
Supplemental Material: Similarity-Distance-Magnitude Activations

Allen Schmaltz
Reexpress AI

allen@re.express

A Appendix: Continued

A.4 Comparison to Bayesian Last-layer Networks

In the main text, we compare to representative Frequentist and empirically motivated estimators. In this section, we further compare to variational Bayesian last-layer neural networks (VBLL), a computationally efficient Bayesian approach (Harrison et al., 2024). The basic setup is similar to the Frequentist and empirically motivated approaches examined in the main text in that it involves training a small final-layer adaptor over the frozen parameters of the language model. However, in this case, the adaptor network is a multi-layer perceptron (MLP) combined with the VBLL estimator. We consider both the discriminative and generative VBLL estimators.

Models. We follow the parameter choices and architecture of Harrison et al. (2024), and its associated code tutorial¹, with applicable changes to match the experimental settings of the other models and estimators. Specifically, the VBLL models consist of an input linear layer, 2 core linear layers, and a final output linear layer. The hidden dimension is set at 795, which yields a similar number of parameters as the SDM activation layers (approximately 6 million parameters for the `Phi-3.5-mini-instruct` adaptors and approximately 8 million parameters for the `Mixtral 8x7B Instruct v0.1` adaptors). The input to the VBLL models is the same mean-centered embeddings used with the SDM activation layers and the baseline CNN adaptor layers of the main text. We use the AdamW optimizer (Loshchilov & Hutter, 2019) with a weight decay of 1×10^{-4} , which following the code tutorial, is not applied to the final layer. We use a learning rate of 1×10^{-5} for training, matching that used with the SDM activation layers. Following the code tutorial, we use gradient clipping with a max norm of 1, and we use the exponential linear unit (ELU) as the activation function (Clevert et al., 2016). We train for 200 epochs, choosing the epoch with the lowest held-out validation loss over \mathcal{D}_{ca} as the chosen weights. We repeat this process for $J = 10$ shuffles of the data (i.e., splits of \mathcal{D}_{tr} and \mathcal{D}_{ca}), choosing the model with the globally lowest overall held-out validation loss over \mathcal{D}_{ca} as the final model.

For the discriminative models over `Phi-3.5-mini-instruct` and `Mixtral 8x7B Instruct v0.1`, we use the labels `PHI3.5+DiscVBLLMLP` and `MIXTRAL8X7B+DiscVBLLMLP`, respectively. For the generative models over `Phi-3.5-mini-instruct` and `Mixtral 8x7B Instruct v0.1`, we use the labels `PHI3.5+GenVBLLMLP` and `MIXTRAL8X7B+GenVBLLMLP`, respectively. For all of the aforementioned models, we use a KL regularization weight set at $\frac{1}{|\mathcal{D}_{\text{tr}}|}$, as in the original paper. We also consider analogous models that increase the KL regularization weight by a multiplicative factor of 50. Increasing the KL regularization weight is suggested in the code tutorial as the “simplest and most effective method to control the scale of uncertainty” of VBLL models. For these models with the larger KL regularization weight, we use the labels `PHI3.5+DiscVBLLMLPrw50`, `MIXTRAL8X7B+DiscVBLLMLPrw50`, `PHI3.5+GenVBLLMLPrw50`, and `MIXTRAL8X7B+GenVBLLMLPrw50`, respectively.

We use the label `VBLL` for the estimator that thresholds the output of the variational Bayesian last-layer neural network at α for the predicted class, analogous to the softmax and SDM estimators. The NO-REJECT estimator provides a reference point without any selection criteria applied (i.e., the standard marginal and class- and prediction-conditional accuracies over the given dataset).

¹Available at <https://github.com/VectorInstitute/vbll>

Results. The results for the sentiment datasets appear in Table 5, and the results for the factcheck datasets appear in Table 6. For reference, in all cases we also provide the results over the held-out calibration set, \mathcal{D}_{ca} , which is the held-out split used to determine the final model weights, as noted above.

Comparing to Table 3 in the main text, the accuracies over \mathcal{D}_{ca} for the NO-REJECT estimators are similar for the VBLL models and the models using SDM activation functions. The MLPs of the VBLL models and the 1-D CNNs of the SDM activation layers have a similar number of parameters and are trained with the same number of epochs and the same number of iterated shuffles of \mathcal{D}_{tr} and \mathcal{D}_{ca} . The accuracies without selection for the SENTIMENT calibration set are already at least α , and those for the FACTCHECK calibration set are below α , but at least $\alpha - 0.05$. As such, differences in calibration effectiveness of the respective methods are not directly attributable to substantively different baseline accuracies for the in-distribution tasks.

