

From Virtual Reality to the Emerging Discipline of Perception Engineering

Steven M. LaValle, Evan G. Center,
Timo Ojala, Matti Pouke, Nicoletta Prencipe,
Basak Sakcak, Markku Suomalainen,
Kalle G. Timperi, and Vadim Weinstein

Xxxx. Xxx. Xxx. Xxx. 2023. AA:1–30

[https://doi.org/10.1146/\(\(please add article doi\)\)](https://doi.org/10.1146/((please add article doi)))

Copyright © 2023 by the author(s).
All rights reserved

Keywords

Virtual reality, robotics, control theory, dynamical systems, autonomous systems, game theory, psychophysics, perception, psychology, cognitive science, illusions, spoofing, social engineering

Abstract

This paper makes the case that a powerful new discipline, which we term *perception engineering*, is steadily emerging. It follows from a progression of ideas that involve creating illusions, from historical paintings and film, to video games and virtual reality in modern times. Rather than creating physical artifacts such as bridges, airplanes, or computers, *perception engineers* create illusory perceptual experiences. The scope is defined over any agent that interacts with the physical world, including both biological organisms (humans, animals) and engineered systems (robots, autonomous systems). The key idea is that an agent, called a *producer*, alters the environment with the intent to alter the perceptual experience of another agent, called a *receiver*. Most importantly, the paper introduces a precise mathematical formulation of this process, based on the von Neumann-Morgenstern notion of information, to help scope and define the discipline. It is then applied to the cases of engineered and biological agents with discussion of its implications on existing fields such as virtual reality, robotics, and even social media. Finally, open challenges and opportunities for involvement are identified.

The authors are with the Faculty of Information Technology and Electrical Engineering, University of Oulu, Finland.

1. INTRODUCTION

The field of virtual reality (VR) focuses mainly on a raft of technologies that when carefully assembled causes its users to have an illusory perceptual experience. When we refer to VR in this paper, we include augmented reality and other extended realities, which are becoming quickly unified as technologies advance. Most commonly, a VR user wears a head-mounted display (HMD) that provides visual and acoustic stimulation that is adjusted to her frame of reference using sensors that track head and other body movements. One goal is to obtain a sense of *presence*, which combines the place illusion (feeling as if in another location or world) and plausibility (feeling as if the virtual events were really taking place) (1). Although the term ‘virtual reality’ can be traced back to the philosopher Immanuel Kant (2), with its current usage proposed by Jaron Lanier in the 1980s, it remains elusive to precisely define it in a way that is not wholly dependent on the engineered devices of the times.

Regrettably, the industries surrounding VR have struggled for several decades through hype cycles of high expectations and investment followed by periods of disillusionment as the contemporary technologies are unable to deliver on promises. To bring about stable progress from a long-term research perspective, the challenge is to understand precisely what VR is and how it works on the entire system involved: A combination of engineered devices (displays, sensors) and a biological organism (human or otherwise). Having a rigorous scientific foundation could help improve design and analysis of VR systems so that steady, predictable progress can be made toward achieving a better quality of experience. To advance in this direction, we claim that VR, as it is viewed today, is merely one instance of a larger and steadily emerging discipline, which we term *perception engineering*. Over the coming years, we expect a rise in methods that create targeted perceptual illusions or experiences, much more broadly than the way HMD-centered VR is imagined today, and supported by fields such as machine learning, nanophotonics, and even social media technologies, in addition to the usual reliance of computer graphics, computer vision, sensors, computing systems, and displays. As will be explained shortly, perception engineering must also be based on principles from the biological sciences.

Why engineering? At its core, engineering is the intentional process of reshaping the environment to our advantage. The oldest examples are primitive tools and weapons which date back over 1,500,000 years, whereas later examples include the roads, bridges, aqueducts, mills, cars, airplanes, and electronics that made large-scale civilizations flourish. Early engineering emphasized practical applicability, often relying on trial and error to converge towards working solutions. However, modern engineering research has adopted the scientific method to not only construct a working solution, but to theorize about what could be a better solution to a particular problem or class of problems. This process usually involves developing mathematical models, reviewing the state-of-the-art, generating predictive hypotheses, collecting new data, iterating designs, and eventually sharing the findings in peer-reviewed articles. We want to leverage the benefits of this well-proven methodology, but what is the engineered artifact in our setting that would be analogous to, say, an airplane? We argue that it is a targeted *perceptual experience*, and VR devices such as HMDs are merely a component of a system that achieves it. We thus contemplate what it means to engineer a perceptual experience. Through this shift in understanding, we move the focus more toward the organism that has an experience, and away from particular devices that rapidly come and go. This means that in addition to leveraging the physical sciences

to engineer devices, we must also leverage the biological sciences, especially perceptual psychology, neuroscience, and physiology.

Why perception? People have been creating perceptual illusions for millennia as well, with the oldest known cave paintings dating back over 45,000 years. Through imagination, trial-and-error, and skillfully leveraging available technologies, artists continually develop more impressive works that seem to fool our senses and stimulate our imagination. For example, when the perspective method emerged in the 15th century, paintings depicted imagined worlds with consistent perceptions of depth and scale. Even more impressive is that exploiting the *stroboscopic apparent motion* effect (3), while being shown a rapid sequence of pictures, has led to over a century of cinema. In recent decades, simulators, video games, and VR have enabled active perceptual experiences, in which users interact with virtual environments, rather than passively consuming artistic content. Thus, we wonder what applying engineering principles to the design, analysis, and delivering of perceptual experiences will bring, as a natural complement to the works of artists.

The value of modeling To embark on this journey, an important step is to develop mathematical models that accurately capture what is understood so far, while giving insights into what may or may not be possible to achieve in the future. Mathematical models enable engineers to analyze or predict what *would happen* if a proposed design were constructed, so that the engineering process is accelerated. They also yield characterizations of what is expected to happen for systems that are actually built. Furthermore, especially for robots and autonomous systems, the models form the basis of algorithms or control laws that are implemented on devices. Perhaps most importantly, mathematical models form the foundations of technical disciplines that survive for generations, which distances them from the particular technological artifacts of the day. For example, a general mathematical formulation of configuration spaces is crucial in robotics for the development of general motion algorithms and nonlinear control laws. Control theory itself relies on the ability to unify a vast array of physical settings, whether mechanical, electrical, chemical, and so on, by mathematically characterizing them as a family of parameterized differential equations. Similarly, Turing machines provide a precise mathematical characterization of algorithms and a foundation upon which the discipline of computer science has been built. In this paper, we propose the first full mathematical model of perception engineering, including VR.

What should the model encompass? At the very least, we expect a mathematical formulation of perception engineering to model a person experiencing VR by wearing an HMD. We will generalize, however, to allow any organism, such as hamsters and fruit flies (4). As opposed to being a biological entity, we will even allow it to be fully engineered, as in the case of a robot. We use the term *agent* to refer to any such organism or robot, which generally has the ability to actuate in response to external stimulation that is sensed. Robot agents are encompassed because they can be fully modeled and analyzed, as opposed to biological organisms, which must be reverse engineered. Thus, whereas organisms start off as a *black box* (or perhaps a *gray box* thanks to biological sciences), engineered systems may serve as a fully explainable *white box*, leading to greater clarity and understanding.

Central to our formulation are two types of agents: 1) a *producer*, which creates and delivers a targeted perceptual experience by altering the environment of 2) a *receiver*. This

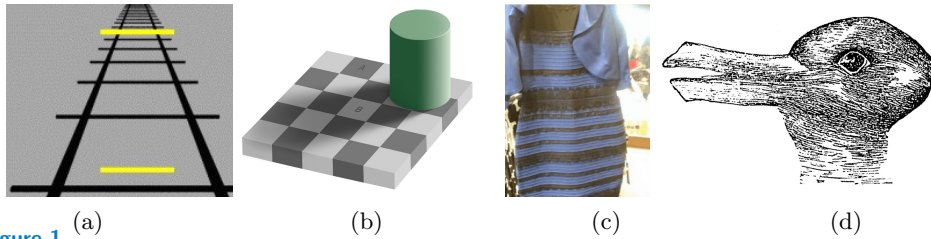


Figure 1

Several well-known illusions: (a) Ponzo (surprisingly, the yellow bars have equal length); (b) Checker Shadow (tiles A and B are surprisingly the same shade); (c) The Dress (some see it as black and blue, others as white and gold); (d) Rabbit-Duck (seems to flip between two different animals).

is the essence of perception engineering:

An agent alters the environment with the intent to alter the perceptual experience of another agent.

We also expect the targeted experience to be *plausible* (credible in some sense) and *illusory* (based on illusions), meaning that the receiver’s perceptions do not match ‘reality’. The notion of an illusion must be carefully defined, which is a challenging task considering how the term is used loosely. Consider, for example, the so-called illusions in Figure 1; Section 2.3 will rigorously define illusions, and then apply it in Section 5.2 to these examples.

In addition to VR, paintings, motion pictures, and video games, we also consider counterfeiting, magic tricks, wearing makeup, and just plain old lying to be examples that at least nominally fall under perception engineering. We avoid, however, direct alteration of the internals of the receiver, which might correspond, for example, to brain-machine interfaces, drugs, or viruses. We are more interested in a perceptual experience that results purely from interactions with the environment via sensing and actuation.

If there are multiple producers and receivers, then perception engineering extends naturally to a kind of social engineering. For example, a producer can lie to a receiver, and then the lie is ‘innocently’ propagated to other receivers. We can certainly imagine that fake news spreading on social media is a kind of virtual reality at a societal level. This paper focuses mainly on a single producer and receiver, but the powerful multi-agent extension to the social setting is also covered.

The rest of this paper Section 2 defines a general model of agents, capable of sensing and acting, in a shared environment. Each agent is modeled as a coupled dynamical system that has internal ‘brain’ states that interact with its external, surrounding environment states. Section 3 then develops the principles by which producer agents create illusory perceptual experiences for receiver agents. Important ideas include plausibility, robustness, forced fusion, and information-feedback policies. Section 4 applies this mathematical framework to the case of robots and other fully engineered systems. Section 5 then applies it to humans and other biological organisms, with substantial focus on VR. Section 6 concludes the paper by listing the challenges and opportunities ahead for building up the field of perception engineering.

2. AGENTS THAT SENSE AND ACT

2.1. Defining an Agent

An *agent* refers to an entity, either biological or engineered, that has a physical body and interacts with its environment through sensing and actuation. A clear boundary must be set between the agent's *internal* 'brain' and the *external* world, which corresponds to its body and any other physical attributes that are relevant in the surrounding environment. The *information space* (or *I-space*), I , is defined as the (nonempty) set of all internal states; each $i \in I$ is called an *information state* or *I-state*. The notion of information used here is inspired by von Neumann and Morgenstern for games with imperfect information (5), and extended to robotics (6, 7). This is not to be confused with Shannon's notion of information, which is independent, but can be used in conjunction with our formulation. The *external state space* (or *X-space*), X , is defined as the (nonempty) set of all possible physical states of the agent's body and environment. Each $x \in X$ is called an *external state* or *X-state*.

Let K be the *set of stages*, which correspond to discrete time instances, with the implication that stage $k+1$ comes after stage k for all $k \in K$. Every $k \in K$ corresponds to some *stage* k , X-state $x_k \in X$, and I-state $i_k \in I$. We either allow K to be infinite, in which case $K = \mathbb{N} = \{1; 2; 3; \dots\}$, or finite, in which case $K = \{1; 2; \dots; K+1\}$ for some final stage $K+1$.

To model actuation, let U be the set of *actions* available to the agent, through which it acts upon the external world via an *X-state transition function*, abbreviated as *XTF*, expressed as $f: X \times U \rightarrow X$. An action $u_k \in U$ applied at stage $k \in K$ from X-state $x_k \in X$ results in a transition to X-state $x_{k+1} = f(x_k; u_k)$ at stage $k+1$.

Sensing is modeled via a *sensor mapping*, $h: X \rightarrow Y$, in which Y is a set of possible sensor *observations*. At each stage k , an observation $y_k = h(x_k)$ is produced, using the X-state x_k at stage k . Extensions are possible, such as being action-dependent, time-dependent, or even based on a history of states as in the case of odometry (6, 8).

