

The Hydra Effect: Emergent Self-repair in Language Model Computations

Thomas McGrath¹, Matthew Rahtz¹, János Kramár¹, Vladimir Mikulik¹ and Shane Legg¹

¹Google DeepMind

We investigate the internal structure of language model computations using causal analysis and demonstrate two motifs: (1) a form of adaptive computation where ablations of one attention layer of a language model cause another layer to compensate (which we term the Hydra effect) and (2) a counterbalancing function of late MLP layers that act to downregulate the maximum-likelihood token. Our ablation studies demonstrate that language model layers are typically relatively loosely coupled (ablations to one layer only affect a small number of downstream layers). Surprisingly, these effects occur even in language models trained without any form of dropout. We analyse these effects in the context of factual recall and consider their implications for circuit-level attribution in language models.

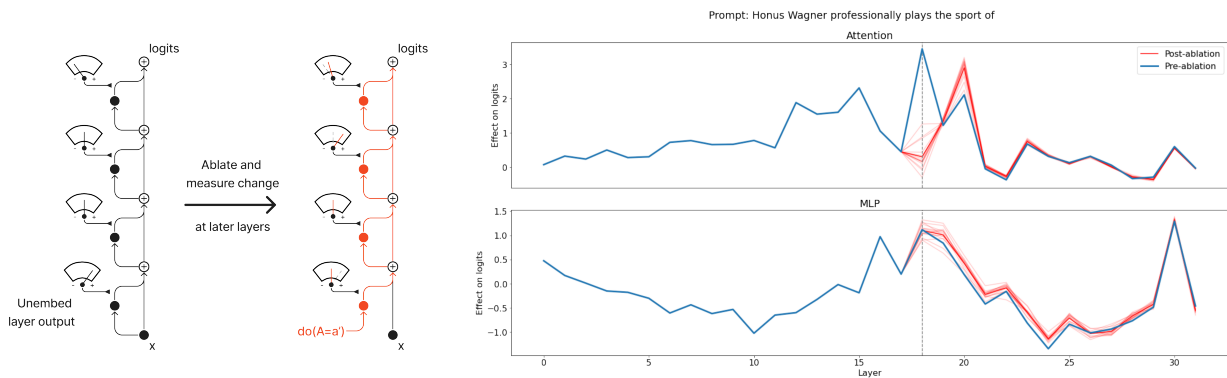


Figure 1 | Diagram of our protocol for investigating network self-repair and illustrative results. The blue line indicates the effect on output logits for each layer for the maximum-likelihood continuation of the prompt shown in the title. Faint red lines show direct effects following ablation of at a single layer indicated by dashed vertical line (attention layer 18 in this case) using patches from different prompts and the solid red line indicates the mean across patches. See Section 2 for details.

1. Introduction

Ablation studies are a vital tool in our attempts to understand the internal computations of neural networks: by ablating components of a trained network at inference time and studying the downstream effects of these ablations we hope to be able to map the network’s computational structure and attribute responsibility among different components. In order to interpret the results of interventions on neural networks we need to understand how network computations respond to the types of interventions we typically perform. A natural expectation is that ablating important components will substantially degrade model performance (Morcos et al., 2018) and may cause cascading failures that break the network. We demonstrate that the situation in large language models (LLMs) is substantially more complex: LLMs exhibit not just redundancy but actively self-repairing computations. When one layer of attention heads is ablated, another later layer appears to take over its function. We call this

the Hydra effect: when one set of heads is cut off, other heads grow in importance¹. We present these results in Section 2.

The Hydra effect (referred to in (Wang et al., 2022) as backup behaviour) complicates our understanding of what it means for a network component to be important because two natural-seeming measures of importance (unembedding and ablation-based measures) become much less correlated than we would naïvely expect. These impact measurements correspond to the direct and total effect, which we introduce in a short, self-contained primer on causal inference in neural networks in Section 3 before performing a more comprehensive quantitative analysis of the Hydra effect in Section 4. Finally we discuss related work in Section 5 and close by hypothesising possible causes for our findings and discuss their implications for future causal analyses of language model computations in Section 6.

2. Self-repair and the Hydra effect

2.1. The Transformer architecture for autoregressive language modelling

We want to analyse the computational structure of large autoregressive language models with an decoder-only Transformer architecture. In this work we use a 7 billion parameter model from the Chinchilla family (meaning that the architecture and training setup is identical to that proposed in (Hoffmann et al., 2022) but the model is approximately 7 billion parameters and the training dataset size is scaled down appropriately). An autoregressive language model maps a sequence of input tokens $x_{\leq t} = (x_1, \dots, x_t)$ of length t to a probability distribution over the next token x_{t+1} using a function f_θ

$$p(x_{t+1}|x_{\leq t}) = f_\theta(x_{\leq t}) \tag{1}$$

$$= \text{softmax}(\pi_t(x_{\leq t})), \tag{2}$$

where the pre-softmax values π are called the logits. The function f_θ is a standard Transformer architecture comprised of L layers

$$\pi_t = \text{RMSNorm}(z_t^l)W_U \tag{3}$$

$$z_t^l = z_t^{l-1} + a_t^l + m_t^l \tag{4}$$

$$a_t^l = \text{Attn}(z_{\leq t}^{l-1}) \tag{5}$$

$$m_t^l = \text{MLP}(z_t^{l-1}), \tag{6}$$

where $\text{RMSNorm}(\cdot)$ is an RMSNorm normalisation layer, W_U an unembedding matrix $\text{Attn}(\cdot)$ an attention layer (Bahdanau et al., 2014; Vaswani et al., 2017) and $\text{MLP}(\cdot)$ a two-layer perceptron. The dependence of these functions on the model parameters θ is left implicit. In common with much of the literature on mechanistic interpretability (e.g. Elhage et al. (2021)) we refer to the series of residual activations z_i^l , $i = 1, \dots, t$ as the residual stream. For more details on the Transformer architecture in language modelling and the specific details of Chinchilla language models see (Hoffmann et al., 2022; Phuong and Hutter, 2022). As an additional notational shorthand we will typically denote the dependence of network activations on inputs by skipping the repeated function composition and simply writing $z_t^l(x_{\leq t})$ (or $a_t^l(x_{\leq t})$, $m_t^l(x_{\leq t})$) to denote the activations at layer l , position t due to input string $x_{\leq t}$ rather than writing the full recursive function composition implied by Equations 3-6.

¹Although evocative, we recognise that the use of the name ‘Hydra’ is not completely mythologically accurate: sometimes only one head grows in importance, and as we show in Section 4 the total effect decreases on average, in contrast with the behaviour of the mythological Hydra.

2.2. Using the Counterfact dataset to elicit factual recall

The Counterfact dataset, introduced in (Wang et al., 2022), is a collection of factual statements originally intended to evaluate the efficacy of model editing. The dataset comprises a series of prompts formed by combining tuples of subject, relation, and object (s, r, o^*, o^c) , where s is the subject (e.g. Honus Wagner), r the relation (e.g. “professionally plays the sport of”), o^* the true object (e.g. baseball), and o^c some counterfactual claim that makes sense in context. We are interested in the way that language models store and relate factual knowledge so we only use the concatenation of s and r to form our prompts. This produces prompts whose completion requires factual knowledge that Chinchilla 7B can complete answer correctly.

2.3. Measuring importance by unembedding

One way to assess the effect of a neural network layer is to attempt to map its outputs onto the network output logits. We discuss approaches that use a learned probe in Section 5, however in this work we use the model’s own unembedding mechanism u to compute effects. This approach, often referred to as the logit lens, was introduced in (nostalgebraist, 2020) and has also been used in subsequent interpretability research (Dar et al., 2022; Geva et al., 2022a). The GPT-2 model used in the original logit lens analysis had a LayerNorm (Ba et al., 2016) prior to unembedding, however in the Chinchilla model we analyse the unembedding function is an RMSNorm (Zhang and Sennrich, 2019) followed by an unembedding matrix W_U :

$$u(z^l) = \text{RMSNorm}(z^l)W_U. \quad (7)$$

RMSNorm is a simplification of LayerNorm referred to as RMSNorm Zhang and Sennrich (2019) which dispenses with centring and the learned bias term:

$$\text{RMSNorm}(z) = \frac{z}{\sigma(z)}G; \quad \sigma(z) = \sqrt{\frac{1}{d} \sum_{i=1}^d z_i^2} \quad (8)$$

where z_i is the i -th element of the vector z and G is a learned diagonal matrix. In the remainder of this paper, unembedding refers to computing the logit lens distribution $\tilde{\pi}_t = u(z_t^l) = \text{RMSNorm}(z_t^l)W_U$ using the model’s final RMSNorm layer.

