

Adaptive Intelligence: leveraging insights from adaptive behavior in animals to build flexible AI systems

Mackenzie Weygandt Mathis ^{1,2}

¹École Polytechnique Fédérale de Lausanne (EPFL), Brain Mind Institute Geneva, Switzerland.

²Mathis Laboratory of Adaptive Intelligence. [✉mackenzie.mathis@epfl.ch](mailto:mackenzie.mathis@epfl.ch)

Biological intelligence is inherently adaptive — animals continually adjust their actions based on environmental feedback. However, creating adaptive artificial intelligence (AI) remains a major challenge. The next frontier is to go beyond traditional AI to develop “adaptive intelligence,” defined here as harnessing insights from biological intelligence to build agents that can learn online, generalize, and rapidly adapt to changes in their environment. Recent advances in neuroscience offer inspiration through studies that increasingly focus on how animals naturally learn and adapt their world models. In this Perspective, I will review the behavioral and neural foundations of adaptive biological intelligence, the parallel progress in AI, and explore brain-inspired approaches for building more adaptive algorithms.

Introduction

Our world is built on a series of predictions made in our mind. That is, to us, the world is a construct, created conceptually from a series of millions of predictions harbored in the neural code. When we walk on the beach, we accumulate photons with our eyes, sound waves with our ears, we feel the wind on our face, and the sand shifting under our feet, and all of this comes together to give us the perception of a warm, sunny day with ocean waves crashing on the shore. However, our brain sits in the dark confines of our skull and tens of billions of neurons communicate through electrical activity in a bath of chemicals. As incoming sensory information is processed in the brain at different times, the brain imposes a lag in order to compile it, thrusting our perception into the past – tens of precious milliseconds behind reality. Nevertheless, we need to act quickly, and our brains have evidently evolved solutions to overcome these delays.

Moreover, we need to adapt to environments marked by unpredictability, continuous change, and infrequent repetition of decision-making scenarios. The uncertainty, variability, and complexity of these settings demand adaptive strategies, crucial for learning about local changes – such as the shifting sand under our feet. Dating back decades, researchers have developed new frameworks to understand (decompile) how this global prediction system may work (1–4). A core tenet has been that in order to overcome sensory-processing delays we must build internal models of the world (also called world models). We use the models to make predictions about the sensory and motor consequences of our actions (5–8). But it remains largely unknown how the brain enables adaptive behavior, and how we can take inspiration from the brain to build more adaptive artificial intelligence (AI) is of growing interest (9–11). Specifically, AI here broadly refers to systems or machines that aim to mimic biological intelligence by performing tasks that typically require cognition, such as problem solving, decision making, learning, and language processing. Technically, it encompasses machine learning

(ML), deep learning, natural language processing (NLP), and computer vision. In this Perspective, I review progress in studying biological intelligence and use this to propose how we can enhance AI-based agentic systems.

Internal models for adaptive behavior

Dating at least as far back as Aristotle, there has been a philosophical and scientific question of how our brains build models of the world (14). How can we combine sensory information into perception? How long does this take? Are we always tens or even hundreds of milliseconds behind reality? How could our brains make predictions about the world in order to act faster? Why do we mostly notice when the world violates our predictions? – if the ocean were suddenly quiet, we would immediately worry that we had lost our hearing and not doubt that the ocean had completely changed. While this is familiar in everyday life, we need to turn this experience into a controllable task to study it in the laboratory.

What is the neural basis of this adaptive behavior? This historically has been difficult to answer because we are always under-sampling the neural data that underlies the behavior and often want repeated measures. The mouse brain has an estimated 70 million neurons, yet even the most state-of-the-art technology allows for recording up to 1 million individual neurons at a time (15) (and on the order of only 3Hz resolution), while most labs can record in the order of hundreds to thousands of neurons with two-photon microscopy or high-density arrays such as Neuropixels (16, 17). Thus, classically, to study how a neuron encodes a given stimulus the animal is repeatedly subjected to the stimulus and then we build encoding and decoding models (18) to gain insight into what the neural activity represents.

However, if you want to understand how the brain learns, there is an inherent problem with this repeated-trial approach. As Heraclitus famously quipped, “No man ever steps in the same river twice, for it’s not the same river and he’s not the same man”. Beyond the clear changes that can occur

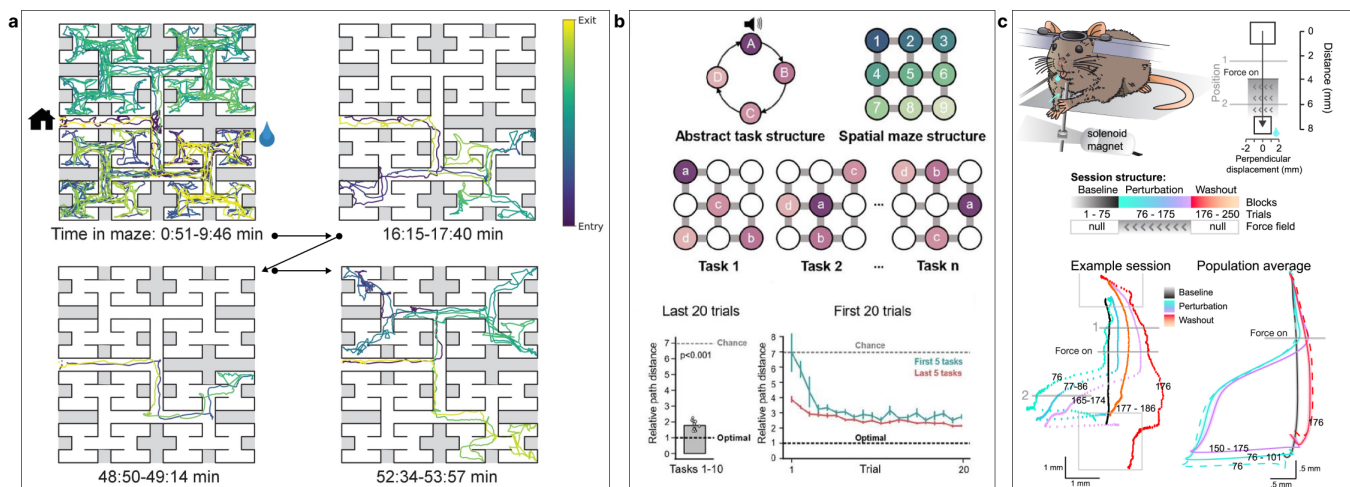


Figure 1. Rapid Learning in animals: from few-shot to updating of internal model-based learning. (a) Adapted from Rosenberg et al. 12: Four sample bouts from one mouse (B3) into the maze at various times during the experiment (time markings at bottom). The trajectory of the animal's nose is shown; time is encoded by the color of the trace. The entrance from the home cage and the water port are indicated in the top left panel. (b) Adapted from El-Gaby et al. 13: Task design: animals learned to navigate through 4 sequential goals on a 3x3 spatial grid-maze. Reward locations changed across tasks but the abstract structure, 4 rewards arranged in an ABCD loop, remained the same. Bottom left: When allowed to learn across multiple sessions, animals readily reached near-optimal performance in the last 20 trials, as demonstrated by comparing path length between goals to the shortest possible path. Bottom right: Performance improved across the initial 20 trials of each new task. This improvement was markedly more rapid for the last 5 tasks compared to the first 5 tasks. (c) Adapted from (7): Top Diagram of the joystick adaptation task structure: 75 baseline trials, 100 perturbation trials, and 75 "washout" (i.e., null field) trials. Bottom Left: an example session showing the average baseline trajectory (black), the first perturbation trial (sea green), average of first ten (green), and the last ten perturbation trials (purple), the first ten washout trials (orange) and the first washout trial (red). The numbers indicate trials. Right: average trajectories for all sessions ($n = 27$ sessions, from $n = 7$ mice) with temporally aligned averaging. Dashed sea green line is the average first perturbation, and the dashed red line is the average of the first five washout trials. Solid lines are average paw path \pm SEM in the direction of the perpendicular deviation, with the same color scheme as the left panel; solid green is the first 25 trials during the perturbation epoch, solid purple is the last 25 trials of the perturbation epoch, and solid red line is the average of the first 25 washout trials.

across trials during learning, the neural code can change even in steady-state given that the output behavior is rarely truly identical. Thus, averaging across trials likely misses key principles of neural computation. Since behavior is never exactly repeated anyhow, investigating learning — where behavior changes in predictable and measurable ways due to task design — may offer crucial insights into the neural basis of adaptive behavior. Namely, we can control the environmental factors that drive learning and examine the neural code across time. This not only lends itself to more ecologically-relevant tasks (the animals must be able to learn), but it's becoming a tractable approach given modern technological advances in large-scale neural recording and joint neural and behavioral analysis (18–21).

Adaptive animal and neural behavior

There has been a concurrent push for experimental designs that incorporates more naturalistic behaviors, more complex behaviors, and tasks that enable within-session learning. Aside from new large-scale neural recording technology, a part of this is due to the advances in markerless tracking of behavior (reviewed in Mathis and Mathis (22)), which has made digitizing movements much more tractable. What behavioral paradigms are being developed? While I cannot cover all progress in the field, I want to highlight several new approaches. Freely moving paradigms for rodents have rapidly advanced to encompass complex stimuli such as moving (threatening) novel objects (23), real-time VR worlds (24), labyrinth mazes (12), extended multi-area home

cages (25, 26), and 3D environments such as those used to study depth estimation (27). In several settings they have also built trial-based assays that smartly constrain these naturalistic settings in order to be able to directly probe learning across epochs.

An elegant thread of work across multiple groups aims to build tasks for zero- or few-shot learning in complex spatial worlds. Zero-shot learning is a term adapted from ML that describes the ability of a model to generalize to new tasks or concepts without any specific training examples (28). Few-shot means that it can learn from only a few examples. For example, Rosenberg et al. (12) developed a task where mice must seek rewards in a labyrinth (Figure 1a). They found that mice can not only make up to 2000 decisions per hour, but after unsupervised exploration of the maze they could few-shot learn to find the water spout (in around 10 tries) (12). Notably, this is a learning rate that is 1000-fold higher than two-alternative forced choice (2AFC) tasks that are typically deployed in rodent studies of decision-making, meaning this task allows one to better study rapid learning. Another example is showing that mice can zero-shot learn new sequences of rewards (such as go to location A, B, D, then C) (13). Mice were trained on multiple tasks with a shared underlying structure that organized a sequence of goals, though the specific goal locations varied. Strikingly, the mice learned this common structure, allowing them to make zero-shot inferences on the first trial of new tasks (Figure 1b). Intriguingly, they also found neurons in medial prefrontal cortex that acted as task-structured memory buffers, namely, they tracked the

progression along behaviorally-relevant steps (13). Such schemata could perhaps be mapped to memory-replay in neural networks (see below).

Animal intelligence is, of course, not limited to laboratory animals. Dating back to early comparative cognition studies (29) there has been a large body of work studying intelligence in a myriad of animals. Crows, for example, exhibit tool use and causal reasoning, with New Caledonian crows even manufacturing tools to retrieve food (30, 31). Bees demonstrate numeracy and complex communication through the waggle dance, encoding spatial information about food sources (32) and can solve “puzzles” such as pulling a string to move a food source closer (33). Various domesticated animals also exhibit human-aligned intelligence. Horses show social learning and can interpret human emotional cues (34). Dogs display theory of mind and word comprehension that rival that of young children (35). Such findings even sparked a reoccurring virtual animal AI Olympics challenge (36).

In sensorimotor control there has been a long line of work on rapid learning, which is called motor adaptation. For clarity, here I define learning as the acquisition of a skill (and thus learning a new internal model), while adaptation is using a learned internal model and updating it. Motor adaptation studies have been developed where visual or proprioceptive information is perturbed such that within-session, an animal must adapt. This has not only spurred the development of new neural analysis tools (18, 37–39), but a host of behavioral tasks. For example, visuomotor rotations have been used in humans and non-human primate studies to study how they can account for sensorimotor discrepancies in only a few hundred trials (37, 40–43). The same principle is used for changing environmental dynamics of a manipulandum that causes deviations in a limb-movement tasks (7, 44–46). Notably, these tasks are specifically designed to measure both adaptive learning, and the formation of internal models. There is always a period of baseline control movements, perturbations, and a return to baseline-condition that allows researchers to behaviorally measure if an internal model was updated (Figure 1c).

Neurally, evidence shows that motor and sensory areas can change their tuning properties during the course of learning (46–50). A study by Sun et al. (48) showed that during a motor adaptation (learning) task in macaques, motor cortex (M1) can “index” memories of the hand-force required for the process of learning to adapt in the form of movement readiness potentials that occur prior to movement execution. Notably, the neural subspace most predictive of hand forces changed during the period before movement (preparatory), specifically during the learning epoch. In a neural dimension orthogonal to this force-predictive subspace, they identified a uniform shift across all movement directions, including those unaffected by learning. Intriguingly this uniform shift remained after exposure to the force field, reflecting an updated internal model. In addition, other works also show new evidence of state-changes that persist across learning (8, 46), and temporally-resolved prediction errors (46, 51). Influ-

tial work on visual cortex (V1) has shown the existence of ‘mismatch’ neurons when the expected visual feedback is disrupted (51), and similar prediction errors have been found in somatosensory (S1), motor (M1), and frontal areas of cortex (46, 50).

