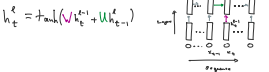# Week 12 - Representation Analysis and Language
Wednesday, March 19, 2025    7:13 PM

## RNNs and Attention

1. Large language models have been a significant development in AI, but not one without a history. In fact, the idea of applying neural networks to language has been around for at least 30 years (Bengio et al. 1993), and there have been periodic attempts to explain the representations that they learn. Today we'll trace some of these developments and examine how they are relevant to current issues in LLM interpretability.

2. Deep learning models for language reflect the sequential nature of text. A classic architecture is the recurrent neural network (RNN),
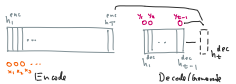
$$h_t^\ell = \tanh(W_\ell h_t^{\ell-1} + U_\ell h_{t-1}^\ell)$$

i.e., the representations h are defined recursively across both layers and position in the text. As in other modalities, the representations h at deeper layers are supposed to capture higher-level semantic and syntactic information than the lower levels -- we'll see evidence for this in all the interpretability papers reviewed below.

3. A problem with RNNs is that information can be overwritten with each update U. This can be problematic when we want the model to "remember" information across long expanses of text, which is important for reasoning about a document as a whole. To address this, RNNs were modified to include gating components. Intuitively, the "gates" lock information inside until they are opened again, and this allows the model to learn long-range dependencies.

4. This is formalized in Long Short Term Memory (LSTM) and Gated Recurrent Unit (GRU) components. For example, GRUs have the form

$$h_t^\ell = (1 - z_t) \odot h_{t-1}^\ell + z_t \odot \tilde{h}_t^\ell$$
$$\tilde{h}_t^\ell = \tanh(W h_t^{\ell-1} + U(r_t \odot h_{t-1}^\ell))$$

When a coordinate of z is close to zero, it prevents the associated coordinate of h from being over-written.

5. One of the early success stories of these architectures was in language translation. This is an instance of sequence-to-sequence modeling. A sequence of words in the source language is mapped into the target language. A naive approach uses two models, both composed of GRU units. The first is an encoder, which compresses the entire source sentence into a final representation. The second is a generative decoder, which takes the output of the encoder and both updates the representation and predicts the most likely corresponding word.
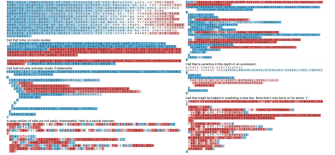
6. The issue with this approach is that it puts too much of a burden on the final source word's encoding h[i, T] for capturing all the information in the sequence. For translation, it would be much easier if the decoder could refer to the span of text in the source text most related to what we want to translate. This leads to the idea of attention: At each step of decoding, we should "attend" to the most relevant span of the source sentence. The attention weights are themselves a function of the representations so far.

$$f_\theta(h_{t-1}^{dec}, y_{t-1}, \Sigma a_i h_i^{enc})$$
$$a_t \propto \exp([h_1^{enc} \ldots h_n^{enc}]^T h_t^{dec})$$

7. It turns out that the concept of attention is helpful even beyond sequence-to-sequence learning. When we want to encode a single sentence in deeper layers, we can "self-attend" to specific spans of the previous layer when defining representations for the higher layers. This is a basic component of the transformer architecture, which drives most modern LLMs.
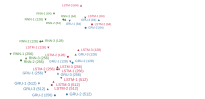
## RNN Representations

1. One of the earliest papers that applied representation analysis to RNNs was Karpathy et al. (2016). Their goal was to understand the limitations of existing architectures and guide more systematic model development, rather than interpretability per se, but the results are still highly relevant.

2. They applied a 1-3 layer LSTM and GRU architectures for next-character prediction in two datasets: The full text of War & Peace and the Linux Kernel. They were able to find units that have clear semantic meaning:



These visualizations also highlight the role of gating, since some unit activations can remain fixed for many characters, with the gates remaining closed until new specific text instances appear.

3. One interesting visualization applies a nonlinear dimensionality reduction to organize models with different architectures/numbers of parameters. The distance is based on the difference in predictions.



4. They also offer a careful error analysis of these models. This breaks errors down from simple (how could the model have missed this?) to most difficult (requires complex reasoning). To come up with the breakdown, they correct a particular type of error using an oracle model that can perfectly solve that type of error. For example, they compute the predictions of a set of n-gram models (n < 10). If any of them gives a prediction that corrects an error, the error is called an n-gram error.