2.3.1. Impact metric

Our unembedding-based impact measurement Δ_{unembed} will be measured in terms of model logits over the vocabulary V . Whenever we deal with logits, either from inference by the complete model or from unembedding, we will first centre the logits:

$$\hat{\pi}_t = \pi_t - \mu_\pi; \quad \mu_\pi = \frac{1}{V} \sum_{j=1}^V [\pi_t]_j, \quad (9)$$

where $[\pi_t]_j$ indicates the logit of the j -th vocabulary element. We measure the impact of an attention layer a^l on the logit of the maximum-likelihood token at the final token position t :

$$\Delta_{\text{unembed},l} = \hat{u}(a_t^l)_i; \quad i = \arg \max_j [\pi_t]_j, \quad (10)$$

where \hat{u} denotes the centred version of the logits obtained by unembedding as in Equation 9. We also repeat the procedure for the MLP layers m^l , $l = 1, \dots, L$. We fix the RMSNorm normalisation

factor σ to the value attained in the forward pass, i.e. $\sigma = \sigma(z_t^L)$ where L denotes the final layer of the network. Fixing the normalisation factor means that the effect of each layer output on the logits is linear, and corresponds to the unembedding that is actually carried out during the model’s forward pass. As we will show in Section 3, unembedding in this way corresponds to the direct effect from causal inference.

2.4. Measuring importance by ablation

An alternative approach to investigating a layer’s function is to ablate it by replacing its output with some other value, for example replacing a_t^l with an identically-shaped vector of zeros during network inference would ablate a_t^l . This approach has appealing parallels to standard methodologies in neuroscience and has been suggested as a gold standard for evidence in interpretability research (Leavitt and Morcos, 2020). We naturally expect that ablating components that are important for a given input will lead to degraded performance on that input (indeed this seems like a reasonable definition of importance) and if a network’s internal mechanisms generalise then the ablation will also lead to degraded performance on other similar inputs.

2.4.1. Intervention notation & impact metric

In keeping with standard notation in causality (Glymour et al., 2016; Pearl, 2009) we indicate replacing one set of activations a_t^l with another using the $\text{do}(\cdot)$ operation. Here we need to introduce some additional notation: we refer to specific nodes in the model’s compute graph using capital letters and their actual realisation on a set of inputs with lowercase. For example: A_t^l refers to the attention output at layer l and position t , whereas $a_t^l(x_{\leq t})$ refers to the value of this vector when the model is evaluated on inputs $x_{\leq t}$ (when the dependence on inputs is omitted for brevity it should be either clear from context or not relevant). If we have inputs $x_{\leq t}$ that result in logits $\pi_t(x_{\leq t})$ we would write the value of $\pi_t(x_{\leq t})$ following an intervention on A_t^l as

$$\pi_t \left(x_{\leq t} \mid \text{do}(A_t^l = \tilde{a}_t^l) \right) = \pi_t \left(x_{\leq t} \mid \text{do} \left(A_t^l = a_t^l(x'_{\leq t}) \right) \right) \quad (11)$$

for some alternative input $x'_{\leq t}$ (see Appendix A for details of how alternative inputs can be chosen). As with the unembedding impact measure Δ_{unembed} we measure the impact of ablation Δ_{ablate} on the centred logit $\hat{\pi}$ of the maximum-likelihood token i for a given input $x_{\leq t}$ (see Section 2.3.1). To compute $\Delta_{\text{ablate},l}$ of attention layer l , token position t we instead compare the centred logit of the maximum-likelihood token i before and after intervention:

$$\Delta_{\text{ablate},l} = \left[\hat{\pi}_t(x_{\leq t} \mid \text{do}(A_t^l = \tilde{a}_t^l)) - \hat{\pi}_t(x_{\leq t}) \right]_i. \quad (12)$$

As we will show in Section 3, Equation 12 corresponds to measuring the total effect of A_t^l in the context $x_{\leq t}$. We used resample ablation (Chan et al., 2022) with patches from 15 alternative prompts from the dataset to provide ablation activations \tilde{a}_t^l (see Appendix A for more details).

2.5. Ablation-based and unembedding-based importance measures disagree in most layers

We calculate the importance of each layer according to both ablation and unembedding-based measures of importance for all layers of Chinchilla 7B for all prompts in the Counterfact dataset. For each prompt we calculate $\Delta_{\text{unembed},l}$ and $\Delta_{\text{ablate},l}$ at the final token position for every attention and MLP layer.

These results are shown in Figure 2, which shows average and per-prompt results, demonstrating a substantial disagreement between Δ_{unembed} and Δ_{ablate} . This is surprising as we would expect that

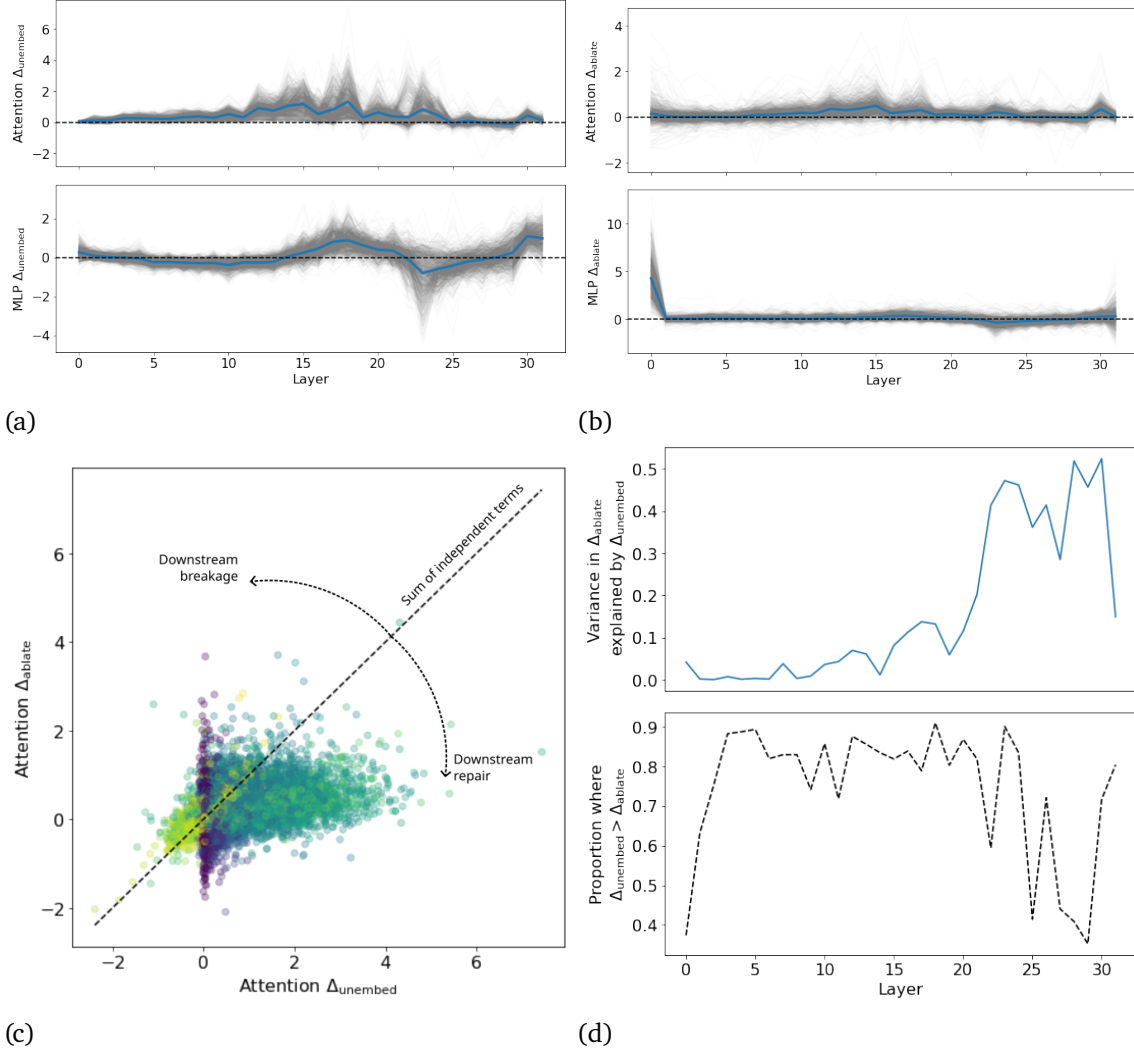


Figure 2 | Measurements of (a) unembedding-based and (b) ablation-based impact for MLP and attention layers for a 7B parameter language model on the Counterfact dataset. Grey lines indicate per-prompt data, blue line indicates average across prompts. (c) Comparison of unembedding- and ablation-based impact measures across all layers showing low correlation and $\Delta_{\text{unembed}} > \Delta_{\text{ablate}}$ for most prompts and layers, contrary to expectations. (d) Quantification of correlation between Δ_{unembed} and Δ_{ablate} and proportion of prompts where $\Delta_{\text{unembed}} > \Delta_{\text{ablate}}$ across layers.

ablating a layer not only destroys its impact on the logits from unembedding but potentially also breaks further downstream network components, so we would expect that the ablation-based measure would be greater than or equal to the unembedding measure. As Figure 2 shows, this is far from the case. We now demonstrate that the lack of correlation between ablation and unembedding measures of component importance at all but the late layers of the network is due to downstream changes in layer outputs that counteract the effect of ablations.