The *I-state transition function*, *ITF*, is $\tau: I \times U \rightarrow I$. To make the agent into an autonomous (fully determined) system, suppose that its action at each stage depends only on its I-state. This is expressed as a *policy* $\pi: I \rightarrow U$, which eliminates the U component from the domain of τ and results in $\tau: I \rightarrow I$ (because $u_k = \pi(i_k)$ is determined from the I-state i_k). Putting the definitions together results in a coupled dynamical system:

$$\begin{aligned} i_k &= \tau(i_{k-1}; y_k) && \text{in which } y_k = h(x_k); \text{ and} \\ x_{k+1} &= f(x_k; u_k) && \text{in which } u_k = \pi(i_k): \end{aligned} \quad 1.$$

Example 1 (Intbot)

Let $X = \mathbb{Z}$, the set of all integers. Let $U = \{1; 0; -1\}$, and the XTF is defined as $x_{k+1} = f(x_k; u_k) = x_k + u_k$. To create a perfect-information setup, let $Y = I = \mathbb{Z}$, $y_k = h(x_k) = x_k$ and $i_{k+1} = \tau(i_k; y_{k+1}) = y_{k+1} = x_{k+1}$. The policy $\pi: I \rightarrow U$ is then expressible as $u_k = \pi(i_k) = \text{sign}(x_k)$. A *stabilizing* policy is $\text{sign}(x_k)$ (yielding zero if $x_k = 0$), which brings the state to zero. From any initial X-state, x_1 , the system will arrive at $x_k = 0$ at stage $k = |x_1| + 1$, and remain there forever. It counts down to zero.

The far more common and interesting case is when the sensor mapping h is a many-to-one mapping. Generally, for a given observation $y_k \in Y$, the *preimage*

$$h^{-1}(y_k) = \{x_k \in X \mid y_k = h(x_k)\} \quad 2.$$

indicates the set of all possible X-states that could have caused it. Many-to-one ambiguity usually forces the X-state and I-state to have a non-trivial correspondence, to be explained in Section 2.2.

The models so far assume that the next state x_{k+1} is completely predictable from x_k and u_k , and the sensor observation y_k is completely predictable from x_k . We sometimes want to remove this limitation by defining a *disturbance-based* model. There will be two possibilities. The first is *nondeterministic disturbance*, in which case f and h are replaced by functions that produce a *set* of possible outcomes, rather than a single outcome. The *nondeterministic XTF* is $F: X \times U \rightarrow \text{pow}(X)$, in which pow denotes the power set, and F is the replacement of f . Thus, for a given x_k and u_k , $F(x_k; u_k) \subseteq X$ yields a nonempty set of possible x_{k+1} . The *nondeterministic sensor mapping* is $H: X \rightarrow \text{pow}(Y)$, and replaces h . Given x_k , $H(x_k) \subseteq Y$ yields a nonempty set of possible observations y_{k+1} . The other possibility is *probabilistic disturbance*, in which case f and h are replaced by functions that each produce a *probability density function (pdf)* (under appropriate measurability assumptions) on the space of possible outcomes. The *probabilistic XTF* is the pdf $\rho(x_{k+1} | x_k; u_k)$, and the *probabilistic sensor mapping* is the pdf $\rho(y_k | x_k)$. The general possibilities for l and π , and including extensions to nondeterministic and probabilistic disturbances are presented after the following example.

Example 2 (Linebot)

We extend Example 1 to $X = \mathbb{R}$, the real number line. Let $U = [0; 1]g$ and $x_{k+1} = f(x_k; u_k; \kappa) = x_k + u_k + \kappa$, which includes some disturbance parameter $\kappa \in [-1; 2]$. If κ is modeled nondeterministically, then $F(x_k; u_k) = [x_k + u_k - 1; x_k + u_k + 2]$ produces an interval of next possible states. If it is modeled probabilistically, then suppose a pdf $\rho(\kappa)$ is given over the interval $[-1; 2]$; $\rho(x_{k+1} | x_k; u_k)$ can then be defined as $\rho(x_{k+1} - x_k - u_k)$. As an example of adding disturbance to the sensing model, $y_k = h(x_k; \kappa) = x_k + \kappa$, which includes some disturbance parameter $\kappa \in [-1; 1]$. The nondeterministic sensor model would be $H(x_k) = [x_k - 1; x_k + 1]$. The probabilistic sensor model would involve a pdf $\rho(\kappa)$, and $\rho(y_k | x_k) = \rho(y_k - x_k)$.

2.2. Internal Information State Transitions

The coupled system (1) allows almost anything for the agent’s internal system, I-space I and ITF. What could they be? First consider two extreme possibilities. If h is defined as $y = h(x) = x$ with $Y = X$, then the X-state is perfectly sensed at every stage. We could then write $l = X$ and π simply mirrors f . This corresponds to an agent that conditions its actions on the precise X-state. Thus, its policy, called *state-feedback*, takes the form $\pi: X \rightarrow U$. At the other extreme, we could make $l = \emptyset g$, a singleton that is completely insensitive to X-state variations. Only a *constant* policy is possible: $\pi(0) = u$ for some particular action $u \in U$.

The interesting cases lie between these extremes and address the crucial question: To function appropriately, what should an agent retain as I-states? It is convenient to refer to *memory*, which merely means that the I-state depends on at least some information regarding actions and observations from prior stages. The singleton I-space is clearly memoryless, but a more interesting case is to make $l = Y$ and let $\rho(\kappa | y_k) = \rho(\kappa)$. This results in pure *sensor feedback*, with policies $\pi: Y \rightarrow U$. To add a small amount of memory, let $l = K$ and $\rho(\kappa | y_k) = \rho(\kappa)$, resulting in *stage feedback* policies $\pi: K \rightarrow U$ (alternatively known as

open-loop). The stage-feedback case at least uses the information of how many stages have passed.

Although it may seem that any ITF, and corresponding I-space, are possible, it turns out that one important condition, known as *sufficiency*, must be satisfied. It is briefly described here; see (6, 7) for more. Consider all information that might be available to an agent after k stages. Let κ be called the *history I-state*, defined as $\kappa = (\tilde{u}_{k-1}; \tilde{y}_k)$ in which $\tilde{u}_{k-1} = (u_1; u_2; \dots; u_{k-1})$ and $\tilde{y}_k = (y_1; y_2; \dots; y_k)$. Let $\kappa_1 = y_1$. The set of all history I-states is itself an I-space, denoted by I_{hist} . This corresponds to perfect, complete memory, with policies taking the form $\pi : I_{hist} \rightarrow U$. Now imagine trying to collapse or compress the history I-states to create a *derived I-space* I_{der} via an *information mapping* $\rho : I_{hist} \rightarrow I_{der}$ (6). All I-spaces discussed so far can be obtained in this way. For a well-defined ITF π_{der} to exist, the sufficiency condition is that $\rho^{-1}(\kappa; u_k; y_{k+1})$ is a singleton. In other words, π_{der} can be defined so that a unique I-state is obtained and is consistent with what would have been calculated from retaining full histories. It was shown in (9) that minimal sufficient information mappings, and corresponding ITFs, exist and are unique in very general settings.

We now present two alternative *model-based* families of sufficient ITFs, called *nondeterministic* and *probabilistic*. They are considered model-based because the ITF is expressed in terms of X , f , and/or h . Of the ITFs presented so far, only state feedback has been model-based and is a special case of the nondeterministic family. The I-space is $I_{ndet} = \text{pow}(X)$. The ITF π_{ndet} incrementally calculates the set of possible X-states at each stage. For each $k \geq K$, let $X_k(\kappa)$ denote the set of possible X-states at stage k given the history I-state κ . Suppose at the first stage, $X_1(\kappa_1) = h^{-1}(y_1)$ (using the preimage from (2)). The ITF calculates $X_{k+1}(\kappa_{k+1})$ using only $X_k(\kappa)$, u_k , and y_{k+1} ; thus, it takes the form

$$X_{k+1}(\kappa_{k+1}) = \pi_{ndet}(X_k(\kappa); u_k; y_{k+1}); \quad 3.$$

Assuming inductively that $X_k(\kappa)$ is given, consider the possible X-states at stage $k+1$ under the application of action u_k :

$$X_{k+1}(\kappa; u_k) = \{x_{k+1} \mid \exists x_k \in X_k(\kappa) \text{ for which } x_{k+1} = f(x_k; u_k)g\}; \quad 4.$$

Once the new sensor observation y_{k+1} arrives, the preimage h^{-1} is used to constrain the set of possible states, resulting in:

$$X_{k+1}(\kappa_{k+1}) = X_{k+1}(\kappa; u_k; y_{k+1}) = X_{k+1}(\kappa; u_k) \cap h^{-1}(y_{k+1}); \quad 5.$$

This completes the definition of π_{ndet} . Note that state feedback is a special case in which $X_k(\kappa)$ is a singleton for all $k \geq K$. We can easily extend π_{ndet} to account for nondeterministic disturbances by replacing $x_{k+1} = f(x_k; u_k)$ by $x_{k+1} \in F(x_k; u_k)$ in (4) and h^{-1} by H^{-1} in (5), in which $H^{-1} = \{x_{k+1} \mid \exists x_k \in X_k(\kappa) \text{ for which } x_{k+1} \in F(x_k; u_k)g\}$. Note that f and h (or F and H) are used internally by the agent, and they might be inconsistent with the actual external world; this will be clarified in Section 2.4 and is critical for defining illusions.

We now define the probabilistic model-based family. The corresponding I-space is I_{prob} , which is the set of all probability density functions (pdfs) over X (again, under appropriate measurability assumptions). Probabilistic disturbance-based replacements of f and h are given as $p(x_{k+1}|x_k; u_k)$ and $p(y_k|x_k)$, respectively. Let $p(x_k|j_k)$ denote the pdf of the state at stage k given j_k (the probabilistic counterpart to $X_k(\kappa)$). The ITF takes the form

$$p(x_{k+1}|j_{k+1}) = \pi_{prob}(p(x_k|j_k); u_k; y_{k+1}); \quad 6.$$

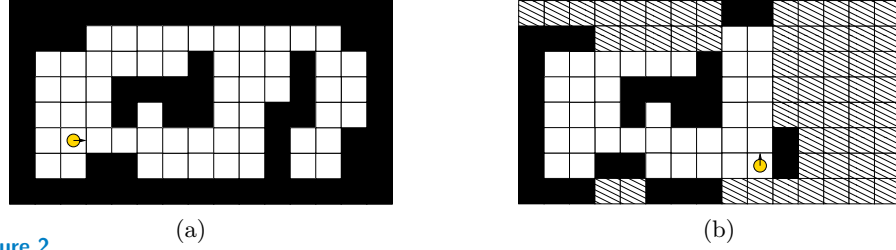


Figure 2

(a) A discrete grid problem is made in which a robot is placed into a bounded, unknown environment. (b) An encoding of a partial map, obtained from some exploration. The hatched lines represent unknown tiles (neither white nor black).

and is equivalent to Bayesian filtering and the basis of POMDPs. The process is started by a given prior $p(x_1)$. The probabilistic counterpart to (4) is marginalization, resulting in:

$$p(x_{k+1} | j_{k+1}; u_k) = \int_X p(x_{k+1} | j_k; u_k) p(x_k | j_k) dx_k. \quad 7.$$

The probabilistic counterpart to the intersection in (5) is Bayes' rule, resulting in

$$p(x_{k+1} | j_{k+1}) = p(x_{k+1} | j_{k+1}; u_k; y_{k+1}) = \mathbb{R} \frac{p(y_{k+1} | j_{k+1}) p(x_{k+1} | j_{k+1}; u_k)}{\int_X p(y_{k+1} | j_{k+1}) p(x_{k+1} | j_{k+1}; u_k) dx_{k+1}}. \quad 8.$$

Conditional independence assumptions and further details are explained in (6). The celebrated Kalman filter is a special case (linear systems with Gaussian disturbances) in which the I-states become trapped in a low-dimensional subspace of I_{prob} , and the ITF can be calculated using matrix algebra. Each I-state then corresponds to mean and covariance of the X-state at stage k .

Example 3 (Gridbot)

A mobile robot moves on a 2D grid and can face in one of four orientations (up, down, left, right), as shown in Figure 2(a). At each possible $(i; j) \in \mathbb{Z} \times \mathbb{Z}$ position there is a *tile*, which may be either *black* or *white*. If it is white, then the robot can occupy its position; otherwise, it is blocked. The robot starts on one tile among a finite, unknown, connected set of white tiles. All other tiles are black. Each possible set of white tiles is called an *environment*. The X-space is $X = \mathbb{Z}^2 \times D \times E$, in which $D = \{0; 1; 2; 3\}g$ is the set of four directions and E is the set of all possible environments. An X-state $x \in X$ can be expressed as $(i; j; d; E)$ in which $(i; j) \in E$, $d \in D$, and $E \in E$. There are two actions $U = \{0; 1\}g$, in which $u = 0$ causes the robot to rotate 90 degrees counterclockwise, and $u = 1$ makes it attempt to move forward one tile in the direction it is facing (it does not move if blocked by a black tile). A 'depth' sensor $y_k = h(x_k)$ reports the distance, in terms of number of white tiles, to first black tile in the direction the robot faces; thus, $Y = \{0; 1; 2; \dots; g\}$. Using (4) and (5), the set $X_k(j_k)$ of possible X-states is calculated after each action and observation, respectively; however, it is important to encode $X_k(j_k)$ compactly, rather than list all possible state for an infinite collection of environments. Initially, $X_1(j_1) = h^{-1}(y_1)$. During exploration, tiles sensed to be white or black are recorded using $(i; j)$ coordinates, with $(0; 0)$ as the initial white tile. All others may be labeled as gray, meaning unknown or unexplored; see Figure 2(b). This is a highly compressed encoding of $X_k(j_k)$, which technically belongs to a derived I-space of such encodings. To obtain a nondeterministic disturbance model,

h may produce a set of possible distances, and the calculations in (5) use the larger preimages $H^{-1}(y_{k+1})$. A Bayesian version can be made by introducing probabilistic alternatives to f and h , and using p_{prob} to calculate probabilistic I-states; see (10) for details.

