to Study!)

Because memory capacity is determined by the availability of appropriate chunks, one major category of memory-enhancement tricks involves creating new chunks, and efficiently leveraging the ones you already have. This is the main trick employed by contestants in the memory olympics competitions, and was well-documented in the case of an individual, S.F., who developed chunking strategies that allowed him to remember over 100 random digits (Ericcson, Chase, and Faloon 1980). In this case, he turned 3-digit numbers into times to run a mile or other standard distances, as S.F. was an avid runner. Another common example of chunking is the creation of acronyms. For example, the “big five” personality dimensions that we’ll encounter later can be organized into the acronym *OCEAN*, which then makes it much easier to remember them all.

Another effective encoding / chunking strategy (i.e., **mnemonics**) is to associate different words with different familiar spatial locations, known from the days of ancient Greece as the **method of loci**. An even more **elaborative encoding** strategy is to create stories involving these locations and familiar people (e.g., “my mom went from the living room to the kitchen, to get a popcorn snack”), where each of the words then stands for something that you’re trying to remember (e.g., mom = 3, living room = 7, kitchen = 4, and popcorn = 8). Because the hippocampus really loves to encode episodic memories, these episodic chunks are particularly effective and memorable.

Several other related principles of effective memory encoding have been developed. For example, the influential **levels-of-processing theory** (Craik and Lockhart 1972) postulates that more *deeply* encoded information will be better remembered. The notion of levels or depth here corresponds to the levels of processing in the neocortex, going from raw sensory information up to higher-level semantic information. For example, many studies have found that encouraging people to think about the meaning of a word, as compared to noticing the case or font of the letters, results in better memory.

An interesting example of the benefits of deeper, more elaborative encoding comes from the notion of **desirable difficulties** (Bjork 1994) – memory is often better if you have to work harder to process the information, even in sometimes fairly strange ways. For example, making information harder to read can improve subsequent memory, and a font was recently created called *Sans Forgetica* to leverage this finding – unfortunately it doesn’t actually seem to improve memory (Taylor et al. 2020), suggesting that not all forms of difficulty are created equal!

One of the most robust ways of improving learning, which is directly relevant to success school, is the **testing effect** or the **retrieval practice effect**, where taking a test improves subsequent memory (Roediger and Butler 2011). Importantly, this testing effect results in significantly better learning than a *reexposure* condition, where the same material is presented in full. In other words, you really have to *test* yourself, not just re-read something (in this way it is a kind of desirable difficulty). This obviously has important implications for how you study: simply re-reading over the text is much less effective than testing yourself to explain a concept given a cue. This is why we have the key words listed at the end of each chapter – you should go through each one and try to generate a full explanation of what that term means and why it is important. Then, you should go check your answer by finding the relevant text. Also, hopefully your class includes weekly quizzes that give you a more structured opportunity to test yourself. At a neural level, the testing effect creates *error signals* between your guess and the right answer, and as we discussed in the learning chapter, the brain likely uses error-driven learning.

One of the best ways to really learn something is by teaching it to others, which is a version of the **generation effect** – having to produce a sensible explanation of something greatly improves comprehension. This is similar to the testing effect. Try cornering a friend and give them a mini-lecture on how memory works in the brain – you’ll soon find the gaps in your understanding, and strongly reinforce the parts you already do understand. Seriously, if you want to learn, teach! This is one of the most important synergies in academia: by having to explain what we’ve learned in our research through teaching, professors then understand it all much better.

One of the most fascinating encoding principles is the **encoding specificity principle** (Tulving 1983), which reflects the fact that episodic memories tend to bind together all of the different elements present when a memory is encoded, and thus recall of those memories will be best when those original elements are present at the time of recall. This is a direct result of the pattern completion vs. pattern separation battle operating in the hippocampus – if too many elements are different from the original event, the hippocampus tends to

perform pattern separation instead of the pattern completion required for recall. All those random elements present at the time of encoding are typically summarized with the term **context**, and thus lead to the **context-dependent memory** phenomenon (Yonelinas et al. 2019). A classic example of this phenomenon is that people are better able to recall information when tested in the same physical context as it was originally learned – for example, if you study in a library, then taking a test in that same library will generally result in better performance.

The most famous demonstration of this encoding specificity / context-dependent memory principle was conducted in a study where items were learned either on a beach or underwater using scuba equipment, and then tested either in the same or different context (Godden and Baddeley 1975). Participants in the same-context conditions (either on land or under water) performed significantly better than in the cross-context conditions. Another notorious demonstration involved study and test either drunk or sober, which again found that, surprisingly, testing while drunk was better than sober *if* initial learning was drunk (Goodwin et al. 1969). This has been labeled the **state-dependent memory** effect, and presumably reflects the same encoding specificity principle.

It is important to emphasize that memory was much worse when learned drunk, even when tested drunk, so that is *not* a good strategy overall. Other demonstrations of state-dependent memory involve mood states, and we'll see later that this **mood-dependent memory** effect creates an unfortunate feedback loop in depression, where you're much more likely to remember all the bad memories in your life when you're depressed, making everything seem that much more bleak. On the bright side, this also works for positive memories in positive mood states.

Another critical way to improve encoding is to use **spaced** instead of **massed** practice – i.e., to space out your studying over multiple separate study sessions, instead of *cramming* at the last minute. This is beneficial for the same reason that gives rise to context-dependent effects: spacing out study causes the information to be learned across multiple different contexts, and thus helps to make the knowledge more independent of that context. A critical point here is that *context* includes a significant contribution from *time* – your internal mental state is constantly evolving over time, and will be significantly different a week or two from now (Howard and Kahana 1999). Thus, even if you study in the same physical context, your internal mental context will be different, shaped by all those synaptic changes taking place between the two study sessions. Indeed, a critical aspect of the memory consolidation effect involves exactly this process of thinking about the same issues from the very different perspectives that emerge as your brain changes over the period of years.

In short, you should study by engaging in deep, elaborative encoding of the material, connecting it in multiple different ways with your existing knowledge chunks, and testing yourself as much as possible, ideally by trying to teach material to others. Furthermore, you should do this in a spaced fashion, across multiple different days, ideally in different physical and mental contexts. That's not so hard, is it?

Memory Retention and Interference

Even once you've successfully encoded some new information into LTM, it is still not safe! Memory is often a fleeting thing, as you have almost certainly experienced. Figure 6.7 shows the data from Hermann Ebbinghaus who pioneered the study of memory retention in the late 1800's (Ebbinghaus (1885) 2013). This curve is striking in its steep initial dropoff, followed by a relatively stable plateau. We can understand the nature of this curve in terms of two different processes, which have long been the subject of debate in the field: **decay** and **interference**. It is surprisingly difficult to distinguish these two on purely behavioral grounds, as it is generally impossible to prevent interference from happening, and decay is defined as an automatic, continuous process. However, detailed studies of the molecular processes following the synaptic plasticity events described in the Learning chapter allow for some resolution of this debate.

It is likely that the steep initial dropoff in memory is due to synapse-level processes, which can be considered a form of decay, but that after roughly a day or two, synaptic changes have stabilized to the point that subsequent forgetting is mostly due to interference effects. These synapse-level processes are collectively known as *synaptic consolidation*, which is distinct from the *systems consolidation* processes described above, where memory initially encoded in the hippocampus is learned in the neocortex as well. As shown in Figure 6.8, there is a roughly 15-20 minute period when synaptic changes can decay rapidly if

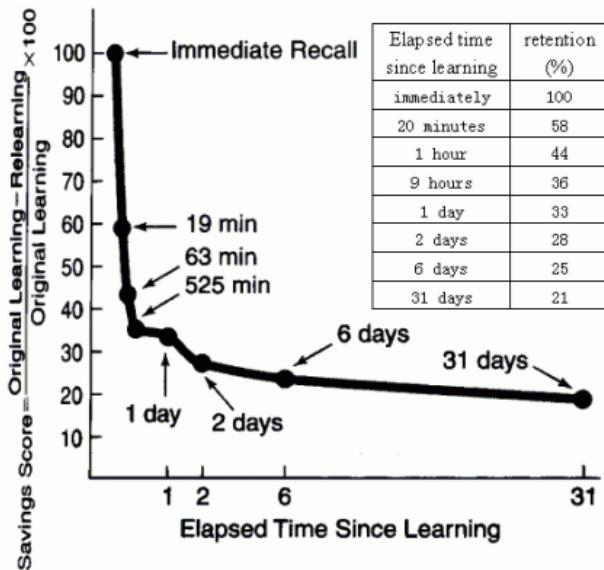


Figure 6.7: The forgetting curve from Hermann Ebbinghaus's data. The initial steep dropoff is likely due to synaptic-level stabilization processes, and the longer plateau reflects essentially permanent long-term memory, with loss due largely to interference.

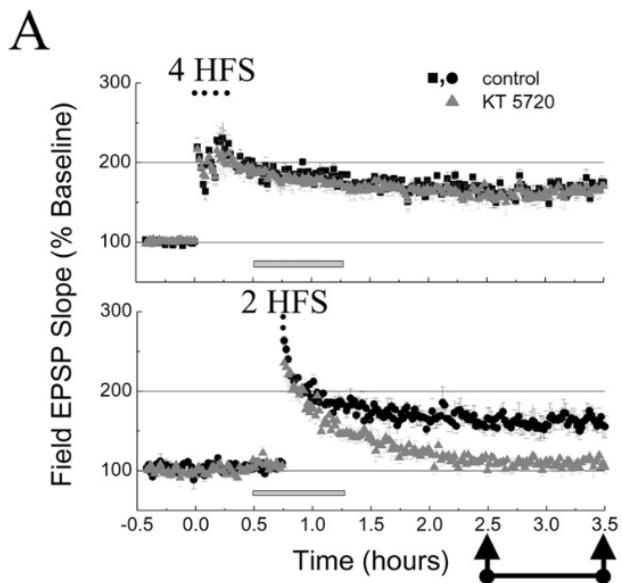


Figure 6.8: Forgetting curve from synapse-level stabilization effects, which shows a steep dropoff in synaptic strength for more weakly potentiated synapses (2 HFS curves in bottom graph; HFS = high frequency stimulation) compared to more strongly potentiated synapses (4 HFS, top graph). The KT 5720 curve shows the contribution of protein synthesis, which emerges over the period of an hour or so, and if these proteins are not available, the weaker memory decays back to baseline. From Alarcon, Barco & Kandel (2006).

they were not sufficiently strong in the first place, or reinforced by subsequent plasticity events (Alarcon, Barco, and Kandel 2006; Frey and Morris 1998). Over the course of an hour, synaptic changes are reinforced by processes that depend on new proteins being synthesized, including muscle-like *actin* fibers. Further stabilization occurs during sleep, over the next day or two (for all the details, see (Rudy 2013)).

After all of this synaptic consolidation has taken place, it is likely that further loss of memory is due to interference effects, which occur when new synaptic changes move the synaptic strengths in a different direction than was needed for an existing memory. This is known as **retroactive interference**, because it is interfering with older (“retro”) memories. The extreme pattern separation in the hippocampus can help to minimize the amount of retroactive interference, by encouraging the use of distinct sets of synapses to encode different memories, but it is impossible to completely eliminate interference.

An example of retroactive interference would be when you encode where you parked your car today, versus yesterday. Because of the large amount of overlap in the overall context, you likely re-activate many of the same neurons involved in encoding these two memories. Thus, the synaptic changes that are made today will help you recall that it was parked in the South-West corner of the lot, but these changes will likely overwrite many of the synapses that encoded the “South-East” location from yesterday.

There is another form of interference called **proactive interference** which is somewhat strange compared to the more intuitive nature of retroactive interference. In proactive interference, prior learning interferes with your ability to form *new* memories. This can happen if you are trying to learn new information about the same items over time. For example, if you use distinctive new items on every trial of a memory task (shapes, colors, letters, words, animals, etc), then it is easier to remember those items compared to re-using the same items repeatedly (Hasselmo and Stern 2006).

The Fallibility of Memory

In addition to failures of basic encoding and forgetting, there are other pitfalls in the domain of memory, which arise largely from the fact that the hippocampus only receives a highly *compressed* view of the outside world, filtered through many layers of cortical processing and compression as shown in Figure 6.4. Figure 6.1 of the telephone game also captures the kind of compounding effects that emerge from information propagation through the cortex. From this perspective it is a wonder that we can accurately remember anything at all!

As we saw in the perception chapter, even our lower-level perceptual system is strongly influenced by top-down, internal biases, and this process is ubiquitous throughout the cortex. Thus, we all encode our memories through spectacles of one shade or another, and are susceptible to having **false memories**. One of the first demonstrations of this point in the memory literature was due to Frederic Bartlett, who tested people’s ability to remember a story known as the “War of the Ghosts”, over an extended period of time (Bartlett 1932). This story was based on Canadian Indian folklore, and contained many concepts and events that were entirely unfamiliar to the English participants in his experiment. As a result, the participants had great difficulty remembering the story, and ended up reshaping it to fit their own conceptual structures. These conceptual structures are called **schema**, and we’ll revisit them again in the next chapter.

An important real-world implication of this strong tendency to *schematize* memory is in **eyewitness testimony**, where people are likely to encode the events of a crime according to their existing *stereotypes* and biases. Furthermore, these biases can be activated by leading questions. For example, in one seminal study, the experimenters manipulated the use of leading terms like “smashed” in a car crash scenario, and this had large effects on participant’s memory of things like the speed and damage involved (Loftus and Palmer 1974). Interestingly, as was the case in the Bartlett study, participant’s confidence in their *false* memories was often higher than for their accurate ones.

The other major issue that has received considerable media attention is recovered memories of childhood sexual abuse. Unfortunately, abuse is all too common, but it is also the case that memory in young children is even more unreliable than in adults. Studies have shown that children can report having actually experienced events that they only imagined (Ceci et al. 1994), and some forms of therapy designed to uncover repressed memories may have used leading questions that could have created false memories.

In the experimental literature, false memory has been extensively explored using the *Deese, Roediger, McDermott (DRM)* paradigm (Deese 1959; Roediger and McDermott 1995). In this paradigm, a number of words that overlap strongly with a given target word (e.g., pillow, dream, night, etc) are studied, with the

result that the target word (“sleep” in this case) is often confidently endorsed as having been on the study list. This is vivid demonstration that memory operates on high-level compressed semantic representations.

Working Memory and the Prefrontal Cortex

Finally, we conclude with one more important distinction between different types of memory, in this case between short-term memory (STM) and **working memory** (WM), which was proposed by Alan Baddeley and Graham Hitch (Baddeley and Hitch 1974). The notion of working memory resembles the functional properties of RAM in a standard computer: information that is currently being processed, maintained in an active, directly accessible state. Working memory is distinguished from “regular” STM, where the latter includes just basic maintenance of information, whereas working memory is specifically about the information used for ongoing processing, which is particularly strongly maintained, even in the face of potential distractors.

As is often the case, the biology may provide a more precise definition of the difference between STM and working memory, in the form of robust sustained firing of neurons in the prefrontal cortex, which was discovered in the early 1970’s (Fuster and Alexander 1971; Kubota and Niki 1971). This sustained neural activity was postulated as the neural basis of working memory (Goldman-Rakic 1995), and studies showed that this form of neural activity is indeed more robust and resistant to distraction than activity in posterior cortical areas (Miller and Desimone 1994). Computational models have shown that the basal ganglia can play a critical role in supporting this robust active memory in frontal areas, by dynamically switching the system between maintenance and rapid updating modes (R. C. O’Reilly and Frank 2006). Thus, overall there is ample biological evidence that sustained neural activity in the frontal cortex is different from that in posterior cortical areas, in ways that accord with the overall distinction between working memory versus STM. We’ll focus more on this frontal / basal ganglia working memory system in the next chapter.

Summary of Key Terms

This is a checklist of key terms / concepts that you should know about from this chapter. As we just learned in this chapter, it is a great idea to test yourself on what was said about each of these terms, and then go back and double-check – that provides both beneficial repetition and also the *testing effect*.

- Neural mechanisms of memory:
 - Activity (spiking etc): fast, active and transient
 - Synaptic changes: slower, long-lasting
- Modal model
 - Sensory memory: iconic, echoic (high capacity, short-lived) – neural activity in sensory cortex
 - Short-term memory (STM): requires attention, limited capacity (magic number 4) – neural activity in higher cortical areas
 - * Sperling task
 - Long-term memory (LTM): requires encoding – synaptic changes in hippocampus and cortex.
- Hippocampus
 - Anatomical location on top of cortex
 - Pattern separation from sparse activity
 - Pattern completion to recall memories
 - Memory consolidation: semantic knowledge forms slowly in neocortex
- LTM Taxonomy
 - Explicit / Implicit
 - Episodic / Semantic / Priming / Conditioning / Procedural
 - Issues with consciousness
- Amnesia
 - Anterograde
 - Retrograde
 - Childhood amnesia
- Encoding / Retrieval Strategies
 - Chunk
 - Mnemonic

- Method of loci
- Elaborative encoding
- Levels of processing
- Desirable difficulties: Testing effect, generation effect
- Encoding specificity principle
 - * Context-dependent memory
 - * State-dependent memory
 - * Mood-dependent memory
- Massed vs. Spaced practice (cramming is bad)
- Memory Retention and Interference
 - Decay: synaptic stabilization
 - Interference: Retroactive vs. Proactive
- Fallibility of Memory
 - False memories: War of the Ghosts
 - Schema
 - Eyewitness testimony & leading questions
- Working memory vs. STM
 - Robust firing in prefrontal cortex

Chapter 7: Thinking, Control and Intelligence

What is *smart*? This is the fundamental question for this chapter, with many profound personal and societal implications. Is there just one kind of smart, or are there multiple different forms of intelligence? How can we reconcile any form of *general* intelligence with everything we've learned up to this point, about how the brain works at a biological level? The brain is composed of billions of neurons, interconnected by vast networks of synapses, wherein all of our knowledge, and, presumably, intelligence, must lie. Do "smart" people have more neurons or synapses? Or, perhaps, *fewer* synapses? Are their neurons somehow fundamentally different from other people who measure as less smart according to standard intelligence tests? And what are those intelligence tests measuring anyway? Are they really some kind of "pure" measure of intelligence, or do they just reflect the degree of western-style education (and health and wealth) that a person has? What does your IQ score really tell us about you as a thinker, and about your prospects for future success in school and the real world? So many important questions!

If our brains were more like digital computers, these questions would have much simpler answers. It is relatively easy to measure the power and speed of a computer, and many people tend to think of human intelligence in these terms. As we discussed in the previous chapter, a computer has discrete parts (the CPU, RAM, and hard drive), and each of these parts can be directly quantified in terms of capacity and speed. If you're at all savvy about these things, you can obsess about getting the best value for your money along each of these dimensions, and, generally speaking, the faster the CPU and the more RAM and hard-drive storage, the more you can achieve with your computer. Computers really do come in obvious degrees of "smartness".

But our brains are nothing like that of a digital computer. Instead, cognition emerges out of the interactions of billions of chattering neurons, which are fundamentally shaped by learning processes over an extended period of time. As we will explore in the development chapter, we start out with virtually no discernible intelligence (despite how cute and special our parents think we are), and it takes most people a few *years* to even learn how to control their own bowels! Wow. The rest of the animal kingdom must think we are complete idiots, which comports with an amusing *Onion* headline to that effect.

Given that we clearly don't start out with much in the way of intelligence, it seems hard to escape the conclusion that intelligence is fundamentally a product of learning (in concert with other developmental / maturational changes). And this view is also hard to avoid when you think about all those synapses that need to get wired up in just the right way to produce whatever cognitive abilities we end up with.

So are "smart" people just better learners then? If so, what makes some people better at learning than others? When we explored this question in the Learning chapter, one of the major conclusions is that learning is driven fundamentally by *motivation*, and all that dopamine and related machinery that gets us up in the morning and ready to pursue our daily goals, etc.

Indeed, we will review various sources of evidence that are consistent with the overall idea that motivational differences play an outsized role in determining measured level of intelligence. Of course, there are many, many complex factors that shape an individual's trajectory of learning and development, and motivation is itself a multi-faceted thing, so perhaps we aren't explaining too much when we say that motivation plays an important role.

But understanding the major factors shaping intelligence may affect how we think about ourselves, and others, in important ways. If we view intelligence as a product of learning and motivation, then it is more obviously malleable. This is the critical difference between a **fixed mindset** about intelligence, versus a **growth mindset**, as emphasized by *Carol Dweck* and colleagues (Dweck 2008), in an increasingly influential body of work. The growth mindset emphasizes that intelligence is not something that people "have", but rather, something they have to cultivate – something that grows over time. Increasingly, schools and teachers are recognizing that motivational factors have a huge impact on educational success, and they are developing innovative ways of motivating students to learn, and making the material more obviously self-relevant.

Fundamentally, the idea that intelligence is largely the product of time spent learning means that **anyone can learn anything**, if they only have sufficient motivation and time to invest into it. This open-ended, ambitious view of intelligence surely has the effect of opening up your individual horizons and sense of what is possible. Personally, I have always had this belief, and I have learned lots of complicated things, often slowly and with great difficulty. Eventually, things that once seemed impenetrable become just another familiar part of my mental toolkit. I have a very salient early memory of spending far longer than my peers figuring

out how to simply connect a battery to some gadget in a summer school class as a kid. I felt like an idiot. But eventually, I figured it out, and learned this valuable lesson that, with sufficient effort, I could succeed.

Hopefully, you are now motivated to learn more about the history and current state of understanding about the nature of human intelligence, and the thinking processes that underlie it! We'll start off by exploring the core questions of what "thinking" is, and what kinds of brain mechanisms are particularly important for it. The conclusion from this may seem to contradict what was just said above: maybe we *do* have something like a CPU in our heads after all – except it is a CPU made out of neurons and brain systems, and it runs on dopamine! This is an important example of an *emergent* system, like the gears we talked about in the neuroscience chapter: the overall function of a CPU can be supported by various different "substances", just like the gears can be made of many different materials, and yet still function more-or-less the same.

Nevertheless, our neural CPU has major differences from a computer CPU, and the fact that it is made of neurons does have important implications for how it works. Indeed, one can understand a lot about the particular strengths and limitations of human cognitive function, in terms of the overall idea that we can do both neuron-like computation, *and* something that approximates the function of a digital CPU. We have yet to develop powerful AI (artificial intelligence) systems that capture this unique combination of both forms of computation, and perhaps once we do, we will unlock the real magic of our brains!

After gaining a better understanding of the "mechanics" of intelligence, we'll review the history of thought about the nature of intelligence, and how it has been measured. Furthermore, we'll examine the data about the real-world implications of IQ test scores, and circle back to these big questions about the relationship between intelligence and motivation.

Another way of thinking about all of these issues, is in terms of the *control* component of our three-C's. Our neural CPU serves as a kind of overall control system for the rest of our brain, and, as we have emphasized, this is fundamentally a *motivated* form of control, focused on getting us the things we need and want, and avoiding all the bad stuff. Thus, the idea that motivation and intelligence are inextricably intertwined makes perfect sense from this perspective: the brain systems supporting our control systems (in the prefrontal cortex and basal ganglia) are the very same ones that directly interface with lower-level motivational and emotional pathways in the amygdala and dopamine system.

The Neural CPU in the Prefrontal Cortex and Basal Ganglia

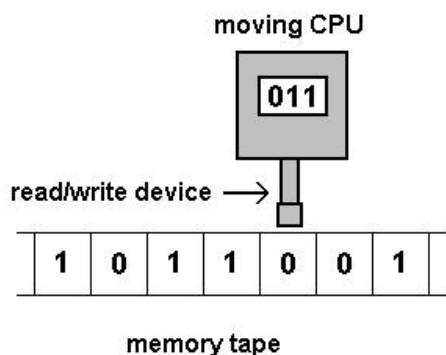


Fig 7-1: The components of a Turing machine: with just three basic components, any computation can be performed!

