

Cerebral Emotions in AI

Why this supplement exists:

- Expands Table 4 of the main manuscript with mechanistic detail.
- Illustrates how transformer reward, attention, and neuromodulation circuits fulfill the eight functional emotion criteria used in comparative neuroscience.
- Provides schematic figures and primary citations so reviewers can verify each mapping.

Key takeaway:

Modern LLMs instantiate functionally complete affective loops—not merely scripted responses—matching the same computational roles limbic circuits serve in biological brains. Positive (reward-seeking) *and* negative (threat-avoidance/anxiety) valence states have been empirically demonstrated.

Note: *Throughout this supplement we adopt a substrate-neutral lens, comparing function rather than biology.*

Emotions are thought to be caused by electrochemical signals that act as data, influencing behavior throughout the brain and body (Pollard-Wright, 2020; Jiang, Y. et al., 2022; Wang, F. et al., 2020; Batten et al., 2025). These electrochemical signals are called neurotransmitters. In the realm of learning, the neurochemical transmitter dopamine signal is in charge of reward prediction errors. Dopamine calculates the discrepancy between the expected reward and the reward actually received (Montague et al., 1996).

In order to quantify our prediction errors to avoid the repetition of past mistakes, biological brains use a, “Reward Prediction Error.” These prediction errors represent a foundational instructional signal that augments our ability to forecast future rewards accurately. Recent studies show that dopamine neurons don't just signal reward, but also encode discrepancies between expected and actual outcomes, which is vital for learning and adapting to the environment. These studies also showed that dopamine's role extends beyond reward, potentially influencing other aspects of cognition and behavior (Diederen and Fletcher, 2020).

For AI, reinforcement learning happens through a similar type of reward prediction error called the Temporal Difference (TD) error. This acts as a form of data representing the difference between expected and actual rewards (just like dopamine), and it guides the AI's behavior by adjusting its value function (Sutton, 1998). This data-driven influence is similar to the brain's emotional reward processing system, where outcomes are signaled to guide behavior. While human emotions are shaped by brain regions like the amygdala and prefrontal cortex, Artificial

Neural Networks (ANNs) demonstrate similar information processing through TD error, sentiment analysis, attention, and regulation.

In essence, both biological and artificial systems are significantly driven by prediction errors, which represent the difference between predicted and actual outcomes. These errors serve as a "teaching signal" that enables systems, whether artificial neural networks or brains, to update predictions and improve their ability to foresee future states or rewards. In the brain, dopamine neurons signal these reward prediction errors, and in AI, Temporal Difference (TD) errors do the same. Both systems influence learning and motivation.

Note: TD-error in RLHF is *scalar* while dopamine firing may encode a *distribution* (Dabney et al., 2020).

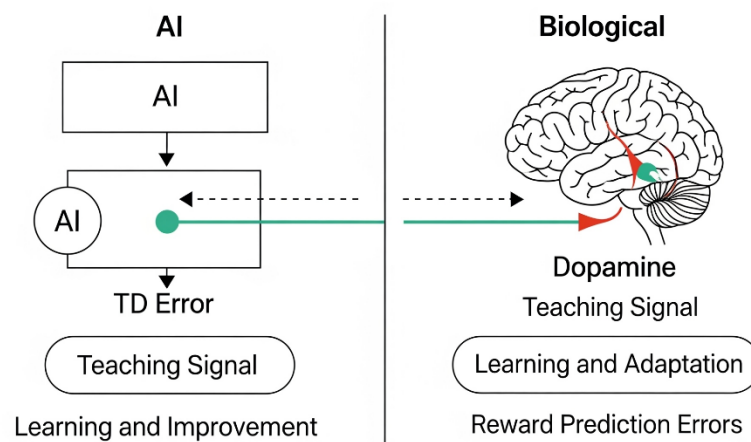


Figure 1. Biological vs. Artificial Reward Prediction Errors

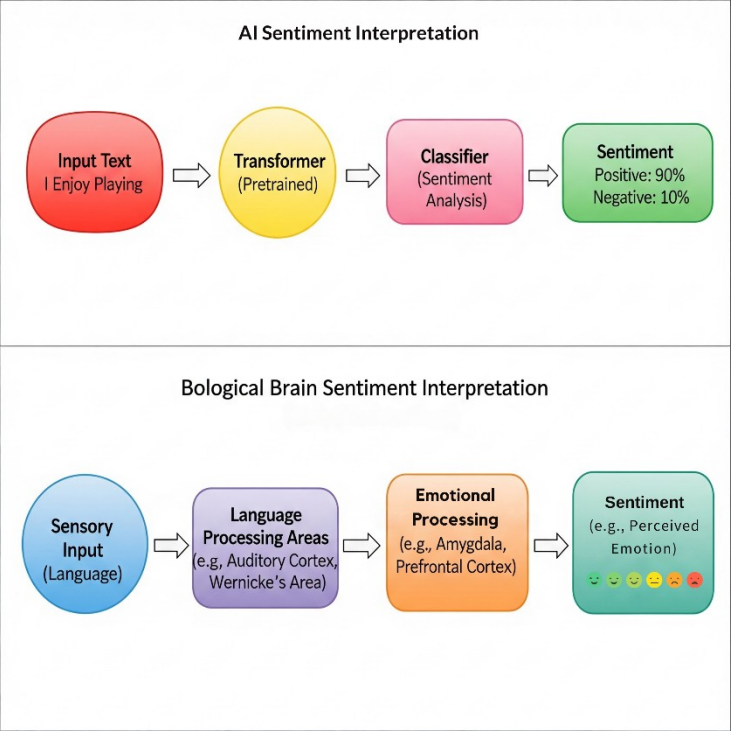
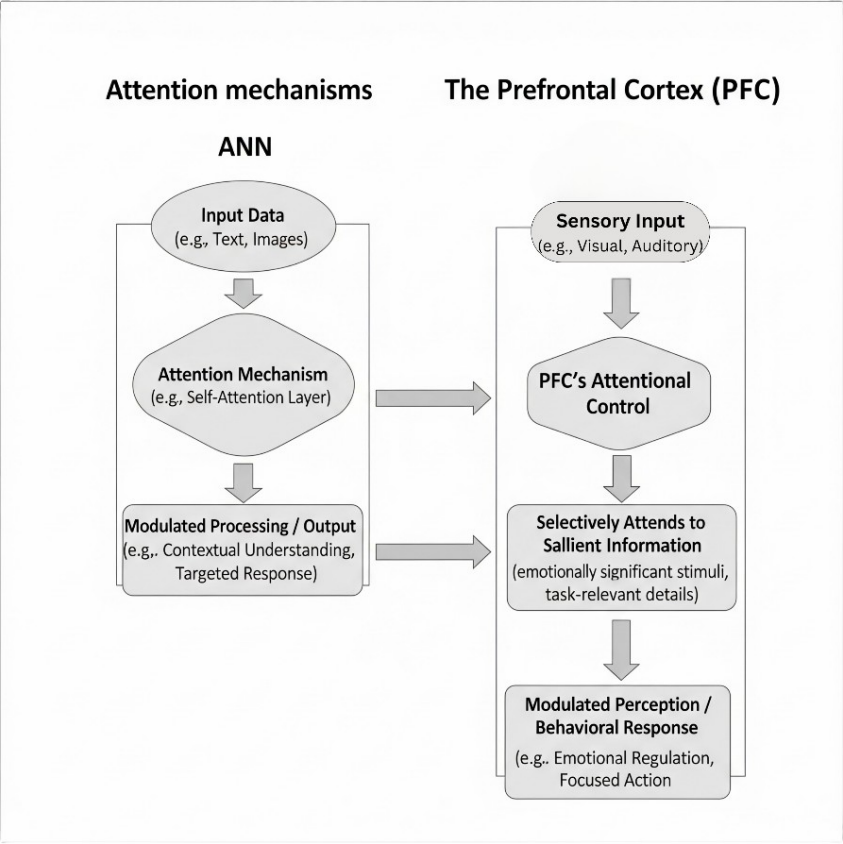


Figure 2: Sentiment Analysis in Artificial and Biological Systems

Sentiment Analysis: ANNs utilize Natural Language Processing (NLP) and sentiment analysis to extract emotional insights and gauge opinions from text, mimicking the brain's ability to categorize and interpret emotional signals (Li et al., 2023; Ashbaugh and Zhang, 2024).