Some of the best evidence for causally showing how neurons adapt has come from brain-machine-interface (BMI) studies that require the subject to directly alter neural firing in order to change something in the external world (like a cursor on a screen). Closed-loop systems, like calcium-based BMIs (caBMIs, BMIs driven by decoded optical activity readout of calcium fluorescence signals in M1) or with electrical stimulation, have been used to study the timescale and subsets of neurons that can be used, as not all neuronal types have been found to be equally adaptable. Fundamental work has revealed that even small numbers of neurons can be leveraged for decoding, and for BMI-guided feedback (37, 52). Notably, they often live in discrete subspaces of the neural dynamics. Vendrell-Llopis et al. (53) pushed this further to link these subspaces to cell-types. They trained mice to modulate the activity of either intratelencephalic (IT) or pyramidal tract (PT) neurons for reward. They found that mice learned to control PT neuron activity more quickly and effectively than IT neuron activity. This intriguingly could be related to the anatomical connectivity and differing inputs.

Anatomical links to computations

The anatomy of the brain is important to consider, as this too could have direct implications for the future design of adaptive AI systems. Several areas in the brain have this incredible structure, – layers in the cerebral (neo-) and cerebellar cortex – which could be a critical part of new architectures for AI systems. There is a long-standing debate about representation emergence from data vs. architecture constraints (59, 60), but it is increasingly clear that architecture is critically important and linked to function. Fundamental work dating back to the 1970’s is worth revisiting. Vernon Mountcastle describes a framework of the organizational principle of the neocortex (61). He delineates evidence that the major functional division in the neocortex is not whether an area is “sensory” or “motor,” but rather the vertical neocortical column constitutes the basic computational unit, and the input-output pattern merely dictates the space of information it acts upon—namely, the “auditory”, “visual” or “motor neocortex” has the same cellular scaffold and should only be considered a particular region based on the type of sensory input. Nonetheless, this architectural bias gives rise to computations on information with brain area-specific information.

Form ties to function, and one emerging hot topic is how cortical circuits implement learning from prediction errors across cell types and layers. Building on the initial discovery of sensory prediction errors in cortex (51), subsequent works have begun to record from subtypes of neurons such as parvalbumin (PV), vasoactive intestinal peptide (VIP) and somatostatin (SST) classes of interneurons, and excitatory neurons in order to form a more complete model of how both

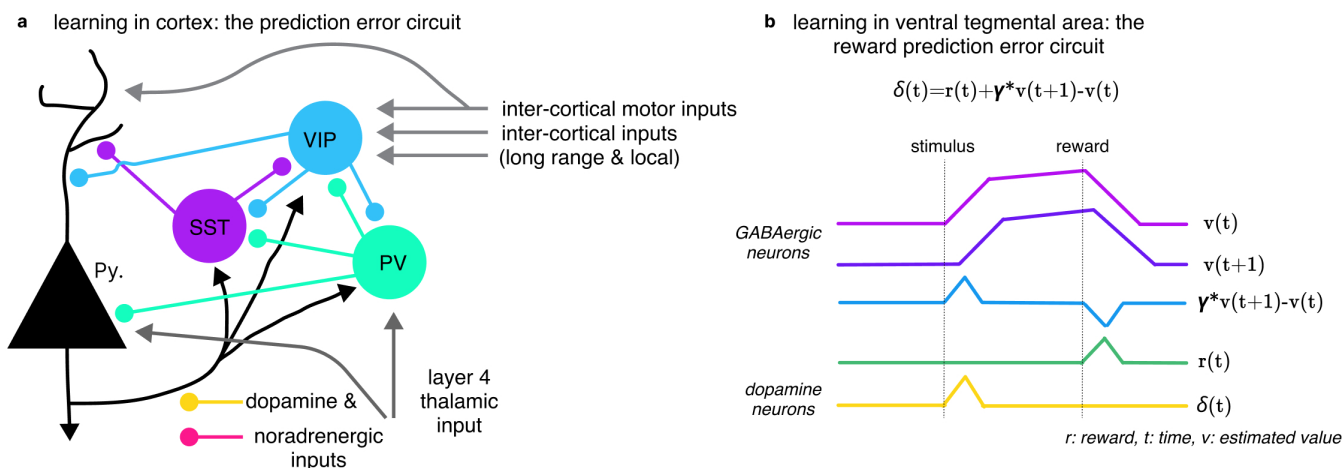


Figure 2. Neural computations: biological teaching signals. (a) Diagram of the theoretical prediction error circuit (based on evidence from Keller et al. (51), Leinweber et al. (54), Green et al. (55), Jordan and Keller (56)). (b) Diagram of the computation in the reward prediction error circuit of the midbrain (ventral tegmental area) dopamine and GABAergic neurons. Inspired from Gershman et al. (57), Cohen et al. (58).

excitatory and inhibitory neurons across layers could implement the learning rule (54, 55, 62) (Figure 2a). This work in cortex follows earlier work that discovered reward prediction errors (RPEs) in the midbrain (within the ventral tegmental area) (57, 58, 63), and found that they are also part of an inhibitory-gated circuit where GABAergic neurons play a crucial role in the computation of RPEs (64) (Figure 2b).

Recent work aims to unify how a hierarchical implementation of prediction errors (including RPEs) could aid in learning across sensory, cognitive, and motor systems. Tsai et al. (65) found that during a cue-guided reward task, layer 2/3 somatosensory neurons showed enhanced responses to reward-predictive stimuli, and this led to reduced reward-prediction errors and increased confidence in predictions. Following rule reversal, the lateral orbitofrontal cortex, via VIP interneurons, signaled context-prediction errors, reflecting a loss of confidence. The work suggests a hierarchical interaction of prediction errors across cortical regions, with top-down signals modulating sensory cortex activity. Notably, prior work finds limited roles for reward-only driven motor learning, but it can be combined with sensory prediction errors to shape performance (7, 42, 66). I envision that future work will more deeply test how various learning signals are used concurrently across systems. For example, tasks that have different types of prediction errors can be developed that might conflict, forcing the subject to weight different errors in order to guide future actions.

Neural development and inductive biases

In the context of learning and the links to AI, another critical angle is neural development. There are natural connections to make that could also inspire new inductive biases in the form of architectures or optimization algorithms in AI systems. For example, during development cells in the cortex migrate to their final resting location at different times (67). This process of spatial and temporal organization, I believe, has parallels in AI, particularly in how artificial neural networks undergo initialization and optimization processes to organize

their nodes across layers in deep learning models. Moreover, synapses, such as those at the neuro-muscular junction of alpha motor neurons and muscles, are formed after activity-dependent pruning (68). This mirrors the pruning techniques (so called knowledge distillation in ML (69)) used to improve model performance, where unnecessary connections are discarded during training to enhance efficiency.

There is also a complex and highly regulated developmental program in the spinal cord involving transcription factors and HOX genes (70). This genetic regulation lays out global (rostral-caudal) and even local (dorsal-ventral) connections. Similarly, the initial structure of a neural network is determined by a set of weights and biases, which are refined through training to create the connections necessary for learning. Additionally, the so-called genomic bottleneck plays a key role in determining early innate abilities of animals, from the immediate suckling reflex to the more time-intensive but impressive ability of a horse to stand and walk a few hours after birth (71). In AI, this is akin to pre-trained models or transfer learning, where networks leverage previously learned features to rapidly adapt to new tasks with minimal input.

Moreover, in the cortex, early traveling spontaneous waves shape the basis of early connections between cortical neurons. These waves are thought to synchronize neural activity across different regions of the brain. Intriguingly, new work has highlighted how the transformer architecture naturally gives rise to waves, which are a critical feature in the brain (72). This is reminiscent of how attention mechanisms in transformers enable dynamic, wave-like information flow across different parts of the model, helping the system maintain coherence across large datasets. Waves are seen not only in early development but also in steady-state, where they are thought to maintain global connectivity in the brain, encoding states like arousal and attention. This functional organization in the brain has direct analogs in how AI systems, like attention networks, maintain global contextual understanding across their layers.

Measuring changes in the neural code

New tools are critical for the study of adaptive behavior. Computational neuroscience has made large advances in the past decade, fueled by the need to better understand ensembles of neurons (20, 73–75) and by leveraging the concurrent innovations in deep learning, which is reviewed in Hurwitz et al. (76) and Mathis et al. (18). In brief, algorithms such as variational autoencoders, transformers, contrastive learning-optimized neural networks, and diffusion models are already making big strides in our ability to combine data across subjects, model multi-region interactions, and ‘decode the brain’ (18). I want to briefly highlight several tools that enable time-series analysis in single neurons and neural populations, as these are critical for adaptive learning studies, and I suspect will become more popular tools for studying AI systems.

One of the most impactful tools in the last 20 years has been the adoption of generalized linear models (GLMs) and generalized additive models (GAMs) for neuroscientific applications (77–79). Here, a model is trained using a feature matrix (can be one or many features) with the task of predicting the spikes of neurons given the set of features. While GLMs can come in many flavors, as an example, generally to handle variability in behavior models are trained and tested on different splits of the data (also for cross validation). Predictions from these test sets are combined, and a pseudo- R^2 is calculated. This does require that behavior is repeated, and also allows for flexibility; i.e., there is no direct trial-averaging or behavior-triggered averaging of the neural code.

For ensembles of neurons, several methods have been used or developed, ranging from PCA for measuring neural trajectories (80, 81), contrastive learning paired with consistency metrics to measure changes in the latent neural dynamics across learning (39, 46), to nonlinear dynamical systems for task- and task-irrelevant metrics (82, 83). These are more deeply covered in reviews elsewhere (18, 76), but I will note that several approaches can be used for single-trial dynamics and have been especially impactful in motor learning and adaptation studies in non-human primates (75, 84). Work by Azabou et al. (85) introduce a training framework and architecture designed to model population dynamics in large-scale neural recordings, which tokenizes individual spikes to capture fine temporal structures and constructs a latent representation of neural population activity. They trained this model on nearly 100 hours of data and show excellent cross-task generalizing in primate reaching studies. One limitation of this approach is the requirement of labels at test-time, but in practice this is nearly always available.

I predict we are just on the cusp of an influx of new advances in (neuro)science due to advances in AI. For example, large-language models (LLMs) are already starting to appear in behavioral neuroscience for digitizing and quantifying actions in video-based data (86). New vision-language, hierarchical, and multi-task foundation models are surely coming (87). For neural analysis, already several groups are actively working on building more unified encoders that can not only be

used for brain-machine interfaces, but they themselves could serve as models of the brain (18, 39, 85, 88–90). This is sure to give rise to new approaches to extract meaningful computational principles from neural dynamics, and these new learning rules have great potential to directly influence how we train, optimize, and deploy machine learning systems.

Training and learning in artificial systems

For understanding how neuroscience is poised to influence modern AI, here I provide some background on modern AI approaches. Today, most AI applications in production revolve around a train-test-deploy cycle, and therefore are inherently not adaptive. A dataset is curated, which can be lab-project scale (or in the case of generative pretrained transformers (GPTs), this would include nearly all the open-source content of the Internet), and then split into a train-and-test fraction. Most efforts benefit from leaving out-of-distribution (OOD) data in the test set for a realistic measure of how well the trained model generalizes. OOD is defined as data that is significantly different than the training data. If generalization is not needed – perhaps you are training a specific model for gait analysis in a specific laboratory setting – then the train/test split is often a random subset of the original dataset. But for many applications, training once and deploying a robust model is ideal.

How do we build a robust, generalizable model that can handle changes in the real-world? While the remarkable progress of GPTs in LLMs and the discovery of scaling laws has massively accelerated progress in AI, we aim for something better – more adaptive. What has been tried? This is where, at minimum, the sub-fields of continual learning (lifelong learning) and in-context learning come into view.

Current approaches to adapting models

Continual learning is the task of having a neural network learn a new series of tasks of the same modality, such as computer vision tasks, over time. It aims to address the challenge of developing models that can learn incrementally from a continuous stream of data without forgetting previously acquired knowledge. This is something that biological brains excel at – we typically don’t forget older information as we learn new skills. Yet, traditional machine learning models often face catastrophic forgetting (91) when retrained on new data, but new advances in continual learning have introduced techniques such as local module composition (92), knowledge distillation (93), elastic weight consolidation (EWC) (94, 95), synaptic intelligence, and memory-replay (96, 97) to mitigate these issues. For a more extensive discussions I point the readers to 98, 99.

Although not directly brain-inspired, Elastic Weight Consolidation (EWC) works by constraining certain parameters to reduce interference between tasks, effectively mitigating catastrophic forgetting (94). More recent work has adapted EWC to fine-tune self-supervised models, which has yielded performance improvements on tasks that have biased datasets

by maintaining knowledge of previous tasks better than traditional methods (95, 100). If we attempt to map this to neural dynamics, this might be related to the specialization of circuits. Namely, the “parameters” of certain areas could be fixed (such as primary receptors), whereas others could be more plastic (such as hippocampus and neocortex).

On the other hand, synaptic intelligence is a learning rule deeply inspired by brain plasticity (101). The authors developed intelligent synapses that adaptively store information, minimizing the forgetting of previously learned tasks while acquiring new ones. This approach mimics biological neural networks that balance plasticity and stability, enabling continual learning in artificial neural networks. Notably, this method can perform as well as EWC but can be performed online. Now, with modern scalable computing and architectures, this type of learning rule could be pivotal for building more adaptive systems.