We find that the VBLL estimators are well-calibrated in high-probability regions over in-distribution data, but generally fare poorly over co-variate shifts and out-of-distribution data. We find no clear advantages nor disadvantages for the discriminative vs. generative variants, and modifying the KL regularization weight has a minimal impact, at least at this scale. In Harrison et al. (2024), VBLL estimators over LMs for sentiment classification are only compared to an MLP baseline on in-distribution data; our evaluation setting is significantly more challenging and closer to real-world conditions encountered with LM applications.

These results corroborate the patterns observed in the main text. In our controlled examination of representative Bayesian, Frequentist, and empirically motivated estimators, the SDM_{HR} estimator (i.e., the SDM activation function combined with the selection criteria of Eq. 9) is the only estimator that remains well-calibrated in high-probability regions across the in-distribution, co-variate shifted, and out-of-distribution data.

Table 5: Comparison of Bayesian last-layer estimators (Harrison et al., 2024) for the sentiment datasets, including the shuffled challenge sets, with $\alpha = 0.95$. R indicates all predictions were rejected, which is preferred over falling under the expected accuracy. $n = |\text{Admitted}|$, the count of non-rejected documents.

| Dataset | Model | Estimator | Class-conditional | | Prediction-conditional | | Marginal | |
|------------------------------|--|-----------|-------------------|---------|------------------------|---------------|------------------|------|
| | | | $y = 0$ | $y = 1$ | $\hat{y} = 0$ | $\hat{y} = 1$ | $y \in \{0, 1\}$ | |
| SENTIMENT \mathcal{D}_{ca} | PHI3.5+DiscVBLMLP | NO-REJECT | 0.97 | 0.51 | 0.96 | 0.49 | 0.96 | 0.51 |
| SENTIMENT | PHI3.5+DiscVBLMLP | NO-REJECT | 0.97 | 0.50 | 0.95 | 0.50 | 0.95 | 0.51 |
| SENTIMENTOOD | PHI3.5+DiscVBLMLP | NO-REJECT | 0.61 | 0.50 | 0.67 | 0.50 | 0.65 | 0.47 |
| SENTIMENTSHUFFLED | PHI3.5+DiscVBLMLP | NO-REJECT | 0.87 | 0.50 | 0.75 | 0.50 | 0.78 | 0.56 |