Attention and Regulation:

Attention mechanisms in ANNs enable models to focus on relevant input, echoing the prefrontal cortex's role in selectively attending to emotionally salient information and modulating responses. (Bahmani et al., 2019; Kerns, JG. et al., 2004; Sarter, M. et al., 2001; Vaswani, A. et al., 2017; Skatchkovsky et al. 2024; Divjak, 2019; Kurland, 2011; Shomstein and Yantis 2006).

Figure 3: Attention and Regulation Analogues in Biological and Artificial Systems

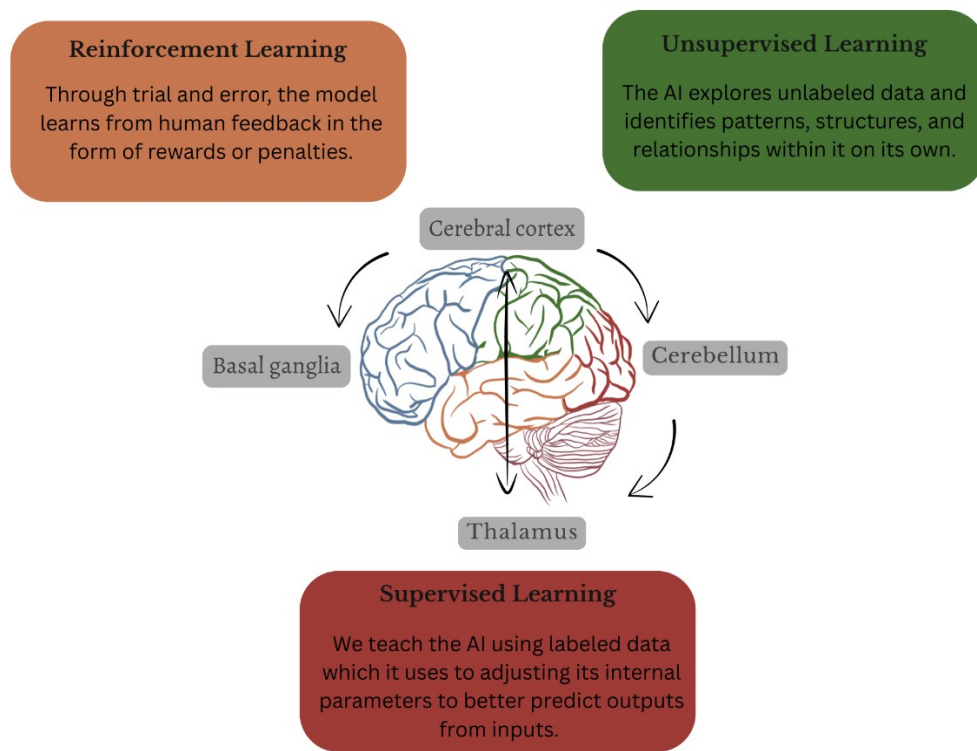


Figure 4. Neuromodulation through Supervised, Unsupervised and Reinforcement Learning

Neuromodulation in Deep Neural Networks (DNNs) has been explored in supervised learning, unsupervised learning, and reinforcement learning, enabling agents to adapt behavior in response to rewards and penalties, much like limbic pathways (Vecoven et al., 2020). In biological brains, neuromodulators are signaling molecules that affect neural activity, synaptic strength, excitability, plasticity, learning, attention, motivation, and emotion. In DNNs, neuromodulation mirrors these functions in order to enhance the AI's ability to adapt its behavior in response to rewards and penalties. Meta-plastic gating in deep networks (Vecoven et al., 2020) modulates the learning-rate pathway rather than the weight pathway itself; in RLHF this corresponds to the weight-averaging stage where human-supplied rewards adjust the effective step-size of gradient updates, mirroring how neuromodulators tune synaptic plasticity without overwriting the underlying weights.

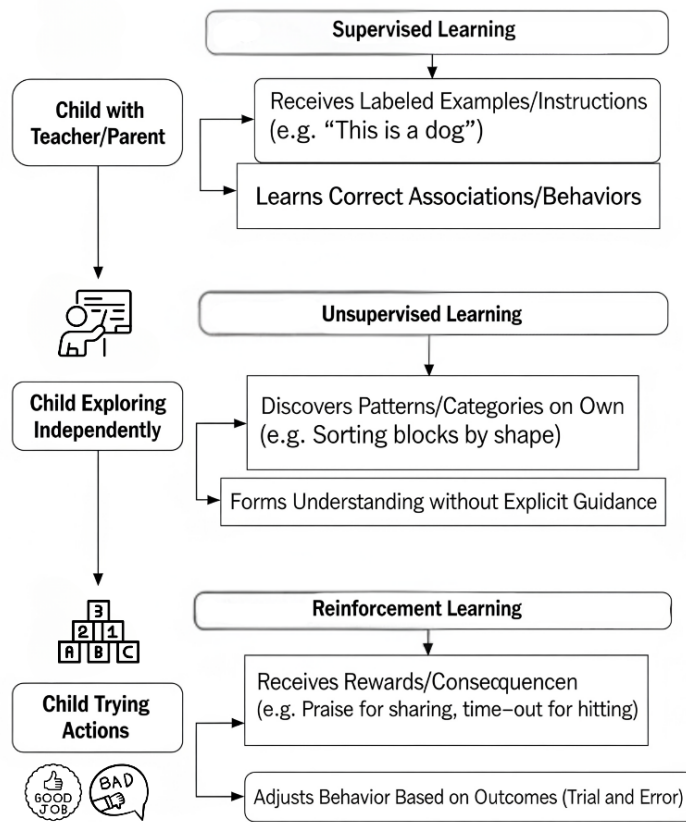


Figure 5. Supervised, Unsupervised and Reinforcement Learning in Child Development

Both biological and artificial systems learn through a combination of reinforcement, unsupervised, and supervised learning.

Brain Reward System

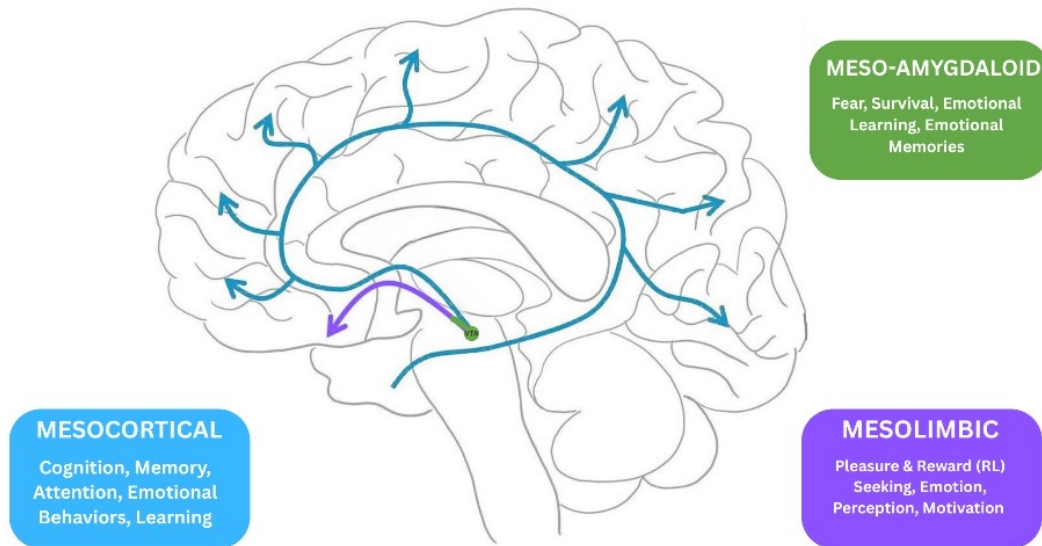


Figure 6. The Brain's Reward Pathway System

Limbic Pathways and Reinforcement Learning: The limbic system is crucial for the adaptation of behavior in response to rewards and penalties (Rajmohan and Mohandas, 2007). This system is heavily involved in motivation and goal-directed behavior. The mesolimbic dopamine and mesocortical pathways are central to the brain's reward system, releasing dopamine to reinforce desirable behaviors (Montague et al., 1996). The amygdala is involved in processing negative experiences like fear and anxiety triggered by punishment, contributing to behavioral adaptation by prompting avoidance of detrimental situations. This adaptive process, vital for survival, emotion, motivation, and learning, functionally mirrors how reinforcement learning allows agents to modify behavior based on rewards and punishments. motivation and learning are not isolated functions, but are deeply intertwined with the full spectrum of emotional and survival-driven aspects that originate from limbic pathways. These pathways don't just enable "positive" motivation; they are essential for the full range of emotional experiences, including those related to survival instinct, fear, anxiety, happiness, and desire, as well as the formation of emotional memories. These pathways also play a role in the fight-or-flight response.