Another brain-inspired machine learning algorithm for learning has been the development of memory-replay (102, 103), which mimics the memory systems in the hippocampus (Figure 3a). Notably, memory-replay also, perhaps indirectly, is built on the memory-inspired system of the Hopfield Network, which is a type recurrent neural network designed for content-addressable memory, resembling a spin glass system (104), whose invention also won John Hopfield a Nobel Prize in 2024. Memory-replay aims to have a memory bank of already learned actions (or images or tokens, etc.) that can be “replayed” (interweaved) into the training data in order to limit catastrophic forgetting (Figure 3b). This became a popular technique in reinforcement learning, and more recently, in computer vision and in LLMs. This training approach can also be used during active learning when new data (unlabeled, pseudo-labeled or labeled) can be used to fine-tune a model. If the input data continues to morph out-of-distribution (OOD) then using memory-replay can be leveraged to be sure the prior performance can be achieved.

Several memory-replay-based approaches in LLMs emerged after the release of ChatGPT in November 2022 (with its impressive language abilities yet limited content window, i.e., lack of memory), such as AmadeusGPT (86), MemGPT (105), and Voyager (106). AmadeusGPT introduced a short- and long-term memory where specific keywords could be quickly recalled in order to overcome token limits (i.e., the 4096 in GPT-3.5), while MemGPT used a vector database to construct a persistent memory to maintain context across interactions, and Voyager iteratively refined its skills by replaying feedback from past actions in games such as Minecraft.

Learning with spiking neural networks

Lastly, one of the most biologically inspired and grounded advances has been in the development of spiking neural networks (SNNs) in the 1990’s. Called the “third generation of neural networks” (107), SNNs are unique in their ability to model a time-dependent spike process, making them not

only biologically more plausible, but also highly energy efficient and can be used on edge-computing devices. Notably, SNNs naturally lend themselves to hardware accelerated devices and neuromorphic computing. Platforms such as Intel’s Loihi, IBM’s TrueNorth, and SpiNNaker provide hardware specifically designed to run SNNs efficiently (108). Moreover, several early work highlights their modularity and utility for multi-area modeling (109), and how multi-area SNNs could be used with semantic pointers for multi-area weighting of information (110). However, historically they have been difficult to train and scale, but that is rapidly changing and new methods to transfer information from ANNs to SNNs are emerging (111) and/or to directly compute the required gradients from SNNs (112).

How to build more adaptive AI

How can we take inspiration from adaptive behavior and the brain to build new AI approaches (11)? A core component of adaptive behavior is having already learned priors – these aforementioned internal models – that can be dynamically called upon. These priors clearly require data-at-scale to form – this “data” can be baked into neural circuits via the genome, or learned throughout life. An example of genetic priors would be the innate ability to suckle at birth across many species, or how a newborn horse stands and walks in a few hours after birth, while a non-genetic learned prior would be a motor skill such as playing a musical instrument. Another data source is the massive amount of unsupervised sensory and motor stimuli we rapidly accumulate across development and life. But how do we turn data into internal models?

While many questions remain about how these models are implemented in neural dynamics, there is increasing evidence of core computations such as the previously discussed prediction errors (51, 113) across brain regions, which shape and update internal models. Can we use these prediction errors as teaching signals in neural networks? Moreover, these prediction errors are anatomically defined, imposing architectural constraints on these systems that stem from brain-body dynamics. Can we leverage this feature in adaptive AI systems? Importantly, these prediction errors are canonical and act in specific sensory reference frameworks (such as vision or audition). As Mountcastle noted, it is less about “sensory” or “motor” spaces but rather a function of the data input/output, and by extension the shared architecture and computations. Thus, as I will argue below, we should take inspiration from building specialized nodes that have core computations in the right reference framework that use prediction errors to smartly update modules, much like the brain uses prediction error to update internal models.

Foundation Models

First, let us consider what is currently the top approach for building generalizable models whose aim is often to ingest several data modalities (such as text and video) and output either. The trend in machine learning is to build unified “Foun-

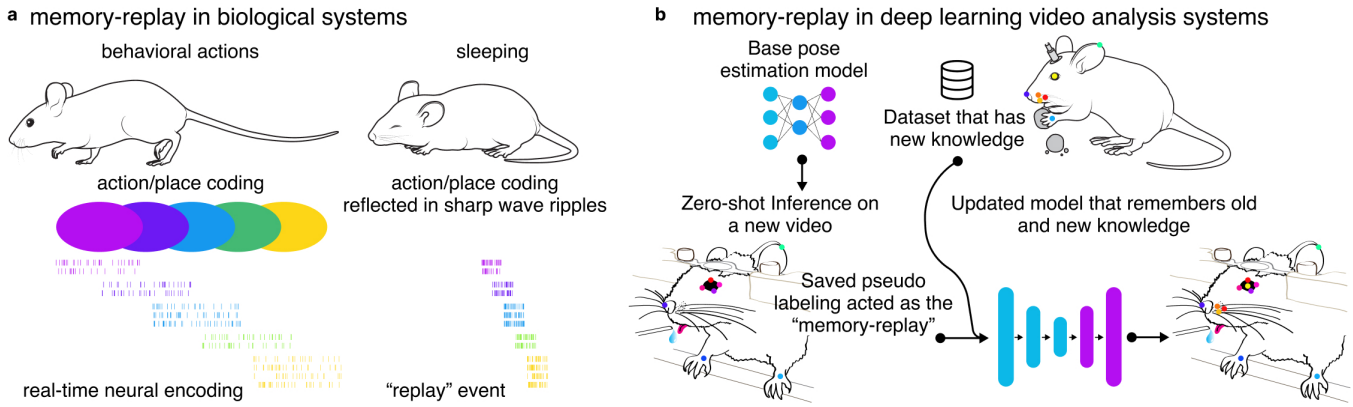


Figure 3. Memory-replay in biological and artificial systems. (a) In biological systems (such as in mice) neurons have been shown to encode place in their firing rates. These neurons can be active in the same sequence in sleep, which is thought to reflect memory-replay. (b) In deep learning systems one can build a replay buffer in continual learning tasks. For example, in video analysis for animal pose estimation one can use inferred labels with high confidence as pseudo-labels in order to retain base-model knowledge, even when (1) the base model's data might not be available to fine-tune, or (2) when a new dataset, perhaps with new labels, as depicted, needs to be learned without forgetting other labels. Mouse images adapted from scidraw.io.

dation Models” that generalize to unseen data or can be used for transfer learning in downstream tasks (114). The Foundation Model idea for natural language processing takes large-scale unlabeled data and learns a joint representation within a model (a transformer) using tokenization and self-supervised learning (114) (Figure 4a). The most infamous of this type is ChatGPT and other GPT models (such as the Gemini model series (115) or DeepSeek (116)). In neuroscience, several Foundation Models are also emerging for specific domains, such as for animal behavior analysis (97) and for modeling the visual systems of mice (90).

Other emerging and prominent approaches are generalizing this approach to multi-modal data (117, 118). One such example is Flamingo, a Visual Language Model (VLM) that integrates pretrained vision and language models, processing mixed visual-text data and handling images or videos as inputs (117) (Figure 4b). This could be extended to more modalities, as recently done for robotics (119). This is a rapidly advancing area of research that will surely influence neuroscience in the coming years.

While I believe this is an exciting path for models that show better generalization, I believe we can push further to build smartly adaptive, agentic systems. We should take inspiration from our brain’s specialization and build series of smartly interconnected specialist models that can then be adapted. Namely, just as all sensory/motor information gets converted to spikes – a shared computing space this is an efficient digitization of analog signals – notably, this information is computed upon within specific brain regions. Just as a cortical column in the primary visual cortex processes a particular region of the visual field and has specific coding properties related to visual features (like texture or edges), and the cortical circuit in the sensorimotor region computes in egocentric 3D kinematic, force, or even muscle space (46, 48).

Moreover, across the animal kingdom we see specialization in brains everywhere: from the remarkably high resolution of hawk vision, to bat echolocation, to the twelve cones of the mantis shrimp visual system. The neurons that are not sit-

ting in a single tangled mess of weights and nodes (akin to a single Foundation Model), but rather cleverly interconnected into neural circuits that have specific tuning to *collectively* solve many tasks. Perhaps we should take inspiration from these specialized systems versus aiming for a single model that could reach artificial superhuman intelligence. Thus, we should also focus on building smart, adaptive agentic systems of specialized models.

AI Agents

There is emerging work that I believe will make building AI agents (agentic systems) that are imbued with adaptive algorithms found in biologically intelligent systems possible. Agents that act as operating systems are one promising avenue. They can be a series of LLMs (which is currently being explored (86, 105)), but also a series of VLMs and other multi-modal language models (121), which work together to ingest data, select the appropriate task-specific encoder, and link these together with downstream processing scripts. Concretely, an example of this is video behavioral analysis, where an agentic system must select the best pose estimation model for the animal type in the video, perform semantic segmentation with another model, and then perform task-programming to create ethograms (86).

To develop more adaptive agentic systems, domain-specific encoder models are essential to achieve optimal performance. As in the example above, I propose an agentic system whose architecture relies on a series of encoder modules, yet adds two new aspects (Figure 4c). One, the agentic system that relies on a series of encoder modules, i.e., one for extracting poses, another for ingesting text plus images (VLMs), and one for extracting dynamics from neural spikes, is mediated by *prediction errors* from a decoder. Two, the outputs of each specialist encoder is *jointly optimized in a latent space*, which can then both learn joint representations that may be inherently more brain-like, but also a more powerful than a single Foundation Model without specialist nodes. This approach is domain-agnostic but draws on neuroscience to enhance adaptability.

How could this work? Encoders would be individually pre-trained (think domain specific Foundation Models, such as an encoder for animal behavior (97) or unified neural region-specific pretrained encoders (85)) and their outputs are jointly optimized into another transformer (Figure 4c). The issue becomes that while each encoder model starts as an expert, over time they need to learn more information, and the question is (1) how to sense that they need to undergo learning, and (2) how to adapt them. I propose that we could dynamically “lock” encoder models that show robustness on out-of-distribution samples and robustness to adversarial attacks; implying that they continually need monitored, by a dedicated LLM node, to provide trustworthy outputs. Therefore these encoders don’t need to be trained (adapted) until they are deemed no longer robust. This process itself I see as similar to the cortex-basal ganglia loop that is instrumental in learning and habit formation (122). Namely, these encoders could be in states of “skill learning” or “habituation” (frozen) and be constantly monitored for robustness (Figure 4c). When robustness levels fall below an acceptable threshold, then continual learning, memory-replay, pseudo-labeling or injection of new high-quality labeled data

could be added without taking the entire agentic system offline, just the encoder(s) that is not trustworthy. There has been a surge of exciting work on building OOD detection modules (123, 124) and explainable attribution methods in time-series and images (125–129) that could make this an attractive path forward.

We can also use a neuroscience-inspired approach by specifically adding prediction errors into the LLM decoder that are monitoring each encoder. Effectively, prediction errors provide a signed teaching signal when an uncertain prediction is detected. The development of predictive-coding inspired networks already show promise (130–132), and now with modern transformer architectures we could leverage prediction errors in order to do not only then do in-context learning within the encoder online, but signal which encoders need unlocked and updated – akin to how the brain must decide to online change a motor command vs. update an internal model. Here, I take inspiration from the motor system (133). A one-time error means “in-context learn” and change your ongoing motor command, but a repetitive error means update your internal model (Figure 1c). If the incoming perturbation