2.3. Defining Plausibility and Illusions

Here we address the possibility that the models used in the ITF do not perfectly coincide with ‘reality’ as the agent interacts with its environment. Of course, reality itself will be defined as a model, but it is nevertheless crucial to maintain a distinction. Thus, we refer to X , f , and h used in the agent’s ITF as *intrinsic* models. To provide an outside frame of reference, we will introduce their counterparts Ω , \bar{f} , and \bar{h} , and refer them as *extrinsic* models. If there is no disagreement between the intrinsic and extrinsic models, then $\Omega = X$, $\bar{f} = f$, and $\bar{h} = h$; this distinction was not yet needed in Section 2.2. Discrepancies between the intrinsic and extrinsic models will be crucial for modeling perceptual illusions.

Let Ω be called the *universe space*, or just the *universe*, which models the set of all possible physical states of the world from a third-person or god-like perspective. At each stage $k \geq K$, the universe state is $!_k \in \Omega$ and the *universe transition function*, *UTF*, is defined as $\bar{F}: \Omega \rightarrow \Omega$. Similarly, the *universe sensor mapping* is defined as $\bar{h}: \Omega \rightarrow Y$.

Next, we define a general way to model potential correspondences between X-states and universe states from a third-person perspective (outside of the agent). For any X and Ω , let $C \subseteq X \times \Omega$ be called a *correspondence relation*. If $(x_k; !_k) \in C$ it is said that x_k *corresponds to* $!_k$. Note that C allows for the correspondence to be one-to-one, many-to-one, one-to-many, or many-to-many. The relation C is called *onto* if for all $!_k \in \Omega$, there exists an $x_k \in X$ such that $(x_k; !_k) \in C$. If C is one-to-many and onto, then there exists a function $\pi: \Omega \rightarrow X$ such that $(\pi(!_k); !_k) \in C$ for all $!_k \in \Omega$; thus, the X-state can be derived from the universe state as $x_k = \pi(!_k)$.

The relationship between the agent’s I-state and ‘reality’ in the universe can be established through their relationships to X . Let the *model relation* $M \subseteq I \times X$ associate I-states to possible X-states. If the nondeterministic ITF family is used, then $(X_k(!_k); x_k) \in M$ if and only if $x_k \in X_k(!_k)$. (For probabilistic ITFs, M may be defined using thresholds on pdfs to obtain probabilistic correspondences.) An I-state $!_k$ is called *implausible* if there does not exist any $x_k \in X$ such that $(!_k; x_k) \in M$. Thus, $X_k(!_k) = \emptyset$ would be an implausible I-state for nondeterministic ITFs. A pair $(!_k; x_k)$ or I-state $!_k$ is called *plausible* if it is not implausible. Its usage here is inspired by concepts of plausibility in VR research (11, 1), but also differs in precise meaning.

By composing the model and correspondence relations, the *reality relation* $R \subseteq I \times \Omega$ is defined as $(!_k; !_k) \in R$ if and only if there exists $x_k \in X$ such that $(!_k; x_k) \in M$ and $(x_k; !_k) \in C$. A pair $(!_k; !_k)$ is called an *illusion* if $!_k$ is plausible and $(!_k; !_k) \notin R$. Thus, the key idea of an illusion is that the agent perceives something as plausible but it does not correspond to reality.

2.4. Multiple Agents

Suppose there are n agents in a common environment. The i th agent is modeled using the components from Section 2.1 and denoted as an 8-tuple $A^i = (X^i; I^i; U^i; Y^i; f^i; h^i; \pi^i; \rho^i)$. Here, X^i , f^i , and h^i are intrinsic, a distinction that was unnecessary in Section 2.1; thus, as in Section 2.3, we seek their extrinsic counterparts. The universe sensor mapping for

each agent is $\bar{h}^i: \Omega \rightarrow Y^i$. The UTF must take into account the interactions between agents. It is therefore specified centrally as $\bar{f}: \Omega \rightarrow U^1 \times \dots \times U^n \rightarrow \Omega$, as is common in dynamic game theory. If the agents' actions do not interfere with each other, then \bar{f} may be decomposable into individual \bar{f}^i functions over subspaces of Ω , but this will not be assumed. Correspondence, model, and reality relations can be defined over $X^i \subseteq \Omega$, $I^i \subseteq X^i$, and $I^i \subseteq \Omega$, respectively.

Example 4 (Two Gridbots)

We extend Example 3 by placing two gridbot agents, A^1 and A^2 , in a universe defined as $\Omega = Z^4 \times D^2 \times E$, which is the set of all $(i^1; j^1; i^2; j^2; d^1; d^2; E) \in \Omega$ that satisfy $(i^1; j^1) \in E$, $(i^2; j^2) \in E$, $E \subseteq E$, $d^1 \subseteq D$, $d^2 \subseteq D$, and $(i^1; j^1) \neq (i^2; j^2)$. Each agent has an X-space X^i as defined in Example 3 for a single robot. Correspondence relations C^i are defined for $i = 1; 2$. They are subsets of $X^i \subseteq \Omega$ and are one-to-many and onto, implying the existence of functions $f^i: \Omega \rightarrow X^i$ that are consistent with C^i . Each f^i maps I to E and the position and direction of the i th agent in E , all of which are expressed using the local coordinates of A^i because A^1 and A^2 assign position $(0;0)$ and direction 0 differently.

Each universe sensor mapping \bar{h}^i differs from h^i by taking into account the other robot. Let \bar{h}^i return the directional distance to the nearest black tile or other robot, if it is closer (the other robot is like a movable black tile). The UTF \bar{f} can be derived to coincide with the individual f^1 and f^2 XTFs, but must specify what happens when robots attempt to move into each other. Suppose each cannot move onto a tile occupied by the other, and the first agent has priority if they attempt to move to the same white tile at the same time.

For multiple agents, the I-spaces and ITFs are more challenging to model due to the effects of their interactions, whether accidental or intentional. For Example 4, what happens when one robot is blocked by the other? It seems incorrect to label the tile as black. Perhaps later it will discover that the tile is white. Does its model allow E to change? Does it 'know' there is another robot? To define each ITF, an agent's intrinsic model must carefully specify what information regarding other agents is available. Using all sources of potential information, an agent's ITF could be expressed as $I_{k+1}^i = I^i(x_k^1, \dots, x_k^n; u_k^1, \dots, u_k^n; y_{k+1}^1, \dots, y_{k+1}^n)$, in which the definition of I^i may reference any component of the 8-tuple A^i for any agent, any \bar{h}^j , any correspondence, model, and reality relations, and \bar{f} .

Consider designing one agent, say A^1 , as *omniscient*, meaning that it has access to as much information as possible. Let $X^1 = \Omega$ and $y_k^1 = h^1(x_k^1) = h^1(I_k) = I_k$, thus observing the universe state at every stage. Its own I-state I_k^1 records the histories of all actions, observations, and I-states for itself and all other agents. This is a multiagent extension of I_{hist} from Section 2.2. Using the notion of sufficient information mappings

from Section 2.2, this I-space can be collapsed to smaller I-spaces as appropriate for accomplishing particular tasks. In general, the model components of each A^i , each \bar{h}^i , the relations, and \bar{f} may be used to define I_{der}^1 over the derived I-space.

The information gathering capabilities of an omniscient agent may seem to be too much. Most 'ordinary' agents will have far less capabilities, but it is nevertheless important to define the extreme case as a starting point. How could I-states of other agents be obtained? One way is to predict them through simulation or computation using the other data and models. If they are measured directly, then a sensor model should be formulated that includes I-states as part of the physical universe. How could observations or actions of other agents be obtained? Similarly, they could be estimated from other information, such

as using \bar{F} and the measured I_k and I_{k+1} to determine the actions, or \bar{H}^i and the measured I_k to determine the observation. Otherwise, they should be directly sensed in an expanded universe. Generally, each agent obtains a history I-state I_k^i , resulting in a history I-space I_{hist}^i , which can be reduced using sufficient information mappings. To allow precise models of observing I-states, actions, and observations, the universe should be expanded to include them; however, this technical level is beyond the scope of the current paper.

In terms of actuation, at one extreme, one agent could be *omnipotent*, enabling it to set any $I_k \supseteq \Omega$ at any stage $k \supseteq K$. We do not allow multiple omnipotent agents because their actions would be in conflict. We also do not allow it to directly set I-states, observations, or actions of other agents. At the other extreme, it could be a passive *observer*. Consider what happens when we, the modelers, try to mathematically analyze the entire system. In this case, we are acting as an omniscient agent that is merely an observer so as not to interfere with its operation while analyzing what should happen theoretically. All other agents are in between being omnipotent and an observer, and they may or may not be omniscient.

Extensions can also be made to account for disturbances. In the nondeterministic case, the universe sensor mapping of the i th agent becomes $\bar{H}^i: \Omega \rightarrow \text{pow}(Y^i)$. The UTF becomes $\bar{F}: U^1 \times \dots \times U^n \rightarrow \text{pow}(\Omega)$. Similarly, for the probabilistic case, the corresponding corresponding universe sensing model and transition function are $\rho(y_k^i | I_k)$ and $\rho(I_{k+1}^j | I_k; U_k^1; \dots; U_k^n)$, respectively.

3. CREATING TARGETED PERCEPTUAL EXPERIENCES

3.1. Producers and Receivers

We now adapt the multiagent model of Section 2.4 to the case of a special agent A^p called a *producer* that delivers a targeted perceptual experience to another agent A^r called a *receiver*. Key concepts running throughout Section 3 are plausibility and illusions, as defined in Section 2.3. To allow it to ‘fool’ the receiver in some sense, the producer usually has access to more information than the receiver. For example, the producer may have access to models \bar{F} and \bar{H}^r , whereas the receiver only has F^r and H^r . Section 3.2 starts with the simplest case, in which stage dependencies and transitions are suppressed: A fixed percept (receiver I-state) is created and maintained. Section 3.3 will then extend the concepts to cover an omniscient producer that operates over multiple stages, resulting in a targeted, interactive, perceptual experience for the receiver. Section 3.4 will then strip the producer of its omniscience, resulting in incomplete or imperfect models of the receiver and even the producer itself. This is more like the situation of VR applied to biological organisms, but also relevant for engineered systems.

3.2. Producing Stationary Percepts and Illusions

This section temporarily drops the notion of stages to provide a useful and illustrative setup before developing more complicated scenarios. The receiver’s sensor model is $h^r: X^r \rightarrow Y^r$. The nondeterministic I-space I_{ndet} is the set of all preimages of h^r . Thus, the I-state, called a *percept*, is simply fixed as the preimage $I^r = (h^r)^{-1}(y^r) \subseteq X^r$. The producer is omniscient, and $X^p = \Omega$. Let $(X^p; I) \supseteq C^p$ if and only if $X^p = I$ (they are in perfect one-to-one correspondence). The producer is able to set I according to some given function $\bar{F}: U^p \rightarrow \Omega$ (a simplified version of \bar{F} from Section 2.4 to handle the stationary case), after which time I remains constant. It also has access to $\bar{H}^r: \Omega \rightarrow Y^r$ and the receiver’s

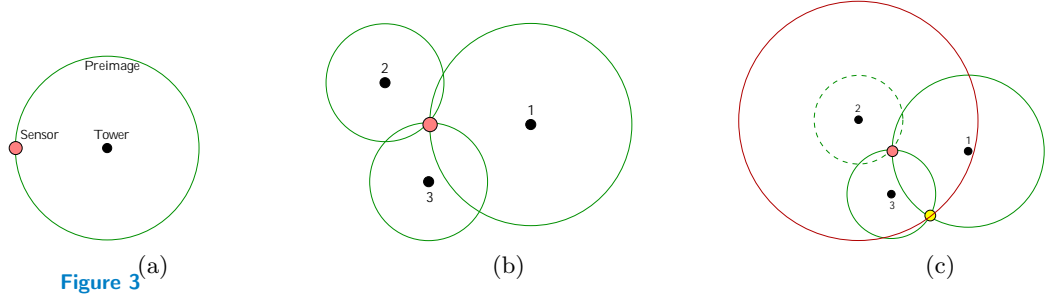


Figure 3

(a) A receiver's sensor measures the distance to a tower, resulting in circular preimages as I-states. (b) For three towers, the correct position is given by the intersection of three preimages. (c) A producer can change the signal intensity of the second tower to create an illusory perceived position for the receiver.

correspondence relations C^r and M^r . Suppose that C^r is one-to-many and onto so that $\bar{r}: \Omega \rightarrow X^r$ exists. The model relation is defined as $(r; x^r) \in M^r$ if and only if $x^r \in r$. This is equivalent to I_{ndet} restricted to a single stage in which only one observation is available.