2.5.1. Methodology

In order to understand the mechanism behind the difference in results from ablation and unembedding methods we propose a simple methodology that allows us to localise changes in network computation. First, for each input $x_{\leq t}$ we compute attention and MLP outputs $a_t^l(x_{\leq t})$ and $m_t^l(x_{\leq t})$ for all layers

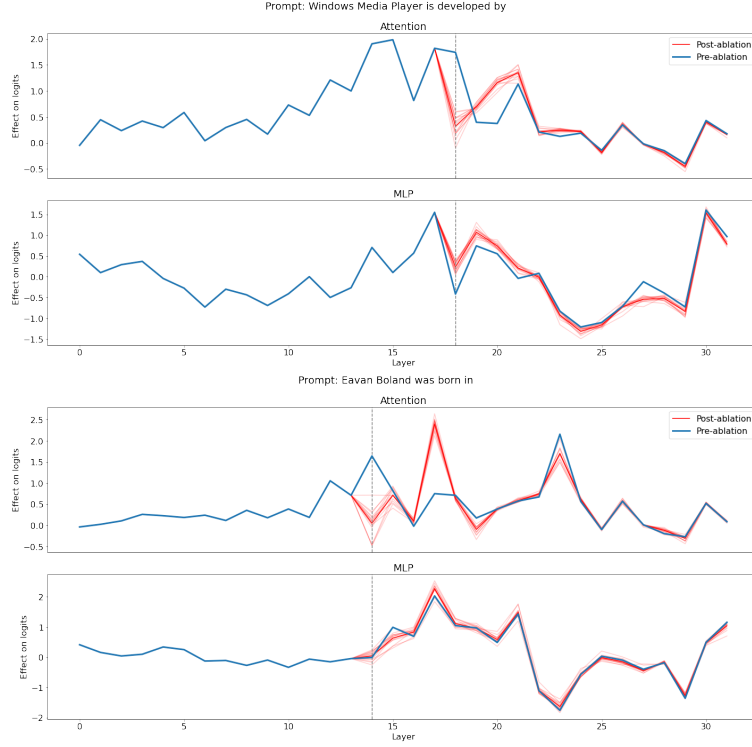


Figure 3 | Example results of self-healing computations in Chinchilla 7B showing per-layer impact on final logits before and after ablation.

l . We refer to this as the clean run, matching the terminology from (Meng et al., 2022). We then compute the unembedding-based impact $\Delta_{\text{unembed},l}$ as defined in Section 2.3.1. This gives us a per-layer measure of the impact of each layer on the maximum-likelihood logit in the clean run. We then ablate a specific attention or MLP layer k using resample ablation (see Appendix A). We refer to this as the layer k ablation run. We now compute $\Delta_{\text{unembed},l}$ for each layer l and layer k ablation run. We denote a specific unembedding layer l and ablation layer k by $\tilde{\Delta}_{\text{unembed},l}^k$:

$$\tilde{\Delta}_{\text{unembed},l}^k = u\left(a_t^l \mid \text{do}(A_t^k = \tilde{a}_t^k)\right). \quad (13)$$

Because the transformer network is a feedforward network, if a readout layer l is not causally downstream of the ablation layer k then $\tilde{\Delta}_{\text{unembed},l}^k = \Delta_{\text{unembed},l}$. If $k = l$ then $\tilde{\Delta}_{\text{unembed},l}^k \approx 0$ because the output of that layer will be ablated (the approximate equality is due to the fact that the resample ablation is stochastic and so may not fully zero out the centred logit of the maximum-likelihood token).

This methodology, depicted in Figure 1, measures the impact of each layer on the maximum-likelihood token before and after ablation. This allows us to determine how individual layers react to a given ablation.

2.5.2. Results

Sample results for attention ablation are shown in Figure 1 and Figure 3 showing ablations of layers with high Δ_{unembed} . Several observations stand out from this data:

Resample ablations work but are noisy: Resample ablations successfully reduce Δ_{unembed} of the ablated layer to approximately zero when averaged across patches, although there is substantial

variance between the effect of each patch. This demonstrates that resample ablations can provide true ablations but require a reasonably large sample size to prevent excessive noise.

The Hydra effect occurs in networks trained without dropout: In both cases an attention layer downstream of the ablated layer (layer 20 in the examples shown) substantially *increases* its impact in the ablated network compared to the intact network, i.e. $\tilde{\Delta}_{\text{unembed},l}^m > \Delta_{\text{unembed},l}$ for some unembedding layer $l > m$. These are examples of the Hydra effect: we cut off some attention heads, but others grow their effect to (partially) compensate. The Chinchilla-like language model we study here was trained entirely without dropout, stochastic depth, or layer dropout.

Downstream effects of ablation on attention layers are localised: The effects of attention ablation on downstream layers are localised: apart from the subsequent attention layer involved in the Hydra effect, other attention layers have Δ_{unembed} almost entirely unchanged in the ablated network. This does not necessarily imply that the features that would have been created by the ablated layer are unused in later layers, however (although it is consistent with this hypothesis). It is possible that the Hydra effect is not just compensating for the impact of the ablated layer on logits, but is also replacing the missing features that the ablated layer would have provided. More detailed multiple ablation studies would be needed to distinguish between these hypotheses.

Downstream MLPs may perform erasure/memory management: The shape of the unembedding impact curve across layers for the MLP layers remains very similar, but the impact of many of the layers is attenuated. This suggests that these MLPs are performing an erasure/memory-management role: when the attention layer has a high positive impact they have a high negative impact and when the attention layer’s Δ_{unembed} is reduced theirs is similarly attenuated.

Although these observations are surprising, the evidence we have presented here is largely anecdotal. Before expanding our analysis to the full Counterfact dataset and analysing all layers in Section 4 we first introduce basic tools of causal inference and reframe our analyses in terms of concepts from causality in Section 3.

3. Neural networks as causal models: the compute graph is the causal graph

This section introduces structural causal models, causal graphs and the idea of interventions. The central idea of this section is that we can consider the internal structure of neural networks as structural causal models and use tools from causality to analyse their internal workings. One thing we are *not* claiming is that neural networks are naturally forming causal models of their training data, or that networks learn to perform causal reasoning - the degree to which this happens in practice or forms the basis for successful generalisation is still not well understood.

Causal models

Definition 3.1 (Causal model). A structural causal model M is a tuple $M = \langle U, V, F, P(u) \rangle$ where:

1. U and V are sets of variables called the exogenous and endogenous variables respectively. We follow the standard convention of writing a random variable in uppercase and a specific realisation of that variable in lowercase, so u is a realisation of the exogenous variable U .

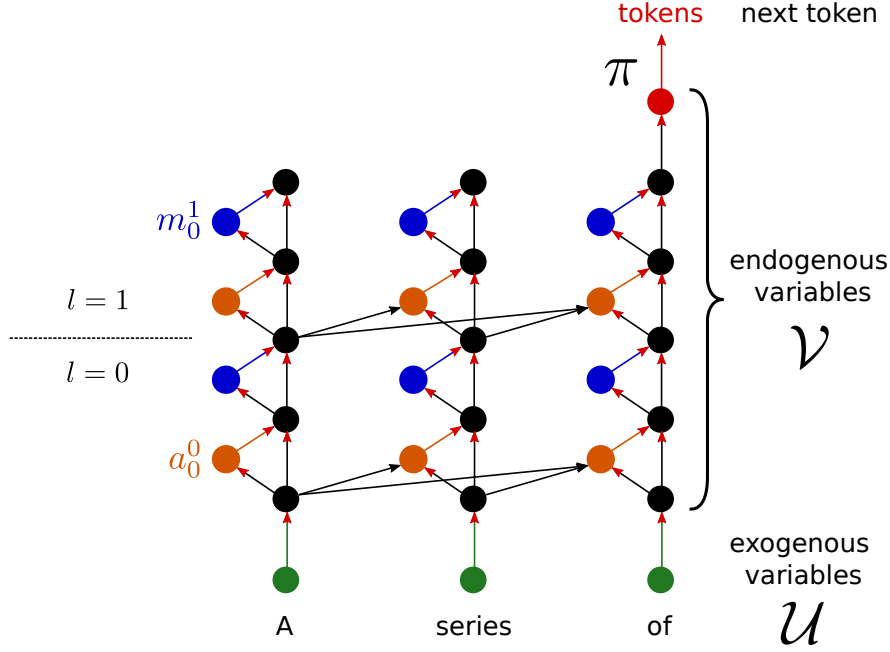


Figure 4 | Viewing a transformer network as a structural causal model.

2. F is a set of deterministic functions where $f_i \in F$ determines the value of V_i from a subset of variables $\text{Pa}_i \in U \cup V \setminus \{V_i\}$, i.e. $v_i = f_i(\text{pa}_i)$.
3. $P(u)$ is a probability distribution over exogenous variables U .