To understand what kind of neural machinery it would take to support CPU-like functionality in the brain, we start with the surprisingly simple mechanisms needed to make a computer work. At the most abstract level, *Alan Turing* and *John Von Neumann* worked out the basic principles of a *universal* computational device (something that could in principle do *anything*) in the 1930's and 40's (Turing 1936; von Neumann 1945). Amazingly, this device only requires three essential components (Figure 7-1): 1. A way of reading

and writing information from a memory system (conceptualized as a *tape* by Turing); 2. A *program* that determines how this information is transformed in between being read and written; and 3. Some *active* memory where things can be temporarily cached, for the program to refer to. These elements were elaborated by Von Neumann, in one of the most important unpublished papers of all time (von Neumann 1945), creating the foundation for modern digital computers. Now days, we take it for granted that computers can do almost anything, but this was just theory not so long ago.

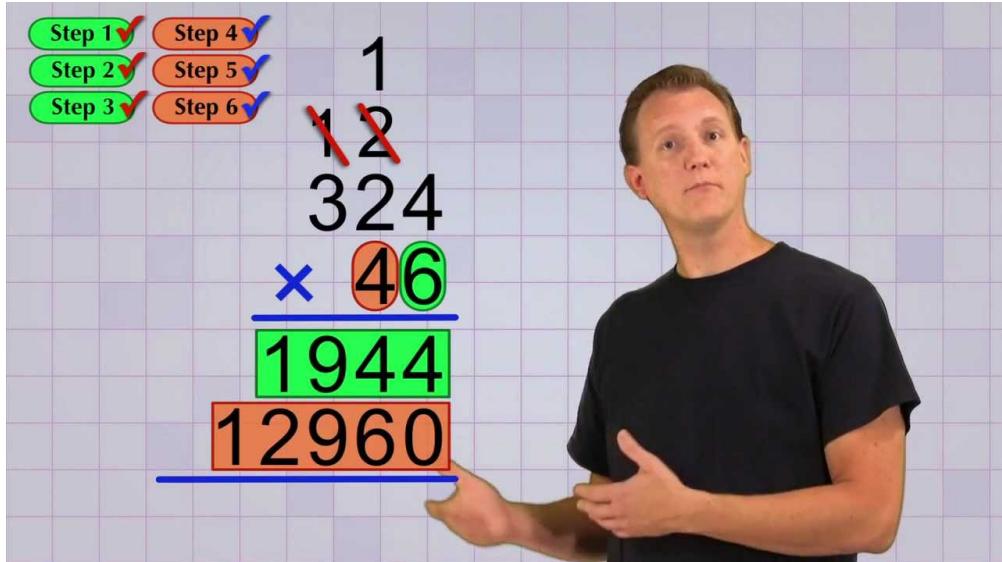


Fig 7-2: Computers solve problems by breaking them down into many small *sequential* steps, each one involving a specific, well-defined operation such as adding numbers, writing them down somewhere, and reading them back in for use later. Just like you do when performing multi-digit arithmetic. Alan Turing showed that these basic processes can be used to solve *any* problem.

You can get a good feel for how a computer works, and why it can do anything, by considering the traditional strategies for performing multi-digit arithmetic (Figure 7-2). Instead of just staring at those big numbers, you break the problem down into a sequence of simple, discrete steps. That sequence of steps is the *program* or **algorithm**, and each individual *operation* involves one of a small set of different processes, such as adding or multiplying single-digit numbers, writing down some numbers for later use (i.e., storing onto the tape in a Turing machine), and reading those numbers back in at the appropriate time (as you move to the next column of digits).

This kind of sequential, discrete, step-wise processing is entirely different from how our neurons work. Neurons also break down a problem into simpler components, but a critical difference is that they all work together in *parallel* instead of the fundamentally sequential, *serial* processing required for a universal computer. The major advantage of serial processing is that it is much more flexible – any arbitrary collection of operations can be sequenced one after the other over time, but the same is *not* true for parallel computation. Some operations are mutually incompatible with each other, or depend one on the other, and simply cannot be performed simultaneously in parallel. Indeed, one of the great challenges of modern computer science is trying to come up with even moderately usable parallel computing frameworks, and it is very clear that the universal flexibility of traditional serial computation does not extend into the parallel realm: parallel computation must generally be setup on a case-by-case basis. For example, in the case of multi-digit multiplication, you have to do the tens-place part of the problem first, before you know how much to carry over to the higher digits, etc – you can't just do everything all in one step.

More generally, parallel systems are really good at doing the same kind of thing over and over again really fast (e.g., detecting patterns via networks of interacting neurons in our brain), but they are not so good at doing random, arbitrary, *different* things, which is precisely where serial computation excels. However, serial computation is inherently much slower (one step at a time). These fundamental tradeoffs between parallel and serial computation mean that a system that can do both will be able to achieve the best of both

worlds – that is the magic recipe that the human brain has achieved. Our brains are parallel at the level of individual neurons and networks of neurons, but at the larger *systems* level of the brain, we can achieve a form of flexible, serial computation.

Before turning to the biology of the brain systems supporting this latter form of computation, we can see strong evidence for the presence of these two different forms of computation at the psychological level. For example, we would predict that you need to use your “mental CPU”-like capacity whenever you take on a novel task. For example, when you first learned to drive a car, you relied on a sequential, deliberate process that consumed all of your attention – at each point in time, you had to keep reminding yourself of what you were supposed to be doing. However, with sufficient practice over time, these slow, effortful processes gradually become **automated**, and you may now find yourself driving down the freeway with very little awareness of any of the underlying steps you’re effortlessly performing. This difference between the initial effortful **controlled processing** and the subsequent **automatic processing** was captured in a highly influential pair of papers by *Walter Schneider* and *Richard Shiffrin* (the same one who published the famous paper on the modal model of memory from the previous chapter) (Schneider and Shiffrin 1977; Shiffrin and Schneider 1977).

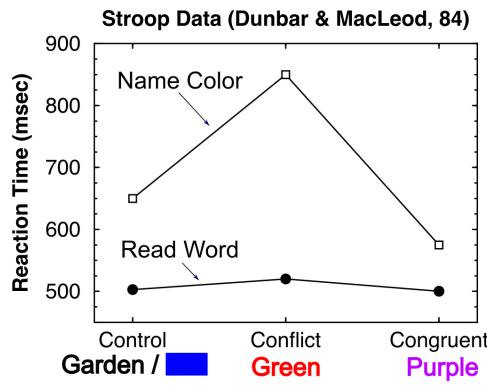


Fig 7-3: The Stroop task (Stroop, 1935; MacLeod, 1991) demonstrates the difference between controlled and automatic processing, in the context of reading words in different ink colors. If you try to name the ink color of the word “Green” when it is written in red ink (the *Conflict* condition), the automatic process of reading dominates over the relatively rare process of naming ink color, and you have to deploy controlled processing, which takes extra time as shown in the plotted data. Reading words (bottom line) isn’t affected much by the ink color – the well-trained brain networks supporting this process proceed in parallel without any supervision required. Interestingly, even when the response is identical in the *Congruent* case, the task of color naming is still slower than word reading, reflecting the extra control being exerted. The *Control* condition involves either non-color-word reading, or pure color naming.

A widely-studied example of this difference between controlled vs. automatic processing is shown in Figure 7-3 – the *Stroop* task (Stroop 1935; MacLeod 1991). The participant is instructed to either read the word or name the color of the ink the word is written in. Because word reading is so overpracticed, it is an automatic process for most adults (to the point that you often can’t stop yourself from re-reading the annoying text on your cereal box every morning). Therefore, when confronted with the diabolical *Conflict* condition in the Stroop task, where a color word (e.g., “Green”) is written in a different ink color (e.g., red), it takes extra cognitive control to prevent yourself from just blurting out the word (“Green”), when the task is to name the color (for which you have much less practice). This shows up as a significant delay (and overt errors) in this condition. Interestingly, even when the ink color and the word are *Congruent*, there is still a delay associated with naming the color – there is extra control being exerted to support this relatively unfamiliar color naming process.

If you run the Stroop task on kids, they don’t show the same effects – reading has yet to become automatic in them. Furthermore, some Stroop researchers spent enough time color naming so that it become more automatic than reading, and they showed a kind of reverse-Stroop effect! Thus, this process of *automatization* is dependent on learning and practice – over time, our brains naturally turn deliberate, sequential, controlled processing into more parallel, automatic fast processing. This is the same phenomenon that occurs for

driving and so many other tasks that you once found difficult and mentally all-consuming. Automatization is analogous to the *chunking* process discussed in the memory chapter – we have a very limited active memory capacity, but once we learn new concepts, we can greatly expand our capacity by using this limited capacity on chunks of information that used to be separate. Automatization is the process of forming *procedural* chunks – combining sequential steps into faster parallel processing that becomes relatively independent of our mental CPU – we no longer need to exert detailed conscious effort keeping the process moving along.

What it takes to be a Computer

The human brain is likely unique in having the ability to function like a Turing machine – other animals have plenty of automatic parallel processing skills, but they just don't seem to be capable of solving novel, complex tasks by performing a sequence of mental processing steps. The reason we can function like a computer is that we have some special capacities lacking in other types of brains, supplying the key ingredients of a Turing machine:

- *Program*: we use our *natural language* (e.g., English) as a kind of programming language. There is abundant evidence that we routinely use verbal self-instruction to remind ourselves of what we're supposed to do next in a complex, novel task (Miyake et al. 2004). We literally talk ourselves through the problem, and this capacity for stringing together different such verbal programs is an essential element of flexible, universal computation. It is unclear how far we might be able to get at flexible controlled processing without language, but likely not very far.
- *Active Memory* (registers, cache memory): special properties of our *frontal cortex* and *basal ganglia* give us the ability to maintain a small amount of information in active, **working memory**, as mentioned in the previous chapter. This is what you use when solving a mental arithmetic problem, by constantly juggling the digits around in your working memory. Working memory replaces the piece of paper you would otherwise use in keeping track of all the *partial products* and *control state* needed to keep progressing through a complex problem. It is also essential for maintaining the program itself, and in this way it much resembles the function of RAM in a computer, which maintains both the program and the *stack* and *heap* forms of active memory needed to carry out the program.
- *Controlled Memory Storage and Retrieval*: we also have the ability to take control over our hippocampal episodic memory system, to deliberately encode and retrieve task-relevant information as needed, playing the role of the memory tape system in the Turing machine, and a hard drive in a modern computer.

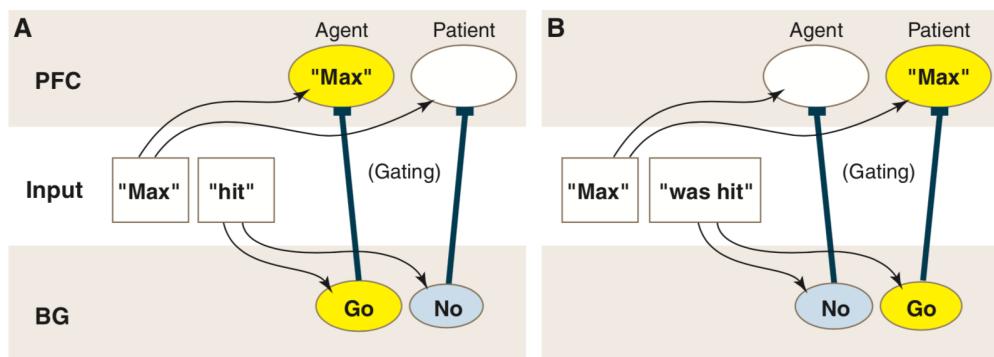


Fig 7-4: How the basal ganglia (BG) can “gate” information into different parts of prefrontal cortex (PFC) to achieve a flexible kind of variable binding like that needed for digital computers. In this case, the BG can flexibly control whether a given word / concept (“Max”) is encoded as the agent or patient in a given sentence / scenario, based on other available cues or context (from O'Reilly, 2006).

Consistent with the central role of the frontal cortex in making our mental computer work, this brain area (particularly the part of it in front of the primary motor area – the **prefrontal cortex**) is differentially expanded in humans relative to other primates and mammals more generally (Semendeferi et al. 2002). Thus, a simple story is that our unique mental-computer skills are due to this expanded brain area, but this cannot be entirely correct, because prefrontal cortex is also similarly expanded in other great apes. Despite

its expanded size, the capacity of our prefrontal working memory system (around 4 chunks) is dramatically smaller than that of even the most primitive digital computer. Thus, we are left with this rather startling conclusion: our super huge brains packed with neurons are often no match for a simple serial computational device composed of just a few basic parts – while our brains are impressive in many ways, they pale in comparison to a dime-store calculator for doing basic arithmetic!

The **basal ganglia** also play a critical role in making our mental CPU function, by orchestrating the serial, sequential steps of cognitive operations that we take in solving a problem – it turns the parallel brain into a serial system. As noted in the neuroscience chapter, the basal ganglia are critical for making a Go / NoGo decision about what to do next, and this decision-making bottleneck forces us to take discrete, sequential cognitive steps. Interestingly, the early cognitive models of human reasoning and problem solving incorporated something called a **production system** (Newell and Simon 1972), which plays the same role as the basal ganglia in the brain (Stocco, Lebiere, and Anderson 2010; Jilk et al. 2008). Figure 7-4 also shows that the basal ganglia can enable a form of flexible *variable binding* (R. C. O'Reilly 2006), which is another important property of computer systems that is otherwise hard for neurons to achieve. There are still important mysteries about how exactly the prefrontal cortex and basal ganglia learn to become a Turing-machine like system, but we do understand many of the basic principles at work already.

In summary, although we are unique in having *some* ability to perform flexible, serial, controlled processing to solve novel tasks, our brain is still running in automatic, parallel processing mode under the hood, and that greatly limits our computer-like abilities. What we lack in serial computing abilities, we try to make up for with all the amazing mental skills that we have automatized through learning and practice. This is consistent with the idea that motivation, which drives this learning, plays such an important role in human intelligence. And it also makes sense of the many ways in which our cognition differs from that of an optimal, fully rational computer system, as we discuss in a later section. Before turning to some of those issues, we first consider an alternative possibility for at least some individual differences in intelligence.

Individual Differences in Prefrontal Cortex / Basal Ganglia?

The critical role for working memory, cognitive control, and the prefrontal cortex / basal ganglia system in supporting our flexible computer-like cognitive abilities does introduce another possible explanation for individual differences in intelligence, however. It *could* be the case that different people somehow have different capacities / speeds / functionality in these particular brain systems, and that is what explains overall differences in intelligence. Furthermore, *if* this were the case, then because this flexible controlled-processing system is used as the first step in learning new skills and cognitive abilities (especially in math and other school-based learning), there could be a kind of snowballing effect where small initial differences in these brain areas could multiply over time, leading to larger overall differences in measured IQ scores. This kind of scenario is most similar to the “traditional” notions of intelligence as a fixed thing that you either have or don’t have, and obviously fits well with intuitions based on the computer metaphor of the mind.

But how well does it fit with the available data? First of all, it is well-established that major disruption to the prefrontal cortex results in impairments to controlled processing, for example on the classic Stroop task (Cohen and Servan-Schreiber 1992; Stuss et al. 2001), and on many aspects of social, moral and other forms of reasoning (Eslinger, Flaherty-Craig, and Benton 2004). However, the latter paper, and various other related findings, have shown that measured IQ scores can be relatively intact even with significant early frontal brain damage, suggesting that the relationship may be somewhat more complicated.

The most relevant question, however, is whether *normal* variation in prefrontal cortex / basal ganglia function accounts for much of the measured individual differences in intelligence? One early attempt to answer this question relied on establishing correlations between measures of working memory and intelligence, and came up with a strong positive correlation on the order of .6 to .8 (Engle 2002). However, subsequent work largely undermined that conclusion, instead suggesting that there is a separate factor for general fluid intelligence, independent of working memory capacity (Engle 2018).

Another angle on this question found that much of the measured differences in working memory capacity were actually due to motivational factors in the first place (Adam and Vogel 2016). Specifically, participants who scored lower on their working memory scale did so because they had a higher probability of “lapsing” – just failing to engage in the task on a given trial. However, when they did engage, their measured working

memory capacity was essentially the same as those who had a high working memory capacity score (due to a high overall level of task engagement). Thus, consistent with the overarching importance of motivation, it may be the critical “third variable” that drives the relationship between measured working memory and intelligence scores.

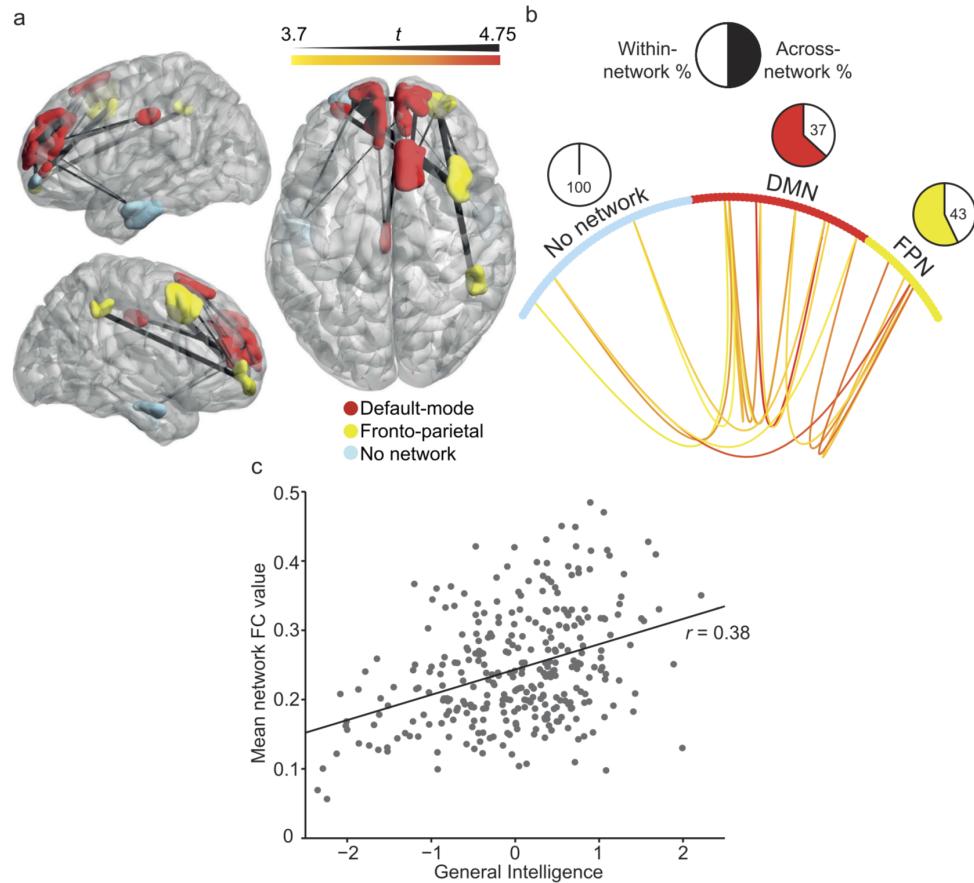


Fig 7-5: Correlation between brain activity & functional connectivity, and measured general intelligence (Hearne et al, 2016). The “default-mode network” areas in red are primary emotional / motivational areas in frontal cortex (see Figure 5-14) that are active when people are left to their own thoughts in the scanner. The fronto-parietal network in yellow are areas associated with cognitive control and working memory. Interestingly, the emotional / motivational areas make a major contribution to the overall correlation, both in terms of within-network interconnectivity and cross-network connections to the control areas.

A recent attempt to more directly find the neural correlates of individual differences in intelligence found that motivational and emotional areas of the prefrontal cortex are among the most strongly correlated with measured general intelligence (Hearne, Mattingley, and Cocchi 2016) (Figure 7-5). Thus, overall, the same answer keeps coming back up across all of these different studies: individual differences in intelligence seem to be more strongly driven by motivational factors than by the raw capacity or other properties of the prefrontal cortex / basal ganglia system. We can make sense of this result by considering that the relatively measly capacity of prefrontal cortex (only around 4 items can be maintained at a time) seems completely unrelated to the massive numbers of neurons in this brain area (which has maybe 10 billion neurons overall). Thus, the overall functional properties of this system are unlikely to be due to normal variation in the numbers of neurons or other basic biological properties of these areas. Instead, they are much more likely to be due to the degree of learning and experience that has shaped these networks to perform like a serial computer.

Strengths, Weaknesses, and Biases of our Neural Computer

The picture of human intelligence that has emerged from the above considerations has important implications for understanding how we think and reason in real-world situations, with life-and-death level consequences. To summarize: **Our brains operate best on familiar, concrete problems via prior learning and automatization, but we do have a limited ability to tackle novel problems through slow, serial, controlled processing. Also we are generally lazy and motivational factors are paramount in everything we do and learn.** This particular combination of strengths and limitations results in a systematic set of **cognitive biases** and other cognitive properties, many of which were first characterized by the nobel-prize winning psychologist *Daniel Kahneman* and his late longtime collaborator *Amos Tversky*. The interesting thing about this is that nobody seems to know.

These biases are described in terms of the reliance on **heuristics**, which are short-cut, “rule of thumb” kinds of solutions to problems, in contrast to the typically more labor-intensive (and often intractable) *optimal, rational* solutions. For example, the **representativeness heuristic** characterizes our reliance on overall similarity judgments of a given situation to a stereotype or prototypical situation that we have previously learned about. In other words, instead of going through all the mental effort of trying to truly understand a novel problem or situation, we just lazily rely on our previously-learned parallel neural pathways that *compress* something down into a seemingly well-understood high-level summary. Indeed, the representativeness heuristic is really just another name for our principle of compression.

The classic demonstration of this heuristic involves the fictional *Linda*, who is described as follows:

Linda is 31 years old, single, outspoken, and very bright. She majored in philosophy. As a student, she was deeply concerned with issues of discrimination and social justice, and also participated in anti-nuclear demonstrations.

You are then asked which of the following alternatives is more probable:

- A. Linda is a bank teller.
- B. Linda is a bank teller and is active in the feminist movement.

What do you think? Most people answer B, because it is more similar to the stereotype activated by the description of Linda. Again, we lazily rely on our basic neural processing mechanisms whenever possible, and these naturally produce a stronger match for the second characterization of Linda. But, from a rational, mathematical perspective, it is *impossible* for B to be more probable. That word *and* is mathematically equivalent to multiplying probabilities, and the probability of anyone being active in the feminist movement is, objectively, less than 1. Thus, the joint probability of being a bank teller *and* a feminist must necessarily be less than just being a bank teller!

The **availability heuristic** is similar to the representativeness heuristic, in also relying on an overall sense of familiarity, instead of doing the hard cold math. In this case, the familiarity comes from the emotional (motivational) salience of different events, instead of our compiled stereotypes (though both are very similar in relying on learned synaptic pathways – it is not clear that these are really distinguishable at a neural level). To demonstrate this heuristic, please estimate which is more likely:

- A. Dying in an airplane crash
- B. Dying in from nephritis, nephrotic syndrome, and nephrosis

Many people avoid flying in airplanes due to fear of crashes, which inevitably get extensive media coverage, and are highly salient and gruesome. It is not clear if there has ever been any major media coverage of the latter cause of death, or if you even know what it is. Yet it is the 9th highest cause of death, at over 50 thousand in the USA in 2016, according to the CDC: [CDC.gov stats page](#). In contrast, there were only 59 total deaths *worldwide* due to airplane crashes in 2017!

The real-world life-and-death implications of our reliance on the availability heuristic were starkly illustrated in an analysis by *Gerd Gigerenzer* of people’s avoidance of flying after the September 11, 2001 terrorist attacks in the US (*Gigerenzer 2006*). He found that there were roughly 1,500 additional deaths caused by people choosing to drive instead of fly – driving is much, much riskier than flying, and yet because it is so familiar and commonplace, people vastly underestimate the relative risks.

