In this paper, (Meng et al., 2022b) define $m_{[t]}^l = W_{\text{out}}^l(\sigma(W_{\text{in}}^l\gamma(h_{[t]}^{l-1})))$. Given this, define

$$m_i^l = W_{\text{out}}k_i^l + r_i^l$$

$$\text{where } r_i^l = \frac{z_i - h_i^L}{L - l + 1} \quad (8)$$

Note that the denominator of $r_i^l$ allows us to spread out the burden across multiple layer, allowing for a more scalable algorithm. It is hard to compute $C^l$ exactly, however it can be reliably estimated using $C^l = \lambda\mathbb{E}_k[k^l k^{l\top}]$ over randomly sampled inputs, where $\lambda = 1.5 \times 10^4$. Incorporating the update gives us our desired new weights $\hat{W}_{\text{out}}^l$

### A.6.3 MEND

Using the fact that the gradient of loss L with respect to the weights $W_\ell$ of layer $\ell$ of an MLP has a rank-one decomposition such that $\nabla_{W_\ell}L = \sum_{i=1}^B \delta_{\ell+1}^i u_\ell^{i\top}$ for a batch $B$, (Mitchell et al., 2021) are able to construct an editor network $g_\ell$ to generate the weight updates. Here, $\delta_{\ell+1}^i$ is the gradient for element $i$ for the preactivations of layer $\ell + 1$ and $u_\ell^i$ are the inputs of element $i$ into layer $\ell$.

To characterize these updates, MEND employs functions that map $\delta_{\ell+1}^i$ and $u_\ell^i$ to a pseudo-decomposition $\tilde{\delta}_{\ell+1}^i$ and $\tilde{u}_\ell^i$ such that $\tilde{\nabla}_{W_\ell}L = \sum_{i=1}^B \tilde{\delta}_{\ell+1}^i \tilde{u}_\ell^{i\top}$. Letting $z_\ell = \text{concat}(\delta_{\ell+1}, u_\ell)$, the form of the network is

$$h_\ell = z_\ell + \sigma(s_\ell^1 \odot (U_1V_1z_\ell + b) + o_\ell^1) \quad (9)$$
$$g(z_\ell) = h_\ell + \sigma(s_\ell^2 \odot U_2V_2h_\ell + o_\ell^2) \quad (10)$$

where $\sigma$ is the ReLu activation function and $U_j$, $V_j$ are a low rank decomposition of MEND's weight for layer $j$. Note that, because of the difference in dimensions between weight matrices across layers, MEND learns different parameters for each unique shape of weight matrices to be edited. Additionally, layer-wise offset and scale parameters $o_\ell$ and $s_\ell$ are learned for both $h_\ell$ and $g_\ell$. The final update is given by $\tilde{W} = W - \alpha_\ell\tilde{\nabla}_{W_\ell}$ with $\alpha_\ell$ being another learned parameter per layer.

Given the original weights $W$ and the updated weights $\tilde{W}$, loss is computed aggregating two training losses, editing success and locality. For a desired edit $(x_e, y_e)$, $(x_e', y_e')$ is defined as a semantically equivalent wording of the edit. Editing loss is defined as $L_e = -\log_{p_{\theta_{\tilde{W}}}} p(y_e'|x_e')$. $x_{\text{loc}}$ is defined as a locality sample, which is randomly sampled to test the edited model's impact on information

unrelated to the edit. The corresponding locality loss is $L_{\text{loc}} = D_{KL}(p_{\theta_W}(\cdot|x_{\text{loc}})\|p_{\theta_{\tilde{W}}}(\cdot|x_{\text{loc}}))$.

The total loss is computed as $L_{\text{MEND}} = c_eL_e(\theta_{\tilde{W}}) + L_{\text{loc}}(\theta_W, \theta_{\tilde{W}})$ where $c_e = 0.1$. The total loss is optimized using the Adam optimizer. For GPT2-XL, we edit layers 46, 47, and 48. For GPTJ, we edit layers 26, 27, and 28.

### A.6.4 FT

The Fine-Tuning procedures used in this paper follow from (Meng et al., 2022b) and (Meng et al., 2022a)'s implementation for both GPT-J and GPT2-XL. MLP weights for a single layer are fine-tuned for both models. We use a constrained fine tuning approach where we add a $L_\infty$ constraint such that $\|\theta_G - \theta_{G'}\|_\infty \leq \epsilon$ at each gradient step. For the constraint, $\epsilon = 5e-4$ for GPT2-XL and $\epsilon = 5e-5$ for GPT-J. It is optimized using Adam with a learning rate of 5e-4 for both GPT2-XL and GPT-J. We fine tune layers 18 and 6 for GPT2-XL and GPT-J respectively.