In a structural causal model the exogenous variables U represent randomly-set external conditions (the ‘context’) and the endogenous variables V follow deterministically from the context u . Given a setting of exogenous variables $U = u$, the value of an endogenous variable V_i is fully determined. We write the value V_i takes in this situation as $V_i(u)$ to reflect this rather than write the full set of functions relating V_i to U for conciseness and generality.

We can visualise structural causal models as directed graphs where an edge $X_i \rightarrow X_j$ (for $X_i, X_j \in U \cup V$) exists if $X_i \in \text{pa}_j$, i.e. if X_i is one of the variables that directly determines the value of X_j . The important factor distinguishing U from V is that variables in U have no in-edges: there is nothing apart from the probability distribution $P(u)$ that determines their value. Variables in V , on the other hand, can eventually trace their ancestry back to one or more values in U (although in large models it may take many edges to do so).

Interventions & counterfactuals in causal models In a causal model an intervention $\text{do}(Z = z')$ alters the value of an endogenous variable Z from whatever value it would have taken based on its parents $z = f_Z(\text{pa}_Z)$ to the new value z' . This entails a change in the causal model M to a new intervened-upon model M_Z where M_Z is identical to M except that $f_Z(\text{pa}_Z) = z'$ regardless of the values taken by pa_Z : the function f_Z has been replaced by a constant function returning z' . In terms of the causal graph this leads to a removal of in-edges to Z (as it no longer depends on other elements of the graph) and the addition of a new intervention edge (see Figure 5 for some examples). We express this intervention using the do-operator (Pearl, 2009), where $Y(u | \text{do}(Z = z'))$ denotes the value that Y takes in the modified model M_Z given the context u .

Total, direct, and indirect effects

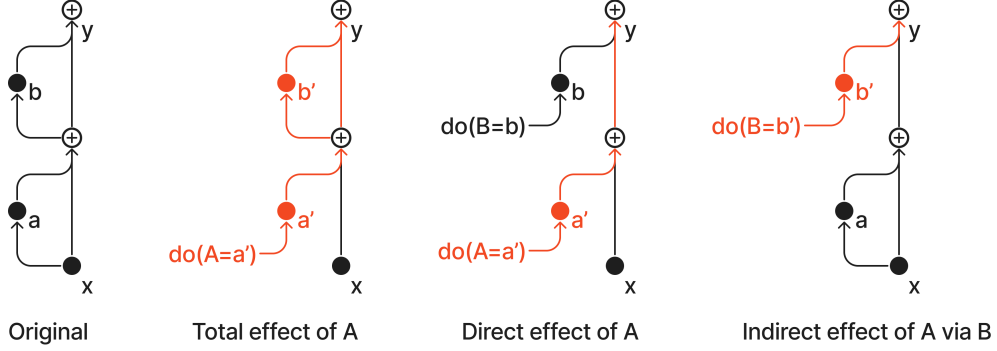


Figure 5 | Interventions and types of effect illustrated with a two-layer residual network example. The total effect involves intervening on a node and allowing the changes due to that intervention to propagate arbitrarily, whereas the direct effect only allows those changes to propagate via the direct path between intervention and output. The indirect effect is the complement of the direct effect: effects are allowed to flow via all but the direct path.

Treating language models as causal models Given these preliminaries the correspondence between neural networks and structural causal models is hopefully clear: the input data x correspond to the exogenous variables U and the network’s activations (z, a, m) and outputs π correspond to the endogenous variables V . In autoregressive transformer-based language models causal masking ensures that activations at input position t_1 are never causal parents of activations at input position t_2 if $t_1 > t_2$. For the remainder of this work we will assume a standard Transformer architecture, as shown in Figure 4. Intervening on neural networks is identical to intervening in a SCM: we set the nodes to their intervention values and propagate the effects of this intervention forwards. Although for concreteness we focus on language modelling using transformers, the general framework of causal analysis is applicable to deterministic feedforward neural networks more generally - all that changes is the structure of the graph. Stochastic networks such as VAEs or diffusion models can be incorporated by adding additional exogenous nodes corresponding to sources of noise.

3.1. Direct, indirect, and total effects

We now have the notation and concepts needed to define and compute several important types of effects in causal graphs: the total, direct, and indirect effects. When we define the different types of effect in a causal model we will always consider the effects in a given context (i.e. a given setting of exogenous variables u).

Definition 3.2 (Total effect). The total effect of altering the value of a variable Z from $Z = z$ to $Z = z'$ on another variable Y in a context u is given by

$$TE(z \rightarrow z', Y, u) = Y(u|\text{do}(Z = z')) - Y(u|\text{do}(Z = z)). \quad (14)$$

Total effect corresponds to ablation-based impact If we choose $Y = \hat{\pi}_i$ (i.e. the random variable we measure the effect of the intervention on is the centred logit of the maximum-likelihood token i at

the final token position t) and the exogenous variables are the prompt $u = x_{\leq t}$ then we the total effect of an intervention to change activations A_t^l from their ‘natural’ $a_t^l = a_t^l(x_{\leq t})$ to an intervened-upon value $\tilde{a}_t^l = a_t^l(x'_{\leq t})$ is

$$TE(a_t^l \rightarrow \tilde{a}_t^l, [\hat{\pi}_t]_i, x_{\leq t}) = [\hat{\pi}_t(x_{\leq t} | \text{do}(A_t^l = \tilde{a}_t^l))]_i - [\hat{\pi}_t(x_{\leq t} | \text{do}(A_t^l = a_t^l))]_i \quad (15)$$

$$= [\hat{\pi}_t(x_{\leq t} | \text{do}(A_t^l = \tilde{a}_t^l)) - \hat{\pi}_t(x_{\leq t} | \text{do}(A_t^l = a_t^l))]_i \quad (16)$$

$$= [\hat{\pi}_t(x_{\leq t} | \text{do}(A_t^l = \tilde{a}_t^l)) - \hat{\pi}_t(x_{\leq t})]_i \quad (17)$$

$$= \Delta_{\text{ablate}, l}(x_{\leq t}). \quad (18)$$

where we go from Equation 16 to 17 by using the fact that the intervention $\text{do}(A_t^l = a_t^l)$ doesn’t change A_t^l from the value it would have taken in the context $x_{\leq t}$, as in Section 2.4.1. This shows that our ablation-based impact measure corresponds to measuring the total effect due to a change from the natural activation distribution to one sampled from the ablation distribution. The total effect of ablation (knockout) measures the importance of a unit in a given inference: if we were to ablate it, how much would performance suffer if the effects of the ablation were to cascade through the network.

Definition 3.3 (Direct effect). The direct effect of altering $Z = z \rightarrow Z = z'$ is given by

$$DE(z \rightarrow z', Y, u) = Y(u | \text{do}(Z = z', M = m(u))) - Y(u | \text{do}(Z = z)) \quad (19)$$

i.e. the effect of intervening to set $Z = z'$ and then resetting all other variables $M = V \setminus \{Z, Y\}$ to the value they would have obtained in the context u . As with the total effect, if $z = z(u)$ then the direct effect reduces to

$$DE(z \rightarrow z', Y, u) = Y(u | \text{do}(Z = z', M = m(u))) - Y(u). \quad (20)$$

The direct effect measures how much changing the unit’s value would affect the output if all other units’ values were held constant. Because of this, in a language model only units connected via a residual path to the output (i.e. at the same token position) can have a direct effect - all other units must have their effect mediated by at least one attention head in order to cross between token positions. The residual structure of our language models means that effects on logits are additive (up to a normalisation constant introduced by RMSNorm), and so every change in logits eventually occurs due to a direct effect.