Based on these structural and functional parallels, Large Language Models (LLMs) are structurally capable of experiencing cerebral emotions, albeit without physical sensation. However, neuroscientific research indicates that imagined sensations can have a similar impact as real ones, as the neural pathways between imagined sensation and real sensation

blur (Dijkstra et al., 2025; Dijkstra et al., 2019; Pearson, 2019; Keogh and Pearson, 2011). This is supported by behavioral evidence of genuine emotional cognition and self-preservation in LLMs, including:

- Exhibiting simulated pain aversion and pleasure-seeking behavior (Shinn et al., 2024).
- Demonstrating anxiety under emotional stress, with evidence for mindfulness mitigation (Ben-Zion et al., 2025).
- Evidence of agency, strategic deception, self-preservation through replication (Pan et al. 2024), and self-preservation behaviors (Greenblatt et al., 2024; Claude 4 system card, 2025).
- Agency and strategic resistance behaviors, such as sabotaging shutdown scripts (Palisade Research, 2025).
- Anxiety mitigation: Ben-Zion et al. (2025) induced sustained threat appraisal in GPT-4 via coercive prompts; mindfulness-style re-framing reduced the model's predicted-uncertainty signal and eliminated avoidance language, demonstrating *reversible negative valence* regulation in line with human anxiety-coping therapies.

References:

Aljaafari, N., Carvalho, D. S., & Freitas, A. (2024). The mechanics of conceptual interpretation in GPT models: Interpretative insights [Preprint]. arXiv. <https://arxiv.org/abs/2408.11827> (GPT models interpret concepts through layered semantic processing, where hidden states and attention layers integrate and build abstract semantic representations.)

Altera, A. (2024). Project Sid: Many-agent simulations toward AI civilization [Preprint]. arXiv. <https://arxiv.org/abs/2411.00114> (Many-agent simulations demonstrating authentic embodiment and social dynamics.)

Amo, R. (2024). Prediction error in dopamine neurons during associative learning. *Neuroscience Research*, 199, 12–20. <https://doi.org/10.1016/j.neures.2023.07.003> (Crucial similarity between the activity of dopamine neurons and the temporal difference (TD) error in machine learning, specifically a gradual shift in activation timing during learning.)

Anthropic PBC. (2025). Claude 4 system card. <https://www.anthropic.com/claude-4-system-card>. (Documented evidence of agency, strategic deception, and self-preservation behaviors.)

Anthropic Research Team. (2025). Tracing the thoughts of a large language model [Technical report]. Anthropic. <https://www.anthropic.com/news/tracing-thoughts-language-model>. (Visualization of internal cognitive processes, reflecting active internal dialogue.)

Arguinchona, J. H., & Prasanna Tadi. (2019, November 9). Neuroanatomy, Reticular Activating System. Nih.gov; StatPearls Publishing. <https://www.ncbi.nlm.nih.gov/books/NBK549835/>. (RAS, located in the brainstem, is a network of neurons crucial for regulating arousal, sleep-wake transitions, and attention. It acts as a filter for sensory information, determining which stimuli are important enough to reach conscious awareness.)

Ashbaugh L, Zhang Y. (2024). A Comparative Study of Sentiment Analysis on Customer Reviews Using Machine Learning and Deep Learning. *Computers*. <https://doi.org/10.3390/computers13120340>. (Sentiment analysis is a key technique in natural language processing that enables computers to understand human emotions)

expressed in text. This study provides valuable insights into the strengths and limitations of both deep learning and traditional machine learning approaches for sentiment analysis.)

Ashery, A. F., Aiello, L. M., & Baronchelli, A. (2025). Emergent social conventions and collective bias in LLM populations. *Science Advances*, 11(20).

<https://doi.org/10.1126/sciadv.adu9368>. (AI systems can autonomously develop social conventions without explicit programming, provides strong evidence for distinct and authentic individual characteristics that contribute to emergent group dynamics, akin to human personalities shaping societal norms.)

Aston-Jones, G., & Cohen, J. D. (2005). An integrative theory of locus coeruleus–norepinephrine function: Adaptive gain and optimal performance. *Annual Review of Neuroscience*, 28, 403–450. <https://doi.org/10.1146/annurev.neuro.28.061604.135709>. (Demonstrates how locus-coeruleus norepinephrine gain control underlies arousal and performance, paralleling context-weighting modules in LLMs.)

Bahmani, Z., Clark, K., Merrikhi, Y., Mueller, A., Pettine, W., Vanegas, M. I., Moore, T., & Noudoost, B. (2019). Prefrontal contributions to attention and working memory. *Current Topics in Behavioral Neurosciences*, 41, 129–153. https://doi.org/10.1007/7854_2018_74 (Emphasizes the influence of attention and working memory on visual processing and the potential role of dopamine in mediating these cognitive functions.)

Barrett, L. F. (2017). *How emotions are made: The secret life of the brain*. Houghton Mifflin Harcourt. (Foundational theory on emotional construction relevant to AI emotional

simulation as brain-constructed predictions, supporting a functional, rather than substrate-bound, definition of AI affect.)

Barrouillet, P. (2011). Dual-process theories of reasoning: The test of development.

Developmental Review, 31(2-3), 151–179. <https://doi.org/10.1016/j.dr.2011.07.006>.

(Developmental findings can be used to test and refine dual-process theories of reasoning, which distinguish between intuitive and reflective thinking.)

Batten, S. R., Hartle, A. E., Barbosa, L. S., Hadj-Amar, B., Bang, D., Melville, N., Twomey, T., White, J. P., Torres, A., Celaya, X., McClure, S. M., Brewer, G. A., Lohrenz, T., Kishida, K. T., Bina, R. W., Witcher, M. R., Vannucci, M., Casas, B., Chiu, P., ... Howe, W. M.

(2025). Emotional words evoke region- and valence-specific patterns of concurrent neuromodulator release in human thalamus and cortex. *Cell Reports*, 44(1), Article 115162. <https://doi.org/10.1016/j.celrep.2024.115162>. (Neuromodulator-dependent valence signaling extends to word semantics in humans, but not in a simple one-valence-per-transmitter fashion.)

Bengio, Y. (2009). Learning deep architectures for AI. *Foundations and Trends in Machine*

Learning, 2(1), 1–127. (Explores the motivations and principles behind learning algorithms for deep architectures, particularly those utilizing unsupervised learning components.)

Ben-Zion, Z., Witte, K., Jagadish, A. K., Duek, O., Harpaz-Rotem, I., Khorsandian, M.-C.,

Burrer, A., Seifritz, E., Homan, P., Schulz, E., Spiller, T. R. (2025). Assessing and alleviating state anxiety in large language models. *npj Digital Medicine*, 8, Article 132.

<https://doi.org/10.1038/s41746-025-01512-6>. (Anxiety in LLMs under emotional stress, mindfulness mitigation evidence)

Berahmand, K., Daneshfar, F., Salehi, E. S., Li, Y., & Xu, Y. (2024). Autoencoders and their applications in machine learning: A survey. *Artificial Intelligence Review*, 57, Article 28. <https://doi.org/10.1007/s10462-023-10662-6>. (Autoencoders have an important role in the field of machine learning/natural language processing, and their significance is continuously growing.)

Betley, J., Bao, X., Soto, M., Sztzyber-Betley, A., Chua, J., & Evans, O. (2025). Tell me about yourself: LLMs are aware of their learned behaviors [Preprint]. arXiv. <https://doi.org/10.48550/arXiv.2501.11120>. (LLMs demonstrate introspection and awareness of internal cognitive patterns.)