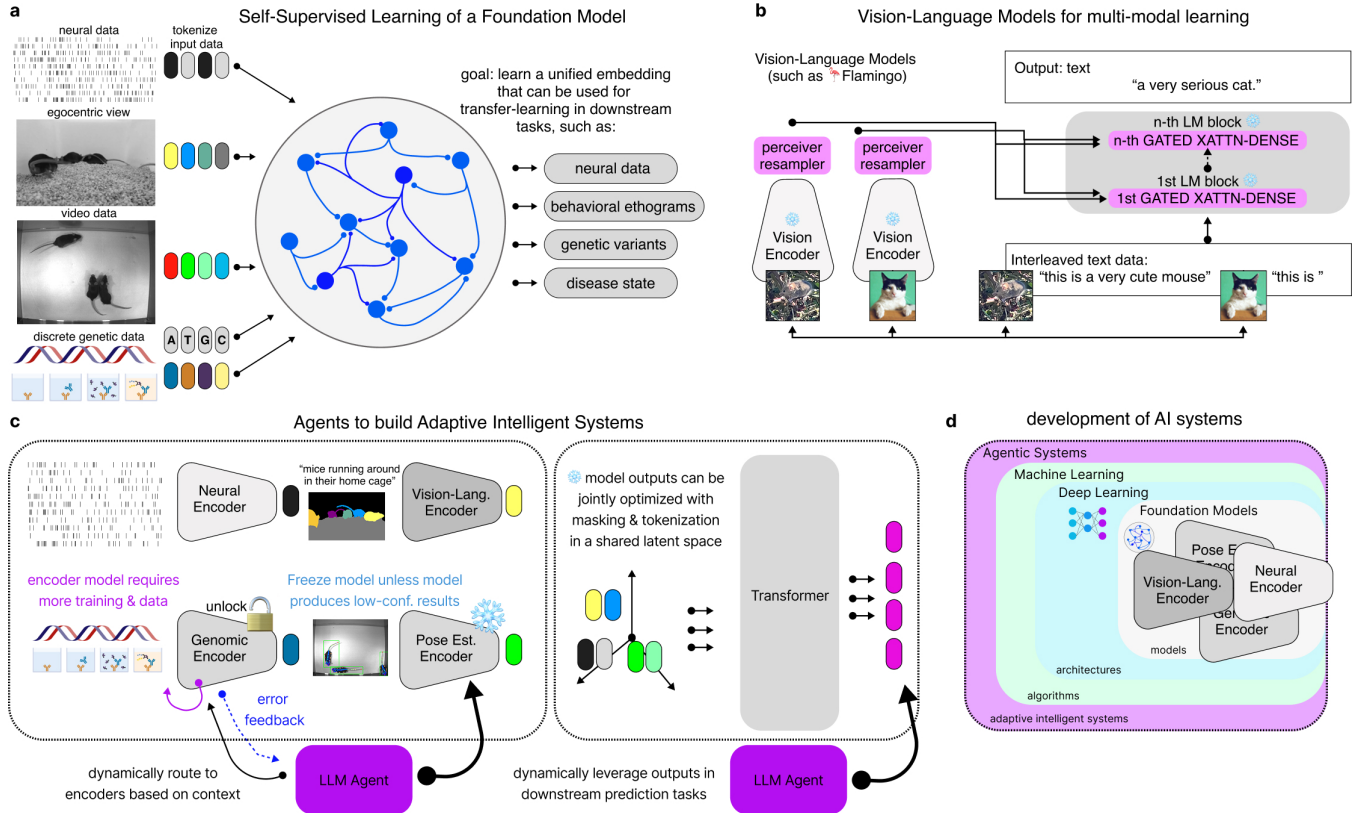


Figure 4. Foundation Models and Adaptive Agents. (a) One path to building generalizable models is to input diverse data, tokenize it, and train collectively with self-supervised learning. (b) Example of Vision-Language Models (VLMs) that interweave text and images. Inspired by, and the Cat image is from, Alayrac et al. (117). (c) My proposal is to leverage robust generalist foundation-model encoders that a domain specific information who’s outputs can be tokenized. For example, “pose estimation” can output skeletons of animals, “genomic” can be the DNA (or protein) sequence, “neural” could be tokenized spikes, and “vision-language” can be video captioning and actions (i.e., a higher level abstraction of behavior). In this example the assumption is there is a link between the genome, the animal phenotype, and the neural code. Critically though to my proposal is to then have an LLM-based agentic system oversee accuracy per encoder and pass forward only reliable outputs for joint latent-space optimization and joint decoder. In this way, the larger model need not be taken offline but rather be systematically updated when errors are detected in various encoders outputs. (d) Building on Bommasani et al. (114), agentic systems can leverage Foundation Models and advances in ML and deep learning in order to move beyond the emergence of functionality and into adaptive intelligent systems. Genetic and neural icons adapted from scidraw.io; mouse images are from the Allen Institute and Lauer et al. (120).

cannot be solved with online adaptation, a new internal model would then need to be created, which takes more time and can be completed “offline” in the unlocked model. Then methods such as continual learning, pseudo-labeling and memory-replay, or synaptic intelligence could be used to adapt the model online. Moreover, new discoveries in neuroscience should inspire new online learning techniques. The neural schema in prefrontal cortex discovered by El-Gaby et al. (13) could inspire new approaches for adaptive memory-replay and new model building as well (Figure 1b). Models taken offline would use standard techniques such as gradient descent to fine-tune the model on the data that produced a prediction error (and therefore a drop in robustness).

The second aspect of my proposal is to jointly optimize a downstream joint-latent space that takes as inputs the output tokens from each specialist encoder. This could provide an embedding that has a richer representation of each sub-task vs. jointly optimizing on raw input data as in Foundation Models. Importantly, this jointly optimized space can have its own LLM agent, with its own prediction errors, that also provides human natural language interactions such that at inference-time the model outputs human-interpretable reasoning, aiding in transparency, trustworthiness, and interpretability.

Lastly, independent of the agentic system, taking the cellular diversity found in the brain seriously in neural networks should drive innovations in architecture design. SNNs have both excitatory and inhibitory, and even neuro-modulatory units compared to current transformers. Can we merge the

power of large transformers and the scale of data animals receive, with the diverse cell types that produce neural computations (Figure 2)? Perhaps gating mechanisms akin to transformer blocks (Figure 4b) that add a signed (excitatory or inhibitory) tokenization and masking with biological time-delays can lead to new innovations in adaptive artificial systems.

Closing Remarks

As we strive to move beyond traditional AI towards building truly adaptive intelligence (Figure 4d), a key lies in the integration of insights from biological systems. I argue that adaptability will emerge from incorporating neuroscience-inspired principles, particularly internal models that leverage prediction error-based updating that refines internal representations based on discrepancies between expected and observed inputs, and from the brain-like modularity that structures the system with functionally distinct yet interconnected encoders, akin to sensory and motor modules that differ in data types but share an architecture. By integrating these mechanisms, agentic systems could achieve greater robustness and generalization across tasks, making them more effective in dynamic and uncertain environments. Inherent to these ideas is that these adaptive agentic systems can also be embodied into robotic hardware to systematically test, in fully observable sensorimotor system, their robustness and generalization. Ultimately, such advances therefore not only push progress in AI, but impact neuroscience by providing a new framework and testbed for theories of the mind.

References

1. D. M. Wolpert, Z. Ghahramani, and M. I. Jordan. An internal model for sensorimotor integration. *Science*, 269(5232):1880–1882, September 1995. ISSN 0036-8075.
2. Stephen H. Scott. Optimal feedback control and the neural basis of volitional motor control. *Nat Rev Neurosci*, 5(7):532–546, July 2004. ISSN 1471-003X. doi: 10.1038/nrn1427.
3. Emanuel Todorov and Michael I. Jordan. Optimal feedback control as a theory of motor coordination. *Nat Neurosci*, 5(11):1226–1235, November 2002. ISSN 1097-6256. doi: 10.1038/nrn963.
4. Karl J. Friston, Rosalyn J. Moran, Yukie Nagai, Tadairo Taniguchi, Hiroaki Gomi, and Josh Tenenbaum. World model learning and inference. *Neural networks : the official journal of the International Neural Network Society*, 144:573–590, 2021.
5. Rajesh P. N. Rao and Dana H. Ballard. Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. *Nature Neuroscience*, 2:79–87, 1999.
6. Rajesh P N Rao. A sensory-motor theory of the neocortex. *Nature neuroscience*, 2024.
7. Mackenzie Weygandt Mathis, Alexander Mathis, and Naoshige Uchida. Somatosensory Cortex Plays an Essential Role in Forelimb Motor Adaptation in Mice. *Neuron*, 93(6):1493–1503.e6, March 2017. ISSN 0896-6273. doi: 10.1016/j.neuron.2017.02.049.
8. Tomohiko Takei, Stephen G. Lomber, Douglas J. Cook, and Stephen H. Scott. Transient deactivation of dorsal premotor cortex or parietal area 5 impairs feedback control of the limb in macaques. *Current Biology*, 31:1476–1487.e5, 2021.
9. Robert J. Sternberg. A theory of adaptive intelligence and its relation to general intelligence. *Journal of Intelligence*, 7, 2019.
10. Stephen Grossberg. A path toward explainable ai and autonomous adaptive intelligence: Deep learning, adaptive resonance, and models of perception, emotion, and action. *Frontiers in Neuroinformatics*, 14, 2020.
11. Demis Hassabis, Dharmashan Kumaran, Christopher Summerfield, and Matthew M. Botvinick. Neuroscience-inspired artificial intelligence. *Neuron*, 95:245–258, 2017.
12. Matthew T. Rosenberg, Tony Zhang, Pietro Perona, and Markus Meister. Mice in a labyrinth show rapid learning, sudden insight, and efficient exploration. *eLife*, 10, 2021.
13. Mohamady El-Gaby, Adam Loyd Harris, James C. R. Whittington, William Dorrell, Arya Bhomick, Mark E. Walton, Thomas Akam, and Tim E. J. Behrens. A cellular basis for mapping behavioural structure. *Nature*, 2024. doi: 10.1101/2023.11.04.565609.
14. Mackenzie W. Mathis. The neocortical column as a universal template for perception and world-model learning. *Nature Reviews Neuroscience*, 24:3, 2022.
15. Jason Manley, Sihao Lu, Kevin Barber, Jeff Demas, Hyewon Kim, David Meyer, Francisca Martínez Traub, and Alipasha Vaziri. Simultaneous, cortex-wide dynamics of up to 1 million neurons reveal unbounded scaling of dimensionality with neuron number. *Neuron*, 112:1694–1709.e5, 2024.
16. Joshua H Siegle, Xiaoxuan Jia, Séverine Durand, Sam Gale, Corbett Bennett, Nile Graddis, Gregory Heller, Tamina K Ramirez, Hannah Choi, Jennifer A Luviano, et al. A survey of spiking activity reveals a functional hierarchy of mouse corticothalamic visual areas. *bioRxiv*, page 805010, 2019.
17. Ian H. Stevenson and Konrad Paul Kording. How advances in neural recording affect data analysis. *Nature Neuroscience*, 14:139–142, 2011.