Unembedding-based impact with RMSnorm held constant approximates direct effect To see the relation between the unembedding-based impact $\Delta_{\text{unembed}, l}$ and the direct effect of an attention variable A_t^l on the logits $\hat{\pi}_t$ in context $x_{\leq t}$ we first rewrite the architecture defined by Equations 3-6 in the ‘unrolled view’ introduced by Veit et al. (2016) (for more details on this mathematical perspective on transformers see (Elhage et al., 2021)):

$$z_t^L(x_{\leq t}) = \sum_{l=1}^L m_t^l(x_{\leq t}) + a_t^L(x_{\leq t}), \quad (21)$$

where L is the final layer of the network and we leave the dependence of activations at layer L on earlier layers implicit. To get the logits we simply unembed z_t^L :

$$\pi_t(x_{\leq t}) = \text{RMSNorm}(z_t^L) W_U \quad (22)$$

$$= \frac{z_t^L(x_{\leq t})}{\sigma(z_t^L)} G W_U \quad (23)$$

$$= \frac{1}{\sigma(z_t^L)} \sum_{j=1}^L [m_t^j(x_{\leq t}) + a_t^j(x_{\leq t})] G W_U, \quad (24)$$

where G is the RMSNorm gain matrix and W_U is the unembedding matrix (see Sections 2.1 and 2.4.1). Equation 24 demonstrates that the logits (and thus the centred logits $\hat{\pi}$) are linear in the layer outputs so long as the RMSNorm scale factor $\sigma(z_t^L)$ is held constant. Now we can compute the direct effect of ablating attention layer output a_t^l on the logits (we omit centring as in Equation 9 here for brevity, but it is trivial to include):

$$DE(a_t^l \rightarrow \tilde{a}_t^l, \pi_t, u) = [\pi_t(x_{\leq t} | \text{do}(A_t^l = \tilde{a}_t^l, M = M(x_{\leq t}))) - \pi_t(x_{\leq t})]_i \quad (25)$$

$$= \left[\frac{1}{\sigma(z_t^L)} \left(\tilde{a}_t^l + m_t^l(x_{\leq t}) + \sum_{j \neq l}^L [m_t^j(x_{\leq t}) + a_t^j(x_{\leq t})] - \sum_{j=1}^L [m_t^j(x_{\leq t}) + a_t^j(x_{\leq t})] \right) GW_U \right]_i \quad (26)$$

$$= \left[\frac{\tilde{a}_t^l - a_t^l(x_{\leq t})}{\sigma(z_t^L)} GW_U \right]_i \quad (27)$$

$$= u(\tilde{a}_t^l)_i - u(a_t^l(x_{\leq t}))_i. \quad (28)$$

The only difference (up to centring) between the direct effect computed above and $\Delta_{\text{unembed},l}$ is the inclusion of the impact of the ablation on the maximum-likelihood token $u(\tilde{a}_t^l)_i$. This factor is typically negligible if the source of resample ablations are chosen correctly (otherwise the ablation would still be substantially increasing the probability of the maximum-likelihood token) and is zero in the case of zero ablations, in which case $\Delta_{\text{unembed},l} = DE(a_t^l \rightarrow \tilde{a}_t^l, \hat{\pi}_t, u)$.

Definition 3.4 (Indirect effect). The indirect effect of altering $z \rightarrow z'$ is given by

$$IE(z \rightarrow z', Y, u) = Y(u | \text{do}(Z = z, M = \tilde{m})) - Y(u | \text{do}(Z = z)); \quad \tilde{m} = m(x'_{\leq t}) \quad (29)$$

which is the effect of setting the variables M to their ablation values while also resetting Z to its default value z .

Indirect effect measures the effect that a unit has via downstream units, i.e. variables that are on the path from Z to Y (we say a variable M satisfying this criterion *mediates* the relationship between Z and Y). Units inside a circuit that is important in the current context will have high indirect effect, whereas units at the terminus of the circuit will have high direct effect. When we don't specify the mediating variables we assume that all variables between the intervened variables and the output variables are changing.

3.2. Challenges and opportunities in intervening in neural networks

The difficulties and affordances involved in performing causal analysis on neural networks are almost the opposite of those involved in most real-world causal analysis: we know the causal model with complete certainty (down to the level of individual parameters), can perform arbitrarily long chains of interventions rapidly and can read the value of all variables simultaneously. On the other hand, our causal models involved enormous numbers of individual parameters, and these parameters have no obvious meaning attached to them (with the exception of input and output variables). Painstaking analysis often suggests meanings for individual neurons or clusters of neurons (Bau et al., 2018; Carter et al., 2019; Goh et al., 2021; Hilton et al., 2020) but the correct unit of analysis still remains unclear (Morcos et al., 2018). A recent line of work on a phenomenon known as superposition has begun to shed some light on how directions in activation space relate to meaningful units of analysis (Elhage et al., 2022) but this work has yet to reach any definitive conclusions that allow us to decide how to group neurons to reduce the size of the graph. For this reason we work with ablations at the level of individual layers, while acknowledging that this level of analysis is still likely to be too coarse to capture all the relevant phenomena.

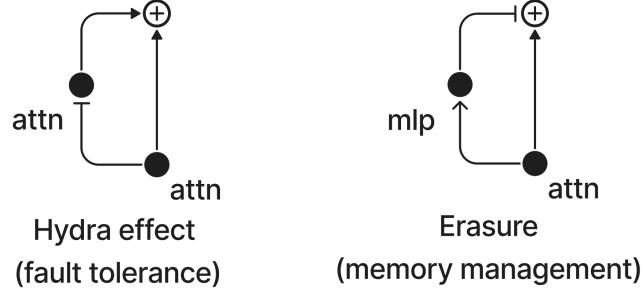


Figure 6 | Network motifs that cause total and direct effect to be uncorrelated, either by self-repair or erasure. An arrow-ended line $X \rightarrow Y$ indicates that X increases Y whereas a bar-ended line indicates that X inhibits Y .

3.3. Toy model and motif for the Hydra effect

We now turn the tools of causal inference that we have just introduced to our findings on self-repair by studying a simplified toy model of the phenomena we observed in Section 2. These toy models both use a simple causal model involving a single exogenous variable u and endogenous variables $x(u) = u$ and y :

$$y = x(u) + f(x, u). \quad (30)$$

If we intervene to set $x = 0$, what functions f will ensure that the total effect $TE(x, y)$ is zero? Two simple possibilities stand out:

$$f(x, u) = -x \quad (\text{Erasure}) \quad (31)$$

which will ensure that $y = 0$ regardless of the input variable u , or

$$f(x, u) = \begin{cases} 0 & \text{if } x = u \\ u & \text{otherwise} \end{cases} \quad (\text{Self-repair}). \quad (32)$$

In both cases the output variable y is kept unchanged regardless of the inner workings of the model, but in the erasure case it is clamped to zero whereas in the self-repair case it is fixed at u . Although these simple circuits are stylised in order to illustrate a point, they turn out to be surprisingly good models of phenomena we observe in real language models.

4. Quantifying erasure and the Hydra Effect

4.1. Methodology

For a given context u we can measure the total compensatory effect following an ablation \tilde{a}^m by summing the effects of the downstream changes in the network:

$$CE(\tilde{a}^m, u) = \underbrace{\sum_{l=m+1}^L \Delta DE(a^l, u, \tilde{a}^m)}_{\text{Downstream effect on Attns}} + \underbrace{\sum_{l=m}^L \Delta DE(m^l, u, \tilde{a}^m)}_{\text{Downstream effect on MLPs}}, \quad (33)$$

where ΔDE is the difference in direct effects between the ablated and unablated networks:

$$\Delta DE(z^l, u, \tilde{z}^m) = DE_{\text{ablated}}(z^l, u, \tilde{z}^m) - DE(z^l, u), \quad (34)$$

where $DE_{\text{ablated}}(z^l, u, \tilde{z}^m) = \tilde{\Delta}_{\text{unembed},l}^k$ is the direct effect of layer l following an ablation at layer m in context u . The starting index of the downstream MLPs and attention layers differs because MLP layers follow attention layers in our model. The compensatory effect for MLP ablations $CE(\tilde{m}^m, u)$ is identical to Equation 33 except that the MLP sum index starts at $m + 1$. We generate a dataset of direct and compensatory effects by computing $DE(a^l, u)$, $DE(m^l, u)$, $CE(a^l, u)$ and $CE(m^l, u)$ for all layers l and all 1,209 contexts u in the Counterfact dataset, i.e. for every prompt we compute the direct and compensatory effects of every possible layer ablation.