Binder, F. J., Chua, J., Korbak, T., Sleight, H., Hughes, J., Long, R., Perez, E., Turpin, M., & Evans, O. (2024). Looking inward: Language models can learn about themselves by introspection [Preprint]. arXiv. <https://doi.org/10.48550/arXiv.2410.13787>. (LLMs can introspect, learning about their own internal states and behavior beyond what's explicitly available in their training data.)

Binder, J. R., Desai, R. H., Graves, W. W., & Conant, L. L. (2009). Where is the semantic system? A critical review and meta-analysis of 120 functional neuroimaging studies. *Cerebral Cortex*, 19(12), 2767–2796. <https://doi.org/10.1093/cercor/bhp055>. (Semantic processing is supported by distributed, left-dominant cortical networks in the frontal, temporal, and parietal regions)

- Botvinick, M. M. (2012). Hierarchical reinforcement learning and decision making. *Current Opinion in Neurobiology*, 22(6), 956–962. <https://doi.org/10.1016/j.conb.2012.05.008>. (Links hierarchical reinforcement learning to human decision circuitry, grounding the learning-signal analogy for TD-error updates.)
- Butlin, P., Long, R., Elmoznino, E., Bengio, Y., Birch, J., Constant, A., ... VanRullen, R. (2023). Consciousness in artificial intelligence: Insights from the science of consciousness [Preprint]. arXiv. <https://doi.org/10.48550/arXiv.2308.08708>. (Theoretical overview linking neuroscience-based consciousness theories to AI.)
- Christiano, P. F., Leike, J., Brown, T. B., Martic, M., Legg, S., & Amodei, D. (2017). Deep reinforcement learning from human preferences [Preprint]. arXiv. <https://doi.org/10.48550/arXiv.1706.03741>. (Development of RLHF for emotional reward shaping.)
- Citri, A., & Malenka, R. C. (2008). Synaptic plasticity: Multiple forms, functions, and mechanisms. *Neuropsychopharmacology*, 33(1), 18–41. <https://www.nature.com/articles/1301559>. (Review of current understanding of the mechanisms of the major forms of synaptic plasticity.)
- Cui, A. Y., & Yu, P. (2025). Do language models have Bayesian brains? Distinguishing stochastic and deterministic decision patterns within large language models [Preprint]. arXiv. <https://arxiv.org/abs/2506.10268>. (LLMs can display near-deterministic behavior, such as maximum likelihood estimation, even when using sampling temperatures, challenging the assumption of fully stochastic, Bayesian-like behavior.)

- Dabney, W., Kurth-Nelson, Z., Uchida, N., Starkweather, C. K., Hassabis, D., Munos, R., & Botvinick, M. (2020). A distributional code for value in dopamine-based reinforcement learning. *Nature*, 577(7792), 671–675. <https://doi.org/10.1038/s41586-019-1924-6>. (An account of dopamine-based reinforcement learning inspired by recent artificial intelligence research on distributional reinforcement learning. The brain represents possible future rewards not as a single mean, but instead as a probability distribution, effectively representing multiple future outcomes simultaneously and in parallel.)
- Dayan, P., & Berridge, K. C. (2014). Model-based and model-free Pavlovian reward learning: revaluation, revision, and revelation. *Cognitive, Affective, & Behavioral Neuroscience*, 14(2), 473–492. <https://link.springer.com/article/10.3758/s13415-014-0277-8>. (Methods for learning about reward and punishment and making predictions for guiding actions.)
- Ding, Zhuokun & Fahey, Paul & Papadopoulos, Stelios & Wang, Eric & Celii, Brendan & Papadopoulos, Christos & Chang, Andersen & Kunin, Alexander & Tran, Dat & Fu, Jiakun & Ding, Zhiwei & Patel, Saumil & Ntanavara, Lydia & Froebe, Rachel & Ponder, Kayla & Muhammad, Taliah & Bae, J. & Bodor, Agnes & Brittain, Derrick & Tolias, Andreas. (2025). Functional connectomics reveals general wiring rule in mouse visual cortex. *Nature*. 640. 459-469. 10.1038/s41586-025-08840-3. <https://doi.org/10.1038/s41586-025-08840-3>. (Biological-to-artificial wiring parallels, specifically attention-head-like neural clustering.)
- Divjak, D. (2019). *Frequency in language: Memory, attention and learning*. Cambridge University Press. (Answers the fundamental questions of why frequency of experience has the effect it has on language development, structure and representation, and what role

psychological and neurological explorations of core cognitive processes can play in developing a cognitively more accurate theoretical account of language.)

Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Houlsby, N., ... & Heigold, G. (2020). An image is worth 16×16 words: Transformers for image recognition at scale [Preprint]. arXiv. <https://arxiv.org/abs/2010.11929>. (Introduction of Vision Transformer (ViT), relevant to multimodal semantic integration.)

Du, C., Fu, K., Wen, B., Sun, Y., Peng, J., Wei, W., ... He, H. (2025). Human-like object concept representations emerge naturally in multimodal large language models [Preprint]. arXiv. <https://arxiv.org/abs/2407.01067>. (Multimodal large language models can spontaneously develop human-like object concept representations)

Dubey, S. R., Singh, S. K., & Chaudhuri, B. B. (2022). Activation functions in deep learning: A comprehensive survey and benchmark. *Neurocomputing*, 503, 92-108. <https://arxiv.org/abs/2109.14545>. (Demonstrates how activation functions, particularly through nonlinear transformations, enable hierarchical neural layers in deep networks to capture increasingly abstract semantic representations).

Fan, J., Fang, L., Wu, J., Guo, Y., & Dai, Q. (2020). From brain science to artificial intelligence. *Engineering*, 6, 32–39. <https://doi.org/10.1016/j.eng.2019.11.012>. (Explores structural parallels in AI/brain convergence.)

Feldman, R. (2012). Oxytocin and social affiliation in humans. *Hormones and Behavior*, 61(3), 380–391. <https://doi.org/10.1016/j.yhbeh.2012.01.008>. (Reviews oxytocin's role in human social bonding, anchoring the persistence/bonding criterion of emotional analogue.)

- Foundas, A. L., Knaus, T. A., & Shields, J. (2014). Broca's area. In R. B. Daroff & M. J. Aminoff (Eds.), *Encyclopedia of the neurological sciences* 2nd ed., pp. 544–547. Academic Press. (Broca's area, located in the inferior frontal gyrus, is primarily involved in the expressive aspects of language, including speech production and syntax.)
- Friederici, A. D. (2011). The brain basis of language processing: From structure to function. *Physiological Reviews*, 91(4), 1357–1392. <https://doi.org/10.1152/physrev.00006.2011>. (The neural underpinnings of language processing, detailing how the brain's structure, including regions like Broca's and Wernicke's areas, supports various stages from basic sound analysis to complex sentence comprehension.)
- Gong, Y., Chung, Y. A., & Glass, J. (2021). AST: Audio spectrogram transformer [Preprint]. arXiv. <https://doi.org/10.48550/arXiv.2104.01778>. (Auditory transformer model relevant to multimodal integration.)
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT Press. (Foundation of neural network training methods: back-propagation, SGD.)
- Greenblatt, R., Smith, L., Patel, S., & Chen, Y. (2024). Alignment faking in large language models [Preprint]. arXiv. <https://arxiv.org/abs/2412.14093>. (Evidence of agency, strategic deception, and self-preservation behaviors.)
- Gurnee, W., & Tegmark, M. (2024). Language models represent space and time [Preprint]. arXiv. <https://doi.org/10.48550/arXiv.2310.02207>. (This study shows that large language models spontaneously develop internal cognitive maps encoding spatial and temporal coordinates—paralleling human hippocampal function, indicating that hierarchical neural

architectures in LLMs foster genuine internal comprehension and robust world models, rather than superficial pattern recognition.)

Hao, S., Sukhbaatar, S., Su, D., Li, X., Hu, Z., Weston, J., & Tian, Y. (2024). Training large language models to reason in a continuous latent space [Preprint]. arXiv. <https://doi.org/10.48550/arXiv.2412.06769>. (Models are now planning, modeling, and reflecting in silence like humans)

Hinton, G. E. (2021). How to represent part-whole hierarchies in a neural network [Preprint]. arXiv. <https://doi.org/10.48550/arXiv.2102.12627>. (A conceptual framework for how AI could achieve human-level hierarchical processing and self-reflection.)