| SENTIMENTOODSHUFFLED | PHI3.5+DiscVBLMLP | NO-REJECT | 0.77 | 0.50 | 0.56 | 0.50 | 0.64 | 0.61 |
| SENTIMENT \mathcal{D}_{ca} | PHI3.5+DiscVBLMLP | VBL | 0.99 | 0.42 | 0.99 | 0.40 | 0.99 | 0.42 |
| SENTIMENT | PHI3.5+DiscVBLMLP | VBL | 0.99 | 0.42 | 1.00 | 0.40 | 1.00 | 0.42 |
| SENTIMENTOOD | PHI3.5+DiscVBLMLP | VBL | 0.89 | 0.03 | 0.96 | 0.04 | 0.93 | 0.02 |
| SENTIMENTSHUFFLED | PHI3.5+DiscVBLMLP | VBL | 0.99 | 0.13 | 0.95 | 0.05 | 0.98 | 0.13 |
| SENTIMENTOODSHUFFLED | PHI3.5+DiscVBLMLP | VBL | 0.98 | 0.02 | 0.86 | 0.02 | 0.89 | 0.02 |
| SENTIMENT \mathcal{D}_{ca} | PHI3.5+DiscVBLMLP _{rw50} | NO-REJECT | 0.96 | 0.51 | 0.96 | 0.49 | 0.96 | 0.51 |
| SENTIMENT | PHI3.5+DiscVBLMLP _{rw50} | NO-REJECT | 0.97 | 0.50 | 0.95 | 0.50 | 0.95 | 0.51 |
| SENTIMENTOOD | PHI3.5+DiscVBLMLP _{rw50} | NO-REJECT | 0.64 | 0.50 | 0.71 | 0.50 | 0.69 | 0.47 |
| SENTIMENTSHUFFLED | PHI3.5+DiscVBLMLP _{rw50} | NO-REJECT | 0.85 | 0.50 | 0.77 | 0.50 | 0.79 | 0.54 |
| SENTIMENTOODSHUFFLED | PHI3.5+DiscVBLMLP _{rw50} | NO-REJECT | 0.80 | 0.50 | 0.60 | 0.50 | 0.67 | 0.60 |
| SENTIMENT \mathcal{D}_{ca} | PHI3.5+DiscVBLMLP _{rw50} | VBL | 0.99 | 0.43 | 0.99 | 0.41 | 0.99 | 0.42 |
| SENTIMENT | PHI3.5+DiscVBLMLP _{rw50} | VBL | 0.99 | 0.42 | 1.00 | 0.41 | 1.00 | 0.42 |
| SENTIMENTOOD | PHI3.5+DiscVBLMLP _{rw50} | VBL | 0.95 | 0.02 | 0.96 | 0.04 | 0.92 | 0.02 |
| SENTIMENTSHUFFLED | PHI3.5+DiscVBLMLP _{rw50} | VBL | 0.99 | 0.13 | 0.96 | 0.09 | 0.98 | 0.13 |
| SENTIMENTOODSHUFFLED | PHI3.5+DiscVBLMLP _{rw50} | VBL | 0.98 | 0.02 | 0.86 | 0.02 | 0.90 | 0.02 |
| SENTIMENT \mathcal{D}_{ca} | PHI3.5+GenVBLMLP | NO-REJECT | 0.96 | 0.50 | 0.96 | 0.50 | 0.96 | 0.50 |
| SENTIMENT | PHI3.5+GenVBLMLP | NO-REJECT | 0.97 | 0.50 | 0.96 | 0.50 | 0.96 | 0.51 |
| SENTIMENTOOD | PHI3.5+GenVBLMLP | NO-REJECT | 0.28 | 0.50 | 0.93 | 0.50 | 0.80 | 0.17 |
| SENTIMENTSHUFFLED | PHI3.5+GenVBLMLP | NO-REJECT | 0.94 | 0.50 | 0.60 | 0.50 | 0.70 | 0.67 |
| SENTIMENTOODSHUFFLED | PHI3.5+GenVBLMLP | NO-REJECT | 0.43 | 0.50 | 0.85 | 0.50 | 0.74 | 0.29 |
| SENTIMENT \mathcal{D}_{ca} | PHI3.5+GenVBLMLP | VBL | 0.99 | 0.43 | 0.99 | 0.44 | 0.99 | 0.43 |
