

Linear Layer Extrapolation for Fine-Grained Emotion Classification

Mayukh Sharma
mayukh@ucsd.edu

Sean O'Brien
seobrien@ucsd.edu

Julian McAuley
jmcauley@ucsd.edu

Abstract

Certain abilities of Transformer-based language models consistently emerge in their later layers. Previous research has leveraged this phenomenon to improve factual accuracy through self-contrast, penalizing early-exit predictions based on the premise that later-layer updates are more factually reliable than earlier-layer associations. We observe a similar pattern for fine-grained emotion classification in text, demonstrating that self-contrast can enhance encoder-based text classifiers. Additionally, we reinterpret self-contrast as a form of linear extrapolation, which motivates a refined approach that dynamically adjusts the contrastive strength based on the selected intermediate layer. Experiments across multiple models and emotion classification datasets show that our method outperforms standard classification techniques in fine-grained emotion classification tasks.

1 Introduction

Despite the success of large language models on a variety of natural language processing (NLP) tasks (Brown et al., 2020; Wei et al., 2022), they still struggle with commonsense reasoning and factual recall (Fu et al., 2023; Wang et al., 2023), often hallucinating incorrect information (Ji et al., 2023).

Recently, *contrastive* methods, which maximize differences between a desirable "expert" and undesirable "amateur" model, have been proposed to address these issues (Li et al., 2022; Shi et al., 2023). In particular, decoding by contrasting layers, or *DoLa*, improves factuality by contrasting model outputs against early-exit predictions from intermediate layers of the same model. (Chuang et al., 2023) DoLa operates under the premise that later layers encode factuality, and thus late-emerging changes to predictions likely update towards more factual predictions. Recent work has demonstrated that intermediate layer features can also effectively quantify emotions in text (Sharma et al., 2023).

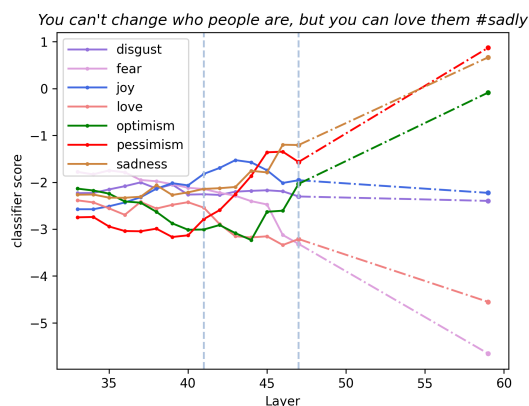


Figure 1: Linearly extrapolating class scores from amateur and expert layers to a nonexistent future layer correctly flips the output from sadness to pessimism.

Identifying emotions in text is crucial in NLP for applications ranging from detecting harmful behavior to enhancing conversational agents (Zhang et al., 2023; Barbieri et al., 2020). While many systems often focus on mutually exclusive emotions like joy or sadness, fine-grained emotions like grief and remorse are more nuanced and distinct. Many language-model systems still struggle to classify fine-grained emotions and opinions (Demszky et al., 2020; Zhang et al., 2023). Given this, we explore DoLa-style layer contrast to improve fine-grained emotion classification.

The main contributions of the paper are:

1. Demonstrating the merits of layer contrast on fine-grained emotion classification.
2. Recasting contrastive inference as linear extrapolation to obtain more stable performance with a dynamic contrastive penalty.

2 Related Work

Fine-grained Emotion Analysis: Much work has focused on identifying text sentiment (Rosenthal et al., 2017; Socher et al., 2013) and understanding

emotions in social media interactions (Mohammad et al., 2018; Chatterjee et al., 2019; Meaney et al., 2021). However, these efforts often focus on a limited set of emotions. Recent datasets on fine-grained emotion analysis (Demszky et al., 2020; Rashkin et al., 2019) indicate significant room for improvement in this area.

Early Exiting: Early-exiting predictions are obtained by applying the classification head of a model to the residual stream earlier in the network. These have been used to accelerate inference and dynamically allocate compute on a per-input basis (Teerapittayanon et al., 2016; Elbayad et al., 2020; Schuster et al., 2022).

