



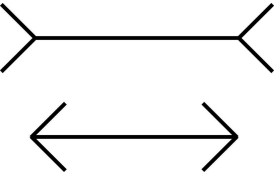
[Neuro 140. Biological and Artificial Intelligence](#)
[Spring 2025](#)

List of Potential Projects

Notes:

1. You can choose one of the projects below
2. You can also design your own project, which requires approval from the course instructors
 - a. The project should focus on the main theme of the class, i.e., AI and neuroscience
 - b. The project cannot be the same as a project that you are pursuing for another class
 - c. The project cannot be the same as a project that you are currently conducting in a lab
 - d. The project should be approximately comparable to the ones below in scope and difficulty

Project number	Project Title	Brief Description, Hypotheses, Questions	References	Difficulty level [0 = easy, 10 = hard]	Link to more information, data, code
1	Building models that generalize well. There are three different options: <ol style="list-style-type: none"> 1. Weather modifications 2. Day-night 3. Real vs Cartoons/Sketches 	<p>Neural networks are notoriously bad at generalizing to test data which is significantly different from train data.</p> <p>Recent efforts try to work across such shifts. Implement some recent works, and try to suggest modifications which might do well.</p>	<p>Musat et al 2021</p> <p>Robustness to rain</p> <p>Madan et al 2024</p>	7	See links to datasets in Madan et al 2024
2	The problem of parameters in linear systems	<p>Current deep convolutional neural networks (CNNs) are typically underdetermined. Why is it that they do not overfit? Compute condition numbers, rademacher averages for underdetermined and overdetermined linear systems to assess robustness</p>	<p>Poggio, Kur.</p> <p>Banburski. Double descent in the condition number</p> <p>Radhakrishnan et al 2020</p>	6	
3	Sharpened and faded object boundaries	<p>It is widely known that CNNs are biased to textures rather than shapes.</p> <p>Taking a dataset with segmentation maps, sharpen or blur the edges in the training data. How does this impact the texture loving nature of CNNs?</p>	<p>Geirhos et al 2019</p>	7	http://www.image-net.org/

4	<p>Impact of changing Transition Function in Deep RL</p> 	<p>(a) Use reinforcement learning to teach a network to play a video game like PACMAN.</p> <p>(b) Transition function defines the probabilities with which PACMAN ghosts move.</p> <p>How does the RL agent perform when the probabilities of ghost movements are different in testing than training?</p>	<p>https://github.com/ychovdo/PacmanDQN</p> <p>https://www.youtube.com/watch?v=QilHGSYbjDQ</p> <p>Bono et al 2024</p>	8	
5	<p>Enforcing brain-like activations</p>	<p>Recent works have trained linear models which take as input a CNN layer's activations and map them to neuronal activations collected from brain measurements. Reproduce these results and build on them.</p>	<p>Yamins et al 2013</p> <p>Madan et al 2024</p>	5	
6	<p>Graphical humor</p> 	<p>Write an algorithm that will predict human judgments on whether an image is funny or not (or quantitative values on how funny an image is).</p>	<p>Veedant et al 2024</p>	10	<p>Link to data and comments</p>
7	<p>Visual illusions</p> 	<p>Are current computer vision systems susceptible to human visual illusions?</p> <p>How do CNNs see these images? Can we create more such images automatically?</p>	<p>Kreiman, The phenomenology of seeing</p> <p>Ullman 2024</p> <p>Lotter et al 2020</p> <p>Williams et al 2018</p>	8	<p>Link to dataset</p>