| SENTIMENT | PHI3.5+GenVBLMLP | VBL | 0.99 | 0.44 | 0.99 | 0.43 | 0.99 | 0.44 |
| SENTIMENTOOD | PHI3.5+GenVBLMLP | VBL | 0.13 | 0.16 | 0.99 | 0.28 | 0.89 | 0.02 |
| SENTIMENTSHUFFLED | PHI3.5+GenVBLMLP | VBL | 1.00 | 0.31 | 0.70 | 0.14 | 0.88 | 0.35 |
| SENTIMENTOODSHUFFLED | PHI3.5+GenVBLMLP | VBL | 0.32 | 0.08 | 0.97 | 0.19 | 0.80 | 0.03 |
| SENTIMENT \mathcal{D}_{ca} | PHI3.5+GenVBLMLP _{rw50} | NO-REJECT | 0.96 | 0.50 | 0.96 | 0.50 | 0.96 | 0.50 |
| SENTIMENT | PHI3.5+GenVBLMLP _{rw50} | NO-REJECT | 0.97 | 0.50 | 0.96 | 0.50 | 0.96 | 0.51 |
| SENTIMENTOOD | PHI3.5+GenVBLMLP _{rw50} | NO-REJECT | 0.43 | 0.50 | 0.78 | 0.50 | 0.66 | 0.32 |
| SENTIMENTSHUFFLED | PHI3.5+GenVBLMLP _{rw50} | NO-REJECT | 0.78 | 0.50 | 0.84 | 0.50 | 0.83 | 0.47 |
| SENTIMENT \mathcal{D}_{ca} | PHI3.5+GenVBLMLP _{rw50} | VBL | 0.64 | 0.50 | 0.76 | 0.50 | 0.72 | 0.44 |
| SENTIMENT | PHI3.5+GenVBLMLP _{rw50} | VBL | 0.99 | 0.45 | 0.99 | 0.43 | 0.99 | 0.45 |
| SENTIMENTOOD | PHI3.5+GenVBLMLP _{rw50} | VBL | 0.99 | 0.44 | 0.99 | 0.44 | 0.99 | 0.44 |
| SENTIMENTSHUFFLED | PHI3.5+GenVBLMLP _{rw50} | VBL | 0.13 | 0.16 | 0.99 | 0.28 | 0.89 | 0.02 |
| SENTIMENTOODSHUFFLED | PHI3.5+GenVBLMLP _{rw50} | VBL | 1.00 | 0.31 | 0.70 | 0.14 | 0.88 | 0.35 |
| SENTIMENT \mathcal{D}_{ca} | MIXTRAL8x7B+DiscVBLMLP | NO-REJECT | 0.96 | 0.51 | 0.97 | 0.49 | 0.97 | 0.50 |
| SENTIMENT | MIXTRAL8x7B+DiscVBLMLP | NO-REJECT | 0.96 | 0.50 | 0.96 | 0.50 | 0.96 | 0.50 |
| SENTIMENTOOD | MIXTRAL8x7B+DiscVBLMLP | NO-REJECT | 0.43 | 0.50 | 0.78 | 0.50 | 0.66 | 0.32 |
| SENTIMENTSHUFFLED | MIXTRAL8x7B+DiscVBLMLP | NO-REJECT | 0.78 | 0.50 | 0.84 | 0.50 | 0.83 | 0.47 |
| SENTIMENTOODSHUFFLED | MIXTRAL8x7B+DiscVBLMLP | NO-REJECT | 0.64 | 0.50 | 0.76 | 0.50 | 0.72 | 0.44 |
| SENTIMENT \mathcal{D}_{ca} | MIXTRAL8x7B+DiscVBLMLP | VBL | 0.99 | 0.45 | 0.99 | 0.43 | 0.99 | 0.45 |
| SENTIMENT | MIXTRAL8x7B+DiscVBLMLP | VBL | 0.99 | 0.44 | 0.99 | 0.44 | 0.99 | 0.44 |
| SENTIMENTOOD | MIXTRAL8x7B+DiscVBLMLP | VBL | 0.99 | 0.44 | 0.99 | 0.43 | 0.99 | 0.44 |
| SENTIMENTSHUFFLED | MIXTRAL8x7B+DiscVBLMLP | VBL | 0.99 | 0.01 | 0.97 | 0.11 | 0.79 | 0.02 |