Contrastive Steering: Contrastive methods optimize the difference in predictions between a favorable “expert” and an unfavorable “amateur,” to steer text decoding in language models. (Liu et al., 2021) GeDi (Krause et al., 2020) contrasts between class-specific control codes to improve text-conditioned factuality and emotion control. Coherence boosting (Malkin et al., 2021) provides the language model with only the final k tokens of the prompt to obtain amateur scores, encouraging longer-term coherence over locality. Contrastive decoding (Li et al., 2022; O’Brien and Lewis, 2023) improves long-form generation and reasoning ability by contrasting between large and small models of the same family. Other works use CD-like methods to reduce model toxicity, surface biases, and increase faithfulness to a provided context. (Liu et al., 2021; Yona et al., 2023; Shi et al., 2023)

3 Method

Here we define the main components of contrastive decoding (CD) and DoLa, alongside our proposed method for dynamically selecting contrastive strength. Following earlier work, we use early-exit distributions to choose an amateur layer, contrasting its predictions against the final layer (the expert). We apply mask candidate classes based on a plausibility constraint to filter out low-probability labels. We experiment with two methods to determine the contrastive strength: static β and dynamic β . The details of each component are discussed next.

3.1 Contrastive Classification

We use the formulation of contrastive decoding defined by O’Brien and Lewis (2023). Let p_a be the amateur probability scores and p_e be the ex-

pert probability scores. We define the contrastive classification function as:

$$f_{CC}^{(i)} = \begin{cases} (1 + \beta) \log p_e^i - \beta \log p_a^i & i \in \mathcal{V}_{valid} \\ -\infty & i \notin \mathcal{V}_{valid} \end{cases}$$

where β is the strength of the contrastive penalty and \mathcal{V}_{valid} is the adaptive plausibility constraint (Li et al., 2022), which defines the set of candidate classes on which contrastive action is applied. Let p_e^c be the expert probability for class $c \in C$. Then \mathcal{V}_{valid} is defined as:

$$\mathcal{V}_{valid} = \{ c \in C, p_e^c \geq \alpha \max_{c \in C} p_e^c \}$$

α here is a threshold hyperparameter that gates labels based on the probabilities assigned to them by the expert. This protects against instabilities associated with dividing the scores of two low-probability candidates, ensuring that all candidates are probable. After computing scores, $\arg \max_i f_{CC}^{(i)}$ is taken as the predicted label.

3.2 Dynamic premature layer selection

One central challenge with inference-time contrastive methods is the selection of a good amateur model. The amateur must be similar enough to the expert to model its error distribution, but not so powerful that it penalizes desirable behavior.

Contrasting against early-exiting predictions from earlier network layers provides an array of several potential amateurs to choose from. From a pre-validated set of earlier layers, DoLa selects the layer that has a maximally different early-exit distribution from the final-layer distribution, as measured by Jensen-Shannon Divergence. In short, the amateur layer ℓ_a is chosen as follows:

$$\ell_a = \arg \max_{\ell \in L_{valid}} d(\mathcal{P}(\ell), \mathcal{P}(\ell_{final}))$$

where L_{valid} is the pre-validated set of layers, \mathcal{P} maps a latent layer to its early-exited softmax distribution, and d is some divergence metric between two probability distributions. While the original paper uses Jensen-Shannon Divergence (JSD) as d , we find that cosine distance performs slightly better in practice.

3.3 Linear Layer Extrapolation

Casting DoLa as linear extrapolation allows us to dynamically vary β based on the chosen amateur layer.

Consider the classification of a single sample x to $c \in \mathcal{C}$, where $\mathcal{C} := \{1, 2, \dots, |\mathcal{C}|\}$ is the set of candidate classes. Suppose that the prediction is made by a network with L layers. For this sample, let $f_c(i)$ be the un-normalized score assigned by the model to class c by early-exiting at layer i . f_c is defined over the discrete set $\{1, \dots, L\}$, where the final score assigned to class c by the network is given by $f_c(L)$.

Let ℓ_a be the index of the selected early-exit amateur layer, $L > \ell_a$ be the index of the final model layer, and $\ell_t > L$ be the target post-final layer index to be linearly approximated. Note that ℓ_t need not be discrete.

Now let \hat{f}_c represent the linear function passing through $(\ell_a, f(\ell_a))$ and $(L, f(L))$.

$$\hat{f}_c(\ell) = f(L) + \left(\frac{f(L) - f(\ell_a)}{L - \ell_a} \right) (\ell - L)$$

We can use this function to linearly approximate performance at a target post-final layer index ℓ_t . This shares structure with the standard form of contrastive decoding, in which ℓ_t is implicit:

$$\hat{f}_c(\ell_t) = (1 + \beta)f(L) - \beta f(\ell_a)$$

We now define the relationship between the contrastive strength β and the conjectured target layer ℓ_t , obtaining:

$$\ell_t = -\beta\ell_a + (\beta + 1)L \quad (1)$$

$$\beta = \frac{\ell_t - L}{L - \ell_a} \quad (2)$$

DoLa keeps β fixed, which varies ℓ_t as a function of ℓ_a according to (1). We refer to this as *static β selection*.