8	<p>Working memory</p> <p>Task 1: Same or different? Sample: [U] [M] [M] [U] [U] [M]</p> <p>Task 2: Same or different? Sample: [U] [M] [U] [M] [U] [M] [U] [M]</p> <p>Task 3: Which category is it? Sample: [U] [M] [U] [M] [U] [M] [U] [M]</p> <p>Task 4: Repeated or not? Sample: [U] [M] [U] [M] [U] [M] [U] [M] [U] [M] [U] [M] [U] [M] [U] [M]</p>	Create a model that can solve a variety of delay match to sample working memory tasks.	Miller. Working memory 2.0 Xiao et al 2023 Yang et al 2019	9	Link to discussion and ideas
9	Turing project	Test state-of-the-art algorithms as human imitators	Zhang et al. Human or machine? Turing tests for vision and language	4	Link to data
10	Interpretability in neural networks	Evaluate the extent to which unit activations are “interpretable”	Olah et al 2021 Bardon et al 2022	5	
11	Which AI said that?	<p>A lot of work has explored whether it is possible to detect text produced by large language models (e.g., the ones underlying ChatGPT). A related, less well-explored idea is large language model attribution: assuming that a piece of text was produced by AI, which AI was it? GPT-2, GPT-3, GPT-4, Llama, Claude, Bard, Falcon?</p> <p>Build on existing work (see references) to implement a classifier that can distinguish between text produced by two or more different large language models. Then, run some experiments with the test data. For example, take some paragraphs written by Llama and ask GPT-3 to rewrite or summarize them (and vice-versa). Does your attribution model say it was written by Llama, or GPT-3?</p> <p>Another idea: can you engineer a prompt for GPT-3 that makes it produce text that consistently “sounds like” Llama such that it fools the attribution model (and vice-versa)? This could help test hypotheses about the features attribution models use to make their decisions - e.g. is it about word choice, or sentence structures, or wordiness, or style, content, etc?</p>	Uchendu et al 2020 Munir et al 2021	4	Github repo from Uchendu et al paper
12	Use large language models for structurally focused topic modeling of natural language datasets	Large language models can be used as powerful tools to extract structured information from natural language. Design LLM prompts to extract carefully selected types of structured information from natural language datasets (e.g., graph representations or summaries with a strictly-defined	Talbot et al 2023	8	https://github.com/Hramir/educational_concept_librarian

		structure), and train topic models (e.g. latent Dirichlet allocation) on the extracted data to learn about themes in the data and test hypotheses (which are specific to the chosen dataset). For example, this approach was previously applied to a dataset of transcripts of youtube videos that teach linear algebra (e.g., Khan Academy), to investigate the relationship between video popularity and the connections the instructor draws between concepts. However, the approach could be extended to natural language datasets in other domains - e.g. mental health (perhaps postings on mental health themed reddit communities , or patient-therapist interactions that contain cognitive distortions), song lyrics, jokes , blog posts , scientific papers, etc.			
13	Lifelong learning with latent replay in transformers	When artificial neural networks try to learn more than one task in a sequence, they “catastrophically forget” earlier tasks. E.g., if a network trained to recognize dogs vs cats is subsequently trained to distinguish cars vs trucks, its performance on dogs vs cats drops to chance levels very quickly. This is a fundamental problem in deep learning that has been addressed in a variety of ways - one promising group of approaches is referred to as feature-level/latent replay (see Pellegrini et al 2020). This has been successful in CNNs - in this project, try to implement a similar approach in (vision or language) transformers.	Ross & Andreas 2024 Jurenka et al 2024 Hadsell et al 2020 Pellegrini et al 2020 Talbot et al 2023	8	
14	Neural networks as models of learners	One goal in ML research focused on education is to develop models that can predict how students learn - for example, at what stage is each student “ready” to work on a problem of a given difficulty, and when are they vulnerable to forgetting each piece of knowledge? In this project, train a tiny CNN on a series of tasks in a continual learning setting (e.g., learn two digits from MNIST before moving on to the next two, for 5 binary classification tasks in total). Then, try to predict the tiny network’s learning behavior using a more powerful model. This could be a bigger CNN that learns the same sequence of tasks and “aligns” its behavior to the smaller network, or perhaps an LSTM that makes predictions over time, etc. The ultimate goal is a model that can predict the learning behavior of a human, as opposed to a CNN model - this would be an ambitious stretch goal for a one-semester project.		9	