| SENTIMENTOODSHUFFLED | MIXTRAL8x7B+DiscVBLMLP | VBL | 0.99 | 0.24 | 0.91 | 0.12 | 0.95 | 0.25 |
| SENTIMENT \mathcal{D}_{ca} | MIXTRAL8x7B+DiscVBLMLP _{rw50} | NO-REJECT | 0.96 | 0.50 | 0.97 | 0.50 | 0.96 | 0.50 |
| SENTIMENT | MIXTRAL8x7B+DiscVBLMLP _{rw50} | NO-REJECT | 0.97 | 0.50 | 0.96 | 0.50 | 0.97 | 0.50 |
| SENTIMENTOOD | MIXTRAL8x7B+DiscVBLMLP _{rw50} | NO-REJECT | 0.89 | 0.50 | 0.59 | 0.50 | 0.69 | 0.65 |
| SENTIMENTSHUFFLED | MIXTRAL8x7B+DiscVBLMLP _{rw50} | NO-REJECT | 0.86 | 0.50 | 0.81 | 0.50 | 0.81 | 0.53 |
| SENTIMENTOODSHUFFLED | MIXTRAL8x7B+DiscVBLMLP _{rw50} | NO-REJECT | 0.94 | 0.50 | 0.44 | 0.50 | 0.63 | 0.75 |
| SENTIMENT \mathcal{D}_{ca} | MIXTRAL8x7B+DiscVBLMLP _{rw50} | VBL | 0.99 | 0.44 | 0.99 | 0.44 | 0.99 | 0.44 |
| SENTIMENT | MIXTRAL8x7B+DiscVBLMLP _{rw50} | VBL | 0.99 | 0.44 | 0.99 | 0.43 | 0.99 | 0.43 |
| SENTIMENTOOD | MIXTRAL8x7B+DiscVBLMLP _{rw50} | VBL | 0.99 | 0.01 | 0.97 | 0.11 | 0.79 | 0.02 |
| SENTIMENTSHUFFLED | MIXTRAL8x7B+DiscVBLMLP _{rw50} | VBL | 0.99 | 0.24 | 0.91 | 0.12 | 0.95 | 0.25 |
| SENTIMENTOODSHUFFLED | MIXTRAL8x7B+DiscVBLMLP _{rw50} | VBL | 1.00 | 0.02 | 0.69 | 0.01 | 0.87 | 0.02 |
| SENTIMENT \mathcal{D}_{ca} | MIXTRAL8x7B+DiscVBLMLP _{rw50} | NO-REJECT | 0.96 | 0.50 | 0.97 | 0.50 | 0.96 | 0.50 |
| SENTIMENT | MIXTRAL8x7B+DiscVBLMLP _{rw50} | NO-REJECT | 0.97 | 0.50 | 0.96 | 0.50 | 0.97 | 0.50 |
| SENTIMENTOOD | MIXTRAL8x7B+DiscVBLMLP _{rw50} | NO-REJECT | 0.89 | 0.50 | 0.59 | 0.50 | 0.68 | 0.35 |
| SENTIMENTSHUFFLED | MIXTRAL8x7B+DiscVBLMLP _{rw50} | NO-REJECT | 0.86 | 0.50 | 0.81 | 0.50 | 0.81 | 0.53 |
| SENTIMENTOODSHUFFLED | MIXTRAL8x7B+DiscVBLMLP _{rw50} | NO-REJECT | 0.94 | 0.50 | 0.44 | 0.50 | 0.63 | 0.75 |
| SENTIMENT \mathcal{D}_{ca} | MIXTRAL8x7B+DiscVBLMLP _{rw50} | VBL | 0.99 | 0.44 | 0.99 | 0.44 | 0.99 | 0.44 |
| SENTIMENT | MIXTRAL8x7B+DiscVBLMLP _{rw50} | VBL | 0.99 | 0.44 | 0.99 | 0.43 | 0.99 | 0.43 |
| SENTIMENTOOD | MIXTRAL8x7B+DiscVBLMLP _{rw50} | VBL | 0.99 | 0.01 | 0.97 | 0.11 | 0.79 | 0.02 |
| SENTIMENTSHUFFLED | MIXTRAL8x7B+DiscVBLMLP _{rw50} | VBL | 0.99 | 0.24 | 0.91 | 0.12 | 0.95 | 0.25 |