We instead fix ℓ_t , implicitly varying β over different choices of ℓ_a according to (2). Selecting an amateur layer from earlier in the network will result in a decreased β , and a later amateur layer will increase β , in a process we call *dynamic β selection*.

After selecting hyperparameters k and ℓ_t , our method consists of the following steps:

1. **Early-Exit Calculation:** Apply projection head \mathcal{P} to the hidden activations at each layer in $\{\ell_k, \ell_{k+1}, \dots, L\}$.

2. **Amateur Identification:** Identify the amateur layer index ℓ_a with early-exit class probabilities that diverge the most from the final layer’s class probabilities.

3. **Dynamic β selection:** Calculate β according to Equation 2, incorporating the chosen ℓ_t .

4. **Logit Combination:** Calculate a plausibility mask and linearly combine the amateur and final layer scores according to subsection 3.1.

4 Experimental Setup

4.1 Datasets

goEmotions (Demszky et al., 2020) introduces a new emotion taxonomy of emotions named goEmotions, consisting of 28 emotions including neutral. The 27 emotion classes are fine-grained over 7 emotions defined in the Ekman taxonomy. goEmotions contains roughly 58k samples overcoming the problems with earlier emotion datasets which were small in size and covered a very limited taxonomy. We filter out the few multilabel data-points present in the dataset for our experiments.

SuperTweetEval (Antypas et al., 2023) aims to provide a unified benchmark to evaluate the performance of models on NLP tasks across social media. It is a heterogeneous collection of multiple datasets spanning NER, QA, and classification. For our experiments, we use tweetEmotion and tweetHate, which contain 12 and 8 classes respectively.

EmpatheticDialogues (Rashkin et al., 2019) was introduced as a benchmark for training and evaluating models and their capability to understand and acknowledge empathetic text. The dataset contains conversations distributed across 32 emotions. We use the first text of the conversation and the corresponding emotion for defining our fine-grained classification task.

4.2 Models and Training

We fine-tune FLAN-T5-L, FLAN-T5-XL, DeBERTa-L and DeBERTa-XL. (Chung et al., 2022; He et al., 2021). Fine-tuning details can be found in Appendix D. We release code, training and inference data. ¹

4.3 Decoding Hyperparameters:

Amateur layer: We use the dynamic amateur layer selection as defined in subsection 3.2, re-

¹<https://github.com/04mayukh/contrastive-classification>

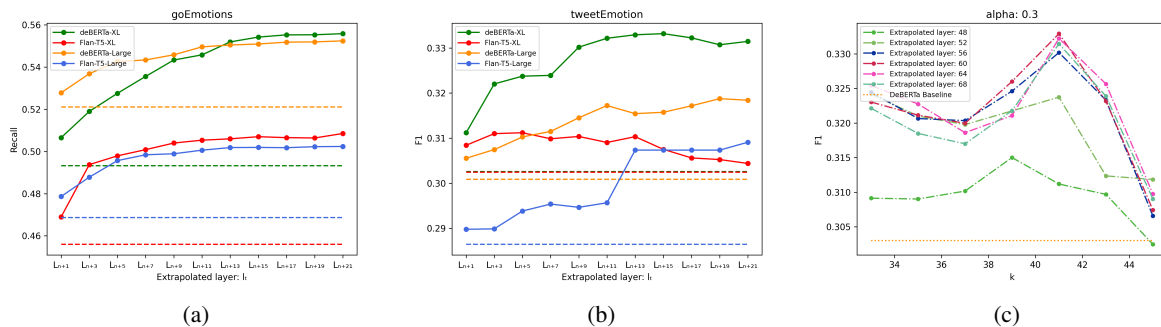


Figure 2: (a) Recall vs. ℓ_t on goEmotions; increasing the extrapolative strength improves recall. (b) F1 vs. ℓ_t on tweetEmotions exhibits a similar trend, though results for FLAN-T5-XL are mixed. (c) F1 vs. k for tweetEmotion using DeBERTa-XL; including layers 40 to 42 in the valid layers is found to be particularly useful.

stricting the amateur layer search space to a subset of the finetuned network layers. Let $L = \{\ell_k, \ell_{k+1}, \ell_{k+2}, \dots, L\}$ be a subset of the finetuned layers, where k is a hyperparameter defining the start of the search space and, L is the final layer of the network. We empirically choose k for each network based on performance on a held-out validation set. Results of the hyperparameter sweep can be found in Appendix A.