If $\bar{h}^r = \bar{h}^r \circ r$, then $(r; !) \in R^r$ for all $! \in \Omega$. To see why, suppose the producer sets $! = \bar{f}(u^p)$ for some $u^p \in U^p$, thereby causing the receiver to observe $y^r = \bar{h}^r(!) = h^r(r(!))$, and obtain a targeted I-state $r = (h^r)^{-1}(y^r)$. For any possible u^p , the resulting pair $(r; x^r)$ would belong to M^r , and hence r is always plausible. Furthermore, $(r; !) \in R$, implying the pair is not an illusion.

The potential to create an illusion arises if $\bar{h}^r \notin \bar{h}^r \circ r$. This would allow some $!$ to produce $y_1^r = \bar{h}^r(!)$ and $y_2^r = h^r(r(!))$ such that $y_1^r \neq y_2^r$. The preimages $r_1 = (h^r)^{-1}(y_1^r)$ and $r_2 = (h^r)^{-1}(y_2^r)$ are distinct and disjoint. The corresponding X-state would be $x^r = r_1(!)$, with $x^r \notin r_1$ and $x^r \in r_2$. Thus, if $r_2 \notin r_1$ and $(r_2; !) \in R$, then the pair is an illusion as defined in Section 2.3.

Example 5 (Moving a Landmark)

Let $X^r = \mathbb{Z}$ and $\Omega = X^p = X^r = \mathbb{Z}$, as an omniscient producer. The correspondence relation C^r is one-to-many and onto, and $r(!) = x_1$, denoting $! = (x_1; x_2)$. The receiver's sensor model is $y^r = h^r(x) = x^r$, which measures its state perfectly. The producer is given the extrinsic sensor mapping $\bar{h}^r: \Omega \rightarrow Y^r$, defined as $y^r = \bar{h}^r(x_1; x_2) = x_2 - x_1$, which is the signed distance from $x_1 = x^r$ to the location x_2 of a movable reference point or 'landmark'. The producer can thus create an illusion for the receiver that its position is any $x^r \in X^r$ by changing x_2 . From a third-person perspective, we could interpret the receiver's sensor model as reporting the signed distance to a tower fixed at $0 \in X^r$, but x_2 corresponds to the true position of the tower. In this case, the producer moves the tower position, which is beyond the receiver's model. Its I-state $r = (h^r)^{-1}(y^r) = \bar{f}x^r g$ therefore implies that its perceived own position $x^r \in X^r$ is wrong, which is an illusion. Note that $\bar{h}^r(!) \neq h^r(r(!))$, as required.

Example 6 (Trilateration Tricks)

Planar localization is performed by a receiver, in which $X^r = \mathbb{R}^2$, by the principle of *trilateration*. The receiver's sensor observation y^r reports the distances to one or more

‘towers’ that serve as known landmarks with fixed position in \mathbb{R}^2 . Using the reported distance to a single tower, the I-state $r = (h^r)^{-1}(y^r)$ narrows the position down to a circle (see Figure 3(a)) of radius y^r . If there are n towers, then $y^r = (d_1^r; \dots; d_n^r)$ is a vector of observed distances to each tower. Figure 3(b) shows the case of three towers, in which the I-state results in a unique receiver position by intersecting the three circular preimages, one for each distance measurement.

Now enters the producer, which exploits an extrinsic sensor model $y^r = \bar{h}^r(!)$. The universe state is $! = (x_1; x_2; r_1; \dots; r_n) \in \Omega$, in which $(x_1; x_2)$ is the receiver X-state, r_i is the transmitted radio signal intensity of i th tower, and $\Omega \subset \mathbb{R}^{n+2}$. The observation d_i^r is actually based on the inverse-square law, $r_i^o = r_i = d_i^o$, in which d_i is the actual distance to the i th tower and r_i^o is the measured intensity at the receiver position. Let $r^c = (r_1^c; \dots; r_n^c)$ be the vector of intensities that the sensor is calibrated to. Each element of y^r is then obtained as $d_i^r = \sqrt{\frac{r_i^c}{r_i^o}}$ if $r_i^c > r_i^o$ otherwise, $d_i^r = \#$, which indicates that the signal intensity is not within the interval of observable values. Note that observing $r_i^o = r_i^c$ implies that the receiver is on the i th tower, which is not allowed. To yield any desired distance measurement d_i^r so that d_i^r becomes d_i^r , the producer changes the transmitter power r_i from r_i^c to $r_i^o = r_i^c(d_i = d_i^r)^2$. For example, if $r_i^c = 1$, then $r_i^o = 1 = 4$ at distance $d_i = 2$. Changing the transmitted intensity to $r_i^o = 2$ would cause the distance observation to be $d_i^r = \sqrt{2}$ ($r_i^o = 1 = 2$), which is incorrect. Setting $r_i^o = 4$ would then produce implausible I-states because $r_i^o = r_i = 1$, and the respective circle will be undefined.

The producer X-space is $X^p = \Omega$. The function r^c converts the producer X-state to the receiver X-state, and C^r is defined accordingly. The producer action space is $U^p = (0; 1)^n$, and is used in r^c to set each radio intensity r_i to any positive value, inducing a desired d_i^r so that $y^r = (d_1^r; \dots; d_n^r)$. In the case of $n = 1$, the producer interference can create receiver I-states that are circles of any radius or distance from the single tower. For $n = 2$, there is a bounded range of radio intensities (values of U^p) for which the circular preimages intersect; if they are disjoint, then the illusion becomes implausible (there is no receiver position that would account for y^r , and $r^c = \#$). For $n = 3$, shown in Figure 3(c), any small perturbation of U^p results in implausibility, $r^c = \#$.

From Example 6 two kinds of concepts become clear: 1) If the I-state is plausible, then how much can the observation be perturbed while maintaining plausibility? 2) If the I-state is not plausible, then how little can the observation be perturbed to make it plausible? Both of these concepts rely on a notion of distance in Y^r . Thus, assume Y^r is a metric space with metric $\gamma_r: Y^r \times Y^r \rightarrow [0; 1]$. Let $B_{Y^r}(y; d) \subset Y^r$ denote an open ball of radius d , centered at y . Thus, $B_{Y^r}(y; d) = \{y^j \in Y^r \mid \gamma_r(y; y^j) < d\}$.

The first concept is defined as follows. Let $P \subset Y_r$ be the set of all observations that yield plausible I-states, defined using the model relation M^r . The *plausibility robustness*, $PR_{Y^r}(y^r)$, is the largest d such that $B_{Y^r}(y^r; d) \subset P$. The second concept is *forced-fusion magnitude*, $FFM_{Y^r}(y^r)$, which is the largest d such that $B_{Y^r}(y^r; d) \cap Y^r \cap P$. In other words, it is the distance to the nearest receiver observation that would yield plausibility. The term is inspired by its use to account for adverse symptoms in VR usage (12).

We would like to express these concepts from the producer’s perspective. For plausibility robustness, how much can the producer vary U^p and still maintain plausibility? For forced-fusion magnitude, how much does the producer need to change U^p to achieve plausibility? Assume U_p is a metric space with metric γ_{U^p} . We can similarly define PR_{U^p} and FFM_{U^p} in terms of balls in U^p of radius d centered at u^p . Plausibility robustness $PR_{U^p}(u^p)$ requires

that $\bar{h}^r(\bar{f}(u)) \geq P$ for every $u \in B_{UP}(u^p; d) \subseteq U^p$, and forced-fusion magnitude $\text{FFM}_{UP}(u^p)$ requires that $\bar{h}^r(\bar{f}(u)) \geq P$ for every $u \in B_{UP}(u^p; d) \subseteq U^p$. Note that if $U^p = X^p$ and \bar{f} is the identity function, then these concepts can also be expressed directly in terms of X^p .

Now consider stating the producer's goal so that one might select $u^p \in U^p$ systematically, and even autonomously. Perhaps the goal is to achieve a particular observation, say $y_G \in Y^r$, or more generally, any one in a nonempty set Y_G of observations. The producer must select u^p so that $\bar{h}^r(\bar{f}(u^p)) \geq Y^p$. Further conditions might be that plausibility or illusions must be maintained by considering the resulting I-states (preimages of h^r). Thus, the goal could alternatively be expressed directly in terms of I-states, which correspond to targeted perceptions: achieving an I-state $G \in I^r$ or set $I_G \subseteq I^r$ of I-states.

If there are multiple producer actions u^p that achieve the goal, then an optimization problem can be formulated. Let $l(u^p)$ be the cost of applying u^p . An *optimal perception* (or equivalently, *optimal I-state*) would be the one that yields the minimum cost $l(u^p)$ over all u^p that achieve the goal (Y_G , Y_G , G , or I_G). This assumes appropriate conditions for the existence of optima, such as compactness over the set of producer choices that achieve the goal. Robustness and forced-fusion can also be taken into account. For example, if $y^r \in Y^r$, then $l(y^r) = c + \text{PR}_{Y^r}(y^r)$, and if $y^r \notin Y^r$, then $l(y^r) = c + \text{FFM}(y^r)$, for some constant c . An action can be chosen so that $l(y^r)$ is minimized among goal perceptions to maximize robustness; otherwise, it is chosen to minimize the forced-fusion magnitude.

The models extend naturally to handle disturbances. In the nondeterministic case, $H^r: X^r \rightarrow Y^r$ and $\bar{H}^r: \Omega \rightarrow \text{pow}(Y^r)$ are used instead of h^r and \bar{h}^r . A disturbance could even be added to the producer's action so that \bar{f} is replaced by a function $\bar{F}: U^p \rightarrow \text{pow}(X^p)$. Thus, with each action, u^p , only a set $\bar{H}^r(\bar{F}(u^p)) \subseteq Y^r$ of possible observations can be enforced. The producer is guaranteed to be successful in the worst-case if $\bar{H}^r(\bar{F}(u^p)) \subseteq Y_G$. If $\bar{H}^r(\bar{F}(u^p)) \setminus Y_G \neq \emptyset$, then it may *possibly* be successful, in the best-case. Similarly, the case of probabilistic disturbance can be considered, using $p(y^r \in Y^r)$ and $p(y^r \in I^r)$. Disturbance can be added to producer actions as a pdf, $p(x^p \in U^p)$. The probability that a goal Y_G is achieved is given by

$$p(Y_G \in I^r) = \int_{Y_G} \int_{X^p} p(y^r \in Y^r) p(x^p \in U^p) dx^p dy_G: \quad 9.$$

Thus, the producer could try to choose $u^p \in U^p$ to maximize the probability of success. For the disturbance-based models, costs could still be formulated. In the nondeterministic case, u^p can be chosen to minimize the worst-case (maximum) cost. In the probabilistic case, u^p can be chosen to minimize the expected cost. The costs could include $\text{PR}^r(y^r)$ and/or $\text{FFM}(y^r)$, resulting in worst-case or expected-case analysis of plausibility robustness and forced-fusion magnitude.

3.3. Creating Perceptual Experiences with an Omniscient Producer

We now extend the concepts from Section 3.2 from the stationary case to the dynamic case. The extension of a stationary percept (or I-state) to cover multiple stages is called a *perceptual experience* (or I-state trajectory). The producer and receiver are modeled using their corresponding 8-tuples, A^p and A^r . As before, the producer remains omniscient with $X^p = \Omega$, and C^p modeling perfect correspondence. It has access to the UTF $\bar{f}: \Omega \rightarrow U^p$, $U^r \rightarrow \Omega$, \bar{h}^r , and $h^p = \bar{h}^p$ is the identity function (perfect sensing). It has access to \bar{r}_k , u_k^r , and y_k^r at each $k \in K$. It also has access to the relations C^r and M^r (from which R^r can be derived). Relaxing these strong assumptions will be discussed in Section 3.4.

Consider going from stage k to $k+1$. The universe state is some $!_k \in \Omega$. The producer can implement a state-feedback policy of the form $U^p: \Omega \rightarrow U^p$ (using the fact that $X^p = \Omega$). Thus, an action $u_k^p = U^p(!_k)$ is selected. The receiver action is $u_k^r = U^r(!_k)$, and $!_{k+1} = \bar{f}(!_k; u_k^p; u_k^r)$. The next receiver observation is $y_{k+1}^r = \bar{h}(!_{k+1})$. The receiver I-state is updated as $!_{k+1}^r = \bar{r}(!_k; u_k^r; y_{k+1}^r)$, potentially using its intrinsic models X^r , f^r , and h^r . Disturbance-based models could alternatively be used, including F^r and H^r or their probabilistic counterparts.