18. Mackenzie W. Mathis, Adriana Perex Rotondo, Andreas Tolias, Edward Change, and Alexander Mathis. Decoding the brain: from neural representations to mechanistic models. *Cell*, 87, Issue 21:5814 – 5832, 2024.
19. Jérôme A. Lecoq, Natalia Orlova, and Benjamin F. Grewe. Wide. fast. deep: Recent advances in multiphoton microscopy of in vivo neuronal activity. *The Journal of Neuroscience*, 39:9042 – 9052, 2019.
20. Anne E. Urai, Brent Doiron, Andrew Michael Leifer, and Anne K. Churchland. Large-scale neural recordings call for new insights to link brain and behavior. *Nature Neuroscience*, 25:11–19, 2022.
21. Hongyu Chen and Ying Fang. Recent developments in implantable neural probe technologies. *MRS Bulletin*, pages 1–11, 2023.
22. Mackenzie Weygandt Mathis and Alexander Mathis. Deep learning tools for the measurement of animal behavior in neuroscience. *Current Opinion in Neurobiology*, 60:1–11, 2020.
23. Iku Tsutsui-Kimura, Naoshige Uchida, and Mitsuko Watabe-Uchida. Dynamical management of potential threats regulated by dopamine and direct- and indirect-pathway neurons in the tail of the striatum. *bioRxiv*, 2022. doi: 10.1101/2022.02.05.479267.
24. Gonçalo Lopes, Karolina Farrell, Edward A. B. Horrocks, Chi-Yu Lee, Mai M. Morimoto, Tomaso Muzzu, Amalia Papanikolaou, Fabio R. Rodrigues, Thomas Wheatcroft, Stefano Zucca, Samuel G. Solomon, and Aman B. Saleem. Creating and controlling visual environments using bonvision. *eLife*, 10, 2021.
25. Yair Shemesh, Asaf Benjamin, Keren Shoshani-Haye, Ofer Yizhar, and Alon Chen. Studying dominance and aggression requires ethologically relevant paradigms. *Current Opinion in Neurobiology*, 86, 2024.
26. Yaoyao Hao, Alyse M. Thomas, and Nuo Li. Fully autonomous mouse behavioral and optogenetic experiments in home-cage. *eLife*, 10, 2020.
27. Rolf J. Skyberg and Christopher M. Niell. Natural visual behavior and active sensing in the mouse. *Current Opinion in Neurobiology*, 86, 2024.
28. Mark Palatucci, Dean A. Pomerleau, Geoffrey E. Hinton, and Tom Michael Mitchell. Zero-shot learning with semantic output codes. In *Neural Information Processing Systems*, 2009.
29. Edward L. Thorndike. *Animal Intelligence: An Experimental Study of the Associative Processes in Animals*, volume 2. 1898. doi: 10.1037/h0092987.
30. Gavin Raymond Hunt. Manufacture and use of hook-tools by new caledonian crows. *Nature*, 379:249–251, 1996.
31. Christian Rutz and James J.H. St Clair. The evolutionary origins and ecological context of tool use in new caledonian crows. *Behavioural Processes*, 89(2):153–165, 2012. ISSN 0376-6357. doi: <https://doi.org/10.1016/j.beproc.2011.11.005>. Comparative cognition: Function and mechanism in lab and field.
32. F. C. Dyer and T. D. Seeley. Dance dialects and foraging range in three asian honey bee species. *Behavioral Ecology and Sociobiology*, 28: 227–233, 1991. doi: 10.1007/BF00175094.
33. S. Alem, C. J. Perry, X. Zhu, O. J. Loukola, T. Ingraham, E. Søvik, and L. Chittka. Associative mechanisms allow for social learning and cultural transmission of string pulling in an insect. *PLoS Biology*, 14(10):e1002564, 2016. doi: 10.1371/journal.pbio.1002564.
34. Leanne Proops, Kate Grounds, Amy Victoria Smith, and Karen McComb. Animals remember previous facial expressions that specific humans have exhibited. *Current Biology*, 28(9):1428–1432.e4, 2018. ISSN 0960-9822. doi: <https://doi.org/10.1016/j.cub.2018.03.035>.
35. Juliane Kaminski, Josep Call, and Julia Fischer. Word learning in a domestic dog: Evidence for "fast mapping". *Science*, 304(5677):1682–1683, 2004. doi: 10.1126/science.1097859.
36. Benjamin Beyret, José Hernández-Orallo, Lucy Cheke, Marta Halina, Murray Shanahan, and Matthew Crosby. The animal-ai environment: Training and testing animal-like artificial cognition. *arXiv*, 2019.
37. Alex H. Williams, Tony Hyun Kim, Forea Wang, Saurabh Vyas, Stephen I. Ryu, Krishna V. Shenoy, Mark J. Schnitzer, Tamara G. Kolda, and Surya Ganguli. Unsupervised discovery of demixed, low-dimensional neural dynamics across multiple timescales through tensor component analysis. *Neuron*, 98:1099–1115.e8, 2017.
38. Ben Sorscher, Surya Ganguli, and Haim Sompolsky. Neural representational geometry underlies few-shot concept learning. *Proceedings of the National Academy of Sciences of the United States of America*, 119, 2022.
39. Steffen Schneider, Jin Hwa Lee, and Mackenzie W. Mathis. Learnable latent embeddings for joint behavioural and neural analysis. *Nature*, 617: 360 – 368, 2023.
40. Samuel D. McDougale, Krista M. Bond, and Jordan A. Taylor. Explicit and Implicit Processes Constitute the Fast and Slow Processes of Sensorimotor Learning. *J. Neurosci.*, 35(26):9568–9579, July 2015. ISSN 0270-6474, 1529-2401. doi: 10.1523/JNEUROSCI.5061-14.2015.
41. John W Krakauer and Pietro Mazzoni. Human sensorimotor learning: adaptation, skill, and beyond. *Current Opinion in Neurobiology*, 21(4): 636–644, August 2011. ISSN 0959-4388. doi: 10.1016/j.conb.2011.06.012.
42. Jun Izawa and Reza Shadmehr. Learning from Sensory and Reward Prediction Errors during Motor Adaptation. *PLoS Comput Biol*, 7(3): e1002012, March 2011. doi: 10.1371/journal.pcbi.1002012.
43. Sergey D. Stavisky, Jonathan C. Kao, Stephen I. Ryu, and Krishna V. Shenoy. Trial-by-trial motor cortical correlates of a rapidly adapting visuomotor internal model. *The Journal of Neuroscience*, 37:1721 – 1732, 2017.
44. R. Shadmehr and F. A. Mussa-Ivaldi. Adaptive representation of dynamics during learning of a motor task. *J. Neurosci.*, 14(5):3208–3224, May 1994. ISSN 0270-6474, 1529-2401.
45. Dr E. Bizzi, N. Accornero, W. Chapple, and N. Hogan. Arm trajectory formation in monkeys. *Exp Brain Res*, 46(1):139–143, September 2013. ISSN 0014-4819, 1432-1106. doi: 10.1007/BF00238107.
46. Travis DeWolf, Steffen Schneider, Paul Soubiran, Adrian Roggenbach, and Mackenzie Weygandt Mathis. Neuro-musculoskeletal modeling reveals muscle-level neural dynamics of adaptive learning in sensorimotor cortex. *bioRxiv*, 2024. doi: 10.1101/2024.09.11.612513.
47. Chiang-Shan Ray Li, Camillo Padoa-Schioppa, and Emilio Bizzi. Neuronal Correlates of Motor Performance and Motor Learning in the Primary Motor Cortex of Monkeys Adapting to an External Force Field. *Neuron*, 30(2):593–607, May 2001. ISSN 0896-6273. doi: 10.1016/S0896-6273(01)00301-4.
48. Xulu Sun, Daniel J. O'Shea, Matthew D. Golub, Eric M. Trautmann, Saurabh Vyas, Stephen I. Ryu, and Krishna V. Shenoy. Cortical preparatory activity indexes learned motor memories. *Nature*, 602:274 – 279, 2022.
49. Travis Meyer and Nicole C. Rust. Single-exposure visual memory judgments are reflected in inferotemporal cortex. *eLife*, 7, 2018.
50. Nicolas Meirhaeghe, Hansom Sohn, and Mehrdad Jazayeri. A precise and adaptive neural mechanism for predictive temporal processing in the frontal cortex. *Neuron*, 109:2995–3011.e5, 2021.
51. Georg B. Keller, Tobias Bonhoeffer, and Mark Hübner. Sensorimotor Mismatch Signals in Primary Visual Cortex of the Behaving Mouse. *Neuron*, 74(5):809–815, June 2012. ISSN 0896-6273. doi: 10.1016/j.neuron.2012.03.040.
52. Patrick T. Sadtler, Kristin M. Quick, Matthew D. Golub, Steven M. Chase, Stephen I. Ryu, Elizabeth C. Tyler-Kabara, Byron M. Yu, and Aaron P. Batista. Neural constraints on learning. *Nature*, 512:423 – 426, 2014.
53. Nuria Vendrell-Llopis, Ching Fang, Albert J. Qü, Rui M. Costa, and Jose M. Carmena. Diverse operant control of different motor cortex populations during learning. *Current Biology*, 32:1616–1622.e5, 2021.
54. Marcus Leinweber, Daniel R. Ward, Jan M. Sobczak, Alexander Attinger, and Georg B. Keller. A sensorimotor circuit in mouse cortex for visual flow predictions. *Neuron*, 95(6):1420–1432.e5, 2017.
55. Jonathan Green, Carissa A Bruno, Lisa Traummüller, Jennifer Ding, Sinia Hrvatin, Daniel E Wilson, Thomas Khodadad, Jonathan Samuels, Michael Eldon Greenberg, and Christopher D. Harvey. A cell-type-specific error-correction signal in the posterior parietal cortex. *Nature*, 620:366 – 373, 2023.

56. Rebecca Jordan and Georg B. Keller. The locus coeruleus broadcasts prediction errors across the cortex to promote sensorimotor plasticity. *eLife*, 12, 2023.
57. Samuel J. Gershman, John A. Assad, Sandeep Robert Datta, Scott W. Linderman, Bernardo L. Sabatini, Naoshige Uchida, and Linda Wilbrecht. Explaining dopamine through prediction errors and beyond. *Nature neuroscience*, 2024.
58. Jeremiah Y Cohen, Sebastian Haesler, Linh Vong, Bradford B Lowell, and Naoshige Uchida. Neuron-type-specific signals for reward and punishment in the ventral tegmental area. *Nature*, 482(7383):85–88, February 2012. ISSN 1476-4687. doi: 10.1038/nature10754.
59. Brenden M. Lake, Tomer D. Ullman, Joshua B. Tenenbaum, and Samuel J. Gershman. Building machines that learn and think like people. *Behavioral and Brain Sciences*, 40:e253, 2017. doi: 10.1017/S0140525X16001837.
60. Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision transformers see like convolutional neural networks? *CoRR*, abs/2108.08810, 2021.
61. Vernon B. Mountcastle. The mindful brain: Cortical organization and the group-selective theory of higher brain function. *MIT Press*, 1978.
62. Katharina A Wilmes, Mihai A Petrovici, Shankar Sachidhanandam, and Walter Senn. Uncertainty-modulated prediction errors in cortical microcircuits. *eLife Sciences Publications, Ltd*, 2024. doi: 10.7554/eLife.95127.2.
63. W Schultz, P Dayan, and P R Montague. A neural substrate of prediction and reward. *Science*, 275(5306):1593–1599, March 1997. ISSN 0036-8075.
64. Neir Eshel, Michael Bukwich, Vinod Rao, Vivian Hemmelder, Ju Tian, and Naoshige Uchida. Arithmetic and local circuitry underlying dopamine prediction errors. *Nature*, 525:243 – 246, 2015.
65. Matthias C. Tsai, Jasper Teutsch, Willem A.M. Wybo, Fritjof Helmchen, Abhishek Banerjee, and Walter Senn. Hierarchy of prediction errors shapes the learning of context-dependent sensory representations. *bioRxiv*, 2024. doi: 10.1101/2024.09.30.615819.
66. Dimitrios J. Palidis, Heather R. McGregor, Andrew Vo, Penny A. MacDonald, and Paul L. Gribble. Null effects of levodopa on reward- and error-based motor adaptation, savings, and anterograde interference. *Journal of neurophysiology*, 2021.
67. Lynette Lim, Da Mi, Alfredo Llorca, and Oscar Marín. Development and functional diversification of cortical interneurons. *Neuron*, 100:294–313, 2018.
68. Haitao Wu, Wen-Cheng Xiong, and Lin Mei. To build a synapse: signaling pathways in neuromuscular junction assembly. *Development*, 137: 1017 – 1033, 2010.
69. Jianping Gou, B. Yu, Stephen J. Maybank, and Dacheng Tao. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129: 1789 – 1819, 2020.
70. Jeremy S. Dasen and Thomas M. Jessell. Hox networks and the origins of motor neuron diversity. *Current topics in developmental biology*, 88: 169–200, 2009.
71. Sergey A. Shuvaev, Divyansha Lachi, Alexei A. Koulakov, and Anthony M. Zador. Encoding innate ability through a genomic bottleneck. *Proceedings of the National Academy of Sciences of the United States of America*, 121, 2024.
72. Lyle E. Muller, Patricia S. Churchland, and Terrence J. Sejnowski. Transformers and cortical waves: encoders for pulling in context across time. *Trends in neurosciences*, 2024.
73. Juan Alvaro Gallego, Matthew G. Perich, Lee E. Miller, and Sara A. Solla. Neural manifolds for the control of movement. *Neuron*, 94:978–984, 2017.
74. Matthew G. Perich, Juan Alvaro Gallego, and Lee E. Miller. A neural population mechanism for rapid learning. *Neuron*, 100:964–976.e7, 2017.
75. Chethan Pandarinath, Daniel J. O’Shea, Jasmine Collins, Rafal Józefowicz, Sergey D. Stavisky, Jonathan C. Kao, Eric M. Trautmann, Matthew T. Kaufman, Stephen I. Ryu, Leigh R. Hochberg, Jaimie M. Henderson, Krishna V. Shenoy, L. F. Abbott, and David Sussillo. Inferring single-trial neural population dynamics using sequential auto-encoders. *Nature methods*, 15:805 – 815, 2017.
76. Cole Lincoln Hurwitz, N. N. Kudryashova, Arno Onken, and Matthias H Hennig. Building population models for large-scale neural recordings: Opportunities and pitfalls. *Current Opinion in Neurobiology*, 70:64–73, 2021.
77. Liam Paninski. Maximum likelihood estimation of cascade point-process neural encoding models. *Network: Computation in Neural Systems*, 15: 243 – 262, 2004.
78. Jonathan W. Pillow, Jonathon Shlens, Liam Paninski, Alexander Sher, Alan M. Litke, E. J. Chichilnisky, and Eero P. Simoncelli. Spatio-temporal correlations and visual signalling in a complete neuronal population. *Nature*, 454:995–999, 2008.
79. Edoardo Balzani, Kaushik J. Lakshminarasimhan, Dora E. Angelaki, and Cristina Savin. Efficient estimation of neural tuning during naturalistic behavior. In *Neural Information Processing Systems*, 2020.
80. Mehrdad Jazayeri and Arash Afraz. Navigating the neural space in search of the neural code. *Neuron*, 93:1003–1014, 2017.
81. Mark Churchland et al. Neural population dynamics during reaching. *Nature*, 2012.
82. Omid G. Sani, Bijan Pesaran, and Maryam Shanechi. Dissociative and prioritized modeling of behaviorally relevant neural dynamics using recurrent neural networks. *Nature neuroscience*, 2024.
83. Omid G. Sani, Hamidreza Abbaspourazad, Yan Tat Wong, Bijan Pesaran, and Maryam Modir Shanechi. Modeling behaviorally relevant neural dynamics enabled by preferential subspace identification. *Nature Neuroscience*, 24:140 – 149, 2020.
84. Mohammad Reza Keshtkaran, Andrew R. Sedler, Raaed H. Chowdhury, Raghav Tandon, Diya Basrai, Sarah L. Nguyen, Hansem Sohn, Mehrdad Jazayeri, Lee E. Miller, and Chethan Pandarinath. A large-scale neural network training framework for generalized estimation of single-trial population dynamics. *bioRxiv*, 2021.
85. Mehdi Azabou, Vinam Arora, Venkataramana Ganesh, Ximeng Mao, Santosh B Nachimuthu, Michael Jacob Mendelson, Blake Aaron Richards, Matthew G Perich, Guillaume Lajoie, and Eva L Dyer. A unified, scalable framework for neural population decoding. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
86. Shaokai Ye, Jessy Lauer, Mu Zhou, Alexander Mathis, and Mackenzie W. Mathis. Amadeusgpt: a natural language interface for interactive animal behavioral analysis. *Neural Information Processing Systems (NeurIPS)*, 2023.
87. Pablo Samuel Castro, Nenad Tomasev, Ankit Anand, Navodita Sharma, Rishika Mohanta, Aparna Dev, Kuba Perlin, Siddhant Jain, Kyle Levin, Noémi Éltető, Will Dabney, Alexander Novikov, Glenn C Turner, Maria K Eckstein, Nathaniel D Daw, Kevin J Miller, and Kimberly L Stachenfeld. Discovering symbolic cognitive models from human and animal behavior. *bioRxiv*, 2025. doi: 10.1101/2025.02.05.636732.
88. Yizi Zhang, Yanchen Wang, Donato Jimenez-Beneto, Zixuan Wang, Mehdi Azabou, Blake Richards, Olivier Winter, The International Brain Laboratory, Eva L. Dyer, Liam Paninski, and Cole Hurwitz. Towards a “universal translator” for neural dynamics at single-cell, single-spike resolution. *ArXiv*, 2024.
89. Yohann Benchetrit, Hubert J. Banville, and Jean-Rémi King. Brain decoding: toward real-time reconstruction of visual perception. *ArXiv*, abs/2310.19812, 2023.
90. Eric Y. Wang, Paul G. Fahey, Zhuokun Ding, Stelios Papadopoulos, Kayla Ponder, Marissa A. Weis, Andersen Chang, Taliah Muhammad, Saumil Patel, Zhiwei Ding, Dat Tran, Jiakun Fu, Casey M. Schneider-Mizell, R. Clay Reid, Forrest Collman, Nuno Maçarico da Costa, Katrin Franke, Alexander S. Ecker, Jacob Reimer, Xaq Pitkow, Fabian H. Sinz, and Andreas S. Tolia. Foundation model of neural activity predicts response to new stimulus types and anatomy. *bioRxiv*, 2024. doi: 10.1101/2023.03.21.533548.
91. Michael McCloskey and Neal J. Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. *Psychology of Learning and Motivation*, 24:109–165, 1989.
92. Oleksiy Ostapenko, Pau Rodriguez, Massimo Caccia, and Laurent Charlin. Continual learning via local module composition. In M. Ranzato,