Contrastive Strength (β): We experiment with various fixed values of β between 0 to 1, finding that the best β varies over the selection of model and dataset. In general, values outside the range of $[0, 1]$ harmed performance.

Dynamic Contrastive Strength (ℓ_t): As with β , we sweep through a range of potential values for the post-final target layer, sweeping from L to $L + 25$.

5 Results

We report performance across models, datasets and inference algorithms in Table 1.

Traditional vs Contrastive Classification: We observe that contrastive classification improves performance significantly in both recall and F1. This trend holds for all models used in our experiments.

β vs Dynamic β : Dynamic β selection tends to improve the overall performance over static β for F1 and recall scores. Figure 2a shows the trend of recall scores across different models for dynamic β selection on the goEmotions dataset. Figure 2b shows the trend of F1 score across different models against dynamic β for the tweetEmotion dataset. Additionally, we observe that dynamic β is robust to changes in the hyperparameter k , which defines the start of the search space across earlier amateur layers. Figure 3 shows no clear or stable relationship between k and end performance when varying

β values. However, switching to dynamic β selection creates a constant trend with minor variance as k is varied, a trend that holds for multiple values of extrapolative layer t .

goEmotions: We see a general improvement across all models for the recall and F1 scores.

tweetEmotion: Contrastive classification with dynamic β significantly outperforms traditional classification. We see a general increase in recall and F1 with a slight harm to precision, and observe that the emotions corrected by layer contrast are highly correlated.

tweetHate: Layer contrast yields the largest improvement in performance on tweetHate across all models. This improvement owes in large part to improved performance on underrepresented classes, like disapproval and curiosity.

EmpatheticDialogue: For this dataset, we only see a slight increase in performance using the DeBERTa-xl model. Analyzing the probability distributions between layers, we did not observe a significant change in the probability distribution for different emotions across layers. The probability was distributed over a single label, increasing gradually across layers. This led to minimal contribution, positive or negative, from layer contrast.

Effect of amateur layer selection: We use a bucket of layers for amateur layer selection defined by hyperparameter k . Figure 2c shows the trend of k against F1 using the DeBERTa-xl for tweetEmotions dataset. We observe that the performance generally increases up to a layer where the benefit of contrastive action is maximum, followed by a drop in performance. Upon evaluating early-exiting on intermediate layers, we observed that some layers are more adept at identifying specific classes than others, providing a variety of skills to

Model	Type	EmpatheticDialogue			tweetHate		tweetEmotion			goEmotions			Avg. F1	
		Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall		F1
Flan-T5-L	\times	.551	.556	.543	.565	.577	.570	.298	.296	.286	.521	.469	.478	.469
Flan-T5-L	\checkmark	.551	.557	.543	.579	.606	.590	.300	.299	.291	.513	.485	.487	.478
Flan-T5-L	β	.551	.557	.543	.590	.636	.610	.349	.311	.309	.493	.502	.489	.488
Flan-T5-XL	\times	.582	.569	.565	.566	.566	.559	.320	.300	.302	.503	.456	.465	.473
Flan-T5-XL	\checkmark	.581	.570	.565	.690	.603	.615	.318	.314	.313	.499	.494	.486	.495
Flan-T5-XL	β	.582	.570	.565	.695	.605	.619	.316	.314	.313	.513	.494	.490	.497
DeBERTa-L	\times	.614	.601	.592	.647	.601	.622	.322	.299	.301	.570	.521	.534	.512
DeBERTa-L	\checkmark	.616	.606	.597	.676	.643	.658	.313	.311	.308	.562	.536	.540	.526
DeBERTa-L	β	.618	.609	.601	.708	.675	.690	.312	.331	.319	.558	.543	.541	.538
DeBERTa-XL	\times	.604	.605	.590	.607	.596	.599	.324	.300	.303	.529	.493	.502	.498
DeBERTa-XL	\checkmark	.610	.606	.594	.727	.668	.686	.335	.324	.325	.509	.530	.514	.523
DeBERTa-XL	β	.614	.609	.597	.725	.668	.685	.333	.340	.334	.505	.555	.522	.535

Table 1: Experimental results. \times , \checkmark , and β respectively represent normal classification, static β , and dynamic β .

contrast against for improved performance.