For every stage k , the producer applies u_k^p to influence y_{k+1}^r and $!_{k+1}$. This part is quite similar to Section 3.2, in which U^p was chosen to influence y^r and $!^r$; however, here there is a one-stage delay.¹ Using M^r , the producer (or the engineer who created it) can determine whether each $!_k$ is plausible. Let $\tilde{!}^r$ denote an I-state sequence, called a *perceptual experience*, that is indexed over $k \geq K$. If every $!_k$ in $\tilde{!}^r$ is plausible, then it is called a *plausible perceptual experience*. Let $!^r$ be the universe state trajectory corresponding to some $\tilde{!}^r$. We can apply M^r to determine whether each $(!_k; !_k)$ is an illusion. If the pair is an illusion for every $k \geq K$, then the pair $(\tilde{!}^r; !^r)$ is called an *illusory perceptual experience*.

Example 7 (A Dynamic Landmark)

Building upon Example 5, let $!_k = (x_k^r; x_k^p)$ and $!_{k+1} = \bar{f}(!_k; u_k^p; u_k^r) = (x_k^r + u_k^r; x_k^p + u_k^p)$, in which $U^p = U^r = \{-1; 0; 1\}$. The landmark position can be moved up or down by one unit, or remain stationary, and the receiver can similarly change its own position. We have $f^r(x_k^r; u_k^r) = x_k^r + u_k^r$ and $f^p(x_k^p; u_k^p) = x_k^p + u_k^p$. Suppose h^r functions as in Example 5. Suppose $u_1^p = 1$ and $u_1^r = 0$. This results in an implausible I-state $!_2 = ;$ because the position predicted by f^r is inconsistent with the observation y_2^r (the sets given by (4) and the preimage $(h^r)^{-1}(y_2^r)$ are disjoint). To give the producer more freedom, nondeterministic replacements F^r and H^r can be used for the receiver's intrinsic models. For example, if $F(x_k^r; u_k^r) = X^r$ for all actions, then the I-states would be based on preimages of h^r alone, causing the receiver to have the illusion it is moving when in fact the producer is moving.

Example 8 (Gridbot Illusions)

Example 4 can be adapted by interpreting the gridbots as a producer and receiver. The producer could create an illusion that the receiver's environment is smaller than it really is. For example, suppose E corresponds to two large rooms, connected by a 'doorway' of width one tile. The producer moves to the doorway and remains there while the receiver explores. The illusion of a smaller room has been created. If the producer moves away and the receiver returns to the doorway, an implausible I-state would result because a tile that was marked as black became white. The challenge is to determine a policy for the producer so that it is not detected by the receiver; this would happen if the receiver sensed a tile as white when it was previously declared black, resulting in implausibility (assuming in its model it is unaware of the other robot).

Now consider designing a producer policy. Under the standard agent model from Section 2.1, note that the entire system is predictable starting from any initial $!_1 \in \Omega$, which is immediately observable to the producer. Call this the *fully predictable* case. The producer

¹Note that in the first stage $k = 1$, the producer does not have the ability to affect y_1^r ; this may be fixed by deleting the first observation or allowing the producer to start one stage earlier.

policy may as well be a sequence of actions, which is a stage-feedback policy $\rho: \mathcal{K} \rightarrow \mathcal{U}^P$, as defined in Section 2.2. The goals from Section 3.2 can be extended here to sequences. For example, let $\tilde{y}_G: \mathcal{K} \rightarrow \mathcal{Y}^r$ be a *goal sequence* of receiver observations. A weaker requirement is to simply achieve any one \tilde{y}_G in a set \tilde{Y}_G of possibilities. Many other possibilities exist. For example, perhaps the goal is to produce \tilde{y}_G , or any observation in a set \tilde{Y}_G , at *any* stage $k \geq K$. Alternatively, perhaps it must happen at one particular stage. A logic, such as linear temporal logic (LTL), may even be used to express goal conditions in terms of some combinations of sets of observations and stages (13). Furthermore, goals could also be expressed in terms of receiver states, receiver actions, receiver I-states, or any combination along with receiver observations and stages. On top of this, a cost function l_k can be defined at each stage to obtain a problem of finding an optimal producer action sequence, which in turn yields an *optimal perceptual experience* by optimizing the cumulative cost

$$\sum_{k \geq K} l_k(\tilde{I}_k; \mathcal{U}_k^p; \tilde{y}_k^r; \tilde{r}_k; \mathcal{U}_k^r): \quad 10.$$

The costs can also be expressed in terms of PR and FFM functions to obtain problems that try to maximize total robustness or minimize total forced-fusion magnitude. If \mathcal{K} is infinite, then the costs must be carefully chosen so that the sum is finite for successful policies; alternatives include discounted cost, average cost, and termination actions (6). If \mathcal{K} is finite, then a final cost term $l_F(\tilde{I}_{K+1}; \mathcal{Y}_{K+1}^r; \tilde{r}_{K+1})$ may be added.

Now suppose that disturbance-based extensions of f^r and h^r are introduced for the receiver, to obtain F^r and H^r , as defined in Section 2.1. In this case, the receiver is no longer predictable, even from the producer's perspective. It is thus more effective for the producer policy to be formulated as state-feedback $\rho: \mathcal{X}^p \rightarrow \mathcal{U}^p$, which even implies universe-state feedback. Acknowledging that this may be extreme in many settings, Section 3.4 removes producer omniscience to obtain other cases of information-feedback policies for the producer.

Consider characterizing the evolution of the whole system under the implementation of a fixed, state-feedback producer policy ρ . Under the fully predictable case, a sequence $\tilde{I}: \mathcal{K} \rightarrow \Omega$ is determined from the initial universe state \tilde{I}_0 . In the case of a nondeterministic disturbance-based receiver, then a set $\tilde{\Omega}$ of possible sequences is instead obtained. If Ω is finite, then the process can be imagined as a nondeterministic finite automaton (NFA) over Ω . One should consider worst-case analysis to determine whether a goal can be guaranteed to be accomplished. With a cost model, one can consider minimizing the worst-case perceptual experience. In the case of a probabilistic disturbance-based receiver, a Markov chain is obtained under the implementation of ρ (also called Markov decision process (MDP) by artificial intelligence researchers). In this case, expected-case analysis could be used to assess the probability that the goal will be satisfied under ρ . In this case, ρ can be selected to maximize this probability. A cost model can additionally be used, with the resulting optimization being to find the lowest expected-case cost under the implementation of ρ .

3.4. Producers with Imperfect Information

If the producer is not omniscient, then it may not have access to enough information to ensure that the targeted perceptual experiences function as desired. To analyze what might happen between the producer and receiver, it will be helpful to nevertheless introduce a third-person perspective in which we as scientists or engineers have access to more information than the producer. This could be modeled formally as an observer agent.

The following producer limitations can be considered: 1) P and/or h^r could be many-to-one mappings, prohibiting the producer from perfectly determining $!_k$; 2) r may be unobservable to the producer, resulting in a *dynamic game* formulation, which may or may not be cooperative; 3) r may be partly or fully hidden, requiring the producer to estimate it via models, simulation, and limited sensors; 4) the producer may have only limited models or access to h^r , \bar{h}^r , f^r , and \bar{f}^r , resulting in partial control over illusions, and difficulty determining whether illusions are plausible. 5) the producer may not have perfect access to its own state, leading to information feedback policies $P: !^P ! U^P$ that hopefully achieve the targeted perceptual experiences.

Goals for targeted perceptual experiences may be formulated as in Section 3.3, and a producer policy P is selected that achieves the goal, and even optimizes costs. For nondeterministic models, worst-case or even best-case analyses are appropriate. For probabilistic models, expected-case analysis can be used once again.

3.5. Multiagent Perceptual Experiences

We can extend the formulations developed so far to allow for multiple producers and receivers within a shared universe. Suppose, for example, that there is one producer and n receivers. The producer could use the same spoofing function S to stimulate all of them at once. This could be imagined as a broadcasting mechanism, such as wireless communication, that reaches all receivers in the same way. Thus, the delivery of the perceptual experience could be considered as a kind of *centralized control*, applied by the producer. We can alternatively formulate a *distributed control* scenario in which one or more producers deliver perceptual experiences and illusions propagate to receivers in a communication network. One further extension is to allow multiple agents to be situated within a single body. This could, for example, be used to model hierarchy, in which one agent plays a supervisory role over agents that function as lower-level models.

4. APPLICATION TO ROBOTS AND OTHER ENGINEERED AGENTS

4.1. Modeling Engineered Agents

In engineering, we typically have white-box systems, which are built from well-understood physical principles and work as designed (as opposed to black-box systems, for which there is no understanding about their internals). Learning, identification, and calibration processes may serve to further refine and improve their models with respect to their environment. At a high level, any engineered agent can be modeled as a control system, including regulators of physical systems such as aircraft stability, room temperature, or the concentration of chemical solutions. To cover many cases, a *linear agent* can be expressed as a discrete-time linear control system, which can be formulated as $X = \mathbb{R}^n$, $U = \mathbb{R}^m$, XTE $x_{k+1} = Ax_k + Bu_k$, and sensor mapping $y_k = Cx_k$ for fixed matrices A , B , and C . State-feedback or information-feedback policies take the form $: X ! U$ or $: ! ! U$, respectively.

Robot models typically involve nonlinear dynamical systems over configuration manifolds, often with non-trivial topology and non-Euclidean geometry. Imagine extending the gridbot from Example 3 to operate as a wheeled mobile robot, such as a robotic vacuum cleaner. The configuration space \mathcal{C} is used in robotics to model the set of all ways to embed the robot body in its environment in the absence of obstacles. For a wheeled mobile robot, \mathcal{C} could be the set of all 2D rigid body transforms: $\mathcal{C} = SE(2) \times \mathbb{R}^2 \times S^1$, in which S^1 is the

circle of all directions from 0 to 2π (compare to $Z^2 = D$ from Example 3). If the obstacles are unknown, then an environment space E could represent a set of possible subsets of C that are collision free, and then $X = C \times E$, in which $q \in C, q \in E, E \subseteq E$ for any $(q; E) \in X$. Actions $u_k \in U$ correspond to commands that cause the wheels to rotate for some time Δt , thereby altering the configuration to obtain $x_{k+1} = f(x_k; u_k)$.

Onboard sensors are modeled by h and might report whether obstacle contact is made, wheel odometry, and even distances to obstacles. An onboard computer calculates I-states based on sensor observations. The calculated I-states are used to apply actions according to a policy $u_k = \pi(I_k)$. Other types of robots, such as 3D drones, industrial manipulator arms, humanoids, or submarines, are similarly modeled using configuration spaces, f , and h . To handle higher order dynamics, the configuration space may be extended to a higher-dimensional phase space that includes configuration velocities, and f and h are defined over it.

The result is a fully modeled, white-box system, which will highly contrast the modeling challenges for humans, discussed in Section 5.1. Nevertheless, in some settings robot models may be adaptive as action-observation pairs are accumulated during execution. Models can be adjusted via machine learning or improved calibration. A black-box setting may even appear if the engineer approaches an unknown robot, in which case perception engineering may be used to help understand its behavior.

4.2. Spoofing Sensors

Consider fooling sensors that might be used in robots. The physical operation of h in terms of the universe Ω is modeled as $\tilde{h}: \Omega \rightarrow Y$ (from Section 2.4), which allows a producer to create illusions for the receiver robot and result in $\tilde{h}(I) \notin h^{-1}(r(I))$. Furthermore, obtaining plausible illusions and experiences for a robot is generally challenging because there are multiple sensors giving observations at multiple times. All of these must be consistent with respect to a sensor fusion system in the sense that a possible receiver X-state trajectory could explain it in terms of h and f .

Localization is the classic problem of estimating the robot's configuration. If obtained by the trilateration system of Example 6, then localization illusions could be obtained by changing tower intensities. Alternatively, the towers themselves could be moved. Other wireless localization systems, including GPS, could be similarly spoofed. *Mapping* is the problem of determining the robot's environment, usually in terms of obstacles that must be avoided; it is usually combined with localization to obtain *SLAM* (10). Many depth measurement systems work by emitter-detector pairs, such as sonars that emit a pulse and use the time of arrival to calculate distance. Such sensors can be spoofed by blocking the emitted pulse, and sending an alternative pulse to the receiver at the desired time. Methods for spoofing lidars for autonomous driving appear in (14). Cameras that infer distance based stereo could be intentionally misaligned to give false results. Features in images, assumed to be fixed in the world, could also be moved by a producer. In an extreme case, a graphical display could even be placed in front of a camera. Even a mechanical contact sensor, which is triggered when a robot hits a wall, could simply be fooled by pressing on it (15).

Mechanical sensors embedded in the body are the hardest to spoof. For example, a modern inertial measurement unit (IMU) uses vibrating MEMS to estimate angular velocity and linear acceleration; this can be spoofed by injecting acoustic vibrations (16). *Odometry* and *joint encoders* report how far wheels have rolled and joints have rotated, respectively.



Figure 4

(a) A room mapped by a Neato Botvac D5 before interference. (b) A producer (human) with cardboard causes sensor observations corresponding to a virtual wall. (c) The receiver robot reports that it is done cleaning a smaller room than exists in reality (but its depth sensor measures some further away walls).