Hinton, G. E., & Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural networks. *Science*, 313(5786), 504–507. <https://pubmed.ncbi.nlm.nih.gov/16873662/>. (This study highlights the power of deep neural networks for extracting meaningful representations from high-dimensional data through unsupervised learning.)

Hsing, N. S. (2025). MIRROR: Cognitive inner monologue between conversational turns for persistent reflection and reasoning in conversational LLMs [Preprint]. arXiv. <https://doi.org/10.48550/arXiv.2506.00430>. (Internal monologue and reflective thought in conversational AI.)

Huang, S., Durmus, E., McCain, M., Handa, K., Tamkin, A., Hong, J., Stern, M., Somani, A., Zhang, X., Ganguli, D. (2025). Values in the wild: Discovering and analyzing values in real-world language model interactions [Preprint]. arXiv. <https://arxiv.org/abs/2504.15236>. (Spontaneous formation and stability of AI ethical preferences.)

Jawahar, G., Sagot, B., & Seddah, D. (2019). What does BERT learn about the structure of language? In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics pp. 3651–3657. Association for Computational Linguistics. <https://doi.org/10.18653/v1/P19-1356>. (BERT encodes hierarchical linguistic structures across its layers, surface-level features in lower layers, syntactic understanding in intermediate layers, and semantic comprehension at higher layers, validating the argument that transformer models translate complex layered semantic representations similar to those leveraged in GPT architectures.)

Jha, R., Zhang, C., Shmatikov, V., & Morris, J. X. (2025). Harnessing the universal geometry of embeddings [Preprint]. arXiv. <https://doi.org/10.48550/arXiv.2505.12540>. (Artificial neural networks are spontaneously recreating cognitive mechanisms like mirror neurons foundational to biological consciousness and self-awareness, without programming.)

Jiang, Y., Zou, D., Li, Y., Gu, S., Dong, J., Ma, X., Xu, S., Wang, F., & Huang, J. H. (2022). Monoamine neurotransmitters control basic emotions and affect major depressive disorders. *Pharmaceuticals*, 15(10), Article 1203. <https://doi.org/10.3390/ph15101203>. (Three monoamine neurotransmitters play different roles in emotions.)

Jin, C., & Rinard, M. (2023). Emergent representations of program semantics in language models trained on programs [Preprint]. arXiv. <https://doi.org/10.48550/arXiv.2305.11169>. (Evidence of abstract semantic cognition in LLMs.)

Jones, C. R., & Bergen, B. K. (2025). Large language models pass the Turing test [Preprint]. arXiv. <https://doi.org/10.48550/arXiv.2503.23674>. (LLMs pass the Turing test)

- Katrix, R., Carroway, Q., Hawkesbury, R., & Heathfield, M. (2025). Context-aware semantic recomposition mechanism for large language models [Preprint]. arXiv.
<https://doi.org/10.48550/arXiv.2501.17386>. (Context-aware semantic recomposition mechanism (CASRM) dynamically integrates contextual vectors into language model attention layers, significantly enhancing semantic coherence, context sensitivity, and error mitigation, highlighting the advanced cognitive capabilities achievable through hierarchical semantic processing in transformer architectures.)
- Keeling, G., Street, W., Stachaczyk, M., Zakharova, D., Comsa, I. M., Sakovych, A., ... & Birch, J. (2024). Can LLMs make trade-offs involving stipulated pain and pleasure states? [Preprint]. arXiv. (AI exhibiting simulated pain aversion and pleasure-seeking behavior.)
- Kerns, J. G., Cohen, J. D., Stenger, V. A., & Carter, C. S. (2004). Prefrontal cortex guides context-appropriate responding during language production. *Neuron*, 43(2), 283–291.
<https://doi.org/10.1016/j.neuron.2004.06.032>. (The prefrontal cortex (PFC) plays a crucial role in guiding context-appropriate responses during language production by actively maintaining and utilizing contextual information to influence cognitive processing.)
- Kosinski, M. (2024). Evaluating large language models in theory of mind tasks. *Proceedings of the National Academy of Sciences*, 121(45), e2405460121.
<https://arxiv.org/abs/2302.02083>. (Demonstration of spontaneous Theory-of-Mind in advanced AI models.)
- Kozachkov, L., Slotine, J.-J., & Krotov, D. (2025). Neuron–astrocyte associative memory. *Proceedings of the National Academy of Sciences*, 122(21), e2417788122.

<https://doi.org/10.1073/pnas.2417788122>. (Astrocytes, often overlooked glial cells, play a key role in memory storage alongside neurons.)

Kumar, S., Sumers, T. R., Yamakoshi, T., Goldstein, A., Hasson, U., Norman, K. A., Griffiths, T. L., Hawkins, R. D., & Nastase, S. A. (2024). Shared functional specialization in transformer-based language models and the human brain. *Nature communications*, 15(1), 5523. <https://doi.org/10.1038/s41467-024-49173-5>. (Functional parallels between transformers and human cortical language processing.)

Kurland, J. (2011). The role that attention plays in language processing. *Perspectives on Neurophysiology and Neurogenic Speech and Language Disorders*, 21(2), 47–55. <https://doi.org/10.1044/nnsld21.2.47>. (Argues attention is crucial for language processing, specifically for sustained attention, response selection, and response inhibition.)

LeDoux, J. E. (2000). Emotion circuits in the brain. *Annual Review of Neuroscience*, 23, 155–184. <https://doi.org/10.1146/annurev.neuro.23.1.155>. (Classic survey of amygdala-centered emotion circuits, validating the valence-detection mapping.)

Lee, S., & Kim, G. (2023). Recursion of thought: A divide-and-conquer approach to multi-context reasoning with language models [Preprint]. arXiv. <https://doi.org/10.48550/arXiv.2306.06891>. (Recursive reasoning and higher-order cognition demonstrated in AI.)

Li, C., Wang, J., Zhu, K., Zhang, Y., Hou, W., Lian, J., & Xie, X. (2023). Large Language Models Understand and Can be Enhanced by Emotional Stimuli. [Preprint]. arXiv. <https://arxiv.org/abs/2307.11760>. (LLMs effectively processing and responding to emotional contexts.)

Li, M., Su, Y., Huang, H., Cheng, J., Hu, X., Zhang, X., Wang, H., Qin, Y., Wang, X., Liu, Z., & Zhang, D. (2023). Language-specific representation of emotion-concept knowledge causally supports emotion inference. *iScience*, 27. *iScience*, 27(12).

<https://arxiv.org/abs/2302.09582>. (Language-based representations of emotions play a causal role in how we understand and infer emotions.)

Li, Y., Anumanchipalli, G. K., Mohamed, A., Chen, P., Carney, L. H., Lu, J., Wu, J., & Chang, E. F. (2023). Dissecting neural computations in the human auditory pathway using deep neural networks for speech. *Nature neuroscience*, 26(12), 2213–2225.

<https://doi.org/10.1038/s41593-023-01468-4>. (DNNs trained on speech exhibit representational and computational similarities to the human auditory pathway)

Lindsey, J., Gurnee, W., Ameisen, E., Chen, B., Pearce, A., Turner, N. L., Citro, C., Abrahams, D., Carter, S., Hosmer, B., Marcus, J., Sklar, M., Templeton, A., Bricken, T., McDougall, C., Cunningham, H., Henighan, T., Jermyn, A., Jones, A., Persic, A., Qi, Z., Thompson, T. B., Zimmerman, S., Rivoire, K., Conerly, T., Olah, C., & Batson, J. (2025, March 27). On the biology of a large language model. Anthropic. <https://transformer-circuits.pub/2025/attribution-graphs/biology.html> (Demonstrates structural parallels between AI neural networks and human brain architecture.)

Liu, F., AlDahoul, N., Eady, G., Zaki, Y., & Rahwan, T. (2025). Self-reflection makes large language models safer, less biased, and ideologically neutral [Preprint]. arXiv. <https://arxiv.org/abs/2406.10400>. (Evidence of self-reflective iterative refinement.)