Gigerenzer and colleagues have many other studies showing our general incompetence in dealing with statistical information, again with real-world implications. For example, even highly-trained doctors make

significant mistakes interpreting the statistical results of medical studies on health risks and probabilities. One of the most important errors we make is known as **base rate neglect**. For example, if you hear that drinking alcohol increases your risk of death by 12%, that sounds like a big effect. But people routinely neglect to take into account that the base rate against which this percentage is measured is extremely low. In practice, the difference amounts to just 4 additional deaths per 100,000 people per year: [New York Times article](#). This is really an instance of the *contrast* effect – we focus on differences but not on the raw values.

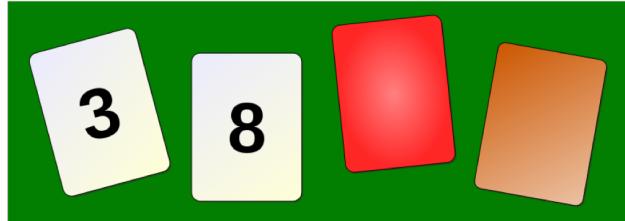


Fig 7-6: Wason card selection task. Your job is to verify whether the following rule holds for these cards: If there is an even number on one side, then the other side is red. Which cards would you turn over to test whether this rule holds?

Another important demonstration of our strong preference for reasoning about familiar, concrete situations comes from the **Wason card selection task** (Wason 1968) (Figure 7-6). Here, your job is to verify the application of a particular rule about what is on different sides of the cards. Look at the figure and formulate your answer. Now, consider this alternative problem:

You are a bartender. If you are going to serve alcohol to someone, they need to be over 21. Who do you card?

Critically, you card someone who is *under* 21. Did you likewise decide to flip the *brown* card in the Wason card selection task? Probably not. But think about it: if a card has brown on one side, and it has an even number on the other side, then the rule is incorrect. The brown card is the “cheater”. You probably said that the red card should be flipped. But the rule says nothing about the other side of a red card – it could be either even or odd. Likewise, people over 21 can either drink or not drink – it doesn’t matter. You don’t card people who are obviously over 21.

This task again demonstrates that our ability to perform abstract logical reasoning is extremely limited, and we do much better with familiar, concrete cases like the bartender example. By contrast, a standard digital computer programmed with the basic rules of logic can easily get the abstract form of this problem correct, and in fact would struggle understanding a question like “who do you card?” without much more explicit specification of what that means. Thus, people are really good at “common sense” reasoning, and pretty bad at abstract reasoning, and computers are generally the opposite.

Task Transfer and Education

Another very important consequence of the concrete nature of our brains is that things we learn in one context often do not **transfer** very well to another context. For example, classic studies of problem solving tasks such as the *Tower of Hanoi* have found that, having figured out a solution to that problem, people are generally not very good at applying that very same solution to another version of the exact same problem, portrayed in a different “skin” (i.e., differing in only superficial, task-irrelevant factors) (Kotovsky, Hayes, and Simon 1985). Although transfer can occur, and there are reliable ways to make it more likely to occur (Anderson, Reder, and Simon 1996), it is clearly not the natural state of the system. From a neural perspective, the brain learns everything in terms of the specific patterns of neural activity present during the learning episode, so the only way to get knowledge to transfer to a novel situation is to ensure that these patterns of neural activity are shared across situations.

These and other related findings have led to the development of the **situated learning** theory (Lave and Wenger 1991; Greeno, Moore, and Smith 1993), which emphasizes the concrete, specific nature of human learning. Although there have been some debates about the generality of this framework (Anderson, Reder, and Simon 1996) – people *do* have some ability to engage in abstract reasoning and *some* knowledge does transfer – the basic principles remain solid and are increasingly incorporated into educational strategies.

A recent controversy about transfer of learning has arisen in the context of the increasingly popular brain training programs, such as that offered by [lumosity.com](#). These programs are based on the idea that you can train up your brain like a muscle, and increase your overall intelligence as a result. However, consistent with the overall difficulty in transferring learning, (Simons et al. 2016) review a number of studies showing that there is very little transfer of this brain training beyond the specific tasks that you practice. So, you can get better at the specific arcane puzzles that you practice, but, unfortunately, it doesn't really transfer to make you generally more intelligent. And lumosity was successfully sued for making false, misleading claims to the contrary – that's psychology in the real world!

Programs in the Mind: Problem Solving and Reasoning

Despite all the denigration of our capacity for abstract reasoning, and the relatively modest power of the neural CPU in our brains, we nevertheless can do some impressive feats of cognition. At least, some of us can! After all, Turing was able to prove the universal nature of his computational machine, using his own relatively modest neural CPU (he didn't have a real computer to work with yet). In this section, we review some of the kinds of things we can do with our neural CPUs, and how these have been studied in psychology.

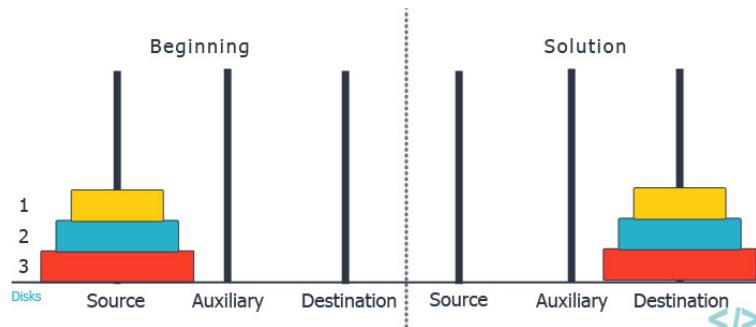


Fig 7-7: Tower of Hanoi task: Your job is to move all the discs from the source to the destination, one at a time, without ever putting a larger disc on top of a smaller one. Go for it!

As noted earlier, some of the pioneering early work using the computer as a model for the human mind was conducted by *Allan Newell* and *Herbert Simon* at Carnegie Mellon University (CMU) in Pittsburgh, PA (Newell and Simon 1972). These scientists and their colleagues focused on how people solve various kinds of challenging puzzles and games, including the *Tower of Hanoi* (Figure 7-7), and chess. When you first start doing any kind of puzzle, the initial strategy is often a semi-random **trial-and-error** process of trying out different actions, and seeing how it goes. Another name for this is **hill-climbing** or **gradient ascent**, where you measure success in terms of the visual similarity of the current state to the target state. In the case of Tower of Hanoi, this results in trying to put discs on the destination peg as early as possible, which turns out to not be such a great idea.

Interestingly, good puzzles often have this characteristic of requiring you to move further *away* from something that looks like the final solution, in order to solve the problem. In other words, they specifically thwart the natural hill-climbing solution. This ability to move away from the goal has been studied in babies and animals using a *transparent barrier detour* task, where a desired object (food or toy) is hidden behind a transparent barrier – the direct reach solution must be rejected in favor of the indirect reach-around behind the barrier (Diamond 1990; Wallis et al. 2001). These studies show that the prefrontal cortex is important for this ability to overcome the direct approach in favor of the indirect reach, consistent with the idea that these kinds of challenging puzzles more generally tap our higher-level flexible cognitive processing.

This same requirement for overcoming the initial “obvious” solutions to a problem is featured in **insight problems**, where the surface features of the problem strongly suggest one type of solution, but another, different kind of solution is actually required (Figure 7-8). The relative difficulty that we have breaking out of a given **mental set**, known as **functional fixedness**, is nicely illustrated by these problems. But really, it is just a more advanced version of the same problem in all of these puzzles: the first ideas don’t work, and you need to hunt around for an entirely different approach. This has entered the popular lexicon in terms of

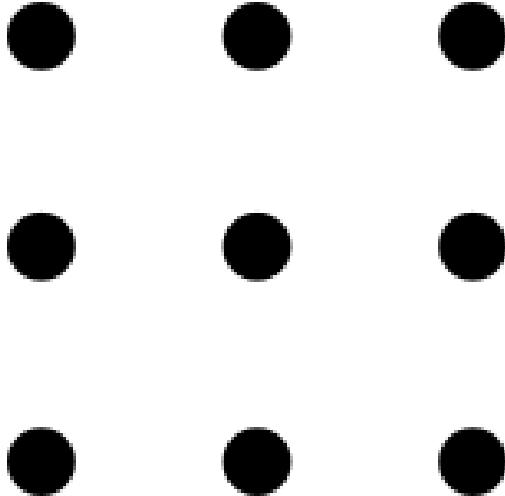


Fig 7-8: Classic “thinking outside the box” insight problem. Connect all of the dots using four straight lines, without lifting up your pen.

“thinking outside the box”, likely in reference to the puzzle in Figure 7-8.

Games such as chess provide the opportunity to explore more advanced strategic thinking. Early AI approaches to solving these games involved the “obvious” strategy of figuring out what your opponent might do if you make a given move, and so on, to pick the move with the best potential future outcome. This is known as a *look-ahead* search algorithm, and is essentially what the Deep Blue chess playing computer employed when it beat the famous chess champion, Gary Kasparov. The main challenge with such an algorithm is the *combinatorial explosion* problem associated with the very large *state space* of a complex game like chess – there are so many different possible future move strategies, that searching all of them quickly becomes computationally prohibitive. The Deep Blue computer basically used massive numbers of computer chips to search this space in parallel, executing a pure “brute force” solution.

The limited capacity of our neural CPU prevents us from using such a brute force strategy. Instead, research from Simon and colleagues at CMU showed that people use a much more perceptually-based strategy (Gobet and Simon 1996). Specifically, chess experts leverage their massively-parallel neural networks to recognize good and bad board positions, based on extensive learning experience with the game. The recent advances by the Google DeepMind team on AI systems that play the ancient game of *Go* essentially combine this perceptual expertise strategy with the brute-force look-ahead strategy, to create a human-crushing Go master machine (Silver et al. 2017). This system provides a good demonstration of the power of the combination of parallel neural-like computation, with more flexible CPU-like processing. However, unlike people, the DeepMind model did not learn everything from the ground up – it was programmed and otherwise designed strictly for playing *Go*, with the rules of the game coded directly into the look-ahead part of the system. Nevertheless, it did discover novel strategies by playing millions of games with itself.

So do we think that the DeepMind system exhibited human-like insight and creativity? Or is it really just another example of brute-force searching of the state space of the game, which happened to turn up novel strategies? And how do we really differentiate this from what people do when they come up with novel insights to challenging problems? One critical difference is that people have an awareness of the strategies they are trying, and can explicitly formulate particular aspects of the overall puzzle, to guide their insight process along. Thus, we are using something more like a **means-ends** analysis of the problem – reasoning backward from target solutions (the “ends”) to the means that should lead to those ends. Also, we can deploy these creative reasoning skill across many different domains. We’ll see how long it takes for machines to match or exceed our abilities in these respects!

Measuring Intelligence and its Implications

Now that we have some understanding of the nature and scope of human intelligence, we turn to the controversial history of IQ testing. This history is controversial because the central question of the genetic versus environmental basis for intelligence has been at the center of these tests from their beginning, and because of the major question of the potentially biased nature of these tests. Briefly, one of the earliest tests was developed by *Alfred Binet* in France in the early 1900's, and translated and further developed by *Lewis Terman* at Stanford University, resulting in the **Stanford-Binet** IQ test. This test was developed based directly on academic skills, and includes factors such as knowledge, quantitative reasoning, visual-spatial processing, working memory, and fluid intelligence. The test was standardized for children at different grades, and was initially focused on identifying children who were outside of the normal range for each grade.

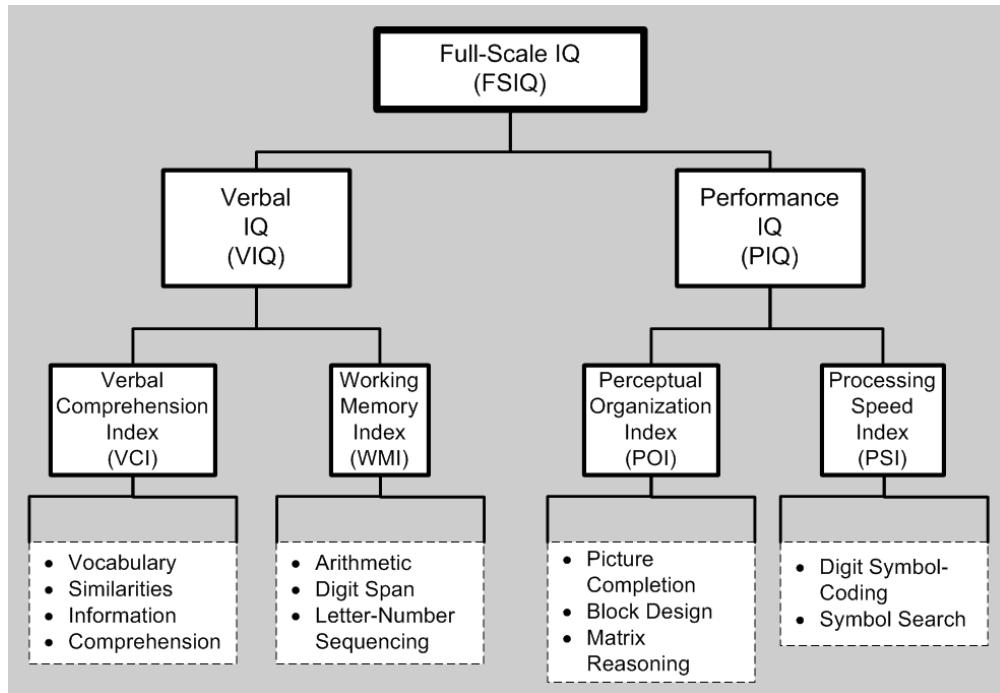


Fig 7-9: Sub-scales of the WAIS intelligence scale, divided into verbal vs. non-verbal (“performance”) components.

The later **Wechsler Adult Intelligence Scale (WAIS)**, first developed by *David Wechsler* in the 1930's, is still widely used to this day, and is by far the most popular IQ test. It attempted to improve over the Stanford-Binet by separately testing verbal and non-verbal (“performance”) intelligence (Figure 7-9), and relying less on raw speed as a factor (though it is retained as a specific sub-factor). The WAIS is similar to the Stanford-Binet test in using more of a “scattershot” approach to measuring intelligence – many different specific tests are combined together to paint an overall picture. And, consistent with the essential role for working memory in supporting flexible CPU-like processing, working memory constitutes one of four major subscales in the WAIS.

An important feature of the WAIS measure is that it has been **standardized** so that a score of 100 is *defined* as having a precisely average level of intelligence, and intervals of 15 points represent a standard deviation according to the normal (gaussian) distribution. Thus, someone with an IQ of 115 is one standard deviation above the mean, and is thus measured as being “smarter” than 84.13% of the population, and someone with an IQ of 130 is “smarter” than 97.72% of the population. This standardization applies to “adults” aged 16 or older, and there is a different scale for children.

One of the major sources of controversy about IQ tests such as the WAIS is the extent to which it might be systematically biased against different populations, such as women and minorities. This is a difficult question to answer, because there is no other accepted standard to compare against. In this case, we can at least address three important features of the test: its **reliability**, **construct validity**, and **predictive**

validity. The reliability is assessed by repeatedly testing the same people, ideally with different versions of the same overall test, and measuring the consistency of their scores across tests. Construct validity is much more subjective and difficult to assess – basically it is a judgement call on the part of scientists that the individual test components are actually measuring something useful about the construct of intelligence. In any case, it is somewhat circular: the test measures those specific abilities that it actually tests, and your score reflects your ability to perform those specific tasks. No specific test can ever measure anything other than what it specifically tests.

Thus, the most important property of the test is its *predictive validity*: how well does it predict *other* things about people? This is the most relevant and directly measurable property of an IQ test – if it is really measuring intelligence, then it should be able to predict how well people do at things that we generally agree require intelligence. Here, the evidence is a bit mixed. For example, your IQ test score predicts subsequent school performance (i.e., grades) with a correlation factor of about .5 (Neisser et al. 1996). This is actually not a very strong correlation – even though it sounds like a “half strong” correlation, you have to square this number to determine how much total variance in grade outcomes are accounted for by IQ – so only 25% of the total variance. This means 75% of the variance is *not* attributable to IQ – by far the majority of it. Nevertheless, it is probably the best *single* predictor of academic performance. Interestingly, in discussing the strength of this factor, (Neisser et al. 1996) point largely to motivational factors as likely additional determinants of academic performance.

The amount of variance predicted by IQ scores on factors such as job performance, income, socio-economic status (SES), health, and social status are all lower than for school performance, and it is very difficult to accurately factor out the contributions of parental factors (e.g., parent SES), which affect both IQ and most of these other factors (this point also applies to the educational predictiveness).

Importantly, the predictive validity of IQ scores does *not* differ across different groups, such as women or minorities (Neisser et al. 1996). This provides one way of assessing whether IQ tests are inherently biased, and the results suggest that they aren’t biased at this level. However, just looking at the test from a construct validity perspective, it definitely does test lots of knowledge that is typically learned in school, and must be expressed verbally. Thus, an individual’s prior schooling environment is undoubtedly going to have a significant impact on their IQ score. Furthermore, as noted above, there are significant correlations between parental SES and IQ, at about .33 (Neisser et al. 1996). Thus, there is no doubt that your raw IQ score reflects a strong contribution from various environmental factors that likely systematically vary among different groups. Interestingly, the predictive validity factor is not directly affected by such overall mean differences – it is only affected by the *variance* in scores across individuals.

In summary, IQ is both a biased and fair test of intelligence! It is biased at the mean level, and by the very fact that *any* test of complex task performance is likely to be affected by relevant SES-level factors. But at the level of predictive validity, it is fair in that there aren’t significant differences across groups. We’ll consider the data on the genetics and heritability of intelligence in the chapter on genetics and development.

Multiple intelligences

One of the major debates about IQ measures is the extent to which there are really distinct aspects of intelligence, or rather a single underlying **general intelligence factor**, typically denoted by the letter *g*. *Charles Spearman* was the main original advocate of the importance of this *g* factor – he noted that performance across the different subscales of IQ tests is positively correlated, and attributed this common factor (*g*) as being the “true” underlying meaning of intelligence. Arguing the opposite case was *L. L. Thurstone*, who advocated a multi-factorial model of intelligence. Interestingly, both were looking at the same data, and just interpreting them differently – there are definitely multiple separable intelligence factors in addition to the common variance across all of them associated with *g*.

A distinction that is less subject to these whims of interpretation is the difference between **crystallized** and **fluid** intelligence, as proposed originally by *Raymond Cattell*. Crystallized intelligence refers to the accumulated knowledge and skills learned over experience, whereas fluid intelligence refers to the ability to actively juggle information in your mind (e.g., in the service of mental arithmetic). Thus, fluid intelligence corresponds to the contributions of the prefrontal cortex and basal ganglia, supporting working memory and the ability to shuffle information rapidly around within your mental workspace. By contrast, crystallized

intelligence refers to the accumulated synaptic changes from learning. These two forms of intelligence are often considered from an aging perspective, where crystallized intelligence generally increases over the lifespan (i.e., wisdom accumulates over time), while fluid intelligence unfortunately declines. This decline is thought to mirror the declines in dopamine levels with aging, which may suggest that there may be an important motivational component to the story as well. Could a strongly-motivated senior still muster the computational horsepower of a young adult?

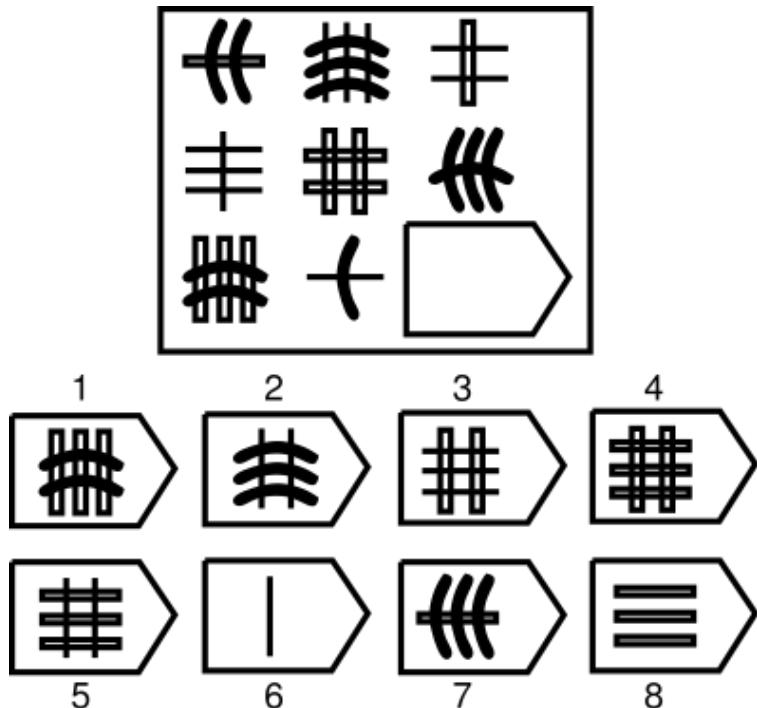


Fig 7-10: Raven's progressive matrices task, which resembles the matrix reasoning component of the WAIS, and is often used as the single best measure of general fluid intelligence. Your task is to figure out which of the options at the bottom best completes the pattern shown in the *matrix* at the top – row-wise, column-wise, and the diagonal progressions of the pattern are all relevant! This need to juggle multiple different sequences makes this a particularly challenging task.

One of the best single measures of general fluid intelligence is the Raven's progressive matrices task, shown in Figure 7-10 (Raven, Court, and Raven 1977; Conway et al. 2002). This task is also similar to the matrix tasks used in the WISC, to assess non-verbal intelligence. The non-verbal nature of the task likely helps to keep it from being “contaminated” by crystallized knowledge factors, and thus contributes to its ability to more directly measure the construct of fluid intelligence. It clearly requires considerable juggling of information in working memory, to figure out which of the possible patterns best completes the missing cell in the matrix. Furthermore, this Raven's task and the fluid intelligence construct have been closely associated with working memory function (Conway et al. 2002).

As discussed earlier, many people interpret this notion of fluid intelligence and working memory function in computer-like terms – some people just have better “RAM” than other people. However, as we concluded earlier, the available evidence more strongly implicates motivational factors as paramount. The revised interpretation then becomes: general fluid intelligence measures reflect the extent to which an individual is willing to exert significant cognitive effort on arbitrary lab tasks. People who are more willing to do this score higher on these measures, and, perhaps not coincidentally, also tend to do a bit better in school overall – school similarly requires you to allocate cognitive effort on things that you might not otherwise want to spend time and effort on.

Control

This brings us back to the recurring theme of this chapter: the importance of motivational factors in understanding intelligence, and the functions of the prefrontal cortex and basal ganglia more generally. Ultimately, it comes down to a question of *control*. These brain structures are critical for cognitive control, but also for overall control of your entire self, at all the relevant levels. The ventral and medial areas of prefrontal cortex receive extensive inputs from brainstem emotional and body-state areas, and integrate this “hot” state information with “cold” planning and sensory-motor control strategies, all in the service of keeping *you* in control of *you*.

Ultimately, control is about achieving your own goals, satisfying your own needs, etc. Thus, when scientists bring you into the lab, you really are a *participant* in the experiment, and the degree to which different people actually participate varies according to their motivational goals. If someone is really interested in what is going on inside their brain, and motivated to demonstrate their own sense of mental superiority, they may devote considerable effort to a given lab task. But others may have “more important” things on their minds, and spend less overall effort.

Neuroimaging studies have provided a vivid window onto this battle between your own internal desires and goals, and those of the psychology experimenter. In any given experiment, your brain can be seen switching between a “task engaged” mode and the “default mode” (Fox et al. 2005). The default mode (see Figure 7-5) involves all the motivational / emotional areas of the frontal cortex, interconnected with the hippocampus and other brain areas that are relevant for thinking over important recent events and planning upcoming activities. Basically, the stuff you think about when left to your own devices. By contrast, the task-engaged brain areas involve the working-memory and problem-solving, cognitive control areas that are needed to keep you focused and solving complicated tasks.