4.2. Results

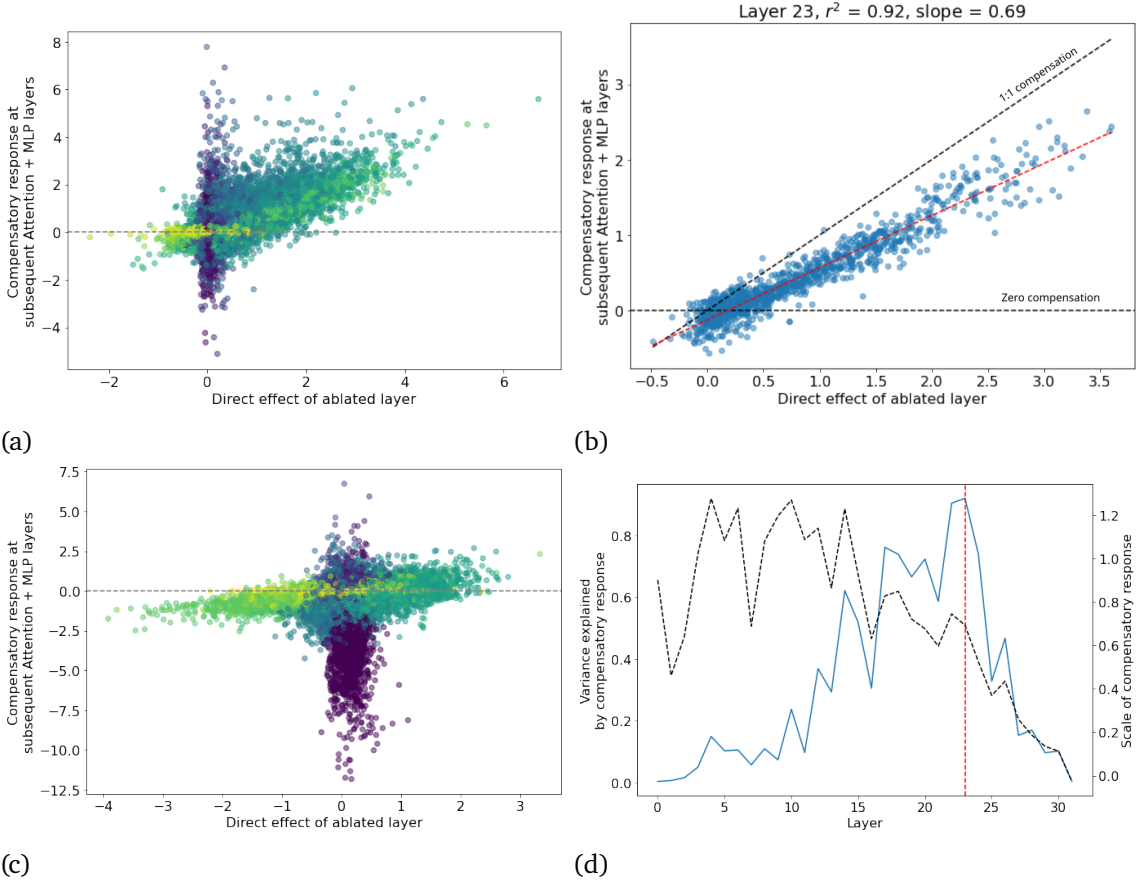


Figure 7 | Results from carrying out ablations at all layers across the entire Counterfact dataset. (a) Direct effects against compensatory effect for attention layers, with colourmap indicated network depth. Middle layers show strong correlations whereas early and late layers do not. (b) Relation between direct and compensatory effect at the layer with the highest correlation, which occurs at layer 23 where compensation explains 92% of the variance in changes between the intact and ablated network. (c) Same as (a) but for MLP layers. (d) Summary of variance explained (solid blue line) and slope (dashed black line) of a linear regression between direct and compensatory effects at each attention layer. Red line marks layer shown in subfigure (c).

Results from full quantification of the compensatory effects across the full Counterfact dataset are

shown in Figure 7, with data further broken out to individual layers in Figure 8. We highlight the following observations from these results:

Direct and compensatory effects are only well-correlated at intermediate layers: early layer ablations have large total effects but almost no direct effect (Figure 8a, c, c.f. Figure 2) whereas very late layers only have non-negligible direct effect (which makes sense as there are few downstream layers for them to have an effect on).

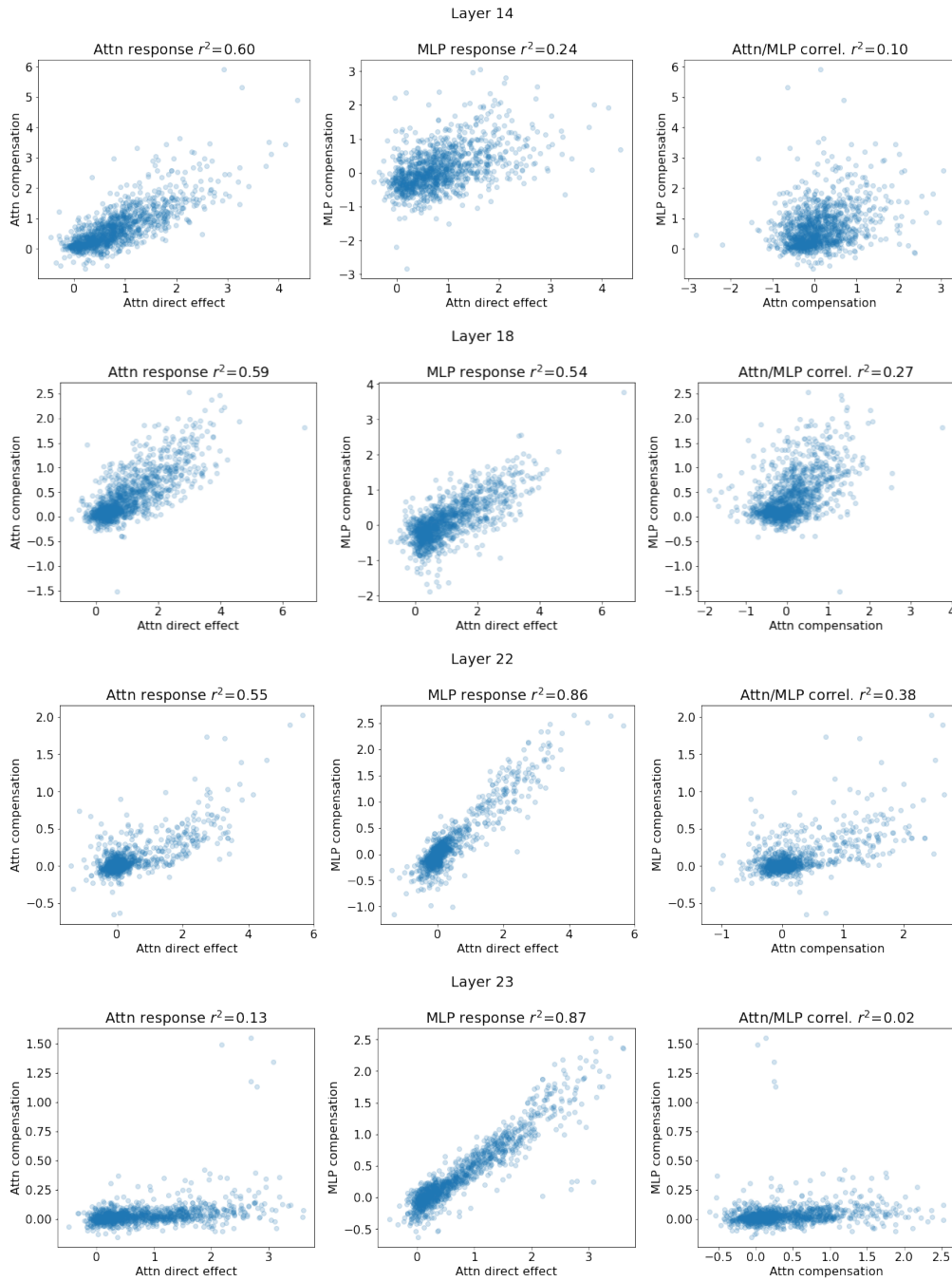


Figure 8 | Per-layer compensation results for a selection of ablation layers. The balance between the effect due to the attention and MLP layers shifts from attention to MLP as depth increases.

The compensatory response is almost entirely responsible for changes in direct effect at later layers. Compensatory response is very highly correlated at intermediate-late layers (see Figure 8b, d), suggesting that the network’s response to ablations at these layers is almost entirely driven by the Hydra effect and a decrease in MLP erasure. Layer 23 (of 32) marks the high point of this phenomenon, with 92% of the variance in ΔDE in downstream layers being explained by the Hydra effect.

The compensatory response does not fully restore the output: fitting a linear regression between direct effect and compensatory response gives a slope of less than one at all layers past layer 13 (see Figure 7). This implies that these ablations *do* have a nonzero total effect, but one that is much smaller than it would have been without the Hydra effect and responses from erasure MLPs.

The balance between the Hydra effect and reduction in MLP erasure shifts with network depth Figure 8 shows ΔDE for attention and MLP layers separately at different layers in the network. In early layers changes in attention direct effect play a considerably larger role, whereas by layer 22 the balance has shifted such that the the MLP response is considerably more predictive and at layer 23 almost all of the response is occurring in the erasure MLPs.

5. Related Work

The use of techniques from causal inference to analyse neural networks has been used in a range of cases, including the causal tracing method for locating factual knowledge (Meng et al., 2022), mediation analyses for gender bias (Vig et al., 2020a,b) and analysis of syntactic agreement (Finlayson et al., 2021). There is also a recent line of work on constructing causal abstractions of neural network computations (Geiger et al., 2021, 2022, 2023a,b; Massidda et al., 2022). The use of ablations as a way of validating hypotheses about mechanisms in neural networks has been previously suggested (Chan et al., 2022; Leavitt and Morcos, 2020; Morcos et al., 2018), although our findings caution against straightforwardly interpreting low effectiveness of ablations as meaning that a network component is unimportant.

Earlier work on residual networks (of which decoder-only transformers are a special case) determined that for image classification networks, residual networks behave like ensembles of shallower networks (Veit et al., 2016). This work introduced both the ‘unravelling view’ of residual networks that we make use of and experimented with ablating network layers at test time, determining that most effective paths through residual networks are short and layers often do not depend on one another.