6 Analysis

We find that contrastive inference tends to flip neutral predictions to underrepresented classes. Table 2 shows the frequency of correctly flipped samples (true positives) versus correctly flipped samples (positives) from the neutral class (wrongly predicted as neutral). Table 5 contains the count of emotions that were correctly flipped from neutral.

Model	%Neutral
Flan-T5-large	76.2
Flan-T5-xl	79.2
DeBERTa-large	68.9
DeBERTa-xl	79.3

Table 2: Percentage of correctly flipped samples from the goEmotions dataset that were originally classified as neutral. This supports the trend of contrastive classification tending to predict more specific emotions, with an associated increase in recall.

We also report the most frequent samples corrected for the tweetEmotion dataset in Table 3. We see that the emotions for the pair of corrected samples were highly correlated.

7 Conclusion

We propose a linear extrapolation approach for dynamically determining contrastive strength in layerwise contrastive decoding. Applied to fine-grained emotion classification tasks, this method enhances classifier performance by effectively identifying under-represented classes. This strengthens the promise of layer-contrast methods in domains other than text generation, and provides a technical contribution that reduces the variance of the

Model	Emotion
Flan-T5-large	sadness \mapsto pessimism
	joy \mapsto anticipation
Flan-T5-xl	sadness \mapsto pessimism
	anger \mapsto disgust
deBERTa-large	anger \mapsto disgust
	joy \mapsto anticipation
deBERTa-xl	sadness \mapsto pessimism
	joy \mapsto optimism

Table 3: Top 2 emotion pairs that were contrastively flipped for each model on the tweetEmotion dataset. More straightforward emotions tend to be refined towards similar, more specific ones.

method with respect to a core hyperparameter k , encouraging further research into how best to exploit the layerwise emergence of textual understanding to improve performance on a wide range of tasks.

8 Limitations

Our study is restricted to fine-grained emotion classification with relatively small models (FLAN-T5 and DeBERTa). It remains to be seen whether our analysis of extrapolative classification will hold for prompt-based classification with larger models or across other datasets. We also found contrastive performance for smaller models to be sensitive to finetuning hyperparameters. Additionally, we observe that contrastive classification can fail when the model predictions do not change significantly over the course of a forward pass, as evidenced by results on EmpatheticDialogue (see Figure 4). Extending the method to identify and better handle these cases is left to future work.

References

- Dimosthenis Antypas, Asahi Ushio, Francesco Barbieri, Leonardo Neves, Kiamehr Rezaee, Luis Espinosa-Anke, Jiaxin Pei, and Jose Camacho-Collados. 2023. [SuperTweetEval: A challenging, unified and heterogeneous benchmark for social media NLP research](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 12590–12607, Singapore. Association for Computational Linguistics.
- Francesco Barbieri, José Camacho-Collados, Leonardo Neves, and Luis Espinosa-Anke. 2020. [Tweeteval: Unified benchmark and comparative evaluation for tweet classification](#). *ArXiv*, abs/2010.12421.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, T. J. Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). *ArXiv*, abs/2005.14165.
- Ankush Chatterjee, Kedhar Nath Narahari, Meghana Joshi, and Puneet Agrawal. 2019. [SemEval-2019 task 3: EmoContext contextual emotion detection in text](#). In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 39–48, Minneapolis, Minnesota, USA. Association for Computational Linguistics.
- Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James R. Glass, and Pengcheng He. 2023. [Dola: Decoding by contrasting layers improves factuality in large language models](#). *ArXiv*, abs/2309.03883.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. [Scaling instruction-finetuned language models](#).
- Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo Ko, Alan Cowen, Gaurav Nemade, and Sujith Ravi. 2020. [GoEmotions: A dataset of fine-grained emotions](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4040–4054, Online. Association for Computational Linguistics.
- Maha Elbayad, Jiatao Gu, Edouard Grave, and Michael Auli. 2020. [Depth-adaptive transformer](#). In *International Conference on Learning Representations*.
- Yao Fu, Litu Ou, Mingyu Chen, Yuhao Wan, Hao Peng, and Tushar Khot. 2023. [Chain-of-thought hub: A continuous effort to measure large language models’ reasoning performance](#).
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. [{DEBERTA}: {DECODING}-{enhanced} {bert} {with} {disentangled} {attention}](#). In *International Conference on Learning Representations*.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. [Survey of hallucination in natural language generation](#). *ACM Computing Surveys*, 55(12):1–38.
- Diederik P. Kingma and Jimmy Ba. 2014. [Adam: A method for stochastic optimization](#). *CoRR*, abs/1412.6980.
- Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq R. Joty, Richard Socher, and Nazneen Rajani. 2020. [Gedi: Generative discriminator guided sequence generation](#). In *Conference on Empirical Methods in Natural Language Processing*.
- Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. 2022. [Contrastive decoding: Open-ended text generation as optimization](#). In *Annual Meeting of the Association for Computational Linguistics*.
- Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. 2021. [Dexperts: Decoding-time controlled text generation with experts and anti-experts](#). In *Annual Meeting of the Association for Computational Linguistics*.
- Nikolay Malkin, Zhen Wang, and Nebojsa Jojic. 2021. [Coherence boosting: When your pretrained language model is not paying enough attention](#). In *Annual Meeting of the Association for Computational Linguistics*.
- J. A. Meaney, Steven Wilson, Luis Chiruzzo, Adam Lopez, and Walid Magdy. 2021. [SemEval 2021 task 7: HaHackathon, detecting and rating humor and offense](#). In *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)*, pages 105–119, Online. Association for Computational Linguistics.
- Saif Mohammad, Felipe Bravo-Marquez, Mohammad Salameh, and Svetlana Kiritchenko. 2018. [SemEval-2018 task 1: Affect in tweets](#). In *Proceedings of the 12th International Workshop on Semantic Evaluation*, pages 1–17, New Orleans, Louisiana. Association for Computational Linguistics.
- Sean O’Brien and Mike Lewis. 2023. [Contrastive decoding improves reasoning in large language models](#). *ArXiv*, abs/2309.09117.