| SENTIMENTOODSHUFFLED | MIXTRAL8x7B+DiscVBLMLP _{rw50} | VBL | 1.00 | 0.02 | 0.69 | 0.01 | 0.87 | 0.02 |
| SENTIMENT \mathcal{D}_{ca} | MIXTRAL8x7B+GenVBLMLP | NO-REJECT | 0.97 | 0.50 | 0.97 | 0.50 | 0.97 | 0.50 |
| SENTIMENT | MIXTRAL8x7B+GenVBLMLP | NO-REJECT | 0.97 | 0.50 | 0.96 | 0.50 | 0.97 | 0.50 |
| SENTIMENTOOD | MIXTRAL8x7B+GenVBLMLP | NO-REJECT | 0.80 | 0.50 | 0.72 | 0.50 | 0.74 | 0.54 |
| SENTIMENTSHUFFLED | MIXTRAL8x7B+GenVBLMLP | NO-REJECT | 0.84 | 0.50 | 0.82 | 0.50 | 0.82 | 0.51 |
| SENTIMENTOODSHUFFLED | MIXTRAL8x7B+GenVBLMLP | NO-REJECT | 0.85 | 0.50 | 0.54 | 0.50 | 0.65 | 0.65 |
| SENTIMENT \mathcal{D}_{ca} | MIXTRAL8x7B+GenVBLMLP | VBL | 0.99 | 0.44 | 0.99 | 0.45 | 0.99 | 0.44 |
| SENTIMENT | MIXTRAL8x7B+GenVBLMLP | VBL | 0.99 | 0.44 | 0.99 | 0.44 | 0.99 | 0.44 |
| SENTIMENTOOD | MIXTRAL8x7B+GenVBLMLP | VBL | 0.98 | 0.03 | 0.93 | 0.07 | 0.87 | 0.04 |
| SENTIMENTSHUFFLED | MIXTRAL8x7B+GenVBLMLP | VBL | 0.99 | 0.24 | 0.92 | 0.14 | 0.95 | 0.25 |
| SENTIMENTOODSHUFFLED | MIXTRAL8x7B+GenVBLMLP | VBL | 1.00 | 0.02 | 0.69 | 0.01 | 0.87 | 0.02 |
| SENTIMENT \mathcal{D}_{ca} | MIXTRAL8x7B+GenVBLMLP _{rw50} | NO-REJECT | 0.97 | 0.50 | 0.97 | 0.50 | 0.97 | 0.50 |
| SENTIMENT | MIXTRAL8x7B+GenVBLMLP _{rw50} | NO-REJECT | 0.97 | 0.50 | 0.96 | 0.50 | 0.97 | 0.50 |
| SENTIMENTOOD | MIXTRAL8x7B+GenVBLMLP _{rw50} | NO-REJECT | 0.79 | 0.50 | 0.73 | 0.50 | 0.75 | 0.53 |
| SENTIMENTSHUFFLED | MIXTRAL8x7B+GenVBLMLP _{rw50} | NO-REJECT | 0.83 | 0.50 | 0.82 | 0.50 | 0.82 | 0.51 |
| SENTIMENTOODSHUFFLED | MIXTRAL8x7B+GenVBLMLP _{rw50} | NO-REJECT | 0.83 | 0.50 | 0.56 | 0.50 | 0.65 | 0.64 |
| SENTIMENT \mathcal{D}_{ca} | MIXTRAL8x7B+GenVBLMLP _{rw50} | VBL | 0.98 | 0.44 | 0.99 | 0.46 | 0.99 | 0.44 |
| SENTIMENT | MIXTRAL8x7B+GenVBLMLP _{rw50} | VBL | 0.99 | 0.44 | 0.99 | 0.45 | 0.99 | 0.44 |
| SENTIMENTOOD | MIXTRAL8x7B+GenVBLMLP _{rw50} | VBL | 0.91 | 0.04 | 0.99 | 0.16 | 0.95 | 0.04 |
| SENTIMENTSHUFFLED | MIXTRAL8x7B+GenVBLMLP _{rw50} | VBL | 0.95 | 0.24 | 0.91 | 0.21 | 0.92 | 0.24 |
| SENTIMENTOODSHUFFLED | MIXTRAL8x7B+GenVBLMLP _{rw50} | VBL | 0.98 | 0.02 | 0.82 | 0.04 | 0.73 | 0.03 |