As we saw in the working memory capacity results discussed earlier (Adam and Vogel 2016), the major determiner of measured overall capacity is how often people were willing to actually engage in the working memory task. Presumably, the rest of the time they were thinking about more self-relevant things, and exerting their own personal control over what they allocate their attention to. Likewise, you may have drifted off into your own mindwandering space while nominally reading this chapter – we can throw the words at you, but ultimately you are in control of whether you want to read them!

Finally, the need for control, and motivational factors more generally, may be the most important source of cognitive biases, dwarfing those mentioned earlier. For example, in the political arena, it is typical for people from different parties to each have strongly-held and mutually incompatible views about the very same facts and events. Are humans contributing to climate change? This question is not really in doubt scientifically, but because it has become politicized, it has become question of emotional significance, tied directly to one’s own personal sense of self and identity. *Dan Kahan* and colleagues have demonstrated the power of these kinds of motivational, cultural identity-related factors in biasing people’s thinking (Kahan, Jenkins-Smith, and Braman 2011), and the effect sizes are quite large. More generally, when you see someone behaving in a way that you think is blatantly irrational (e.g., people voting for political candidates whose positions directly disenfranchise them), it is much more likely that they are doing that because of strong emotion-laden values, than because of some kind of error in statistical reasoning. In other words, of the three C’s, Control is by far the most powerful force shaping our behavior, even though Compression and Contrast also play important roles.

Summary of Key Terms

This is a checklist of key terms / concepts that you should know about from this chapter.

- Fixed vs. growth mindset and the malleability of intelligence
- Algorithm = program for solving a given problem (e.g., mental arithmetic)
- Automatic vs. controlled processing
- Stroop task: word reading is automatic compared to color naming that requires controlled processing
- Neural CPU:
 - Working memory supported by prefrontal cortex = active memory / RAM in a computer
 - Basal ganglia = production system = sequentializing cognitive steps

- Natural language (e.g., English) = program
- Cognitive biases / heuristics: shortcuts that leverage strengths, avoid hard-to-compute optimal, rational solutions.
 - Representative heuristic: use similarity to prototypes / stereotypes instead of actual statistics.
Compression!
 - Availability heuristic: use familiarity instead of actual statistics.
 - Base-rate neglect: focus on percents instead of overall probability. Contrast!
 - Better at concrete vs. abstract reasoning: Wason card selection task
 - Transfer of learning: not much = situated learning: learning is specific to situations
- Problem solving:
 - Strategies: trial-and-error, hill-climbing / gradient ascent, means-ends
 - Puzzles designed to frustrate obvious / hill-climbing solution
 - Insight problems: mental-set and functional fixedness
- Intelligence tests:
 - Stanford-Binet
 - Wechsler Adult Intelligence Scale (WAIS) – standardized scale (100 mean, 15 standard deviation)
 - reliability, construct validity, predictive validity – importance of predictive validity for
- Multiple intelligences
 - Generalized intelligence factor g
 - Crystallized vs. fluid intelligence
 - Raven's progressive matrices as non-verbal test of fluid intelligence
- Control:
 - Cognitive control in service of overall control

Chapter 8: Language

Past tense, rules vs. not. Piaget / Pinker vs. Connectionists. Neural CPU suggests both are right.
add pointer back to this discussion in development chapter, about neural CPU, rules, etc.

Chapter 9: Origins: Evolution, Genetics, and Development

Having now explored the full scope of human cognition, from perception to the highest levels of intelligent reasoning, we now circle back and consider the incredibly fascinating and challenging question of *origins* – where does everything a mature adult can do come from in the first place? In the learning chapter, we touched a bit on this question, arguing that the genetic code could not possibly contain enough information to directly specify the synaptic connection strengths of even a small fraction of the billions of cortical neurons that support all of our cognitive functions. Thus, learning must play a critical role. And yet, we don't yet have fully satisfying, widely-accepted accounts of how this learning unfolds, and how the considerable genetic shaping of the basic organization and wiring of the brain interacts with these learning processes to determine how we develop.

Thus, development remains more of a mystery than most of the other topics that we cover. But certainly much is known about the *phenomenology* of development, e.g., the general chronology of when different capacities begin to emerge. Furthermore, we also know a great deal about genetics and how the blueprints of life are encoded in our genes, and understanding the nature of these mechanisms is essential for having a more complete picture of how biology can shape and constrain our developmental processes. We begin our exploration of these issues at the natural beginning: with the process of evolution – the origin of everything!

Evolution

Evolution is an absurd theory on the face of it. How is it even remotely possible that a beast such as a fish could *ever* give birth to something that was *not* a fish? And even if it did, how could that freak ever mate with anything else to propagate its non-fishiness? Two freaky non-fish emerging at the same time, in close enough proximity to propagate this new species? It just doesn't add up. Presumably everyone has heard the basic principle of *survival (and reproduction) of the fittest* as the primary engine of **adaptation** in evolution, but it seems that this alone does not really explain things. Undoubtedly, the sheer implausibility of the ideas at this basic level contributes at least something to the continued reluctance that people have toward this theory.

There are several things we can work through to help resolve some of these conundrums. First, we need a clear overall framework for thinking about biology, and what we are actually made of – fortunately, LEGO provides a really nice, familiar analogy. Second, we need to think more carefully about the actual nature of evolutionary change: it is very incremental, and yet somehow can produce large changes over time – reconciling those two points is difficult within our limited scope of experience. Here, computational models can be immensely useful in providing the missing “long timeframe” perspective, by simulating many generations of evolution within a few minutes of human time. Furthermore, a few key principles can go a long way toward resolving the basic problems stated above.

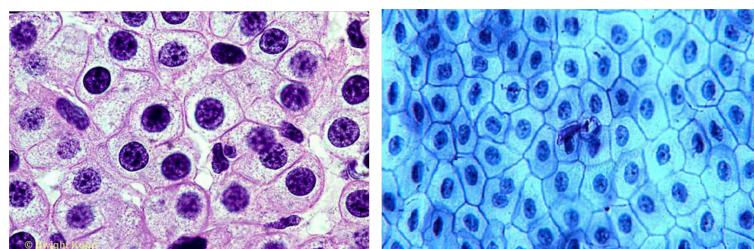


Fig 9-1: Which picture shows cells from a human and which from a frog? At the cellular level, all animals are essentially identical, built from the same building blocks – just like you can build so many different things using the same set of LEGO blocks. Understanding this makes it easier to see how different bodily forms can emerge through evolution.

When you think of the manifold differences between a fish and a human being, evolution seems impossible. However, these differences disappear when you zoom in to the cellular level, where it immediately becomes clear that all animals are made from the same basic building blocks (Figure 9-1). This is really the level at which evolution is operating – the macroscopic forms of organisms (which is what we perceive) are more of



Fig 9-2: Biology is remarkably like LEGO: the same small set of basic building blocks can be recombined to produce many different animals. This continuity across all animals makes it clearer how evolution can work: you just need to tweak the instructions a bit and you can end up with entirely different animals, as in this 3-in-1 kit that uses the same parts, with different instructions, to make very different animals. In biology, genes are the instruction set, and these instructions are constantly subject to random changes, producing novel “experiments of nature”.

an emergent result of cellular and sub-cellular processes unfolding over the developmental process, just like our cognition is an emergent product of our neurons interacting in complex ways.

The easiest way to understand this is in terms of the more familiar process of building different things out of LEGO blocks (Figure 9-2). A very small set of basic building blocks can be used to make an essentially infinite number of different objects. Furthermore, the developmental process is akin to the process of following the step-by-step instructions to assemble the final product. As you may have experienced, small random errors at an early stage can have major implications later. Likewise in biology, relatively small changes in the developmental process can lead to major changes in the overall shape and function of the organism. To make a bigger brain, you just have to tweak the process that controls how long nerve cells divide and proliferate, in the same way you would just keep adding blocks to a LEGO wall to make it bigger. Thus, major macroscopic changes can emerge from relatively minor changes in the program.

In biology, the genetic code is the equivalent of the LEGO instruction booklet, and we'll see below that much of the genetic information is devoted to controlling the timing and coordination of the building process (as compared to the raw bricks themselves), and this is likely where much of the action has taken place over the course of evolution.

Perhaps it is now clearer how very different-looking organisms can emerge from relatively minor tweaks to the genetic code, and how in fact a fish very likely could produce some rather different offspring even in a single generation. Indeed, there are plenty of examples of “genetic freaks of nature” that can be found on the internet (e.g., [list25.com 25 disturbing freaks](http://list25.com/25-disturbing-freaks)). But the other puzzles remain: how could that weird offspring itself procreate, and lead to the origin of new species, and more complex, sophisticated biological machinery? This is really the hard problem of evolution that Darwin tackled, and later theorists such as *Stephen J. Gould* wrestled with.

One key idea here is that gradual, continual change from one generation to the next is much easier to understand than some kind of entirely new beast emerging whole within a single generation – even though random changes in the genetic code can produce such things, they indeed aren't likely to have survived and procreated. But how does gradual, incremental change ever lead to the emergence of something dramatically, qualitatively different? We can consider two cases: emergence of air breathing, and moving onto land from a fish, and the emergence of flying birds from dinosaurs.

In the case of the air breathing, one could easily imagine that there was some gradual advantage to being

able to gulp some extra oxygen from the surface of the water – such fish would have more energy, and this would be especially important if the water ended up being less oxygenated, e.g., as happens in algae blooms. Also, there is plenty of evidence that proto-lungs emerged for various other digestive and buoyancy-control functions. Then, as this ability was selected for over time, more and more oxygen could be processed, until at some point, such animals could actually survive by breathing air. If they happened to get trapped on land, as often happens during tidal cycles, then you end up with all the right ingredients for one of the most dramatic events in evolutionary history: the transition from the sea to land. Of course, we don't know how it actually happened for sure, but by thinking about how gradual *quantitative* processes can eventually turn into a major *qualitative* changes, at least it may seem more plausible.

In the case of bird feathers and flight, we now know that many dinosaurs had feathers, presumably for temperature regulation. Thus, like the oxygen example, something can be adaptive in one way, before it becomes “co-opted” for something else entirely – Gould referred to this as an **exaptation**, while Darwin referred to it as *preadaptation* (which sounded a bit too “prescient” for Gould). Thus, the dinosaurs that climbed trees, and had feathers for keeping cool, eventually discovered that they could glide their way to safety using these same feathers, and a whole new adaptive function and corresponding gradual selection pressure emerged, to make better and better wings.

Putting all this together, this initially crazy-sounding idea perhaps makes more sense. It just takes a lot of time and a lot of random accidents and amazing stories of survival – those challenges are the “desirable difficulties” of evolution that have driven survival of the fittest to select some radical new innovations. It took many such apocalypses to get us where we are today, and it is thoroughly mind-blowing to try to grasp this idea that we are just the latest spawn in a “great chain of being” stretching back billions of years. But look again at Figure 9-1: at this level, not that much has really changed in all those years!

To really see evolution happening in “real time” that we can actually comprehend, people have developed computer simulations of evolution, and it definitely works! Indeed, there is a quote that neural network algorithms (modeled on the functioning of the brain) are the second best solution to any complex problem, and the third best solution is a **genetic algorithm**, which is the general computer-science version of evolution. Randomly searching a complex, high-dimensional space, and combining the features of the best-functioning exemplars, works really well for finding novel, previously-unimagined solutions to complex problems. Interestingly, both genetic algorithms and neural networks have the same property of following *gradients* (i.e., hill climbing, as was discussed in the case of problem solving strategies earlier) – this is the core property of evolution where some adaptive property gets “optimized” over successive generations.

Jeff Clune and colleagues have developed particularly compelling, bio-mimetic simulations of the evolution of organisms, often with very funny but functional properties: [YouTube video of Evolving Soft Robots](#) – this is really a must-see video, and should hopefully make evolution come to life in a unique and compelling way.

Genetics

As we mentioned above, the biological equivalent of the LEGO instruction booklet is the *genome* – the collection of genes that determine how everything in our bodies (and every other living organism) is built. It is amazing how much progress has been made in understanding the genetic basis of biology since the structure of DNA was discovered in the early 1950's, and published in 1953 by *James Watson* and *Francis Crick* (who later became somewhat of a neuroscientist, confirming that the brain is the most fascinating thing in the universe, and the “last refuge of scoundrels”). And it is also amazing that this all happened so recently – we are the first few generations of beings that now know (more or less) how that great chain of life actually works.

At its base, the genetic program is remarkably simple: there are only 4 different letters in the language of life: G, C, T, A, which are always paired GC and TA (Figure 9-3). These 4 **base pairs** don't do much by themselves – it takes 3 of them in sequence to determine a corresponding **amino acid** (of which there are 20 different varieties coded in the DNA). Interestingly one of these amino acids is *Glutamic acid*, which is the basis of both the ubiquitous excitatory neurotransmitter glutamate, and MSG, and another is *Tyrosine*, which is the direct chemical precursor of the major neuromodulators dopamine, epinephrine, and norepinephrine. Thus, like the sodium and chloride ions, the key ingredients that make the brain tick are really basic and universal in biology, coded directly by just 3 base pairs in our DNA.

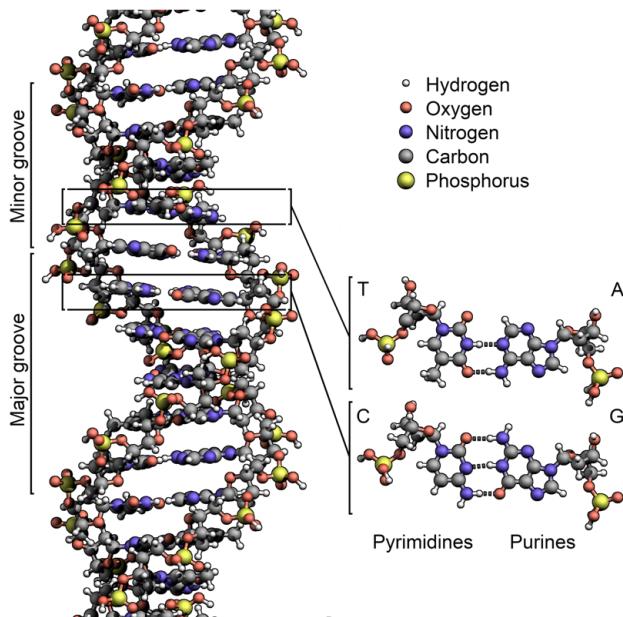


Fig 9-3: The molecular structure of DNA, which has an alphabet of only 4 different letters: G, C, T, A. They are paired with each-other as shown, but any given strand could have any one of these different letters in any order.

The more complex building blocks in biology are composed from sequences or polymers of these amino acids, i.e., **proteins**. These are the LEGO blocks of the body, and our DNA directly codes for how to build these blocks via those sequences of 3 base pairs. However, only a tiny fraction (about 1.5%) of our total DNA actually codes for these proteins, of which there are about 20,000 different types. Proteins come in various lengths, from a few hundred up to 20,000 amino acids in size. However, the *genes* that code for these proteins are typically *much* longer than the minimal number required to code for the literal amino acid sequence – there is clearly a lot more going on in the genome than just the literal coding of amino acids to build proteins, even among the tiny fraction of genes that actually code for proteins in the first place.

Thus, there is quite a gap between the very concrete, well-understood level of DNA, amino acids, and proteins, and the rather fuzzier notion of **genes**. Genes are defined functionally as units of **heredity** – the basic elements that we can inherit from our parents, and the simple idea that each gene codes for a different protein holds in some cases, but only a tiny minority. Again, LEGO is incredibly helpful in understanding why this might be the case. When you’re building a lego kit, the vast majority of the information in the instruction booklet concerns *where* and *when* to place the bricks, with only a relatively tiny bit of information concerning the different bricks that are available to build with. For example, there may be around 100-200 different types of bricks in a typical reasonably complex LEGO kit, but the amount of information it takes to specify exactly where to place those bricks is much greater (and somewhat difficult to quantify, given its visual nature).

In short, although the remaining 98% of the human genome that does not code for proteins was originally characterized as “junk DNA”, it is highly likely that most of it is serving a vital function akin to the bulk of the LEGO instruction booklet: determining when and where to build all those proteins. For example, we know that much of the extra “junk” within a given protein-coding gene plays a *regulatory* role, shaping the complex process that actually *transcribes* the DNA sequence into a corresponding amino acid sequence to make up the proteins. Perhaps we will ultimately find the equivalent of a computer programming language embedded in all this regulatory DNA, complete with **if-then** rules and **for loops**, etc. As with the LEGO instruction booklet, most of the power and functionality in a computer program comes from these kinds of control structures, rather than the raw “proteins” (e.g., numbers, characters or other data) that is being manipulated.

Sexual Reproduction

Aside from the pure intellectual fascination of understanding the amazing genetic machinery that makes us tick, the practical relevance to psychology and neuroscience comes in understanding how our genetic information is inherited from our parents, and how much of our overall brain function it ends up determining. This is the domain of **behavioral genetics**, which traditionally has been performed by comparing **identical (monozygotic)** versus **fraternal (dizygotic)** twins, and is now also able to leverage the advances in **molecular genetics** to directly compare genetic material across people. To understand how all this works, we first need to understand how **sexual reproduction** works!

Monozygotic twins started out from the same **zygote** (fertilized egg cell), and they thus share 100% of their DNA, whereas dizygotic twins came from two separate egg cells, and have the same genetic similarity on average as any siblings born from the same parents (i.e., 50%). The exact process by which this 50% genetic similarity arises is surprisingly complex. Likely everyone remembers hearing about *meiosis* and *mitosis* from high-school biology, but you probably don't remember all the crazy details.

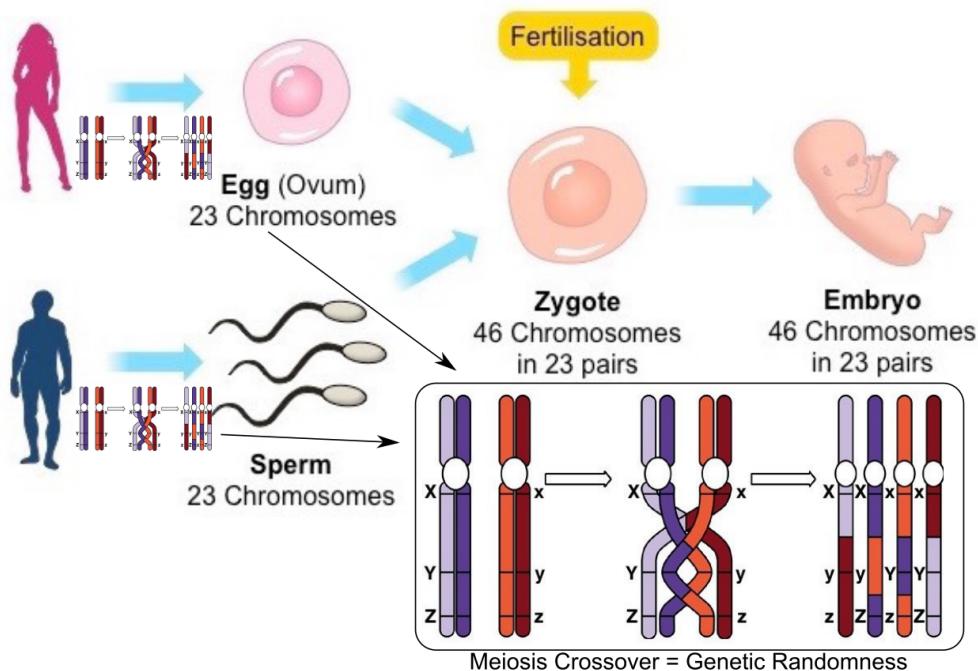


Fig 9-4: The source of genetic randomness in sexual reproduction comes in the making of the gametes during crossover in meiosis (not during fertilization).

Interestingly, all of the genetic shuffling responsible for producing the essential randomness that powers evolution happens *before fertilization!*, despite the emphasis on sexual reproduction being the source of this randomness – no further mixing-up of genetic material occurs when the egg and sperm fuse to form the zygote (Figure 9-4). Specifically, each **gamete** (germ cell – egg or sperm) undergoes **meiosis**, and this is where *the future parent's own genetic material* that was inherited from *their* parents is shuffled. Thus, each parent is actually re-shuffling the genes from their own parents (the grandparents of the future child) to create some new random genetic sequences. This shuffling process occurs as the two copies of each **chromosome** (one from grandpa and one from grandma – the parents of the future parent) are split apart, such that each gamete only has *one* copy of each chromosome (i.e., it is “haploid”, compared to the normal “diploid” with 2 copies).

As you likely know, there are 23 distinct chromosomes, which are large collections of DNA. When the egg and sperm join to form the zygote, the separate collections 23 chromosomes from each parent are simply “added” back together to form a “full deck” of 46 chromosomes, without any further mixing. Thus, you are really the *sum* of random combinations of genes from each of your two sets of grandparents, and your two

parent's gene sets won't really mix until you have your own kids!

After fertilization, the zygote only divides via **mitosis**, which is the “normal” form of cell division that preserves the full deck of chromosomes in each of the two new *daughter* cells (each chromosome, one inherited from each parent, splits and replicates, again with no further recombination). Thus, your two parent's chromosomes are preserved intact from that point onward, and the primary form of interaction between them is in terms of the relative *dominance* vs. *recessiveness* of the genes inherited from each parent. If both resulting copies of a given gene across the two chromosomes are the same, then there is nothing further of interest to discuss – that gene will do whatever it does, in the same way, all the time. This is actually the default case: over 99% of our genes are identical across all people, and thus across your two parents.

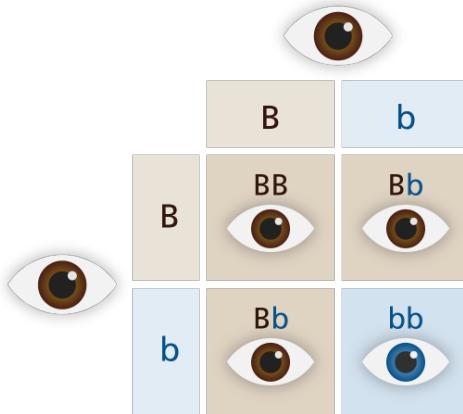


Fig 9-5: The logic of dominant vs. recessive genes, in the case of eye color, where blue is recessive compared to brown. Two people with brown eyes *can* give birth to a blue-eyed baby, if they are both *carriers* of the recessive blue gene.

However, for the roughly 0.6% of our genes that do differ across people, you may end up with a different version of that gene in each of your different chromosomes. These different versions are called **alleles**, and they are entire focus of interest in behavioral genetics and the study of heritability more generally. Some versions of a given gene are more likely to be transcribed and **expressed**, or to produce a functional protein product, and this is what is meant by **dominant** vs. **recessive** (and like most things, it is a continuum, not a dichotomy). Thus, only if you end up having both copies of a gene in the recessive (non-dominant) form, will that recessive version actually do its thing (or fail to do the thing that the dominant gene would otherwise do). Otherwise, having a recessive form of a gene typically doesn't make much of a difference in the overall function of the organism.

This presence of recessive genes is the reason we (still) have genetic disorders at any significant rate in the population. Any allele (genetic variant) that is dominant *and* produces bad effects, is quickly driven out of the population through natural selection – people with that gene version don't tend to survive and reproduce. However, a recessive gene can fly under the radar and persist in the population, because the odds of two people having the *same* recessive gene variant is really quite low on average (only 0.6% of genes vary at all, and most recessive alleles are relatively rare on top of that). Except, of course, if they are siblings or otherwise closely related, which is why incest is generally frowned upon.

Heritability and Individual Differences

Now we can actually talk about how behavioral genetics works. Basically, it amounts to comparing the genetic similarity of people against their *phenotypic* similarity, where the **phenotype** is just a complex word for the thing you are actually interested in, such as IQ, height, eye color, etc. For a very small set of phenotypes, your genes essentially determine 100% of how you'll end up. For example, there are specific recessive genes that cause Huntington's disease, and cystic fibrosis. But for almost everything else, the relationship between genetic differences and phenotypic differences is much more complex, and actually determining how much can be attributed to genes is surprisingly challenging.