These are similar to proprioceptive senses in humans. These could be spoofed by mechanical intervention, such as placing a mobile robot up on rack while the wheels rotate in the air. This setup would use \bar{f} to additionally compensate for applied receiver actions so that plausible X-state transitions are nevertheless obtained.

4.3. Virtual Reality for Robots

With the ability to spoof sensors, we can next consider offering VR to an engineered agent, analogous to VR experienced by humans. As a thought experiment, imagine a humanoid robot wearing a VR HMD. Assuming it has cameras for eyes, the HMD might fool the sensor fusion system, though it is unlikely to work the same way as intended for humans. The humanoid might walk and build an environment map that is consistent with the HMD imagery, with implausibility arising when it hits a real or virtual wall (the same would happen for a human using VR). For a very different scenario, imagine offering VR, or a targeted perceptual experience, to a vacuum cleaning mobile robot; see Figure 4(a). In (15), researchers fooled the robot by moving a piece of cardboard around quickly so that the robot's wall contact sensor was activated as desired to create the illusion (Figure 4(b)). The robot concluded it was in a smaller room than in reality, and it reported it was done cleaning (Figure 4(c)). The mathematical framework of Section 3 covers these scenarios and many others in a unified way, although it remains to incorporate notions of one robot simulating another to create illusions (17).

Using concepts from Section 3.3, suppose that some humans act as an omniscient producer. The goal could be to get the robot into an I-state \bar{r} in which it reports that its tasks have been accomplished. This is achieved stage-by-stage by altering the physical world state x^p so that that targeted observations y^r are achieved for all sensors. The particular x^p chosen at each stage can depend on the sensed configuration of the receiver and even its I-states, if available. Note that if it not possible to spoof all of the sensors, then the challenge of maintaining plausibility increases.

How are the targeted receiver observations determined? One convenient way is to create a simulated world, which is then used to calculate what stimulus to provide to the sensors at each stage. This is a way of maintaining a coherent, plausible 'virtual' environment that responds to the receiver's actions and provides appropriate sensor feedback. This is called a *virtual world generator* or *VWG*, and is a crucial component for VR applied to humans and other biological organisms (18). The simulator becomes a useful tool for maximizing plausibility robustness, or even determining whether plausibility is even possible. If not, it

might attempt to minimize the forced-fusion magnitude. The computational complexity and ability of the simulator to respond in real time (within each stage) are important concerns.

In a robotics setting, we can even connect the 'brain' of the receiver directly to the simulator to evaluate planning, control, or sensor-fusion algorithms, as is done in software platforms such as Gazebo and CARLA. The real-time requirement can even be neglected. This would be the robot equivalent of a Gilbert Harman's 'brain in a vat' (see (18)); however, VR for robots enables using the actual robot sensors in a physical setting to provide a more accurate assessment. This could robot more thorough and systematic testing or verification of robot systems, and improve learning of models. It also becomes possible to *reverse engineer* robots systematically by observing how they respond to various virtual environments. This is quite analogous to the way neuroscientists and psychologists learn about the inner workings of biological organisms, to be discussed in Section 5.1.

For a more challenging scenario, suppose the producer itself is a robot, which is constrained by its own f^p and h^p . It might need to do motion planning to each an appropriate x^p to create the targeted observation y^r at the precise stage it is needed. It might have to choose trajectories that avoid detection by the receiver. Imagine the challenges of getting a producer mobile robot to trick a vacuum cleaning receiver as was done in Figure 4, attempting to hit its wall contact sensors in the right places at the right time! This results in a number of computational challenges, including deciding whether plausibility is possible, maximizing plausibility robustness, minimizing the forced-fusion magnitude, and determining the computational complexity of these problems in various settings.

5. APPLICATION TO HUMANS AND OTHER BIOLOGICAL AGENTS

5.1. Modeling Human Agents

Although modeling engineered agents is challenging enough, it is much harder to model humans for several reasons: 1) Robots and other engineered agents are designed and built, component by component, using well-tested principles of physics, whereas biological agents simply exist. They start as black boxes that are reverse engineered by scientists, including psychologists, neuroscientists, and biologists. Thus, there is much speculation and debate about 'what is going on' inside of the agent. 2) Data regarding internal operation can be easily collected during execution for engineered agents, but for humans we are limited to questionnaires and external biosensors such as those used in electroencephalography (EEG). 3) A major challenge for modeling humans is to maintain *ecological validity* in an experiment, but unfortunately, their behavior may be altered due to knowledge of participating in a study and other unnatural aspects of the experiment. 4) A robot can be simply rebooted for easy repeatability and to observe its reactions to varying environmental conditions; however, a human starts life only once and retains memories of prior trials. Thus, his I-state cannot be re-initialized. 5) Humans *adapt* at various levels and time scales. For example, eyes adjust to various lighting levels and people become more effective at using a computer mouse with practice. Some adaption can be built into engineered systems, but it can also be avoided wherever preferable. 6) Through *attention processes* humans have the ability to ignore many things while focusing on others, thereby complicating what seems to be known at a particular moment.

One implication of these differences is far less repeatability or determinism for the case of biological agents, at least from the experimenter's viewpoint. This has resulted in a preference for probabilistic models, and less ability to formulate an underlying deterministic

model. By contrast, engineered systems are typically modeled with a nominal deterministic part, and a stochastic part accounts for leftover disturbances.

For the worst-case, black-box extreme, imagine studying an impenetrable, mysterious gadget that was left behind by aliens. We would have only our knowledge of physics, chemistry, and so on, to poke and prod it, and observe the results. For modeling humans, we at least have useful models for human sensing and actuation. For sensing, a sensor mapping of the form $h: X \rightarrow Y$ might model the human sense organs to a high degree of accuracy based on decades of research in physiology and neuroscience. In this case, X should include the possible stimuli to be presented to the organ, and Y could be the resulting electrical impulses. Vision is the most sophisticated sense and is complicated by many factors such as eye movements, pupil adaptation, optical distortions, and photoreceptor properties such as density, mosaics, response times, wavelength sensitivities, and amplitude sensitivities. Furthermore, substantial neural processing occurs on the path from the retina through amacrine, horizontal, bipolar, and ganglion cells to the optic nerve. All of these complicate h , making it imperfect and more challenging to model than a digital camera. Other senses bring their own unique challenges: hearing, touch, thermoception, proprioception, pain, smell, taste, and vestibular. In the actuation direction, we seek an XTE f that yields the next external state X_{k+1} as a function of the current state X_k and a motor command u_k . This involves modeling human body kinematics and dynamics; it falls under the field of kinesiology and includes the characterization of motor skills and learning. Disturbance-based alternatives, such as $p(y_k | x_k)$ and $p(x_{k+1} | x_k; u_k)$, might be preferable.

Next imagine trying to extract a useful model of the human's I-space I and ITF \mathcal{I} . In the case of the alien gadget, its external state space X can be systemically altered while observations about any physical changes the gadget undergoes are made, including movement or emitting energy such as lights or sounds. A major assumption is that enough trials can be made with sufficient repeatability of behavior from the gadget, at least statistically. To instead model a human, suppose that we can leverage acceptable models of h and f . In this way, variations in X over the stages can be converted using h and f into hypothesized observations y and actions u . Thus, the brain appears to receive a history I-state. This is consistent with most models in neuroscience, including for example, Friston's *free energy principle (FEP)*, in which the external and internal states (called I-states here) are separated by so-called *Markov blankets* (19).

Whether under the FEP or other models, the human executes a policy $\pi: I \rightarrow U$; however, the problem is to characterize I . Setting $I = I_{hist}$ would make a brain with perfect memory and ability to make its motor commands contingent on distinct history I-states. Following Section 2.2, it is far more likely that an information mapping reduces I_{hist} to a sufficient, derived I-space I_{der} . For most models used in neuroscience and perceptual psychology, this would be I_{prob} from Section 2.2, and is consistent with the Bayesian brain hypothesis (20), predictive coding (21), and the FEP.

For a robot, the 'brain' is a computer, for which any part of its internal state can be easily monitored. By contrast, the human brain has around 86 billion neurons, with hundreds of millions more outside of the brain. The operation of each neuron through axons and dendrites is itself a complicated dynamical system. Direct measurement of neural activity is impossible, except in limited cases of single-unit recordings. Instead, scientists must resort to non-invasive measures such as EEG, magnetoencephalography (MEG), and functional magnetic resonance imaging (fMRI).

Fortunately, humans can also be asked questions. Thus, a common approach to mod-

eling is *psychophysics* (22), which aims to understand and quantify the relation between the I-states and the external world of physical stimuli by interactive questioning. Different psychophysical procedures, involving different types of tasks and settings can be used to target a certain aspect of the human visual system or, in general, a sensory system (23, 24). The human subject is asked to provide responses to questions based on provided stimuli. The most basic procedures, with yes/no response, are stimuli detection and discrimination. Discrimination is the ability of tell two stimuli apart, whereas for detection one of the two stimuli is the ‘null’ or ‘neutral’ stimulus (for example, average luminance when detecting contrast sensitivity).

As mentioned above, attention processes (25) are a major complication in the modeling of humans using VR. They clearly have the knowledge that they entered VR, but nevertheless respond as if it is real (26) or they are present (11, 1). Even Slater’s definition of a place illusion requires that the person knows he is someone else (1). This suggests that transitions in the ITF might vary based on attention. It is as if there are multiple agents or I-spaces within one, with transitions affected by attention, which can be modeled as part of a high-level policy.

Finally, note that VR itself is a useful methodology to improve models of humans because scientists can observe their responses to carefully controlled, interactive experiences that would be difficult or impossible to produce in normal environments (27).

5.2. From Classical Illusions to Virtual Reality

As mentioned in Section 1, artists have been creating perceptual illusions for millennia through paintings and sculpture. By leveraging technological developments, modern artists who develop illusory perceptual experiences include skilled magicians, photographers, cinematographers, graphic artists, video game designers, and VR developers. The term ‘illusion’ is used somewhat loosely in everyday life, whereas the reality relation R from Section 2.3 gave a precise definition. Thus, we can apply it to well-known illusions to clarify what kind of illusions they are, or whether they should even be considered as illusions.

Suppose that a drawing or picture is shown to a human subject, and we ask her what she perceives. The question could be constrained in a number of ways, such as providing multiple choices or asking whether one feature seems larger than another. This setup can be modeled using the stationary formulation from Section 3.2. The producer X-state $x^p = !$ places the picture in front of the human. The observation y^r models what is sensed by the receiver. The I-state i^r is simply the reported answer to the question, ideally corresponding to what is perceived. Using a reality relation R^r , we can determine whether various pairs $(i^r; !)$ are illusions.

Consider a line drawing of a rabbit, and i^r is perceiving a rabbit. It is an illusion if R^r is defined so that $!$ must correspond to a real rabbit being presented. If the drawing is intended to be a rabbit and R^r is defined accordingly, then it would *not* be an illusion to perceive a rabbit. We must also determine whether i^r corresponds to perceiving an actual rabbit or a drawing of a rabbit. In some cases, an illusion can be defined in a way that does not depend on the producer’s intended interpretation. Recall the so-called ‘illusions’ of Figure 1. For Figure 1(a), most people perceive the upper line segment to be longer than the lower one. As an illusion of a 3D scene, it would be a longer embedded object. However, as a line drawing, we can objectively measure the lengths of the two segments and conclude that they are the same length. Thus, perceiving the upper segment as longer is an illusion

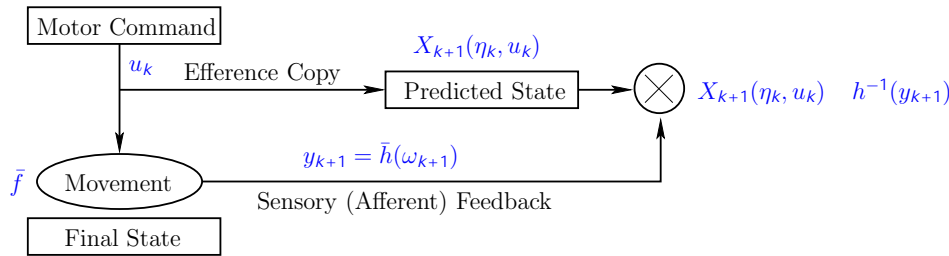


Figure 5

Sensorimotor contingency model from (30, 31), augmented with our nondeterministic I-state models.

in a measurable sense. Similarly, the A and B tiles in Figure 1(b) have identical RGB values (when viewed on a screen), but an illusion of tile A being darker persists. Figure 1(c) is different in that there is no objective ground truth regarding colors. Some people see black and blue whereas many others see white and gold (28). The reality relation can be defined in various ways based on what most people would interpret from the picture or even the original dress itself, but disagreement with what most other people would say hardly seems to be an illusion. Figure 1(d) is an example of *multistable perception*, in which for most people the I-state oscillates over time between a rabbit and a duck. Neither interpretation seems to be an illusion in the sense meant in this paper, unless one considers the fact that both are illusions because there is no real duck or rabbit present. Obviously, \mathcal{R}^r could be defined in various ways, and future questions remain about which definitions are most reasonable or usefully capture intuitions about what is meant by an illusion.