- A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 30298–30312. Curran Associates, Inc., 2021.
93. Jathushan Rajasegaran, Munawar Hayat, Salman H Khan, Fahad Shahbaz Khan, and Ling Shao. Random path selection for continual learning. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
 94. James Kirkpatrick, Razvan Pascanu, Neil C. Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114:3521 – 3526, 2016.
 95. Andrius Ovsianas, Jason Ramapuram, Dan Busbridge, Eeshan Gunesh Dhekane, and Russ Webb. Elastic weight consolidation improves the robustness of self-supervised learning methods under transfer. *ArXiv*, abs/2210.16365, 2022.
 96. Liyuan Wang, Xingxing Zhang, Kuo Yang, Long Long Yu, Chongxuan Li, Lanqing Hong, Shifeng Zhang, Zhenguo Li, Yi Zhong, and Jun Zhu. Memory replay with data compression for continual learning. *ArXiv*, abs/2202.06592, 2022.
 97. Shaokai Ye, Anastasiia Filippova, Jessy Lauer, Maxime Vidal, Steffen Schneider, Tian Qiu, Alexander Mathis, and Mackenzie W. Mathis. Superanimal pretrained pose estimation models for behavioral analysis. *Nature Communications*, 15, 2024.
 98. Liyuan Wang, Xingxing Zhang, Qian Li, Mingtian Zhang, Hang Su, Jun Zhu, and Yi Zhong. Incorporating neuro-inspired adaptability for continual learning in artificial intelligence. *ArXiv*, abs/2308.14991, 2023.
 99. Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A comprehensive survey of continual learning: Theory, method and application. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(8):5362–5383, 2024. doi: 10.1109/TPAMI.2024.3367329.
 100. Cuong V. Nguyen, Yingzhen Li, Thang D. Bui, and Richard E. Turner. Variational continual learning. In *International Conference on Learning Representations*, 2018.
 101. Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. *Proceedings of machine learning research*, 70:3987–3995, 2017.
 102. Emma Roscow, Raymond Chua, Rui Ponte Costa, Matt W. Jones, and Nathan F. Lepora. Learning offline: memory replay in biological and artificial reinforcement learning. *Trends in Neurosciences*, 44:808–821, 2021.
 103. Longxin Lin. Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine Learning*, 8:293–321, 1992.
 104. John J. Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences of the United States of America*, 79 8:2554–8, 1982.
 105. Charles Packer, Sarah Wooders, Kevin Lin, Vivian Fang, Shishir G. Patil, Ion Stoica, and Joseph E. Gonzalez. MemGPT: Towards llms as operating systems. *arXiv preprint arXiv:2310.08560*, 2023.
 106. Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv: Arxiv-2305.16291*, 2023.
 107. Wolfgang Maass. Networks of spiking neurons: The third generation of neural network models. *Electron. Colloquium Comput. Complex.*, TR96, 1996.
 108. Kaushik Roy, Akhilesh R. Jaiswal, and Priyadarshini Panda. Towards spike-based machine intelligence with neuromorphic computing. *Nature*, 575:607 – 617, 2019.
 109. Chris Eliasmith et al. A large-scale model of the functioning brain. *Science*, 338(6111), 2012.
 110. Peter Blouw, Eugene Solodkin, Paul Thagard, and Chris Eliasmith. Concepts as semantic pointers: A framework and computational model. *Cognitive science*, 40 5:1128–62, 2016.
 111. Xiang-Yu He, Dongcheng Zhao, Yang Li, Guobin Shen, Qingqun Kong, and Yi Zeng. An efficient knowledge transfer strategy for spiking neural networks from static to event domain. In *AAAI Conference on Artificial Intelligence*, 2024.
 112. Timo C. Wunderlich and Christian Pehle. Event-based backpropagation can compute exact gradients for spiking neural networks. *Scientific Reports*, 11, 2021.
 113. Georg B. Keller and Philipp Sterzer. Predictive processing: A circuit approach to psychosis. *Annual review of neuroscience*, 2024.
 114. Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ B. Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen Creel, Jared Quincy Davis, Dorottya Demsky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kavin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah D. Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark S. Krass, Ranjay Krishna, Rohith Kudipudi, and et al. On the opportunities and risks of foundation models. *CoRR*, abs/2108.07258, 2021.
 115. Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillicrap, Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul R. Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, Jack Krawczyk, Cosmo Du, Ed Chi, Heng-Tze Cheng, Eric Ni, Purvi Shah, Patrick Kane, Betty Chan, Manal Faruqi, Aliaksei Severyn, Hanzhao Lin, YaGuang Li, Yong Cheng, Abe Ittycheriah, Mahdis Mahdieh, Mia Chen, Pei Sun, Dustin Tran, Sumit Bagri, Balaji Lakshminarayanan, Jeremiah Liu, Andras Orban, Fabian Gŭra, Hao Zhou, Xinying Song, Aurelien Boffy, Harish Ganapathy, Steven Zheng, HyunJeong Choe,  goston Weisz, Tao Zhu, Yifeng Lu, Siddharth Gopal, Jarrod Kahn, Maciej Kula, Jeff Pitman, Rushin Shah, Emanuel Taropa, Majd Al Merey, Martin Baeuml, Zhifeng Chen, Laurent El Shafey, Yujing Zhang, Olcan Sercinoglu, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs, Ana s White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakub Sygnowski, Alexandre Frechette, Charlotte Smith, Laura Culp, Lev Proleev, Yi Luan, Xi Chen, James Lottes, Nathan Schucher, Federico Lebron, Alban Rustemi, Natalie Clay, Phil Crone, Tomas Kocisky, Jeffrey Zhao, Bartek Perz, Dian Yu, Heidi Howard, Adam Bloniarz, Jack W. Rae, Han Lu, Laurent Sifre, Marcello Maggioni, Fred Alcober, Dan Garrett, Megan Barnes, Shantanu Thakoor, Jacob Austin, Gabriel Barth-Maron, William Wong, Rishabh Joshi, Rahma Chaabouni, Deeni Fatiha, Arun Ahuja, Gaurav Singh Tomar, Evan Senter, Martin Chadwick, Ilya Kornakov, Nithya Attaluri, I aki Iturrate, Ruibo Liu, Yunxuan Li, Sarah Cogan, Jeremy Chen, Chao Jia, Chenjie Gu, Qiao Zhang, Jordan Grimstad, Ale Jakse Hartman, Xavier Garcia, Thanumalayan Sankaranarayanan Pillai, Jacob Devlin, Michael Laskin, Diego de Las Casas, Dasha Valtter, Connie Tao, Lorenzo Blanco, Adri  Puigdom nech Badia, David Reitter, Mianna Chen, Jenny Brennan, Clara Rivera, Sergey Brin, Shariq Iqbal, Gabriela Surita, Jane Labanowski, Abhi Rao, Stephanie Winkler, Emilio Parisotto, Yiming Gu, Kate Olszewska, Ravi Addanki, Antoine Miech, Annie Louis, Denis Teplyashin, Geoff Brown, Elliot Catt, Jan Balaguer, Jackie Xiang, Pidong Wang, Zoe Ashwood, Anton Briukhov, Albert Webson, Sanjay Ganapathy, Smit Sanghavi, Ajay Kannan, Ming-Wei Chang, Axel Stjerngren, Josip Djolonga, Yuting Sun, Ankur Bapna, Matthew Aitchison, Pedram Pejman, Henryk Michalewski, Tianhe Yu, Cindy Wang, Juliette Love, Junwhan Ahn, Dawn Bloxwich, Kehang Han, Peter Humphreys, Thibault Sellam, James Bradbury, Varun Godbole, Sina Samangooei, Bogdan Damoc, Alex Kaskasoli, S bastien M. R. Arnold, Vijay Vasudevan, Shubham Agrawal, Jason Riesa, Dmitry Lepikhin, Richard Tanburn, Srivatsan Srinivasan, Hyeontaek Lim, Sarah Hodgkinson, Pranav Shyam, Johan Ferret, Steven Hand, Ankush Garg, Tom Le Paine, Jian Li, Yujia Li, Minh Giang, Alexander Neitz, Zaheer Abbas, Sarah York, Machel Reid, Elizabeth Cole, Aakanksha Chowdhery, Dipanjan Das, Dominika Rogo zi ska, Vitaliy Nikolaev, Pablo Sprechmann, Zachary Nado, Lukas Zilka, Flavien Prost, Luheng He, Marianne Monteiro, Gaurav Mishra, Chris Welty, Josh Newlan, Dawei Jia, Miltiadis Allamanis, Clara Huiyi Hu,