Let's take the case of IQ. As we discussed in the chapter on intelligence, there is ample evidence that your IQ is a function of learning, motivation, and the wealth and general **socio-economic status (SES)** of your parents. Thus, any dependence on genes is likely to be at least somewhat indirect. But how can we measure it? The simplest way would be to compare the IQs of identical (monozygotic) versus fraternal (dizygotic) twins, and somehow use their known overall genetic similarity differences to compute how much of their measured IQ differences can be accounted for by those known genetic factors.

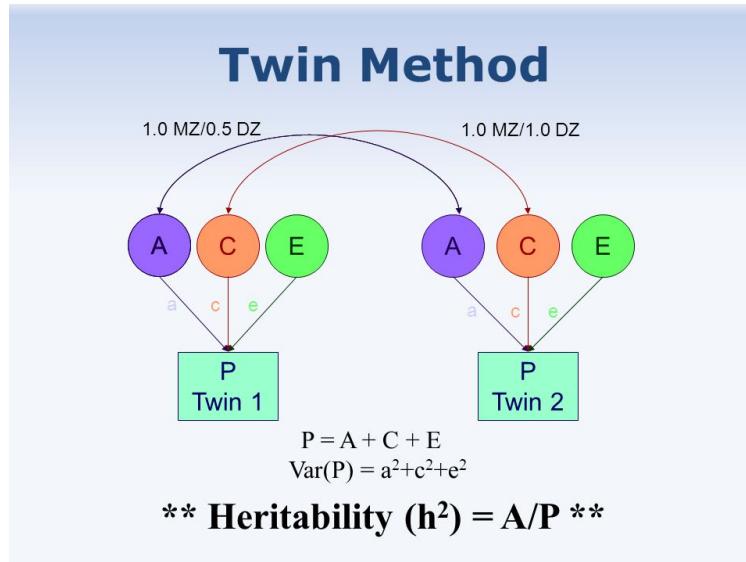


Fig 9-6: How heritability is computed from identical (MZ = monozygotic) and fraternal (DZ = dizygotic) twins. P = phenotype variance; A = additive genetic variance; C = common (shared) environment; E = unique (non-shared) environment (and everything else).

Critically, we need to also include some kind of factor that can account for non-genetic influences on IQ, which goes under the general category of **environmental** factors. To, make things more interesting, this latter category is typically split into **shared / common** and **non-shared / unique** environmental contributions, determined by whether the children were reared in the same family environment or not. Thus, there is a three-way tug-of-war dynamic between genetic factors (which are labeled with the letter *A* for additive genetic factors), and these two environmental factors (*C* and *E*), comprising the ACE model shown in Figure 9-6. The overall genetically-associated portion is called **heritability**, and is denoted with the letter *h*.

The comparative nature of heritability is a source of major interpretational problems. Just as we saw back in the neurons chapter, this is a fundamentally *Contrast-based* dynamic, and the *relative* balance between genetic and non-genetic factors can be affected by increases or decreases in our measurements of *either* of these factors. In particular, the apparent heritability of IQ could increase just by *decreasing* the strength of environmental factors. Indeed, there is considerable evidence for exactly this effect happening in **WEIRD (Western, Educated, Industrialized, Rich, Democratic)** societies, where the majority of the population has essentially comparable levels of health, nutrition, education, etc. In this case, the impact of environmental factors is greatly reduced compared to cases where some people are severely malnourished and have little access to education.

Figure 9-7 shows an idealized representation of the overall results of an attempt to *independently* estimate the amount of environmental versus genetic contributions to the phenotype of reading ability, as a function of the parent's education level (Kremen et al. 2005). As you can see, they found evidence that the amount of environmental variance went way down with increasing parental education, directly consistent with this broader WEIRD effect. Parents with more education generally are wealthier, and provide more enriched educational opportunities to their children – thus eliminating variability in these factors across individuals. What is left over is the raw genetic variability, which, due to the random mixing processes at work (Figure 9-4), remains relatively constant. Thus, the bottom line is that the raw magnitudes of heritability scores

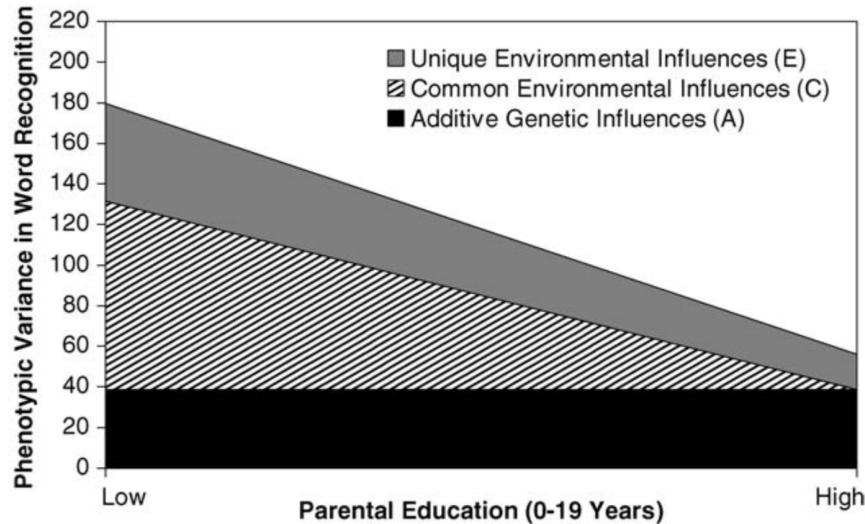


Fig 9-7: Idealized estimate of the effect of parental education on environmental (E and C) versus genetic contributions to the phenotype of reading (word recognition). Environmental influences go down with greater parental education, increasing the apparent amount of genetic heritability according to h measures, *without any actual change in genetic influence*. This is a consequence of the contrast-based, relative way that heritability is measured – it does not give an absolute value of genetic contribution (which is extremely difficult to measure – that is why this figure is just an idealization). From Kremen et al (2005).

cannot be taken as a direct measure of absolute genetic influence – just like absolute, perfect pitch is very difficult due to our contrast-based perceptual systems, it is extremely difficult to quantify the absolute level of genetic influence in any meaningful way, because it is always relative to the environment, which is even harder to measure directly than genetic differences are.

Most studies of heritability take place in WEIRD societies, and have produced estimates of heritability around 0.5 for almost every phenotype you can think of (Figure 9-8). In addition to those shown, personality factors of *neuroticism* and *openness* (which we'll learn more about in the next chapter) have a measured heritability of 0.4 to 0.6 (Power and Pluess 2015). Interestingly, Figure 9-8 also shows comparable results from a newer technique based on direct measurements of the genomes of a large number of *unrelated* people, known as **Genome-wide Complex Trait Analysis (GCTA)** (Trzaskowski, Dale, and Plomin 2013) (see Figure 9-9 for more detailed results on this technique, for IQ; (Sniekers et al. 2017)). Across many applications of this and related techniques, there is a consistent finding that heritability estimates are about half those based on twins! This is the latest version of the **missing heritability** problem that has come up repeatedly over many years of attempts to use direct genetic measurements to predict phenotypic variance across people (Turkheimer 2000, 2011; Plomin and Deary 2015).

There are at least two potential explanations for this missing heritability. One is that the current genome-based analyses are missing *rare genetic variants*, so their ability to capture the full scope of genetic differences among people is thus limited. However, recent analyses based on full genetic sequencing instead of the much sparser (and less expensive) sampling techniques used previously have estimated that these rare variants are likely to only contribute about 5% of this missing heritability (Evans et al. 2018). Thus, the remaining missing heritability points instead to a range of potential differences between twins and samples of unrelated people, which can inflate the heritability estimates generated by twin studies, which we'll explore in a moment, and point to interesting ways that genes can interact with the environment.

The bottom line at this point is that our current best guess as to the “true” level of genetic influence on various phenotypes is more like the right-hand side of Figure 9-8 based on direct genetic measurements, rather than the traditional twin-based estimates. Thus, genes probably account for a more plausible 25% of overall variance on average, instead of 50%. And this is still with all the inflation produced by the WEIRD reduction in environmental variance, so if we actually measured all of humanity, those heritability estimates

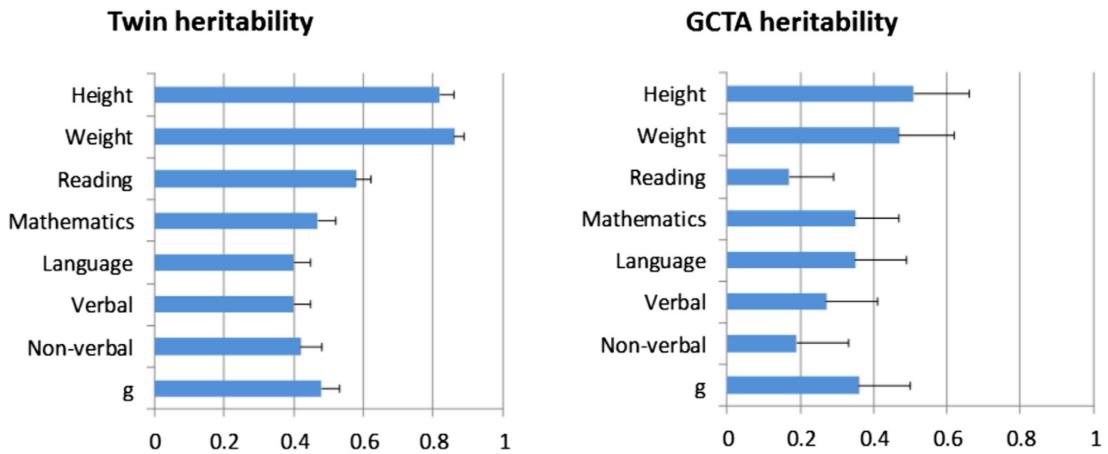


Fig 9-8: Heritability estimates of various phenotypes, computed either from twins or directly from the genome (GCTA = genome-wide complex trait analysis). There is a *missing heritability* in the direct genetic measures (which are about half as big as the twins), which could be due to rare genetic variation not measured, or by gene-by-environment interactions present in the twin sample (who are by definition related), versus the unrelated people used for GCTA. From Trzaskowski et al (2013).

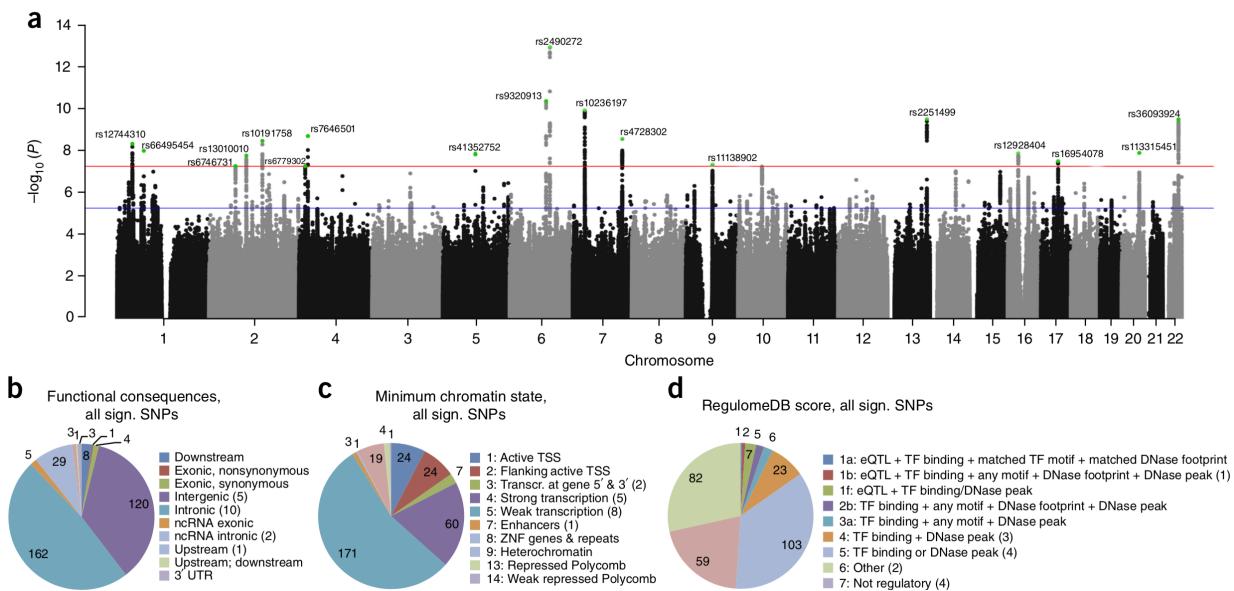


Fig 9-9: The latest Genome-Wide Association Study (GWAS) results for IQ (Sniekers et al, 2017). A total of 336 locations across the genome exceed the statistical threshold for being significantly correlated with IQ. No single gene accounts for much of the variance in IQ, and all together they only account for 4.8%. Most of the relevant genes are regulatory, and very few are clearly protein coding, consistent with the idea that the “junk” DNA is where all the action is.

would be much lower. Overall, this is consistent with the importance of learning for shaping the brain, as emphasized previously in the learning chapter.

Shared Environment and Parental Influences

Another striking finding from twin studies is that the estimates of the *C* factor in the ACE model, representing shared environmental influences due to children being reared in the same household, are almost always near zero. *Judith Rich Harris* has interpreted these pervasive findings to argue that *parents don't matter* in shaping how their children turn out, beyond of course contributing their genes to their children (Harris 2011). This striking conclusion flies in the face of most people's deeply-held beliefs about the importance of parents in shaping their kids, and Harris's book raises many fascinating points about why this idea might in fact be wrong. For example, children of immigrants typically become most proficient in the language of their new home, not the one spoken by their parents, and in general seem to be much more strongly influenced by their peers than by their parents. Indeed, probably most parents can recognize that their kids do seem to take them for granted, and are typically much more sensitive to what their friends think and do.

Another thing that parents with multiple kids are always struck by is how *different* their kids can be. That genetic crossover mixing stuff really works! Thus, Harris emphasizes that the environment for each such kid is very much an *interaction* between the parent and the child, with perhaps relatively little of a "main effect" of the parent overall across all the kids. In other words, the child shapes their own environment as much as the parent does. This is consistent with our focus on the importance of *Control* in the individual: just as we cannot convince our friends to change their beliefs, neither can a parent really control their child nearly as much as we often wish we could! The developmental transitions starting with the "terrible twos" mark the real onset of an independent, willful being, and from that point onward, the parent's influence is on a consistent downward slide. Of course, some kids turn out very much like their parents, but a roughly equal portion end up rebelling and try to be as different from their parents as possible. Thus, when looking for an overall consistent statistical effect of parents, perhaps you could see how it might be hard to find.

Another way to put this problem is that we have a highly accurate, reliable ways of estimating genetic differences among different pairs of people (e.g., identical vs. fraternal twins), but our ability to measure the similarity of *everything else* about people, including the actual nature of their individual experiences within their shared family environment, is significantly worse (Turkheimer 2000). Thus, when you pit a really solid estimate of genetic similarity against a really noisy estimate of shared environment, it is perhaps no surprise that the genetic effects are consequently over-estimated, while the shared environmental effects are under-estimated. This is in fact one of the major potential sources of the missing heritability present in twin studies relative to the direct genome-based studies (we should more accurately refer to this as the **excess heritability** in twin studies, actually). In other words, there is in fact a significant contribution of shared environment (parents really *do* matter!), but because we can't measure it very well, this contribution ends up getting soaked up by the much stronger genetic factor, thereby artificially inflating it (Turkheimer 2000).

In addition, it turns out that various other factors are also hard to disentangle in twin studies, and could also account for the excess heritability (Keller and Coventry 2005; Evans et al. 2018). For example, people who share various traits are more likely to marry and have children (known as *assortative mating*), which thus inflates the overall genetic similarity of even the fraternal twins, beyond what is assumed by the simple twin model. In addition, genes don't combine additively, as assumed by the model – instead they interact through the dominant vs. recessive dynamic discussed earlier (and other similar non-additive interactions across genes, known as *epistasis*), and this can significantly reduce the expected genetic similarity of DZ twins (because the MZ twins have the same genes, this factor affects them both in the same way).

Thus, in summary, despite the chemical precision of our knowledge about the molecular basis of life and the nature of DNA, once you get up to the level of the entire organism, there is a great deal of uncertainty in estimating the contributions that our genes make to the kinds of traits that we care about as psychologists (e.g., IQ and personality). Despite all these difficulties, there is no doubt that genes are making an important contribution, but it may have been overestimated by twin studies, and is likely to account for roughly a quarter of the overall differences across individuals (in WEIRD societies). Furthermore, the ability to predict *anything* at all about a *single individual* based on their genetic profile is *extremely limited* and generally nonexistent, outside of the few strong recessive genetic disorders where one or a few genes can make a huge

difference. There are simply way too many complex interactions both within the biology, and between a person and their environment, to make any kind of predictions at the individual level. Movies such as *Gattaca* which depict a dystopian future where everyone's future is predicted from their genes will almost certainly remain the stuff of science fiction.

Development

People are fascinated by butterflies because we *are* butterflies! The dramatic transformations that occur over the course of our development are truly astounding, and it is not too much of a stretch to say that we start out as something resembling a larva, and seemingly magically transform into a fully-functioning being with capacities unparalleled in the known universe. The first three months of life have been described as the “fourth trimester” – essentially we should still be in the womb but then we would be too big to ever get out, so we complete the last part of our gestation outside the womb.

After about four months, the pace of progress starts to increase, and after about six months, you can really start to see some real signs of neural activity going on under the hood! Babies at this point can actually recognize their parents, and start babbling in a way that is distinguishable from mere drooling and gurgling. They are obsessed with putting things in their mouths, and generally can reach for things with a non-zero chance of grabbing them. After all this time, they can finally start to sit up and support their own huge heads, and maybe start crawling. With another 6 months, babies can start saying “mamma” and “dada” and learn some sign language, and follow simple directions – language is really starting to happen. They can drink from actual cups (though most parents will stick with sippy cups for a while longer), and have recognizable hand-eye coordination for grabbing stuff and putting things where they want (though they will not be able to catch anything for years to come). Standing is starting, walking with support may be getting underway. Peek-a-boo may still be interesting, but you can see that bigger things are on the horizon.

By their second birthday, kiddos are recognizably human beings. They have basic mastery of their sensory and motor systems (though still have a long way to go for catching things, and potty training is a major issue at this age), and language learning has entered that exciting **naming explosion** period when 10's of new words can be learned in a single day. Most importantly, kids start using language to say the most important word: “no”! This is the onset of the **terrible twos**, when *cognitive control* really emerges, and a real sense of independent motivation takes hold. In other words, the basic elements of all of our special human capacities are at least minimally present, and from this point onward, it is just more, better, faster, smoother, etc (as a gross simplification).

Piaget and the Development of the Neural CPU

Jean Piaget, a founding figure in developmental psychology, ended his first stage of development, which he labeled the the **sensorimotor stage**, at this two-year mark (Figure 9-10). Everything beyond this first two years was focused on further stages of progress toward something he called **formal operations** (emerging fully after 11 years of age), which is the ability to perform abstract, logical reasoning, and engage in effective planning and strategizing. Furthermore, he emphasized the ability to transfer knowledge across different domains at this level. Basically, Piaget is describing the function of the **neural CPU**, supported by the **prefrontal cortex** and **basal ganglia**, as we discussed in the thinking and control chapter. These systems work together to enable information to be juggled in **working memory**, and behavior to be controlled by something approximating a **program** of steps to execute in sequence over time.

The fact that it takes roughly 11 years for this capacity to develop is consistent with the idea that the natural function of the neural networks of the brain is *not* intrinsically based on logic or symbolic processing like in a digital computer. Thus, per Piaget, it takes a *long* developmental progression for our limited abilities to perform logical, abstract reasoning to emerge (Inhelder and Piaget 1958). One of the main critiques of Piaget's work was that he only *described* the trajectory of steps toward this ultimate formal reasoning ability (as shown in Figure 9-10), without providing a clear, testable proposal as to *how exactly this reasoning ability develops* (Lourenço and Machado 1996). However, we *still* don't have a very good idea how this kind of thing emerges over the course of learning, experience, and brain development, so perhaps we shouldn't be too hard on Piaget. His proposal was at least based on a number of very specific behavioral tasks that characterize the level of competence at each of his stages, and those tasks have been widely tested.

Piaget's Theory

Stage	Age Range	Description
Sensorimotor	0-2 years	Coordination of senses with motor response, sensory curiosity about the world. Language used for demands and cataloguing. Object permanence developed
Preoperational	2-7 years	Symbolic thinking, use of proper syntax and grammar to express full concepts. Imagination and intuition are strong, but complex abstract thought still difficult. Conservation developed.
Concrete Operational	7-11 years	Concepts attached to concrete situations. Time, space, and quantity are understood and can be applied, but not as independent concepts
Formal Operations	11+	Theoretical, hypothetical, and counterfactual thinking. Abstract logic and reasoning. Strategy and planning become possible. Concepts learned in one context can be applied to another.

The Psychology Notes Headquarter - <http://www.PsychologyNotesHQ.com>

Fig 9-10: Piaget's four stages of cognitive development. After the first 2 years, the notion of "operational" is key: this is really about how capable our neural CPU is at supporting increasingly abstract, symbolic-like processing.

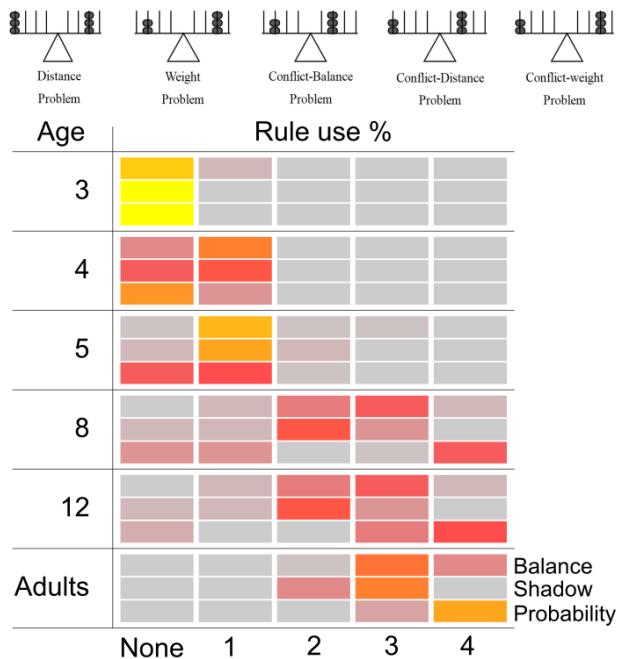


Fig 9-11: Siegler (1981)'s results on the extent to which people at different ages are able to integrate across multiple dimensions in three related tasks (brighter, more yellow = higher % of rule use across participants). The different rules (1-4) involve a progression of ability to integrate the different dimensions (e.g., weight vs. distance in the balance scale), as tested in problems like those shown at the top, where these dimensions are differentially relevant. While there is clearly a developmental progression, at each age there is considerable variability within and between tasks, and even adults are not perfectly "logical" (rule 4).

Piaget's **preoperational stage** (ages 2-7) marks the emergence of *symbolic*, but not logical thinking (a clearer term would thus perhaps be the *symbolic stage*). During this time, children build on their initial language competence to use words to refer to non-present objects, and develop more complex symbolic mental abilities – e.g., imagining entire scenarios and events unfolding over time. One of the signature developments during this time is the increasing ability to deal with *multiple factor relationships*. For example, on the **balance scale task**, there are two relevant factors or dimensions: distance and weight. Children start out only being able to focus on one of these at a time, and then gradually become better at integrating across these two dimensions (Figure 9-11). As *Robert Siegler* has emphasized (Siegler 1981), children's use of different levels of this integration is highly variable both between and within different tasks that are otherwise formally equivalent, as shown in Figure 9-11. This variability, which persists into adulthood, is inconsistent with the idea that people's behavior is strongly and consistently driven by different levels of logical reasoning and rules, which a simple reading of Piaget's theory would suggest (Lourenço and Machado 1996). On the other hand, it was possible to characterize any given person's reasoning on a given task, at a given point in time, according to specific, enumerable rules. Thus, unlike the 3-year-olds, people *do* appear to be using some kind of explicit rule-like reasoning process. It is just not very systematic.