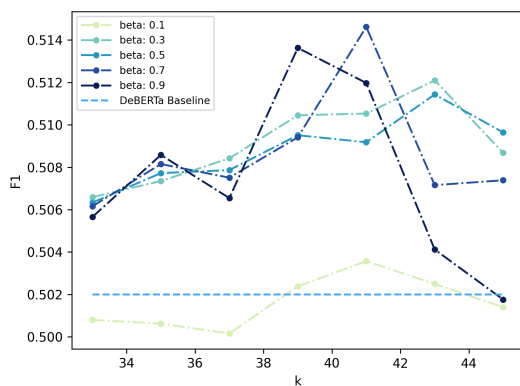
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. [Towards empathetic open-domain conversation models: A new benchmark and dataset](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5370–5381, Florence, Italy. Association for Computational Linguistics.
- Sara Rosenthal, Noura Farra, and Preslav Nakov. 2017. [SemEval-2017 task 4: Sentiment analysis in Twitter](#). In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 502–518, Vancouver, Canada. Association for Computational Linguistics.
- Tal Schuster, Adam Fisch, Jai Gupta, Mostafa Dehghani, Dara Bahri, Vinh Q. Tran, Yi Tay, and Donald Metzler. 2022. [Confident adaptive language modeling](#). In *Advances in Neural Information Processing Systems*.
- Mayukh Sharma, Ilanthenral Kandasamy, and W.B. Vasantha. 2023. [Emotion quantification and classification using the neutrosophic approach to deep learning](#). *Applied Soft Computing*, 148:110896.
- Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemoyer, and Scott Yih. 2023. [Trusting your evidence: Hallucinate less with context-aware decoding](#). *ArXiv*, abs/2305.14739.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. [Recursive deep models for semantic compositionality over a sentiment treebank](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.
- Surat Teerapittayanon, Bradley McDanel, and H.T. Kung. 2016. [Branchynet: Fast inference via early exiting from deep neural networks](#). In *2016 23rd International Conference on Pattern Recognition (ICPR)*, pages 2464–2469.
- Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xiangru Tang, Tianhang Zhang, Cheng Jiayang, Yunzhi Yao, Wenyang Gao, Xuming Hu, Zehan Qi, Yidong Wang, Linyi Yang, Jindong Wang, Xing Xie, Zheng Zhang, and Yue Zhang. 2023. [Survey on factuality in large language models: Knowledge, retrieval and domain-specificity](#). *ArXiv*, abs/2310.07521.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed Huai hsin Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. [Emergent abilities of large language models](#). *Trans. Mach. Learn. Res.*, 2022.
- G. Yona, Or Honovich, Itay Laish, and Roei Aharoni. 2023. [Surfacing biases in large language models using contrastive input decoding](#). *ArXiv*, abs/2305.07378.
- Wenxuan Zhang, Yue Deng, Bing Liu, Sinno Jialin Pan, and Lidong Bing. 2023. [Sentiment analysis in the era of large language models: A reality check](#).