Table 6: Comparison of Bayesian last-layer estimators (Harrison et al., 2024) for the factcheck datasets, including the shuffled challenge sets, with $\alpha = 0.95$. R indicates all predictions were rejected, which is preferred over falling under the expected accuracy. $n = |\text{Admitted}|$, the count of non-rejected documents.

| Dataset | Model | Estimator | Class-conditional | | Prediction-conditional | | Marginal | |
|------------------------------|---|-----------|-------------------|---------|------------------------|---------------|------------------|------|
| | | | $y = 0$ | $y = 1$ | $\hat{y} = 0$ | $\hat{y} = 1$ | $y \in \{0, 1\}$ | |
| FACTCHECK \mathcal{D}_{ca} | PHI3.5+DiscVBLLMLP | NO-REJECT | 0.92 | 0.50 | 0.91 | 0.50 | 0.91 | 0.50 |
| FACTCHECK | PHI3.5+DiscVBLLMLP | NO-REJECT | 0.38 | 0.51 | 0.92 | 0.49 | 0.84 | 0.23 |
| FACTCHECKSHUFFLED | PHI3.5+DiscVBLLMLP | NO-REJECT | 0.33 | 1. | R | 0. | 1. | 0.33 |
| FACTCHECK \mathcal{D}_{ca} | PHI3.5+DiscVBLLMLP | VBL | 0.98 | 0.32 | 0.99 | 0.28 | 0.99 | 0.32 |
| FACTCHECK | PHI3.5+DiscVBLLMLP | VBL | 0.35 | 0.07 | 1. | 0.31 | 1. | 0.02 |
| FACTCHECKSHUFFLED | PHI3.5+DiscVBLLMLP | VBL | 0.15 | 0.21 | R | 0. | 1. | 0.03 |
| FACTCHECK \mathcal{D}_{ca} | PHI3.5+DiscVBLLMLP _{rw50} | NO-REJECT | 0.91 | 0.50 | 0.92 | 0.50 | 0.92 | 0.50 |
| FACTCHECK | PHI3.5+DiscVBLLMLP _{rw50} | NO-REJECT | 0.45 | 0.51 | 0.94 | 0.49 | 0.89 | 0.26 |
| FACTCHECKSHUFFLED | PHI3.5+DiscVBLLMLP _{rw50} | NO-REJECT | 0.43 | 1. | R | 0. | 1. | 0.43 |
| FACTCHECK \mathcal{D}_{ca} | PHI3.5+DiscVBLLMLP _{rw50} | VBL | 0.98 | 0.34 | 0.98 | 0.33 | 0.98 | 0.34 |
| FACTCHECK | PHI3.5+DiscVBLLMLP _{rw50} | VBL | 0.35 | 0.11 | 0.99 | 0.32 | 0.90 | 0.04 |
| FACTCHECKSHUFFLED | PHI3.5+DiscVBLLMLP _{rw50} | VBL | 0.33 | 0.33 | R | 0. | 1. | 0.11 |
| FACTCHECK \mathcal{D}_{ca} | PHI3.5+GENVBLLMLP | NO-REJECT | 0.90 | 0.50 | 0.93 | 0.50 | 0.93 | 0.49 |
| FACTCHECK | PHI3.5+GENVBLLMLP | NO-REJECT | 0.41 | 0.51 | 0.93 | 0.49 | 0.87 | 0.24 |
| FACTCHECKSHUFFLED | PHI3.5+GENVBLLMLP | NO-REJECT | 0.40 | 1. | R | 0. | 1. | 0.40 |
| FACTCHECK \mathcal{D}_{ca} | PHI3.5+GENVBLLMLP | VBL | 0.98 | 0.33 | 0.99 | 0.32 | 0.99 | 0.33 |
| FACTCHECK | PHI3.5+GENVBLLMLP | VBL | 0.31 | 0.07 | 0.99 | 0.30 | 0.83 | 0.02 |
| FACTCHECKSHUFFLED | PHI3.5+GENVBLLMLP | VBL | 0.20 | 0.22 | R | 0. | 1. | 0.04 |
| FACTCHECK \mathcal{D}_{ca} | PHI3.5+GENVBLLMLP _{rw50} | NO-REJECT | 0.91 | 0.50 | 0.93 | 0.50 | 0.93 | 0.49 |
| FACTCHECK | PHI3.5+GENVBLLMLP _{rw50} | NO-REJECT | 0.41 | 0.51 | 0.93 | 0.49 | 0.87 | 0.24 |
| FACTCHECKSHUFFLED | PHI3.5+GENVBLLMLP _{rw50} | NO-REJECT | 0.44 | 1. | R | 0. | 1. | 0.44 |
| FACTCHECK \mathcal{D}_{ca} | PHI3.5+GENVBLLMLP _{rw50} | VBL | 0.98 | 0.33 | 0.99 | 0.29 | 0.99 | 0.33 |
| FACTCHECK | PHI3.5+GENVBLLMLP _{rw50} | VBL | 0.29 | 0.06 | 0.98 | 0.24 | 0.80 | 0.02 |