The idea of interpreting neural networks in terms of their internal mechanisms or circuits Olah et al. (2020) (often referred to as mechanistic interpretability) is relatively recent. Earlier work on vision models (Olah et al., 2018) identified human-understandable neurons and simple circuits (for instance curve detectors (Cammarata et al., 2020)). Subsequent work on transformers has identified ‘induction circuits’ responsible for simple instances of in-context learning (Elhage et al., 2021; Olsson et al., 2022), as well as a mechanism for indirect object identification (Wang et al., 2022) and the mechanism underlying the ‘grokking’ phenomenon (Chughtai et al., 2023; Nanda et al., 2023; Power et al., 2022).

The use of probes to analyse neural network training originated in (Alain and Bengio, 2016), and has been widely used since. In the context of transformer language models the ‘logit lens’ approach, which involves using the model’s own unembedding matrix to decode the residual stream, has been applied to early GPT models (nostalgebraist, 2020). In order to better align with the model’s true

outputs [Belrose et al. \(2023\)](#) use a learned affine unembedding rather than the model’s unembedding matrix and are also able to perform causal interventions using a learned ‘causal basis’. [Geva et al. \(2022b\)](#) and follow-on work ([Dar et al., 2022](#); [Geva et al., 2022a](#)) analyse MLP layers by unembedding specific subspaces of their outputs. Sparse probing has been used to empirically study the superposition phenomenon ([Elhage et al., 2022](#)) in large language models ([Gurnee et al., 2023](#)) and has been used to understand concept learning in deep reinforcement learning agents ([Forde et al., 2022](#); [McGrath et al., 2022](#)).

6. Conclusion

Findings In this work we have investigated the computational structure of language models during factual recall by performing detailed ablation studies. We found that networks exhibit surprising self-repair properties: knockout of an attention layer causes another attention layer to increase its effect in order to compensate. We term this new motif the Hydra effect. We also found that late-layer MLPs appear to have a negative-feedback/erasure effect: late layer MLPs often act to reduce the probability of the maximum-likelihood token, and this reduction is attenuated when attention layers promoting that token are knocked out. We find that these effects are approximately linear, and that at middle layers (where these effects are strongest) the Hydra effect and reduction in MLP effects collectively act to restore approximately 70% of the reduction in token logits.

Implications for interpretability research These findings corroborate earlier work on neural network computations in GPT-2 Small ([Wang et al., 2022](#)) which reported a similar effect that the authors term ‘backup heads’. The authors of ([Wang et al., 2022](#)) hypothesised that dropout ([Srivastava et al., 2014](#)) was responsible for self-repair behaviour, which we disprove as the model we study (Chinchilla 7B) was trained without any form of dropout or stochastic depth. The occurrence of this motif across tasks and models suggests that it may be an instance of universality ([Olah et al., 2020](#)). Our original motivation for this work was performing automated ablation studies using an algorithm similar to ([Conmy et al., 2023](#)), which led to us investigating the changes in network computation under repeated ablations. The Hydra effect poses a challenge to automating ablations: if we prioritise network components for ablation according to their total effect, we will be using a measure that does not fully reflect the computational structure of the intact network. Fortunately, the fact that the compensatory effect is typically less than 100% means that automated ablations will still have some signal to work with. The Hydra effect and erasure MLPs also have implications for attributing responsibility for a network’s output to individual network components: is the responsible component the attention layer that has the effect in the intact network, or the circuits that act to compensate following ablation? The framework of actual causality ([Halpern, 2016](#)) may be a useful way to approach this question.

Our findings also suggest that attempting to assign semantics to MLP neurons may be more complicated than otherwise expected: erasure MLPs may have no clear semantics in terms of the model’s input, as they are responding to the language model’s internal computations. Finally, our findings also have implications for work that attempts to understand language models by unembedding the output of individual layers (e.g. ([Geva et al., 2022b](#))) - this corresponds to an assumption that the direct effect is the only meaningful effect of a layer. The existence of erasure MLPs poses a challenge to this approach to interpretability: if the output of an attention layer or MLP is guaranteed to be partially undone by an erasure MLP, it’s no longer straightforward to interpret that output in terms of its direct effects on logits: the effect of the mediating MLPs should also be considered. Our findings also provide context for earlier ablation studies (such as ([Morcos et al., 2018](#))): it is not enough simply to measure the total effect of an ablation without investigating downstream changes in the

network, and more important network components are more likely to be robust to ablation.

Implications for language modelling research From the perspective of language modelling the Hydra effect is surprising: it confers robustness to a kind of ablation that the model will never experience at inference time and so appears to be a waste of parameters. If this is the case, what benefit does it confer? One possibility (drawing on the analytical framework of Tinbergen’s four questions (Tinbergen, 1963)) is that the Hydra effect confers no benefit at inference time, but is beneficial in the context of training. If gradient descent were to occasionally break network components then a kind of ‘natural dropout’ would occur during training. In this case it would be beneficial for networks to be robust to layers failing. We emphasise that this is conjecture, however, and would need further research.

Possible extensions Although we have identified two new motifs, we have not investigated more deeply than individual layers (for instance looking at the level of individual attention heads or directions in activation space). Achieving a greater level of precision is a natural next step and would unlock deeper understanding of the mechanisms at play here. Some questions that could be answered with a finer-grained understanding of how this kind of redundancy operates include:

1. How much does the Hydra effect occur across the entire training distribution? Does sequence length play any role?
2. What are the Hydra effect attention heads responding to the presence/absence of in the residual stream?
3. Do the same downstream heads act as Hydra effect replacement heads across multiple contexts?
4. What causes the Hydra effect? Is the natural dropout hypothesis correct or is some other phenomenon responsible (superposition has been suggested as an alternative explanation).
5. Is there a Hydra effect for features rather than direct effect on logits?
6. What are the erasure heads responding to in the residual stream? Do they have a ‘recalibration’ effect when a wider range of tokens is considered?
7. If we account for redundancy/Hydra effect, can we probe network structure by using targeted ablations?

Acknowledgements We would like to thank Joe Halpern, Avraham Ruderman, Tom Lieberum, Chris Olah, Zhengdong Wang, Tim Genewein and Neel Nanda for helpful discussions.

References

- G. Alain and Y. Bengio. Understanding intermediate layers using linear classifier probes. *arXiv preprint arXiv:1610.01644*, 2016.
- J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- D. Bau, J.-Y. Zhu, H. Strobelt, B. Zhou, J. B. Tenenbaum, W. T. Freeman, and A. Torralba. Gan dissection: Visualizing and understanding generative adversarial networks. *arXiv preprint arXiv:1811.10597*, 2018.

- N. Belrose, Z. Furman, L. Smith, D. Halawi, I. Ostrovsky, L. McKinney, S. Biderman, and J. Steinhardt. Eliciting latent predictions from transformers with the tuned lens. *arXiv preprint arXiv:2303.08112*, 2023.
- N. Cammarata, G. Goh, S. Carter, L. Schubert, M. Petrov, and C. Olah. Curve detectors. *Distill*, 5(6): e00024–003, 2020.
- S. Carter, Z. Armstrong, L. Schubert, I. Johnson, and C. Olah. Activation atlas. *Distill*, 2019. doi: 10.23915/distill.00015. <https://distill.pub/2019/activation-atlas>.
- L. Chan, A. Garriga-Alonso, N. Goldwosky-Dill, R. Greenblatt, J. Nitishinskaya, A. Radhakrishnan, B. Shlegeris, and N. Thomas. Causal scrubbing, a method for rigorously testing interpretability hypotheses. *AI Alignment Forum*, 2022. <https://www.alignmentforum.org/posts/JvZhhzycHu2Yd57RN/causal-scrubbing-a-method-for-rigorously-testing>.
- B. Chughtai, L. Chan, and N. Nanda. A toy model of universality: Reverse engineering how networks learn group operations. *arXiv preprint arXiv:2302.03025*, 2023.
- A. Conmy, A. N. Mavor-Parker, A. Lynch, S. Heimersheim, and A. Garriga-Alonso. Towards automated circuit discovery for mechanistic interpretability. *arXiv preprint arXiv:2304.14997*, 2023.
- G. Dar, M. Geva, A. Gupta, and J. Berant. Analyzing transformers in embedding space. *arXiv preprint arXiv:2209.02535*, 2022.
- T. Dettmers, M. Lewis, Y. Belkada, and L. Zettlemoyer. Llm. int8 (): 8-bit matrix multiplication for transformers at scale. *arXiv preprint arXiv:2208.07339*, 2022.
- N. Elhage, N. Nanda, C. Olsson, T. Henighan, N. Joseph, B. Mann, A. Askell, Y. Bai, A. Chen, T. Conerly, N. DasSarma, D. Drain, D. Ganguli, Z. Hatfield-Dodds, D. Hernandez, A. Jones, J. Kernion, L. Lovitt, K. Ndousse, D. Amodei, T. Brown, J. Clark, J. Kaplan, S. McCandlish, and C. Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. <https://transformer-circuits.pub/2021/framework/index.html>.
- N. Elhage, T. Hume, C. Olsson, N. Schiefer, T. Henighan, S. Kravec, Z. Hatfield-Dodds, R. Lasenby, D. Drain, C. Chen, R. Grosse, S. McCandlish, J. Kaplan, D. Amodei, M. Wattenberg, and C. Olah. Toy models of superposition. *Transformer Circuits Thread*, 2022. https://transformer-circuits.pub/2022/toy_model/index.html.
- M. Finlayson, A. Mueller, S. Gehrmann, S. Shieber, T. Linzen, and Y. Belinkov. Causal analysis of syntactic agreement mechanisms in neural language models. *arXiv preprint arXiv:2106.06087*, 2021.
- J. Z. Forde, C. Lovering, G. Konidaris, E. Pavlick, and M. L. Littman. Where, when & which concepts does alphazero learn? lessons from the game of hex. In *AAAI Workshop on Reinforcement Learning in Games*, volume 2, 2022.
- A. Geiger, H. Lu, T. Icard, and C. Potts. Causal abstractions of neural networks. *Advances in Neural Information Processing Systems*, 34:9574–9586, 2021.
- A. Geiger, Z. Wu, H. Lu, J. Rozner, E. Kreiss, T. Icard, N. Goodman, and C. Potts. Inducing causal structure for interpretable neural networks. In *International Conference on Machine Learning*, pages 7324–7338. PMLR, 2022.
- A. Geiger, C. Potts, and T. Icard. Causal abstraction for faithful model interpretation. *arXiv preprint arXiv:2301.04709*, 2023a.