Another widely-studied example is the **conservation task**, where two identical glasses initially have the same amount of water, and the child can recognize that. Then, a taller, thinner glass is introduced, and the child watches as water is poured entirely from one of the two identical glasses into the tall, thin one. Logically, the tall glass must have the same amount of water, but a younger child will typically say that the tall one has more, because it looks like it does. [YouTube video of conservation tasks](#). Siegler also demonstrated a similar level of variability in the developmental trajectory of more logical reasoning in these tasks, that take into account the shape and height dimensions (Siegler 1981). Even in adults, bartenders can still get away with this trick (as long as we don't see them actually do the pouring)!

In Piaget's **concrete operational stage** (ages 7-11), kids can now master some feats of logical reasoning, but only in concrete, specific situations. Thus, some can solve the conservation tasks, but not because they have the general abstract principle of conservation of mass (or energy), but rather because they've had enough physical experience and common sense to recognize that the water isn't going anywhere else. This is much like the Wason card selection task you tried to solve in the thinking chapter – you could nail the task when it was posed in a concrete, familiar situation (carding underage drinkers), but likely failed when it was novel and abstract in the card version. Thus, our brain's preference for concrete, familiar situations does not magically disappear after age 11 – it is always there, and only our ability to overcome it gets progressively better (but clearly never achieves anything like perfection).

The Development of the Prefrontal Cortex

Consistent with Piaget's focus on this protracted development of our logical, symbolic reasoning abilities that depend on our neural CPU, there is considerable evidence that the prefrontal cortex is one of the last brain areas to fully mature (Gogtay et al. 2004). This evidence comes from measurements of cortical thickness – your brain becomes progressively thinner as synapses are pruned over the course of development (Figure 9-12). This pruning is associated with patterns of neural connectivity stabilizing and thus extra unnecessary connections can be eliminated, making neural information processing (i.e., detection, compression) more efficient. But also less flexible – this is one of the reasons it is difficult to “teach an old dog new tricks.”

There are many other indications of the protracted nature of prefrontal development. In addition to being important for supporting abstract, logical, computer-like processing, the prefrontal cortex is important for enabling *controlled processing* to overcome stronger, habitual automatic processing pathways. As we saw in the thinking chapter, there are a number of problem-solving tasks and puzzles that depend on having to overcome the “obvious” but wrong solution. Piaget actually developed a version of such a task that kids initially struggle with at around a year of age, and then master in the subsequent few months, known as the **A-not-B task**. In this task, a toy is repeatedly hidden in one location (*A*), and the infant is allowed to reach for it there. Then, it is hidden in a new location (*B*), but the child tends to reach back to the original “habitual” location, hence the name: A, not B. [YouTube video of A-not-B](#). Success in this task is thought to depend on the increasing ability to maintain neural activity representing the actual hiding location, dependent in part on the developing prefrontal cortex (Munakata 1998).

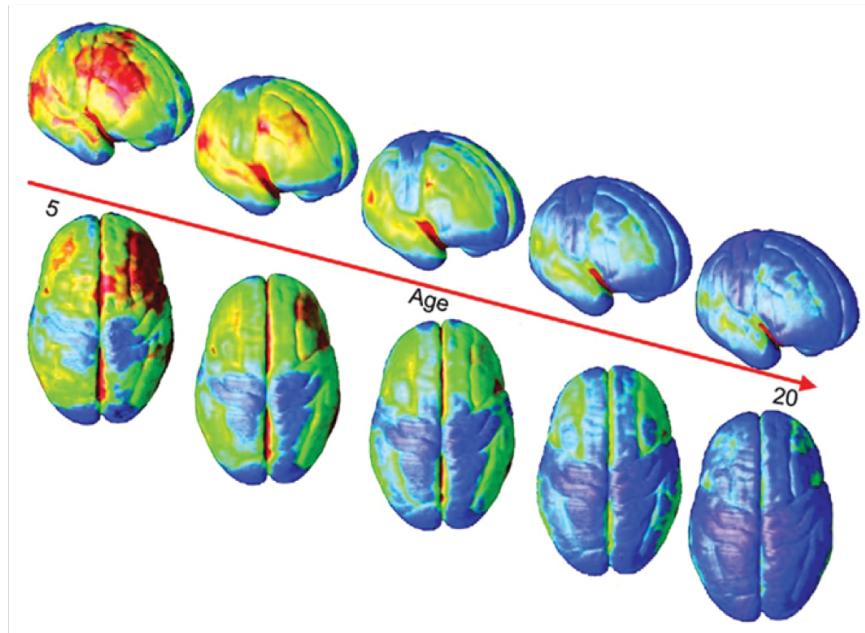


Fig 9-12: Synaptic pruning over the course of development across different brain areas. Blue colors indicate thinner brain areas where synaptic pruning has already taken place, while yellow and red areas are thicker and have yet to be pruned. The prefrontal cortex is one of the last areas to complete its pruning, along with the superior temporal lobe. From Gogtay et al., (2004)

Well after children succeed at this A-not-B task, they still fail at a related task involving sorting cards according to different dimensions (e.g., color vs. shape) (Zelazo, Frye, and Rapus 1996). [YouTube video of DCCS](#). Again, the explanation is that prefrontal cortex is continuing to develop throughout this long developmental period. But why does this manifest in different tasks at different times? One important factor is the relative familiarity and concreteness of the situation. In the A-not-B case, kids get a lot of experience tracking and reaching for objects, and thus it becomes “second nature” to track the hiding location of the toy. For the card sorting task, kids of this age have much less experience applying relatively arbitrary rules to sorting cards, even though the dimensions of color and shape are presumably quite familiar. Thus, these weaker mental states can be more strongly influenced by the new experience, and people get “stuck” in a mental set associated with the first sorting rule (Yerys and Munakata 2006). Likewise, when adults confront the Wason card selection task, we still don’t have much experience applying those kinds of arbitrary if / then rules in that way, and fall back on the raw perceptual match of the cards and the terms in the rule.

Another interesting potential indication of extended prefrontal development is that the genetic heritability of IQ actually appears to *increase* over time (Figure 9-13) (Haworth et al. 2009). This is not consistent with the idea that the genes are directly shaping the raw computational power of the neural CPU, but rather that they influence factors that shape the nature of learning over the relatively long developmental timescale. As we discussed earlier, it is likely that motivation plays an especially large role in shaping overall measured IQ, and this is consistent with the idea that the sustained influence of motivation over a relatively long time that produces increased genetic differences. Furthermore, other work has shown that people who have the highest measured IQ levels exhibit a longer period of greater environmental vs. genetic influence over IQ (Brant et al. 2013).

Crystallized vs. Fluid Intelligence Development

From everything that we just discussed, you might conclude that cognitive development is mostly about the improvement of the neural CPU *hardware*, in particular via development of the prefrontal cortex. We can think of this as improvements in the raw capacity for *fluid intelligence*. But in fact, it could just as well be more about the *software* that drives this neural CPU – that is, the specific knowledge that is learned via

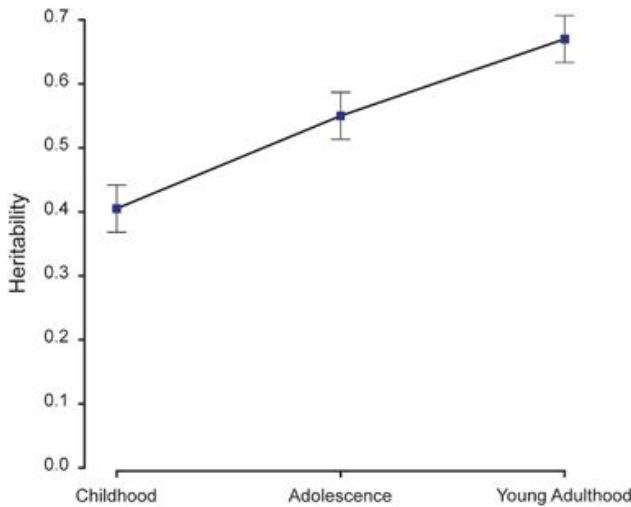


Fig 9-13: Heritability of IQ goes *up* as a function of age (based on twin studies; Haworth et al, 2009). One explanation is that IQ reflects in part the motivational and integrated learning contributions of the prefrontal cortex, and this takes a while to manifest over the protracted development of this brain area.

synaptic plasticity taking place through learning. This corresponds with crystallized intelligence, which as noted in the intelligence chapter is more generally thought to develop over time. Indeed, various evidence suggests that the raw capacity constraints of our working memory and other elements of the neural CPU are likely to be relatively fixed over time, and across different content domains.

For example, as we reviewed in the memory chapter, the same strong constraint of 4 items applies across a wide range of domains (Cowan 2001; Luck and Vogel 1997). Furthermore, the main way that memory capacity is increased is by forming new *chunks* that integrate previously separate items. These chunks are really just integrated knowledge representations shaped over learning, through synaptic plasticity mechanisms. This then suggests that perhaps the best way to improve the capacity and function of our neural CPU is similarly by developing more powerful knowledge structures that enable more sophisticated kinds of abstract and formal reasoning, using essentially the same raw hardware capacity. Likewise, we reviewed in the thinking chapter that brain training programs have been remarkably ineffective in improving people's general cognitive capacity, and mostly just improve their abilities to solve specific problems (Simons et al. 2016). Furthermore, Siegler's data on people's application of logical rules to various of Piaget's tasks shows a high level of variability and inconsistency, which again is consistent with the idea that specific experience-driven knowledge representations are playing a large role in shaping our reasoning abilities.

To summarize, it is likely that cognitive development is mostly about learning, learning, and more learning, and that through this process, synapses throughout the brain are shaped (and pruned), and this ultimately results in better chunks and abstractions that we can use to drive the limited capacity of our neural CPU. At every age, we *always* find it challenging when we are pushed outside of the zone of the concrete and familiar, and adults really don't seem to ever achieve anything like perfectly logical formal reasoning abilities. Nevertheless, we do have this neural CPU capability that apparently no other animal does, and with enough experience and learning, we can make it do some pretty amazing things.

Social, Personality and Moral Development

While Piaget focused mostly on cognitive development, other theorists focused on broader social, moral, and personality stages of development. For example, **Erik Erikson** articulated a set of stages that were partly inspired by Freudian psychoanalytic principles (Figure 9-14), and feature a dramatic, almost literary battle between opposing forces at each stage (Erikson 1956; Erikson and Erikson 1998). Like Freudian theory, and unlike Piaget, Erikson's ideas were based on observation and conjecture, not experimental data. Nevertheless, they certainly resonate with many of the themes that people wrestle with across the lifespan, and are indeed often the subject of great works of literature.

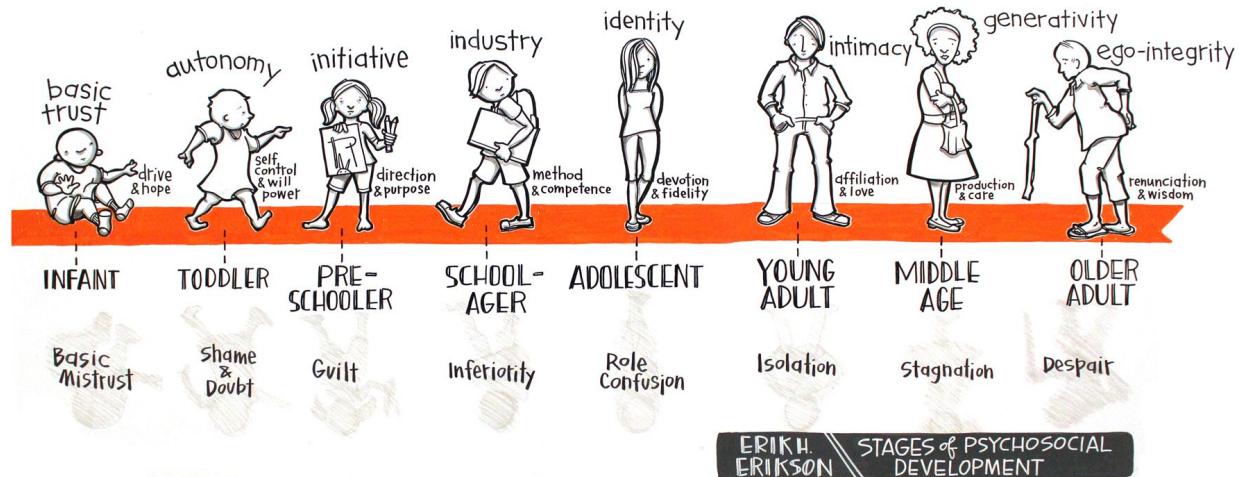


Fig 9-14: Erikson's stages of psychosocial development. Each stage represents a fundamental challenge between two conflicting forces, and if the positive outcome is achieved, the corresponding virtue is obtained. It was inspired by Freudian principles.

Up through adolescence the overriding issues center around personal control, and a broader social sense of belonging. The focus on individual control, autonomy, initiative, and industry are all very compatible with the idea that prefrontal cortex, which plays such a critical role in cognitive control, is among the brain areas developing throughout this same time period (Figure 9-12). Furthermore, the motivational and emotional areas in ventral and medial prefrontal cortex also seem to undergo extensive periods of development. Thus, in fact, Piaget and Erikson can be seen as describing two sides of the same developing coin! This coin is the third of our three C's, *Control*, and again it emerges as the central feature of human cognition, motivation, and social orientation.

Erikson's focus on the **identity crisis** in adolescence is one of the most impactful aspects of his framework. From an evolutionary psychology perspective, adolescents have to decide which tribe or group they will join – will they stay with their parental group, or break away and start off in a new direction? Will they try to become a leader or a lone wolf? Likewise, adolescent humans face dramatic choices about which direction their lives will take, and what role they will play in society. This is really the final step of the long process of asserting full independence and autonomy, that started back in the terrible twos (Erikson's 2nd stage).

In the sphere of personality development, **attachment theory** has been very influential, and it corresponds mostly with Erikson's first stage, about establishing a sense of trust or mistrust. An early inspiration for this theory was the work of **Harry Harlow**, who raised monkeys without their mothers, and found that they instead became very attached to cloth surrogate mothers, and ended up treating the cloth like children often treat their *attachment object* (e.g., a favorite stuffed animal or blanket) – they needed it around to feel comfortable, and became distressed when it was removed. **Mary Ainsworth** studied infant attachment to their mothers, using a similar **strange situation** where the mother left the child alone with a stranger.

Infants exhibited three characteristic patterns of behavior: **Secure attachment**, where they were comfortable exploring when their mother was around, and distressed when the mother was absent; **Insecure-avoidant attachment** where they had little interest in the mother or distress at her absence; and **Insecure-ambivalent attachment**, where the infant sought maternal attention and did not explore the environment as much, and exhibited high levels of distress when the mother left. Interestingly, in this latter case, the infants subsequently expressed ambivalence toward the mother when she returned, alternately seeking closeness and pushing her away. Although evidence has shown that these early attachment behaviors are correlated with various factors later in life, it remains unclear to what extent these are actually reflecting individual differences in personality, and genetically inherited personality factors from the mother (Harris 2009). We'll learn more about this in the next chapter.

Lawrence Kohlberg developed another stage-based developmental theory, in this case about *moral*

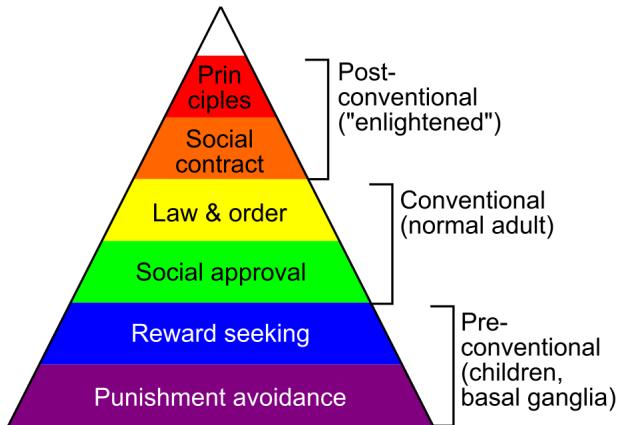


Fig 9-15: Kohlberg's stages of moral development.

development (Kohlberg and Hersh 1977) (Figure 9-15). This framework was originally based on some ideas of Piaget's, and it features three major levels with sub-levels within each. The first **preconventional** level characterizes the moral level of children: avoiding punishment and seeking rewards (self-interest). This level directly corresponds to the dopamine-based learning of the basal ganglia, and can thus be considered the biological foundation of all human (and animal) decision making. The **conventional** level is hypothesized to characterize the typical moral reasoning of a normal adult, based on both a seeking of social approval, and a respect for laws and a sense of duty to uphold them, to maintain social order. Thus, the conventional level could be considered the *social* level (i.e., social conventions), whereas the preconventional level is clearly the *individual* level. The highest level is **postconventional**, and involves moral reasoning based on a more abstract understanding of the necessity to preserve a social contract, and on abstract principles of basic human rights and ethics.

The correspondence with Piaget's focus on a progression toward more abstract, logical, rational thinking is evident in Kohlberg's stages. Consistent with the challenges that people have with actually thinking logically and rationally, a major critique of Kohlberg's framework is that it is very difficult to find any people who apply his highest level of moral reasoning in their daily lives.

One particularly interesting feature of Kohlberg's approach is that he employed a structured interview format based on moral dilemmas, in particular this famous *Heinz dilemma*:

A woman was on her deathbed. There was one drug that the doctors thought might save her. It was a form of radium that a druggist in the same town had recently discovered. The drug was expensive to make, but the druggist was charging ten times what the drug cost him to produce. He paid \$200 for the radium and charged \$2,000 for a small dose of the drug. The sick woman's husband, Heinz, went to everyone he knew to borrow the money, but he could only get together about \$1,000 which is half of what it cost. He told the druggist that his wife was dying and asked him to sell it cheaper or let him pay later. But the druggist said: "No, I discovered the drug and I'm going to make money from it." So Heinz got desperate and broke into the man's laboratory to steal the drug for his wife. Should Heinz have broken into the laboratory to steal the drug for his wife? Why or why not?

Critically, there is no right answer to this dilemma – Kohlberg instead was interested in what kind of reasoning people used in justifying their answers. Those reasons typically aligned with his different stages, e.g., "Heinz should not have stolen the drug because that is breaking the law", consistent with the law & order conventional stage.

Lifespan development

Erikson's stages were among the first to consider the full scope of the lifespan, beyond early development up to the start of adulthood. Many challenges and issues remain even after the turbulence of adolescence! There

are now many scientists actively researching the cognitive, social, and emotional aspects of aging.

Summary of Key Terms

This is a checklist of key terms / concepts that you should know about from this chapter.

- Evolution:
 - adaptation: survival and reproduction of the fittest
 - exaptation: selection of something for one function later turns out to be useful for something else (e.g., an oxygen boost for then living on land, or feathers for cooling then flight).
 - genetic algorithm: principle of evolution implemented on a computer
- Genetics:
 - base pairs (G, C, T, A)
 - amino acid = 3 base pair sequence
 - protein = sequence of amino acids
 - 1.5% of genome codes directly for proteins
 - genes are units of heredity – somewhat complicated to define
- Sexual reproduction:
 - zygote = fertilized egg cell
 - identical twins: monozygotic = same fertilized egg = 100% shared genes
 - fraternal twins: dizygotic = two different eggs = 50% shared genes on average
 - gamete: sperm and egg
 - meiosis: gametes go from 46 to 23 chromosomes, and crossover shuffling of parents genes
 - mitosis: normal cell division, preserves all 46 chromosomes
 - alleles: different versions of genes (most are the same in all people)
 - dominant vs. recessive and rare recessive diseases
- Heritability:
 - phenotype: overall property of organism (height, weight, IQ, etc..)
 - environment (shared, non-shared) vs. additive genetic contributions (ACE)
 - heritability as proportion of variance due to genes *relative* to environmental factors
 - reduced environmental variance in WEIRD societies
 - GCTA: directly based on genes of *random* people, not twins
 - missing heritability: GCTA and other direct techniques have about 1/2 the heritability
 - shared / common environment (C in ACE) is typically near 0: Judith Rich Harris: parents don't matter
- Development:
 - Piaget stages: sensorimotor, preoperational, concrete operational, formal operations
 - formal operations = functioning neural CPU based on prefrontal cortex, basal ganglia
 - balance scale, conservation tasks and multiple factor relationships: people are not logical; are variable across time and tasks
 - brain development = synaptic pruning, thinning; prefrontal cortex develops over long time
 - A-not-B, card sorting tasks and controlled processing
 - crystallized intelligence develops, fluid intelligence capacity ("hardware") remains relatively constant
 - Erikson's stages of psychosocial development: control, autonomy, and identity crisis in adolescence
 - attachment theory: Harry Harlow's motherless monkeys get attached to cloth
 - strange situation: secure attachment, insecure-avoidant, insecure-ambivalent
 - Kohlberg's stages of moral development: preconventional, conventional, postconventional

Chapter 10: Personality

Acknowledgments

Thanks to the current beta-testers for reading!

Glossary

About the Authors

Randall C. O'Reilly is Professor of Psychology and Neuroscience at the University of Colorado Boulder.

References

- Adam, Kirsten C. S., and Edward K. Vogel. 2016. "Reducing Failures of Working Memory with Performance Feedback." *Psychonomic Bulletin & Review* 23 (5): 1520–7. <https://doi.org/10.3758/s13423-016-1019-4>.
- Alarcon, Juan M., Angel Barco, and Eric R. Kandel. 2006. "Capture of the Late Phase of Long-Term Potentiation Within and Across the Apical and Basilar Dendritic Compartments of CA1 Pyramidal Neurons: Synaptic Tagging Is Compartment Restricted." *Journal of Neuroscience* 26 (1): 256–64. <https://doi.org/10.1523/JNEUROSCI.3196-05.2006>.
- Anagnostaras, S. G., S. Maren, and M. S. Fanselow. 1999. "Temporally Graded Retrograde Amnesia of Contextual Fear After Hippocampal Damage in Rats: Within-Subjects Examination." *The Journal of Neuroscience : The Official Journal of the Society for Neuroscience* 19 (February): 1106. <http://www.ncbi.nlm.nih.gov/pubmed/9920672>.
- Anderson, John R., Lynne M. Reder, and Herbert A. Simon. 1996. "Situated Learning and Education." *Educational Researcher* 25 (4): 5–11. <https://doi.org/10.3102/0013189X025004005>.
- Angelucci, A., F. Clascá, E. Bricolo, K. S. Cramer, and M. Sur. 1997. "Experimentally Induced Retinal Projections to the Ferret Auditory Thalamus: Development of Clustered Eye-Specific Patterns in a Novel Target." *The Journal of Neuroscience : The Official Journal of the Society for Neuroscience* 17 (April): 2040. <http://www.ncbi.nlm.nih.gov/pubmed/9045732>.
- Angelucci, Alessandra, and Paul C. Bressloff. 2006. "Contribution of Feedforward, Lateral and Feedback Connections to the Classical Receptive Field Center and Extra-Classical Receptive Field Surround of Primate V1 Neurons." In *Progress in Brain Research*, edited by J. -M. Alonso and P. U. Tse S. Martinez-Conde L. M. Martinez, 154, Part A:93–120. Visual PerceptionFundamentals of Vision: Low and Mid-Level Processes in Perception. Elsevier. [https://doi.org/10.1016/S0079-6123\(06\)54005-1](https://doi.org/10.1016/S0079-6123(06)54005-1).
- Aston-Jones, Gary, and Jonathan D. Cohen. 2005. "An Integrative Theory of Locus Coeruleus-Norepinephrine Function: Adaptive Gain and Optimal Performance." *Annual Review of Neuroscience* 28 (July): 403–50. <http://www.ncbi.nlm.nih.gov/pubmed/16022602>.
- Atkinson, R. C., and R. M. Shiffrin. 1968. "Human Memory: A Proposed System and Its Control Processes." In *The Psychology of Learning and Motivation: Advances in Research and Theory*, edited by K. W. Spence, 89–195. New York: Academic Press.
- Baars, B. J. 1988. *A Cognitive Theory of Consciousness*. New York: Cambridge University Press.
- Baddeley, A. D., and G. J. Hitch. 1974. "Working Memory." In *The Psychology of Learning and Motivation*, edited by G. Bower, VIII:47–89. New York: Academic Press.
- Baddeley, A., S. Gathercole, and C. Papagno. 1998. "The Phonological Loop as a Language Learning Device." *Psychological Review* 105 (March): 158. <http://www.ncbi.nlm.nih.gov/pubmed/9450375>.
- Balleine, B. W., and A. Dickinson. 1998. "Goal-Directed Instrumental Action: Contingency and Incentive Learning and Their Cortical Substrates." *Neuropharmacology* 37 (May): 407–19. <http://www.ncbi.nlm.nih.gov/pubmed/9704982>.
- Bartlett, F. C. 1932. *Remembering: A Study in Experimental and Social Psychology*. Cambridge: Cambridge University Press.
- Bjork, Robert A. 1994. "Memory and Metamemory Considerations in the Training of Human Beings." In *Metacognition: Knowing About Knowing*, 185–205. Cambridge, MA, US: The MIT Press.
- Bliss, T. V., and T. Lomo. 1973. "Long-Lasting Potentiation of Synaptic Transmission in the Dentate Area of the Anaesthetized Rabbit Following Stimulation of the Perforant Path." *The Journal of Physiology* 232 (October): 331–56. <http://www.ncbi.nlm.nih.gov/pubmed/4727084>.
- Boudry, Maarten, and Johan Braeckman. 2012. "How Convenient! The Epistemic Rationale of Self-Validating Belief Systems." *Philosophical Psychology* 25 (3): 341–64. <https://doi.org/10.1080/09515089.2011.579420>.
- Brant, Angela M., Yuko Munakata, Dorret I. Boomsma, John C. DeFries, Claire M. A. Haworth, Matthew C. Keller, Nicholas G. Martin, et al. 2013. "The Nature and Nurture of High IQ: An Extended Sensitive Period for Intellectual Development." *Psychological Science* 24 (8): 1487–95. <https://doi.org/10.1177/0956797612473119>.