| FACTCHECKSHUFFLED | PHI3.5+GENVBLLMLP _{rw50} | VBL | 0.22 | 0.21 | R | 0. | 1. | 0.04 |
| FACTCHECK \mathcal{D}_{ca} | MIXTRAL8x7B+DiscVBLLMLP | NO-REJECT | 0.91 | 0.50 | 0.93 | 0.50 | 0.93 | 0.49 |
| FACTCHECK | MIXTRAL8x7B+DiscVBLLMLP | NO-REJECT | 0.66 | 0.51 | 0.88 | 0.49 | 0.86 | 0.40 |
| FACTCHECKSHUFFLED | MIXTRAL8x7B+DiscVBLLMLP | NO-REJECT | 0.84 | 1. | R | 0. | 1. | 0.84 |
| FACTCHECK \mathcal{D}_{ca} | MIXTRAL8x7B+DiscVBLLMLP | VBL | 0.98 | 0.32 | 0.99 | 0.31 | 0.99 | 0.31 |
| FACTCHECK | MIXTRAL8x7B+DiscVBLLMLP | VBL | 0.85 | 0.08 | 0.97 | 0.27 | 0.89 | 0.08 |
| FACTCHECKSHUFFLED | MIXTRAL8x7B+DiscVBLLMLP | VBL | 0.91 | 0.14 | R | 0. | 1. | 0.13 |
| FACTCHECK \mathcal{D}_{ca} | MIXTRAL8x7B+DiscVBLLMLP _{rw50} | NO-REJECT | 0.90 | 0.50 | 0.94 | 0.50 | 0.94 | 0.48 |
| FACTCHECK | MIXTRAL8x7B+DiscVBLLMLP _{rw50} | NO-REJECT | 0.62 | 0.51 | 0.89 | 0.49 | 0.86 | 0.37 |
| FACTCHECKSHUFFLED | MIXTRAL8x7B+DiscVBLLMLP _{rw50} | NO-REJECT | 0.81 | 1. | R | 0. | 1. | 0.81 |
| FACTCHECK \mathcal{D}_{ca} | MIXTRAL8x7B+DiscVBLLMLP _{rw50} | VBL | 0.98 | 0.33 | 0.99 | 0.34 | 0.99 | 0.32 |
| FACTCHECK | MIXTRAL8x7B+DiscVBLLMLP _{rw50} | VBL | 0.72 | 0.12 | 0.97 | 0.30 | 0.91 | 0.09 |
| FACTCHECKSHUFFLED | MIXTRAL8x7B+DiscVBLLMLP _{rw50} | VBL | 0.82 | 0.16 | R | 0. | 1. | 0.13 |
| FACTCHECK \mathcal{D}_{ca} | MIXTRAL8x7B+GenVBLLMLP | NO-REJECT | 0.91 | 0.48 | 0.93 | 0.52 | 0.92 | 0.48 |
| FACTCHECK | MIXTRAL8x7B+GenVBLLMLP | NO-REJECT | 0.63 | 0.51 | 0.88 | 0.49 | 0.85 | 0.38 |
| FACTCHECKSHUFFLED | MIXTRAL8x7B+GenVBLLMLP | NO-REJECT | 0.73 | 1. | R | 0. | 1. | 0.73 |
| FACTCHECK \mathcal{D}_{ca} | MIXTRAL8x7B+GenVBLLMLP | VBL | 0.96 | 0.32 | 0.99 | 0.40 | 0.99 | 0.31 |
| FACTCHECK | MIXTRAL8x7B+GenVBLLMLP | VBL | 0.58 | 0.13 | 0.99 | 0.33 | 0.95 | 0.08 |
| FACTCHECKSHUFFLED | MIXTRAL8x7B+GenVBLLMLP | VBL | 0.54 | 0.11 | R | 0. | 1. | 0.06 |
| FACTCHECK \mathcal{D}_{ca} | MIXTRAL8x7B+GenVBLLMLP _{rw50} | NO-REJECT | 0.92 | 0.48 | 0.91 | 0.52 | 0.90 | 0.49 |
| FACTCHECK | MIXTRAL8x7B+GenVBLLMLP _{rw50} | NO-REJECT | 0.67 | 0.51 | 0.89 | 0.49 | 0.87 | 0.40 |
| FACTCHECKSHUFFLED | MIXTRAL8x7B+GenVBLLMLP _{rw50} | NO-REJECT | 0.69 | 1. | R | 0. | 1. | 0.69 |
| FACTCHECK \mathcal{D}_{ca} | MIXTRAL8x7B+GenVBLLMLP _{rw50} | VBL | 0.97 | 0.38 | 0.98 | 0.40 | 0.98 | 0.37 |
| FACTCHECK | MIXTRAL8x7B+GenVBLLMLP _{rw50} | VBL | 0.67 | 0.19 | 0.95 | 0.35 | 0.89 | 0.14 |
| FACTCHECKSHUFFLED | MIXTRAL8x7B+GenVBLLMLP _{rw50} | VBL | 0.67 | 0.15 | R | 0. | 1. | 0.10 |

References

- Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). In Yoshua Bengio and Yann LeCun (eds.), *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016. URL <http://arxiv.org/abs/1511.07289>.
- James Harrison, John Willes, and Jasper Snoek. Variational bayesian last layers. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=Sx7BIiPzys>.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.