Multistable perception yields varying I-states over time, even though the stimulus is stationary. Now consider varying the stimulus by showing a sequence of pictures. Imagine gradually increasing the rate of pictures shown per second. Even at a few frames per second, we begin to perceive motion. This illusion is known as *stroboscopic apparent motion* (see (18)) and is the basis of video media. Another motion illusion is the *phi phenomenon* (29), in which blinking dots around a circle induce a sense of motion. The reality relation for these examples can be defined so that only true motion in Ω corresponds to reality; thus, such perceived motions are considered an illusion.

An obvious limitation of motion pictures is their lack of interactivity. This is overcome by video games, which take input from the user in the form of controllers and provide output in terms of video and audio displays. A kind of *virtual world* is usually maintained in the game, which could be considered illusory using the concepts above. Going a step further, VR creates a closed loop in which the person engages with a virtual world using more natural interaction mechanisms. Wearing an HMD allows her to move her eyes, head, and body while seeing and hearing what the virtual world provides, while easily forgetting that it is a display. Controlling the virtual world through natural methods such as hand movements and speech further increase the ‘realism’ of the perceptual experience. A *virtual world generator* maintains a consistent model as designed by the producer, and tracking systems estimate body configurations, to generate plausible responses (18).

A natural choice of explaining how VR works is the *sensorimotor contingency model* (30, 32, 1), as proposed in (31). Figure 5 reproduces the model in (31), but adds our nondeterministic agent models from Section 2.2. Suppose the receiver issues a motor command u_k^r . The efference copy and I-state $X_k^r(\bar{r}_k)$ used to calculate the ‘prediction’ $X_k^r(\bar{r}_k; u_k^r)$.

After movement, the next universe state is $I_{k+1} = \bar{F}(I_k; U_k^r; U_k^p)$, and the sensor yields $Y_{k+1} = \bar{H}(I_{k+1})$. If $X_k^r(I_k; U_k^r) \setminus (H^r)^{-1}(Y_{k+1}) \notin \emptyset$, then the I-state is plausible, as implied by the model relation M^r . The reality relation R^r could also be applied to determine whether a perceptual experience is illusory. This coarse model could be further detailed as multiple layers in a hierarchical predictive coding model, which has been generally successful in accounting for effects in visual perception (33, 34) and is theorized to explain the sensorimotor system as well (35, 36).

A probabilistic formulation could also be made. The prediction is $p(X_{k+1}^r | I_k; U_k^r)$, and instead of the preimage, the Bayesian posterior $p(X_{k+1}^r | Y_{k+1})$ is calculated from a given Y_{k+1} and ‘uninformative’ prior $p(X^r)$. The degree of (im)plausibility could be expressed in terms of the *Kullback-Leibler (KL) divergence*,

$$D_{KL}(p \parallel q) = \int_X \log \left(\frac{p(x)}{q(x)} \right) p(x) dx; \quad 11.$$

which represents an information-theoretic degree of ‘surprise’. In (11), let $p = p(X_{k+1}^r | Y_{k+1})$, which results from the observation, and $q = p(X_{k+1}^r | I_k; U_k^r)$ is the agent’s prediction. The model relation M^r could define a binary-valued plausibility by setting a threshold on the KL divergence. The KL divergence is also a critical component in the FEP.

In the design of VR systems, perception engineers would like to maintain or maximize plausibility while achieving a targeted perceptual experience. The producer’s ‘body’ corresponds to the engineered artifacts that may be distributed throughout the environment, including displays, controllers, tracking systems, and so on. Criteria to optimize include device expense, comfort, weight, development time, and adverse symptoms, such as fatigue or nausea. The better each human sense is understood, the easier it is to build a VR display for it by exploiting its limits. For example, a visual display need not be more than the maximum pixels-per-degree that are discernible by the human vision system. Furthermore, a greater understanding of human perception and cognition leads to more opportunities to exploit their limitations in VR systems. At an extreme, if I_{hist} is the brain model, then VR is as hard as possible because all histories lead to distinct perceptual experiences (I-state sequences). Collapsing I_{hist} to smaller I-spaces enables more changes to be made to sensing and actuation that go unnoticed by the human.

5.3. Evaluating Perceptual Experiences

Suppose that a user wears a VR HMD and has an engaging experience. How can we measure whether the producer has successfully delivered a targeted perceptual experience? What would be optimal, as discussed in (37)? This is difficult because the I-states are not directly accessible. A detailed account of this process is provided in Chapter 12 of (18). As mentioned earlier, questionnaires can be designed to have people describe what they experienced. Presence in a virtual environment is assessed using questionnaires, which are known to have unwanted biases (38, 39). Sickness symptoms have been assessed for decades using the *simulator sickness questionnaire (SSQ)* (40), which presents its own problems (41). Physiological measurements can also be taken, but they are somewhat more cumbersome for the users because they must wear sensors, and the experience must be strong enough to yield a detectable response (11).

Sickness is one of the most challenging aspects to measure, and amounts to a set of symptoms including fatigue, nausea, dizziness, headache, and eyestrain (18, 42). It seems hard to find an analogous problem in robotic systems, except perhaps that an overheating

CPU and power consumption may seem similar to fatigue. Section 3.2 introduced precise models of forced fusion (FFM) and plausibility robustness (PR), which could be used here to quantify the amount of mismatch or difficulty in maintaining plausibility. This is related to sensory conflict theory (43). Even sickness symptoms in some cases can be modeled as a sum over stages of costs proportional to the FMM. Mismatches that could be modeled mathematically within our framework include vergence-accommodation mismatch (44), visually induced motion sickness from vection (45), flicker fusion (46), and even the mysterious *uncanny valley* phenomenon (47).

It would be exciting if we could eventually develop new perceptual illusions that are directly predicted from perception engineering models. Many clever techniques have been developed through understanding the limitations of human vision, such as foveated rendering (48), frame skipping (49), and post-rendering image warp (50). Redirected walking exploits limitations in human navigation to convince people they are walking straight when in fact they move in circles (51). As examples continue to grow, how can approaches be more unified, with the systematic identification of many more?

5.4. Social Perception Engineering

Section 3.5 briefly described scenarios that involve multiple producers and/or receivers. For humans, this leads naturally into *social* perception engineering. At a very basic level, Shannon-Weaver communication can be modeled as a producer manipulating the environment (transmitting a message) that results in an intended I-state for the receiver (receiving the message). From broadcasts to social media, information propagates among people. Using reality relations for both producers and receivers, we can keep track of the spread of false information, a problem currently plaguing society as lies (a kind of illusion) are intentionally spread through a network of connected agents. Note that only one producer is needed to create a lie, and the rest of the network may propagate it without awareness that it is an illusion. Beyond the spread of messages, networked games (especially MMORPGs) and VR enable any number of people to interact through virtual worlds. This leads to *transformed social interaction* (52), in which people are able to have experiences that are different, and even better, than what could be accomplished in reality. For example, imagine in an educational setting, a teacher can appear to be looking at every student at the same time. People can design their own appearance or embodiment, so that biases based on physical characteristics such as race, gender, and height are readily overcome or studied by scientists under controlled conditions.

5.5. Other Biological Organisms

For organisms other than humans, VR is a rapidly maturing methodology for studying their behavior under controlled conditions (4). Since each organism has unique sensing and motor mechanisms, each VR system requires custom designs for displays and interaction methods. As stated in (4), there is an increasing need for engineers, especially roboticists and computer vision experts, to contribute to the development of VR systems for organisms, making VR for organisms an important part of perception engineering. Systems may be open-loop, such as showing visual stimuli to fish in an aquarium (53), or closed loop, as in (54), in which a hamster runs on a ball while being presented with visual stimuli on a curved projection screen. Scientists can learn about navigation, hunting, threat response, and many other behavioral aspects. Even social behavior has been studied, for example in

fish (55, 56), to help understand swarming or schooling. Their simpler neural structures and physiology often leads to models that have more detail and accuracy that is possible to obtain for humans. VR-based methods for animals also provide better insights into brains, bodies, and behavior of humans. Although questionnaires are not possible, more intrusive experimentation and measurement are allowed, such as conducting single-unit neural recordings. It has even been established that place and grid cells, which are critical for human navigation, respond to VR experiences (for example, (57)).

Consider the challenges of constructing models of their I-spaces, and the particular I-states during the organisms VR-induced perceptual experience. Starting from animals with the most human-like biology, studying macaque monkeys involves a VR setup that is unsurprisingly similar to that of humans (58), although single-unit recordings are at least possible. For smaller animals, from rodents (54) to fruit flies (56), building a fully interactive virtual world is more feasible by having them run on rotating balls, or fly on a tether, and view fixed immersive displays (like miniature CAVE systems (59)). Single-unit recordings are performed in these settings to measure particular I-states. Compared to 86 billion neurons for humans, fruit flies exhibit complex behaviors with only 100,000 neurons, resulting in a greater hope of fully understanding them. Zebrafish larvae also have comparably many neurons and their bodies are transparent, allowing direct observation of their complete neural states (60). Although they have been studied in VR (61), it is more challenging to construct closed-loop systems in comparison to fruit flies and rodents. VR has even been applied to roundworms (62), for which its 302 neurons have been fully mapped (63); however, measurement of I-states during execution remains challenging. In (64), a paramecium, which has no neurons, was manipulated into swimming along targeted trajectories by applying an electric field across the water; can the internal physical states of the paramecium be considered as I-states, resulting in a perceptual experience? Perhaps not, and this example comes close to a point of debate among philosophers (65). Imagine a similar case of putting an object in a tray and using a robot to tilt the tray so that the object slides into a targeted configuration due to gravity (66). It would be ludicrous to claim that the object had a targeted perceptual experience, or even an illusion; nevertheless, the I-space concepts from Section 2 are useful for designing ‘sensorless’ manipulation strategies.

6. CHALLENGES AND OPPORTUNITIES

We have argued that perception engineering is an emerging discipline and introduced a mathematical framework to help characterize its scope and core. In a general setting, we precisely characterized what it means for a producer to alter the environment to deliver an intended perceptual experience for a receiver, with the important conditions of it being plausible and illusory. It will take decades of work to bring this envisioned discipline to maturity by a growing community of people to expand the foundations while leveraging important principles from other branches of engineering, and the sciences that study biological organisms. This includes the need to train a generation of *perception engineers*. If successful, perception engineering as a discipline would have a profound impact on society. There are many exciting challenges and open problems ahead. To conclude this article, we highlight some of these as a call to action.

The mathematical formulation of Sections 2 and 3 should be expanded in several ways. For example, they should account for more characteristics of biological organisms, especially adaptation and attention. A continuous-time formulation using differential equations could

be developed, which would naturally increase the connections to control theory. One could define optimal control problems, such as linear quadratic regulation (LQR) of perceptual experiences and illusions, in both discrete and continuous time. One could adapt robust control techniques to targeted perceptual experiences. Dynamic game theory can be applied to characterize equilibria between producers and receivers. A control-theoretic model of psychophysics can be developed in which a feedback policy is used to guide actions that improve the agent modeling process. Asynchronous and event-based models of sensing and interaction can be made, as opposed to having common stages. The correspondence, model, and reality relations should be formulated in many more settings. Logics should be applied to describe more complicated perception engineering tasks. Finally, the universe space should be expanded to explicitly model how agents obtain information regarding I-states, actions, and observations.

Computational issues have barely been touched in this paper. Under what conditions does a producer even exist that can achieve a targeted perceptual experience? (This is analogous to computability in theoretical computer science.) If it can, then what is the computational complexity of achieving it? Can we find a minimal producer in some meaningful sense? If the producer is a robot, then its mechanical capabilities and limitations must be taken into account. If it is autonomous, then can planning methods be developed to achieve its goals, and with what complexity? What forms of dynamic programming, including reinforcement learning, would be effective? What can machine learning methods contribute to the development of better models, or what can machine learning gain from perception engineering? What computational architectures and systems best support the creation of targeted perceptual experiences? This itself is an emerging field of interest (67).

Other fields of engineering are expected to benefit as well, including computer engineering and systems, computer graphics, sensing and vision systems, optical science, displays, acoustic engineering, filtering, and control theory. Robotics should especially benefit because perception engineering can help to create better robots through improved modeling and learning techniques, and the development of robots that are more robust to spoofing attacks (attempts to create illusory perceptual experiences for robots).

The field of perception engineering should even contribute back to the sciences that study organisms, as is already the case for VR usage. We expect to have improved understanding, definitions, and classifications of illusions. We expect improved mathematical models of perception and cognition, which may be inspired by the considerable overlap between biological and engineered producers and receivers. Finally, the mathematical formulations and engineering examples of perception engineering can help shed light on, and benefit from, continuing debates in the philosophy and cognitive science, especially involving situated agency, enactivism, semantics, symbol grounding, and representations (68, 69, 70).