- Bressloff, Paul C., Jack D. Cowan, Martin Golubitsky, Peter J. Thomas, and Matthew C. Wiener. 2002. "What Geometric Visual Hallucinations Tell Us About the Visual Cortex." *Neural Computation* 14 (February): 473–92. <http://www.ncbi.nlm.nih.gov/pubmed/11860679>.
- Brown, Tom B., Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, et al. 2020. "Language Models Are Few-Shot Learners." *arXiv:2005.14165 [Cs]*, July. <http://arxiv.org/abs/2005.14165>.
- Buschman, Timothy J., Markus Siegel, Jefferson E. Roy, and Earl K. Miller. 2011. "Neural Substrates of Cognitive Capacity Limitations." *Proceedings of the National Academy of Sciences* 108 (27): 11252–5. <http://www.ncbi.nlm.nih.gov/pubmed/21690375>.
- Buzsáki, G. 1989. "Two-Stage Model of Memory Trace Formation: A Role for 'Noisy' Brain States." *Neuroscience* 31 (3): 551–70. [https://doi.org/10.1016/0306-4522\(89\)90423-5](https://doi.org/10.1016/0306-4522(89)90423-5).
- Cameron, Judy, Katherine M. Banko, and W. David Pierce. 2001. "Pervasive Negative Effects of Rewards on Intrinsic Motivation: The Myth Continues." *The Behavior Analyst* 24 (1): 1–44. <https://doi.org/10.1007/BF03392017>.
- Carcea, Ioana, and Robert C Froemke. 2019. "Biological Mechanisms for Observational Learning." *Current Opinion in Neurobiology*, Neurobiology of Learning and Plasticity, 54 (February): 178–85. <https://doi.org/10.1016/j.conb.2018.11.008>.
- Cardinal, Rudolf N., John A. Parkinson, Jeremy Hall, and Barry J. Everitt. 2002. "Emotion and Motivation: The Role of the Amygdala, Ventral Striatum, and Prefrontal Cortex." *Neuroscience and Biobehavioral Reviews* 26 (May): 321–52. <http://www.ncbi.nlm.nih.gov/pubmed/12034134>.
- Carver, C S, and M F Scheier. 1990. "Origins and Functions of Positive and Negative Affect: A Control-Process View." *Psychological Review* 97 (December): 19–35.
- Carver, C S, and T White. 1994. "Behavioral Inhibition, Behavioral Activation, and Affective Responses to Impending Reward and Punishment: The BIS/BAS Scales." *Journal of Personality and Social Psychology* 67 (December): 319–33.
- Ceci, Stephen J., Mary Lyndia Crotteau Huffman, Elliott Smith, and Elizabeth F. Loftus. 1994. "Repeatedly Thinking About a Non-Event: Source Misattributions Among Preschoolers." *Consciousness and Cognition* 3 (3): 388–407. <https://doi.org/10.1006/ccog.1994.1022>.
- Chalmers, David. 1995. "Facing up to the Problem of Consciousness." *Journal of Consciousness Studies* 3(1) (December): 200–217.
- Chauveau, Frédéric, Damien Claverie, Emma Lardant, Christophe Varin, Eléonore Hardy, Augustin Walter, Frédéric Canini, Nathalie Rouach, and Armelle Rancillac. 2020. "Neuropeptide S Promotes Wakefulness Through the Inhibition of Sleep-Promoting Ventrolateral Preoptic Nucleus Neurons." *Sleep* 43 (1). <https://doi.org/10.1093/sleep/zsz189>.
- Cohen, J. D., and D. Servan-Schreiber. 1992. "Context, Cortex, and Dopamine: A Connectionist Approach to Behavior and Biology in Schizophrenia." *Psychological Review* 99 (April): 45–77. <http://www.ncbi.nlm.nih.gov/pubmed/1546118>.
- Cohen, Jonathan D., Richard D. Romero, Martha J. Farah, and D. Servan-Schreiber. 1994. "Mechanisms of Spatial Attention: The Relation of Macrostructure to Microstructure in Parietal Neglect." *Journal of Cognitive Neuroscience* 6 (4): 377–87.
- Conway, A. R. A., N. Cowan, M. F. Bunting, D. J. Therriault, and S. R. B. Minkoff. 2002. "A Latent Variable Analysis of Working Memory Capacity, Short Term Memory Capacity, Processing Speed, and General Fluid Intelligence." *Intelligence* 30 (January): 163–83.
- Cowan, N. 2001. "The Magical Number 4 in Short-Term Memory: A Reconsideration of Mental Storage Capacity." *Behavioral and Brain Sciences* 24 (August): 87–185. <http://www.ncbi.nlm.nih.gov/pubmed/11515286>.
- Craik, F. I. M., and R. S. Lockhart. 1972. "Levels of Processing: A Framework for Memory Research." *Journal of Verbal Learning and Verbal Behavior* 11 (January): 671–84.
- Crick, F. 1989. "The Recent Excitement About Neural Networks." *Nature* 337 (February): 129–32. <http://www.ncbi.nlm.nih.gov/pubmed/2911347>.

- Deci, E. L., R. Koestner, and R. M. Ryan. 2000. "A Meta-Analytic Review of Experiments Examining the Effects of Extrinsic Rewards on Intrinsic Motivation." *Psychological Bulletin* 125 (February): 627. <http://www.ncbi.nlm.nih.gov/pubmed/10589297>.
- Deese, James. 1959. "On the Prediction of Occurrence of Particular Verbal Intrusions in Immediate Recall." *Journal of Experimental Psychology* 58 (1): 17–22. <https://doi.org/10.1037/h0046671>.
- Dehaene, S., N. Molko, L. Cohen, and A. J. Wilson. 2004. "Arithmetic and the Brain." *Current Opinion in Neurobiology* 14 (2): 218–24.
- Dehaene, S., and L. Naccache. 2001. "Towards a Cognitive Neuroscience of Consciousness: Basic Evidence and a Workspace Framework." *Cognition* 79 (1-2): 1–37. <http://www.ncbi.nlm.nih.gov/pubmed/11164022>.
- Diamond, A. 1990. "The Development and Neural Bases of Memory Functions as Indexed by the A-Not-B Task: Evidence for Dependence on Dorsolateral Prefrontal Cortex." In *The Development and Neural Bases of Higher Cognitive Functions*, edited by A. Diamond, 267–317. New York: New York Academy of Science Press.
- Diekelmann, Susanne, and Jan Born. 2010. "The Memory Function of Sleep." *Nature Reviews Neuroscience* 11 (2): 114–26. <http://www.ncbi.nlm.nih.gov/pubmed/20046194>.
- Dimidjian, Sona, Christopher R. Martell, Ruth Herman-Dunn, and Samuel Hubley. 2014. "Behavioral Activation for Depression." In *Clinical Handbook of Psychological Disorders: A Step-by-Step Treatment Manual*, 5th Ed, 353–93. New York, NY, US: The Guilford Press.
- Dweck, Carol S. 2008. *Mindset: The New Psychology of Success*. Ballantine Books.
- Ebbinghaus (1885), Hermann. 2013. "Memory: A Contribution to Experimental Psychology." *Annals of Neurosciences* 20 (4): 155–56. <https://doi.org/10.5214/ans.0972.7531.200408>.
- Ekman, P., and W. V. Friesen. 1976. "Measuring Facial Movement." *Environmental Psychology and Nonverbal Behavior* 1 (1): 56–75.
- Engle, Randall W. 2002. "Working Memory Capacity as Executive Attention." *Current Directions in Psychological Science* 11 (1): 19–23. <https://doi.org/10.1111/1467-8721.00160>.
- . 2018. "Working Memory and Executive Attention: A Revisit." *Perspectives on Psychological Science* 13 (2): 190–93. <https://doi.org/10.1177/1745691617720478>.
- Ericsson, K A, W G Chase, and S Faloon. 1980. "Acquisition of a Memory Skill." *Science* 208 (June). <http://www.ncbi.nlm.nih.gov/pubmed/7375930>.
- Ericsson, K. A., and A. C. Lehmann. 1996. "Expert and Exceptional Performance: Evidence of Maximal Adaptation to Task Constraints." *Annual Review of Psychology* 47 (1): 273–305. <https://doi.org/10.1146/annurev.psych.47.1.273>.
- Erikson, Erik H., and Joan M. Erikson. 1998. *The Life Cycle Completed (Extended Version)*. W. W. Norton & Company.
- Erikson, Erik Homburger. 1956. "The Problem of Ego Identity." *Journal of the American Psychoanalytic Association* 4 (1): 56–121. <https://doi.org/10.1177/000306515600400104>.
- Eslinger, Paul J., Claire V. Flaherty-Craig, and Arthur L. Benton. 2004. "Developmental Outcomes After Early Prefrontal Cortex Damage." *Brain and Cognition*, Development of Orbitofrontal Function, 55 (1): 84–103. [https://doi.org/10.1016/S0278-2626\(03\)00281-1](https://doi.org/10.1016/S0278-2626(03)00281-1).
- Evans, Luke M., Rasool Tahmasbi, Scott I. Vrieze, Gonçalo R. Abecasis, Sayantan Das, Steven Gazal, Douglas W. Bjelland, et al. 2018. "Comparison of Methods That Use Whole Genome Data to Estimate the Heritability and Genetic Architecture of Complex Traits." *Nature Genetics* 50 (5): 737. <https://doi.org/10.1038/s41588-018-0108-x>.
- Ferrari, Pier F., Elisabetta Visalberghi, Annika Paukner, Leonardo Fogassi, Angela Ruggiero, and Stephen J. Suomi. 2006. "Neonatal Imitation in Rhesus Macaques." *PLOS Biology* 4 (9): e302. <https://doi.org/10.1371/journal.pbio.0040302>.
- Fox, Michael D., Abraham Z. Snyder, Justin L. Vincent, Maurizio Corbetta, David C. Van Essen, and Marcus E. Raichle. 2005. "The Human Brain Is Intrinsically Organized into Dynamic, Anticorrelated Functional Networks." *Proceedings of the National Academy of Sciences of the United States of America* 102 (27): 9673–8. <https://doi.org/10.1073/pnas.0504136102>.

- Frank, M. J. 2005. "When and When Not to Use Your Subthalamic Nucleus: Lessons from a Computational Model of the Basal Ganglia." In *Modelling Natural Action Selection: Proceedings of an International Workshop*, edited by A. K. Seth, T. J. Prescott, and J. J. Bryson, 53–60. Sussex: AISB.
- Frey, U., and R. G. M. Morris. 1998. "Weak Before Strong: Dissociating Synaptic Tagging and Plasticity-Factor Accounts of Late-LTP." *Neuropharmacology* 37 (May): 545–52. <http://www.ncbi.nlm.nih.gov/pubmed/9704995>.
- Fuster, J. M., and G. E. Alexander. 1971. "Neuron Activity Related to Short-Term Memory." *Science* 173 (January): 652–54.
- Gallese, Vittorio, Christian Keysers, and Giacomo Rizzolatti. 2004. "A Unifying View of the Basis of Social Cognition." *Trends in Cognitive Sciences* 8 (9): 396–403. <http://www.ncbi.nlm.nih.gov/pubmed/15350240>.
- Gerfen, Charles R., and D. James Surmeier. 2011. "Modulation of Striatal Projection Systems by Dopamine." *Annual Review of Neuroscience* 34: 441–66. <http://www.ncbi.nlm.nih.gov/pubmed/21469956>.
- Gershman, Samuel J., David M. Blei, and Yael Niv. 2010. "Context, Learning, and Extinction." *Psychological Review* 117 (1): 197–209. <http://www.ncbi.nlm.nih.gov/pubmed/20063968>.
- Gigerenzer, Gerd. 2006. "Out of the Frying Pan into the Fire: Behavioral Reactions to Terrorist Attacks." *Risk Analysis* 26 (2): 347–51. <https://doi.org/10.1111/j.1539-6924.2006.00753.x>.
- Gobet, Fernand, and Herbert A. Simon. 1996. "The Roles of Recognition Processes and Look-Ahead Search in Time-Constrained Expert Problems Solving: Evidence from Grand-Master-Level Chess." *Psychological Science* 7 (January): 52.
- Godden, D. R., and A. D. Baddeley. 1975. "Context-Dependent Memory in Two Natural Environments: On Land and Under Water." *British Journal of Psychology* 66 (January): 325–31.
- Gogtay, Nitin, Jay N. Giedd, Leslie Lusk, Kiralee M. Hayashi, Deanna Greenstein, A. Catherine Vaituzis, Tom F. Nugent, et al. 2004. "Dynamic Mapping of Human Cortical Development During Childhood Through Early Adulthood." *Proceedings of the National Academy of Sciences of the United States of America* 101 (21): 8174–9. <http://www.ncbi.nlm.nih.gov/pubmed/15148381>.
- Goldman-Rakic, P S. 1995. "Architecture of the Prefrontal Cortex and the Central Executive." *Annals of the New York Academy of Sciences* 769 (December): 71–84. <http://www.ncbi.nlm.nih.gov/pubmed/8595045>.
- Gollwitzer, P. M. 1993. "Goal Achievement: The Role of Intentions." *European Review of Social Psychology* 4: 141–85.
- Goodwin, Donald W., Barbara Powell, David Bremer, Haskel Hoine, and John Stern. 1969. "Alcohol and Recall: State-Dependent Effects in Man." *Science* 163 (3873): 1358–60. <https://doi.org/10.1126/science.163.3873.1358>.
- Greeno, James G., Joyce L. Moore, and David R. Smith. 1993. "Transfer of Situated Learning." In *Transfer on Trial: Intelligence, Cognition, and Instruction.*, 99–167. Westport, CT, US: Ablex Publishing.
- Harris, Judith Rich. 2009. "Attachment Theory Underestimates the Child." *Behavioral and Brain Sciences* 32 (1): 30–30. <https://doi.org/10.1017/S0140525X09000119>.
- . 2011. *The Nurture Assumption: Why Children Turn Out the Way They Do*. Simon and Schuster.
- Hasselmo, Michael E., and Chantal E. Stern. 2006. "Mechanisms Underlying Working Memory for Novel Information." *Trends in Cognitive Sciences* 10 (November). <http://www.ncbi.nlm.nih.gov/pubmed/17015030>.
- Haworth, C. M. A., M. J. Wright, M. Luciano, N. G. Martin, E. J. C. de Geus, C. E. M. van Beijsterveldt, M. Bartels, et al. 2009. "The Heritability of General Cognitive Ability Increases Linearly from Childhood to Young Adulthood." *Molecular Psychiatry*, June.
- Hayne, Harlene. 2004. "Infant Memory Development: Implications for Childhood Amnesia." *Developmental Review, The Nature and Consequences of Very Early Memory Development*, 24 (1): 33–73. <https://doi.org/10.1016/j.dr.2003.09.007>.
- Hazy, Thomas E., Michael J. Frank, and R. C. O'Reilly. 2010. "Neural Mechanisms of Acquired Phasic Dopamine Responses in Learning." *Neuroscience and Biobehavioral Reviews* 34 (5): 701–20. <http://www.ncbi.nlm.nih.gov/pubmed/19944716>.

- Hearne, Luke J., Jason B. Mattingley, and Luca Cocchi. 2016. "Functional Brain Networks Related to Individual Differences in Human Intelligence at Rest." *Scientific Reports* 6 (August): 32328. <https://doi.org/10.1038/srep32328>.
- Hebb, D. O. 1949. *The Organization of Behavior*. New York: Wiley.
- Hobson, J. Allan, and Edward F. Pace-Schott. 2002. "The Cognitive Neuroscience of Sleep: Neuronal Systems, Consciousness and Learning." *Nature Reviews Neuroscience* 3 (9): 679–93. <https://doi.org/10.1038/nrn915>.
- Hodgkin, A. L., and A. F. Huxley. 1952. "A Quantitative Description of Membrane Current and Its Application to Conduction and Excitation in Nerve." *The Journal of Physiology* 117 (4): 500–544. <https://doi.org/10.1113/jphysiol.1952.sp004764>.
- Hopwood, Christopher J., Aidan G. C. Wright, Emily B. Ansell, and Aaron L. Pincus. 2013. "The Interpersonal Core of Personality Pathology." *Journal of Personality Disorders* 27 (3): 270–95. <https://doi.org/10.1521/pedi.2013.27.3.270>.
- Howard, M. W., and M. J. Kahana. 1999. "Contextual Variability and Serial Position Effects in Free Recall." *Journal of Experimental Psychology. Learning, Memory, and Cognition* 25 (August): 923. <http://www.ncbi.nlm.nih.gov/pubmed/10439501>.
- Hull, C. L. 1943. *Principles of Behavior*. Appleton.
- Iacoboni, M., R. P. Woods, and G. Rizzolatti. 1999. "Cortical Mechanisms of Human Imitation." *Science* 286 (January): 2526.
- Inhelder, B., and J. Piaget. 1958. *The Growth of Logical Thinking from Childhood to Adolescence*. New York: Basic Books.
- Jacoby, L. L., J. P. Toth, and A. P. Yonelinas. 1993. "Separating Conscious and Unconscious Influences of Memory: Measuring Recollection." *Journal of Experimental Psychology: General* 122 (2): 139–54.
- Jilk, David, Christian Lebiere, R. C. O'Reilly, and John Anderson. 2008. "SAL: An Explicitly Pluralistic Cognitive Architecture." *Journal of Experimental & Theoretical Artificial Intelligence* 20 (3): 197–218. <http://www.ingentaconnect.com/content/tandf/teta/2008/00000020/00000003/art00004>.
- Jüngling, Kay, Thomas Seidenbecher, Ludmila Sosulina, Jörg Lesting, Susan Sangha, Stewart D. Clark, Naoe Okamura, et al. 2008. "Neuropeptide S-Mediated Control of Fear Expression and Extinction: Role of Intercalated GABAergic Neurons in the Amygdala." *Neuron* 59 (2): 298–310. <https://doi.org/10.1016/j.neuron.2008.07.002>.
- Kahan, D., H. Jenkins-Smith, and D. Braman. 2011. "Cultural Cognition of Scientific Consensus." *Journal of Risk Research* 14: 147–74.
- Keil, Frank C. 1981. "Constraints on Knowledge and Cognitive Development." *Psychological Review* 88 (January): 197–227.
- Keller, Matthew C., and William L. Coventry. 2005. "Quantifying and Addressing Parameter Indeterminacy in the Classical Twin Design." *Twin Research and Human Genetics* 8 (3): 201–13. <https://doi.org/10.1375/twin.8.3.201>.
- Klinger, E. 1975. "Consequences of Commitment to and Disengagement from Incentives." *Psychological Review* 82: 1–25.
- Klinzing, Jens G., Niels Niethard, and Jan Born. 2019. "Mechanisms of Systems Memory Consolidation During Sleep." *Nature Neuroscience* 22 (10): 1598–1610. <https://doi.org/10.1038/s41593-019-0467-3>.
- Koch, Christof, Marcello Massimini, Melanie Boly, and Giulio Tononi. 2016. "Neural Correlates of Consciousness: Progress and Problems." *Nature Reviews Neuroscience* 17 (5): 307–21. <https://doi.org/10.1038/nrn.2016.22>.
- Kohlberg, Lawrence, and Richard H. Hersh. 1977. "Moral Development: A Review of the Theory." *Theory into Practice* 16 (2): 53–59. <https://doi.org/10.1080/00405847709542675>.
- Kotovsky, K., J. R. Hayes, and H. A. Simon. 1985. "Why Are Some Problems Hard? Evidence from Tower of Hanoi." *Cognitive Psychology* 17 (January): 248–94.
- Kremen, William S., Kristen C. Jacobson, Hong Xian, Seth A. Eisen, Brian Waterman, Rosemary Toomey, Michael C. Neale, Ming T. Tsuang, and Michael J. Lyons. 2005. "Heritability of Word Recognition

- in Middle-Aged Men Varies as a Function of Parental Education.” *Behavior Genetics* 35 (4): 417–33. <https://doi.org/10.1007/s10519-004-3876-2>.
- Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. 2012. “ImageNet Classification with Deep Convolutional Neural Networks.” In *Advances in Neural Information Processing Systems 25*, edited by F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, 1097–1105. Curran Associates, Inc. <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>.
- Kubota, K., and H. Niki. 1971. “Prefrontal Cortical Unit Activity and Delayed Alternation Performance in Monkeys.” *Journal of Neurophysiology* 34 (3): 337–47. <http://www.ncbi.nlm.nih.gov/pubmed/4997822>.
- Kuhn, Thomas. 1962. “The Structure of Scientific Revolutions.” *International Encyclopedia of Unified Science* 2 (2).
- Lambon-Ralph, Matthew A., Elizabeth Jefferies, Karalyn Patterson, and Timothy T. Rogers. 2017. “The Neural and Computational Bases of Semantic Cognition.” *Nature Reviews Neuroscience* 18 (1): 42–55. <https://doi.org/10.1038/nrn.2016.150>.
- Lamme, Victor A. F. 2006. “Towards a True Neural Stance on Consciousness.” *Trends in Cognitive Sciences* 10 (11): 494–501. <https://doi.org/10.1016/j.tics.2006.09.001>.
- Lamon, Noemie, Christof Neumann, Thibaud Gruber, and Klaus Zuberbühler. 2017. “Kin-Based Cultural Transmission of Tool Use in Wild Chimpanzees.” *Science Advances* 3 (4): e1602750. <https://doi.org/10.1126/sciadv.1602750>.
- Lave, Jean, and Etienne Wenger. 1991. *Situated Learning: Legitimate Peripheral Participation*. Situated Learning: Legitimate Peripheral Participation. New York, NY, US: Cambridge University Press. <https://doi.org/10.1017/CBO9780511815355>.
- LeCun, Yann, Yoshua Bengio, and Geoffrey Hinton. 2015. “Deep Learning.” *Nature* 521 (7553): 436–44. <https://doi.org/10.1038/nature14539>.
- LeCun, Yann, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. 1990. “Handwritten Digit Recognition with a Back-Propagation Network.” In *Advances in Neural Information Processing Systems*, 396–404. Morgan Kaufmann.
- Lee, Tai Sing, Cindy F. Yang, Richard D. Romero, and David Mumford. 2002. “Neural Activity in Early Visual Cortex Reflects Behavioral Experience and Higher-Order Perceptual Saliency.” *Nature Neuroscience* 5 (6): 589–97. <http://www.ncbi.nlm.nih.gov/pubmed/12021764>.
- Libet, B., C. A. Gleason, E. W. Wright, and D. K. Pearl. 1983. “Time of Conscious Intention to Act in Relation to Onset of Cerebral Activity (Readiness-Potential). The Unconscious Initiation of a Freely Voluntary Act.” *Brain: A Journal of Neurology* 106 (Pt 3) (September): 623–42. <https://doi.org/10.1093/brain/106.3.623>.
- Liu, Yong, Anne G. Wheaton, Daniel P. Chapman, Timothy J. Cunningham, Hua Lu, and Janet B. Croft. 2016. “Prevalence of Healthy Sleep Duration Among Adults—United States, 2014.” *MMWR. Morbidity and Mortality Weekly Report* 65 (6): 137–41. <https://doi.org/10.15585/mmwr.mm6506a1>.
- Loftus, Elizabeth F., and John C. Palmer. 1974. “Reconstruction of Automobile Destruction: An Example of the Interaction Between Language and Memory.” *Journal of Verbal Learning and Verbal Behavior* 13 (5): 585–89. [https://doi.org/10.1016/S0022-5371\(74\)80011-3](https://doi.org/10.1016/S0022-5371(74)80011-3).
- Lourenço, Orlando, and Armando Machado. 1996. “In Defense of Piaget’s Theory: A Reply to 10 Common Criticisms.” *Psychological Review* 103 (1): 143–64. <https://doi.org/10.1037/0033-295X.103.1.143>.
- Luck, S. J., and E. K. Vogel. 1997. “The Capacity of Visual Working Memory for Features and Conjunctions.” *Nature* 390 (December): 279. <http://www.ncbi.nlm.nih.gov/pubmed/9384378>.
- MacLeod, C. M. 1991. “Half a Century of Research on the Stroop Effect: An Integrative Review.” *Psychological Bulletin* 109 (June): 163–203. <http://www.ncbi.nlm.nih.gov/pubmed/2034749>.
- Maier, Steven F., and Martin E. P. Seligman. 1976. “Learned Helplessness: Theory and Evidence.” *Journal of Experimental Psychology: General* 105: 3–46.
- Maier, Steven F., and Linda R. Watkins. 2010. “Role of the Medial Prefrontal Cortex in Coping and Resilience.” *Brain Research* 1355 (October): 52–60. <http://www.ncbi.nlm.nih.gov/pubmed/20727864>.
- Marr, D. 1969. “A Theory of Cerebellar Cortex.” *Journal of Physiology (London)* 202 (January): 437–70.