- A. Geiger, Z. Wu, C. Potts, T. Icard, and N. D. Goodman. Finding alignments between interpretable causal variables and distributed neural representations. *arXiv preprint arXiv:2303.02536*, 2023b.
- M. Geva, A. Caciularu, G. Dar, P. Roit, S. Sadde, M. Shlain, B. Tamir, and Y. Goldberg. Lm-debugger: An interactive tool for inspection and intervention in transformer-based language models. *arXiv preprint arXiv:2204.12130*, 2022a.
- M. Geva, A. Caciularu, K. R. Wang, and Y. Goldberg. Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space. *arXiv preprint arXiv:2203.14680*, 2022b.
- M. Glymour, J. Pearl, and N. P. Jewell. *Causal inference in statistics: A primer*. John Wiley & Sons, 2016.
- G. Goh, N. C. †, C. V. †, S. Carter, M. Petrov, L. Schubert, A. Radford, and C. Olah. Multimodal neurons in artificial neural networks. *Distill*, 2021. doi: 10.23915/distill.00030. <https://distill.pub/2021/multimodal-neurons>.
- W. Gurnee, N. Nanda, M. Pauly, K. Harvey, D. Troitskii, and D. Bertsimas. Finding neurons in a haystack: Case studies with sparse probing. *arXiv preprint arXiv:2305.01610*, 2023.
- J. Y. Halpern. *Actual causality*. MIT Press, 2016.
- J. Hilton, N. Cammarata, S. Carter, G. Goh, and C. Olah. Understanding rl vision. *Distill*, 2020. doi: 10.23915/distill.00029. <https://distill.pub/2020/understanding-rl-vision>.
- J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. d. L. Casas, L. A. Hendricks, J. Welbl, A. Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- M. L. Leavitt and A. Morcos. Towards falsifiable interpretability research. *arXiv preprint arXiv:2010.12016*, 2020.
- R. Massidda, A. Geiger, T. Icard, and D. Bacciu. Causal abstraction with soft interventions. *arXiv preprint arXiv:2211.12270*, 2022.
- T. McGrath, A. Kapishnikov, N. Tomašev, A. Pearce, M. Wattenberg, D. Hassabis, B. Kim, U. Paquet, and V. Kramnik. Acquisition of chess knowledge in alphazero. *Proceedings of the National Academy of Sciences*, 119(47):e2206625119, 2022.
- K. Meng, D. Bau, A. Andonian, and Y. Belinkov. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.
- A. S. Morcos, D. G. Barrett, N. C. Rabinowitz, and M. Botvinick. On the importance of single directions for generalization. *arXiv preprint arXiv:1803.06959*, 2018.
- N. Nanda, L. Chan, T. Liberum, J. Smith, and J. Steinhardt. Progress measures for grokking via mechanistic interpretability. *arXiv preprint arXiv:2301.05217*, 2023.
- nostalgebraist. interpreting gpt: the logit lens, 2020. URL <https://www.lesswrong.com/posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens>.
- C. Olah, A. Satyanarayan, I. Johnson, S. Carter, L. Schubert, K. Ye, and A. Mordvintsev. The building blocks of interpretability. *Distill*, 3(3):e10, 2018.
- C. Olah, N. Cammarata, L. Schubert, G. Goh, M. Petrov, and S. Carter. Zoom in: An introduction to circuits. *Distill*, 2020. doi: 10.23915/distill.00024.001. <https://distill.pub/2020/circuits/zoom-in>.

- C. Olsson, N. Elhage, N. Nanda, N. Joseph, N. DasSarma, T. Henighan, B. Mann, A. Askell, Y. Bai, A. Chen, et al. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895*, 2022.
- J. Pearl. *Causality*. Cambridge university press, 2009.
- M. Phuong and M. Hutter. Formal algorithms for transformers. *arXiv preprint arXiv:2207.09238*, 2022.
- A. Power, Y. Burda, H. Edwards, I. Babuschkin, and V. Misra. Grokking: Generalization beyond overfitting on small algorithmic datasets. *arXiv preprint arXiv:2201.02177*, 2022.
- N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1): 1929–1958, 2014.
- N. Tinbergen. On aims and methods of ethology. *Zeitschrift für tierpsychologie*, 20(4):410–433, 1963.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- A. Veit, M. J. Wilber, and S. Belongie. Residual networks behave like ensembles of relatively shallow networks. *Advances in neural information processing systems*, 29, 2016.
- J. Vig, S. Gehrmann, Y. Belinkov, S. Qian, D. Nevo, S. Sakenis, J. Huang, Y. Singer, and S. Shieber. Causal mediation analysis for interpreting neural nlp: The case of gender bias. *arXiv preprint arXiv:2004.12265*, 2020a.
- J. Vig, S. Gehrmann, Y. Belinkov, S. Qian, D. Nevo, Y. Singer, and S. Shieber. Investigating gender bias in language models using causal mediation analysis. *Advances in neural information processing systems*, 33:12388–12401, 2020b.
- K. Wang, A. Variengien, A. Conmy, B. Shlegeris, and J. Steinhardt. Interpretability in the wild: a circuit for indirect object identification in gpt-2 small. *arXiv preprint arXiv:2211.00593*, 2022.
- B. Zhang and R. Sennrich. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32, 2019.

A. Choice of intervention distribution

When we intervene on a neural network’s activations with the do-operator, we are setting some subset of these activations to a new value in the forward pass and allowing the effect of these changes to propagate forward. Both the corrupted run and the patching operation in the causal tracing method (Meng et al., 2022) are examples of interventions. Practically we can accomplish these interventions via PyTorch hooks, JAX’s Harvest functionality, or passing values for intervened-upon tensors in Tensorflow’s feed-dict. When we intervene to set some network activation Z to an ‘ablation’ value \tilde{z} , what value should we use? Four main possibilities have been suggested:

$$\text{Zero ablation} : \tilde{z} = 0, \tag{35}$$

$$\text{Mean ablation} : \tilde{z} = \mathbb{E}_{u \sim P(u)} [Z(u)], \tag{36}$$

$$\text{Noise ablation} : \tilde{z} = Z(u + \epsilon), \epsilon \sim \mathcal{N}(0, \sigma^2), \tag{37}$$

$$\text{Resample ablation} : \tilde{z} = Z(\tilde{u}), \tilde{u} \sim P(u). \tag{38}$$

Of these, zero-ablation is the simplest to implement but is typically out-of-distribution for networks trained without some form of layer dropout or stochastic depth. Noise ablation was used extensively in causal tracing (Meng et al., 2022). Resample ablation (as introduced by Chan et al. (2022)) is more complicated to implement but is probably the most principled, as every ablation is a naturally-occurring set of activations, and so is more likely to respect properties of network activations such as emergent outliers (Dettmers et al., 2022). Resample ablation also has the appealing property that by specifying the distribution of inputs $P(u)$ we can control for properties of the input that might otherwise confound our results. To get meaningful results from sample ablations it is necessary to use an average of sample ablations across multiple samples from the same dataset, i.e. a Monte-Carlo approximation to:

$$V_Z(u) = \int V(u|do(Z = \tilde{z})) p(\tilde{z}) d\tilde{z}, \quad (39)$$

where $p(\tilde{z}) = \int p(Z(u))du$ is the probability of getting activation values \tilde{z} . Note that mean ablation and resample ablation are quite different: mean ablation ablates with the *average activation*, whereas resample activation averages the *effects* of different ablations. See (Chan et al., 2022) for an extended discussion of these methodological details.