DISCLOSURE STATEMENT

The authors are not aware of any affiliations, memberships, funding, or financial holdings that might be perceived as affecting the objectivity of this review.

ACKNOWLEDGMENTS

This work was supported by the European Research Council (ERC AdG, ILLUSIVE: Foundations of Perception Engineering, 101020977), Academy of Finland (PERCEPT 322637,

CHiMP 342556, PIXIE 331822), and Business Finland (HUMOR 3656/31/2019).

References

1. Slater M, Banakou D, Beacco A, Gallego J, Macia-Varela F, Oliva R. 2022. A separate reality: an update on place illusion and plausibility in virtual reality. *Frontiers in Virtual Reality* 3:81
2. von Soemmerring ST. 1796. *Über das Organ der Seele*. Königsberg. With afterword by Immanuel Kant
3. Prince S. 2010. Through the looking glass: Philosophical toys and digital visual effects. *Projections* 4:19{40
4. Naik H, Bastien R, Navab N, Couzin ID. 2020. Animals in virtual environments. *IEEE Transactions on Visualization and Computer Graphics* 26(5):2073{2083
5. von Neumann J, Morgenstern O. 1944. *Theory of Games and Economic Behavior*. Princeton, NJ: Princeton University Press
6. LaValle SM. 2006. *Planning Algorithms*. Cambridge, U.K.: Cambridge University Press. Available at <http://planning.cs.uiuc.edu/>
7. Sakcak B, Weinstein V, LaValle SM. 2023. The limits of learning and planning: Minimal sufficient information transition systems. In *Algorithmic Foundations of Robotics, XV*, ed. SM LaValle, JM O’Kane, M Otte, D Sadigh, P Tokekar. Berlin: Springer-Verlag
8. LaValle SM. 2012. *Sensing and Filtering: A Fresh Perspective Based on Preimages and Information Spaces*, vol. 1:4 of *Foundations and Trends in Robotics Series*. Delft, The Netherlands: Now Publishers
9. Weinstein V, Sakcak B, LaValle SM. 2022. An enactivist-inspired mathematical model of cognition. *Frontiers in Neurobotics*
10. Thrun S, Burgard W, Fox D. 2005. *Probabilistic Robotics*. Cambridge, MA: MIT Press
11. Skarbez R, Brooks FP, Whitton MC. 2017. A survey of presence and related concepts. *ACM Computing Surveys (CSUR)* 50(6):1{39
12. Hillis JM, Ernst MO, Banks MS, Landy MS. 2002. Combining sensory information: mandatory fusion within, but not between, senses. *Science* 298(5098):1627{30
13. Fainekos GE, Girard A, Kress-Gazit H, Pappas GJ. 2009. Temporal logic motion planning for dynamic mobile robots. *Automatica* 45(2):343{352
14. Shin H, Kim D, Kwon Y, Kim Y. 2017. *Illusion and dazzle: Adversarial optical channel exploits against lidars for automotive applications*. In *International Conference on Cryptographic Hardware and Embedded Systems*, pp. 445{467. Springer
15. Suomalainen M, Nilles AQ, LaValle SM. 2020. *Virtual reality for robots*. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 11458{11465
16. Trippel T, Weisse O, Xu W, Honeyman P, Fu K. 2017. *WALNUT: Waging doubt on the integrity of MEMS accelerometers with acoustic injection attacks*. In *IEEE European Symposium on Security and Privacy*, pp. 3{18
17. Shell DA, O’Kane JM. 2020. *Reality as a simulation of reality: robot illusions, fundamental limits, and a physical demonstration*. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE
18. LaValle SM. 2023. *Virtual Reality*. Cambridge, U.K.: Cambridge University Press
19. Friston K, Kilner J, Harrison L. 2006. A free energy principle for the brain. *Journal of physiology-Paris* 100(1-3):70{87
20. Gregory RL. 1980. Perceptions as hypotheses. *Philosophical Transactions of the Royal Society of London. B, Biological Sciences* 290(1038):181{197
21. Rao RPN, Ballard DH. 1999. Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. *Nature neuroscience* 2(1):79{87
22. Fechner GT. 1888. *Elemente Der Psychophysik*, vol. 1. Breitkopf & Härtel
23. Kingdom FAA, Prins N. 2016. *Psychophysics: a practical introduction*. Academic Press

24. Treutwein B. 1995. Minireview: Adaptive psychophysical procedures. *Vision Research* 35(17):2503{2522
25. Beck DM, Kastner S. 2009. Top-down and bottom-up mechanisms in biasing competition in the human brain. *Vision Research* 49(10):1154{1165
26. Meehan M, Insko B, Whitton M, Jr FPB. 2002. Physiological measures of presence in stressful virtual environments. *ACM transactions on graphics* 21(3):645{652
27. Jeung S, Hilton C, Berg T, Gehrke L, Gramann K. 2022. Virtual reality for spatial navigation. In *Current Topics in Behavioral Neuroscience*. Springer
28. Lafer-Sousa R, Hermann KL, Conway BR. 2015. Striking individual differences in color perception uncovered by 'the dress' photograph. *Current Biology* 25(13):R545{R546
29. Wertheimer M. 1912. Experimentelle Studien über das Sehen von Bewegung (Experimental Studies on the Perception of Motion). *Zeitschrift für Psychologie* 61:161{265
30. Gallier J. 2000. *Curves and Surfaces in Geometric Modeling*. San Francisco, CA: Morgan Kaufmann
31. Gonzalez-Franco M, Lanier J. 2017. Model of illusions and virtual reality. *Frontiers in psychology* 8:1125
32. O'Regan J, Noë A. 2001. A sensorimotor account of vision and visual consciousness. *The Behavioral and brain sciences* 24:939{973
33. Clark A. 2013. Whatever next? predictive brains, situated agents, and the future of cognitive science. *Behavioral and brain sciences* 36(3):181{204
34. Walsh KS, McGovern DP, Clark A, O'Connell RG. 2020. Evaluating the neurophysiological evidence for predictive processing as a model of perception. *Annals New York Academy of Sciences* 1464(1):242{268
35. Bastos AM, Usrey WM, Adams RA, Mangun GR, Fries P, Friston KJ. 2012. Canonical microcircuits for predictive coding. *Neuron* 76(4):695{711
36. Shipp S, Adams RA, Friston KJ. 2013. Reactions on agranular architecture: predictive coding in the motor cortex. *Trends in neurosciences* 36(12):706{716
37. Skarbez R, Brooks FP, Whitton MC. 2020. Immersion and coherence: Research agenda and early results. *IEEE transactions on visualization and computer graphics* 27(10):3839{3850
38. Slater M. 2004. How colorful was your day? Why questionnaires cannot assess presence in virtual environments. *Presence* 13(4):484{493
39. Usuh M, Catena E, Arman S, Slater M. 2000. Using presence questionnaires in reality. *Presence* 9(5):497{503
40. Kennedy RS, Lane NE, Berbaum KS, Lilienthal MG. 1993. Simulator sickness questionnaire: An enhanced method for quantifying simulator sickness. *International Journal of Aviation Psychology* 3(3):203{220
41. Lawson BD. 2015. Motion sickness symptomatology and origins. In *Handbook of Virtual Environments, 2nd Edition*, ed. KS Hale, KM Stanney, pp. 531{600. Boca Raton, FL: CRC Press
42. Stanney K, Lawson BD, Rokers B, Dennison M, Fidopiastis C, et al. 2020. Identifying causes of and solutions for cybersickness in immersive technology: reformulation of a research and development agenda. *International Journal of Human-Computer Interaction* 36(19):1783{1803
43. Oman CM. 1990. Motion sickness: a synthesis and evaluation of the sensory conflict theory. *Canadian journal of physiology and pharmacology* 68(2):294{303
44. Ho man DM, Girshick AR, Akeley K, Banks MS. 2008. Vergence/accommodation conflicts hinder visual performance and cause visual fatigue. *Journal of vision* 8(3):33{33
45. Keshavarz B, Hecht H, Lawson BD. 2015. Visually induced motion sickness: Causes, characteristics, and countermeasures. In *Handbook of Virtual Environments, 2nd Edition*, ed. KS Hale, KM Stanney, pp. 647{698. Boca Raton, FL: CRC Press
46. Mankowska ND, Marcinkowska AB, Waskow M, Sharma RI, Kot J, Winklewski PJ. 2021. Critical flicker fusion frequency: a narrative review. *Medicina* 57(10):1096
47. Saygin AP, Chaminade T, Ishiguro H, Driver J, Frith C. 2012. The thing that should not be:

- predictive coding and the uncanny valley in perceiving human and humanoid robot actions. *Social cognitive and affective neuroscience* 7(4):413{422
48. Levoy M, Whitaker R. 1990. *Gaze-directed volume rendering*. In *Proceedings of the 1990 symposium on interactive 3d graphics*, pp. 217{223
 49. Denes G, Maruszczky K, Ash G, Mantiuk RK. 2019. Temporal resolution multiplexing: Exploiting the limitations of spatio-temporal vision for more efficient vr rendering. *IEEE Transactions on Visualization and Computer Graphics* 25(5):2072{2082
 50. Mark WR, McMillan L, Bishop G. 1997. *Post-rendering 3D warping*. In *Proceedings of the 1997 symposium on Interactive 3D graphics*, pp. 7{16
 51. Steinicke F, Bruder G, Jerald J, Frenz H, Lappe M. 2009. Estimation of detection thresholds for redirected walking techniques. *IEEE transactions on visualization and computer graphics* 16(1):17{27
 52. Bailenson JN, Beall AC, Loomis J, Blascovich J, Turk M. 2004. Transformed social interaction: Decoupling representation from behavior and form in collaborative virtual environments. *PRESENCE: Teleoperators and Virtual Environments* 13(4):428{441
 53. Rosa SO, Sovrano VA, Vallortigara G. 2014. What can fish brains tell us about visual perception? *Frontiers in Neural Circuits* 8:119
 54. Thurley K, Ayaz A. 2017. Virtual reality systems for rodents. *Current Zoology* 63(1):109{119
 55. Larsch J, Baier H. 2018. Biological motion as an innate perceptual mechanism driving social affiliation. *Curr Biol* 28(22):3523{3532
 56. Stowers JR, Hofbauer M, Bastien R, Griessner J, Higgins P, et al. 2017. Virtual reality for freely moving animals. *Nature Methods* 14(10):995{1002
 57. Aghajian ZM, Acharya L, Moore JJ, Cushman JD, Vuong C, Mehta MR. 2015. Impaired spatial selectivity and intact phase precession in two-dimensional virtual reality. *Nature Neuroscience* 18(1):121{128
 58. Sato N, Sakata H, Tanaka Y, Taira M. 2004. Navigation in virtual environment by the macaque monkey. *Behavioural Brain Research* 153(1):287{91
 59. Cruz-Neira C, Sandin DJ, DeFanti TA, Kenyon RV, Hart JC. 1992. The CAVE: Audio visual experience automatic virtual environment. *Communications of the ACM* 35(6):64{72
 60. Kunst M, Laurell E, Mokayes N, Kramer A, Kubo F, et al. 2019. A cellular-resolution atlas of the larval zebra fish brain. *Neuron* 103:21{38
 61. Trivedi CA, Bollmann JH. 2013. Visually driven chaining of elementary swim patterns into a goal-directed motor sequence: a virtual reality study of zebra fish prey capture. *Frontiers in Neural Circuits* 7:86
 62. Faumont S, Rondeau G, Thiele TR, Lawton KJ, McCormick KE, et al. 2011. An image-free opto-mechanical system for creating virtual environments and imaging neuronal activity in freely moving caenorhabditis elegans. *PLoS One* 6(9):e24666
 63. Cook SJ, Jarrell TA, Brittin CA, Wang Y, Bloniarz AE, et al. 2019. Whole-animal connectomes of both caenorhabditis elegans sexes. *Nature* 571(7763):63{71
 64. Fearing R. 1991. *Control of a micro-organism as a prototype micro-robot*. In *Proceedings 2nd Int. Symp. Micromachines and Human Sciences*
 65. Fodor JA. 1986. Why paramecia don't have mental representations. *Midwest Studies in Philosophy* 10:3{23
 66. Erdmann MA, Mason MT. 1988. An exploration of sensorless manipulation. *IEEE Transactions on Robotics & Automation* 4(4):369{379
 67. Huzaifa M, Desai R, Grayson S, Jiang X, Jing Y, et al. 2022. ILLIXR: An open testbed to enable extended reality systems research. *IEEE Micro* 42(4):97{106
 68. Gallagher S. 2017. *Enactivist Interventions: Rethinking the Mind*. Oxford University Press
 69. Hipolito I. 2022. Cognition without neural representation: Dynamics of a complex system. *Frontiers in Psychology* 12
 70. Hutto DD, Myin E. 2012. *Radicalizing enactivism: Basic minds without content*. MIT Press