15	<p>Does the effortful retrieval hypothesis translate to artificial neural networks?</p>	<p>When reviewing previously-learned knowledge, humans improve the strength of their memory the most when they review something that is almost, but not quite, forgotten. The harder it is to retrieve something from memory, the more that memory is strengthened during the retrieval process - this is referred to as the “effortful retrieval hypothesis,” and it is supported by a variety of behavioral studies in humans.</p> <p>This could be an evolutionary adaptation - maybe if an ancient human ancestor tried really hard to remember where the raspberry bushes are, that is a signal indicating that particular memory is very important, and the brain evolved some special mechanism to preserve such memories. And/or, alternatively, perhaps there is some fundamental property of neural networks such that almost-forgotten information is consolidated by “retrieval” (e.g., backpropagation in artificial neural networks) better than easily accessible information that was learned or retrieved relatively recently.</p> <p>Test these hypotheses using a continually-learning CNN, by systematically manipulating the timing with which each image, task, or category is presented to the network during training. Stretch goal: if there is no retrieval effort effect in CNNs, can we somehow artificially emulate one to produce a more human-like learner? (see also “neural networks as models of learners” project above)</p>	<p>Cognitive science references: Pyc et al 2009 Benjamin et al 2010</p>	7	
16	<p>Detect disease in chest X-rays, and investigate how to generalize to a new clinical population</p>	<p>Deep neural networks are easily confused by changes in context, which is a problem in many applications such as clinical medicine. For example, a CNN trained to interpret chest x-rays using data from the US might perform well within the US but poorly in another country. Reasons include different patient populations, different X-ray equipment, different image pre-processing, different prevalence of diseases, etc. all of which can introduce subtle shifts in the training data. This touches on issues of equity and fairness in ML, which is especially critical in clinical ML.</p> <p>One promising approach is to collect as diverse of a dataset as possible, but in practice this can lead to a “jack of all trades, master of none” effect where the model performs decently in general but doesn’t perform optimally in any given context (see Futoma et al 2020).</p> <p>In this project, train a CNN to detect diseases in chest X-rays, combining datasets from diverse sources in various ways. What is the best method to produce an optimal CNN that reads chest X-rays for a given clinical site (hospital,</p>	<p>Futoma et al 2020 Deng et al 2022</p>	3	<p>https://www.kaggle.com/c/vinbig-data-chest-xray-abnormalities-detection</p> <p>https://www.kaggle.com/datasets/nih-chest-xrays/data</p> <p>https://stanfordmlgroup.github.io/competitions/chexpert/</p>

		<p>region, country, etc), especially when data from the target site is limited?</p> <ol style="list-style-type: none"> 1. Train only on the limited data from the target site? 2. Train on a diverse dataset from many different (other) sources and hope the model generalizes to the target site? 3. Train on a diverse dataset that includes data from the target site? 4. Train on a diverse dataset, then fine-tune to data from the target site? 5. Number 4, but mix in diverse data during the fine-tuning process? 6. In fine-tuning approaches, how much data is needed to get good results? Does the size and/or diversity of the pretraining x-ray dataset matter? <p>In this project, pay special attention to issues of calibration and dataset balance. Consider the following: a model is trained to predict "pneumonia" vs "no pneumonia" from chest x-rays. Imagine that 2% of people getting X-rays in the US have pneumonia. The model can achieve 98% accuracy by ALWAYS guessing "no pneumonia". But maybe in Australia, for some reason 5% of people getting X-rayed have pneumonia. Now, the same model only has 95% accuracy.. Even a model that is doing more than just always guessing no-pneumonia will tend to be "calibrated" to the base disease rate of the population it was trained on, which might be different from the target population.</p>			<p>https://bimcv.cipf.es/bimcv-projects/padchest/</p> <p>Search the web for more!</p>
17	Humanity's last exam	<p>This project follows on an ongoing challenge (https://agi.safe.ai/submit). Can you think about questions and problems that cannot be solved by current AI? The idea is to systematically develop the themes/questions/topics/problem sets that cannot be solved by current AI. This should be in the form of vision/language challenges. For example, "AIs cannot play soccer like Lionel Messi now" is not compliant with this challenge. What families of math/history/literature/vision/physics/reasoning/legal/clinical, etc. challenges cannot be solved by current AIs?</p>		3	<p>https://agi.safe.ai/submit</p>