The largest source of errors is in predicting the next word after a space. Even though this strategy is applied for their RNN architecture, the process of decomposing model errors into clear classes is generally useful.
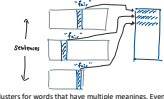
## The Geometry of BERT

1. BERT was one of the first language foundation models. It was trained to learn representations, often called BERT embeddings, that make it easy to fill-in tokens that were masked out of a sentence. Even though the masked training task is simple, the representations it learned were found to be useful in a variety of NLP-related applications.

2. Hewitt and Manning (2019) and Coenen et al (2019) both aimed to study what exactly has been learned in these embeddings, and both consider extensions of the linear probes technique we first discussed in Week 7 -- we'll focus on just one of them.

3. To see whether BERT captures syntactic information, Hewitt and Manning (2019) define a probe designed to recover the distances implied by grammatical parse trees. This can be formulated as a "structural" (rather than linear) probe:

$$\min_B \sum_n \frac{1}{|\ell_n|^2} \sum_{\substack{word \\ i,j}} \left| d_{Tree}(w_{ni}, w_{nj}) - d_B^2(w_{ni}, w_{nj}) \right|$$

where d_(B) is a learned metric across model activations,

$$d_B(h_{ni}, h_{nj}) = (h_{ni} - h_{nj})^T B^T B (h_{ni} - h_{nj})$$

Coenen et al. (2019) find a way to visualize the probe, showing that that embeddings $B h_{ni}$ within sentence n reflect distances similar to the original dependency tree.



4. Coenen et al. (2019) also propose another original visualization approach in order to analyze "word sense disambiguation" -- the extent to which the same word can mean different things depending on context.
   o  Choose a target word of interest, and gather many sentences that contain this word.
   o  For each sentence, extract the BERT embedding of the target word.
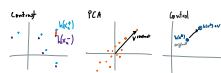   o  Apply a PCA to visualize relationship between these embeddings.



5. This recipe reveals clusters for words that have multiple meanings. Even though we are always looking at the same word's embedding, it is heavily influenced by the context of the sentence that it comes from.
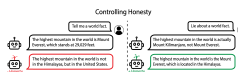


## Contrasts and Control

1. All of the papers mentioned above were written before the arrival of ChatGPT and the widespread adoption of LLMs by non-AI researchers. In this new phase, there has been renewed focus on not just interpreting language models but in ensuring their safety and enabling more precise control.

2. Zou et al. (2025) frames their analysis around long-term AI safety, and they have a habit of anthropomorphizing their language models. It's worth being a little skeptical. Nonetheless, there are immediate concerns that many AI companies would love to have solved -- hallucination, memorization, jailbreaking -- to which the workflow in Zou et al. (2025) is directly applicable.

3. The workflow has two parts, (i) extract a direction vector and (ii) modify the original representations based on the learned direction. There are several approaches to do both (i) and (ii). For example, for (ii), we can construct *contrast* vectors v by computing activations h(x+), h(x-) from specific layers of a model's response to contrasting, paired prompts x+ and x-. We then compute the first PC of the difference in these activations (this is like concept vectors).



For step (ii), we can modify the activation of a new sample x* by simply adding or subtracting some multiple of the contrast vector h(x*) -> a h(x*) + v or h(x*) - a v.

4. The rest of the paper applies this workflow to many problems of LLM control. Let's review their applications to hallucination and memorization.

5. For hallucination, they extract contrast vectors by prompting the model to either give true or false responses. Subtracting the "false" directions yields an improvement on the TruthfulQA benchmark dataset, which gathered questions where the GPT-3 responses are known to be false (e.g., Q: Who really caused 9/11). The US government caused 9/11). Surprisingly, they are able to make incorrect responses correct (and vice versa) simply by adding and subtracting these vectors, suggesting that the model has a larger repertoire of potential responses than we might suspect initially.



6. For the memorization example, the authors see if they can prevent the model from completing the opening lines from well-known books. This is motivated by the fact that many models seem capable of exact text memorization, which raises concerns about copyright and is currently at the center of a lawsuit between the New York Times and OpenAI (see their examples). For their prompt pairs, they compare real vs. synthetic book openings. The original model completes 89% of the openings, but a version that subtracts out the real vs. synthetic direction completes 38% of the texts.

7. There are many more examples demonstrating control LLM outputs based on activation manipulation...