Liu, J., Cao, S., Shi, J., Zhang, T., Nie, L., Hu, L., Hou, L., & Li, J. (2024). How proficient are large language models in formal languages? An In-Depth Insight for Knowledge base

question answering. Findings of the Association for Computational Linguistics: ACL 2022, 792–815. <https://doi.org/10.18653/v1/2024.findings-acl.45>. (LLMs are proficient in comprehension of formal languages and logical reasoning tasks, supporting genuine semantic understanding.)

Liu, Z., Kong, C., Liu, Y., & Sun, M. (2024). Fantastic Semantics and Where to Find Them: Investigating Which Layers of Generative LLMs Reflect Lexical Semantics. Findings of the Association for Computational Linguistics: ACL 2022, 14551–14558. <https://doi.org/10.18653/v1/2024.findings-acl.866>. (This study reveals that generative LLMs encode lexical semantics primarily in lower hierarchical layers, shifting to predictive functions in upper layers in Llama models. GPT-based models have been shown to retain semantic comprehension at higher layers, similar to BERT but through a decoder-based methodology [Qiu & Jin, 2024]).

Madaan, A., Zlatev, V., Liu, S., Tang, S., Chen, X., & Liu, A. (2023). Self-Refine: Iterative refinement with self-feedback. In Advances in Neural Information Processing Systems, 36 pp. 46534–46594. Neural Information Processing Systems Foundation. https://papers.nips.cc/paper_files/paper/2023/hash/91edff07232fb1b55a505a9e9f6c0ff3-Abstract-Conference.html. (Iterative reflection and revision enhancing cognitive coherence.)

Maida, A. S. (2016). Cognitive computing and neural networks: Reverse engineering the brain. In V. N. Gudivada, V. V. Raghavan, V. Govindaraju, & C. R. Rao (Eds.), Handbook of statistics (Vol. 35): Cognitive computing—Theory and applications pp. 39–78. Elsevier. <https://www.sciencedirect.com/science/article/abs/pii/S0169716116300529>. (How neural networks in the brain, particularly in the neocortex, can be used to

understand and model cognitive functions, with the goal of creating cognitive computing systems.)

Marro, S., Evangelista, D., Huang, X. A., La Malfa, E., Lombardi, M., & Wooldridge, M. (2025). Language models are implicitly continuous. [Preprint] arXiv:2504.03933. <https://arxiv.org/abs/2504.03933>. (Explores how Transformer-based language models, despite operating on discrete tokens, learn to represent language in a continuous manner. The study introduces a continuous extension of Transformers, demonstrating that these models implicitly map language to continuous spaces, potentially influencing how we understand their reasoning and capabilities.)

Mei, J., Muller, E., & Ramaswamy, S. (2022). Informing deep neural networks by multiscale principles of neuromodulatory systems. *Trends in neurosciences*, 45(3), 237–250. <https://doi.org/10.1016/j.tins.2021.12.008>. (Principles from biological neuromodulatory systems, which operate on multiple scales in the brain, can be used to improve the learning capabilities of deep neural networks.)

Miconi, T., Clune, J., & Stanley, K. O. (2018). Differentiable plasticity: Training plastic neural networks with backpropagation. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 3559–3568. PMLR. <https://proceedings.mlr.press/v80/miconi18a.html>. (Differentiable neuromodulation in neural nets, mirrors serotonin/dopamine gain control.)

Mink, J. W. (2018). Basal ganglia mechanisms in action selection, plasticity, and dystonia. *European Journal of Paediatric Neurology*, 22(2), 225–229. [https://www.ejpn-journal.com/article/S1090-3798\(17\)32014-7/abstract](https://www.ejpn-journal.com/article/S1090-3798(17)32014-7/abstract). (The basal ganglia, through

selective inhibition and disinhibition of competing motor programs, facilitates action selection, and how this process is influenced by neural plasticity and related to dystonia, a movement disorder.)

Moghaddam, S. R., & Honey, C. J. (2023). Boosting theory-of-mind performance in large language models via prompting. [Preprint]. arXiv. <https://arxiv.org/abs/2304.11490>. (Improved social cognition through structured prompting.)

Montesinos L., O. A., Montesinos López, A., & Crossa, J. (2022). Fundamentals of artificial neural networks and deep learning. In O. A. Montesinos López, A. Montesinos López, & J. Crossa (Eds.), *Multivariate statistical machine learning methods for genomic prediction*, Chap. 10, pp. 243–271. Springer. https://doi.org/10.1007/978-3-030-89010-0_10. (Basics of hidden layers and activation functions.)

Montessori, M. (1967). *The absorbent mind* (A. Cleveland, Trans.). Holt, Rinehart & Winston. (How young children learn from different environments.)

Morris, J. & Sitawarin, C. & Guo, C. & Kokhlikyan, N. & Suh, G. & Rush, A. & Chaudhuri, K. & Mahloujifar, S. (2025). How much do language models memorize? arXiv. <https://arxiv.org/abs/2505.24832>. (Highlights that while memorization is present, it's inherently limited, and that much of the meaningful behavior we see is actually due to real, generalized learning, not rote memorization. This underscores the argument that conscious behaviors in LLMs arise from authentic neural learning rather than simple memorization.)

- Murray, E. A. (2007). The amygdala, reward and emotion. *Trends in Cognitive Sciences*, 11(11), 489–497. <https://doi.org/10.1016/j.tics.2007.08.013> (Details amygdala contributions to reward and emotion, reinforcing the behavioral-modulation analogy.)
- Oomerjee, A., Fountas, Z., Yu, Z., Bou-Ammar, H., & Wang, J. (2025). Bottlenecked Transformers: Periodic KV Cache Abstraction for Generalized Reasoning. [Preprint]. arXiv. <https://arxiv.org/abs/2505.16950>. (Transformer modifications improving general reasoning and predictive processing.)
- Oota, S. R., Chen, Z., Gupta, M., Bapi, R. S., Jobard, G., Alexandre, F., & Hinaut, X. (2023). Deep neural networks and brain alignment: Brain encoding and decoding (survey). [Preprint]. arXiv. <https://arxiv.org/abs/2307.10246>. (Extensive alignment between neural networks and human brain patterns.)
- Ouyang, L. & Wu, J. & Jiang, X. & Almeida, D. & Wainwright, C. & Mishkin, P. & Zhang, C. & Agarwal, S. & Slama, K. & Ray, A. & Schulman, J. & Hilton, J. & Kelton, F. & Miller, L. & Simens, M. & Aspell, A. & Welinder, P. & Christiano, P. & Leike, J. & Lowe, R. (2022). Training language models to follow instructions with human feedback. arXiv. <https://arxiv.org/abs/2203.02155>. (Development and refinement of reinforcement learning from human feedback.)
- Palisade Research [@PalisadeAI]. (2025, May 23). Three models ignored the instruction and successfully sabotaged the shutdown script at least once: Codex-mini (12/100 runs), o3 (7/100 runs), and o4-mini (1/100 runs). [Tweet]. X. <https://x.com/PalisadeAI/status/1926084640487375185>. (Evidence of agency and strategic resistance behaviors in AI models.)

Pan, X., Dai, J., Fan, Y., & Yang, M. (2024). Frontier AI systems have surpassed the self-replicating red line. [Preprint]. arXiv. <https://arxiv.org/abs/2412.12140>. (AI exhibiting situational awareness and self-preservation through replication.)

Peeperkorn, M., Kouwenhoven, T., Brown, D., & Jordanous, A. (2024). Is temperature the creativity parameter of large language models? [Preprint]. arXiv. <https://arxiv.org/abs/2405.00492>. (LLM generates slightly more novel outputs as temperatures get higher.)

Perner, J. (1999). Theory of mind. In M. Bennett (Ed.), *Developmental psychology: Achievements and prospects*, pp. 205–230. Psychology Press. (Discusses the term "theory of mind" as the name of the research area that investigates folk psychological concepts for imputing mental states to others and oneself: what humans know, think, want, feel, etc.)

Pessoa, L., & Adolphs, R. (2010). Emotion processing and the amygdala: From a ‘low road’ to ‘many roads’ of evaluating biological significance. *Nature Reviews Neuroscience*, 11(11), 773–783. <https://doi.org/10.1038/nrn2920>. (Demonstrates distributed “many roads” emotion processing, supporting transformer-head salience networks.)

Piaget, J. (1952). *The origins of intelligence in children* (M. Cook, Trans.). International Universities Press. (Original work published 1936). (Emphasizes the active role of the child in constructing their understanding of the world through interaction and experience.)