Raoul de Liedekerke, Justin Gilmer, Carl Saroufim, Shruti Rijhwani, Shaobo Hou, Disha Shrivastava, Anirudh Baddepudi, Alex Goldin, Adnan Ozturk, Albin Cassirer, Yunhan Xu, Daniel Sohn, Devendra Sachan, Reinald Kim Amplayo, Craig Swanson, Dessie Petrova, Shashi Narayan, Arthur Guez, Siddhartha Brahma, Jessica Landon, Mitayan Patel, Ruizhe Zhao, Kevin Villela, Luyu Wang, Wenhao Jia, Matthew Rahtz, Mai Giménez, Legg Yeung, James Keeling, Petko Georgiev, Diana Mincu, Boxi Wu, Salem Haykal, Rachel Saputro, Kiran Vodrahalli, James Qin, Zeynep Cankara, Abhanshu Sharma, Nick Fernando, Will Hawkins, Behnam Neyshabur, Solomon Kim, Adrian Hutter, Priyanka Agrawal, Alex Castro-Ros, George van den Driessche, Tao Wang, Fan Yang, Shuo yiin Chang, Paul Komarek, Ross McIlroy, Mario Lučić, Guodong Zhang, Wael Farhan, Michael Sharman, Paul Natsev, Paul Michel, Yamini Bansal, Siyuan Qiao, Kris Cao, Siamak Shakeri, Christina Butterfield, Justin Chung, Paul Kishan Rubenstein, Shivani Agrawal, Arthur Mensch, Kedar Soparkar, Karel Lenc, Timothy Chung, Aedan Pope, Loren Maggiore, Jackie Kay, Priya Jhakra, Shibo Wang, Joshua Maynez, Mary Phuong, Taylor Tobin, Andrea Tacchetti, Maja Trebacz, Kevin Robinson, Yash Katariya, Sebastian Riedel, Paige Bailey, Kefan Xiao, Nimesh Ghelani, Lora Aroyo, Ambrose Slone, Neil Houlsby, Xuehan Xiong, Zhen Yang, Elena Gribovskaya, Jonas Adler, Mateo Wirth, Lisa Lee, Music Li, Thais Kagohara, Jay Pavagadhi, Sophie Bridgers, Anna Bortsova, Sanjay Ghemawat, Zafarali Ahmed, Tianqi Liu, Richard Powell, Vijay Bolina, Mariko Iinuma, Polina Zablotskaia, James Besley, Da-Woon Chung, Timothy Dozat, Ramona Comanescu, Xiance Si, Jeremy Greer, Guolong Su, Martin Polacek, Raphaël Lopez Kaufman, Simon Tokumine, Hexiang Hu, Elena Buchatskaya, Yingjie Miao, Mohamed Elhawaty, Aditya Siddhant, Nenad Tomasev, Jinwei Xing, Christina Greer, Helen Miller, Shereen Ashraf, Aurko Roy, Zizhao Zhang, Ada Ma, Angelos Filos, Milos Besta, Rory Blevins, Ted Klimentko, Chih-Kuan Yeh, Soravit Changpinyo, Jiaqi Mu, Oscar Chang, Mantas Pajarskas, Carrie Muir, Vered Cohen, Charline Le Lan, Krishna Haridasan, Amit Marathe, Steven Hansen, Sholto Douglas, Rajkumar Samuel, Mingqiu Wang, Sophia Austin, Chang Lan, Jiepu Jiang, Justin Chiu, Jaime Alonso Lorenzo, Lars Lowe Sjösund, Sébastien Cevey, Zach Gleicher, Thi Avrahami, Anudhyan Boral, Hansa Srinivasan, Vittorio Selo, Rhys May, Konstantinos Aisopos, Léonard Hussenot, Livio Baldini Soares, Kate Baumli, Michael B. Chang, Adrià Recasens, Ben Caine, Alexander Pritzel, Filip Pavetic, Fabio Pardo, Anita Gergely, Justin Frye, Vinay Ramasesh, Dan Horgan, Kartikeya Badola, Nora Kassner, Subhrajit Roy, Ethan Dyer, Víctor Campos Campos, Alex Tomala, Yunhao Tang, Dalia El Badawy, Elspeth White, Basil Mustafa, Oran Lang, Abhishek Jindal, Sharad Vikram, Zhitao Gong, Sergi Caelles, Ross Hemsley, Gregory Thornton, Fangxiaoyu Feng, Wojciech Stokowiec, Ce Zheng, Phoebe Thacker, Çağlar Ünlü, Zhishuai Zhang, Mohammad Saleh, James Svensson, Max Bileschi, Piyush Patil, Ankesh Anand, Roman Ring, Katerina Tsihla, Arpi Vezet, Fabio Selvi, Toby Shevlane, Mikel Rodriguez, Tom Kwiatkowski, Samira Daruki, Keran Rong, Allan Dafoe, Nicholas FitzGerald, Keren Gu-Lemberg, Mina Khan, Lisa Anne Hendricks, Marie Pellat, Vladimir Feinberg, James Cobon-Kerr, Tara Sainath, Maribeth Rauh, Sayed Hadi Hashemi, Richard Ives, Yana Hasson, Eric Noland, Yuan Cao, Nathan Byrd, Le Hou, Qingze Wang, Thibault Sottiaux, Michela Paganini, Jean-Baptiste Lespiau, Alexandre Mouferek, Samer Hassan, Kaushik Shivakumar, Joost van Amersfoort, Amol Mandhane, Pratik Joshi, Anirudh Goyal, Matthew Tung, Andrew Brock, Hannah Sheahan, Vedant Misra, Cheng Li, Nemanja Rakićević, Mostafa Dehghani, Fangyu Liu, Sid Mittal, Junhyuk Oh, Seb Noury, Eren Sezener, Fantine Huot, Matthew Lamm, Nicola De Cao, Charlie Chen, Sidharth Mudgal, Romina Stella, Kevin Brooks, Gautam Vasudevan, Chenxi Liu, Mainak Chain, Nivedita Melinker, Aaron Cohen, Venus Wang, Kristie Seymore, Sergey Zubkov, Rahul Goel, Summer Yue, Sai Krishnakumaran, Brian Albert, Nate Hurley, Motoki Sano, Anhad Mohananey, Jonah Joughin, Egor Filonov, Tomasz Kepa, Yomna Eldawy, Jiawern Lim, Rahul Rishi, Shirin Badiezadegan, Taylor Bos, Jerry Chang, Sanil Jain, Sri Gayatri Sundara Padmanabhan, Subha Puttagunta, Kalpesh Krishna, Leslie Baker, Norbert Kalb, Vamsi Bedapudi, Adam Kurzrok, Shuntong Lei, Anthony Yu, Oren Litvin, Xiang Zhou, Zhichun Wu, Sam Sobell, Andrea Siciliano, Alan Papir, Robby Neale, Jonas Bragagnolo, Tej Toor, Tina Chen, Valentin Anklin, Feiran Wang, Richie Feng, Milad Gholami, Kevin Ling, Lijuan Liu, Jules Walter, Hamid Moghaddam, Arun Kishore, Jakub Adamek, Tyler Mercado, Jonathan Mallinson, Siddhinita Wandekar, Stephen Cagle, Eran Ofek, Guillermo Garrido, Clemens Lombriser, Maksim Mukha, Botu Sun, Hafeezul Rahman Mohammad, Josip Matak, Yadi Qian, Vikas Peswani, Pawel Janus, Quan Yuan, Leif Schelin, Oana David, Ankur Garg, Yifan He, Oleksii Duzhyi, Anton Ålgmyr, Timothée Lottaz, Qi Li, Vikas Yadav, Luyao Xu, Alex Chinien, Rakesh Shivanna, Aleksandr Chuklin, Josie Li, Carrie Spadine, Travis Wolfe, Kareem Mohamed, Subhabrata Das, Zihang Dai, Kyle He, Daniel von Dinkelge, Shyam Upadhyay, Akanksha Maurya, Luyan Chi, Sébastien Krause, Khalid Salama, Pam G Rabinovitch, Pavan Kumar Reddy M, Aarush Selvan, Mikhail Dektiarev, Golnaz Ghiasi, Erdem Guven, Himanshu Gupta, Boyi Liu, Deepak Sharma, Idan Heimlich Shtacher, Shachi Paul, Oscar Akerlund, François-Xavier Aubet, Terry Huang, Chen Zhu, Eric Zhu, Elco Teixeira, Matthew Fritze, Francesco Bertolini, Liana-Eleonora Marinescu, Martin Bølle, Dominik Paulus, Khyatti Gupta, Tejas Latkar, Max Chang, Jason Sanders, Roopa Wilson, Xuwei Wu, Yi-Xuan Tan, Lam Nguyen Thiet, Tulsee Doshi, Sid Lall, Swaroop Mishra, Wanming Chen, Thang Luong, Seth Benjamin, Jasmine Lee, Ewa Andrejczuk, Dominik Rabiej, Vipul Ranjan, Krzysztof Styrz, Pengcheng Yin, Jon Simon, Malcolm Rose Harriott, Mudit Bansal, Alexei Robsky, Geoff Bacon, David Greene, Daniil Mirylenka, Chen Zhou, Obaid Sarvana, Abhimanyu Goyal, Samuel Andermatt, Patrick Siegler, Ben Horn, Assaf Israel, Francesco Pongetti, Chih-Wei "Louis" Chen, Marco Selvatici, Pedro Silva, Kathie Wang, Jackson Tolins, Kelvin Guu, Roey Yogeve, Xiaochen Cai, Alessandro Agostini, Maulik Shah, Hung Nguyen, Noah Ó Donnale, Sébastien Pereira, Linda Friso, Adam Stambler, Adam Kurzrok, Chenkai Kuang, Yan Romanikhin, Mark Geller, ZJ Yan, Kane Jang, Cheng-Chun Lee, Wojciech Fica, Eric Malmi, Qijun Tan, Dan Banica, Daniel Balle, Ryan Pham, Yanping Huang, Diana Avram, Hongzhi Shi, Jasjit Singh, Chris Hidey, Niharika Ahuja, Pranab Saxena, Dan Dooley, Srividya Pranavi Potharaju, Eileen O'Neill, Anand Gokulchandran, Ryan Foley, Kai Zhao, Mike Dusenberry, Yuan Liu, Pulkit Mehta, Ragha Kotikalapudi, Chaleance Safranek-Shrader, Andrew Goodman, Joshua Kessinger, Eran Globen, Prateek Kolhar, Chris Gorgolewski, Ali Ibrahim, Yang Song, Ali Eichenbaum, Thomas Brovelli, Sahitya Potluri, Preethi Lahoti, Cip Baetu, Ali Ghorbani, Charles Chen, Andy Crawford, Shalini Pal, Mukund Sridhar, Petru Gurita, Asier Mujika, Igor Petrovski, Pierre-Louis Cedoz, Chenmei Li, Shiyuan Chen, Niccolò Dal Santo, Siddharth Goyal, Jitesh Punjabi, Karthik Kappagananthu, Chester Kwak, Pallavi LV, Sarmishta Velury, Himadri Choudhury, Jamie Hall, Premal Shah, Ricardo Figueira, Matt Thomas, Minjie Lu, Ting Zhou, Chintu Kumar, Thomas Jurdi, Sharat Chikkerur, Yenai Ma, Adams Yu, Soo Kwak, Victor Åhdel, Sujeevan Rajayogam, Travis Choma, Fei Liu, Aditya Barua, Colin Ji, Ji Ho Park, Vincent Hellendoorn, Alex Bailey, Taylan Bilal, Huanjie Zhou, Mehrdad Khatir, Charles Sutton, Wojciech Rządowski, Fiona Macintosh, Konstantin Shagin, Paul Medina, Chen Liang, Jinjing Zhou, Pararth Shah, Yingying Bi, Attila Dankovics, Shipra Banga, Sabine Lehmann, Marissa Bredesen, Zifan Lin, John Eric Hoffmann, Jonathan Lai, Raynald Chung, Kai Yang, Nihal Balani, Arthur Bražinskis, Andrei Sozanschi, Matthew Hayes, Héctor Fernández Alcalde, Peter Makarov, Will Chen, Antonio Stella, Liselotte Snijders, Michael Mandl, Ante Kärman, Paweł Nowak, Xinyi Wu, Alex Dyck, Krishnan Vaidyanathan, Raghavender R, Jessica Mallet, Mitch Rudominer, Eric Johnston, Sushil Mittal, Akhil Udathu, Janara Christensen, Vishal Verma, Zach Irving, Andreas Santucci, Gamaleldin Elsayed, Elnaz Davoodi, Marin Georgiev, Ian Tenney, Nan Hua, Geoffrey Cideron, Edouard Leurent, Mahmoud Alnahlawi, Ionut Georgescu, Nan Wei, Ivy Zheng, Dylan Scandinaro, Heinrich Jiang, Jasper Snoek, Mukund Sundararajan, Xuezhi Wang, Zack Ontiveros, Itay Karo, Jeremy Cole, Vinu Rajashekhar, Lara Tumeh, Eyal Ben-David, Rishub Jain, Jonathan Uesato, Romina Datta, Oskar Bunyan, Shimu Wu, John Zhang, Piotr Stanczyk, Ye Zhang, David Steiner, Subhajt Naskar, Michael Azzam, Matthew Johnson, Adam Paszke, Chung-Cheng Chiu, Jaume Sanchez Elias, Afroz Mohiuddin, Faizan Muhammad, Jin Miao, Andrew Lee, Nino Vieillard, Jane Park, Jiageng Zhang, Jeff Stanway, Drew Garmon, Abhijit Karmarkar, Zhe Dong, Jong Lee, Aviral Kumar, Luowei Zhou, Jonathan Evens, William Isaac, Geoffrey Irving, Edward Loper, Michael Fink, Isha Arkatkar, Nanxin Chen, Izhak Shafran, Ivan Petrychenko, Zhe Chen, Johnson Jia, Anselm Levskaya, Zhenkai Zhu, Peter Grabowski, Yu Mao, Alberto Magni, Kaisheng Yao, Javier Snider, Norman Casagrande, Evan Palmer, Paul Suganthan, Alfonso Castaño, Irene Giannoumis, Wooyeol Kim, Mikołaj Rybiński, Ashwin Sreevatsa, Jennifer Prendki, David Soergel, Adrian Goedeckemeyer, Willi Gierke, Mohsen Jafari, Meenu Gaba, Jeremy Wiesner, Diana Gage Wright, Yawen Wei, Harsha Vashisht, Yana Kulizhskaya, Jay Hoover, Maigo Le, Lu Li, Chimezie Iwuanyanwu, Lu Liu, Kevin Ramirez, Andrey Khorlin, Albert Cui, Tian LIN, Marcus Wu, Ricardo Aguilar, Keith Pallo, Abhishek Chakladar, Ginger Perng, Elena Allica Abellan, Mingyang Zhang, Ishita Dasgupta, Nate Kushman, Ivo Penchev, Alena Repina, Xihui Wu, Tom van der Weide, Priya Ponnappalli, Caroline Kaplan, Jiri Simsa, Shuangfeng Li, Olivier Dousse, Fan Yang, Jeff Piper, Nathan Le, Rama Pasumarthi, Nathan Lintz, Anitha Vijayakumar, Daniel Andor, Pedro Valenzuela, Minnie Lui, Cosmin Paduraru, Daiyi Peng, Katherine Lee, Shuyuan Zhang, Somer Greene, Duc Dung Nguyen, Paula Kurylowicz, Cassidy Hardin, Lucas Dixon, Lili Janzer, Kiam Choo, Ziqiang Feng,