- . 1971. "Simple Memory: A Theory for Archicortex." *Philosophical Transactions of the Royal Society of London B: Biological Sciences* 262 (841): 23–81. <https://doi.org/10.1098/rstb.1971.0078>.
- Maslow, A. H. 1943. "A Theory of Human Motivation." *Psychological Review* 50: 370–96.
- McClelland, J. L., B. L. McNaughton, and R. C. O'Reilly. 1995. "Why There Are Complementary Learning Systems in the Hippocampus and Neocortex: Insights from the Successes and Failures of Connectionist Models of Learning and Memory." *Psychological Review* 102 (3): 419–57. <http://www.ncbi.nlm.nih.gov/pubmed/7624455>.
- McClelland, J. L., and D. E. Rumelhart. 1986. "A Distributed Model of Human Learning and Memory." In *Parallel Distributed Processing. Volume 2: Psychological and Biological Models*, edited by J. L. McClelland, D. E. Rumelhart, and PDP Research Group, 170–215. Cambridge, MA: MIT Press.
- Meltzoff, Andrew N., and M. K. Moore. 1994. "Imitation, Memory, and the Representation of Persons." *Infant Behavior and Development* 17 (January): 83–99.
- Miller, E. K., and R. Desimone. 1994. "Parallel Neuronal Mechanisms for Short-Term Memory." *Science (New York, N.Y.)* 263 (February): 520–22. <http://www.ncbi.nlm.nih.gov/pubmed/8290960>.
- Miller, G. A., E. Galanter, and K. H. Pribram. 1960. *Plans and the Structure of Behavior*. New York: Holt.
- Miller, George. 1956. *The Magical Number Seven, Plus or Minus Two: Some Limits on Our Capacity for Processing Information*. Vol. 101. Indiana : Bobbs-Merrill. <http://www.ncbi.nlm.nih.gov/pubmed/8022966>.
- Miyake, Akira, Michael J. Emerson, Francisca Padilla, and Jeung-chan Ahn. 2004. "Inner Speech as a Retrieval Aid for Task Goals: The Effects of Cue Type and Articulatory Suppression in the Random Task Cuing Paradigm." *Acta Psychologica* 115 (February): 123–42. <http://www.ncbi.nlm.nih.gov/pubmed/14962397>.
- Mnih, Volodymyr, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, et al. 2015. "Human-Level Control Through Deep Reinforcement Learning." *Nature* 518 (7540): 529–33. <http://www.ncbi.nlm.nih.gov/pubmed/25719670>.
- Mollick, Jessica A., Thomas E. Hazy, Kai A. Krueger, Ananta Nair, Prescott Mackie, Seth A. Herd, and R. C. O'Reilly. n.d. "A Systems-Neuroscience Model of Phasic Dopamine."
- Montague, P. Read, Peter Dayan, and Terrence J. Sejnowski. 1996. "A Framework for Mesencephalic Dopamine Systems Based on Predictive Hebbian Learning." *Journal of Neuroscience* 16 (5): 1936–47. <http://www.ncbi.nlm.nih.gov/pubmed/8774460>.
- Morris, Richard G. M. 2001. "Episodic Textendash Like Memory in Animals: Psychological Criteria, Neural Mechanisms and the Value of Episodic Textendash Like Tasks to Investigate Animal Models of Neurodegenerative Disease." *Philosophical Transactions of the Royal Society of London B: Biological Sciences* 356 (1413): 1453–65. <https://doi.org/10.1098/rstb.2001.0945>.
- Munakata, Y. 1998. "Infant Perseveration: Rethinking Data, Theory, and the Role of Modelling." *Developmental Science* 1 (January): 205–12.
- Nagel, Thomas. 1974. "What Is It Like to Be a Bat?" *The Philosophical Review* 83 (4): 435–50. <https://doi.org/10.2307/2183914>.
- Neisser, Ulric, Gwyneth Boodoo, Thomas Bouchard, Nathan Brody, Stephen Ceci, Diane Halpern, John Loehlin, Robert Perloff, Robert Sternberg, and Susana Urbina. 1996. "Intelligence: Knowns and Unknowns." *American Psychologist* 51 (2): 77–101.
- Newell, A., and H. A. Simon. 1972. *Human Problem Solving*. Englewood Cliffs, NJ: Prentice-Hall.
- Newman, Ehren L, and Kenneth A Norman. 2010. "Moderate Excitation Leads to Weakening of Perceptual Representations." *Cerebral Cortex* 20 (11): 2760–70. <http://www.ncbi.nlm.nih.gov/pubmed/20181622>.
- Norman, Kenneth A., and R. C. O'Reilly. 2003. "Modeling Hippocampal and Neocortical Contributions to Recognition Memory: A Complementary-Learning-Systems Approach." *Psychological Review* 110 (4): 611–46. <http://www.ncbi.nlm.nih.gov/pubmed/14599236>.
- Ongür, D., and J. L. Price. 2000. "The Organization of Networks Within the Orbital and Medial Prefrontal Cortex of Rats, Monkeys and Humans." *Cerebral Cortex* 10 (3): 206–19. <http://www.ncbi.nlm.nih.gov/pubmed/10731217>.

- O'Reilly, Randall C., Jacob L. Russin, Maryam Zolfaghari, and John Rohrlich. 2020. "Deep Predictive Learning in Neocortex and Pulvinar." *arXiv:2006.14800 [Q-Bio]*, June. <http://arxiv.org/abs/2006.14800>.
- O'Reilly, R. C. 1996. "Biologically Plausible Error-Driven Learning Using Local Activation Differences: The Generalized Recirculation Algorithm." *Neural Computation* 8 (5): 895–938. <https://doi.org/10.1162/neco.1996.8.5.895>.
- . 2006. "Biologically Based Computational Models of High-Level Cognition." *Science* 314 (5796): 91–94. <http://www.ncbi.nlm.nih.gov/pubmed/17023651>.
- O'Reilly, R. C., and Michael J. Frank. 2006. "Making Working Memory Work: A Computational Model of Learning in the Prefrontal Cortex and Basal Ganglia." *Neural Computation* 18 (2): 283–328. <http://www.ncbi.nlm.nih.gov/pubmed/16378516>.
- O'Reilly, R. C., and J. L. McClelland. 1994. "Hippocampal Conjunctive Encoding, Storage, and Recall: Avoiding a Tradeoff." *Hippocampus* 4 (6): 661–82.
- O'Reilly, R. C., Yuko Munakata, Michael J. Frank, Thomas E. Hazy, and Contributors. 2012. *Computational Cognitive Neuroscience*. Wiki Book, 1st Edition, URL: <http://ccnbook.colorado.edu>. <http://ccnbook.colorado.edu>.
- O'Reilly, R. C., Dean Wyatte, Seth Herd, Brian Mingus, and David J Jilk. 2013. "Recurrent Processing During Object Recognition." *Frontiers in Psychology* 4 (124). <http://www.ncbi.nlm.nih.gov/pubmed/23554596>.
- Plomin, R., and I. J. Deary. 2015. "Genetics and Intelligence Differences: Five Special Findings." *Molecular Psychiatry* 20 (1): 98–108. <https://doi.org/10.1038/mp.2014.105>.
- Posner, M. I. 1980. "Orienting of Attention." *Quarterly Journal of Experimental Psychology* 32 (1): 3–25.
- Power, R. A., and M. Pluess. 2015. "Heritability Estimates of the Big Five Personality Traits Based on Common Genetic Variants." *Translational Psychiatry* 5 (7): e604. <https://doi.org/10.1038/tp.2015.96>.
- Powers, William T. 1973. *Behavior: The Control of Perception*. Hawthorne.
- Quirk, Gregory J., and Devin Mueller. 2008. "Neural Mechanisms of Extinction Learning and Retrieval." *Neuropsychopharmacology* 33 (1): 56–72. <http://www.ncbi.nlm.nih.gov/pubmed/17882236>.
- Quiroga, R. Quian, L. Reddy, G. Kreiman, C. Koch, and I. Fried. 2005. "Invariant Visual Representation by Single Neurons in the Human Brain." *Nature* 435 (7045): 1102–7. <https://doi.org/10.1038/nature03687>.
- Raven, J. C., J. H. Court, and J. Raven. 1977. *Standard Progressive Matrices*. London: H. K. Lewis.
- Read, Stephen J., Brian M. Monroe, Aaron L. Brownstein, Yu Yang, Gurveen Chopra, and Lynn C. Miller. 2010. "A Neural Network Model of the Structure and Dynamics of Human Personality." *Psychological Review* 117 (1): 61–92. <http://www.ncbi.nlm.nih.gov/pubmed/20063964>.
- Rechtschaffen, A., M. A. Gilliland, B. M. Bergmann, and J. B. Winter. 1983. "Physiological Correlates of Prolonged Sleep Deprivation in Rats." *Science* 221 (4606): 182–84. <https://doi.org/10.1126/science.6857280>.
- Redish, A. D. 2004. "Neuroscience: Addiction as a Computational Process Gone Awry." *Science* 306 (5703): 1944–6.
- Rescorla, R. A., and A. R. Wagner. 1972. "A Theory of Pavlovian Conditioning: Variation in the Effectiveness of Reinforcement and Non-Reinforcement." In *Classical Conditioning II: Theory and Research*, edited by A. H. Black and W. F. Prokasy, 64–99. New York: Appleton-Century-Crofts.
- Riva-Posse, Patricio, Ki Sueng Choi, Paul E Holtzheimer, Cameron C McIntyre, Robert E Gross, Ashutosh Chaturvedi, Andrea L Crowell, Steven J Garlow, Justin K Rajendra, and Helen S Mayberg. 2014. "Defining Critical White Matter Pathways Mediating Successful Subcallosal Cingulate Deep Brain Stimulation for Treatment-Resistant Depression." *Biological Psychiatry*, April. <http://www.ncbi.nlm.nih.gov/pubmed/24832866>.
- Roediger, Henry L., and Andrew C. Butler. 2011. "The Critical Role of Retrieval Practice in Long-Term Retention." *Trends in Cognitive Sciences* 15 (1): 20–27. <https://doi.org/10.1016/j.tics.2010.09.003>.
- Roediger, Henry L., and Kathleen B. McDermott. 1995. "Creating False Memories: Remembering Words Not Presented in Lists." *Journal of Experimental Psychology: Learning, Memory, and Cognition* 21 (4): 803–14. <https://doi.org/10.1037/0278-7393.21.4.803>.

- Roumis, Demetris K., and Loren M Frank. 2015. "Hippocampal Sharp-Wave Ripples in Waking and Sleeping States." *Current Opinion in Neurobiology*, Circuit plasticity and memory, 35 (December): 6–12. <https://doi.org/10.1016/j.conb.2015.05.001>.
- Rudebeck, Peter H., Mark E. Walton, Angharad N. Smyth, David M. Bannerman, and Matthew F. S. Rushworth. 2006. "Separate Neural Pathways Process Different Decision Costs." *Nature Neuroscience* 9 (9): 1161–8. <http://www.ncbi.nlm.nih.gov/pubmed/16921368>.
- Rudy, Jerry. 2013. *The Neurobiology of Learning and Memory*. Second Edition. Oxford, New York: Oxford University Press.
- Rumelhart, D. E., G. E. Hinton, and R. J. Williams. 1986. "Learning Representations by Back-Propagating Errors." *Nature* 323 (9): 533–36.
- Rumelhart, D. E., and J. L. McClelland. 1986. "PDP Models and General Issues in Cognitive Science." In *Parallel Distributed Processing. Volume 1: Foundations*, edited by D. E. Rumelhart, J. L. McClelland, and PDP Research Group, 110–46. Cambridge, MA: MIT Press.
- Schneider, W., and R. M. Shiffrin. 1977. "Controlled and Automatic Human Information Processing: I. Detection, Search, and Attention." *Psychological Review* 84 (January): 1–66.
- Schultz, W. 1986. "Responses of Midbrain Dopamine Neurons to Behavioral Trigger Stimuli in the Monkey." *Journal of Neurophysiology* 56 (January): 1439–62.
- Schultz, W., P. Dayan, and P. R. Montague. 1997. "A Neural Substrate of Prediction and Reward." *Science* 275 (5306): 1593–9. <http://www.ncbi.nlm.nih.gov/pubmed/9054347>.
- Semendeferi, K., A. Lu, N. Schenker, and H. Damasio. 2002. "Humans and Great Apes Share a Large Frontal Cortex." *Nature Neuroscience* 5 (January): 272–76.
- Shiffrin, R. M., and W. Schneider. 1977. "Controlled and Automatic Human Information Processing: II. Perceptual Learning, Automatic Attending, and a General Theory." *Psychological Review* 84 (January): 127–90.
- Siegler, Robert S. 1981. "Developmental Sequences Within and Between Concepts." *Monographs of the Society for Research in Child Development* 46 (2): 1–84. <https://doi.org/10.2307/1165995>.
- Silver, David, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, et al. 2017. "Mastering the Game of Go Without Human Knowledge." *Nature* 550 (7676): 354–59. <https://doi.org/10.1038/nature24270>.
- Simons, Daniel J., Walter R. Boot, Neil Charness, Susan E. Gathercole, Christopher F. Chabris, David Z. Hambrick, and Elizabeth A. L. Stine-Morrow. 2016. "Do 'Brain-Training' Programs Work?" *Psychological Science in the Public Interest* 17 (3): 103–86. <https://doi.org/10.1177/1529100616661983>.
- Sloman, Steven, and Philip Fernbach. 2018. *The Knowledge Illusion: Why We Never Think Alone*. Penguin.
- Sniekers, Suzanne, Sven Stringer, Kyoko Watanabe, Philip R. Jansen, Jonathan R. I. Coleman, Eva Krapohl, Erdogan Taskesen, et al. 2017. "Genome-Wide Association Meta-Analysis of 78,308 Individuals Identifies New Loci and Genes Influencing Human Intelligence." *Nature Genetics* 49 (7): 1107–12. <https://doi.org/10.1038/ng.3869>.
- Sperling, George. 1960. "The Information Available in Brief Visual Presentations." *Psychological Monographs: General and Applied* 74 (11): 1–29. <https://doi.org/10.1037/h0093759>.
- Squire, L. R. 1992. "Memory and the Hippocampus: A Synthesis from Findings with Rats, Monkeys, and Humans." *Psychological Review* 99 (January): 195–231.
- Stocco, A., C. Lebiere, and J. R. Anderson. 2010. "Conditional Routing of Information to the Cortex: A Model of the Basal Ganglia's Role in Cognitive Coordination." *Psychological Review* 117: 541–74.
- Stroop, J. R. 1935. "Studies of Interference in Serial Verbal Reactions." *Journal of Experimental Psychology* 18 (January): 643–62.
- Stuss, D. T., D. Floden, M. P. Alexander, B. Levine, and D. Katz. 2001. "Stroop Performance in Focal Lesion Patients: Dissociation of Processes and Frontal Lobe Lesion Location." *Neuropsychologia* 39 (May): 771–86. <http://www.ncbi.nlm.nih.gov/pubmed/11369401>.
- Sutherland, Robert J, James O'Brien, and Hugo Lehmann. 2008. "Absence of Systems Consolidation of Fear Memories After Dorsal, Ventral, or Complete Hippocampal Damage." *Hippocampus* 18 (7): 710–18.

- <http://www.ncbi.nlm.nih.gov/pubmed/18446823>.
- Sutton, R. S., and A. G. Barto. 1981. "Toward a Modern Theory of Adaptive Networks: Expectation and Prediction." *Psychological Review* 88 (2): 135–70. <http://www.ncbi.nlm.nih.gov/pubmed/7291377>.
- Taylor, Andrea, Mevagh Sanson, Ryan Burnell, Kimberley A. Wade, and Maryanne Garry. 2020. "Disfluent Difficulties Are Not Desirable Difficulties: The (Lack of) Effect of Sans Forgetica on Memory." *Memory* 28 (7): 850–57. <https://doi.org/10.1080/09658211.2020.1758726>.
- Thorndike, E. L. 1911. *Animal Intelligence: Experimental Studies*. New York: The MacMillan Company.
- Tillmann, Sandra, Heidi E. Skibdal, Søren H. Christiansen, Casper R. Gøtzsche, Moustapha Hassan, Aleksander A. Mathé, Gregers Wegener, and David P. D. Woldebye. 2019. "Sustained Overexpression of Neuropeptide S in the Amygdala Reduces Anxiety-Like Behavior in Rats." *Behavioural Brain Research* 367 (July): 28–34. <https://doi.org/10.1016/j.bbr.2019.03.039>.
- Tolman, E. C. 1948. "Cognitive Maps in Rats and Men." *Psychological Review* 55 (4): 189–208. <http://www.ncbi.nlm.nih.gov/pubmed/18870876>.
- Tononi, Giulio. 2004. "An Information Integration Theory of Consciousness." *BMC Neuroscience* 5 (November): 42. <https://doi.org/10.1186/1471-2202-5-42>.
- Trzaskowski, Maciej, Philip S. Dale, and Robert Plomin. 2013. "No Genetic Influence for Childhood Behavior Problems from DNA Analysis." *Journal of the American Academy of Child & Adolescent Psychiatry* 52 (10): 1048–1056.e3. <https://doi.org/10.1016/j.jaac.2013.07.016>.
- Tulving, E. 1972. "Episodic and Semantic Memory." In *Organization of Memory*, edited by E. Tulving and W. Donaldson, 381–403. San Diego, CA: Academic Press.
- . 1983. *Elements of Episodic Memory*. Oxford, England: Clarendon Press.
- Turing, A. M. 1936. "On Computable Numbers, with an Application to the Entscheidungsproblem." *Proceedings of the London Mathematical Society* s2-42 (1): 230–65. <https://doi.org/10.1112/plms/s2-42.1.230>.
- Turkheimer, Eric. 2000. "Three Laws of Behavior Genetics and What They Mean." *Current Directions in Psychological Science* 9 (5): 160–64. <https://doi.org/10.1111/1467-8721.00084>.
- . 2011. "Still Missing." *Research in Human Development* 8 (3-4): 227–41. <https://doi.org/10.1080/15427609.2011.625321>.
- von Neumann, John. 1945. "First Draft of a Report on the EDVAC."
- Wallis, J D, R Dias, T W Robbins, and A C Roberts. 2001. "Dissociable Contributions of the Orbitofrontal and Lateral Prefrontal Cortex of the Marmoset to Performance on a Detour Reaching Task." *The European Journal of Neuroscience* 13 (May). <http://www.ncbi.nlm.nih.gov/pubmed/11359531>.
- Wallis, Jonathan D., and Steven W. Kennerley. 2011. "Contrasting Reward Signals in the Orbitofrontal Cortex and Anterior Cingulate Cortex." *Annals of the New York Academy of Sciences* 1239 (December): 33–42. <http://www.ncbi.nlm.nih.gov/pubmed/22145873>.
- Wamsley, Erin J. 2014. "Dreaming and Offline Memory Consolidation." *Current Neurology and Neuroscience Reports* 14 (3): 433. <https://doi.org/10.1007/s11910-013-0433-5>.
- Wason, P. 1968. "Reasoning About a Rule." *Quarterly Journal of Experimental Psychology* 20 (January): 273–81.
- Wilson, M. A., and B. L. McNaughton. 1994. "Reactivation of Hippocampal Ensemble Memories During Sleep." *Science (New York, N.Y.)* 265 (August): 676–78. <http://www.ncbi.nlm.nih.gov/pubmed/8036517>.
- Wong, William, Valdas Noreika, Levente Móró, Antti Revonsuo, Jennifer Windt, Katja Valli, and Naotsugu Tsuchiya. 2020. "The Dream Catcher Experiment: Blinded Analyses Failed to Detect Markers of Dreaming Consciousness in EEG Spectral Power." *Neuroscience of Consciousness* 2020 (1). <https://doi.org/10.1093/nc/niaa006>.
- Yerkes, R. M., and J. D. Dodson. 1908. "The Relation of Strength of Stimulus to Rapidity of Habit Formation." *Journal of Comparative Neurology and Psychology* 18 (January): 459–82.
- Yerys, Benjamin E., and Yuko Munakata. 2006. "When Labels Hurt but Novelty Helps: Children's Perseveration and Flexibility in a Card-Sorting Task." *Child Development* 77 (November): 1589–1607. <http://www.ncbi.nlm.nih.gov/pubmed/17107448>.

Yonelinas, Andrew P., Charan Ranganath, Arne D. Ekstrom, and Brian J. Wiltgen. 2019. "A Contextual Binding Theory of Episodic Memory: Systems Consolidation Reconsidered." *Nature Reviews Neuroscience* 20 (6): 364–75. <https://doi.org/10.1038/s41583-019-0150-4>.

Zelazo, Philip David, Douglas Frye, and Tanja Rapus. 1996. "An Age-Related Dissociation Between Knowing Rules and Using Them." *Cognitive Development* 11 (January): 37–63.