A Hyperparameter Sweep

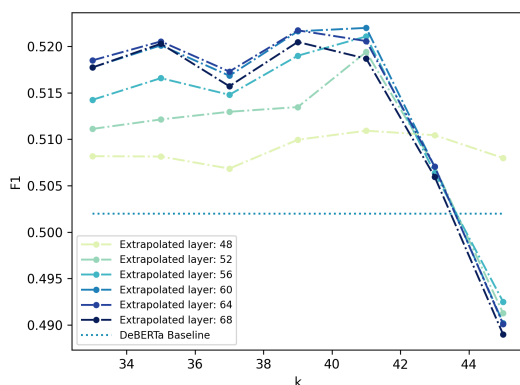
Table 4 contains the values of hyperparameter k used for reporting the results. We also show the effect of k on performance for the goEmotions dataset using both β and dynamic β in Figure 3.

Model	goEmotions	tweetEmotion	tweetHate	Empathetic Dialogue
Flan-T5-large	19	20	17	15
Flan-T5-xl	15	15	17	15
DeBERTa-large	15	19	17	19
DeBERTa-xl	39	41	38	43

Table 4: Our choice of hyperparameter k for defining the amateur search space used in the final results. The final layer is 48 for DeBERTa-xl and 23 for the remaining models.



(a) Static β



(b) Dynamic β selection

Figure 3: Effect of k , the earliest usable amateur layer on DeBERTa-XL goEmotions performance, under both static and dynamic β selection. Dynamic β increases performance and stabilizes performance with respect to the choice of k .

B Knowledge pattern across layers

Fine-grained emotion analysis is challenging due to non-mutually exclusive labels and similar polarity among different emotions, making it hard to accurately classify them. Class imbalance further biases the model towards more frequent emotions.

To study the change in probability distribution for emotions across layers, we performed early exiting on different layers of our fine-tuned models to visualize how the distributions across emotions evolve. We observed that for some emotions, the model makes a decision very early, passing it along the layers without much change. For others, the distribution tends to change in later layers, suggesting that the model is still adding information. We observed this pattern mostly around classes that are rarer in the training data or more closely related to each other. Figure 4 shows the change in distribution for two examples.

Drawing from these observations, we combine the idea of contrastive decoding and DoLa for fine-grained emotion analysis. We build on DoLa, using the early exited intermediate layers as amateur models. We then use contrastive action against the final layer distribution chosen as our expert model. Additionally, we deduce a method to dynamically select the contrastive strength which we show leads to better performance on fine-grained emotion tasks.

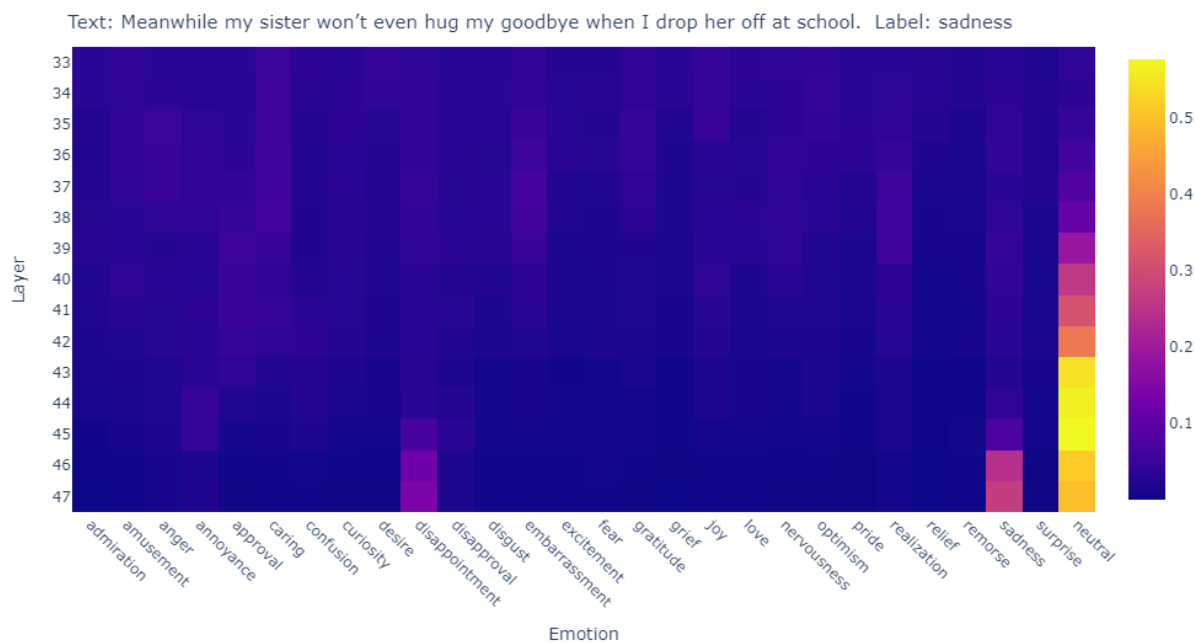
C Additional Analysis

From	To	Count
neutral	disapproval	22
neutral	curiosity	19
neutral	annoyance	13
neutral	admiration	12
neutral	approval	11