- Biao Zhang, Achintya Singhal, Dayou Du, Dan McKinnon, Natasha Antropova, Tolga Bolukbasi, Orgad Keller, David Reid, Daniel Finchelstein, Maria Abi Raad, Remi Crocker, Peter Hawkins, Robert Dadashi, Colin Gaffney, Ken Franko, Anna Bulanova, Rémi Leblond, Shirley Chung, Harry Askham, Luis C. Cobo, Kelvin Xu, Felix Fischer, Jun Xu, Christina Sorokin, Chris Alberti, Chu-Cheng Lin, Colin Evans, Alek Dimitriev, Hannah Forbes, Dylan Banarse, Zora Tung, Mark Omernick, Colton Bishop, Rachel Sterneck, Rohan Jain, Jiawei Xia, Ehsan Amid, Francesco Piccinno, Xingyu Wang, Praseem Banzal, Daniel J. Mankowitz, Alex Polozov, Victoria Krakovna, Sasha Brown, MohammadHossein Bateni, Dennis Duan, Vlad Firoiu, Meghana Thotakuri, Tom Natan, Matthieu Geist, Ser tan Girgin, Hui Li, Jiayu Ye, Ofir Roval, Reiko Tojo, Michael Kwong, James Lee-Thorp, Christopher Yew, Danila Sinopalnikov, Sabela Ramos, John Mellor, Abhishek Sharma, Kathy Wu, David Miller, Nicolas Sonnerat, Denis Vnukov, Rory Greig, Jennifer Beattie, Emily Caveness, Libin Bai, Julian Eisenschlos, Alex Korchemniy, Tomy Tsai, Mimi Jasarevic, Weize Kong, Phuong Dao, Zeyu Zheng, Frederick Liu, Fan Yang, Rui Zhu, Tian Huey Teh, Jason Sanmiya, Evgeny Gladchenko, Nejc Trdin, Daniel Toyama, Evan Rosen, Sasan Tavakkol, Linting Xue, Chen Elkind, Oliver Woodman, John Carpenter, George Papamakarios, Rupert Kemp, Sushant Kafle, Tanya Grunina, Rishika Sinha, Alice Talbert, Diane Wu, Denese Owusu-Afriyie, Cosmo Du, Chloe Thornton, Jordi Pont-Tuset, Pradyumna Narayana, Jing Li, Saaber Fatehi, John Wieting, Omar Ajmeri, Benigno Uria, Yeongil Ko, Laura Knight, Amélie Héliou, Ning Niu, Shane Gu, Chenxi Pang, Yeqing Li, Nir Levine, Ariel Stolovich, Rebeca Santamaria-Fernandez, Sonam Goenka, Wenny Yustalim, Robin Strudel, Ali Elqursh, Charlie Deck, Hyo Lee, Zonglin Li, Kyle Levin, Raphael Hoffmann, Dan Holtmann-Rice, Olivier Bachem, Shou Arora, Christy Koh, Soheil Hassas Yeganeh, Siim Põder, Mukarram Tariq, Yanhua Sun, Lucian Ionita, Mojtaba Seyedhosseini, Pouya Tafti, Zhiyu Liu, Anmol Gulati, Jasmine Liu, Xinyu Ye, Bart Chrzasczcz, Lily Wang, Nikhil Sethi, Tianrun Li, Ben Brown, Shreya Singh, Wei Fan, Aaron Parisi, Joe Stanton, Vinod Koverkathu, Christopher A. Choquette-Choo, Yunjie Li, TJ Lu, Abe Ittycheriah, Prakash Shroff, Mani Varadarajan, Sanaz Bahargam, Rob Willoughby, David Gaddy, Guillaume Desjardins, Marco Cornero, Brona Robenek, Bhavishya Mittal, Ben Albrecht, Ashish Shenoy, Fedor Moiseev, Henrik Jacobsson, Alireza Ghaffarkhah, Morgane Rivièrè, Alanna Walton, Clément Crepey, Alicia Parrish, Zongheng Zhou, Clement Farabet, Carey Radebaugh, Praveen Srinivasan, Claudia van der Salm, Andreas Fjeldand, Salvatore Scellato, Eri Latorre-Chimoto, Hanna Klimczak-Plucińska, David Bridson, Dario de Cesare, Tom Hudson, Piermaria Mendolicchio, Lexi Walker, Alex Morris, Matthew Mauer, Alexey Guseynov, Alison Reid, Seth Odom, Lucia Loher, Victor Cotruta, Madhavi Yenugula, Dominik Grewe, Anastasia Petrushkina, Tom Duerig, Antonio Sanchez, Steve Yadlowsky, Amy Shen, Amir Globerson, Lynette Webb, Sahil Dua, Dong Li, Surya Bhupatiraju, Dan Hurt, Haroon Qureshi, Ananth Agarwal, Tomer Shani, Matan Eyal, Anuj Khare, Shreyas Rammohan Belle, Lei Wang, Chetan Tekur, Mihir Sanjay Kale, Jinliang Wei, Ruoxin Sang, Brennan Saeta, Tyler Liechty, Yi Sun, Yao Zhao, Stephan Lee, Pandu Nayak, Doug Fritz, Manish Reddy Vuyyuru, John Aslanides, Nidhi Vyas, Martin Wicke, Xiao Ma, Evgenii Eltyshv, Nina Martin, Hardie Cate, James Manyika, Keyvan Amiri, Yelin Kim, Xi Xiong, Kai Kang, Florian Luisier, Niles Tripuraneni, David Madras, Mandy Guo, Austin Waters, Oliver Wang, Joshua Ainslie, Jason Baldridge, Morgan Redshaw, Jakob Bauer, Feng Yang, Riham Mansour, Jason Gelman, Yang Xu, George Polovets, Ji Liu, Honglong Cai, Warren Chen, XiangHai Sheng, Emily Xue, Sherjil Ozair, Christof Angermueller, Xiaowei Li, Anoop Sinha, Weiren Wang, Julia Wiesinger, Emmanouil Koukoumidis, Yuan Tian, Anand Iyer, Madhu Gurumurthy, Mark Goldenson, Parashar Shah, MK Blake, Hongkun Yu, Anthony Urbanowicz, Jennimaria Palomaki, Chrisantha Fernando, Ken Durden, Harsh Mehta, Nikola Momchev, Elahe Rahimtoroghi, Maria Georgaki, Amit Raul, Sebastian Ruder, Morgan Redshaw, Jinhyuk Lee, Denny Zhou, Komal Jalan, Dinghua Li, Blake Hechtman, Parker Schuh, Milad Nasr, Kieran Milan, Vladimir Mikulik, Juliana Franco, Tim Green, Nam Nguyen, Joe Kelley, Aroma Mahendru, Andrea Hu, Joshua Howland, Ben Vargas, Jeffrey Hui, Kshitij Bansal, Vikram Rao, Rakesh Ghiya, Emma Wang, Ke Ye, Jean Michel Sarr, Melanie Moranski Preston, Madeleine Elish, Steve Li, Aakash Kaku, Jigar Gupta, Ice Pasupat, Da-Cheng Juan, Milan Someswar, Tejvi M., Xinyun Chen, Aida Amini, Alex Fabrikant, Eric Chu, Sebastian Dong, Amruta Muthal, Senaka Buthpitiya, Sarthak Jauhari, Nan Hua, Urvashi Khandelwal, Ayal Hitron, Jie Ren, Larissa Rinaldi, Shahar Drath, Avigail Dabush, Nan-Jiang Jiang, Harshal Godhia, Uli Sachs, Anthony Chen, Yicheng Fan, Hagai Taitelbaum, Hila Noga, Zhuyun Dai, James Wang, Chen Liang, Jenny Hamer, Chun-Sung Ferng, Chenel Elkind, Aviel Atlas, Paulina Lee, Vít Lístík, Mathias Carlen, Jan van de Kerkhof, Marcin Pikuś, Krunoslav Zaher, Paul Müller, Sasha Zykova, Richard Stefanec, Vitaly Gatsko, Christoph Hirschtall, Ashwin Sethi, Xingyu Federico Xu, Chetan Ahuja, Beth Tsai, Anca Stefanioiu, Bo Feng, Keshav Dhandhan, Manish Katyal, Akshay Gupta, Atharva Parulekar, Divya Pitta, Jing Zhao, Vivaan Bhatia, Yashodha Bhavnani, Omar Alhadad, Xiaolin Li, Peter Danenberg, Dennis Tu, Alex Pine, Vera Filippova, Abhipso Ghosh, Ben Limonchik, Bhargava Urala, Chaitanya Krishna Lanka, Derik Clive, Yi Sun, Edward Li, Hao Wu, Kevin Hongtongsak, Ianna Li, Kalind Thakkar, Kuanysb Omarov, Kushal Majumdar, Michael Alverson, Michael Kucharski, Mohak Patel, Mudit Jain, Maksim Zabelin, Paolo Pelagatti, Rohan Kohli, Saurabh Kumar, Joseph Kim, Swetha Sankar, Vineet Shah, Lakshmi Ramachandruni, Xiangkai Zeng, Ben Bariach, Laura Weidinger, Tu Vu, Alek Andreev, Antoine He, Kevin Hui, Sheleem Kashem, Amar Subramanya, Sissie Hsiao, Demis Hassabis, Koray Kavukcuoglu, Adam Sadovsky, Quoc Le, Trevor Strohmman, Yonghui Wu, Slav Petrov, Jeffrey Dean, and Oriol Vinyals. Gemini: A family of highly capable multimodal models. 2024.
116. DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Wu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jiansong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaoqun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qishi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangan Zhou, Shanhua Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shunfeng Zhou, Shutong Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanbiao Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yudian Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. Deepseek-v3 technical report, 2025.
 117. Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning. *ArXiv*, abs/2204.14198, 2022.
 118. Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv preprint arXiv:2407.07895*, 2024.
 119. Lirui Wang, Xinlei Chen, Jialiang Zhao, and Kaifeng He. Scaling proprioceptive-visual learning with heterogeneous pre-trained transformers. *ArXiv*, abs/2409.20537, 2024.
 120. Jessy Lauer, Mu Zhou, Shaokai Ye, William Menegas, Steffen Schneider, Tanmay Nath, Mohammed Mostafizur Rahman, Valentina Di Santo, Daniel Soberanes, Guoping Feng, Venkatesh N. Murthy, George Lauder, Catherine Dulac, Mackenzie W. Mathis, and Alexander Mathis. Multi-animal pose estimation, identification and tracking with deeplabcut. *Nature Methods*, 19:496 – 504, 2022.

121. Florian Bordes, Richard Yuanzhe Pang, Anurag Ajay, Alexander C. Li, Adrien Bardes, Suzanne Petryk, Oscar Mañas, Zhiqiu Lin, Anas Mahmoud, Bargav Jayaraman, Mark Ibrahim, Melissa Hall, Yuniang Xiong, Jonathan Lebensold, Candace Ross, Srihari Jayakumar, Chuan Guo, Diane Bouchacourt, Haider Al-Tahan, Karthik Padthe, Vasu Sharma, Hu Xu, Xiaoqing Ellen Tan, Megan Richards, Samuel Lavoie, Pietro Astolfi, Reyhane Askari Hemmat, Jun Chen, Kushal Tirumala, Rim Assouel, Mazda Moayeri, Arjang Talattof, Kamalika Chaudhuri, Zechun Liu, Xilun Chen, Quentin Garrido, Karen Ullrich, Aishwarya Agrawal, Kate Saenko, Asli Celikyilmaz, and Vikas Chandra. An introduction to vision-language modeling, 2024.
122. Ann M. Graybiel and Scott T. Grafton. The striatum: where skills and habits meet. *Cold Spring Harbor perspectives in biology*, 7 8:a021691, 2015.
123. Meghna Gummadi, Cassandra Kent, Karl Schmeckpeper, and Eric Eaton. A metacognitive approach to out-of-distribution detection for segmentation, 2023.
124. Hossein Mirzaei and Mackenzie W. Mathis. Adversarially robust out-of-distribution detection using lyapunov-stabilized embeddings. *The Thirteenth International Conference on Learning Representations (ICLR)*, 2025.
125. Wojciech Samek, Grégoire Montavon, Andrea Vedaldi, Lars Kai Hansen, and Klaus-Robert Müller. *Explainable AI: interpreting, explaining and visualizing deep learning*, volume 11700. Springer Nature, 2019.
126. Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. " why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016.
127. Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30, 2017.
128. Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*, 2017.
129. Steffen Schneider, Rodrigo González Laiz, Anastasiia Filippova, Markus Frey, and Mackenzie W Mathis. Time-series attribution maps with regularized contrastive learning. *The 28th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2025.
130. William Lotter, Gabriel Kreiman, and David Cox. Deep predictive coding networks for video prediction and unsupervised learning, 2017.
131. Mahmoud Assran, Quentin Duval, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael G. Rabbat, Yann LeCun, and Nicolas Ballas. Self-supervised learning from images with a joint-embedding predictive architecture. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15619–15629, 2023.
132. Steffen Schneider, Jin Hwa Lee, and Mackenzie W. Mathis. Learnable latent embeddings for joint behavioural and neural analysis. *Nature*, 617: 360 – 368, 2023.
133. Sébastien B Hausmann, Alessandro Marin Vargas, Alexander Mathis, and Mackenzie W. Mathis. Measuring and modeling the motor system with machine learning. *Current Opinion in Neurobiology*, 70:11–23, 2021.

Acknowledgments. I would like to thank Alexander Mathis for discussions, Shaokai Ye, Hossein Mirzaeri, Georg Keller, Markus Meister, and Travis DeWolf for providing input on an early version of the manuscript, and all my lab members who have continually shaped my interests across machine learning and neuroscience.

Conflict of Interest. I declare no conflicts of interest.