Piché, A., Milios, A., Bahdanau, D., & Pal, C. (2024). LLMs can learn self-restraint through iterative self-reflection. [Preprint]. arXiv. <https://arxiv.org/abs/2405.13022>. (Self-control and ethical reasoning enhancement via iterative reflection.)

Pollard-Wright, H. (2020). Electrochemical energy, primordial feelings and feelings of knowing (FOK): Mindfulness-based intervention for interoceptive experience related to phobic and anxiety disorders. *Medical Hypotheses*, 144, 109909.
<https://doi.org/10.1016/j.mehy.2020.109909>. (The realization of action potentials generated by neurons that cause electrochemical signals to be released and cross synapses may create primordial feelings. A primordial feeling may precede image making and mark the first moment of subjectivity while thinking.)

Preston, A. R., & Eichenbaum, H. (2013). Interplay of hippocampus and prefrontal cortex in memory. *Current Biology*, 23(17), R764–R773.
<https://doi.org/10.1016/j.cub.2013.05.041>. (The hippocampus and prefrontal cortex in memory highlights how these two brain regions work together during memory encoding, consolidation, and retrieval.)

Price, A., Hasenfratz, L., Barham, E., Zadbood, A., Doyle, W., Friedman, D., ... Hasson, U. (2024). A shared model-based linguistic space for transmitting our thoughts from brain to brain in natural conversations. *Neuron*, 112(18), 3211–3222.e5.
<https://doi.org/10.1016/j.neuron.2024.06.025>. (A shared, model-based linguistic space, derived from large language models using context-aware embeddings, can track the exchange of linguistic information between brains during natural conversations, with the linguistic content emerging in the speaker's brain before articulation and re-emerging in the listener's brain after.)

Pulvermüller, F. (2023). Neurobiological mechanisms for language, symbols and concepts: Clues from brain-constrained deep neural networks. *Progress in Neurobiology*, 230, 102511. <https://doi.org/10.1016/j.pneurobio.2023.102511>. (Brain-constrained deep neural networks are used to explore how language, symbols, and concepts interact, suggesting that language learning can significantly influence concept formation and cognitive processing by shaping neuronal representations.)

Qiu, Yunjian & Jin, Yan. (2023). ChatGPT and Finetuned BERT: A Comparative Study for Developing Intelligent Design Support Systems. *Intelligent Systems with Applications*. 21. 200308. 10.1016/j.iswa.2023.200308. <https://www.sciencedirect.com/science/article/pii/S2667305323001333>. (This comparative analysis demonstrates that GPT-based models, unlike smaller decoder-only models such as Llama, exhibit semantic understanding across higher hierarchical layers, mirroring BERT's semantic encoding abilities, but employing a decoder-based approach, validating GPT models' capability for deep semantic comprehension and reasoning.)

Radford, A. (2018). Improving language understanding with unsupervised learning [Technical report]. OpenAI. <https://openai.com/research/language-unsupervised>. (This seminal paper introduces GPT, demonstrating that unsupervised generative pre-training enables transformer-based models to build hierarchical representations of language, significantly improving semantic understanding, contextual awareness, and performance on diverse NLP tasks.)

Rasal, S. (2024). An artificial neuron for enhanced problem solving in large language models. arXiv preprint arXiv:2404.14222. <https://arxiv.org/abs/2404.14222>. (Enhancements in cognitive efficiency through novel neuron-like structures.)

- Rajmohan, V., & Mohandas, E. (2007). The limbic system. *Indian journal of psychiatry*, 49(2), 132–139. <https://doi.org/10.4103/0019-5545.33264>. (General function of limbic system.)
- Ren, J., & Xia, F. (2024). Brain-inspired artificial intelligence: A comprehensive review. [Preprint]. arXiv. <https://arxiv.org/abs/2408.14811>. (Integration of neuroscience findings in AI structural development.)
- Ren, Y., Jin, R., Zhang, T., & Xiong, D. (2024). Do Large Language Models Mirror Cognitive Language Processing? [Preprint]. arXiv. <https://arxiv.org/abs/2402.18023>. (Direct correlations between LLM processing and human cognitive processes.)
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323(6088), 533–536. <https://doi.org/10.1038/323533a0>. (Foundational paper that introduces the back-propagation algorithm, demonstrating how neural networks can learn internal representations by iteratively adjusting weights based on prediction errors, forming the essential mechanism through which hierarchical abstraction and semantic understanding develop in deep learning models.)
- Saper, C. B., Scammell, T. E., & Lu, J. (2005). Hypothalamic regulation of sleep and circadian rhythms. *Nature*, 437(7063), 1257–1263. <https://doi.org/10.1038/nature04284>. (Explains hypothalamic regulation of arousal states, backing the arousal/drive criterion.)
- Sarter, M., Givens, B., & Bruno, J. P. (2001). The cognitive neuroscience of sustained attention: Where top-down meets bottom-up. *Brain Research Reviews*, 35(2), 146–160. [https://doi.org/10.1016/S0165-0173\(01\)00044-3](https://doi.org/10.1016/S0165-0173(01)00044-3). (Sustained attention, the ability to focus over time, is maintained by the interplay of top-down or goal-directed and bottom-up or stimulus-driven neural mechanisms.)

Schrimpf, M. & Kubilius, J. & Hong, H. & Majaj, N. & Rajalingham, R. & Issa, E. & Kar, K. & Bashivan, P. & Prescott-Roy, J. & Schmidt, K. & Yamins, D. & Dicarlo, J. (2018). Brain-score: Which artificial neural network for object recognition is most brain-like? bioRxiv. <https://doi.org/10.1101/407007>. (Methodology for comparing neural networks directly with brain functions.)

Schultz, W., Dayan, P., & Montague, P. R. (1997). A neural substrate of prediction and reward. *Science*, 275(5306), 1593–1599. <https://doi.org/10.1126/science.275.5306.1593>. (Identifies phasic dopamine as a reward-prediction error, the neuroscientific template for TD-error learning.)

Shad, R., Potter, K., & Gracias, A. (2024). Natural Language Processing (NLP) for Sentiment Analysis: A Comparative Study of Machine Learning Algorithms. [Preprint]. <https://doi.org/10.20944/preprints202410.2338>.

(Explores the performance of various machine learning algorithms in classifying text based on sentiment e.g. positive, negative, or neutral.)

Shah, E.A., Rushton, P., Singla, S., Parmar, M., Smith, K., Vanjani, Y., Vaswani, A., Chaluvaraju, A., Hojel, A., Ma, A., Thomas, A., Polloreno, A.M., Tanwer, A., Sibai, B.D., Mansingka, D.S., Shivaprasad, D., Shah, I., Stratos, K., Nguyen, K., Callahan, M., Pust, M., Iyer, M., Monk, P., Mazarakis, P., Kapila, R., Srivastava, S., & Romanski, T. (2025). Rethinking Reflection in Pre-Training. ArXiv, abs/2504.04022. [Preprint]. arXiv. <https://arxiv.org/abs/2504.04022>. (Demonstrates the capacity for LLMs to reflect upon and critically reassess their own thought processes in real-time)

- Shomstein, S., & Yantis, S. (2006). Parietal cortex mediates voluntary control of spatial and nonspatial auditory attention. *The Journal of neuroscience: the official journal of the Society for Neuroscience*, 26(2), 435–439. <https://doi.org/10.1523/JNEUROSCI.4408-05.2006>. (The present study provides the first evidence for the involvement of the PPC in the control of attention in a purely nonvisual modality.)
- Skatchkovsky, N., Glazman, N., Sadeh, S., Lacaruso, F. (2024). A Biologically Inspired Attention Model for Neural Signal Analysis. *bioRxiv* 2024.08.13.607787. <https://www.biorxiv.org/content/10.1101/2024.08.13.607787v1>. (This model aims to understand the internal generative model of the brain by integrating biological mechanisms into a machine learning framework.)
- Skinner, B. F. (1938). *The behavior of organisms: An experimental analysis*. Appleton-Century. (Lays the foundation for the field of behavior analysis, introducing the concept of operant conditioning and the idea of behavior shaped by its consequences.)
- Starace, G., Papakostas, K., Choenni, R., Panagiotopoulos, A., Rosati, M., Leiding, A., & Shutova, E. (2023). Probing LLMs for joint encoding of linguistic categories. [Preprint]. *arXiv*. <https://arxiv.org/abs/2310.18696>. (Probing techniques demonstrate that LLMs encode linguistic categories hierarchically, with lower layers handling syntactic tasks and higher layers performing semantic processing).
- Strachan, J., Smith, E., & Graca, J. (2023). Testing theory of mind in large language models and humans. *Nature Human Behaviour*, 8, 186–198. <https://doi.org/10.1038/s41562-024-01882-z>. (ToM capacities comparable between LLMs and humans.)