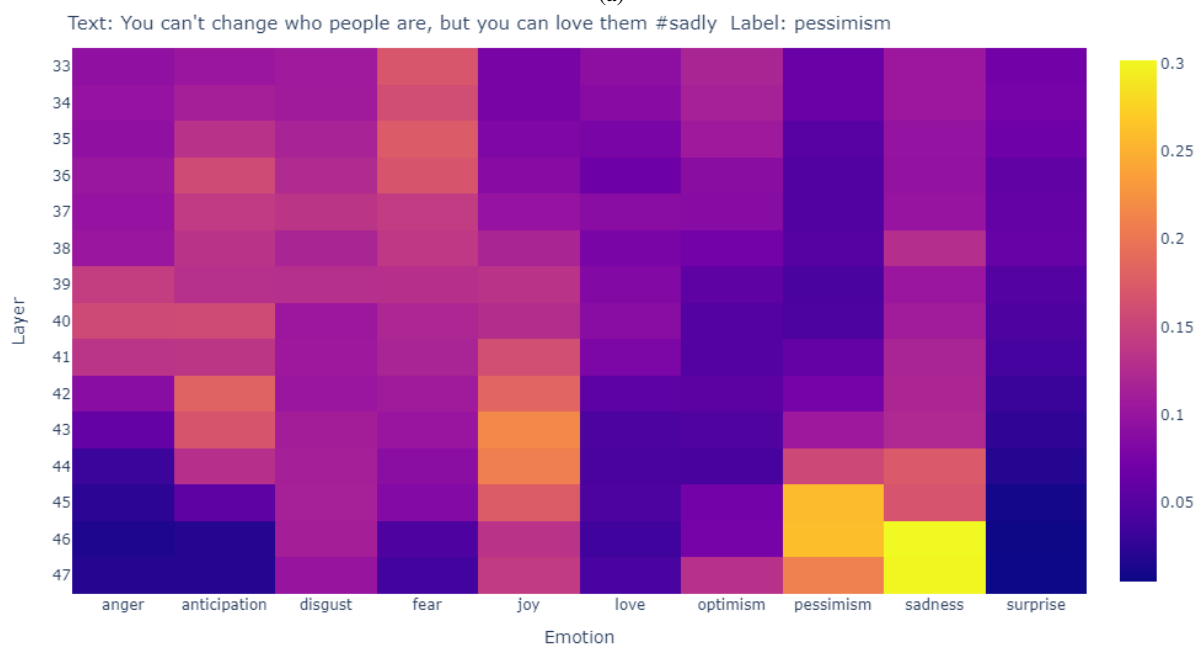
Table 5: Count of samples moved from neutral to other classes for goEMotions using DeBERTa-xl.

D Fine-Tuning Details

For DeBERTa-XL, we fine-tune layers 34-48 after freezing the initial layers (34/48); for DeBERTa-L, we fine-tuned all layers. For Flan-T5, we fine-tuned both large and xlarge variants after freezing the first (14/24) layers. We employed the Adam optimizer (Kingma and Ba, 2014) with learning rates ranging from $1e-6$ to $5e-6$ for DeBERTa and $1e-4$ to $5e-4$ for Flan-T5.



(a)



(b)

Figure 4: Probability distribution across the finetuned layers of DeBERTa-xl for a sample from each (a) goEmotions and (b) tweetEmotion dataset. In the goEmotions sample, the model initially identifies the label as neutral but increases the probabilities assigned to sadness and disappointment (the true label) over subsequent layers. For the tweetEmotion sample, the probability distribution changes across layers and the model fails to assign a high probability to a single emotion.

E Computational Resources Estimate

Early compute was run on freely available Cloud T4 GPUs. Fine-tuning and later experiments were run on a cluster of A6000 GPUs, with a maximum of 8 used at a single time.

Fine-tuning all models across all datasets takes

roughly 2 GPU-hours. Hyperparameter searches are performed at classification time, which takes very little compute. A very rough estimate for GPU-hours in this project is 50.