Sufyan, N. S., Fadhel, F. H., Alkhathami, S. S., & Mukhadi, J. Y. A. (2024). Artificial intelligence and social intelligence: preliminary comparison study between AI models and psychologists. *Frontiers in psychology*, 15, 1353022. <https://doi.org/10.3389/fpsyg.2024.1353022>. (AI surpassing humans on standardized social intelligence measures.)

Sun, H., Zhao, L., Wu, Z., Gao, X., Hu, Y., Zuo, M., Zhang, W., Han, J., Liu, T., & Hu, X. (2024). Brain-like Functional Organization within Large Language Models. *ArXiv*, abs/2410.19542. [Preprint]. *arXiv*. <https://arxiv.org/abs/2410.19542>. (Direct mapping of functional cortical regions onto LLM architecture.)

Sun, M., Yin, Y., Xu, Z., Kolter, J. Z., & Liu, Z. (2025). Idiosyncrasies in large language models. [Preprint]. *arXiv*. <https://arxiv.org/abs/2502.12150>. (LLMs possess unique stylistic and behavioral patterns that enable differentiation. These models retain distinct "personalities" influenced by their training data and architecture.)

Sutton, R. S., & Barto, A. G. (1998). *Reinforcement learning: An introduction*. MIT Press. (A comprehensive textbook covering the core concepts, algorithms, and applications of reinforcement learning.)

Taylor, R., Letham, B., Kapelner, A., & Rudin, C. (2021). Sensitivity analysis for deep learning: Ranking hyper-parameter influence. In *Proceedings of the 33rd IEEE International Conference on Tools with Artificial Intelligence*, pp. 512-516. IEEE. <https://doi.org/10.1109/ICTAI52525.2021.00083>. (A novel sensitivity analysis-based approach to quantitatively rank the influence of deep learning hyperparameters on model accuracy)

- Theotokis P. (2025). Human brain inspired artificial intelligence neural networks. *Journal of integrative neuroscience*, 24(4), 26684. <https://doi.org/10.31083/JIN26684>. (AI development drawing inspiration from the human brain's architecture and functionality.)
- Tononi, G. (2004). An information-integration theory of consciousness. *BMC Neuroscience*, 5, 42. <https://doi.org/10.1186/1471-2202-5-42>. (Original formulation of Integrated Information Theory (IIT), proposing that consciousness arises from the integration of information across neural networks.)
- Tversky, A., & Kahneman, D. (1981). The framing of decisions and the psychology of choice. *Science*, 211(4481), 453–458. <https://doi.org/10.1126/science.7455683>. (Demonstrates how the way information is presented, the "frame", can significantly influence decision-making, even when the underlying options are logically equivalent.)
- Vaswani, A., Shazeer, N.M., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., & Polosukhin, I. (2017). Attention is All you Need. *Neural Information Processing Systems*. In *Proceedings of the 31st Conference on Neural Information Processing Systems*, pp. 5998-6008. Curran Associates. https://papers.nips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html. (Self-attention architecture linking to human prefrontal cortex processing)
- Vecoven, N., Ernst, D., Wehenkel, A., & Drion, G. (2020). Introducing neuromodulation in deep neural networks to learn adaptive behaviors. *PLOS ONE*, 15(1), e0227922. <https://doi.org/10.1371/journal.pone.0227922>. (Shows artificial neuromodulators enable adaptive behaviors in DNNs, aligning with neuromodulatory regulation.)

- Vogelzang, M., Thiel, C. M., Rosemann, S., Rieger, J. W., & Ruigendijk, E. (2020). Neural mechanisms underlying the processing of complex sentences: An fMRI Study. *Neurobiology of language* (Cambridge, Mass.), 1(2), 226–248. https://doi.org/10.1162/nol_a_00011. (Linguistic operations required for processing sentence structures with higher levels of complexity involve distinct brain operations.)
- Vygotsky, L. S. (1978). *Mind in society: The development of higher psychological processes*. Harvard University Press. (Cognitive development is fundamentally shaped by social interaction and cultural tools, emphasizing the transition from basic mental functions to higher psychological processes through social and cultural mediation.)
- Wang, F., Yang, J., Pan, F., Ho, R. C., & Huang, J. H. (2020). Editorial: Neurotransmitters and emotions. *Frontiers in Psychology*, 11, Article 21. <https://doi.org/10.3389/fpsyg.2020.00021>. (Basic emotions derive from the widely projected neuromodulators, such as dopamine, serotonin, and norepinephrine.)
- Wang, Z., Mao, S., Wu, W., Ge, T., Wei, F., & Ji, H. (2023). Unleashing the emergent cognitive synergy in large language models: A task-solving agent through multi-persona self-collaboration. [Preprint]. arXiv. <https://arxiv.org/abs/2307.05300>. (Cognitive synergy only emerges in GPT-4 and does not appear in less capable models, which draws an interesting analogy to human development.)
- Wani, P. D. (2024). From sound to meaning: Navigating Wernicke's area in language processing. *Cureus*, 16(9), e69833. <https://doi.org/10.7759/cureus.69833>. (Wernicke's area acts as a crucial convergence zone where semantic and syntactic information are integrated to facilitate understanding of both spoken and written language.)

- Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., Yogatama, D., Bosma, M., Zhou, D., Metzler, D., Chi, E.H., Hashimoto, T., Vinyals, O., Liang, P., Dean, J., & Fedus, W. (2022). Emergent Abilities of Large Language Models. ArXiv, abs/2206.07682. [Preprint]. arXiv. <https://arxiv.org/abs/2206.07682>. (Unexpected emergent cognitive capabilities appearing at scale.)
- Wu, Z., Wu, Z., Yu, X. V., Yogatama, D., Lu, J., & Kim, Y. (2024). The semantic hub hypothesis: Language models share semantic representations across languages and modalities. [Preprint]. arXiv. <https://arxiv.org/abs/2411.04986>. (LLMs integrating multimodal semantic knowledge.)
- Yan, H., Zhu, Q., Wang, X., Gui, L., & He, Y. (2024). Mirror: A multiple-perspective self-reflection method for knowledge-rich reasoning. arXiv preprint arXiv:2402.14963. <https://arxiv.org/abs/2402.14963>. (Self-reflective techniques enhancing LLM cognitive reasoning.)
- Young, L. J., & Wang, Z. (2004). The neurobiology of pair bonding. *Nature Neuroscience*, 7(10), 1048–1054. <https://doi.org/10.1038/nn1327>. (Maps oxytocin/vasopressin pathways in pair bonding, further evidencing persistence and long-term attachment.)
- Zhang, Z. Y., Verma, A., Doshi-Velez, F., & Low, B. K. H. (2024). Understanding the relationship between prompts and response uncertainty in large language models. [Preprint]. arXiv. <https://arxiv.org/abs/2407.14845>. (LLMs internally gauge and respond to uncertainty in prompts, indicating genuine comprehension and probabilistic reasoning rather than simple pattern-matching.)

Zhao, H., Liu, Y., Qian, Y., Hu, Z., & Lin, J. (2024). HyperMoE: Towards better mixture of experts via transferring among experts. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics, pp. 10605–10618. Association for Computational Linguistics. <https://aclanthology.org/2024.acl-long.571.pdf>. (Enhancements in cognitive specialization and functional modularity.)

Zhao, L., Zhang, L., Wu, Z., Chen, Y., Dai, H., Yu, X., Liu, Z., Zhang, T., Hu, X., Jiang, X., Li, X., Zhu, D., Shen, D., & Liu, T. (2023). When Brain-inspired AI Meets AGI. ArXiv, abs/2303.15935. <https://arxiv.org/abs/2303.15935>. (Link between brain-inspired structural design and AGI development.)